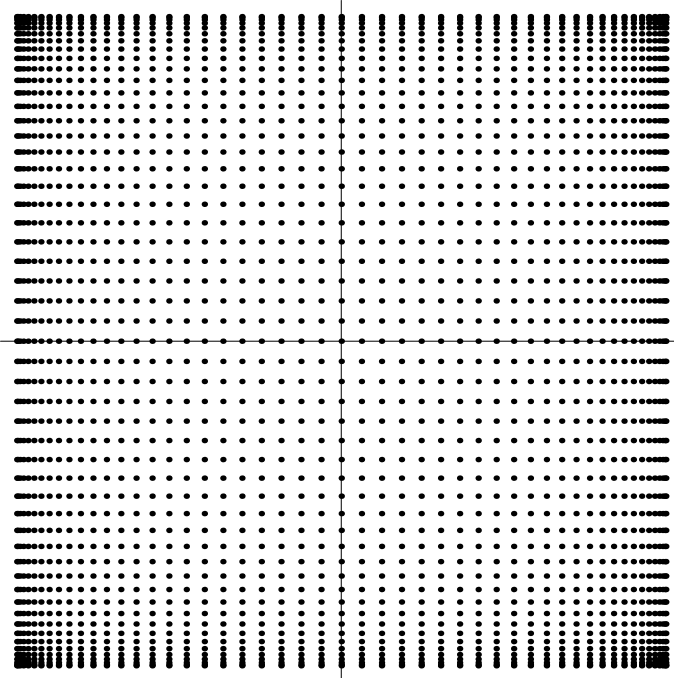


EWARYST RAFAJŁOWICZ

OPTYMALIZACJA
EKSPERYMENTU



Z ZASTOSOWANIAM I
W MONITOROWANIU
JAKOŚCI PRODUKCJI

EWARYST RAFAJŁOWICZ

**OPTYMALIZACJA
EKSPERYMENTU**

**z zastosowaniami
w monitorowaniu jakości produkcji**



OFICYNA WYDAWNICZA POLITECHNIKI WROCŁAWSKIEJ
WROCŁAW 2005

Recenzent
Dariusz Uciński

Opracowanie redakcyjne i korekta
Hanna Jurek

Skład komputerowy
Ewaryst Rafajłowicz

Projekt okładki
Ewaryst Rafajłowicz

© Copyright by Oficyna Wydawnicza Politechniki Wrocławskiej, Wrocław 2005

OFICYNA WYDAWNICZA POLITECHNIKI WROCŁAWSKIEJ
50-370 Wrocław, Wybrzeże Wyspiańskiego 27

ISBN 83-7085-913-5

Drukarnia Oficyny Wydawniczej Politechniki Wrocławskiej. Zam. 2005.

Spis treści

Wstęp	5
I. Klasyczne zagadnienia optymalizacji eksperymentu	
1. Modele liniowe i ich estymacja	9
1.1. Model regresji liniowej	9
1.2. Metoda najmniejszych kwadratów	14
1.3. Dokładność estymatora	16
2. Plan eksperymentu	20
2.1. Jakościowe wymagania stawiane eksperymentom	20
2.2. Ilościowa ocena jakości planów eksperymentu	22
2.3. Definicje planów eksperymentu	26
3. Geometria zbioru macierzy informacyjnych	30
3.1. Elementarne własności macierzy informacyjnych	30
3.2. Klasa realizowalnych macierzy informacyjnych	31
4. Optymalne plany eksperymentu	35
4.1. Dalsze uwagi na temat oceny jakości planu	35
4.2. Własności kryteriów D- i G-optymalności	38
4.3. Równoważność planów D- i G-optymalnych	40
4.4. Uogólnienie twierdzenia Kiefera i Wolfowitza	42
4.5. Wybrane optymalne plany eksperymentów	43
II. Plany optymalne dla modeli wielowymiarowych	
5. Numeryczne poszukiwanie planów optymalnych	51
5.1. Ogólny algorytm gradientowy dla planów D-optymalnych	51
5.2. Algorytm selektywnych poszukiwań losowych	56
5.3. Optymalna alokacja pomiarów	67
6. Plany dla modeli o zmiennych zblokowanych	72
6.1. Iloczyn planów zależnych od zblokowanych zmiennych	72
6.2. Planowanie dla modeli addytywnych względem zblokowanych zmiennych	74
6.3. Plany dla modeli z pełnym zestawem interakcji zblokowanych zmiennych	78
6.4. Analityczne wyznaczanie planów	85
7. Zalety planów produktowych	93
7.1. Poszukiwanie i realizacja planów produktowych	93
7.2. Zyski obliczeniowe w estymacji z zastosowaniem planów produktowych .	98

III. Eksperyment w testowaniu jakości wyrobów

8.	Diagnozowanie i poprawa odporności wyrobów	107
8.1.	Bezpośredni model odporności wyrobów na warunki eksploatacji	107
8.2.	Plany produktowe w modelu Taguchi	109
9.	Plany o minimalnym koszcie i zadanej jakości	115
9.1.	Minimalizacja kosztu przy zadanej macierzy informacyjnej	116
9.2.	Charakteryzacje planów optymalnych	119
9.3.	Czy klasyczne eksperymenty są planami o minimalnym koszcie?	123
10.	Sekwencje planów nadążających za zmianami otoczenia	128
10.1.	Sformułowanie problemu	128
10.2.	Warunki optymalności sekwencji planów	133
10.3.	Zastosowanie – sterowanie ruchomymi czujnikami	139
11.	Pokrewne zadania planowania	144
11.1.	Eksperyment w zadaniach estymacji – wybrane aspekty	144
11.2.	Składniki eksperymentu w estymacji systemów dynamicznych	147
IV. Eksperyment w diagnostyce procesów		
12.	Próbkowanie funkcyjnych charakterystyk wyrobów	151
12.1.	Ocena jakości funkcyjnej charakterystyki wyrobu	151
12.2.	Dobór planu i pomiary charakterystyki	153
13.	Próbkowanie obrazów do celów diagnostycznych	156
13.1.	Próbkowanie zmienności obrazu. Szybki algorytm okonturowywania	157
13.2.	Próbkowanie wymiaru fraktalnego obrazu jako wskaźnika diagnostycznego	165
14.	Wykrywanie zmian jakości w sekwencjach obrazów	173
14.1.	Nieparametryczna karta kontrolna	173
14.2.	Dostrajanie karty i wyniki porównań	181
14.3.	Metodyka wykrywania zmian w sekwencjach obrazów	189
V. Dodatek		
15.	Iloczyn Kroneckera i jego własności	197
15.1.	Definicja i podstawowe własności	197
15.2.	Wartości własne iloczynu Kroneckera macierzy	198
Literatura		200
Skorowidz		211

Wstęp

W prezentowanej Czytelnikowi monografii omawiane są problemy optymalizacji eksperymentów, których celem jest zebranie danych dla estymacji funkcji regresji. Przedstawiono też propozycje potencjalnych zastosowań technik planowania eksperymentu i próbkowania pól losowych (obrazów) w wybranych zagadnieniach monitorowania jakości procesów wytwórczych.

Monografia podzielona została na cztery części. W części I zebrano podstawowe fakty z matematycznej teorii planowania eksperymentów optymalnych. Rezultaty te znaleźć można w kilku książkach (por. [3], [35], [110], [36], [138], [202]). Powyżej i w całej monografii stosujemy konwencję przytaczania numerów cytowanych pozycji w kolejności ich związków z omawianym zagadnieniem.

W części II zebrano rezultaty badań autora na temat poszukiwania planów optymalnych w zagadnieniach wielowymiarowych. Algorytm selektywnych poszukiwań losowych przedstawiono na tle znanego algorytmu Wynna-Fedorova, który był inspiracją dla autora przy konstruowaniu algorytmu omawianego w rozdziale 5. W pozostałych dwóch rozdziałach części II omawiane są rezultaty, uzyskane przez autora i dr. Wojciecha Myszkę, na temat możliwości komponowania planów optymalnych dla zagadnień wielowymiarowych na podstawie planów dla modeli o mniejszej liczbie zmiennych. Przedstawiono także algorytm obliczeniowy wykorzystujący strukturę takich planów.

Części III i IV poświęcone są takim zagadnieniom planowania eksperymentu i próbkowania, które mogą mieć zastosowania w zagadnieniach projektowania wyrobów oraz monitorowania i oceny jakości w procesach wytwórczych. Rezultaty zawarte w tych częściach albo nie były dotąd publikowane, albo zawierają nowe spojrzenie i uogólnienia wcześniejszych wyników autora. Badania przedstawione w częściach III i IV finansowane były z grantu KBN Nr 4 T11A 025 23 oraz z subsydium Fundacji na Rzecz Nauki Polskiej.

Niniejsza monografia przeznaczona jest dla osób zajmujących się problematyką planowania eksperymentu oraz dla tych Czytelników, którzy zechcą podjąć próbę stosowania technik planowania eksperymentu w swoich pracach badawczych lub wdrożeniowych. Można mieć nadzieję, że książka ta, a zwłaszcza jej część I, będzie przydatna także dla słuchaczy studiów doktoranckich.

Autor wyraża serdeczne podziękowania profesorowi Dariuszowi Ucińskiemu, recenzentowi tej monografii, za wiele cennych uwag oraz Wojciechowi Rafajłowiczowi za pomoc w programowaniu badań symulacyjnych karty kontrolnej.

Wrocław, 15 listopada 2005

CZĘŚĆ I

**Klasyczne zagadnienia
optymalizacji eksperymentu**

1. Modele liniowe i ich estymacja

Termin „model liniowy” odnosi się jedynie do liniowości względem nieznanych parametrów modelu, natomiast wpływ innych wielkości może nie być liniowy.

1.1. Model regresji liniowej

Przyjmijmy, że na podstawie zebranych wcześniej informacji potrafimy wskazać pewne wielkości fizyczne, które wpływają na przebieg procesu. Będziemy zakładać, że potrafimy mierzyć ich wartości. Wyniki jednokrotnego pomiaru s wielkości oddziałujących na proces oznaczają będziemy przez

$$x = [x^{(1)}, x^{(2)}, \dots, x^{(s)}]^T$$

i nazywać wielkościami wejściowymi lub krótko – wejściami; $x \in R^s$, powyżej T oznacza transpozycję.

Pomiary wejść w kolejnych eksperymentach oznaczają będziemy przez

$$x_i, \quad i = 1, 2, \dots, N.$$

Zakładamy, że wybrano pewne wielkości charakteryzujące przebieg procesu, nazywane dalej wielkościami wyjściowymi (lub krótko – wyjściami) oraz, że potrafimy je mierzyć. Dalej zakładamy, że interesuje nas zależność jednej tylko spośród wielkości wyjściowych od x . Wartości tego wybranego wyjścia oznaczają będziemy przez y lub

$$y_1, y_2, \dots, y_N,$$

jeśli mamy zestaw pomiarów.

Problem doboru wielkości wejściowych i wyjściowych pozostaje poza zakresem tej monografii. Jest on zresztą dość rzadko dyskutowany w literaturze. Pewne rezultaty na ten temat zawarto w pracy [143].

Przyjęte przez nas założenie o oddzielnym rozpatrywaniu zależności poszczególnych wyjść od x jest głęboko zakorzenione w literaturze. Warto jednak zauważyć, że nie powinno być ono przyjmowane całkiem bezkrytycznie (szerzej temat ten omówiono w [138]).

Problem konstrukcji modelu empirycznego badanego zjawiska można w ogólnym zarysie sformułować następująco. Dysponując zestawem pomiarów wejść i odpowiadającego im wyjścia

$$(x_i, y_i), \quad i = 1, 2, \dots, N$$

chcemy znaleźć pewien opis, matematyczny lub algorytmiczny, zależności y od x . Opis ten powinien spełniać następujące wymagania:

- Być prostszy w stosowaniu niż same obserwacje, a jednocześnie dostatecznie dokładny w okolicy punktów, gdzie dokonywano pomiarów.
- Zapewniać możliwości predykcji wartości wyjścia dla takich wartości wejść, dla których nie dokonywano pomiarów.
- Zachowywać cechy jakościowe badanego procesu. Przykładowo, powinien zachowywać monotoniczność zależności y od x , o ile cechę tę ma badany proces.

To ostatnie wymaganie jawnie formułowane jest dopiero w ostatnich latach.

W praktyce rzadko występuje sytuacja, w której otrzymuje się tę samą wartość y każdorazowo, gdy na wejściu pojawi się ustalona wartość x . Gdy własność ta zachodzi w całym zakresie interesujących nas wartości x , naturalne jest poszukiwanie funkcyjnej zależności y od x .

Częściej spotykamy się z przypadkiem, gdy każdorazowe wystąpienie na wejściu tego samego x daje w wyniku inną wartość mierzonego wyjścia y .

Zakładać będziemy, że te wartości y -ków są niezależnymi obserwacjami pewnej zmiennej losowej o rozkładzie prawdopodobieństwa zależnym od x .

Najpełniejszą formą opisu zachowania y -ków jest dystrybuanta rozkładu wyjścia przy ustalonym x . Jeśli jest ona różniczkowalna dla każdego z x -ów, to posłużyć się można rodziną gęstości rozkładów y -ków sparametryzowaną przez x . Gęstość tę oznaczamy będziemy jako $f(y; x)$. Średnik w tym oznaczeniu zastosowany został po to, by odróżnić je od warunkowej gęstości y względem x .

Posługiwanie się opisem w postaci $f(y; x)$ nie jest zbyt łatwe, a często nie jest nawet możliwe, gdyż funkcja ta nie jest znana, a do jej estymacji mamy zwykle zbyt mało danych. Pewnym uproszczeniem tego opisu są zależności funkcyjne postaci $\tilde{y} = \bar{y}(x)$, gdzie \bar{y} dobieramy tak, by \tilde{y} dobrze, w wybranym przez nas sensie, opisywało zmienność y -ków przy ustalonym x . Jedną z najczęściej używanych form opisu w postaci zależności funkcyjnej jest

$$\tilde{y} = \bar{y}^*(x) = \int_{-\infty}^{\infty} y f(y; x) dy, \quad (1.1.1)$$

czyli wartość oczekiwana y liczona dla każdego ustalonego x . Łatwo sprawdzić, że \bar{y}^* minimalizuje, względem $\psi \in R$, następujący wskaźnik jakości:

$$q(\psi; x) = \int_{-\infty}^{\infty} (y - \psi)^2 f(y; x) dy, \quad (1.1.2)$$

będący średniokwadratowym błędem przybliżenia dla ustalonego x .

Przypuśćmy, że posiadamy obserwacje (x_i, y_i) , $i = 1, 2, \dots, N$, przy czym y_i ma rozkład o gęstości $f(y; x_i)$, wówczas można je przedstawić w postaci:

$$y_i = \bar{y}^*(x_i) + \epsilon_i, \quad i = 1, 2, \dots, N, \quad (1.1.3)$$

gdzie ϵ_i są pewnymi zmiennymi losowymi o średniej zero. Można je interpretować jako losowe zakłócenia pomiaru $\bar{y}^*(x_i)$.

Zauważmy, że zwykle $\bar{y}^*(x)$ nie jest znane, gdyż nie znamy rozkładu y dla poszczególnych x .

Estymacją $\bar{y}^*(x)$ w sytuacji, gdy posiadamy mało informacji o tej funkcji, zajmuje się teoria nieparametrycznej estymacji funkcji regresji.

Gdy liczba posiadanych lub planowanych obserwacji nie jest dostatecznie duża, pozostaje nam możliwość uzupełnienia procedury estymacji o posiadaną (lub założoną hipotetycznie) wiedzę aprioryczną o postaci funkcyjnej zależności y od wejść x .

Aż do odwołania, obowiązywać będą następujące założenia.

1. Zależność funkcyjna służąca przybliżaniu obserwacji ma często postać kombinacji liniowej wybranych funkcji z nieznanymi współczynnikami. Gdy na obiekt eksperymentu nie oddziałują losowe zakłócenia, to zależność obserwacji wyjścia od wejść x jest postaci:

$$\bar{y}(x, a) = a^T v(x) = \sum_{k=1}^r a^{(k)} v^{(k)}(x), \quad (1.1.4)$$

gdzie T oznacza transpozycję, natomiast:

$a = [a^{(1)}, a^{(2)}, \dots, a^{(r)}]^T$ jest wektorem nieznanymi parametrów,

$v(x) = [v^{(1)}(x), v^{(2)}(x), \dots, v^{(r)}(x)]^T$ – to wektor znanych funkcji, zadawanych przez nas na podstawie wiedzy o badanym zjawisku.

2. O funkcjach $v^{(1)}(x), v^{(2)}(x), \dots, v^{(r)}(x)$ zakładamy, że są one liniowo niezależne w pewnym obszarze $X \subset R^s$, z którego pochodzą obserwacje x_i .
3. Zakłócenia ϵ_i , $i = 1, 2, \dots, N$ są zmiennymi losowymi o wartości oczekiwanej zero i skończonych wariancjach. Ponadto dla $i \neq j$ zmienne losowe ϵ_i oraz ϵ_j są nieskorelowane, $i, j = 1, 2, \dots, N$. Zakładamy też, że zakłócenia te oddziałują addytywnie na obiekt badań. Dostępne pomiary (x_i, y_i) związane są zależnością:

$$y_i = a^T v(x_i) + \epsilon_i, \quad i = 1, 2, \dots, N, \quad (1.1.5)$$

dla pewnego wektora parametrów $a \in R^r$. Składowe tego wektora traktowane są jak „prawdziwe” wartości nieznanymi parametrów. Równanie (1.1.5) zapisać można następująco:

$$\mathbf{y} = V_N^T \cdot a + \epsilon, \quad (1.1.6)$$

gdzie $\mathbf{y} = [y_1, y_2, \dots, y_N]^T$, $\epsilon = [\epsilon_1, \epsilon_2, \dots, \epsilon_N]^T$ są wektorami kolumnowymi, natomiast $V_N = [v(x_1), v(x_2), \dots, v(x_N)]$ jest macierzą o N kolumnach i r wierszach.

4. O wariancji zakłóceń przyjmować będziemy jedno z następujących założeń.
- Wariancje $\text{var}(\epsilon_i) = \sigma^2$, $i = 1, 2, \dots, N$ są jednakowe, a wartość σ nie jest znana.
 - Wariancje $\text{var}(\epsilon_i) = \sigma^2 w_i$ mogą być różne, przy czym $\sigma > 0$ nie jest znane, natomiast ciąg $w_i > 0$, $i = 1, 2, \dots, N$ jest znany. W tym przypadku przyjmujemy:

$$\text{var}(\epsilon_i) = \sigma^2(x_i) = \sigma^2 w_i, \quad (1.1.7)$$

gdzie $\sigma^2(x)$ jest pewną znaną funkcją opisującą względną dokładność obserwacji w poszczególnych punktach x . W zadaniach planowania eksperymentu wygodnie będzie posługiwać się funkcją $w(x) = \sigma^2(x)/\sigma^2$, o której będziemy zakładać, że jest znana. Zachodzi oczywiście zależność: $w(x_i) = w_i$, $i = 1, 2, \dots, N$.

5. Wartości x_i , $i = 1, 2, \dots, N$ są znane dokładnie, to znaczy bez błędów pomiarowych, niezależnie od tego czy pochodzą z obserwacji, czy też są wynikiem świadomie zaplanowanego eksperymentu.
6. Punkty x_i , $i = 1, 2, \dots, N$ rozmieszczone są tak, że

$$\text{rząd } [v(x_1), v(x_2), \dots, v(x_N)] = r, \quad (1.1.8)$$

gdzie r oznacza liczbę nieznanymi parametrów a i równocześnie liczbę elementów wektora $v(x)$.

A oto kilka uwag o powyższych założeniach.

- Macierz $V_N = [v(x_1), v(x_2), \dots, v(x_N)]$ ma wymiary $r \times n$. Warunkiem koniecznym dla (1.1.8) jest $n \geq r$, co oznacza, że liczba pomiarów nie może być mniejsza od liczby estymowanych na ich podstawie parametrów.
- Warto zwrócić uwagę, że postać warunku (1.1.8) nie zależy od jakichkolwiek założeń o zakłóceniami. Jeśli warunek ten nie jest spełniony, to wskazać można takie zestawy parametrów, powiedzmy, a_1 i a_2 , dla których $V_N^T a_1 = V_N^T a_2$, mimo że $a_1 \neq a_2$. Oznacza to nierozróżnialność a_1 i a_2 nawet wówczas, gdyby obserwacje dokonywane były bez zakłóceń.
- Zależność (1.1.4) nazywa się zwykle liniową. Termin ten odnosi się do liniowej zależności \bar{y} od parametrów a . Nie wymagamy, by \bar{y} zależała liniowo od x .
- Liniowa niezależność składowych $v(x)$ oznacza, że żadna z funkcji $v^{(i)}$ nie może być przedstawiona jako liniowa kombinacja pozostałych funkcji tworzących wektor $v(x)$. Liniowa zależność między składowymi $v(x)$ może prowadzić do osobliwości macierzy układu równań normalnych, z którym mamy do czynienia, szacując wektor a metodą najmniejszych kwadratów (por. rozdz. 1.2).

- W literaturze rozważa się również ogólniejszy przypadek zakłóceń skorelowanych (por. [145]). Zakłada się wówczas, że znana jest dodatnio określona macierz kowariancji zakłóceń, powiedzmy W . Oznaczmy przez $\mathcal{V} = W^{1/2}$ (por. [77], gdzie znaleźć można definicję podnoszenia macierzy do potęg ułamkowych). Dokonajmy teraz zamiany zmiennych: $\bar{\mathbf{y}} = \mathcal{V}^{-1} \cdot \mathbf{y}$, $\bar{V}_N^T = \mathcal{V}^{-1} \cdot V_N^T$, $\bar{\epsilon} = \mathcal{V}^{-1} \cdot \epsilon$, to otrzymamy model $\bar{\mathbf{y}} = \bar{V}_N^T \cdot a + \bar{\epsilon}$, w którym zakłócenia są nieskorelowane.
- W ramach ogólnej teorii modeli liniowych dopuszcza się niespełnienie założenia 6 (por. [145]). Bada się wówczas możliwości estymacji liniowych kombinacji składowych wektora a .
- W rozdziale tym przyjmujemy, że punkty x_i , $i = 1, 2, \dots, N$, w których dokonywane były pomiary, mogły pochodzić z zaplanowanego wcześniej eksperymentu lub były wynikiem biernej obserwacji (np. reprezentują chwile, w których rejestrowano y_i).

Definicja 1.1. Składowe wektora $v(x)$ nazywać będziemy funkcjami rozpinającymi regresję lub model liniowy.

Definicja 1.2. Model regresji liniowej oznaczać będziemy przez $(X, v(x), \sigma^2(x))$. Na opis modelu składają się:

- obszar jego określoności X ,
- zestaw funkcji $v(x)$, rozpinających model $\bar{y} = a^T v(x)$,
- funkcja (lub stała), opisująca wariancję zakłóceń.

Termin model używany będzie także w znaczeniu węższym, obejmującym tylko $\bar{y} = a^T v(x)$ i domyślnie obszar X , bez specyfikowania opisu wariancji zakłóceń.

Poniżej podano przykłady często stosowanych zestawów $v(x)$.

1. Model wielomianowy: $v^{(k)}(x) = x^{k-1}$, $k = 1, 2, \dots, r$.
2. Model trygonometryczny: $v^{(1)}(x) = 1$, $v^{(2k)}(x) = \sin(2kx)$, $v^{(2k-1)}(x) = \cos((2k-1)x)$, $k = 1, 2, \dots$.
3. Model liniowy: dla $x \in R^s$ $v^{(1)}(x) = 1$, $v^{(k+1)}(x) = x^{(k)}$, $k = 1, 2, \dots, s$.
4. Wielomiany Bernsteina: $v^{(k)}(x) = \binom{N}{k} x^k (1-x)^{N-k}$, $x \in [0, 1]$, $k = 0, 1, \dots, N$. Z twierdzenia Bernsteina o jednostajnej, wielomianowej aproksymacji wynika, że dowolnie dokładne przybliżenie funkcji $f \in C([0, 1])$ uzyskamy, wybierając dostatecznie duże N i kładąc $a_k = f(k/N)$ we wzorze

$$\bar{y}(x, a) = \sum_{k=0}^N a^{(k)} v^{(k)}(x). \quad (1.1.9)$$

5. Przypuśćmy, że chcemy uzyskać model, który jest w stanie opisać ekstremum względem zmiennych $x^{(1)}$, $x^{(2)}$, co wymaga użycia modelu stopnia co najmniej

drugiego. Jednocześnie, dla każdej wartości $x^{(1)}$, $x^{(2)}$ model ma być funkcją liniową względem $x^{(3)}$. Wymagania te prowadzą do przyjęcia modelu o postaci:

$$\begin{aligned}\bar{y}(x, a) = & a^{(1)} + a^{(2)}x^{(1)} + a^{(3)}x^{(2)} + a^{(4)}x^{(1)}x^{(2)} + a^{(5)}(x^{(1)})^2 \\ & + a^{(6)}(x^{(2)})^2 + a^{(7)}x^{(3)} + a^{(8)}x^{(1)}x^{(3)} + a^{(9)}x^{(2)}x^{(3)} \\ & + a^{(10)}x^{(1)}x^{(2)}x^{(3)} + a^{(11)}(x^{(1)})^2x^{(3)} + a^{(12)}(x^{(2)})^2x^{(3)}.\end{aligned}$$

W modelu tym wektor rozpinający ma postać:

$$v(x) = \left[1, x^{(1)}, x^{(2)}, x^{(1)}x^{(2)}, (x^{(1)})^2, (x^{(2)})^2, x^{(3)}, x^{(1)}x^{(3)}, x^{(2)}x^{(3)}, \right. \\ \left. x^{(1)}x^{(2)}x^{(3)}, (x^{(1)})^2x^{(3)}, (x^{(2)})^2x^{(3)} \right]^T,$$

a wektor a ma dwanaście elementów.

1.2. Metoda najmniejszych kwadratów

Metoda najmniejszych kwadratów należy do klasycznych metod opracowywania obserwacji, a jej bibliografia liczy setki pozycji (por. [62], [145], [172] i cytowaną tam literaturę). Ograniczymy się do zestawienia tylko podstawowych faktów, mających wpływ na sformułowanie zagadnień planowania eksperymentu.

Niech spełnione będą założenia 1, 3, 5, 6 oraz założenie 4a). Przyjmijmy, że zebraliśmy pomiary (x_i, y_i) , $i = 1, 2, \dots, N$.

Definicja 1.3. Wektor \hat{a} , który minimalizuje błąd średniokwadratowy:

$$Q(a) = \sum_{i=1}^N \left[y_i - a^T v(x_i) \right]^2 \quad (1.2.10)$$

względem wszystkich $a \in R^s$ nazywać będziemy estymatorem otrzymanym metodą najmniejszych kwadratów.

Właściwie należałoby mówić: *estymator uzyskany metodą minimalizacji sumy kwadratów błędów*, lecz przyjęło się używanie jeszcze krótszego terminu: *estymator MNK*.

Zauważmy, że sposób oceny błędów w (1.2.10) bierze pod uwagę jedynie różnice między obserwacjami odpowiedzi obiektu i modelu liniowego. Stosując takie podejście, przyjmować musimy założenie 5.

Jeśli zamiast założenia 4a) spełniona jest jego wersja b), to naturalne jest przypisanie różnych wag do błędów $\left[y_i - a^T v(x_i) \right]^2$, w zależności od dokładności

obserwacji w danym punkcie. Prowadzi to do minimalizacji, względem $a \in R^s$, ważonej sumy kwadratów danej wzorem:

$$Q_u(a) = \sum_{i=1}^N u_i \left[y_i - a^T v(x_i) \right]^2, \quad (1.2.11)$$

gdzie $u_i \geq 0$ oznaczają wagi. Naturalne jest przyjęcie

$$u_i = \frac{1}{w(x_i)}, \quad i = 1, 2, \dots, N, \quad (1.2.12)$$

gdyż przypisuje się małe wagi punktom, w których wariancja zakłóceń jest relatywnie duża.

Jeśli założymy na chwilę, że rozkład prawdopodobieństwa każdego spośród ϵ_i , $i = 1, 2, \dots, N$ jest $N(0, \sigma^2 w(x_i))$ (rozkład normalny o średniej zero i wariancji $\sigma^2 w(x_i)$), to przyjęcie (1.2.11) jako wskaźnika jakości estymacji można uzasadnić na gruncie metody największej wiarygodności, zakładając nieskorelowanie ϵ_i oraz ϵ_j dla $i \neq j$. W dalszej części rozdziału nie zakładamy normalności rozkładu zakłóceń.

Łatwo sprawdzić, że jeśli spełnione jest założenie 6 oraz $u_i > 0$, $i = 1, 2, \dots, N$, to funkcja $Q_u(a)$, będąca formą kwadratową wektora a , jest ściśle wypukła. Dlatego znalezienie jej minimum sprowadza się do rozwiązania układu równań, powstałego z przyrównania gradientu $Q_u(a)$ do zera. Powstaje w ten sposób układ równań liniowych względem a , opisanych w następującej definicji.

Definicja 1.4. *Następujący układ równań liniowych względem $a \in R^r$ nazywa się układem równań normalnych:*

$$M_N a = \sum_{i=1}^N u_i y_i v(x_i), \quad (1.2.13)$$

gdzie $r \times r$ macierz M_N dana jest wzorem:

$$M_N = \sum_{i=1}^N u_i v(x_i) v^T(x_i). \quad (1.2.14)$$

Uwaga 1.1. *Dla macierzy układu równań normalnych M_N stosujemy to samo oznaczenie, które używane będzie w dalszych rozdziałach dla macierzy informacyjnej, gdyż różnią się one jedynie stałym mnożnikiem. Tam gdzie mnożnik ten nie jest istotny, również (1.2.14) nazywać będziemy macierzą informacyjną.*

Można wykazać, że układ równań (1.2.13) zawsze ma co najmniej jedno rozwiązanie. Jest tak nawet wtedy, gdy macierz M_N jest osobliwa. Jeśli wagi u_i są niezerowe, to założenie 6 gwarantuje nam, że rozwiązanie układu równań normalnych jest jedyne.

Twierdzenie 1.1. *Jeśli spełnione jest założenie 6, to \hat{a} określone wzorem*

$$\hat{a} = M_N^{-1} b_N, \quad b_N = \sum_{i=1}^N u_i y_i v(x_i) \quad (1.2.15)$$

jest estymatorem MNK, to znaczy minimalizuje błąd średniokwadratowy. Ponadto \hat{a} jest liniową funkcją obserwacji y_1, y_2, \dots, y_N .

Zwracamy uwagę, że ta postać estymatora MNK służy do celów teoretycznych, natomiast w praktycznych zastosowaniach rozwiązać należy układ równań normalnych (1.2.13). Jest to liniowy układ równań algebraicznych, a liczba niewiadomych $r = \dim(a)$ nie przekracza zwykle kilkudziesięciu. Mimo to, rozwiązanie go w sposób numerycznie poprawny nie zawsze jest łatwe [71].

Powyższe twierdzenie pokazuje postać estymatora MNK parametrów funkcji regresji. Jeśli natomiast celem naszym jest estymacja wartości tejże funkcji w wybranych punktach, to jako estymator MNK wartości funkcji regresji w punkcie $x \in X$ wybieramy

$$\hat{y}(x) = \hat{a}^T v(x). \quad (1.2.16)$$

W następnych rozdziałach podamy uzasadnienie teoretyczne takiego właśnie sposobu estymacji wartości funkcji regresji.

Jeżeli wszystkie wagi u_i są jednakowe, to

$$\hat{a} = M_N^{-1} b_N, \quad b_N = \sum_{i=1}^N y_i v(x_i), \quad M_N = \sum_{i=1}^N v(x_i) v^T(x_i). \quad (1.2.17)$$

1.3. Dokładność estymatora

Własności estymatorów MNK należą do klasycznych rezultatów i dlatego przedstawimy je krótko.

Własność 1. *Niech spełnione będą założenia 1–6. Estymator \hat{a} parametrów modelu liniowego jest estymatorem nieobciążonym, tzn., $\mathbf{E}\hat{a} = a$, niezależnie od tego, jaką wartość przyjmuje wektor nieznanych parametrów a .*

Nieobciążoność estymatora zinterpretować można jako brak systematycznego – w sensie statystycznym – błędu oceny parametrów, niezależnie od tego, jakie są wartości tych parametrów.

Własność 2. Niech spełnione będą założenia 1, 2, 3, 4b), 5, 6. Ponadto, zakładamy, że wagi dobrane są w sposób uzgodniony z wariancjami zakłóceń, tzn. $u_i = 1/(\sigma^2 w(x_i))$. Wówczas macierz kowariancji estymatora \hat{a}

$$\text{cov}(\hat{a}) = \mathbf{E} \left[(\hat{a} - \mathbf{E}\hat{a}) (\hat{a} - \mathbf{E}\hat{a})^T \right]$$

ma postać:

$$\text{cov}(\hat{a}) = \sigma^2 \left[\sum_{i=1}^N w^{-1}(x_i) v(x_i) v^T(x_i) \right]^{-1}. \quad (1.3.18)$$

Ze wzoru $\text{cov}(\hat{a}) = \mathbf{E} \left[(\hat{a} - \mathbf{E}\hat{a}) \cdot (\hat{a} - \mathbf{E}\hat{a})^T \right]$ wynikają natychmiast następujące własności macierzy kowariancji:

- Elementy diagonalne tej macierzy są równe wariancom ocen poszczególnych elementów wektora a .
- Macierz $\text{cov}(\hat{a})$ jest symetryczna.
- Jej elementy leżące pod i ponad przekątną są proporcjonalne do współczynników korelacji między ocenami poszczególnych parametrów (współczynnikami proporcjonalności są iloczyny dyspersji ocen tychże parametrów).

Przyjęliśmy, że funkcja $a^T v(x)$ estymowana będzie przez $\hat{y}(x) \stackrel{\text{def}}{=} \hat{a}^T v(x)$. Własności 1 oraz 2 pozwalają uzyskać natychmiast informacje na temat dokładności oceny $a^T v(x)$ przez $\hat{a}^T v(x)$.

Własność 3. Niech spełnione będą założenia własności 2. Wówczas $\hat{y}(x)$ jest nieobciążonym estymatorem $a^T v(x)$.

$$\begin{aligned} \text{var}(\hat{y}(x)) &= v^T(x) \text{cov}(\hat{a}) v(x) \\ &= \sigma^2 v^T(x) \left[V_N W_N^{-1} V_N^T \right]^{-1} v(x) \end{aligned} \quad (1.3.19)$$

gdzie $r \times N$ macierz ma postać $V_N = [v(x_1), v(x_2), \dots, v(x_N)]$, natomiast W_N jest $N \times N$ macierzą diagonalną o elementach na przekątnej równych $w(x_i)$, $i = 1, 2, \dots, N$.

Definicja 1.5. Ciągami reszt (ang. residuals) nazywamy

$$\hat{\epsilon}_i = y_i - \hat{y}(x_i), \quad i = 1, 2, \dots, N. \quad (1.3.20)$$

Ciąg ten ma dużą wartość diagnostyczną dla oceny dokładności, gdyż daje ocenę wartości zakłóceń i dlatego w jego oznaczeniu użyto symbolu $\hat{\epsilon}$. Stosuje się również unormowane oceny reszt (por. [2]). Wobec nieobciążoności $\hat{y}(x_i)$ jako estymatora $a^T v(x_i)$, nietrudno zauważyć, że $\hat{\epsilon}_i$ jest nieobciążonym estymatorem dla ϵ_i , to znaczy, $\mathbf{E} \hat{\epsilon}_i = 0$, $i = 1, 2, \dots, N$.

Powyższe wzory uproszczą się wtedy, gdy $w(x_i) = 1$, $i = 1, 2, \dots, N$. Rzeczywiście, wówczas

$$\text{cov}(\hat{a}) = \sigma^2 \left(V_N V_N^T \right)^{-1}, \quad \text{cov}(\hat{y}(x)) = \sigma^2 v^T(x) \left(V_N V_N^T \right)^{-1} v(x),$$

natomiast ciąg reszt pozwala oszacować wariancję zakłóceń następująco:

$$\hat{\sigma}^2 = (N - r)^{-1} \sum_{i=1}^N \hat{\epsilon}_i^2. \quad (1.3.21)$$

Można wykazać, że $\hat{\sigma}^2$ jest nieobciążonym estymatorem σ^2 .

Własność 4. *Jeśli dodatkowo założymy, że zakłócenia pomiarowe ϵ_i mają rozkład normalny $\mathcal{N}(0, \sigma^2)$, to:*

- estymator parametrów \hat{a} ma rozkład $\mathcal{N} \left(a, \sigma^2 \left(V_N V_N^T \right)^{-1} \right)$,
- $\hat{y}(x)$ ma rozkład $\mathcal{N} \left(a^T v(x), \text{cov}(\hat{y}(x)) \right)$.

Własności te są podstawą testów służących do badania hipotez o zerowości pewnych parametrów regresji oraz procedur doboru struktury regresji i testowania jej adekwatności. Procedury te i testy należą do klasycznych rezultatów analizy regresji i są szeroko opisanie w literaturze (por. [62], [172], [145]).

Jak już wspomniano, estymator MNK jest estymatorem liniowym, czyli liniową funkcją obserwacji y -ków. Gdy oznaczymy przez L macierz o r wierszach i N kolumnach, to $L y$ jest ogólną postacią estymatorów liniowych. Estymator \hat{a} jest też estymatorem nieobciążonym, więc sensowne jest porównywanie go z innymi estymatorami liniowymi i nieobciążonymi parametrów a w modelu regresji. Estymatory te będziemy oznaczać skrótem ELN.

Niech A i B będą dwiema kwadratowymi $r \times r$ macierzami nieujemnie określonymi. Jeśli macierz $A - B$ jest nieujemnie określona, to będziemy pisać $A \geq B$ lub $A - B \geq 0$.

Nie każde dwie nieujemnie określone macierze o tych samych wymiarach są porównywalne w sensie powyższej definicji. Można jednak udowodnić, że macierze kowariancji estymatorów z klasy ELN są zawsze porównywalne w sensie powyższego określenia. Jeśli $L_1 y$ i $L_2 y$ są dwoma estymatorami z klasy ELN, to przyjmujemy, że estymator $L_1 y$ jest nie gorszy niż $L_2 y$, gdy $\text{cov}(L_1 y) \leq \text{cov}(L_2 y)$. Wówczas sens ma pytanie o estymator posiadający „najmniejszą” macierz kowariancji. Odpowiedź na to pytanie zawarta jest w następującym twierdzeniu.

Twierdzenie 1.2. (Gaussa–Markova) *Niech spełnione będą założenia wniosku 2.*

- Estymator MNK jest estymatorem najlepszym w klasie ELN, w sensie jaki wyznacza uporządkowanie macierzy kowariancji.

- Wśród wszystkich liniowych i nieobciążonych estymatorów wartości funkcji $a^T v(x)$ estymator $\hat{y}(x)$ ma minimalną wariancję.
- Jeśli dodatkowo rozkład zakłóceń ϵ_i , $i = 1, 2, \dots, N$ jest normalny $N(0, \sigma)$, to wymienione własności zachodzą w klasie wszystkich (nie tylko liniowych) estymatorów nieobciążonych.

Dowód tego twierdzenia w powyższej wersji znaleźć można w [17], str. 362 (por. także [145]). Warto zauważyć, że twierdzenie to zachodzi dla dowolnego skończonego N .

Założenie o porównywaniu estymatora MNK wyłącznie z estymatorami liniowymi i nieobciążonymi jest bardzo istotne. Jeśli dopuścimy do konkurencji również estymatory obciążone, to sama macierz kowariancji nie w pełni oddaje popełniane błędy. Gdy zdecydujemy się porównywać błędy średniokwadratowe $\mathbf{E} \left[(\hat{a} - a)^T (\hat{a} - a) \right]$, to okaże się, że zarówno estymator grzebieniowy (ang. *ridge regression*), jak i estymator Jamesa–Steina o postaci $\gamma \hat{a}$ mogą dawać mniejszy błąd średniokwadratowy pod warunkiem, że ich parametry t i γ są właściwie dobrane. Zwracamy uwagę, że ten właściwy dobór w praktyce polega na uzależnieniu t i γ od \mathbf{y} , co czyni te estymatory nieliniowymi (por. [40]).

2. Plan eksperymentu

W rozdziale tym dokonamy przeglądu jakościowych i ilościowych wymagań stawianych planom eksperymentów. Na razie, pod pojęciem planu eksperymentu rozumiemy zestaw wejść $\xi^{(N)} = \{x_1, x_2, \dots, x_N\}$, wybieranych z obszaru X , dla których przeprowadzane będą pomiary reakcji (wyjścia) badanego obiektu.

2.1. Jakościowe wymagania stawiane eksperymentom

Przytoczymy zestaw wymagań jakościowych, które stawiane bywają eksperymentom, stanowi jednocześnie słownik podstawowych terminów tej dziedziny. Nie trzeba dodawać, że nie istnieją plany zdolne sprostać wszystkim wymienionym wymaganiom.

ORTOGONALNOŚĆ PLANU. Plan nazywać będziemy ortogonalnym, gdy kolumny macierzy pomiarów V_N są względem siebie ortogonalne jako wektory w przestrzeni R^r . Plany te mają wiele zalet z punktu widzenia statystycznego. W szczególności, pominięcie w modelu pewnych członów nie powoduje konieczności przeliczania oszacowań pozostałych jego parametrów, o ile tylko pomiary wykonywane były zgodnie z planem ortogonalnym dla tego modelu. Dla wielu modeli i wskaźników jakości planowania, plany ortogonalne są optymalne.

Termin „ortogonalność planu” nie jest precyzyjny, gdyż faktycznie wymaganie to odnosi się do macierzy informacyjnej, która z kolei zależy zarówno od modelu, jak i od planu. Termin ten przyjął się jednak w literaturze i dlatego dalej też będzie używany.

WALORY NUMERYCZNE. Od planu eksperymentu wymagać możemy tego, by jego użycie prowadziło do redukcji błędów numerycznych powstających przy obliczaniu oszacowań parametrów modelu (np. wówczas, gdy do obliczeń stosujemy metodę najmniejszych kwadratów).

SYMETRIA OBROTOWA PLANU. Symetria obrotowa planu to wymaganie stałości wariancji oszacowania wyjścia modelu w stałej odległości od centrum planu. Centrum planu to punkt w przestrzeni wejść, w którego otoczeniu tworzony jest model matematyczny procesu. Symetria obrotowa planu ma zatem na celu zapewnienie, by dokładność oszacowania wartości wyjścia modelu nie preferowała żadnego kierunku. Dalsze rezultaty na temat planów o symetrii obrotowej i ich znaczenia znaleźć można w [62], [188]. Dawniej w literaturze polskiej plany te nazywano planami rotatabilnymi.

OPTIMALNOŚĆ PLANU. Wymaganie optymalności planu oznacza, że przyjęty został pewien wskaźnik mierzący jakość (np. dokładność oszacowania parametrów modelu) różnych planów, a plan optymalny to taki, który zapewnia największą możliwą do osiągnięcia w danych warunkach wartość tego wskaźnika.

PLANY UWZGLĘDNIAJĄCE NIEPRAWIDŁOWĄ SPECYFIKACJĘ MODELU. Badania planów optymalnych opierają się na założeniu, że struktura modelu jest znana i poprawna. Wiedzę tę musimy mieć jeszcze przed doświadczeniem. Od dość dawna zdawano sobie sprawę z ograniczającej roli tego założenia [70], [78], [95], lecz pierwsze istotne wyniki dotyczące odporności planów produktowych uzyskano niedawno [167], [165].

KOSZT I CZAS TRWANIA EKSPERYMENTU. Minimalizacja kosztów i czasu eksperymentu są wymaganiami równie oczywistymi, co trudnymi do spełnienia. Zwykle wymaganie to uwzględnia się tylko pośrednio, dążąc do minimalizacji liczby przeprowadzanych eksperymentów. W [93] zawarte są rezultaty na temat doboru liczby pomiarów w sytuacji, gdy uwzględnia się ich koszt. Inne podejście omówimy w rozdziale 9.

PLANY UWZGLĘDNIAJĄCE KORELACJĘ ZAKŁÓCEŃ. Założenie o braku korelacji zakłóceń między kolejnymi pomiarami nawet wtedy, gdy pomiary wykonywane są w tym samym punkcie przestrzeni eksperymentu, jest jednym z podstawowych warunków nakładanych w całej klasycznej teorii planowania eksperymentu. Będziemy je przyjmować także w tej książce, gdyż wydaje się, że w wielu problemach praktycznych jest ono do utrzymania, zwłaszcza tam, gdzie źródłem błędów losowych są zakłócenia pomiarowe. Zagadnieniom planowania w sytuacji, gdy przyjmuje się założenie o skorelowaniu zakłóceń poświęcona jest relatywnie bogata literatura [157], [158], [159], [182], [183], [181], [13], [14], [18], [30], [33], [50], [94], [161].

PLANY WYKORZYSTUJĄCE INFORMACJĘ A PRIORI O PARAMETRACH. Planowanie takie mieści się w obszarze tzw. podejścia bayesowskiego. Jego zastosowanie wymaga założenia, że nieznanne wartości parametrów modelu są wynikiem losowania ze znanego nam rozkładu prawdopodobieństwa. Gdy informacja taka jest dostępna, to do planowania eksperymentu zastosować można podejścia omawiane w [106], [107].

PLANOWANIE SEKWENCYJNE. Przy czynionych wcześniej założeniach możliwe było jednorazowe rozdysponowanie wszystkich obserwacji, jakie można w danych warunkach wykonać. W sytuacji, gdy wyjście modelu zależy nieliniowo od nieznanymi parametrów, celowe jest etapowe podejmowanie decyzji o alokacji eksperymentów. W pracach [35], [171], [166] zawarte są interesujące rezultaty na temat sekwencyjnego planowania, w którym decyzje o rozmieszczeniu kolejnych obserwacji podejmuje się na podstawie już zdobytych informacji.

PLANY UWZGLĘDNIAJĄCE CZYNNIKI ILOŚCIOWE I JAKOŚCIOWE. Zmienne wejściowe mogą mieć charakter ilościowy lub jakościowy. Zadania planowania

z czynnikami obu rodzajów są od dłuższego czasu badane, por. [76], [204], [205], [163], [5]. Specyfika tego rodzaju planowania polega na tym, że wpływ czynników jakościowych uwzględnić można poprzez dopuszczenie zmian parametrów modelu, w zależności od wartości czynników jakościowych.

2.2. Ilościowa ocena jakości planów eksperymentu

Rozważmy dwa plany $\xi_1^{(N)} = \{x'_1, x'_2, \dots, x'_N\}$, $\xi_2^{(N)} = \{x''_1, x''_2, \dots, x''_N\}$, przeznaczone do estymacji parametrów tego samego modelu. Chcemy ocenić, który z nich jest lepszy z punktu widzenia dokładności estymacji. Poprzednio w twierdzeniu Gaussa–Markowa porównywaliśmy liniowe sposoby oceny parametrów modelu liniowego za pomocą relacji \leq między macierzami kowariancji ocen parametrów, które w pełni opisują jakość liniowych nieobciążonych estymatorów. Twierdzenie Gaussa–Markowa daje nam dobre podstawy do porównywania planów i planowania eksperymentu, gdyż zapewnia ono, że na drodze ulepszania metody przetwarzania wyników pomiarów nie uzyskamy poprawy dokładności, jeśli na metodę przetwarzania nakładamy wymóg liniowości i nieobciążoności odpowiadającego jej estymatora. Dalszej poprawy możemy szukać na drodze ekstensywnej lub intensywnej. Przez drogę ekstensywną rozumiemy tu wykonanie jednej lub kilku czynności:

1. *Zmniejszenie wariancji zakłóceń $\sigma^2(x)$.* Zabieg ten wymaga albo zmiany urządzeń pomiarowych, albo wielokrotnego powtarzania tego samego pomiaru i uśredniania wyników.
2. *Zwiększenie liczby pomiarów N .* Wiąże się to ze wzrostem kosztów i/lub czasu eksperymentu.
3. *Zwiększenie obszaru eksperymentu X lub zmiana jego kształtu.* Zwiększenie obszaru eksperymentu wymaga zwykle nakładów na urządzenia wykonawcze realizujące wymuszenia. Ponadto, zwiększając nadmiernie obszar planowania, ryzykujemy, że badany model przestanie być adekwatny.

Jeśli zrealizowane zostały wszystkie powyższe proste sposoby zwiększania dokładności estymacji, to pozostaje jedynie droga intensywna, polegająca na doborze planu eksperymentu $\xi^{(N)}$, przy założeniu, że N i X są ustalone, a wariancja pomiarów nie może ulec zmniejszeniu.

Jeśli z góry przyjmiemy, że do estymacji użyjemy metody najmniejszych kwadratów, to do porównywania planów skorzystać będzie można z macierzy kowariancji $\text{cov}(\hat{a})$ lub macierzy informacyjnej, będącej jej odwrotnością. Powstaje więc problem porównywania ze sobą macierzy $\text{cov}(\hat{a})$ odpowiadających różnym planom. Macierze takie na ogół nie są ze sobą porównywalne w sensie wprowadzonej w poprzednim rozdziale relacji \geq między macierzami symetrycznymi. Z tego

względem pozostaje nam posługiwanie się skalarnymi¹ funkcjami macierzy $\text{cov}(\hat{a})$ (lub macierzy informacyjnej) tak wybranymi, by miały interpretację statystyczną.

Zakładamy, że spełnione są założenia rozdziału 1.1, a oszacowania parametrów \hat{a} obliczono metodą najmniejszych kwadratów. W celach interpretacyjnych zakładamy, ale tylko w tym podrozdziale, normalność rozkładu zakłóceń. Przy tych założeniach udowodnić można, że elipsoida ufności dla oszacowań parametrów regresji ma postać:

$$(a - \hat{a})^T M_N(\xi^{(N)})(a - \hat{a}) \leq c, \quad a \in R^r, \quad (2.2.1)$$

gdzie $c > 0$ jest stałą zależną od liczby obserwacji, liczby estymowanych parametrów, poziomu ufności $0 < \beta < 1$ i oszacowania wariancji zakłóceń. Stała ta nie zależy natomiast ani od a , ani od $\xi^{(N)}$. Centrum elipsoidy postaci (2.2.1) położone jest w punkcie \hat{a} . Elipsoida ta pokrywa wektor nieznanych parametrów z prawdopodobieństwem $0 < \beta < 1$.

Długości poszczególnych osi elipsoidy (2.2.1) równe są $2/\sqrt{\lambda_i(M_N(\xi^{(N)}))}$, gdzie $\lambda_i(M_N(\xi^{(N)}))$, $i = 1, 2, \dots, r$ oznaczają wartości własne macierzy $M_N(\xi^{(N)})$.

Interpretacja sposobów porównywania jakości planów jest bardziej przejrzysta, gdy wyrazimy długości osi elipsoidy ufności w równoważnej postaci:

$$2/\sqrt{\lambda_i(M_N(\xi^{(N)}))} = 2\sqrt{\lambda_i(M_N^{-1}(\xi^{(N)}))},$$

czyli w terminach wartości własnych macierzy kowariancji ocen parametrów.

Interesować nas będzie wpływ doboru planu na kształt elipsoidy ufności. Intuicyjnie jest jasne, że ten z planów jest lepszy (zapewnia większą dokładność estymacji parametrów), którego elipsoida ufności „jest mniejsza”.

Zmiany planu wpływają mogą zarówno na zmianę orientacji głównej osi elipsoidy, jak i na proporcje między długościami poszczególnych osi. Dlatego nie ma jednego sposobu przypisania elipsoidom – a poprzez nie także planom – wartości liczbowych, określających ich użyteczność w zadaniu estymacji parametrów.

Dalej przytaczamy najczęściej stosowane sposoby porządkowania planów na podstawie różnych sposobów mierzenia rozmiarów elipsoidy ufności. Używać będziemy następującej symboliki: jeśli plan $\xi_1^{(N)}$ nie jest lepszy od planu $\xi_2^{(N)}$, w sensie określonym na następnych stronach, to będziemy pisać $\xi_1^{(N)} \ll \xi_2^{(N)}$. Określenie, że plan *nie jest lepszy* oznacza, że jest on równie dobry lub gorszy w sensie wartości wybranego wskaźnika jakości.

¹ Przegląd podejść do definiowania planów „uniwersalnie” optymalnych znaleźć można w [105], [110]. Okazało się jednak, że plany takie istnieją jedynie w bardzo specyficznych i rzadkich przypadkach.

OCENA NA PODSTAWIE OBJĘTOŚCI ELIPSOIDY UFNOŚCI. Objętość elipsoidy (2.2.1) jest proporcjonalna do $[\det M_N^{-1}(\xi^{(N)})]^{1/2}$. Dlatego plan $\xi_1^{(N)}$ uznajemy za nie lepszy od planu $\xi_2^{(N)}$ na podstawie wyznacznika macierzy kowariancji, gdy

$$\xi_1^{(N)} \stackrel{D}{\ll} \xi_2^{(N)} \Leftrightarrow \det [M_N^{-1}(\xi_1^{(N)})] \geq \det [M_N^{-1}(\xi_2^{(N)})]. \quad (2.2.2)$$

Litera D w symbolu $\stackrel{D}{\ll}$ nawiązuje do angielskiej nazwy wyznacznika *determinant*.

OCENA NA PODSTAWIE ŚREDNIEJ DŁUGOŚCI OSI ELIPSOIDY UFNOŚCI. Średnia długość osi elipsoidy (2.2.1) jest proporcjonalna do śladu macierzy $M_N^{-1}(\xi^{(N)})$. Ślad macierzy oznaczamy przez $\text{tr}[\cdot]$. Jeśli będziemy oceniać plany na tej podstawie, to

$$\xi_1^{(N)} \stackrel{A}{\ll} \xi_2^{(N)} \Leftrightarrow \text{tr} [M_N^{-1}(\xi_1^{(N)})] \geq \text{tr} [M_N^{-1}(\xi_2^{(N)})]. \quad (2.2.3)$$

Litera A w symbolu $\stackrel{A}{\ll}$ nawiązuje do angielskiego słowa *average*.

OCENA NA PODSTAWIE ŚREDNIEJ p -TEGO RZĘDU. Ta metoda oceny planów jest naturalnym uogólnieniem poprzedniej i polega na następującym sposobie liczenia średniej długości osi elipsoidy ufności:

$$\xi_1^{(N)} \stackrel{L_p}{\ll} \xi_2^{(N)} \Leftrightarrow \left\{ \text{tr} [M_N^{-p}(\xi_1^{(N)})] \right\}^{1/p} \geq \left\{ \text{tr} [M_N^{-p}(\xi_2^{(N)})] \right\}^{1/p}, \quad (2.2.4)$$

gdzie $p > 0$ jest wybranym parametrem. Związki tej metody oceny jakości planów z innymi sposobami ich uporządkowania omawiamy dalej (bezpośrednio po podaniu Definicji 4.5).

OCENA NA PODSTAWIE LINIOWO WAŻONEJ ŚREDNIEJ. Innym uogólnieniem oceny planów na podstawie średniej długości osi elipsoidy ufności jest liniowo ważona średnia

$$\xi_1^{(N)} \stackrel{L}{\ll} \xi_2^{(N)} \Leftrightarrow \text{tr} [AM_N^{-1}(\xi_1^{(N)})] \geq \text{tr} [AM_N^{-1}(\xi_2^{(N)})], \quad (2.2.5)$$

gdzie A jest nieujemnie określoną macierzą wybieraną przez eksperymentatora. Warto nadmienić, że ten sposób ważenia uwzględniać może także pozadiagonalne elementy macierzy kowariancji ocen.

OCENA NA PODSTAWIE MAKSYMALNEJ DŁUGOŚCI OSI ELIPSOIDY UFNOŚCI. Ocena ta mierzy „wielkość” elipsoidy ufności za pomocą maksymalnej długości osi tej elipsoidy, co prowadzi do następującego uporządkowania planów:

$$\xi_1^{(N)} \stackrel{E}{\ll} \xi_2^{(N)} \Leftrightarrow \max_i \sqrt{\lambda_i(M_N^{-1}(\xi_1^{(N)}))} \geq \max_i \sqrt{\lambda_i(M_N^{-1}(\xi_2^{(N)}))}. \quad (2.2.6)$$

Litera E pochodzi od angielskiego *eigenvalue*. Wobec monotoniczności funkcji \sqrt{t} dla $t > 0$, można pominąć pierwiastki w poprzednim wzorze – wpisano je ze względów interpretacyjnych.

Wymienione sposoby porównania jakości planów brały pod uwagę dokładność oszacowań parametrów regresji. Alternatywnym, lecz w pewnych ważnych przypadkach równoważnym, spojrzeniem jest porównanie planów pod względem dokładności szacowania wartości funkcji regresji (wyjścia modelu). Spojrzenie takie jest ważne wtedy, gdy regresja używana jest dla dokonania predykcji lub oszacowania wartości wyjścia pomiędzy punktami, w których dokonywano pomiarów. Zgodnie z własnością 3, jeśli dopuszczamy tylko nieobciążone estymatory wyjścia, to dokładność estymacji regresji w punkcie x mierzymy za pomocą wariancji $\text{var}(\hat{y}(x))$. Aby podkreślić zależność wariancji od planu, wprowadzimy oznaczenie $\phi_N(x, \xi_N) = \text{var}(\hat{y}(x))$, a następujący wzór jawnie pokazuje tę zależność

$$\varphi_N(x, \xi^{(N)}) = \sigma^2 v^T(x) M_N^{-1}(\xi^{(N)}) v(x). \quad (2.2.7)$$

Porównywanie planów na tej podstawie wymaga jeszcze przyjęcia sposobu porównywania dwóch funkcji. Najczęściej używane w teorii planowania eksperymentu funkcjonały przytaczamy poniżej. Ponownie ograniczymy się tylko do takich planów, które zapewniają estymowalność wszystkich parametrów regresji.

EKSTRAPOŁACJA W PUNKCIE x_0 . Niech x_0 będzie z góry wybranym punktem, w którym interesuje nas możliwie dokładna estymacja funkcji regresji. Wówczas

$$\xi_1^{(N)} \overset{v(x_0)}{\ll} \xi_2^{(N)} \Leftrightarrow \varphi_N(x_0, \xi_1^{(N)}) \geq \phi_N(x_0, \xi_2^{(N)}). \quad (2.2.8)$$

Symbol $\overset{v(x_0)}{\ll}$ nie jest w literaturze powszechnie przyjęty, lecz wskazuje on na związek tego sposobu oceny planów z (2.2.5), w którym jako macierz A wybrano $v(x_0)v^T(x_0)$.

PORÓWNANIE ŚREDNICH WARIANCJI WYJŚĆ. Ponieważ wariancja jest funkcją nieujemną to całka z niej poprawnie opisuje uśrednione jej zachowanie, co prowadzi do uporządkowania

$$\xi_1^{(N)} \overset{Q}{\ll} \xi_2^{(N)} \Leftrightarrow \int_X \varphi_N(x, \xi_1^{(N)}) dx \geq \int_X \varphi_N(x, \xi_2^{(N)}) dx. \quad (2.2.9)$$

Symbol $\overset{Q}{\ll}$ jest w literaturze używany, lecz zwracamy uwagę, że powyższe uporządkowanie jest zgodne z (2.2.5), jeśli jako A wybierzemy $\int_X v(x)v^T(x) dx$.

PORÓWNANIE MAKSYMALNYCH WARIANCJI WYJŚĆ. Istotą tego sposobu porównywania planów jest przyjęcie, że jakość planu mierzymy wariancją odpowiedzi w tym punkcie, w którym jest ona największa, co prowadzi do

$$\xi_1^{(N)} \overset{G}{\ll} \xi_2^{(N)} \Leftrightarrow \sup_{x \in X} \varphi_N(x, \xi_1^{(N)}) \geq \sup_{x \in X} \varphi_N(x, \xi_2^{(N)}). \quad (2.2.10)$$

Oznaczenie $\overset{G}{\ll}$ jest powszechnie przyjęte. Związki tego sposobu uporządkowania z innymi omawiane będą w dalszych rozdziałach.

Oprócz wprowadzonych uporządkowań planów warto wprowadzić pojęcie równoważności planów. Każde z powyższych uporządkowań takie pojęcie zawierało. Przykładowo, jeśli posługujemy się relacją porządkującą $\overset{D}{\ll}$ (por. (2.2.2)), to $\xi_1^{(N)} \overset{D}{\ll} \xi_2^{(N)}$ oraz $\xi_2^{(N)} \overset{D}{\ll} \xi_1^{(N)}$ musi implikować równoważność planów $\xi_1^{(N)}$ i $\xi_2^{(N)}$. Dla relacji $\overset{D}{\ll}$ równoważność planów oznacza równość wyznaczników macierzy kowariancji obu planów (a zatem, także równość wyznaczników macierzy informacyjnych). Analogiczne rozważania przeprowadzić można dla każdego typu wymienionych wyżej relacji porządkujących.

Warto jednak wprowadzić także mocniejsze pojęcie równoważności planów.

Definicja 2.1. Plany $\xi_1^{(N)}$ i $\xi_2^{(N)}$ nazwiemy równoważnymi w zadaniu estymacji parametrów modelu $(X, v(x), \sigma^2(x))$, i będziemy pisać $\xi_1^{(N)} \equiv \xi_2^{(N)}$, jeśli

$$\xi_1^{(N)} \equiv \xi_2^{(N)} \Leftrightarrow M_N(\xi_1^{(N)}) = M_N(\xi_2^{(N)}). \quad (2.2.11)$$

Określenie to stosowane będzie także w odniesieniu do planów ciągłych (por. z Definicją 2.3).

Jest oczywiste, że plany równoważne w sensie tej definicji są także równoważne w sensie implikowanym przez każde z omawianych wcześniej uporządkowań.

2.3. Definicje planów eksperymentu

Załóżmy, że w planie $\xi^{(N)}$ ciąg x_1, x_2, \dots, x_N został uporządkowany w ten sposób, że pierwszych $m \leq N$ jego elementów jest różnych i są one reprezentantami pozostałych. Bardziej precyzyjnie, podciąg ten wybieramy zgodnie z następującymi regułami:

- Wśród punktów x_1, x_2, \dots, x_m , $m \leq N$ nie ma takich, które się powtarzają, tzn. $x_i \neq x_j$ dla $i \neq j$, $i, j = 1, 2, \dots, m$.
- Dla dowolnego punktu x_j , $j > m$ istnieje punkt x_i , $i \leq m$ taki, że $x_i = x_j$.

Definicja 2.2. Unormowaną wersją planu $\xi^{(N)}$, lub krótko – planem unormowanym, nazywać będziemy tablicę

$$\begin{bmatrix} x_1 & x_2 & \dots & x_m \\ p_1 & p_2 & \dots & p_m \end{bmatrix}, \quad (2.3.12)$$

gdzie $p_i \stackrel{\text{def}}{=} n_i/N$, $i = 1, 2, \dots, m$, natomiast n_i jest krotnością, z jaką punkt x_i występuje w ciągu x_1, x_2, \dots, x_N .

Ze względu na równoważność obu postaci planu, również tablicę (2.3.12) oznaczamy będziemy przez $\xi^{(N)}$.

Zachodzą proste zależności $n_i \geq 0$, $\sum_{i=1}^m n_i = N$, oraz

$$p_i \geq 0, \quad \sum_{i=1}^m p_i = 1. \quad (2.3.13)$$

Zauważmy, że jeśli $p_i = n_i/N$, to oczywiście $N \cdot p_i$ są liczbami naturalnymi (włączając 0). Odwrotnie, pewien zestaw par $\{x_i, p_i\}$, $i = 1, 2, \dots, m$ oraz pewna liczba naturalna N stanowią plan unormowany tylko wówczas, gdy spełnione jest $\sum_{i=1}^m p_i = 1$ oraz warunek

$$N p_i \quad \text{są liczbami naturalnymi,} \quad i = 1, 2, \dots, m. \quad (2.3.14)$$

Jeśli pominiemy ten warunek, zachowując pozostałe wymagania zawarte w definicji planu unormowanego, to uzyskamy użyteczne poszerzenie pojęcia planu, a mianowicie tak zwany plan ciągły, skupiony w skończonej liczbie punktów, który zdefiniowany jest następująco.

Definicja 2.3. *Planem ciągłym ξ , skupionym w skończonej liczbie punktów, nazywamy tablicę*

$$\xi = \begin{bmatrix} x_1 & x_2 & \dots & x_m \\ p_1 & p_2 & \dots & p_m \end{bmatrix}, \quad (2.3.15)$$

której elementami są punkty planu $x_i \in X$, $i = 1, 2, \dots, m$ oraz wagi p_i , spełniające warunki

$$p_i \geq 0, \quad \sum_{i=1}^m p_i = 1. \quad (2.3.16)$$

Odnotujmy, że termin „ciągły” w powyższej definicji odnosi się do pominięcia wymagania (2.3.14). Tego ważnego dla teorii i rozumienia istoty planów optymalnych uogólnienia dokonali Kiefer i Wolfowitz [68]. W literaturze angielskojęzycznej coraz częściej w ostatnich latach stosuje się termin *plany aproksymacyjne* na określenie planów ciągłych, skupionych w skończonej liczbie punktów. Termin ten dobrze oddaje istotę tych planów. Pozostaniemy jednak przy terminie stosowanym dotychczas w polskiej literaturze.

Warunki (2.3.16) pozwalają traktować wagi p_i , $i = 1, 2, \dots, m$ jako rozkład prawdopodobieństwa skupiony w skończonej liczbie punktów x_1, x_2, \dots, x_m . Innymi słowy, plan ξ to miara prawdopodobieństwa, która punktom x_1, x_2, \dots, x_m przypisuje p_i , $i = 1, 2, \dots, m$.

W teorii planowania eksperymentu pojęcie planu uogólniane jest jeszcze bardziej (por. [64], [105]).

Definicja 2.4. Ciągłym planem eksperymentu nazywa się dowolną miarę probabilistyczną μ zadaną na σ -ciele zbiorów borelowskich w X .

Aby ocenić praktyczną przydatność powyższej definicji, przytoczymy (bez technicznie zaawansowanego dowodu) następujący rezultat. Użyte w nim macierze informacyjne dla planów ciągłych są naturalnymi uogólnieniami klasycznych macierzy informacyjnych (formalne definicje podamy w następnym rozdziale).

Twierdzenie 2.1. Niech obszar planowania $X \subset R^s$ będzie zbiorem zwartym. Załóżmy, że funkcja $\sigma(x) > 0$ oraz funkcje tworzące wektor $v(x)$ o r składowych są ciągłe w X . Niech μ będzie dowolnie wybraną miarą probabilistyczną określoną na X i taką, dla której całka w (2.3.17) istnieje i jest skończona. Wówczas istnieje odpowiadający mierze μ plan ciągły ξ postaci (2.3.12) skupiony w co najwyżej $r(r+1)/2 + 1$ punktach i taki, że

$$\int_X \sigma^{-2}(x) v(x) v^T(x) \mu(dx) = \sum_{i=1}^m \sigma^{-2}(x_i) p_i v(x_i) v^T(x_i). \quad (2.3.17)$$

Równość (2.3.17) oznacza, że $\mu \equiv \xi$ w sensie Definicji 2.1.

Całka po lewej stronie równości (2.3.17) dotyczy każdego elementu macierzy $\sigma^{-2}(x)v(x)v^T(x)$, a dla każdego z nich interpretowana jest w sensie Lebesgue'a (por. [12]).

Poprzedni rezultat pozwala ograniczyć się do badania planów ciągłych, skupionych w skończonej liczbie punktów, bez straty jakości estymacji. Ograniczenie to obowiązywać będzie we wszystkich dalszych rozdziałach, w których zajmować się będziemy estymacją modelu regresji liniowej.

Definicja 2.5. Zbiór wszystkich planów ciągłych, skupionych w skończonej liczbie punktów ustalonego zbioru X , nazywać będziemy zbiorem planów dopuszczalnych i oznaczać będziemy przez Ξ lub przez $\Xi(X)$, jeśli będziemy chcieli jawnie wskazać obszar planowania X .

Plany z klasy $\Xi(X)$ podają względny rozkład częstości pomiarów w punktach x_1, x_2, \dots, x_m . Aby zastosować taki plan w praktyce, musimy przetworzyć te częstości w liczby eksperymentów, które należy wykonać w każdym z tych punktów. Poniżej przedstawiamy kilka uwag na temat realizacji planów ciągłych, skupionych w skończonej liczbie punktów.

1. Realizacja planu opisanego w Definicji 2.3 wymaga:
 - wybrania sumarycznej liczby obserwacji $N > 0$,
 - obliczenia liczb $n'_i = N p_i$, $i = 1, 2, \dots, m$,
 - wykonania operacji zaokrąglania, która polega na:
 - obliczeniu $n''_i = \lfloor N p_i \rfloor$, $i = 1, 2, \dots, m$ ($\lfloor \alpha \rfloor$ to największa liczba całkowita, nie większa niż α)

- obliczeniu $N_\rho = N - \sum_{i=1}^m n_i''$,
 - rozmieszczeniu pozostałych N_ρ pomiarów (np. losując N_ρ -krotnie (z powtórzeniami) elementy ze zbioru x_1, x_2, \dots, x_m).
2. Powyższe kroki dają w wyniku pary (x_i, n_i) , $i = 1, 2, \dots, m$. Należy jeszcze zdecydować o kolejności podawania poszczególnych „kopi” x_i na wejście badanego obiektu. Kolejność tę można wybrać na co najmniej dwa sposoby:
- losowy – zgodnie z sugestiami Fishera z początków XX wieku; zasada randomizacji w realizacji eksperymentu powinna redukować wpływ czynników, które w modelu nie zostały uwzględnione,
 - systematyczny – stosowany wtedy, gdy losowa kolejność nie może być zrealizowana, lub wtedy, gdy systematyczny wybór kolejności może przynieść dodatkowe korzyści (uproszczenie obliczeń, zmniejszenie kosztów itp.).

Bardzo szczegółową analizę dokładności zaokrąglania przeprowadzono w [110].

3. Geometria zbioru macierzy informacyjnych

Badanie geometrii zbioru wszystkich osiągalnych macierzy informacyjnych jest ważnym etapem poszukiwania planów optymalnych. Wnioski z własności geometrycznych i topologicznych tego zbioru posłużą nam, między innymi, do podania warunków dostatecznych istnienia rozwiązań optymalnych.

3.1. Elementarne własności macierzy informacyjnych

Rozpocznijmy od określenia, czy też raczej rozszerzenia, pojęcia macierzy informacyjnej na wszystkie plany ciągle skupione w skończonej liczbie punktów. Przypomnijmy, że zbiór wszystkich takich planów oznaczyliśmy przez Ξ (lub przez $\Xi(X)$).

Unormowany plan $\xi^{(N)}$ należy do klasy planów Ξ . Jednocześnie, ponieważ jest on tylko innym zapisem planu wyrażonego w „naturalnej postaci” $\{x_1, x_2, \dots, x_N\}$, potrafimy dla niego obliczyć macierz kowariancji i macierz informacyjną (por. rozdział 1.3).

Przypomnijmy, że jeśli zastosujemy taki plan i oszacujemy parametry za pomocą ważonej MNK, to macierz kowariancji $\text{cov}(\hat{a}) = [M_N(\xi^{(N)})]^{-1}$, gdzie

$$M_N(\xi^{(N)}) = N \sum_{i=1}^m \sigma^{-2}(x_i) p_i v(x_i) v^T(x_i), \quad (3.1.1)$$

zakładając, że $M_N(\xi^{(N)})$ jest nieosobliwa.

Jeśli pominiemy współczynnik proporcjonalności N i dopuścimy, by $N p_i$ nie były liczbami naturalnymi, to otrzymamy unormowaną macierz informacyjną także dla planu ciągłego. Rozszerzenie to prowadzi do następującej definicji.

Definicja 3.1. *Macierzą informacyjną (unormowaną) planu $\xi \in \Xi(X)$ zadania estymacji parametrów modelu $(X, v(x), \sigma^2(x))$ nazywamy*

$$M(\xi) = \sum_{i=1}^m \sigma^{-2}(x_i) p_i v(x_i) v^T(x_i), \quad (3.1.2)$$

przy założeniu, że dla każdego $x \in X$ zachodzi $\sigma(x) > 0$.

Naturalnym uogólnieniem definicji unormowanej macierzy informacyjnej, na przypadek gdy pojęcie planu traktujemy ogólnie, jest następujące określenie.

Definicja 3.2. *Jeśli plan ξ jest miarą probabilistyczną na X i jeśli poniższa całka istnieje i jest skończona, to jako unormowaną macierz informacyjną tego planu przyjmujemy*

$$M(\xi) = \int_X \sigma^{-2}(x) v(x) v^T(x) \xi(dx), \quad (3.1.3)$$

gdzie całka rozumiana jest w sensie Lebesgue'a, a operacja całkowania odnosi się do wszystkich elementów macierzy $\sigma^{-2}(x)v(x)v^T(x)$.

A oto podstawowe własności macierzy informacyjnych. Pomijamy ich elementarne dowody.

1. Dla dowolnego $\xi \in \Xi$, $M(\xi)$ jest symetryczną i nieujemnie określoną macierzą o r wierszach i kolumnach, gdzie $r = \dim v(x)$.
2. Macierz informacyjna $M(\xi)$ planu

$$\xi = \begin{bmatrix} x_1 & x_2 & \dots & x_m \\ p_1 & p_2 & \dots & p_m \end{bmatrix}, \quad (3.1.4)$$

- którego liczba punktów nośnika m jest mniejsza niż r , jest macierzą osobliwą.
3. Niech $M(\xi)$ będzie macierzą informacyjną modelu $(X, v(x), \sigma^2(x))$ dla planu $\xi \in \Xi(X)$, przy założeniu, że $\sigma(x) > 0$, $x \in X$. Wówczas $M(\xi)$ jest również macierzą informacyjną dla tego samego planu, lecz formalnie w innym modelu: $(X, v(x) \sigma^{-1}(x), 1)$.

Innymi słowy, modyfikując funkcje $v(x)$ (rozpinające regresję) przez podzielenie ich przez pierwiastek z wariancji zakłóceń i przyjmując formalnie stałą wariancję równą 1, otrzymujemy tę samą macierz informacyjną.

4. Niech $M(\xi)$ będzie macierzą informacyjną modelu $(X, v(x), \sigma^2(x))$ dla planu $\xi \in \Xi(X)$. Niech A , $\dim A = r \times r$ będzie zadaną macierzą. Macierz informacyjna w modelu $(X, Av(x), \sigma^2(x))$ ma postać: $A \cdot M(\xi) \cdot A^T$.

Jeśli funkcja $a^T v(x)$ poddana zostanie zmianie parametryzacji $a^T = b^T \cdot A$ (np. zmianie skali dla poszczególnych parametrów), to problem estymacji parametrów a przekształca się w zadanie estymacji wektora b w funkcji regresji rozpiętej na funkcjach $A \cdot v(x)$, a macierz informacyjna po reparametryzacji ma postać $A \cdot M(\xi) \cdot A^T$.

3.2. Klasa realizowalnych macierzy informacyjnych

W zbiorze $\Xi(X)$ określić można kombinację wypukłą planów ciągłych, korzystając z tego, że plany te można utożsamiać z odpowiednimi dyskretnymi rozkładami prawdopodobieństw, a dla rozkładów potrafimy określić ich kombinację wypukłą, zwaną często mieszaniną rozkładów.

Definicja 3.3. Dla $0 \leq \alpha \leq 1$ kombinację wypukłą planów $\xi' \in \Xi$ i $\xi'' \in \Xi$ oznaczamy następująco $(1-\alpha)\xi' + \alpha\xi''$ i definiujemy jako mieszaninę odpowiadających im rozkładów prawdopodobieństwa.

W celu uświadomienia sobie jak powstaje kombinacja wypukła planów, zapiszmy proces jej tworzenia w postaci algorytmu.

Rozważmy dwa plany o postaci:

$$\xi' = \begin{bmatrix} x'_1 & x'_2 & \cdots & x'_{m1} \\ p'_1 & p'_2 & \cdots & p'_{m1} \end{bmatrix}, \quad \xi'' = \begin{bmatrix} x''_1 & x''_2 & \cdots & x''_{m2} \\ p''_1 & p''_2 & \cdots & p''_{m2} \end{bmatrix}. \quad (3.2.5)$$

Ich kombinacja wypukła $(1-\alpha)\xi' + \alpha\xi''$ jest tablicą o tej samej postaci jak te w (3.2.5). Elementy nowej tablicy powstają następująco.

Krok 1. Jako punkty pierwszego wiersza należy wziąć

$$\{x'_1, x'_2, \dots, x'_{m1}\} \cup \{x''_1, x''_2, \dots, x''_{m2}\}.$$

Elementy tej sumy będą oznaczone przez x_1, x_2, \dots, x_m , gdzie $m \leq m1 + m2$.

Krok 2a. Plany ξ' i ξ'' poszerzamy tak, by w pierwszym wierszu obu planów pojawiły się punkty x_1, x_2, \dots, x_m .

Krok 2b. Punktom nośnika, które w danym planie pojawiły się na skutek poszerzenia, przypisujemy wagi o wartościach zero. Wagi punktów po tym poszerzeniu oznaczać nadal będziemy, odpowiednio, przez p'_i i p''_i .

Zauważmy, że mimo poszerzenia, wagi te nadal sumują się do jedności.

Krok 3. Wagi wynikowe powstają zgodnie ze wzorem $(1-\alpha)p'_i + \alpha p''_i$, $i = 1, 2, \dots, m$.

Z kroków tych wynika natychmiast, że kombinacja wypukła planów $\xi' \in \Xi(X)$ i $\xi'' \in \Xi(X)$ jest również ciągłym planem eksperymentu, tzn.

$$[(1-\alpha)\xi' + \alpha\xi''] \in \Xi(X),$$

dla dowolnego $0 < \alpha < 1$. Innymi słowy:

Własność 5. Zbiór $\Xi(X)$ wszystkich planów ciągłych na ustalonym zbiorze X jest zbiorem wypukłym.

Kombinację wypukłą planu ξ' i planu jednopunktowego skupionego w x oznaczać będziemy przez $(1-\alpha)\xi' + \alpha\delta_x$.

Własność 6. Niech ξ'' będzie planem skupionym w jednym punkcie x z wagą równą 1 i niech

$$\xi' = \begin{bmatrix} x'_1 & x'_2 & \cdots & x'_{m1} \\ p'_1 & p'_2 & \cdots & p'_{m1} \end{bmatrix} \quad (3.2.6)$$

Załóżmy, że $x \neq x_i$, $i = 1, 2, \dots, x_{m1}$. Wówczas, dla $0 < \alpha < 1$, plan $\xi = (1 - \alpha)\xi' + \alpha\xi''$ ma postać

$$\xi = \begin{bmatrix} x'_1, & x'_2, & \dots, & x'_{m1}, & x \\ (1 - \alpha)p'_1, & (1 - \alpha)p'_2, & \dots, & (1 - \alpha)p'_{m1}, & \alpha \end{bmatrix}. \quad (3.2.7)$$

Własność ta będzie często wykorzystywana w algorytmach numerycznego poszukiwania planów optymalnych. Jest ona także podstawowa dla badania geometrii zbioru osiągalnych macierzy informacyjnych, czyli takich, które da się wygenerować poprzez zmienianie wszystkich ciągłych planów eksperymentu.

Definicja 3.4. Zbiorem osiągalnych macierzy informacyjnych $\mathcal{M}(X, v, \sigma)$ w zadaniu estymacji $(X, v(x), \sigma^2(x))$ nazywamy

$$\mathcal{M}(X, v, \sigma) = \{M(\xi) : \xi \in \Xi(X)\}. \quad (3.2.8)$$

Będziemy pomijać argumenty \mathcal{M} , gdy postać zadania estymacji wynikać będzie z kontekstu.

Zauważmy, że jeśli $M_1, M_2 \in \mathcal{M}$, to także

$$(\alpha M_1 + (1 - \alpha) M_2) \in \mathcal{M}, \quad (3.2.9)$$

dla dowolnego $0 \leq \alpha \leq 1$. Jest tak, gdyż M_1, M_2 mogą być elementami \mathcal{M} tylko wówczas, gdy istnieją takie plany $\xi_1, \xi_2 \in \Xi$, dla których $M_1 = M(\xi_1)$ oraz $M_2 = M(\xi_2)$. Tak więc

$$\alpha M_1 + (1 - \alpha) M_2 = M[\alpha \xi_1 + (1 - \alpha) \xi_2]. \quad (3.2.10)$$

Wobec wypukłości zbioru planów ciągłych mamy: $(\alpha \xi_1 + (1 - \alpha) \xi_2) \in \Xi$, co implikuje $M[\alpha \xi_1 + (1 - \alpha) \xi_2] \in \mathcal{M}$. W ten sposób udowodniliśmy:

Własność 7. $\mathcal{M}(X, v, \sigma)$ jest zbiorem wypukłym.

Następna własność zbioru \mathcal{M} jest ważna, gdyż zapewnia istnienie rozwiązań zadań planowania optymalnego.

Twierdzenie 3.1. Jeżeli obszar planowania X jest zbiorem zwartym, a funkcje $v(x)$ oraz $\sigma(x)$ są ciągle w X , to zbiór $\mathcal{M}(X, v, \sigma)$ też jest zwarty.

Dowód tego twierdzenia pominiemy (patrz [105]). Następna własność jest wnioskiem z Twierdzenia 3.1 i znanego twierdzenia Caratheodory'ego (por. [105]).

Własność 8. Niech spełnione będą założenia Twierdzenia 3.1.

1. Dla każdego planu $\xi \in \Xi(X)$ istnieje plan $\xi' \in \Xi(X)$ skupiony w co najwyżej $r(r + 1)/2 + 1$ punktach i taki, że $M(\xi) = M(\xi')$.
2. Jeśli ponadto macierz $M(\xi)$ leży na brzegu zbioru \mathcal{M} , to plan ξ' można wybrać tak, by jego nośnik zawierał co najwyżej $r(r + 1)/2$ punktów i jednocześnie, by równoważność planów była zachowana.

Zbiorowi \mathcal{M} można nadać interpretację geometryczną. Wobec symetrii macierzy informacyjnych, możemy utożsamiać je z wektorami $r(r+1)/2$ wymiarowymi, utworzonymi z uporządkowanych elementów macierzy leżących na jej przekątnej lub ponad nią. \mathcal{M} da się zatem zobrazować w przestrzeni $R^{r(r+1)/2}$. Potencjalna konstrukcja geometrycznego kształtu zbioru \mathcal{M} rozpoczyna się od zobrazowania zbioru $\{\sigma^{-2}(x)v(x)v^T(x) : x \in X\}$, a następnie znalezienia jego powłoki wypukłej. Rzeczywista możliwość narysowania zbioru \mathcal{M} ogranicza się oczywiście do przypadku $r = 2$. Gdy $r = 3$ i pierwszym elementem wektora $v(x)$ jest 1, szkicowany bywa rzut prostopadły zbioru \mathcal{M} , powstający przez pominięcie składowej stałej.

4. Optymalne plany eksperymentu

Rozdział ten zajmuje centralne miejsce w pierwszej części tej monografii. Zestawiamy w nim znane kryteria optymalności planów i szczegółowo omawiamy klasyczne twierdzenia Kiefera–Wolfowitza, podające warunki konieczne i dostateczne optymalności. W ostatnim podrozdziale omawiamy te klasy zadań planowania, które dzięki temu twierdzeniu rozwiązać można analitycznie.

4.1. Dalsze uwagi na temat oceny jakości planu

Każdy ze sposobów uporządkowania planów, które były omawiane w rozdziale 2.2 można natychmiast rozszerzyć na plany ciągłe o nieosobliwej macierzy informacyjnej. Każda z tych relacji porządkujących pozwala też zdefiniować odpowiadające jej zadanie poszukiwania optymalnego planu. Przegląd sformułowań takich zadań jest celem niniejszego podrozdziału. Pokażemy je na tle ogólnego sformułowania problemu planowania optymalnego.

Podobnie jak przy definiowaniu relacji porządkujących, posługiwać się można funkcjami zdefiniowanymi na zbiorze macierzy kowariancji lub macierzy informacyjnych. Przyjmijmy konwencję wyrażania ich w terminach $M(\xi)$. Porównywać będziemy jakość jedynie takich planów ciągłych, których macierz informacyjna jest nieosobliwa. Zbiór wszystkich takich planów oznaczamy będziemy przez $\hat{\Xi} \stackrel{\text{def}}{=} \{\xi \in \Xi : \det M(\xi) > 0\}$. Zauważmy, że zbiór planów Ξ zależy tylko od obszaru planowania X , natomiast jego podzbiór $\hat{\Xi}$ wyznaczany jest także przez przyjęty model $(X, v(x), \sigma^2(x))$. Nie będziemy jawnie wskazywać na tę zależność, aby nie komplikować notacji.

Funkcje oceniające jakość planu i służące do zdefiniowania planów optymalnych przyjęto nazywać kryteriami planowania eksperymentu, choć nazwa „wskaźnik jakości planowania” byłaby bardziej właściwa. Dalej podajemy zestaw ogólnych wymagań nakładanych zwykle na funkcje kryterialne.

Niech Φ będzie funkcją, która odwzorowuje elementy zbioru symetrycznych $r \times r$ macierzy nieujemnie określonych w zbiór liczb rzeczywistych. Zbiór wszystkich symetrycznych i nieujemnie określonych macierzy o r wierszach i kolumnach oznaczamy dalej przez \mathcal{A}_r . Zakładamy przy tym, że jeśli argument Φ jest macierzą osobliwą, to Φ przyjmuje wartość $-\infty$. Prócz tego przyjmujemy następujące założenia:

JEDNORODNOŚĆ. Dla dowolnego $\alpha > 0$ i macierzy $A \in \mathcal{A}_r$ $\Phi(\alpha A) = s(\alpha)\Phi(A)$ dla pewnej skalarnej, nieujemnej i rosnącej funkcji s .

MONOTONICZNOŚĆ. Dla dowolnych macierzy $A, B \in \mathcal{A}_r$ ale takich, że $A - B \geq 0$, zachodzi $\Phi(A) \geq \Phi(B)$.

WKŁĘŚŁOŚĆ. Dla dowolnego $0 < \alpha < 1$ oraz $A, B \in \mathcal{A}_r$ zachodzi nierówność:

$$\Phi((1 - \alpha)A + \alpha B) \geq (1 - \alpha)\Phi(A) + \alpha\Phi(B).$$

RÓŻNICZKOWALNOŚĆ. Dla nieosobliwych macierzy $A \in \mathcal{A}_r$ istnieje $r \times r$ macierz

$$F(A) = \left[\frac{d\Phi(A)}{da_{ij}} \right], \quad i, j = 1, 2, \dots, r.$$

Definicja 4.1. Plan $\xi^* \in \Xi$ nazywamy Φ -optymalnym, jeśli

$$\max_{\xi \in \Xi} \Phi(M(\xi)) = \Phi(M(\xi^*)). \quad (4.1.1)$$

Przyjęto tu konwencję maksymalizacji funkcji Φ . Jeśli dalej pojawi się zadanie minimalizacji kryterium planowania, to wystarczy za Φ podstawić $-\Phi$, by dopasować takie zadanie do powyższego ogólnego schematu.

Jak już wspominaliśmy, porównywanie jakości planów za pośrednictwem funkcji $\det[M^{-1}(\xi)]$ należy do najbardziej rozpowszechnionych. Ponieważ

$$\det M^{-1}(\xi) = 1/\det M(\xi),$$

to przyjęła się następująca wersja definiowania zadania D-optymalnego planowania.

Definicja 4.2. Plan ξ^* nazywa się D-optymalnym, jeżeli

$$\det M(\xi^*) = \max_{\xi \in \Xi} \det M(\xi). \quad (4.1.2)$$

Wobec ścisłej monotoniczności funkcji logarytm, planów D-optymalnych poszukiwać możemy maksymalizując funkcjonal $\Phi(\xi) = \ln \det(M(\xi))$.

Definicja 4.3. Niech A będzie wybraną $r \times r$ macierzą symetryczną. Plan $\tilde{\xi}$ nazywamy L-optymalnym, jeżeli

$$\text{tr} [A M^{-1}(\tilde{\xi})] = \min_{\xi \in \Xi} \text{tr} [A M^{-1}(\xi)] \quad (4.1.3)$$

Szczególnymi przypadkami tego kryterium jest ważne kryterium A-ptymalności, które otrzymuje się, gdy $A = I_r$. Podobnie, przyjęcie $A = cc^T$, gdzie $c \in R^r$ jest wybranym wektorem, prowadzi do kryterium c-ptymalności. Jeśli, z kolei jako c przyjmiemy wektor $v(x_0)$, to otrzymamy kryterium ekstrapolacji w (wybranym przez nas) punkcie x_0 . Nazwa ta ma następujące uzasadnienie. Jeśli $A = v(x_0)v^T(x_0)$, to kryterium L-ptymalności sprowadza się do minimalizacji $v^T(x_0)M^{-1}(\xi)v(x_0)$, a wyrażenie to jest proporcjonalne do wariancji predykcji w punkcie x_0 .

Definicja 4.4. Plan $\tilde{\xi}$ nazywamy E-ptymalnym, jeżeli

$$\lambda_{\max} [M^{-1}(\tilde{\xi})] = \min_{\xi \in \Xi} \lambda_{\max} [M^{-1}(\xi)], \quad (4.1.4)$$

gdzie $\lambda_{\max}[\cdot]$ oznacza maksymalną wartość własną macierzy w nawiasach kwadratowych.

Definicja 4.5. Niech $p > 0$ będzie wybraną liczbą naturalną. Plan $\tilde{\xi}$ nazywamy L_p -ptymalnym, jeżeli

$$\left\{ \text{tr} [M^{-p}(\tilde{\xi})] \right\}^{1/p} = \min_{\xi \in \Xi} \left\{ \text{tr} [M^{-p}(\xi)] \right\}^{1/p} \quad (4.1.5)$$

Gdy $p \rightarrow 0+$, to $\left\{ \text{tr} [M^{-p}(\xi)] \right\}^{1/p}$ jest zbieżne do $\det(M^{-1}(\xi))$. Natomiast, gdy $p \rightarrow \infty$, to $\left\{ \text{tr} [M^{-p}(\xi)] \right\}^{1/p} \rightarrow \lambda_{\max} [M^{-1}(\xi)]$. Dla $p = 1$ otrzymujemy oczywiście kryterium A-ptymalności. Gama kryteriów L_p pokrywa szeroki zakres kryteriów od D-ptymalności do E-ptymalności.

Omówione poprzednio kryteria planowania odpowiadały uporządkowaniom, w których plan oceniany był z punktu widzenia dokładności oceny parametrów. Omawiane dalej kryteria są odpowiednikami uporządkowań (2.2.9) i (2.2.10), a więc optymalność widziana jest tu pod kątem jakości estymacji funkcji regresji.

Odpowiednikiem wariancji odpowiedzi (2.2.7) dla planów ciągłych $\xi \in \Xi$ jest

$$\phi(x, \xi) = v^T(x) M^{-1}(\xi) v(x). \quad (4.1.6)$$

Uwaga 4.1. W powyższym wzorze i dalej w całej książce przyjęto $\sigma^2(x) \equiv 1$. Nie ogranicza to ogólności rozważań, gdyż podstawienie $\sigma^{-1}(x)v(x)$ w miejsce $v(x)$ redukuje problem do takiego, w którym wariancja równa jest jeden. Zwracamy też uwagę Czytelnika na fakt, że funkcje $\phi(x, \xi)$ i (2.2.7) są do siebie proporcjonalne tylko wówczas, gdy wariancja zakłóceń jest stała (niezależna od x) w całym obszarze planowania. W przeciwnym razie, $\phi(x, \xi)$ traktować należy jako wygodny skrócony zapis dla prawej strony wzoru (4.1.6). Dla zaznaczenia tej różnicy wprowadziliśmy inny krój czcionki (ϕ zamiast dotychczasowego φ) i zamieniliśmy kolejność argumentów.

Definicja 4.6. Plan $\hat{\xi}$ nazywamy *G- optymalnym*, jeżeli

$$\max_{x \in X} v^T(x) M^{-1}(\hat{\xi}) v(x) = \min_{\xi \in \Xi} \left\{ \max_{x \in X} v^T(x) M^{-1}(\xi) v(x) \right\}. \quad (4.1.7)$$

Temu ważnemu kryterium poświęcamy dużo uwagi w następujących rozdziałach.

Gdy oceniamy wariancję predykcji, po uśrednieniu jej w pewnym obszarze $X_0 \subseteq X$, to plan, który minimalizuje funkcjonal $\int_{X_0} \phi(\xi, x) dx$ nazywa się Q- optymalnym. Kryterium Q- optymalności jest szczególnym przypadkiem kryterium L- optymalności.

W literaturze (por. [105], [28], [79], [25]) rozpatruje się uogólnienie kryterium D- optymalności o postaci $\det(A M(\xi) A^T)$, gdzie A jest $s \times r$ macierzą. Jeśli $s \neq r$, to otrzymamy plany inne niż D- optymalne. Kryterium to jest przykładem tzw. częściowego kryterium optymalności (por. [105], [175], [110]). Kryteria te biorą pod uwagę dokładność estymacji tylko części parametrów i nie będą dalej omawiane, gdyż prowadzić mogą do planów o osobliwych macierzach informacyjnych.

4.2. Własności kryteriów D- i G- optymalności

W rozdziale tym funkcja $\Phi(\cdot) \stackrel{\text{def}}{=} \ln \det(\cdot)$. Zestawimy jej podstawowe własności jako funkcji macierzy i jako kryterium D- optymalności.

Przypomnijmy, że wklęsłość (wypukłość w górę) oznacza, że dla dowolnych $M_1, M_2 \in \mathcal{M}$ oraz dowolnego $\alpha \in (0,1)$

$$(1 - \alpha) \Phi(M_1) + \alpha \Phi(M_2) \leq \Phi[(1 - \alpha) M_1 + \alpha M_2].$$

A oto zapowiadany zestaw własności logarytmu wyznacznika macierzy. W ich opisie ξ^* oznacza plan D- optymalny w pewnym ustalonym modelu $(X, v(x), \sigma^2(x))$ o nieosobliwej macierzy informacyjnej.

1. Funkcja $\ln \det(\cdot)$ rozpatrywana na zbiorze \mathcal{M} jest funkcją wklęsłą. Jest ona ściśle wklęsła na $\hat{\mathcal{M}}$ (dowód [105]).
2. Niech $A(t)$ będzie pewną funkcją skalarne argumentu t , której wartościami są symetryczne i dodatnio określone macierze o wymiarach $r \times r$. Załóżmy, że funkcja ta jest różniczkowalna w otoczeniu punktu t . Wówczas

$$\frac{d}{dt} \Phi(A(t)) = \text{tr} \left[A^{-1}(t) \cdot \frac{dA(t)}{dt} \right], \quad (4.2.8)$$

gdzie tr oznacza ślad macierzy (dowód [35]).

3. Niech $\xi_1 \in \hat{\Xi}$ i $\xi_2 \in \Xi$ będą dwoma planami w modelu $(X, v(x), \sigma^2(x))$. Wówczas

$$\frac{d}{d\alpha} \Phi [(1 - \alpha) \xi_1 + \alpha \xi_2] \Big|_{\alpha=0+} = -r + \text{tr} [M^{-1}(\xi_1) \cdot M(\xi_2)]. \quad (4.2.9)$$

Jeżeli ponadto plan ξ_2 skupiony jest tylko w jednym punkcie $x \in X$, to

$$\frac{d}{d\alpha} \Phi [(1 - \alpha) \xi_1 + \alpha \xi_2] \Big|_{\alpha=0+} = -r + \sigma^{-2}(x) v^T(x) M^{-1}(\xi_1) v(x). \quad (4.2.10)$$

Dowód – bezpośrednio zastosowanie (4.2.8).

4. Niech macierz informacyjna D- optymalnego planu ξ^* , będzie nieosobliwa. Jeśli plan ξ^{**} jest również D- optymalny w tym samym problemie estymacji, to $M(\xi^*) = M(\xi^{**})$.

Dowód – natychmiastowa konsekwencja własności 1.

5. Niech A , $\dim A = r \times r$ będzie zadaną macierzą nieosobliwą. Rozważmy zadanie D- optymalnego planowania dla modelu, w którym wektor $v(x)$ poddano transformacji $(X, A \cdot v(x), \sigma^2(x))$. Wówczas plan ξ^* jest równocześnie optymalny w modelu przekształconym i odwrotnie.

Dowód wynika z tego, że wyznacznik iloczynu macierzy kwadratowych jest iloczynem ich wyznaczników. (Uwaga – założenie, że A jest nieosobliwa jest istotne.)

6. Dla dowolnego planu $\xi \in \hat{\Xi}$ wariancja $\phi(\xi, x) = v^T(x) M^{-1}(\xi) v(x)$ estymacji wyjścia modelu $(X, v(x), 1)$ spełnia następujące zależności:

$$\sum_{i=1}^m \phi(\xi, x_i) p_i = r \quad (4.2.11)$$

$$\max_{x \in X} \phi(\xi, x) \geq r, \quad (4.2.12)$$

gdzie r jest liczbą estymowanych parametrów.

Dowody powyższych własności wynikają z następującego ciągu zależności:

$$r = \text{tr} [M^{-1}(\xi) \cdot M(\xi)] = \sum_{i=1}^m \text{tr} [M^{-1}(\xi) v(x_i) v^T(x_i)] p_i = \quad (4.2.13)$$

$$= \sum_{i=1}^m \phi(\xi, x_i) p_i \leq \max_{i=1, \dots, m} \phi(\xi, x_i) \leq \max_{x \in X} \phi(\xi, x).$$

Własność (4.2.11) oznacza, że ważona średnia wariancji predykcji w punktach planu nie zależy od użytego planu. Nie znaczy to oczywiście, że nie warto planować eksperymentu. Jak wynika z (4.2.12), w punktach gdzie niedokładność estymacji jest największa, wariancja wyrażenia $\hat{a}^T v(x)$ jest zawsze nie mniejsza niż r . Przez odpowiedni dobór planu można się zatem starać o to, by wariancja ta nie przekraczała r w żadnym punkcie obszaru X . Jak zobaczymy, własność ta jest charakterystyczna dla planów D- i G- optymalnych.

4.3. Równoważność planów D- i G- optymalnych

Twierdzenie Kiefera i Wolfowitza (por. [69] oraz [68], [65], [66], [67]) było i nadal jest jednym z najważniejszych rezultatów teorii planowania eksperymentu. Odgrywa ono również zasadniczą rolę w konstruowaniu numerycznych algorytmów planowania eksperymentu.

Twierdzenie 4.1. (Kiefer i Wolfowitz 1960) *Niech w modelu liniowym funkcje $v(x)$ będą ciągle na zwartym zbiorze X . Załóżmy, że w zadaniu estymacji parametrów modelu $(X, v(x), 1)$ istnieje w $\Xi(X)$ plan o nieosobliwej macierzy informacyjnej, wówczas następujące stwierdzenia, dotyczące planów z $\Xi(X)$, są równoważne:*

1. plan ξ^* jest D-optymalny w estymacji parametrów modelu $(X, v(x), 1)$,
2. plan ξ^* jest G-optymalny w estymacji parametrów modelu $(X, v(x), 1)$,
3. plan ξ^* spełnia warunek

$$\max_{x \in X} \phi(\xi^*, x) = r, \quad (4.3.14)$$

gdzie $\phi(\xi^*, x) = v(x)^T M^{-1}(\xi^*) v(x)$. Ponadto, we wszystkich punktach nośnika $x_1^*, x_2^*, \dots, x_m^*$ planu ξ^* osiągnięte jest maksimum w wyrażeniu (4.3.14), tzn. $\phi(\xi^*, x_i^*) = r$, $i = 1, 2, \dots, m$.

Ponieważ przy poczynionych założeniach funkcja $\phi(\xi^*, x)$ jest ciągła na domkniętym i ograniczonym zbiorze X , to warunek (4.3.14) zastąpić można następującym

$$\sup_{x \in X} \phi(\xi^*, x) = r. \quad (4.3.15)$$

Dowód tego twierdzenia, na podstawie różniczkowych warunków optymalności funkcji wypukłych, znaleźć można w wielu monografiach [105], [35], [175]. My ograniczymy się do wykazania, że tezy 1) i 3) są równoważne.

Rozważmy plan o postaci: $\xi_\alpha = (1 - \alpha)\xi^* + \alpha\delta_x$, $\alpha \in (0, 1)$, gdzie δ_x oznacza plan skupiony w jednym punkcie $x \in X$. D-optymalność planu ξ^* implikuje, że

$$\left. \frac{d\Phi(\xi_\alpha)}{d\alpha} \right|_{\alpha=0} \leq 0,$$

co, na mocy własności 3), prowadzi do nierówności:

$$v^T(x) M^{-1}(\xi^*) v(x) \leq r, \quad x \in X. \quad (4.3.16)$$

Stąd mamy natychmiast $\max_{x \in X} \phi(\xi^*, x) \leq r$. Nierówność ta po zestawieniu jej z (4.2.12) kończy dowód stwierdzenia 3), o ile zachodzi 1).

Dowód implikacji odwrotnej przeprowadzimy przez sprowadzenie do sprzeczności. Załóżmy, że (4.3.14) zachodzi, ale ξ^* nie jest D-optymalny. Jednocześnie

wiemy, że istnieje pewien plan D- optymalny, który oznaczmy przez $\bar{\xi}$. A skoro ξ^* nie jest optymalny, to $\Phi(\bar{\xi}) > \Phi(\xi^*)$. Ścisła wklęsłość funkcji Φ implikuje wtedy, że również plany postaci:

$$\bar{\xi}_\alpha = (1 - \alpha)\xi^* + \alpha\bar{\xi}, \quad \alpha \in (0,1),$$

będą lepsze niż ξ^* . Zatem,

$$\frac{d}{d\alpha} \Phi(\bar{\xi}_\alpha) \Big|_{\alpha=0} > 0.$$

Nierówność ta, po skorzystaniu z własności 3 i (4.2.9)), daje w rezultacie

$$r < \text{tr} \left[M^{-1}(\xi^*) M(\bar{\xi}) \right] = \sum_{i=1}^m \bar{p}_i \phi(\xi^*, \bar{x}_i) \leq \max_{i=1, \dots, m} \phi(\xi^*, \bar{x}_i), \quad (4.3.17)$$

gdzie \bar{p}_i i \bar{x}_i to, odpowiednio, wagi i punkty nośnika planu $\bar{\xi}$. Nierówność ta jest sprzeczna z warunkiem (4.3.14), którego prawdziwość wcześniej założono, co kończy dowód D- optymalności planu ξ^* .

Poniżej podajemy kilka uwag na temat interpretacji i wniosków z twierdzenia Kiefera–Wolfowitza.

1. Rezultat ten można wysłowić następująco: plan D- optymalny jest równocześnie G- optymalny. Plan o możliwie dużej dokładności estymacji parametrów modelu zapewnia zatem najdokładniejszą estymację samej funkcji regresji w sensie minimaxowym, omówionym dalej.

Kolejność operacji min i max we wzorze (4.1.7) wybrana jest tak, że dla każdego ustalonego planu $\xi \in \Xi$ obliczana jest maksymalna w X wartość funkcji wariancji (maksimum osiągane może być w różnych punktach). Minimalizacja z kolei realizowana jest po wszystkich planach $\xi \in \Xi$ w celu znalezienia takiego planu, dla którego odpowiadająca mu maksymalna wartość wariancji w X jest najmniejsza.

2. Można udowodnić następujący wniosek z twierdzenia Kiefera–Wolfowitza (por. [35], [36]).

Niech $\xi_0 \in \hat{\Xi}$ będzie pewnym planem dla estymacji modelu $(X, v(x), \sigma^2)$ o nie- osobliwej macierzy informacyjnej, a ξ^* jest planem D- optymalnym. Oznaczmy $\phi_0 \stackrel{\text{def}}{=} \sup_{x \in X} \phi(\xi_0, x)$. Spełniona jest następująca nierówność:

$$\phi_0 - r \geq \ln \det M(\xi^*) - \ln \det M(\xi_0) \geq 0. \quad (4.3.18)$$

Nierówność (4.3.18) pozwala ocenić stratę dokładności, jaką poniesiemy stosując plan ξ_0 , zamiast planu ξ^* . Do dokonania tej oceny nie musimy znać planu ξ^* , wystarczy nam wartość $e = \phi_0 - r$; jest ona dodatnia dla każdego planu, który nie jest D- optymalny.

3. Znana jest wersja tego twierdzenia dla wariancji zakłóceń zależnej od czynników eksperymentu (por. uwaga 4.1).

4.4. Uogólnienie twierdzenia Kiefera i Wolfowitza

W podrozdziale tym Φ oznacza dowolną funkcję macierzy informacyjnej, która spełnia wymagania, nałożone na kryteria planowania na początku tego rozdziału.

Uogólnieniem twierdzenia Kiefera–Wolfowitza na bardzo szeroką klasę kryteriów, prowadzących do planów o nieosobliwej macierzy informacyjnej, jest następujący rezultat [67], [105], [110].

Twierdzenie 4.2. *Załóżmy, że kryterium planowania Φ jest funkcją jednorodną, monotoniczną, wklęsłą i różniczkowalną w sensie omówionym w rozdziale 4.1.*

W zadaniu estymacji parametrów modelu $(X, v(x), \sigma^2(x))$ plan ξ^ jest Φ -optymalny wtedy i tylko wtedy, gdy spełniony jest warunek*

$$\max_{x \in X} \Psi(\xi^*, x) = \sum_{i=1}^m p_i^* \Psi(\xi^*, x_i^*), \quad (4.4.19)$$

gdzie: $\Psi(\xi^*, x) = \sigma^{-2}(x) v^T(x) F(M(\xi^*)) v(x)$, $x \in X$, a macierz $F(M(\xi^*))$ obliczana jest następująco, najpierw obliczany jest gradient

$$F(A) = \left[\frac{d\Phi(A)}{da_{ij}} \right], \quad i, j = 1, 2, \dots, r,$$

a następnie w miejsce macierzy A wstawiana jest macierz $M(\xi^*)$.

Warto zauważyć, że funkcja $\Psi(\xi^*, x)$ pełni podobną rolę jak funkcja $\phi(\xi^*, x)$ dla kryterium D-optymalności, a warunek (4.4.19) jest odpowiednikiem (4.3.14).

Efektywne korzystanie z tego twierdzenia wymaga znajomości pochodnych kryteriów planowania względem elementów macierzy informacyjnej (dalej nazywać będziemy je gradientami kryteriów planowania mimo, że nazwa ta nie jest powszechnie przyjęta). Poniżej przytaczamy za [105] zestaw gradientów najczęściej używanych kryteriów planowania.

Lemat 4.1. *Zakładamy, że obliczane poniżej gradienty liczone są dla nieosobliwej macierzy M .*

Gradient kryterium D-optymalności. *W zadaniu maksymalizacji*

$$\Phi(M) = \log \det(M) \quad \text{otrzymujemy: } F(M) = M^{-1}.$$

Gradient kryterium A-optymalności. *W problemie minimalizacji $\text{tr}[M^{-1}]$, czyli maksymalizacji*

$$\Phi(M) = -\text{tr}[M^{-1}] \quad \text{dostaniemy: } F(M) = M^{-2}.$$

Gradient kryterium L- optymalności. W problemie minimalizacji $\text{tr}[A M^{-1}]$, czyli maksymalizacji

$$\Phi(M) = -\text{tr}[A M^{-1}] \quad \text{dostaniemy:} \quad F(M) = M^{-1} A M^{-1}.$$

Jako przykład zastosowania powyższego twierdzenia rozpatrzmy kryterium A- optymalności. Łatwo sprawdzić jednorodność, monotoniczność i różniczkowalność funkcji $\text{tr}(M^{-1})$. Odnotujmy bez dowodu (por. [105]), że funkcja ta jest ściśle wypukła na $\hat{\mathcal{M}}$. Z warunku (4.4.19) natychmiast otrzymamy wniosek.

Wniosek 4.1. Załóżmy, że $\sigma(x) \equiv 1$. Plan ξ^* jest A- optymalny wtedy i tylko wtedy, gdy

$$\max_{x \in X} v^T(x) M^{-2}(\xi^*) v(x) = \sum_{i=1}^m p_i^* v^T(x_i^*) M^{-2}(\xi^*) v(x_i^*). \quad (4.4.20)$$

4.5. Wybrane optymalne plany eksperymentów

Analityczne znajdowanie optymalnych planów może polegać na „odgadnięciu” pewnego planu i sprawdzeniu jego optymalności, korzystając z (4.3.14). Można również próbować wybrać dostatecznie ogólną strukturę planu, a jego parametry tak „dostroić”, by spełniony został warunek optymalności. Autorzy rezultatów o optymalności poszczególnych klas planów zwykle nie podają sposobu, w jaki zostały one otrzymane, lecz można się domyślać, że – wobec braku ogólnej metodologii – postacie planów optymalnych były otrzymywane w jeden z opisanych już sposobów. W rozdziale tym dokonamy przeglądu planów optymalnych (głównie D- optymalnych) dla estymacji funkcji regresji (głównie) o jednej zmiennej. Dalej plany te posłużą nam jako „klocki” do konstruowania planów dla szerokiej klasy funkcji wielu zmiennych w systematyczny sposób.

Analityczne znajdowanie planów D- optymalnych znacznie się upraszcza, gdy liczba punktów nośnika planu jest równa minimalnej możliwej liczbie r , przy której macierz informacyjna może być nieosobliwa.

Rzeczywiście, jeśli plan D- optymalny w modelu $(X, v(x), 1)$ skupiony jest dokładnie w r punktach, gdzie $r = \dim v(x)$, to wszystkie wagi tego planu są równe i wynoszą $1/r$.

Tak być musi, gdyż wyznacznik macierzy informacyjnej takiego planu ma postać

$$(\det V_r)^2 \prod_{i=1}^r p_i, \quad (4.5.21)$$

gdzie $r \times r$ macierz $V_r \stackrel{\text{def}}{=} [v(x_1^*), v(x_2^*), \dots, v(x_r^*)]$, natomiast x_i^* , $i = 1, 2, \dots, r$, to punkty nośnika planu D- optymalnego skupionego w r punktach. Wystarczy

teraz znaleźć maksimum wyrażenia $\prod_{i=1}^r p_i$ względem $p_i \geq 0$, przy ograniczeniu $\sum_{i=1}^r p_i = 1$. Jak wiadomo, osiągnięte jest ono dla $p_i = 1/r$, $i = 1, 2, \dots, r$.

Jeśli podejrzewamy, że istnieje D-optimalny plan r -punktowy, to wystarczy teraz optymalizować położenie punktów jego nośnika.

Jako prosty przykład wykażemy D-optimalność planu

$$\xi^* = \begin{bmatrix} -1 & 1 \\ 1/2 & 1/2 \end{bmatrix}$$

w estymacji parametrów modelu liniowego $y = a^{(1)} + a^{(2)}x$, gdy obszarem planowania jest $X = [-1, 1]$. Macierz informacyjna tego planu ma postać $M(\xi^*) = I_2$, gdzie I_2 oznacza macierz jednostkową 2×2 , co implikuje: $\phi(\xi^*, x) = 1 + x^2$. Funkcja ta osiąga maksimum równe 2 w punktach ± 1 , co dowodzi D-optimalności tego planu, zgodnie z twierdzeniem Kiefera–Wolfowitza.

Korzystając z warunku (4.3.14), można sprawdzić, że plan

$$\xi^* = \begin{bmatrix} -1 & 0 & 1 \\ 1/3 & 1/3 & 1/3 \end{bmatrix} \quad (4.5.22)$$

jest D-optimalny dla estymacji regresji $a^{(1)} + a^{(2)}x + a^{(3)}x^2$ na odcinku $[-1, 1]$.

Podobnie, korzystając z warunku (4.4.20), można sprawdzić A-optimalność planu

$$\tilde{\xi} = \begin{bmatrix} -1 & 0 & 1 \\ 1/4 & 1/2 & 1/4 \end{bmatrix} \quad (4.5.23)$$

w estymacji funkcji regresji $a^{(1)} + a^{(2)}x + a^{(3)}x^2$, $X = [-1, 1]$, przy założeniu, że dyspersja zakłóceń jest stała $\sigma(x) = 1$, $x \in X$.

W pracach [64], [105] znaleźć można dowody następujących rezultatów.

Twierdzenie 4.3. *Niech funkcja regresji będzie wielomianem stopnia $r - 1$ o postaci*

$$\bar{y}(x) = \sum_{k=1}^r a^{(k)} x^{k-1} = a^T v(x), \quad (4.5.24)$$

gdzie $v(x) = [1, x, \dots, x^{r-1}]^T$. Parametry $a = [a^{(1)}, a^{(2)}, \dots, a^{(r)}]^T$ są estymowane na podstawie obserwacji w $X = [-1, 1]$.

1. Jeśli ponadto, $\sigma(x) \equiv 1$, to plan D-optimalny jest jednoznacznie wyznaczony, a charakteryzuje się on tym, że skupiony jest w r punktach, które są pierwiastkami następującego równania

$$(1 - x^2) \cdot \frac{dP_{r-1}(x)}{dx} = 0, \quad (4.5.25)$$

gdzie $P_k(x)$ jest wielomianem Legendre'a stopnia k . Wszystkie wagi planu są równe $1/r$.

2. Jeśli $\sigma^{-2}(x) = (1-x)^{\alpha+1}(1+x)^{\beta+1}$, gdzie $\alpha > -1$, $\beta > -1$, to plan D-optimalny również jest jedyny i skupiony w r punktach z wagami $1/r$. Punktami nośnika planu są pierwiastki wielomianu Jacobiego stopnia r -tego o parametrach α i β .

W monografii [64] opisano plany D-optimalne dla regresji wielomianowej z wykładniczymi funkcjami dyspersji $\sigma(x)$ i obszarami planowania $X = [0, \infty)$ oraz $X = (-\infty, \infty)$.

Czytelnik zauważy pokrewieństwo tezy 1) powyższego twierdzenia i następującego rezultatu (por. [64]).

Twierdzenie 4.4. Niech estymowana na odcinku $X = [0, \pi]$ regresja trygonometryczna ma postać

$$\bar{y}(x) = a_0 + \sum_{j=1}^l a_j \cos(jx). \quad (4.5.26)$$

Przez τ_j oznaczmy pierwiastki równania

$$(1-x^2) \frac{dP_l(x)}{dx} = 0,$$

w którym $P_l(x)$ są wielomianami Legendre'a stopnia l -tego. Plan D-optimalny dla estymacji regresji (4.5.26) skupiony jest w $r = l + 1$ punktach

$$x_j^* = \arccos(\tau_j), \quad j = 1, 2, \dots, r,$$

którym przypisuje się jednakowe wagi $1/r$.

Nietrudno zauważyć, że omawiane w twierdzeniu 4.5 plany D-optimalne dla regresji trygonometrycznej zawierającej zarówno składniki sinusowe i kosinusowe jest istotnie różne od planów dla regresji (4.5.26). Powodem jest bliski związek tej ostatniej z wielomianami Czebyszewa. Dowód tego twierdzenia podano w [105].

Twierdzenie 4.5. Niech trygonometryczna funkcja regresji, określona na odcinku $X = [0, 2\pi]$, ma postać:

$$\bar{y}(x) = \alpha_0 + \sum_{j=1}^l \alpha_j \cos(jx) + \sum_{j=1}^l \beta_j \sin(jx). \quad (4.5.27)$$

Plany D-optimalne dla estymacji $r = (2l + 1)$ -wymiarowego wektora parametrów $\alpha_0, \alpha_1, \dots, \beta_1, \dots, \beta_l$ charakteryzują następujące stwierdzenia.

1. Planem D-optimalnym jest (nierealizowalny) plan ciągły $\xi^*(dx) = \frac{1}{2\pi} dx$, odpowiadający mierze o gęstości $1/(2\pi)$ na $[0, 2\pi]$.

2. Jeśli wybierzemy liczbę punktów planu $m \geq 2l + 1$ i punkty nośnika rozmieszczone będą tak, że spełniony jest warunek

$$x_{j+1} - x_j = \frac{2\pi}{m}, \quad j = 1, 2, \dots, m-1, \quad (4.5.28)$$

to plan, który przypisuje tym punktom oraz dowolnemu punktowi $x_0 \in X$ wagi $1/m$ jest planem D -optymalnym dla estymacji parametrów funkcji (4.5.27).

3. Plan opisany w punkcie 1 jest także D -optymalny dla estymacji każdej funkcji postaci (4.5.27) z $l' < l$ liczbą harmonicznych.

Własność 3 D -optymalnych planów dla regresji trygonometrycznej ma znaczenie wówczas, gdy nie znamy dokładnie liczby harmonicznych potrzebnych do dostatecznie dokładnej aproksymacji nieznannej funkcji. Można wówczas zastosować plan dla regresji o większej liczbie harmonicznych, a następnie zredukować złożoność modelu bez utraty optymalności planu. Należy jednak zwrócić uwagę na to, że własność ta nie zachodzi dla modeli innych niż trygonometryczne, a wśród modeli rozpinanych przez funkcje trygonometryczne własność tę mają tylko te o postaci (4.5.27).

Poniżej zestawimy znane plany optymalne w kostce wielowymiarowej dla regresji liniowej wielu zmiennych o postaci

$$\bar{y}(x) = a^{(0)} + \sum_{k=1}^s a^{(k)} x^{(k)}, \quad (4.5.29)$$

gdzie

$$X = \underbrace{[-1, 1] \times \dots \times [-1, 1]}_{s\text{-razy}}, \quad s \geq 1. \quad (4.5.30)$$

Potrzebne nam będzie klasyczne pojęcie planu czynnikowego (por. [62]).

Definicja 4.7. Pełnym planem czynnikowym na dwóch poziomach nazywamy plan, który nakazuje wykonać eksperymenty we wszystkich punktach wierzchołkowych kostki (4.5.30).

Dalej pełnym planem czynnikowym na dwóch poziomach nazywać będziemy także plan $\xi^* \in \Xi(X)$, który przypisuje jednakowe wagi $1/2^s$ wszystkim wierzchołkom kostki X .

Twierdzenie 4.6. Załóżmy stałość wariancji zakłóceń w X .

1. Pełen plan czynnikowy ξ^* jest planem D -optymalnym w zadaniu estymacji parametrów modelu (4.5.29).
2. Macierz informacyjna tego planu jest macierzą diagonalną.
3. Plan ξ^* jest także A - i E -optymalny (por. [105]).

Tezę 2 sprawdza się bezpośrednim rachunkiem. Dowód tezy 1 sprowadza się zatem do maksymalizacji funkcji $\phi(\xi^*, x)$ na kostce, co jest bardzo łatwe, gdyż dla diagonalnej macierzy informacyjnej funkcja ta jest sumą kwadratów poszczególnych zmiennych.

CZĘŚĆ II

Plany optymalne dla modeli wielowymiarowych

5. Numeryczne poszukiwanie planów optymalnych

Zagadnienie numerycznego poszukiwania planów optymalnych przedstawimy posługując się kryterium $\Phi(\cdot) \stackrel{\text{def}}{=} \ln \det(\cdot)$, czyli D-optymalności, gdyż metody dla innych kryteriów są dość podobne.

5.1. Ogólny algorytm gradientowy dla planów D-optymalnych

W podrozdziale tym opiszemy klasyczną – pochodzącą z lat sześćdziesiątych XX wieku – metodę Wynna–Fedorowa.

Motywacje konstrukcji algorytmu

Przypuśćmy, że wybraliśmy plan ξ_0 , który chcemy poprawić, w tym sensie, że chcemy znaleźć plan ξ_1 , który poprawi wartość Φ . Rozważmy rodzinę planów $(1 - \alpha)\xi_0 + \alpha\xi_{\text{popr}}$, $0 \leq \alpha \leq 1$, gdzie ξ_{popr} jest pewnym planem na X , o którym sądzimy, że poprawi wartość kryterium Φ . Planom tej rodziny odpowiada funkcja jednej zmiennej α

$$\varpi(\alpha) \stackrel{\text{def}}{=} \Phi((1 - \alpha)\xi_0 + \alpha\xi_{\text{popr}}). \quad (5.1.1)$$

Pochodna tej funkcji względem α , obliczona dla $\alpha = 0+$, ma postać (por. (4.2.9))

$$\frac{d}{d\alpha} \Phi[(1 - \alpha)\xi_0 + \alpha\xi_{\text{popr}}] \Big|_{\alpha=0+} = -r + \text{tr} \left[M^{-1}(\xi_0) M(\xi_{\text{popr}}) \right]. \quad (5.1.2)$$

Wartość powyższa wskazuje jak szybko zmienia się Φ w otoczeniu ξ_0 , jeśli zmiana zachodzi w kierunku planu ξ_{popr} . Próba znalezienia „najlepszego” planu ξ_{popr} byłaby co najmniej równie złożona jak wyjściowe zadanie. Jednakże, możemy ograniczyć klasę planów ξ_{popr} do planów skupionych w jednym punkcie, który oznaczmy przez $x_{\text{popr}} \in X$. Wówczas, z (5.1.2) otrzymamy

$$\frac{d}{d\alpha} \Phi[(1 - \alpha)\xi_0 + \alpha\xi_{\text{popr}}] \Big|_{\alpha=0+} = -r + \sigma^{-2}(x_{\text{popr}}) v^T(x_{\text{popr}}) M^{-1}(\xi_0) v(x_{\text{popr}}).$$

Możemy teraz szukać takiego punktu $x_{\text{popr}} \in X$, który wprowadzony do planu ξ_0 najbardziej zwiększy $\sigma^{-2}(x_{\text{popr}}) v^T(x_{\text{popr}}) M^{-1}(\xi_0) v(x_{\text{popr}})$, co prowadzi do lokalnie największej poprawy wartości kryterium Φ . Rozważania te prowadzą do metody Wynna–Fedorowa, którą opisano na następnych stronach. Jak się okazuje, operacja znajdowania najlepszej wartości x_{popr} jest najbardziej czasochłonnym

elementem tej metody w zastosowaniu do problemów wielowymiarowych. W następnych podrozdziałach pokażemy, że znajdowanie dokładnej wartości x_{popr} , która maksymalizuje $\sigma^{-2}(x_{\text{popr}}) v^T(x_{\text{popr}}) M^{-1}(\xi_1) v(x_{\text{popr}})$, nie jest konieczne – wystarczy znaleźć „zadowalające” x_{popr} .

Opis metody

Metoda Wynna–Fedorowa ulepszania planu polega na dodawaniu do jego nośnika punktów z odpowiednio dobraną wagą oraz modyfikacji wag punktów wcześniej włączonych do planu.

Zakładamy, że dany jest model $(X, v(x), \sigma^2 w(x))$, a zadanie polega na znalezieniu przybliżenia planu D- optymalnego. W algorytmie Wynna-Fedorowa punktem startowym jest wybór planu $\xi_0 \in \Xi$, którego macierz informacyjna jest nieosobliwa. Wybieramy także dokładność $\epsilon > 0$ dopuszczalnego niespełnienia warunku optymalności po zatrzymaniu algorytmu.

Przyjmijmy, że indeks k zlicza iteracje poprawiania planu początkowego (na początku kładziemy $k = 0$). Obliczamy funkcję ϕ dla planu ξ_k

$$\phi(\xi_k, x) = (\sigma^2 w(x))^{-1} v^T(x) M^{-1}(\xi_k) v(x) \quad (5.1.3)$$

i znajdujemy punkt

$$x_k = \arg \max_{x \in X} \phi(\xi_k, x). \quad (5.1.4)$$

Użyta powyżej operacja $\arg \max$ oznacza, że szukamy tej wartości argumentu funkcji, dla której osiągnane jest maksimum. Wobec założonej zwartości zbioru X i ciągłości funkcji ϕ istnienie takiego argumentu jest zagwarantowane.

ξ_k uznajemy za dostatecznie dokładną aproksymację planu optymalnego, jeśli

$$\phi(\xi_k, x_k) / r < 1 + \epsilon. \quad (5.1.5)$$

W przeciwnym razie, obliczamy wagę α_k przypisywaną punktowi x_k zgodnie ze wzorem

$$\alpha_k = \frac{\phi(\xi_k, x_k) - r}{(\phi(\xi_k, x_k) - 1) r}. \quad (5.1.6)$$

Tak obliczamy $\alpha_k \in (0, 1)$, gdy warunek stopu (5.1.5) nie jest spełniony. Następnie poprawiamy przybliżenie planu optymalnego następująco:

$$\xi_{k+1} = (1 - \alpha_k) \xi_k + \alpha_k \delta(x_k), \quad (5.1.7)$$

gdzie $\delta(x_k)$ jest planem skupionym w jednym punkcie x_k . W planie $\delta(x_k)$ punktowi x_k formalnie przypisujemy wagę 1, lecz po wykonaniu operacji (5.1.7) punkt ten otrzymuje wagę α_k . Teraz numer iteracji k zwiększany jest o jeden i ponownie wyznaczany jest punkt wprowadzany do planu zgodnie z (5.1.4).

Zbieżność ciągu wartości kryterium

Można wykazać (por. [210], [35], [105]), że jeżeli w omawianej metodzie pominąć sprawdzanie warunku (5.1.5), to generuje ona nieskończony ciąg planów ξ_k , dla którego zachodzi

$$\lim_{k \rightarrow \infty} \det[M(\xi_k)] = \det[M(\xi^*)]. \quad (5.1.8)$$

Dowodu tego rezultatu nie przedstawiamy, gdyż w następnych podrozdziałach przedstawimy dowód zbieżności w nieco ogólniejszej sytuacji.

Warto zaznaczyć, że szybkość zbieżności algorytmu Wynna–Fedorova nie jest zbyt duża. W istocie należy on do tzw. algorytmów I rzędu, tzn. takich, które wykorzystują informacje tylko o gradiencie optymalizowanej funkcji. W praktyce algorytm ten dość szybko znajduje nośnik planu zbliżonego do optymalnego, natomiast precyzyjny dobór wag jest długim procesem, który można przyspieszyć włączając w każdej iteracji procedurę optymalizacji wag. Procedurę tę opiszemy w ostatniej części niniejszego rozdziału.

Bez istotnej zmiany asymptotycznych własności metody Wynna–Fedorova wagę punktu wprowadzanego do planu można obliczać zgodnie ze wzorem

$$\alpha_k = \frac{1}{k+1}, \quad (5.1.9)$$

zamiast stosować (5.1.6). Dopuszczalny jest nawet taki wybór dowolnego ciągu $0 < \alpha_k < 1$, który spełnia warunki

$$\sum_{k=0}^{\infty} \alpha_k = \infty, \quad \sum_{k=0}^{\infty} \alpha_k^2 < \infty. \quad (5.1.10)$$

Wybór zgodnie z równaniem (5.1.6) ma walor lokalnej optymalności, gdyż wartość ta maksymalizuje

$$\det M[(1-\alpha)\xi_k + \alpha\xi(x_k)]$$

względem α .

Algorytm Wynna–Fedorova formalnie nie zmienia się w przypadku poszukiwania planów dla estymacji regresji o wielu zmiennych, jednakże nakład obliczeń drastycznie rośnie wraz ze wzrostem wymiaru przestrzeni planowania. Powodem tego zjawiska jest wybór punktu wprowadzanego do planu zgodnie z (5.1.4)

Metoda Wynna–Fedorova jako algorytm optymalizacji wag planu skupionego w zadanych punktach

Jednym ze sposobów zmniejszenia trudności obliczeniowych związanych z maksymalizacją globalną w (5.1.4) jest wytypowanie skończonego zbioru, powiedzmy

$\Omega = \{\omega_1, \omega_2, \dots, \omega_m\}$, $m > 1$ punktów przestrzeni, w których może znajdować się nośnik planu optymalnego,¹ a następnie optymalizacji wag przypisanych tym punktom. Otrzymany w ten sposób plan nie będzie na ogół planem optymalnym, chyba że uda nam się dobrać zbiór Ω tak, by rzeczywiście zawierał nośnik planu optymalnego. W przeciwnym razie dostaniemy plan optymalny, ale tylko na tym zbiorze punktów. Zauważmy, że do planu takiego odnosi się twierdzenie Kiefera–Wolfowitza i podane w nim warunki optymalności, jeśli jako zbiór planowania X wstawimy Ω .

Opiszemy teraz w skrócie wersję algorytmu Wynna–Fedorova dostosowaną do skończonego zbioru Ω . Wynikiem jego działania będą wagi, które przypisać należy punktom z Ω . Inną procedurę o podobnym zakresie zastosowań przedstawimy w podrozdziale 5.3.

Krok 0. Przypisujemy punktom ω_i wagi początkowe $p_i^{(0)} \geq 0$, $i = 1, 2, \dots, m$, $\sum_{i=1}^m p_i^{(0)} = 1$ i na ich podstawie tworzymy plan początkowy ξ_0 . Ustawiamy licznik iteracji $k = 0$. Wybieramy żadaną dokładność $\epsilon > 0$ spełnienia warunku optymalności.

Krok 1. Dla $j = 1, 2, \dots, m$ obliczamy wartości

$$\zeta_j^{(k)} \stackrel{\text{def}}{=} \phi(\xi_k, \omega_j) = (\sigma^2 w(\omega_j))^{-1} v^T(\omega_j) M^{-1}(\xi_k) v(\omega_j) \quad (5.1.11)$$

i znajdujemy indeks j^* , dla którego

$$j^* = \arg \max_{j \in \{1, 2, \dots, m\}} [\zeta_j^{(k)}]. \quad (5.1.12)$$

Uwaga: w kolejnych iteracjach indeks j^* będzie na ogół różny. Nie zaznaczamy tego faktu w notacji, aby nie komplikować zapisów.

Krok 2. Sprawdzamy warunek

$$\zeta_{j^*}^{(k)} / r < 1 + \epsilon. \quad (5.1.13)$$

Gdy jest on spełniony, to zatrzymujemy obliczenia i ξ_k uznajemy za dostatecznie dokładną aproksymację planu optymalnego.

Krok 3. Jeśli warunek (5.1.13) nie jest spełniony, to obliczamy wagę α_k przypisywaną punktowi ω_{j^*} zgodnie ze wzorem

$$\alpha_k = \frac{\zeta_{j^*}^{(k)} - r}{(\zeta_{j^*}^{(k)} - 1) r}. \quad (5.1.14)$$

¹ W najprostszym przypadku wytypowanie takie może polegać na zastosowaniu siatki równoodległych punktów w przestrzeni planowania.

Krok 4. Obliczamy nowe przybliżenie planu

$$\xi_{k+1} = (1 - \alpha_k) \xi_k + \alpha_k \delta(\omega_{j^*}), \quad (5.1.15)$$

gdzie $\delta(\omega_{j^*})$ jest planem skupionym w jednym punkcie ω_{j^*} . Jeśli numer iteracji k nie przekroczył maksymalnej, założonej wcześniej wartości, to zwiększamy k o jeden i wracamy do kroku 1.

Algorytm ten generuje ciąg planów o rosnących wartościach wyznacznika macierzy informacyjnej. Ciąg ten jest zbieżny, lecz $\lim_{k \rightarrow \infty} \det M(\xi_k) \leq \det M(\xi^*)$. W istocie ciąg ten jest zbieżny do wartości wyznacznika macierzy informacyjnej planu, który jest najlepszy wśród planów o nośnikach zawartych w zbiorze Ω .

Uwagi na temat modyfikacji metody Wynna–Fedorova

W literaturze (por. [6], [29]) zaproponowano kilka modyfikacji algorytmu Wynna–Fedorova, mających na celu przyspieszenie zbieżności. Opisane poniżej modyfikacje odnoszą się także do opisanej w poprzednim podrozdziale wersji algorytmu poprawy wag z tym tylko, że zamiast $\phi(\xi_k, x)$, $x \in X$ podstawiać należy $\zeta_j^{(k)}$.

1. Dołączanie do planu w danej iteracji wszystkich punktów, dla których zachodzi $\max_{x \in X} \phi(\xi_k, x)$. Zabieg taki formalnie zmniejsza liczbę iteracji potrzebnych do osiągnięcia zadanej dokładności ϵ , lecz znalezienie wszystkich (większości) punktów, w których osiągane jest maksimum nie jest zadaniem łatwym, nawet w przypadku funkcji jednej zmiennej.
2. Usuwanie z bieżącego planu punktów, dla których osiągnęte jest

$$\min_{x \in \text{supp}(\xi_k)} \phi(\xi_k, x),$$

gdzie przez $\text{supp}(\xi_k)$ oznaczono nośnik planu ξ_k . Zauważmy, że sytuacja jest tutaj inna niż w poprzednim komentarzu, gdyż do obliczenia i porównania mamy tylko wartości funkcji $\phi(\xi_k, x)$ w znanej i skończonej liczbie punktów. Dlatego warto tę modyfikację stosować. Nie należy się przy tym obawiać, że bezpowrotnie utracimy punkt, który powinien się znaleźć w nośniku planu, gdyż po ewentualnym usunięciu takiego punktu zostanie on i tak w dalszych iteracjach wprowadzony do planu. Teoretycznie istnieje możliwość tzw. pętlenia się algorytmu, tzn. usuwania i wprowadzania do planu tego samego punktu co kilka iteracji. Przed możliwością tą można się zabezpieczyć stosując nierówność wyprowadzoną w pracy [109]. Nierówność ta pozwala stwierdzić, który punkt nośnika aktualnego planu na pewno nie należy do nośnika planu optymalnego (mimo, że tego ostatniego nie znamy). Takie bezpieczne postępowanie polegałoby więc na usuwaniu punktu, dla którego osiągnęte jest $\min_{x \in \text{supp}(\xi_k)} \phi(\xi_k, x)$, ale tylko wówczas, gdy wspomniana nierówność wskazuje, że punkt ten nie należy do optymalnego planu.

Uwagi bibliograficzne

Na przełomie lat sześćdziesiątych i siedemdziesiątych opublikowane zostały pierwsze algorytmy numeryczne znajdowania przybliżeń planów D- optymalnych (por. [210], [35]). W literaturze znane są one pod nazwą metody Wynna i Fedorova. Następne lata przyniosły wiele ulepszeń tego algorytmu i uogólnienia na inne kryteria (por. [211], [105], [6], [27]).

Znane są szybsze algorytmy wykorzystujące macierz drugich pochodnych logarytmu wyznacznika macierzy informacyjnej, czyli procedury II rzędu (por. [211], [28]). Nie przedstawiamy ich tutaj, gdyż są one znacznie bardziej kłopotliwe w implementacji, a uzyskiwane przyspieszenie szybkości zbieżności dotyczy głównie poprawiania wag planu. Warto wskazać również na podejście oparte na dualności zadań programowania wypukłego zaprezentowane w monografii [16] (s. 384–391), które może być stosowane także do nieróżniczkowalnych kryteriów, takich jak wskaźnik E- optymalności.

Jak wynika z doświadczeń autora, trudności w numerycznym znajdowaniu planów optymalnych polegają głównie na znalezieniu lokalizacji punktów nośnika planu, a pod tym względem algorytm Wynna–Fedorova jest bardzo dobry – w większości przypadków lokalizuje on punkty nośnika w kilkunastu pierwszych iteracjach. Spostrzeżenie to jest podstawą rozważań przedstawionych w następnym podrozdziale.

5.2. Algorytm selektywnych poszukiwań losowych

W pracy autora [139] zaproponowano metodę selektywnych poszukiwań losowych D- optymalnych planów eksperymentu. Zarys tej metody poprzedzimy przypomnieniem znanej metody odrzucania (*rejection method*) generowania realizacji zmiennej losowej o zadanej gęstości.

Przekleństwo wymiarowości i wieloekstremalności

Najtrudniejszym w praktycznej implementacji elementem metody Wynna–Fedorova jest (5.1.4). Na realizację tej operacji składa się:

- obliczenie macierzy informacyjnej dla aktualnego planu,
- znalezienie jej odwrotności, po to, by móc uformować funkcję ϕ wariacji oceny regresji,
- znalezienie punktu (lub punktów, w przypadku gdy stosowana jest modyfikacja opisana w komentarzu 5.1), dla którego osiągnane jest $\max_{x \in X} \phi(\xi_k, x)$.

Pierwsze dwie operacje można znacznie przyspieszyć, stosując odpowiednie wzory algebry na odwracanie macierzy o specjalnej strukturze. Natomiast operacja maksymalizacji ϕ z ograniczeniem $x \in X$ jest najbardziej pracochłonnym podproblemem tego algorytmu, gdyż funkcja ϕ jest z natury swej wielomodal-

na. Jak wiadomo, numeryczne znajdowanie globalnego maksimum funkcji wielu zmiennych jest jednym z najtrudniejszych problemów, czy wręcz wyzwania intelektualnych. Z tego powodu przedstawiamy tutaj inne podejście, natomiast szeroką klasę problemów wielowymiarowych proponujemy rozwiązywać dekomponując je do modeli addytywnych lub modeli z pełnym zestawem interakcji. W następnym podrozdziale przedstawiamy, zaproponowaną przez autora [139], metodę obchodzącą ten problem w ten sposób, że zamiast poszukiwać $\max_{x \in X} \phi(\xi_k, x)$, znajdować będziemy (poprzez losowanie z odrzucaniem) punkt o dostatecznie dużej wartości $\phi(\xi_k, x)$ i to on będzie wprowadzany do planu.

Opis metody selektywnych poszukiwań losowych

Podstawą metody jest proste spostrzeżenie, że funkcja $\phi(\xi_k, x)$ nieujemna i całkowalna na X dla dowolnego planu ξ_k , którego macierz informacyjna jest nieosobliwa. Wynika stąd, że funkcja

$$f_k(x) \stackrel{\text{def}}{=} \frac{\phi(\xi_k, x)}{\int_X \phi(\xi_k, x) dx} \quad (5.2.16)$$

może być interpretowana jako gęstość prawdopodobieństwa pewnej zmiennej losowej. Łatwo wykazać, że funkcja ta może być przedstawiona w postaci

$$f_k(x) = \frac{v^T(x) M^{-1}(\xi_k) v(x)}{\text{tr}[M^{-1}(\xi_k) V]}, \quad (5.2.17)$$

gdzie $r \times r$ macierz V zdefiniowana jest następująco: $V = \int_X v(x) v^T(x) dx$. W celu zapewnienia efektywności proponowanego algorytmu macierz tę warto wstępnie obliczyć (analitycznie lub numerycznie).

W celu wygenerowania zmiennej losowej x'_k o gęstości rozkładu prawdopodobieństwa f_k użyjemy metodę odrzucania [26]). Załóżmy, że znamy ograniczenie od góry $0 < \mu_k < \infty$ funkcji $f_k(x) \leq \mu_k$, $x \in X$. Podkreślmy, że μ_k nie musi być maksimum $f_k(x)$ w zbiorze X , lecz im jest ono bliższe tej wartości, tym większa jest efektywność algorytmu odrzucania.

Algorytm generowania liczb losowych metodą odrzucania

Krok 1. Niech $(x', t') \in X \times [0, \mu_k]$ będzie $(s + 1)$ -wymiarowym wektorem wylosowanym zgodnie z rozkładem równomiernym w zbiorze $X \times [0, \mu_k]$.

Krok 2. Jeśli $f_k(x') \leq t'$, to jako liczbę losową o rozkładzie f_k przyjmij $x'_k = x'$.

W przeciwnym razie odrzuć parę (x', t') i powtórz krok 1.

Zgodnie z twierdzeniem von Neumanna (por. [26]), tak wygenerowane x'_k można interpretować jako wektor losowy o gęstości rozkładu prawdopodobieństwa f_k .

Opis algorytmu selektywnych losowych poszukiwań planów D-optimalnych

Krok 0. Obliczyć macierz $V = \int_X v(x)v^T(x)dx$ i wybrać plan początkowy ξ_0 o nieosobliwej macierzy informacyjnej. Ustawić licznik iteracji $k = 0$.

Krok 1. Obliczyć f_k zgodnie ze wzorem (5.2.17) i wylosować x'_k o rozkładzie f_k (stosując, np. algorytm odrzucania).

Krok 2. Jeśli $\phi(\xi_k, x'_k) > r$, podstaw $x_k = x'_k$ i przejdź do kroku 3. W przeciwnym razie odrzuć x'_k i powtórz losowanie opisane w kroku 1.

Krok 3. Utwórz nowy plan

$$\xi_{k+1} = (1 - \alpha_k) \xi_k + \alpha_k \delta(x_k), \quad (5.2.18)$$

w którym α_k wybrane jest następująco:

$$\alpha_k = \frac{\phi(\xi_k, x_k) - r}{(\phi(\xi_k, x_k) - 1)r} \quad (5.2.19)$$

podstaw $k + 1$ w miejsce k i przejdź do kroku 1.

Praktyczne aspekty implementacji powyższego algorytmu omawiane będą w dalszych podrozdziałach.

Odnotujmy, że powyższy algorytm różni się od algorytmów losowego szukania maksimum funkcji $\phi(\xi_k, x)$. Podstawowa różnica polega na tym, że x'_k losowane jest raz w każdej iteracji, a losowanie odbywa się zgodnie z gęstością $f_k(x)$, która jest proporcjonalna do $\phi(\xi_k, x)$. x'_k losowane jest zatem z większym prawdopodobieństwem w obszarach, gdzie wartości $\phi(\xi_k, x)$ są większe. Po wylosowaniu x'_k następuje ponowna selekcja, gdyż punkt ten wprowadzany jest do planu tylko wówczas, gdy spełniony jest warunek $\phi(\xi_k, x_k) > r$. Wymienione fakty uzasadniają nazwę: algorytm selektywnych poszukiwań losowych.

Dowód zbieżności algorytmu

Przejdziemy teraz do sformułowania podstawowych własności omawianego algorytmu. Dowód następnego lematu wzorowany jest na odpowiednim rezultacie z [35].

Lemat 5.1. *Niech ξ_k będzie planem o nieosobliwej macierzy informacyjnej. Załóżmy, że punkt $x_k \in X$ wybrano tak, by spełniony był warunek $\phi(\xi_k, x_k) > r$. Niech długość kroku poszukiwań α_k wybrana będzie zgodnie z (5.2.19), wówczas dla następującego planu*

$$\xi_{k+1} = (1 - \alpha_k) \xi_k + \alpha_k \delta(x_k)$$

zachodzi nierówność $\det(M(\xi_{k+1})) > \det M(\xi_k)$.

Dowód. Potrzebny nam będzie następujący wzór (por. [35] lub [36], wzór (3.1.10) s. 47)

$$\det M(\xi_{k+1}) = \left[\frac{\phi(\xi_k, x_k)}{r} \right]^r \cdot \left[\frac{r-1}{\phi(\xi_k, x_k) - 1} \right]^{r-1} \det M(\xi_k). \quad (5.2.20)$$

Zwracamy uwagę na to, że wzór (5.2.20) zachodzi wówczas, gdy α_k obliczane jest zgodnie ze wzorem (5.2.19) i spełniony jest warunek $\phi(\xi_k, x_k) > r$. Zdefiniujmy funkcję $\omega(t) = \left(\frac{t}{r}\right)^r \cdot \left(\frac{r-1}{t-1}\right)^{r-1}$. Łatwo sprawdzić, że $\omega(r) = 1$ oraz

$$\omega'(t) = t^{-1}(r-1)^{-1} \cdot (t-r) \left(\frac{r-1}{t-1}\right)^r \left(\frac{t}{r}\right)^r \quad (5.2.21)$$

Można też sprawdzić, że gdy $t > r$ to $\omega'(t) > 0$. Dla $t > r$ funkcja $\omega(t)$ jest ściśle rosnąca. Własność ta oraz (5.2.20) implikują $\det M(\xi_{k+1}) > \det M(\xi_k)$, co kończy dowód. •

Możemy teraz sformułować podstawowy rezultat dotyczący algorytmu selektywnych poszukiwań losowych.

Twierdzenie 5.1. *Dla dowolnego planu startowego ξ_0 o nieosobliwej macierzy informacyjnej ciąg planów ξ_k , $k = 1, 2, \dots$ generowany przez algorytm selektywnych poszukiwań losowych spełnia, z prawdopodobieństwem jeden, następującą zależność $\lim_{k \rightarrow \infty} \det M(\xi_k) = \det M(\xi^*)$, gdzie ξ^* jest planem D- optymalnym.*

Dowód. Niech ξ_k będzie bieżącym planem wygenerowanym przez algorytm selektywnych poszukiwań losowych. Rozważymy dwa przypadki, odnoszące się do zachowania algorytmu w kroku 2.

Przypadek I. Przypuśćmy, że w kroku 2) wygenerowany został punkt x_k taki, że $\phi(\xi_k, x_k) > r$. Wówczas, zgodnie z Lematem 5.1, zachodzi $\det M(\xi_{k+1}) > \det M(\xi_k)$. W przypadku I generowany jest zatem ściśle rosnący ciąg liczbowy $\det M(\xi_k)$, $k = 1, 2, \dots$. Ciąg ten jest ograniczony od góry przez wskaźnik jakości planu optymalnego $\det M(\xi^*)$. Dla takiego ciągu istnieje granica skończona

$$\lim_{k \rightarrow \infty} \det M(\xi_k).$$

Pozostaje pokazać, że ciąg ten nie osiąga granicy niższej niż wartość kryterium dla planu optymalnego, czyli że zachodzi $\lim_{k \rightarrow \infty} \det M(\xi_k) = \det M(\xi^*)$.

Jeśli przypuścimy, że tak nie jest, to zachodzić musi

$$\lim_{k \rightarrow \infty} \det M(\xi_k) < \det M(\xi^*).$$

Powtarzając teraz odpowiedni fragment dowodu z pracy [210] dochodzimy do wniosku, że przypuszczenie $\lim_{k \rightarrow \infty} \det M(\xi_k) < \det M(\xi^*)$ implikuje

$$\sup_{x \in X} \phi(\xi_k, x) > r. \quad (5.2.22)$$

Z założonej ciągłości funkcji $v(x)$ wynika ciągłość $\phi(\xi_k, x)$ na zwartym (bo domkniętym i ograniczonym) zbiorze X . Wynika stąd, że supremum w wyrażeniu po lewej stronie (5.2.22) osiągane jest w pewnym punkcie $\hat{x} \in X$ i w punkcie tym $\phi(\xi_k, \hat{x}) > r$. Ponadto, ciągłość funkcji $\phi(\xi_k, x)$ implikuje, że nierówność $\phi(\xi_k, x) - r > 0$ zachodzi także w pewnym otoczeniu \hat{x} , które ma niezerową miarę Lebesgue'a. Oznaczmy to otoczenie przez $U(\hat{x})$. Z nierówności $\phi(\xi_k, x) - r > 0$ wynika, że na zbiorze $U(\hat{x})$ także gęstość $f_k(x)$ jest dodatnia, a co za tym idzie dodatnie jest także prawdopodobieństwo wylosowania punktu z $U(\hat{x})$ w kroku 1 algorytmu. Prawdopodobieństwo to rośnie do jedności gdy algorytm cyrkuluje między krokami 1 i 2, więc z prawdopodobieństwem dążącym do jeden wybrany będzie punkt z $U(\hat{x})$, co prowadzi do sprzeczności z przypuszczeniem, że $\lim_{n \rightarrow \infty} \det M(\xi_n) < \det M(\xi^*)$.

Przypadek II. Algorytm tworzy nieskończoną pętlę między krokami 1 i 2. Takie zachowanie oznacza, że algorytm nie jest w stanie znaleźć punktu, w którym $\phi(\xi_k, x) > r$ i własność ta zachodzi z prawdopodobieństwem jeden, względem miary o gęstości $f_k(x)$. Oznacza to, że zbiór $B_k \stackrel{\text{def}}{=} \{x : \phi(\xi_k, x) > r\}$ ma miarę Lebesgue'a równą zero. Można wykazać ponadto, że zbiór ten jest pusty. Gdyby tak nie było, wówczas w pewnym punkcie $\tilde{x} \in B_k$ zachodziłoby $\phi(\xi_k, \tilde{x}) > r$. Co więcej nierówność ta zachodziłaby nie tylko w punkcie \tilde{x} , ale także w pewnym jego otoczeniu o dodatniej mierze Lebesgue'a (ponownie na skutek ciągłości funkcji $\phi(\xi_k, \cdot)$). Ale stoi to w sprzeczności z tym, że zbiór B_k ma zerową miarę Lebesgue'a, co kończy dowód, gdyż $B_k = \emptyset$ pozwala zastosować twierdzenie Kiefera–Wolfowitza i wywnioskować optymalność planu ξ_k . •

Zanim przejdziemy do omawiania szczegółów implementacyjnych, przedstawimy przykład ilustrujący ideę algorytmu.

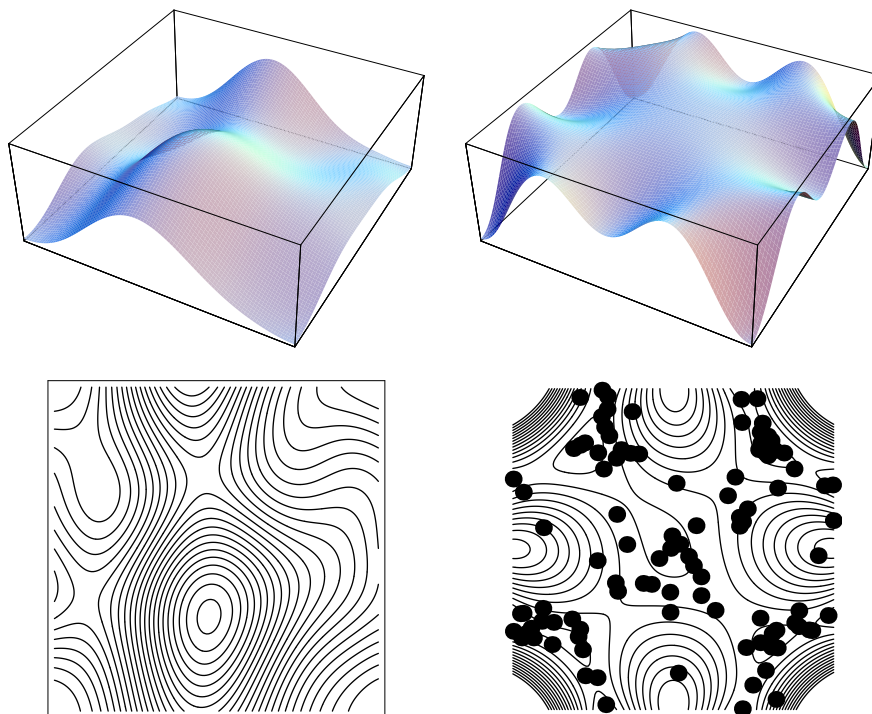
Niech funkcja regresji rozpięta będzie przez zestaw funkcji

$$v(x^{(1)}, x^{(2)}) = \left[1, \sin(\pi x^{(1)}), \sin(\pi x^{(2)}), \cos(\pi x^{(1)}), \cos(\pi x^{(2)}), \sin(\pi x^{(2)} x^{(1)}) \right]^T,$$

określonych w $X = [-1, 1] \times [-1, 1]$. Poszukiwania rozpoczynamy od planu ξ_0 o postaci

$$\left[\begin{array}{ccccccc} \left[\begin{array}{c} -\frac{3}{4} \\ -\frac{3}{4} \\ \frac{1}{8} \end{array} \right] & \left[\begin{array}{c} \frac{1}{4} \\ \frac{7}{10} \\ \frac{5}{64} \end{array} \right] & \left[\begin{array}{c} -\frac{9}{10} \\ 0 \\ \frac{3}{64} \end{array} \right] & \left[\begin{array}{c} \frac{3}{4} \\ -\frac{3}{4} \\ \frac{1}{8} \end{array} \right] & \left[\begin{array}{c} \frac{3}{4} \\ \frac{3}{4} \\ \frac{1}{8} \end{array} \right] & \left[\begin{array}{c} -\frac{7}{10} \\ \frac{1}{4} \\ \frac{1}{8} \end{array} \right] & \left[\begin{array}{c} -\frac{3}{4} \\ \frac{3}{4} \\ \frac{1}{8} \end{array} \right] \end{array} \right],$$

który uzupełniono o punkty $[1, -1]$ i $[-1, 1]$ oba z wagami $1/8$. Planowi temu odpowiada gęstość prawdopodobieństwa $f_0(x)$ pokazana w lewym górnym rogu rysunku 5.1. W lewym dolnym rogu tego rysunku pokazano warstwicę tej gęstości. Następnie wykonano 100 iteracji algorytmu. W prawym górnym rogu rysunku 5.1



Rys. 5.1. Przykład: selektywne poszukiwanie – zmiany gęstości rozkładu poszukiwań

pokazano gęstość $f_{100}(x)$, a poniżej odpowiadające jej warstwy, na które nałożono punkty wprowadzane do planu ξ_{100} w poprzedzających 100 iteracjach. Na rysunku 5.2 zilustrowano przebieg zmian wyznacznika macierzy informacyjnej w kolejnych iteracjach. Dla zachowania przejrzystości przykładu obliczenia przerwano po 100 iteracjach, mimo że plan bliski optymalnemu nie został znaleziony. Jednakże rysunek 5.2 jest dość charakterystyczny dla omawianego algorytmu – w początkowej fazie następuje bardzo szybka poprawa planu startowego.

Warunek stopu

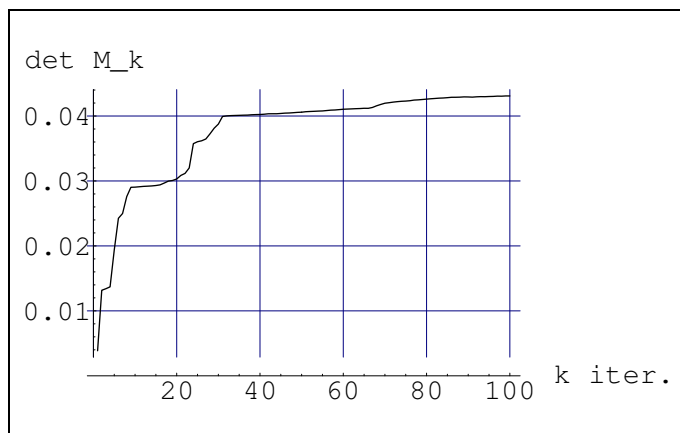
Analizując dowód Twierdzenia 5.1 stwierdzamy, że gdy ξ_k jest planem optymalnym, to algorytm wykonuje nieskończoną pętlę między krokiem 1 i krokiem 2 (przypadek II w dowodzie).

Spostrzeżenie to pozwala zaproponować następujący warunek zatrzymania algorytmu selektywnych poszukiwań losowych.

Warunek stopu:

Wybrać dostatecznie dużą liczbę naturalną N .

Jeśli w bezpośrednio po sobie następujących iteracjach N -krotnie w kroku 2 odrzu-



Rys. 5.2. Przykład: selektywne poszukiwanie planu – przebieg wyznacznika

cone będą punkty $x'_k, x'_{k+1} \dots$ i nastąpi powrót do kroku 1, to zatrzymaj algorytm, a ostatnio obliczony plan ξ_k przyjmij za przybliżenie planu optymalnego.

Teoretycznego uzasadnienia zaproponowanego warunku stopu dostarcza następujący lemat.

Lemat 5.2. *Jeśli ξ_k nie jest planem D -optymalnym i w algorytmie selektywnych poszukiwań losowych stosowany jest powyższy warunek stopu, to prawdopodobieństwo przedwczesnego zatrzymania algorytmu można uczynić dowolnie małym, jeśli wybierze się N dostatecznie duże.*

Przez przedwczesne zatrzymanie rozumiemy w teorii zatrzymanie w przypadku gdy ξ_k nie jest planem optymalnym. Tak też rozumiany jest ten termin w następującym dowodzie. Jednakże w praktyce N musi mieć skończoną wartość i dlatego po zatrzymaniu algorytmu z planem ξ_k możemy spodziewać się, że $\det M(\xi_k)$ jest bliskie $\det M(\xi^*)$, lecz nie możemy gwarantować równości $\det M(\xi_k) = \det M(\xi^*)$ po skończonej liczbie iteracji.

Dowód. Jeśli ξ_k nie jest planem optymalnym, to zbiór $B_k = \{x : \phi(\xi_k, x) > r\}$ nie jest pusty, jego miara Lebesgue'a jest zatem dodatnia (jako konsekwencja ciągłości $v(x), x \in X$). Oznaczmy $c_k = \int_X \phi(\xi_k, x) dx$. Ponieważ x'_k ma gęstość $f_k(x) = \phi(\xi_k, x)/c_k$, zatem

$$P \{x'_k \in B_k\} = c_k^{-1} \int_{B_k} \phi(\xi_k, x) dx > c_k^{-1} r \lambda(B_k), \quad (5.2.23)$$

gdzie $\lambda(\cdot)$ oznacza miarę Lebesgue'a zbioru. Poniższe prawdopodobieństwo szacujemy w schemacie Bernoulliego

$$P \{x'_k \notin B_k \text{ w każdej z } N \text{ kolejnych prób}\} \quad (5.2.24)$$

$$= (1 - P\{x'_k \in B_k\})^N < (1 - c_k^{-1} r \lambda(B_k))^N.$$

Wyrażenie po prawej stronie (5.2.24) możemy uczynić dowolnie małym przez wybór dostatecznie dużego N . Chcemy wykazać, że $c_k^{-1} r \lambda(B_k) < 1$. Aby udowodnić tę nierówność, wystarczy przeanalizować ciąg oszacowań:

$$c_k = \int_X \phi(\xi_k, x) dx \geq \int_{B_k} \phi(\xi_k, x) dx > r \lambda(B_k). \quad (5.2.25)$$

Mając to oszacowanie, możemy uczynić dowolnie małym wyrażenie po prawej stronie (5.2.24) przez wybór dostatecznie dużego N . •

Możliwe są i, zdaniem autora, na dalsze badania zasługują następujące modyfikacje algorytmu selektywnych poszukiwań losowych. Pierwsza możliwość to uzupełnienie algorytmu o deterministyczne, lokalne poprawianie punktów kolejno wprowadzanych do planu. Druga modyfikacja może polegać na wyostreniu lokalnych maksimum funkcji f_n przez zastąpienie jej wyrażeniem

$$\frac{\phi^\gamma(\xi_n, x)}{\int_X \phi^\gamma(\xi_n, x) dx},$$

gdzie $\gamma > 1$ jest parametrem sterującym wyostreniem. W dalszym etapie modyfikacji można uzmiennić parameter γ , uzależniając go od numeru iteracji i zwiększając jego wartość stopniowo, w sposób, odwrotny do stosowanego w algorytmach symulowanego wyżarzania.

Działanie algorytmu – przykłady symulacyjne

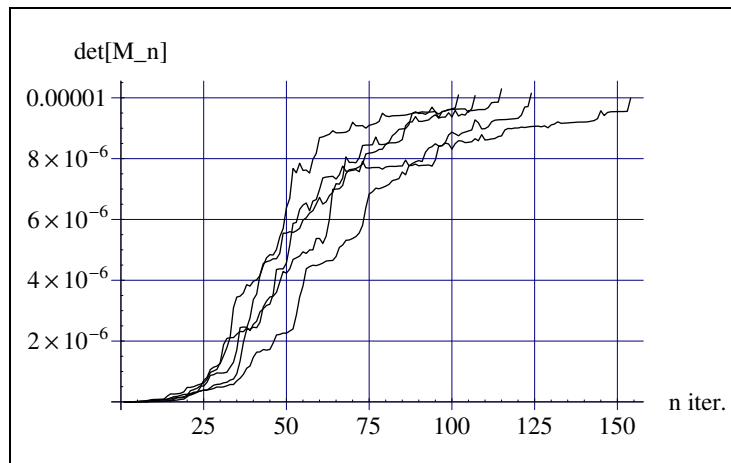
W badaniach symulacyjnych przedstawionych w tym podrozdziale zastosowano dwa ważne usprawnienia algorytmu:

- a) optymalizację wag aktualnego planu w każdej iteracji i odrzucanie punktów o bardzo małych wagach,
- b) „sklejanie” w jeden tych punktów planu, które położone są blisko siebie.

Modyfikacje te omówimy w następnych podrozdziałach, gdyż zakres ich zastosowań obejmuje także inne algorytmy.

Badania przeprowadzono na przykładach, których rozwiązania optymalne są znane, gdyż wtedy możliwa jest ocena jakości uzyskanego planu przybliżającego plan optymalny.

Symulacje 5.1. *Na rysunku 5.3 pokazano cztery przebiegi omawianego algorytmu. Jako przykład testowy wybrano zadanie planowania D-optymalnego dla estymacji parametrów funkcji regresji o dwóch zmiennych, oznaczonych x i t , rozpiętej na zestawie funkcji $\{1, t, t^2, x, tx, t^2x, x^2, tx^2, t^2x^2\}$. Jako obszar planowania wybrano $X = [-1, 1] \times [-1, 1]$. W każdym z czterech przebiegów poszukiwania rozpoczynano od tego samego planu startowego, a różnice przebiegów widoczne*



Rys. 5.3. Przykładowe przebiegi algorytmu selektywnych poszukiwań losowych
– opis w przykładzie symulacji 5.1

na wykresie są rezultatem losowej natury algorytmu. Zastosowano następujące parametry algorytmu:

$N = 10^4$ w regule zatrzymania algorytmu,

$\epsilon = \frac{25}{1000}$ w warunku odrzucania punktów o zbyt małych wagach,

$L = 250$ jako liczbę wewnętrznych iteracji optymalizacji wag (znaczenie tego parametru stanie się jasne po przejrzaniu następnego podrozdziału, gdzie algorytm optymalizacji wag został szczegółowo opisany),

$\eta = 0.05$ jako promień otoczeń, wewnątrz których punkty planu traktowane były jako nierozróżnialne i „sklejane” w jeden.

Jeden z planów znalezionych w tych poszukiwaniach pokazano w tabeli 5.1. Zauważyć można, że plan ten zawiera dwa punkty położone bardzo blisko siebie i punktu $(0, 1)$, który należy do nośnika planu optymalnego. Zwiększenie parametru η spowodowałoby utworzenie jednego punktu planu z dwóch wyżej wymienionych.

Chociaż naszym celem była prezentacja ogólnego zachowania się opisywanego algorytmu, a nie znalezienie planu optymalnego, to podajemy porównanie jakości jednego z otrzymanych planów z planem optymalnym. Wyznacznik macierzy informacyjnej dla znalezionej planu wyniósł 0.0000102939, natomiast wyznacznik ten dla planu optymalnego ma wartość 0.0000105725, co oznacza 97.4% efektywności względnej. Najgorszy plan uzyskany w tych czterech przebiegach miał efektywność 95%.

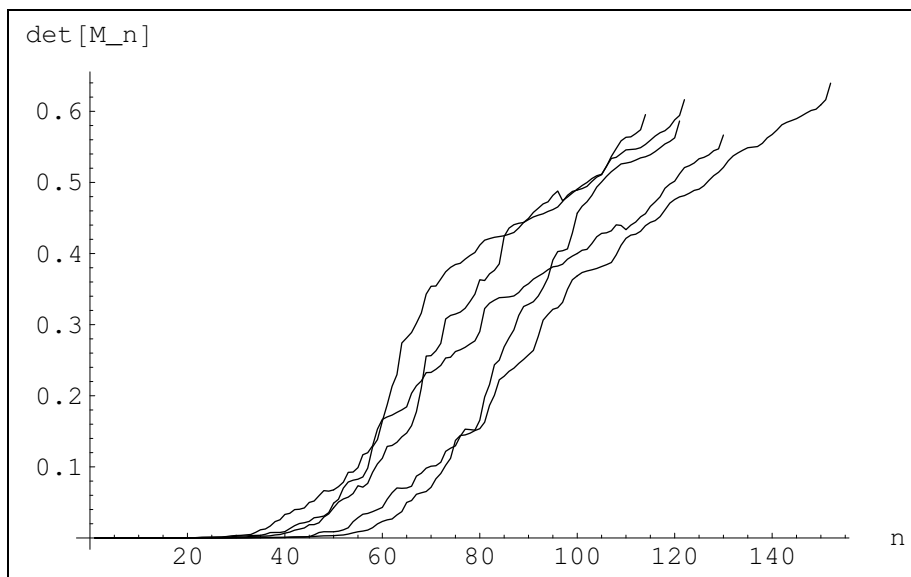
Punkty	p_i	Punkty	p_i
(+1, - 1)	0.11111	(+1, - 1)	1/9
(-1, - 1)	0.11111	(-1, - 1)	1/9
(-0.9995, + 0.0039)	0.111111	(+1, + 1)	1/9
(+0.0002, + 0.0008)	0.111111	(-1, + 1)	1/9
(+0.0159, - 0.9995)	0.111113	(0, - 1)	1/9
(+0.9996, + 0.0047)	0.111111	(0, + 1)	1/9
(+0.9974, + 0.9973)	0.111091	(-1, 0)	1/9
(+0.0236, + 0.9992)	0.059695	(+1, 0)	1/9
(-0.0305, + 0.9994)	0.051457	(0, 0)	1/9
(-0.9998, + 0.9993)	0.111091		

Tabela 5.1.

Tabela po lewej stronie – przybliżenie D-optimalnego planu eksperymentu, otrzymane za pomocą algorytmu selektywnych poszukiwań losowych, w przypadku estymacji kwadratowej funkcji regresji o dwóch zmiennych (szczegółowy opis – w przykładzie symulacji 5.1). Tabela po prawej stronie – dla porównania – plan D-optimalny

Punkty	p_i
(+0.99, - 0.99, - 0.97)	0.125
(-0.98, + 0.98, + 0.98)	0.125
(-0.97, + 0.98, - 0.99)	0.125
(+0.99, + 0.99, + 0.97)	0.125
(-0.98, - 0.99, - 0.99)	0.125
(+0.98, - 0.98, + 0.98)	0.125
(+0.97, + 0.98, - 0.98)	0.125
(-0.99, - 0.98, + 0.98)	0.125

Tabela 5.2. Jeden z planów przybliżonych dla estymacji parametrów regresji (5.2.26), znaleziony w przykładzie symulacji 5.2



Rys. 5.4. Przykładowe przebiegi algorytmu selektywnych poszukiwań losowych – opis w przykładzie symulacji 5.2

Symulacja 5.2. Rozważmy model liniowy z pełnym zestawem interakcji rzędu pierwszego, który można zapisać następująco:

$$\bar{y}(x) = a^T \left\{ \begin{bmatrix} 1 \\ x^{(1)} \end{bmatrix} \otimes \begin{bmatrix} 1 \\ x^{(2)} \end{bmatrix} \otimes \begin{bmatrix} 1 \\ x^{(3)} \end{bmatrix} \right\}. \quad (5.2.26)$$

Użyty w tym wzorze symbol \otimes oznacza iloczyn Kroneckera wektorów. Jego znaczenie dla zapisu modeli omawiamy szczegółowo w następnym rozdziale, a definicję i własności przedstawiamy w dodatku.

Jak wiemy, plan D -optymalny na kostce $X = [-1, 1]^3$ skupiony jest w wierzchołkach tej kostki, a wagi im przypisane wynoszą $1/8$. Na rysunku 5.4 pokazano początkową fazę (ok. 150 iteracji) przykładowych przebiegów algorytmu selektywnych poszukiwań losowych, a tabela 5.2 zawiera jeden ze znalezionych planów przybliżonych. W warunku stopu algorytmu zastosowano wartość $N = 50\,000$. Pozostałe parametry algorytmu były takie jak w poprzednich przebiegach symulacyjnych. Zauważmy, że zarówno punkty, jak i wagi planu z tabeli 5.2 są bardzo bliskie odpowiednim wartościom dla planu optymalnego.

Z drugiej jednak strony, wyznacznik macierzy informacyjnej uzyskanego planu ma wartość około 0.65, podczas gdy dla planu optymalnego wartość ta wynosi 1. Zauważmy, że w poprzednim przykładzie bez trudu uzyskiwaliśmy ponad 95 procentową efektywność planu znalezionej numerycznie, podczas gdy tutaj uzyskujemy

65% efektywności. Jest to skutek wzrostu wymiaru przestrzeni poszukiwań. Uzyskanie lepszej dokładności wymagałoby znacznego wydłużenia czasu poszukiwań. Zalety i wady algorytmu selektywnych poszukiwań losowych podsumować można następująco:

- Algorytm znajduje przybliżenie optymalnego planu po stosunkowo małej liczbie iteracji (rzędu kilkuset), nawet w wielowymiarowych zadaniach. Jednakże osiągnięcie dużej precyzji przybliżenia znacznie wydłuża czas obliczeń.
- W algorytmie unika się bezpośredniego i wielokrotnego poszukiwania globalnego maksimum funkcji wielu zmiennych. Maksymalizacja prowadzona jest pośrednio poprzez losowe znajdowanie „dostatecznie dobrych” punktów w krajobrazie tworzonym przez funkcję wariancji, która zmienia się wraz z przebiegiem obliczeń. W rezultacie, algorytm ten jest łatwy do zaprogramowania.
- Algorytm nie produkuje zbędnych punktów nośnika planu z bardzo małymi wagami, gdyż w każdej iteracji włączana jest procedura optymalizacji wag, a po jej zakończeniu punkty takie uzyskują wagi tak znacząco mniejsze od innych, że ich zidentyfikowanie i usunięcie staje się oczywiste.
- Dalszego zmniejszenia nakładów obliczeniowych można oczekiwać w wyniku właściwego doboru liczby iteracji procedury optymalizacji wag.

5.3. Optymalna alokacja pomiarów

W podrozdziale 5.1 przedstawiliśmy wersję algorytmu Wynna–Fedorowa, przeznaczoną do optymalizacji wag planu o zadanym nośniku. Przedstawimy teraz, pochodzącą od Fellmana i Silveya, Titteringtona i Torsneya (por. uwagi bibliograficzne na końcu podrozdziału) metodę, która także ogranicza się do optymalizacji jedynie wag planu, podczas gdy nośnik planu jest (na czas optymalizacji wag lub na stałe) ustalony i ma postać: x_1, x_2, \dots, x_d , $d \geq 1$. Na takim zbiorze każdy plan $\xi \in \Xi(X)$ wyznaczony jest jednoznacznie przez wagi $p_i \geq 0$, $i = 1, 2, \dots, d$, $\sum_{i=1}^d p_i = 1$.

Metodę optymalizacji wag traktować będziemy jako uzupełnienie algorytmów szukania planów optymalnych metodą Wynna–Fedorowa lub selektywnych poszukiwań losowych. Opisany niżej algorytm można włączyć do tych metod jako jeden z kroków, realizowany po wygenerowaniu nowego punktu wprowadzanego do planu. Inne sposoby korzystania z metody optymalizacji wag planu omówione są na końcu tego podrozdziału.

Opis algorytmu optymalizacji wag

Zdefiniujmy odwzorowanie T , które przeprowadza plan ξ w plan $\xi' = T(\xi)$, w ten sposób, że modyfikujemy wagi p_i , $i = 1, 2, \dots, d$ planu ξ następująco:

$$p'_i = p_i \frac{\phi(\xi, x_i)}{r}, \quad i = 1, 2, \dots, d, \quad (5.3.27)$$

gdzie p'_i to wagi nowego planu ξ' . Nośniki obu planów ξ i ξ' pozostają bez zmian. Zauważmy, że twierdzenie Kiefera–Wolfowitza, w szczególności warunek (4.3.14), pozostają w mocy w przypadku skończonego zbioru X i – zgodnie z Twierdzeniem 4.1 – dla planu optymalnego

$$\frac{\phi(\xi^*, x_i)}{r} = 1, \quad i = 1, 2, \dots, d.$$

W takim przypadku modyfikacja wag w (5.3.27) nie następuje. Wektor wag planu D-optymalnego ξ^* jest punktem stałym odwzorowania T , tzn. $T(\xi^*) = \xi^*$.

Iteracyjne zastosowanie odwzorowania T prowadzi do następującego algorytmu optymalizacji proporcji rozdziału pomiarów między punkty nośnika planu.

- 1) Wybieramy plan ξ^0 o nieosobliwej macierzy informacyjnej i taki, że wagi wszystkich punktów ze skończonego zbioru X są dodatnie $p_i^{(0)} > 0, i = 1, 2, \dots, d$.
- 2) Obliczamy plan $\xi^{(n+1)}$ (modyfikując jedynie wagi planu $\xi^{(n)}$) zgodnie ze wzorem

$$p_i^{(n+1)} = p_i^{(n)} \frac{\phi(\xi^{(n)}, x_i)}{r}, \quad i = 1, 2, \dots, d. \quad (5.3.28)$$

Korzystając z własności 6 można wykazać, że po wykonaniu (5.3.28) spełniony jest warunek unormowania $\sum_{i=1}^d p_i^{(n+1)} = 1$ jeśli tylko $\sum_{i=1}^d p_i^{(n)} = 1$.

- 3) Sprawdzamy warunek:

$$\frac{\phi(\xi^{(n)}, x_i)}{r} < 1 + \epsilon, \quad i = 1, 2, \dots, d, \quad (5.3.29)$$

gdzie $\epsilon > 0$ oznacza wybraną wcześniej dokładność planowania. Jeśli nie jest on spełniony, to zwiększamy n o jeden i ponownie obliczamy (5.3.28). W przeciwnym razie, uznajemy $\xi^{(n)}$ za przybliżenie planu D-optymalnego.

Można wykazać (por. [105]), że ciąg $\det(M(\xi^{(n)}))$, $n = 0, 1, \dots$ jest niemalejący. Ponadto, jeśli warunek (5.3.29) nie jest sprawdzany, to generowany jest nieskończony ciąg taki, dla którego zachodzi

$$\lim_{n \rightarrow \infty} \det(M(\xi^{(n)})) = \det(M(\xi^*)),$$

gdzie ξ^* jest planem D-optymalnym na skończonym zbiorze X .

Jeśli plan ξ nie jest D-optymalny, to w pewnych punktach jego nośnika musi zachodzić $\phi(\xi, x_i) > r$. Po zastosowaniu (5.3.28), wagi tych punktów ulegną zwiększeniu, co prowadzi do zmniejszenia wariancji estymatora wyjścia modelu w tychże punktach. Zwiększenie wag w punktach, gdzie $\phi(\xi, x_i) > r$ musi odbyć

się kosztem zmniejszenia wag tych punktów, w których $\phi(\xi, x_i) < r$, gdyż w nich wariancja predykcji była względnie mała.

Wymaganie, by wybrany plan początkowy przypisywał dodatnie wagi wszystkim punktom zbioru X jest istotne, gdyż – zgodnie z powyższym algorytmem – punktowi, którego waga równa jest zero nigdy nie zostanie przypisana waga dodatnia.

Porównanie algorytmów optymalizacji wag

O ile autorowi wiadomo, w literaturze nie ma teoretycznych rezultatów na temat porównania szybkości zbieżności algorytmu Wynna-Fedorova w wersji służącej tylko do optymalizacji wag (por. podrozdz. 5.1) i opisanego wyżej algorytmu multiplikatywnego.

W podrozdziale tym przytaczamy wyniki wycinkowego, symulacyjnego porównania tych algorytmów przeprowadzone przez autora. Badania przeprowadzono na zbiorze 21 punktów równomiernie rozmieszczonych w przedziale $[-1, 1]$ w taki sposób, że do zbioru tego wchodziły punkty ± 1 oraz 0. Na zbiorze tym szukano wag planów dla estymacji parametrów wielomianowych funkcji regresji od stopnia pierwszego do siódmego. Zauważmy, że zbiór dyskretnych punktów zawierał wszystkie punkty nośnika planów optymalnych na $[-1, 1]$ dla regresji stopnia pierwszego i drugiego, natomiast dla wielomianów wyższych stopni miał jedynie część punktów wspólnych z nośnikami planów optymalnych na $[-1, 1]$.

Ze względu na różnice w konstrukcji i warunkach stopu obu algorytmów liczba iteracji nie jest dobrym miernikiem do porównywania ich jakości. Z tego względu zdecydowano się na porównania bezpośrednie, polegające na pomiarze czasów obliczeń, które oba algorytmy potrzebowały na znalezienie planu o tej samej (i bliskiej maksymalnej) wartości wyznacznika macierzy informacyjnej. We wszystkich przypadkach jako wagi startowe wybierano rozkład równomierny na opisanym wyżej zbiorze punktów. Obliczeń dokonywano na tym samym komputerze, a procedury napisano w języku środowiska Mathematica w taki sposób, by mogły służyć jako elementy składowe programów do szukania planów optymalnych, na przykład metodą selektywnych poszukiwań losowych. Wyniki porównania obu algorytmów przedstawiono w tabeli 5.3. Algorytm Wynna-Fedorova oznaczono w niej skrótem W-F, a algorytm multiplikatywny oznaczono jako mult. Porównań dokonywano dla modeli wielomianowych stopnia od jeden do siedem (kolumna oznaczona jako Deg). Wyniki zamieszczone w ostatniej kolumnie tabeli 5.3 wskazują na znaczną przewagę algorytmu multiplikatywnego poprawiania wag.

Klasteryzacja i usuwanie punktów

Algorytm Wynna-Fedorova, algorytm selektywnych poszukiwań losowych i algorytm optymalizacji wag stopniowo zmniejszają wagi tych punktów, które nie należą do planu optymalnego. Ponadto, dwa pierwsze z nich dodają do nośnika

Deg	Algorytm				T_1/T_2
	W-F		Mult.		
	Det	Czas T_1	Det	Czas T_2	
1	0.997	0.109	0.995	0.062	1.8
2	0.146	0.234	0.145	0.062	3.8
3	0.0048	0.359	0.0048	0.063	5.7
4	$3.8 \cdot 10^{-5}$	0.406	$3.8 \cdot 10^{-5}$	0.062	6.5
5	$7.7 \cdot 10^{-8}$	1.734	$7.7 \cdot 10^{-8}$	0.078	22
6	$3.7 \cdot 10^{-11}$	2.9	$3.7 \cdot 10^{-11}$	0.094	30
7	$2.8 \cdot 10^{-15}$	2.735	$2.8 \cdot 10^{-15}$	0.094	30

Tabela 5.3. Wyniki porównania czasów obliczeń algorytmu poprawy wag w wersji Wynna–Fedorova i algorytmu multiplikatywnego (objaśnienia w tekście)

planu nowe punkty – znajdowane w wyniku numerycznej maksymalizacji funkcji ϕ na X . Oba te działania prowadzą do tego, że na etapach pośrednich nośnik planu zawiera wiele punktów zbędnych. Te zbędne podzielić można na dwa rodzaje.

1. Punkty, których wagi są bardzo małe. Określenie to jest nieprecyzyjne i wymaga podania progu $\epsilon > 0$, poniżej którego punkt uznajemy za zbędny. Punktem odniesienia dla wyboru ϵ może być liczba $1/N$, gdzie N jest całkowitą liczbą pomiarów, które mają być wykonane przy realizacji zaplanowanego eksperymentu. Jeśli jako ϵ wybierzemy liczbę mniejszą niż $1/N$, to w punktach mających wagi mniejsze niż tak wybrane ϵ pomiary nie będą wykonywane.
2. Punkty planu, których nośniki są sobie równe dokładnie lub w przybliżeniu. Jest to wynik skończonej dokładności znajdowania maksimum ϕ . Często zdarza się bowiem, że nowo włączany do planu punkt powinien mieć nośnik pokrywający się z punktem już w planie istniejącym.

Kierując się tymi spostrzeżeniami, warto dolożyć do wszystkich omawianych wyżej algorytmów proste procedury usuwania punktów o zbyt małych wagach i grupowania w jeden punktów planu leżących blisko siebie. Stosując takie grupowanie należy dodać wagi wszystkich punktów grupowanych w jeden, a usuwanie punktów o zbyt małych wagach wykonywać dopiero po grupowaniu. Te proste zabiegi implementacyjne w praktyce bardzo przyspieszają zbliżanie się do planu optymalnego. Należy je jednak stosować z pewną ostrożnością, gdyż nie trudno wskazać „patologiczne” przykłady sytuacji, w których odrzucanie punktów o zbyt małych wagach lub grupowanie prowadzić będzie do niebezpieczeństwa algorytmu poszukiwania planu optymalnego. Przez wybór zbyt dużych progów odrzucania i gru-

powania punktów łatwo wywołać efekt wprowadzania i odrzucania w kolejnych iteracjach tych samych punktów. W pracy [109] udowodniono nierówność, która jest podstawą bezpiecznego testu służącego do usuwania punktów z planu, ale tylko takich, które na pewno nie należą do nośnika planu optymalnego.

Uwagi bibliograficzne i inne

Prezentowana tu procedura optymalizacji wag, tak by otrzymać wagi planu D-optymalnego ma dość długą historię (por. [37], [174], [190], [105], [191] i zawarte tam bibliografie). Znane są także jej uogólnienia na inne kryteria [38].

Poświęcamy jej sporo uwagi, gdyż – w odróżnieniu od większości procedur optymalizacji w ogóle, a planów w szczególności – procedura ta dokonuje poprawy multiplikatywnie. Daje to w efekcie znaczną szybkość zbieżności. Zastosowanie podobnego podejścia do przetwarzania obrazów podjęto w pracy [190].

A oto możliwe sposoby wykorzystania procedury optymalizacji wag w problemach planowania.

1. Jeśli zbiór dopuszczalnych punktów nośnika planu jest skończony (dyskretny zbiór X), to omawiana procedura znajdować będzie plan bliski optymalnemu z zadaną z góry dokładnością.
2. Jeśli zbiór X jest pewnym obszarem w R^s , to można w nim wybrać pewien skończony podzbiór punktów X_d , który dostatecznie równomiernie pokrywa X . Zbiór ten można następnie wykorzystać tak jak w punkcie 1. W tym przypadku procedura gwarantuje dobór wag bliski optymalnemu, ale na zbiorze X_d zamiast na X . Ocenę jakości tak uzyskanego planu przeprowadzić należy dodatkowo, korzystając z warunków optymalności.
3. Jeśli zbiór X jest pewnym obszarem w R^s , to można próbować wyselekcjonować pewien „mały” zbiór punktów $X_d \subset X$, który zawiera punkty nośnika planu optymalnego. Następnie, stosując omawianą w tym podrozdziale procedurę, dobieramy wagi planu. Selekcja punktów należących do nośnika planu optymalnego może odbywać się metodą Wynn–Fedorova, metodą selektywnych poszukiwań losowych lub metodą zaproponowaną w pracy [176], w której do znajdowania kolejnych punktów wprowadzanych do planu zaproponowano wykorzystanie krzywych wypełniających takich jak krzywa Hilberta, Peano czy Sierpińskiego (por. [177] oraz [178], [179], gdzie opisano krzywe wypełniające i algorytmy ich generowania).

6. Plany dla modeli o zmiennych zblokowanych

Idea rozmieszczania obserwacji w blokach należy do klasycznych pojęć teorii planowania eksperymentu (por. [96], [108], [201]). Potrzeba rozpatrywania planów blokowych pojawiła się najwcześniej w eksperymentach rolniczych, gdyż tam potomstwo pochodzące z jednego miotu czy też poletka wydzielone z jednego pola były – w naturalny sposób – zgrupowane, gdyż posiadały wspólne cechy. W naukach technicznych grupowanie eksperymentów może również następować naturalnie, w trzech co najmniej sytuacjach:

- gdy ten sam zestaw próbek pewnego materiału poddawany jest badaniom o różnej naturze, na przykład wytrzymałości, oporności elektrycznej, przewodności cieplnej,
- gdy ustalony zestaw badań wykonywany jest na partiach próbek z różnych materiałów,
- gdy występują trudności ze zmienianiem wartości pewnej podgrupy zmiennych, na przykład, można je zmieniać w tylko w warunkach laboratoryjnych, ale nie wówczas, gdy badany obiekt pracuje w warunkach normalnej eksploatacji.

Biorąc pod uwagę wymienione sytuacje, w rozdziale tym łączenie zmiennych w grupy traktować będziemy czysto formalnie – kierując się wygodą obliczania optymalnych planów eksperymentu. Jak się okaże, dopatrzenie się w modelu odpowiedniej struktury związków między zmiennymi wejściowymi (morfologii) znacznie ułatwia problem planowania.

6.1. Iloczyn planów zależnych od zblokowanych zmiennych

Założmy, że zmienne wejściowe zostały ponumerowane w ten sposób, że $\mathbf{x}^{(1)} = [x^{(1)}, \dots, x^{(s_1)}]^T$ oraz $\mathbf{x}^{(2)} = [x^{(s_1+1)}, \dots, x^{(s)}]^T$. $\mathbf{x}^{(1)}$ i $\mathbf{x}^{(2)}$ zestawione razem tworzą wektor x . Wektory $\mathbf{x}^{(1)}$ i $\mathbf{x}^{(2)}$ nazywać będziemy blokami zmiennych.

Przypuśćmy, że dla każdego zestawu zblokowanych zmiennych dany jest plan ciągły, skupiony w skończonej liczbie punktów

$$\xi^{(1)} = \begin{bmatrix} \mathbf{x}_1^{(1)} & \mathbf{x}_2^{(1)} & \dots & \mathbf{x}_{K_1}^{(1)} \\ p_1^{(1)} & p_2^{(1)} & \dots & p_{K_1}^{(1)} \end{bmatrix}, \quad \xi^{(2)} = \begin{bmatrix} \mathbf{x}_1^{(2)} & \mathbf{x}_2^{(2)} & \dots & \mathbf{x}_{K_2}^{(2)} \\ p_1^{(2)} & p_2^{(2)} & \dots & p_{K_2}^{(2)} \end{bmatrix}.$$

Punkty nośników tych planów są elementami zbiorów $\mathbf{X}_1 \subset \mathbf{R}^{s_1}$ i $\mathbf{X}_2 \subset \mathbf{R}^{s_2}$, odpowiednio, gdzie $s_2 \stackrel{\text{def}}{=} s - s_1$. Ponadto spełnione są warunki:

$$p_i^{(1)} \geq 0, \quad \sum_{i=1}^{K_1} p_i^{(1)} = 1, \quad p_j^{(2)} \geq 0, \quad \sum_{j=1}^{K_2} p_j^{(2)} = 1. \quad (6.1.1)$$

Uwaga 6.1. W poprzednich wzorach i dalej pogrubioną czcionką oznaczamy bloki zmiennych i odpowiadające im podzbiory i przestrzenie. To odstępstwo od notacji stosowanej wcześniej ma pomóc w odróżnianiu modelu i jego części oraz odpowiadających im planów.

Mając dwa takie plany dla bloków zmiennych, możemy stworzyć plan dla wszystkich s zmiennych wejściowych w następujący sposób.

Definicja 6.1. *Produktem planów $\xi^{(1)}$ i $\xi^{(2)}$ nazywamy plan skupiony w $K_1 \cdot K_2$ punktach, skonstruowany następująco:*

- punkty $(\mathbf{x}_i^{(1)}, \mathbf{x}_j^{(2)})$ stanowią nośnik planu,
- odpowiadające im wagi p_{ij} mają postać $p_{ij} = p_i^{(1)} p_j^{(2)}$, $i = 1, 2, \dots, K_1$, $j = 1, 2, \dots, K_2$.

Definicja 6.2. *Ciągłym planem produktowym skupionym w skończonej liczbie punktów nazywamy plan powstały jako produkt planów $\xi^{(1)}$ i $\xi^{(2)}$. Oznaczamy go przez $\xi^{(1)} \otimes \xi^{(2)}$ i nazywamy w skrócie planem produktowym.*

Plany $\xi^{(1)}$ i $\xi^{(2)}$ nazywać będziemy planami dla bloków zmiennych lub planami częściowymi planu $\xi^{(1)} \otimes \xi^{(2)}$. Zauważmy, że plan częściowy $\xi^{(1)}$ interpretować można jako fragmenty planu $\xi^{(1)} \otimes \xi^{(2)}$ powstałe w wyniku ustalenia pewnej wartości zmiennych \mathbf{x}_1 , a więc ustalenia wpływu jednego z bloków.

Nośnik planu $\xi^{(1)} \otimes \xi^{(2)}$ zawarty jest w iloczynie kartezjańskim $\mathbf{X}_1 \times \mathbf{X}_2$. Łatwo sprawdzić, że wagi $p_{ij} \geq 0$ spełniają wymaganie $\sum_{i=1}^{K_1} \sum_{j=1}^{K_2} p_{ij} = 1$. Z obu tych własności wynika, że jeśli

$$\xi^{(1)} \in \Xi(\mathbf{X}_1) \quad \text{i} \quad \xi^{(2)} \in \Xi(\mathbf{X}_2), \quad \text{to} \quad \xi^{(1)} \otimes \xi^{(2)} \in \Xi(\mathbf{X}_1 \times \mathbf{X}_2),$$

co oznacza, że produkt planów ciągłych skupionych w skończonych liczbach punktów jest również planem ciągłym skupionym w skończonej liczbie punktów. Jeśli jednak rozważymy wszystkie możliwe plany produktowe

$$\left\{ \xi^{(1)} \otimes \xi^{(2)} : \xi^{(1)} \in \Xi(\mathbf{X}_1), \xi^{(2)} \in \Xi(\mathbf{X}_2) \right\}, \quad (6.1.2)$$

to okaże się, że klasa planów produktowych na $\mathbf{X}_1 \times \mathbf{X}_2$ nie wypełnia całego zbioru $\Xi(\mathbf{X}_1 \times \mathbf{X}_2)$ wszystkich planów ciągłych na $\mathbf{X}_1 \times \mathbf{X}_2$. Dalej pokażemy jednak, że klasa planów (6.1.2) jest dostatecznie bogata, by zawierać w sobie plany optymalne dla szerokiej klasy funkcji regresji wielu zmiennych.

Z formalnego punktu widzenia operacja \otimes nie jest przemienne, ponieważ $\xi^{(1)} \otimes \xi^{(2)} \neq \xi^{(2)} \otimes \xi^{(1)}$. Jednakże oba te plany produktowe różnią się między sobą kolejnością numeracji zmiennych wejściowych, co pozwala je w praktyce utożsamiać.

Operacja \otimes jest łączna w tym sensie, że jeśli mamy trzy plany $\xi^{(1)} \in \mathbf{X}_1$, $\xi^{(2)} \in \mathbf{X}_2$, $\xi^{(3)} \in \mathbf{X}_3$, to

$$\left(\xi^{(1)} \otimes \xi^{(2)}\right) \otimes \xi^{(3)} = \xi^{(1)} \otimes \left(\xi^{(2)} \otimes \xi^{(3)}\right). \quad (6.1.3)$$

Własność ta pozwala poprawnie zdefiniować produkt dowolnej skończonej liczby planów dla poszczególnych bloków

$$\prod_{i=1}^N \xi^{(i)} = \left[\prod_{i=1}^{N-1} \xi^{(i)} \right] \otimes \xi^{(N)}, \quad N = 2, 3, \dots \quad (6.1.4)$$

Jeśli jako $\xi^{(1)}$ i $\xi^{(2)}$ przyjmiemy plany skupione w punktach ± 1 z wagami $1/2$, to plan $\xi^{(1)} \otimes \xi^{(2)}$ skupiony będzie we wszystkich wierzchołkach kwadratu $[-1, 1] \times [-1, 1]$ z wagami $1/4$. Jeśli również plan $\xi^{(3)}$ ma taką samą postać jak $\xi^{(1)}$, to plan $\xi^{(1)} \otimes \xi^{(2)} \otimes \xi^{(3)}$ ma nośnik we wszystkich wierzchołkach kostki $[-1, 1]^3$ z wagami $1/8$. Nie należy jednak sądzić, że za pomocą operacji \otimes otrzymać można jedynie plany o geometrycznym kształcie siatki wielowymiarowej. Jeśli np. plan $\xi^{(1)}$ skupiony jest na okręgu, a plan $\xi^{(2)}$ ma nośnik na odcinku, to punkty planu $\xi^{(1)} \otimes \xi^{(2)}$ położone będą na poboczniczy walca.

6.2. Planowanie dla modeli addytywnych względem zblokowanych zmiennych

Wygodnie jest zapisywać modele addytywne korzystając z następującej definicji sumy prostej wektorów (por. [82] i [147], gdzie znaleźć można także definicję iloczynu Kroneckera wektorów, która przydatna będzie w dalszych rozważaniach).

Definicja 6.3. Sumą prostą $a \oplus b$ kolumnowych wektorów $a \in \mathbf{R}^{d_1}$ i $b \in \mathbf{R}^{d_2}$, oznaczaną jako $a \oplus b$, nazywa się wektor $[a^T, b^T]^T$ o $d_1 + d_2$ składowych.

Odnotujemy, że $a \oplus b \neq b \oplus a$. Operacja ta jest natomiast łączna.

Podobnie jak we poprzednim podrozdziale, założmy, że zmienne wejściowe podzielone na bloki zostały ponumerowane w ten sposób, że $\mathbf{x}^{(1)} = [x^{(1)}, \dots, x^{(s_0)}]^T$ oraz $\mathbf{x}^{(2)} = [x^{(s_0+1)}, \dots, x^{(s)}]^T$ i tworzą one wektor $\mathbf{x} \stackrel{\text{def}}{=} \mathbf{x}^{(1)} \oplus \mathbf{x}^{(2)}$.

Definicja 6.4. Modelem addytywnym względem bloków zmiennych $\mathbf{x}^{(1)}$, $\mathbf{x}^{(2)}$ nazywamy

$$\bar{y}(x) = \alpha^T \mathbf{v}_1(\mathbf{x}^{(1)}) + \beta^T \mathbf{v}_2(\mathbf{x}^{(2)}) = \alpha^T \left[\mathbf{v}_1(\mathbf{x}^{(1)}) \oplus \mathbf{v}_2(\mathbf{x}^{(2)}) \right], \quad (6.2.5)$$

gdzie $a \stackrel{\text{def}}{=} \alpha \oplus \beta$ tworzą wektor a estymowanych parametrów, natomiast $v(x) \stackrel{\text{def}}{=} \mathbf{v}_1(\mathbf{x}^{(1)}) \oplus \mathbf{v}_2(\mathbf{x}^{(2)})$ tworzą bazę, na której rozpięta jest funkcja regresji. Zakładamy przy tym, że $\dim(\mathbf{v}_1(\mathbf{x}^{(1)})) = \dim(\alpha)$ oraz $\dim(\mathbf{v}_2(\mathbf{x}^{(2)})) = \dim(\beta)$.

Funkcje $\alpha^T \mathbf{v}_1(\mathbf{x}^{(1)})$, $\beta^T \mathbf{v}_2(\mathbf{x}^{(2)})$ nazywać będziemy modelami częściowymi modelu (6.2.5).

Jeśli w $\mathbf{x}^{(1)}$ i/lub w $\mathbf{x}^{(2)}$ wyróżnić można podgrupy (podbloki) zmiennych w sposób podobny do opisanego w Definicji 6.4, to celowa jest dalsza dekompozycja modelu. W skrajnym przypadku otrzymamy model (w pełni) addytywny o postaci:

$$\bar{y}(x) = a^{(0)} + \sum_{k=1}^s a^{(k)} v^{(k)}(x^{(k)}). \quad (6.2.6)$$

Notacja we wzorze (6.2.6) nie jest w pełni zgodna z tą, która była użyta w Definicji 6.4, gdyż, w celu zachowania zgodności z notacją tradycyjną, funkcje $v^{(k)}$ i ich argumenty $x^{(k)}$ nie zostały złożone czcionką pogrubioną.

Przykład modelu, w którym nie jest możliwa pełna dekompozycja, to

$$\bar{y}(x) = \alpha^{(0)} + \alpha^{(1)}x^{(1)} + \alpha^{(2)}x^{(2)} + \alpha^{(3)}x^{(1)}x^{(2)} + \beta^{(1)}x^{(3)} + \beta^{(2)}x^{(4)}. \quad (6.2.7)$$

Model ten zapisać można w postaci (6.2.5), wprowadzając zmienne

$$\mathbf{x}^{(1)} = [x^{(1)}, x^{(2)}]^T, \quad \mathbf{x}^{(2)} = [x^{(3)}, x^{(4)}]^T$$

oraz wektory funkcji rozpinających

$$\mathbf{v}_1(\mathbf{x}^{(1)}) = [1, x^{(1)}, x^{(2)}, x^{(1)}x^{(2)}]^T, \quad \mathbf{v}_2(\mathbf{x}^{(2)}) = [x^{(3)}, x^{(4)}]^T.$$

Łatwo zauważyć, że ta część modelu, która rozpinana jest przez $\mathbf{v}_2(\mathbf{x}^{(2)})$ może być zdekomponowana do poziomu poszczególnych zmiennych. Trzy modele częściowe dla pełnego modelu (6.2.7) rozpięte są zatem przez $[1, x^{(1)}, x^{(2)}, x^{(1)}x^{(2)}]^T$, $x^{(3)}$ oraz $x^{(4)}$.

Korzystając z tego przykładu, warto zwrócić uwagę na to, że funkcja tożsamościowo równa 1 wystąpić może w zestawie funkcji rozpinających co najwyżej jeden model częściowy (w przeciwnym razie funkcje te nie będą liniowo niezależne).

Rozpatrując modele częściowe (6.2.6) zakładamy, że obszar planowania X da się przedstawić w postaci iloczynu kartezjańskiego

$$X = \mathbf{X}_1 \times \mathbf{X}_2, \quad \mathbf{X}_1 \subset \mathbf{R}^{\dim(\mathbf{x}^{(1)})}, \quad \mathbf{X}_2 \subset \mathbf{R}^{\dim(\mathbf{x}^{(2)})} \quad (6.2.8)$$

dwóch obszarów planowania, odpowiadających zmiennym w poszczególnych blokach. Jeśli, tak jak w podanym przykładzie, możliwa jest dalsza dekompozycja modelu, to zakładamy, że również \mathbf{X}_1 i/lub \mathbf{X}_2 dadzą się przedstawić jako iloczyny kartezjańskie odpowiednich obszarów planowania.

Dalsze obliczenia znacznie łatwiej się zapisuje, jeśli wprowadzi się pojęcie sumy prostej macierzy (por. [82]).

Definicja 6.5. Niech M_1 i M_2 będą odpowiednio macierzami kwadratowymi $r_1 \times r_1$ i $r_2 \times r_2$. Sumą prostą macierzy, oznaczaną przez \oplus , nazywa się macierz blokową:

$$M_1 \oplus M_2 = \begin{bmatrix} M_1 & \mathbf{0} \\ \mathbf{0} & M_2 \end{bmatrix}, \quad (6.2.9)$$

gdzie $\mathbf{0}$ oznaczają macierze złożone z elementów zerowych. Macierz $M_1 \oplus M_2$ jest macierzą kwadratową o $r_1 + r_2$ wierszach i kolumnach.

Uwaga 6.2. Zauważmy, że dziedzinę określoności operacji \oplus ograniczyliśmy do macierzy kwadratowych. Stosowanie jej do macierzy prostokątnych doprowadziłoby do niezgodności rezultatu sumy prostej wektorów (por. z Definicją 6.3) i sumy tychże wektorów traktowanych jako macierze $r_1 \times 1$ i $r_2 \times 1$.

Obliczmy macierz informacyjną dla pewnego planu produktowego $\xi^{(1)} \otimes \xi^{(2)}$, $\xi^{(1)} \in \Xi(\mathbf{X}_1)$, $\xi^{(2)} \in \Xi(\mathbf{X}_2)$, który stosujemy do estymacji parametrów modelu podanego w Definicji 6.4. Korzystając z definicji sumy prostej wektorów, otrzymamy

$$M(\xi^{(1)} \otimes \xi^{(2)}) = \begin{bmatrix} M_1(\xi^{(1)}) & M_{12}(\xi^{(1)}, \xi^{(2)}) \\ M_{21}(\xi^{(1)}, \xi^{(2)}) & M_2(\xi^{(2)}) \end{bmatrix}, \quad (6.2.10)$$

gdzie

$$M_1(\xi^{(1)}) \stackrel{\text{def}}{=} \int_{\mathbf{X}_1} \mathbf{v}_1(\mathbf{x}^{(1)}) \mathbf{v}_1(\mathbf{x}^{(1)})^T \xi^{(1)}(d\mathbf{x}^{(1)}), \quad (6.2.11)$$

$$M_2(\xi^{(2)}) \stackrel{\text{def}}{=} \int_{\mathbf{X}_2} \mathbf{v}_2(\mathbf{x}^{(2)}) \mathbf{v}_2(\mathbf{x}^{(2)})^T \xi^{(2)}(d\mathbf{x}^{(2)}), \quad (6.2.12)$$

$$M_{12}(\xi^{(1)}, \xi^{(2)}) \stackrel{\text{def}}{=} \int_{\mathbf{X}_1} \mathbf{v}_1(\mathbf{x}^{(1)}) \xi^{(1)}(d\mathbf{x}^{(1)}) \int_{\mathbf{X}_2} \mathbf{v}_2^T(\mathbf{x}^{(2)}) \xi^{(2)}(d\mathbf{x}^{(2)}) \quad (6.2.13)$$

oraz $M_{21}(\xi^{(1)}, \xi^{(2)}) \stackrel{\text{def}}{=} M_{12}^T(\xi^{(1)}, \xi^{(2)})$.

Definicja 6.6. Plan $\xi^{(i)} \in \Xi(\mathbf{X}_i)$ nazywamy planem symetrycznym względem $\mathbf{v}_i(\mathbf{x}^{(i)})$, wówczas gdy

$$\int_{\mathbf{X}_i} \mathbf{v}_i(\mathbf{x}^{(i)}) \xi^{(i)}(d\mathbf{x}^{(i)}) = \mathbf{0}. \quad (6.2.14)$$

Odnosząc to określenie do modelu opisanego w Definicji 6.4 możemy mieć do czynienia z planami dla bloków zmiennych, które są symetryczne względem $\mathbf{v}_1(\mathbf{x}^{(1)})$ i/lub $\mathbf{v}_2(\mathbf{x}^{(2)})$. Jednakże, jeśli $\mathbf{v}_1(\mathbf{x}^{(1)})$ zawiera funkcję stałą (model zawiera wyraz wolny), to nie istnieje plan symetryczny względem tego zestawu funkcji. Wówczas istnieje możliwość znalezienia planu symetrycznego względem drugiego z tych zestawów funkcji.

Lemat 6.1. *Rozważmy model addytywny opisany w Definicji 6.4 i plan produktowy $\xi^{(1)} \otimes \xi^{(2)}$, $\xi^{(1)} \in \Xi(\mathbf{X}_1)$, $\xi^{(2)} \in \Xi(\mathbf{X}_2)$. Załóżmy, że spełniony jest co najmniej jeden z warunków:*

- plan częściowy $\xi^{(1)}$ jest symetryczny względem $\mathbf{v}_1(\mathbf{x}^{(1)})$,
- plan częściowy $\xi^{(2)}$ jest symetryczny względem $\mathbf{v}_2(\mathbf{x}^{(2)})$.

Wówczas

$$M(\xi^{(1)} \otimes \xi^{(2)}) = M_1(\xi^{(1)}) \oplus M_2(\xi^{(2)}). \quad (6.2.15)$$

Równość (6.2.15) wynika z (6.2.10) i z faktu, że gdy plan(y) częściowy(e) ma własność symetrii, to $M_{21}(\xi^{(1)}, \xi^{(2)}) = 0$ i $M_{12}(\xi^{(1)}, \xi^{(2)}) = 0$.

Niech $\Phi(M(\xi); r)$ będzie kryterium planowania eksperymentu dla funkcji regresji o r parametrach. W oznaczeniu tym jawnie figuruje liczba parametrów r , gdyż rozważane będą zadania modele częściowe funkcji regresji, różniące się zestawem funkcji rozpinających, obszarem planowania i liczbą parametrów.

Zakładamy, że Φ jest funkcją wklęsłą, różniczkowalną i taką, że w rozważanym zadaniu planowania istnieje rozwiązanie o nieosobliwej macierzy informacyjnej. Pojęcie gradientu kryterium planowania wprowadziliśmy już w poprzednim rozdziale. Przytaczamy tę definicję ponownie, by zaakcentować zależność gradientu od r .

Definicja 6.7. *Niech M będzie $r \times r$ macierzą dodatnio określoną. Gradientem kryterium planowania $\Phi(M; r)$, oznaczanym przez $F(M; r)$, nazywamy macierz o elementach postaci*

$$[F(M, r)]_{ij} = \frac{\partial \Phi(M; r)}{\partial m_{ij}}, \quad i, j = 1, 2, \dots, r,$$

gdzie m_{ij} oznaczają elementy macierzy M .

Twierdzenie 6.1. *Załóżmy, że kryterium planowania Φ jest funkcją jednorodną, monotoniczną, wklęsłą i różniczkowalną. Zakładamy także, że dla dowolnej nieosobliwej i osiągalnej w danym problemie planowania macierzy informacyjnej M , macierz $F(M; r)$ jest nieujemnie określona, a ponadto macierz ta spełnia warunek:*

$$F(M_1 \oplus M_2, r_1 + r_2) = F(M_1, r_1) \oplus F(M_2, r_2) \quad (6.2.16)$$

dla dowolnych nieosobliwych $r_1 \times r_1$ macierzy M_1 i $r_2 \times r_2$ macierzy M_2 .

Niech $\hat{\xi}^{(1)} \in \Xi(\mathbf{X}_1)$ i $\hat{\xi}^{(2)} \in \Xi(\mathbf{X}_2)$ będą planami $\Phi(\cdot; r_1)$ i $\Phi(\cdot; r_2)$ optymalnymi dla modeli częściowych opisanych w Definicji 6.4, tzn., dla modeli $\alpha^T \mathbf{v}_1(\mathbf{x}^{(1)})$ i $\beta^T \mathbf{v}_2(\mathbf{x}^{(2)})$, odpowiednio. Załóżmy, że dla planów tych zachodzi

$$M(\hat{\xi}^{(1)} \otimes \hat{\xi}^{(2)}) = M_1(\hat{\xi}^{(1)}) \oplus M_2(\hat{\xi}^{(2)}) \quad (6.2.17)$$

a macierz $M(\hat{\xi}^{(1)} \otimes \hat{\xi}^{(2)})$ jest nieosobliwa. Przy tych założeniach plan $\hat{\xi}^{(1)} \otimes \hat{\xi}^{(2)}$, będący produktem planów dla bloków zmiennych, jest $\Phi(\cdot; r_1 + r_2)$ optymalny w zadaniu estymacji parametrów modelu opisanego w Definicji 6.4.

Rezultat ten wynika z Twierdzenia 4.2, a szczegóły dowodu można znaleźć w [136]. Warunek (6.2.16) spełniają kryteria D-, A-, L_p -optymalności oraz kryterium L-optymalności, przy założeniu, że macierz wag A w tym kryterium ma postać $A = A_1 \oplus A_2$, gdzie A_1 i A_2 są macierzami nieujemnie określonymi o rozmiarach $r_1 \times r_1$ i $r_2 \times r_2$, odpowiednio. Weryfikacja warunku (6.2.17) wymaga wcześniejszego obliczenia optymalnych planów dla modeli częściowych. Pełnej weryfikacji tego warunku można uniknąć, jeżeli spełnione są założenia Lematu 6.1, gdyż są one warunkami dostatecznymi dla (6.2.17).

Twierdzenie 6.1 łatwo uogólnić na przypadek modeli addytywnych względem większej liczby grup zblokowanych zmiennych niezależnych. Jako elementarny wniosek z tego twierdzenia otrzymamy następujący klasyczny rezultat.

Wniosek 6.1. *Rozważmy problem planowania dla modelu*

$$\bar{y}(x) = a^{(0)} + \sum_{k=1}^{r-1} a^{(k)} x^{(k)} \quad (6.2.18)$$

w obszarze $X = [-1, 1]^{r-1}$. Plan D- i A-optymalny dla tego modelu skupiony jest we wszystkich wierzchołkach kostki X z wagami $1/2^{r-1}$.

Dla dowodu zauważmy, że A- i D-optymalne plany dla modeli częściowych rozpiętych przez $[1, x^{(1)}]^T, x^{(j)}, j = 2, 3, \dots, r-1$ skupione są w punktach ± 1 z wagami $1/2$. Teraz wystarczy zauważyć, że plany te są symetryczne względem funkcji $x^{(j)}, j = 2, 3, \dots, r-1$ rozpinających modele częściowe. Warunek (6.2.17) jest zatem spełniony na mocy Lematu 6.1.

6.3. Plany dla modeli z pełnym zestawem interakcji zblokowanych zmiennych

Podobnie jak w poprzednim rozdziale, przyjmiemy, że zmienne wejściowe zostały podzielone na bloki i ponumerowane w ten sposób, by $\mathbf{x}^{(1)} = [x^{(1)}, \dots, x^{(s_0)}]^T$

oraz $\mathbf{x}^{(2)} = [x^{(s_0+1)}, \dots, x^{(s)}]^T$ tworzyły razem wektor x . Modelem z pełnym zestawem interakcji dla bloków nazywać będziemy zależność postaci:

$$\bar{y}(x) = a^T [\mathbf{v}_1(\mathbf{x}^{(1)}) \otimes \mathbf{v}_2(\mathbf{x}^{(2)})], \quad (6.3.19)$$

gdzie a jest kolumnowym wektorem nieznanymi parametrów,

$$\dim(a) = \dim(\mathbf{v}_1(\mathbf{x}^{(1)})) \cdot \dim(\mathbf{v}_2(\mathbf{x}^{(2)})).$$

\otimes oznacza iloczyn Kroneckera wektorów (Definicja 15.1 w Dodatku). W (6.3.19) dopuszczamy też większą liczbę czynników iloczynu Kroneckera.

Uwaga 6.3. *Czytelnik zauważył, że ten sam symbol \otimes stosujemy zarówno dla iloczynu Kroneckera, jak i dla oznaczenia operacji tworzenia planów produktowych. Ta zbieżność oznaczeń wynika z tradycji stosowania tego symbolu w algebrze oraz w teorii miary i funkcji. Zbieżność ta nie jest przypadkowa. Wystarczy porównać zapis dla funkcji $f \otimes g \stackrel{\text{def}}{=} f(x)g(t)$ i zapis*

$$a \otimes b \stackrel{\text{def}}{=} [a_1 b, a_2 b, \dots, a_m b]^T$$

$$= [a_1 b_1, a_1 b_2, \dots, a_1 b_n, a_2 b_1, a_2 b_2, \dots, a_2 b_n, \dots, a_m b_1, a_m b_2, \dots, a_m b_n]^T$$

dla wektorów $a \in R^m$, $b \in R^n$, traktowanych jako funkcje dyskretnych argumentów a_i oraz b_j .

Zgodnie z tradycją teorii planowania eksperymentu, składniki modelu będące iloczynami funkcji zależnych od zmiennych wejściowych nazywamy interakcjami. Iloczyny takie mają za zadanie modelowanie łącznego wpływu grup zmiennych na wyjście modelu.

Model (6.3.19) nazywamy modelem z pełnym zestawem interakcji dla bloków, gdyż wystąpią w nim wszystkie możliwe iloczyny składowych wektorów $\mathbf{v}_1(\mathbf{x}^{(1)})$ oraz $\mathbf{v}_2(\mathbf{x}^{(2)})$.

Rozważmy dwa bloki zmiennych wejściowych $\mathbf{x}^{(1)} = x^{(1)}$, $\mathbf{x}^{(2)} = x^{(2)} \cdot x^{(3)}$ oraz bloki zmiennych przekształconych $\mathbf{v}_1(\mathbf{x}^{(1)}) = [1, x^{(1)}]^T$, $\mathbf{v}_2(\mathbf{x}^{(2)}) = [1, x^{(2)} \cdot x^{(3)}]^T$. Model z pełnym zestawem interakcji tych dwóch bloków zmiennych przekształconych ma postać:

$$\bar{y}(x) = a^{(1)} + a^{(2)}x^{(2)}x^{(3)} + a^{(3)}x^{(1)} + a^{(4)}x^{(1)}x^{(2)}x^{(3)}. \quad (6.3.20)$$

Zauważmy, że w modelu tym nie występują wszystkie możliwe iloczyny zmiennych $x^{(1)}$, $x^{(2)}$, $x^{(3)}$. Wszystkie interakcje między pierwszym blokiem zmiennych $x^{(1)}$ i blokiem drugim $x^{(2)}x^{(3)}$ sprowadzają się do jednego (ostatniego) składnika w (6.3.20).

Model nazywać będziemy modelem z pełnym zestawem interakcji dla zmiennych, jeżeli możliwe jest przedstawienie wszystkich czynników tak, by każdy z nich

zależał jedynie od pojedynczej zmiennej niezależnej. Słowo czynnik oznacza taki składnik modelu, przy którym występuje nieznaną parametr. Dokładniej, jeśli $x = [x^{(1)}, x^{(2)}, \dots, x^{(s)}]^T$ jest wektorem wszystkich wielkości wejściowych, to model z pełnym zestawem interakcji wszystkich zmiennych ma postać

$$\bar{y}(x) = a^T \prod_{k=1}^s \mathbf{v}_k(x^{(k)}), \quad (6.3.21)$$

gdzie $\mathbf{v}_k(x^{(k)})$ są wektorami kolumnowymi funkcji zależnych jedynie od k -tej wielkości wejściowej $x^{(k)}$. Jeśli $\dim(\mathbf{v}_k(x^{(k)})) = r_k$, to liczba wszystkich składników modelu, a zatem także $r = \dim(a)$ wynosi $\prod_{k=1}^s r_k$. Można powiedzieć, że model z pełnym zestawem interakcji dla zmiennych, to model z pełnym zestawem interakcji dla bloków, z tym, że teraz każdy z bloków składa się z jednej zmiennej niezależnej.

Przykłady te ilustrują różnice między wprowadzonymi klasami modeli:

$$\bar{y}(x) = a^{(0)} + a^{(1)}x^{(1)} + a^{(2)}x^{(2)} + a^{(3)}x^{(1)}x^{(2)}, \quad (6.3.22)$$

$$\bar{y}(x) = a^{(0)} + a^{(1)}x^{(1)} + a^{(2)}x^{(2)} + a^{(3)}x^{(1)}x^{(2)} + a^{(4)}(x^{(1)})^2, \quad (6.3.23)$$

$$\bar{y}(x) = a^{(0)} + a^{(1)}x^{(1)} + a^{(2)}x^{(2)} + a^{(3)}x^{(3)} + a^{(4)}x^{(1)}x^{(2)} + a^{(5)}x^{(1)}x^{(3)}. \quad (6.3.24)$$

Model (6.3.22) ma pełen zestaw interakcji, gdyż jest rozpięty przez

$$[1, x^{(1)}]^T \otimes [1, x^{(2)}]^T,$$

natomiast (6.3.23) nie spełnia wymagań nałożonych przez (6.3.21). Rzeczywiście, w modelu tym wystąpił czynnik $(x^{(1)})^2$, ale nie wystąpił składnik będący jego iloczynem z $x^{(2)}$. Podobnie, w modelu (6.3.24) brakuje członu interakcji $x^{(1)}x^{(2)}x^{(3)}$.

Ważnym przypadkiem szczególnym jest sytuacja, gdy wektory rozpinające funkcje regresji zawierają funkcję tożsamościowo równą 1, tzn.

$$\mathbf{v}_j(x^{(j)}) = [1, v_j^{(2)}(x^{(j)}), \dots, v_j^{(r_j)}(x^{(j)})]^T, \quad j = 1, 2. \quad (6.3.25)$$

Wtedy model z pełnym zestawem interakcji dwóch zmiennych można zapisać tak:

$$\begin{aligned} \bar{y}(x) = & \alpha^0 + \sum_{k=2}^{r_1} \alpha_k v_1^{(k)}(x^{(1)}) \\ & + \sum_{k=2}^{r_2} \beta_k v_2^{(k)}(x^{(2)}) + \sum_{k=2}^{r_1} \sum_{l=2}^{r_2} \gamma_{lk} v_1^{(k)}(x^{(1)}) v_2^{(l)}(x^{(2)}), \end{aligned} \quad (6.3.26)$$

gdzie, uporządkowane, współczynniki α_k , β_k , γ_{lk} tworzą wektor a o $r_1 \cdot r_2$ składowych. Zauważmy, że w rozważanym przypadku model z pełnym zestawem interakcji dla zmiennych zawiera w sobie wszystkie składniki modelu w pełni addytywnego względem wszystkich zmiennych.

Definicja 6.8. Modelem z pełnym zestawem interakcji względem zblokowanych zmiennych $\mathbf{x}^{(1)}$ i $\mathbf{x}^{(2)}$ nazywamy

$$\bar{y}(x) = a^T [\mathbf{v}_1(\mathbf{x}^{(1)}) \otimes \mathbf{v}_2(\mathbf{x}^{(2)})], \quad (6.3.27)$$

gdzie a jest kolumnowym wektorem nieznanymi parametrów,

$$\dim(a) = \dim(\mathbf{v}_1(\mathbf{x}^{(1)})) \dim(\mathbf{v}_2(\mathbf{x}^{(2)})).$$

Modelami częściowymi dla funkcji (6.3.27) nazywamy wyrażenia $\alpha^T \mathbf{v}_1(\mathbf{x}^{(1)})$ oraz $\beta^T \mathbf{v}_2(\mathbf{x}^{(2)})$, określone odpowiednio na X_1 oraz X_2 . W modelach tych $\alpha \in \mathbf{R}^{\dim(\mathbf{v}_1(\mathbf{x}^{(1)}))}$ oraz $\beta \in \mathbf{R}^{\dim(\mathbf{v}_2(\mathbf{x}^{(2)}))}$ są kolumnowymi wektorami stałych parametrów.

Zadanie planowania dla (6.3.27) rozpatrywane będzie na zbiorze $X_1 \times X_2$, gdzie $\mathbf{x}^{(1)} \in X_1$, $\mathbf{x}^{(2)} \in X_2$.

Modelom częściowym modelu (6.3.27) można nadać następującą interpretację: $\alpha^T \cdot \mathbf{v}_1(\mathbf{x}^{(1)})$ to zależność wyjścia od zblokowanych wielkości wejściowych $\mathbf{x}^{(1)}$, przy ustalonych wartościach bloku wejść $\mathbf{x}^{(2)}$. Modele częściowe można także traktować zupełnie formalnie – jako narzędzie znajdowania planów optymalnych.

Rozważmy prosty przykład:

$$\bar{y}(x) = a^T \left\{ \begin{bmatrix} 1 \\ x^{(1)} \end{bmatrix} \otimes \begin{bmatrix} 1 \\ x^{(2)} \end{bmatrix} \otimes \begin{bmatrix} 1 \\ x^{(3)} \end{bmatrix} \right\}. \quad (6.3.28)$$

Można w nim wyróżnić dwie różne pary modeli częściowych. Pierwsza para rozpięta jest przez funkcje

$$\begin{bmatrix} 1 \\ x^{(1)} \end{bmatrix} \otimes \begin{bmatrix} 1 \\ x^{(2)} \end{bmatrix} \quad \text{oraz} \quad \begin{bmatrix} 1 \\ x^{(3)} \end{bmatrix}.$$

Para druga rozpinana jest przez zestawy

$$\begin{bmatrix} 1 \\ x^{(1)} \end{bmatrix} \quad \text{oraz} \quad \begin{bmatrix} 1 \\ x^{(2)} \end{bmatrix} \otimes \begin{bmatrix} 1 \\ x^{(3)} \end{bmatrix}.$$

Podany przykład pozwala zauważyć, że wyróżnianie modeli częściowych nie jest procedurą o jednoznacznym wyniku. W praktyce dążyć powinniśmy do możliwie daleko idącego rozbicia modelu na modele częściowe. W omawianym przykładzie korzystne jest wyróżnienie trzech modeli częściowych, gdyż wówczas problem znajdowania planu optymalnego upraszcza się najbardziej.

Lemat 6.2. *Rozważmy model z pełnym zestawem interakcji zblokowanych zmiennych*

$$\bar{y}(x) = a^T [\mathbf{v}_1(\mathbf{x}^{(1)}) \otimes \mathbf{v}_2(\mathbf{x}^{(2)})], \quad \mathbf{x}^{(1)} \in X_1, \mathbf{x}^{(2)} \in X_2 \quad (6.3.29)$$

i dowolny plan produktowy $\xi^{(1)} \otimes \xi^{(2)}$ na $X_1 \times X_2$, gdzie plany $\xi^{(1)}, \xi^{(2)}$ są postaci

$$\xi^{(l)} = \begin{bmatrix} \mathbf{x}_1^{(l)} & \mathbf{x}_2^{(l)} & \dots & \mathbf{x}_{m_1}^{(l)} \\ p_1^{(l)} & p_2^{(l)} & \dots & p_{m_1}^{(l)} \end{bmatrix}, \quad l = 1, 2. \quad (6.3.30)$$

Wówczas macierz informacyjna planu produktowego dana jest wzorem

$$M(\xi^{(1)} \otimes \xi^{(2)}) = M(\xi^{(1)}) \otimes M(\xi^{(2)}), \quad (6.3.31)$$

gdzie macierze informacyjne dla modeli częściowych mają postać

$$M^{(l)}(\xi^{(l)}) = \sum_{i=1}^{m_1} p_i^{(l)} \mathbf{v}^{(l)}(\mathbf{x}_i^{(l)}) \left(\mathbf{v}^{(l)}(\mathbf{x}_i^{(l)}) \right)^T, \quad l = 1, 2. \quad (6.3.32)$$

Dowód tego lematu, jako czysto algebraiczny, pomijamy.

Podamy teraz sposób konstruowania optymalnego planu produktowego dla omawianego modelu z pełnym zestawem interakcji względem zblokowanych zmiennych (6.3.29). Po pierwsze, należy znaleźć plan D-optymalny $\hat{\xi}^{(1)} \in \Xi(X_1)$ dla modelu częściowego $\alpha^T \cdot \mathbf{v}_1(\mathbf{x}^{(1)})$. W tym celu należy rozwiązać zadanie:

$$\max_{\xi^{(1)} \in \Xi(X_1)} \det M^{(1)}(\xi^{(1)}) = \det M^{(1)}(\hat{\xi}^{(1)}), \quad (6.3.33)$$

gdzie dla każdego planu $\xi^{(1)} \in \Xi(X_1)$

$$\xi^{(1)} = \begin{bmatrix} \mathbf{x}_1^{(1)} & \mathbf{x}_2^{(1)} & \dots & \mathbf{x}_{m_1}^{(1)} \\ p_1^{(1)} & p_2^{(1)} & \dots & p_{m_1}^{(1)} \end{bmatrix} \quad (6.3.34)$$

macierz informacyjna dana jest wzorem

$$M^{(1)}(\xi^{(1)}) = \sum_{i=1}^{m_1} p_i^{(1)} \mathbf{v}^{(1)}(\mathbf{x}_i^{(1)}) \left(\mathbf{v}^{(1)}(\mathbf{x}_i^{(1)}) \right)^T. \quad (6.3.35)$$

W drugim kroku należy znaleźć plan D-optymalny $\hat{\xi}^{(2)} \in \Xi(X_2)$ dla modelu częściowego $\beta^T \mathbf{v}_2(\mathbf{x}^{(2)})$. W tym celu należy rozwiązać zadanie analogiczne do (6.3.33). W ostatnim kroku konstruujemy plan produktowy $\xi^* = \hat{\xi}^{(1)} \otimes \hat{\xi}^{(2)}$.

Twierdzenie 6.2. *Uzyskany w opisany sposób plan $\xi^* = \hat{\xi}^{(1)} \otimes \hat{\xi}^{(2)}$, jest planem D- optymalnym dla estymacji wektora parametrów a modelu (6.3.29) w obszarze planowania $X_1 \times X_2$.*

Dowód. Dowód przedstawiamy przy założeniu, że wariancja zakłóceń jest stała w całym obszarze planowania. Dowód w nieco ogólniejszym przypadku, gdy wariancja jest faktoryzowalną funkcją zmiennych przestrzennych, znaleźć można w [116].

Plany $\hat{\xi}^{(1)}$ i $\hat{\xi}^{(2)}$, jako D- optymalne dla modeli częściowych, spełniają następujące warunki (wynikające z twierdzenia Kiefera i Wolfowitza):

$$\sup_{\mathbf{x}^{(1)} \in X_1} \mathbf{v}_1^T(\mathbf{x}^{(1)}) [M^{(1)}(\hat{\xi}^{(1)})]^{-1} \mathbf{v}_1(\mathbf{x}^{(1)}) = r_1 \quad (6.3.36)$$

$$\sup_{\mathbf{x}^{(2)} \in X_2} \mathbf{v}_2^T(\mathbf{x}^{(2)}) [M^{(2)}(\hat{\xi}^{(2)})]^{-1} \mathbf{v}_2(\mathbf{x}^{(2)}) = r_2. \quad (6.3.37)$$

Jednocześnie, dla planu produktowego $\xi^* = \hat{\xi}^{(1)} \otimes \hat{\xi}^{(2)}$ zachodzi równość $M(\xi^*) = M^{(1)}(\hat{\xi}^{(1)}) \otimes M^{(2)}(\hat{\xi}^{(2)})$. Plany $\hat{\xi}^{(1)}$, $\hat{\xi}^{(2)}$ prowadzą do nieosobliwych macierzy informacyjnych $M^{(1)}(\hat{\xi}^{(1)})$ i $M^{(2)}(\hat{\xi}^{(2)})$, gdyż są one planami D- optymalnymi dla modeli częściowych. Zatem, zgodnie z własnościami iloczynu Kroneckera macierzy (por. Dodatek), nieosobliwa jest także macierz $M(\xi^*)$, co prowadzi do

$$\begin{aligned} \sup_{x \in X_1 \times X_2} v^T(x) M^{-1}(\hat{\xi}) v(x) &= \sup_{\mathbf{x}^{(1)} \in X_1} \mathbf{v}_1^T(\mathbf{x}^{(1)}) [M^{(1)}(\hat{\xi}^{(1)})]^{-1} \mathbf{v}_1(\mathbf{x}^{(1)}) \quad (6.3.38) \\ &\times \sup_{\mathbf{x}^{(2)} \in X_2} \mathbf{v}_2^T(\mathbf{x}^{(2)}) [M^{(2)}(\hat{\xi}^{(2)})]^{-1} \mathbf{v}_2(\mathbf{x}^{(2)}) = r_1 r_2 = r. \end{aligned}$$

Wystarczy teraz ponownie odwołać się do twierdzenia Kiefera–Wolfowitza, aby z warunku (6.3.38) odczytać D- optymalność planu ξ^* w modelu (6.3.27). •

Warto dodać, że pierwsze wskazówki, że plany optymalne mogą mieć postać taką jak w Twierdzeniu 6.2 znaleźć można w pracy [55], lecz tam rozważano szczególnie przypadek modeli trygonometrycznych.

Za pracą [116] przedstawimy uogólnienie rezultatu zawartego w Twierdzeniu 6.2 na inne kryteria planowania. Wymagania stawiane kryteriom planowania omówiono w rozdziale 4.1. Tutaj także wymagać będziemy, aby kryterium Φ było funkcją wklęsłą i różniczkowalną w zbiorze wszystkich nieosobliwych i osiągalnych macierzy informacyjnych. Ograniczamy się tu do tzw. problemów regularnych, to znaczy takich, dla których w rozważanym zadaniu planowania istnieje rozwiązanie Φ -optymalne o nieosobliwej macierzy informacyjnej.

Przypomnijmy, że przez $F(M; r)$ oznaczamy macierz gradientu kryterium Φ , której elementy są postaci

$$[F(M, r)]_{ij} = \frac{\partial \Phi(M; r)}{\partial m_{ij}}, \quad i, j = 1, 2, \dots, r,$$

gdzie m_{ij} oznaczają elementy macierzy M . Od macierzy $F(M; r)$ wymagamy, by **F1)** dla dowolnych nieosobliwych M_1 ($\dim M_1 = r_1 \times r_1$) oraz M_2 ($\dim M_2 = r_2 \times r_2$) zachodziło

$$F(M_1 \otimes M_2; r_1 r_2) = F(M_1; r_1) \otimes F(M_2; r_2). \quad (6.3.39)$$

Nażłone wcześniej (s. 36) wymagania monotoniczności i różniczkowalności kryterium Φ implikują:

F2) dla dowolnej nieosobliwej macierzy informacyjnej M , macierz $F(M; r)$ jest nieujemnie określona.

Precyzyjny dowód tej własności podano w [199] (Lemat 3.4, s. 44). Można także wykazać (por. [116]), że warunek F1) (a zatem także F2)) spełnia wiele kryteriów planowania, a w szczególności: kryteria D-, L_p -, A-, Q-optymalności, i ekstrapolacji w zadanym punkcie. Jak wspomniano wcześniej, ostatnie trzy kryteria są szczególnymi przypadkami kryterium L-optymalności. W używanej tu notacji, ta klasa kryteriów ma postać $\Phi(M(\xi), r) = -\text{tr}[WM^{-1}(\xi)]$, gdzie W jest $r \times r$ macierzą wagową. Aby warunek F1) był spełniony, macierz wag musi mieć odpowiednią strukturę, a mianowicie powinna dać się przedstawić jako iloczyn Kroneckera $W_1 \otimes W_2$ macierzy wagowych odpowiadających modelom częściowym. Jeśli przyjmiemy model (6.3.42), to faktoryzacja taka zachodzi między innymi dla kryteriów A-, Q-optymalności i ekstrapolacji w zadanym punkcie.

Twierdzenie 6.3. *Załóżmy, że kryterium planowania Φ jest funkcją jednorodną, monotoniczną, wklęsłą i różniczkowalną. Niech spełnione będzie założenia F1). Załóżmy, że plany $\hat{\xi}^{(1)} \in \Xi(X_1)$ oraz $\hat{\xi}^{(2)} \in \Xi(X_2)$ są Φ -optymalne dla modeli częściowych, to znaczy zachodzi*

$$\max_{\xi_1 \in \Xi(X_1)} \Phi(M(\xi_1); r_1) = \Phi(M(\hat{\xi}^{(1)}); r_1) \quad (6.3.40)$$

oraz

$$\max_{\xi_2 \in \Xi(X_2)} \Phi(M(\xi_2); r_2) = \Phi(M(\hat{\xi}^{(2)}); r_2). \quad (6.3.41)$$

Wówczas plan $\hat{\xi}^{(1)} \otimes \hat{\xi}^{(2)}$ jest planem Φ -optymalnym w zadaniu estymacji parametrów modelu

$$\bar{y}(x) = a^T [\mathbf{v}_1(\mathbf{x}^{(1)}) \otimes \mathbf{v}_2(\mathbf{x}^{(2)})]. \quad (6.3.42)$$

Uogólnienie tego rezultatu na większą liczbę modeli częściowych jest natychmiastowe. Twierdzenie to udowodniono w [116].

Twierdzenia 6.3 oraz 6.1 pozwalają na poszukiwanie optymalnych planów dla modeli, które nie są ani modelami w pełni addytywnymi, ani modelami z pełnym zestawem interakcji względem zblokowanych zmiennych, lecz da się w nich wydzieleć modele częściowe o takich cechach. Następujący rezultat ilustruje zarysowaną poprzednio ogólną ideę.

Rozważmy model o postaci

$$\bar{y} = a^T \{[g(x) \oplus h(u)] \otimes R(z)\}, \quad (6.3.43)$$

gdzie $g(x)$, $h(u)$ i $R(z)$ są wektorami liniowo niezależnych funkcji skalarnych argumentów, określonymi, odpowiednio, na odcinkach domkniętych X , U , Z .

Φ -optymalne plany częściowe, odpowiadające modelom częściowym rozpinanym przez $g(x)$, $h(u)$ i $R(z)$, oznaczamy, odpowiednio, przez $\hat{\xi}_x$, $\hat{\xi}_u$ i $\hat{\xi}_z$.

Twierdzenie 6.4. *Niech spełnione będą założenia Twierdzenia 6.1 odnoszące się do wskaźnika optymalności planowania. Ponadto, załóżmy, że przynajmniej jeden z planów $\hat{\xi}_x$ lub $\hat{\xi}_u$ jest symetryczny, odpowiednio, względem $h(x)$ lub $g(u)$. Wówczas plan $\hat{\xi}_x \otimes \hat{\xi}_u \otimes \hat{\xi}_z$ jest Φ -optymalny dla estymacji a w modelu (6.3.43) w obszarze planowania $X \times U \times Z$.*

Dowód pomijamy. Dalsze rezultaty na temat planów produktowych znaleźć można w pracach [168], [169] i [170].

6.4. Analityczne wyznaczanie planów

W podrozdziale tym przedstawiamy zastosowania Twierdzeń 6.1–6.4 do znajdowania planów dla modeli addytywnych, modeli z pełnym zestawem interakcji i modeli z niepełnym zestawem interakcji, które jednak mają postać addytywną względem zablokowanych zmiennych.

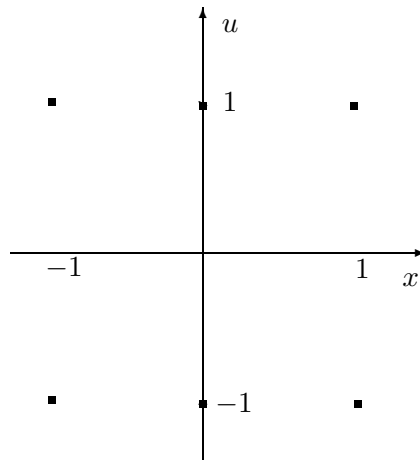
Przykłady dla modeli addytywnych

Formalne uogólnienie Twierdzenia 6.1 na przypadek wielowymiarowego modelu addytywnego względem każdej ze zmiennych wejściowych jest natychmiastowe. W tym przypadku, uogólnienie Lematu 6.1 sprowadza się do wymagania, by w każdej parze planów częściowych, które występują w danym problemie, co najmniej jeden z planów był symetryczny względem stowarzyszonego z nim modelu częściowego.

Optymalność pełnych planów czynnikowych

Poprzednie spostrzeżenie pozwala natychmiast udowodnić klasyczny rezultat, a mianowicie D- i A-optymalność planu o wagach $1/2^s$ przypisanych wszystkim wierzchołkom kostki $[-1,1]^s$ w problemie estymacji parametrów modelu

$$\bar{y} = a_0 + \sum_{i=1}^s a_i x^{(i)}. \quad (6.4.44)$$



Rys. 6.1. Nośnik planu D-optimalnego dla estymacji parametrów modelu (6.4.45)

Plan dla prostego modelu addytywnego

Plan D-optimalny dla modelu

$$\bar{y} = a_0 + a_1 x + a_2 x^2 + a_3 u, \quad (x, u) \in [-1, 1] \times [-1, 1] \quad (6.4.45)$$

ma nośnik pokazany na rysunku 6.1, z wagami $1/6$ przypisanymi każdemu z tych punktów. Łatwo też wykazać, że plan A-optimalny dla modelu (6.4.45) ma nośnik pokazany na rysunku 6.1 z wagami $1/8$ przypisanymi punktom wierzchołkowym i $1/4$ w obu punktach $(0, 1)$ i $(0, -1)$.

Przykłady dla modeli z interakcjami

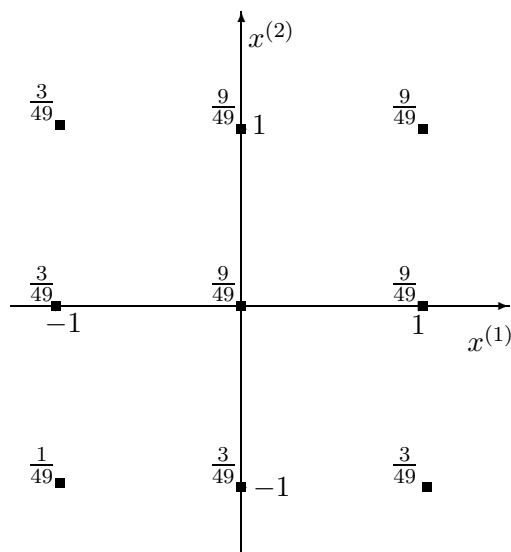
Przedstawiamy przykłady analitycznego znajdowania planów dla estymacji parametrów w modelach z interakcjami, zarówno z pełnym zestawem interakcji, jak i takich, w których – mimo braku pewnych interakcji – potrafimy przedstawić model w postaci addytywnej względem grup zmiennych.

Model z pojedynczą interakcją

Rozważmy model o postaci

$$\bar{y}(x) = a^T \left\{ \begin{bmatrix} 1 \\ x^{(1)} \end{bmatrix} \otimes \begin{bmatrix} 1 \\ x^{(2)} \end{bmatrix} \right\}, \quad (6.4.46)$$

który w klasycznym zapisie ma postać: $\bar{y}(x) = a^{(1)} + a^{(2)}x^{(1)} + a^{(3)}x^{(2)} + a^{(4)}x^{(1)}x^{(2)}$. W obszarze $X = [-1, 1] \times [-1, 1]$ szukamy D-optimalnego planu dla estymacji

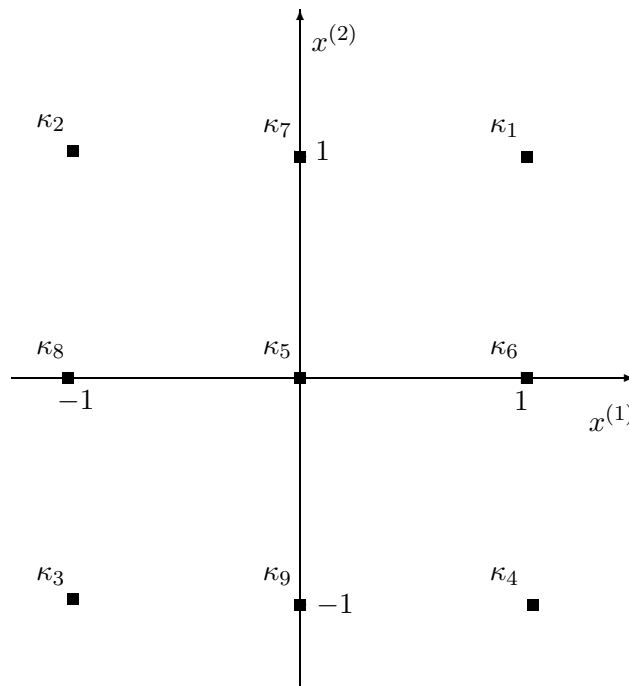


Rys. 6.2. Plan optymalny dla regresji kwadratowej względem każdej ze zmiennych $x^{(1)}$ i $x^{(2)}$ przy ekstrapolacji w punkcie $(2, 2)$

wektora a . Model częściowy dla pierwszego wejścia ma postać: $\alpha^{(1)} + \alpha^{(2)}x^{(1)}$. Plan D-optymalny dla tego modelu na $[-1, 1]$ jest znany $\hat{\xi}^{(1)} = \begin{bmatrix} -1 & 1 \\ 1/2 & 1/2 \end{bmatrix}$. Wobec symetrii, zupełnie analogicznie wygląda drugi model częściowy i odpowiadający mu plan optymalny plan częściowy. W rezultacie zastosowania Twierdzenia 6.2 stwierdzamy, że następujący plan

$$\xi^* = \begin{bmatrix} (-1, -1) & (-1, 1) & (1, -1) & (1, 1) \\ 1/4 & 1/4 & 1/4 & 1/4 \end{bmatrix} \quad (6.4.47)$$

jest D-optymalny dla modelu (6.4.46). Fakt ten jest od dawna znany, lecz tu uzyskuje się go elementarnie. Zauważmy, że plan (6.4.47) jest także D-optymalny dla estymacji parametrów modelu $\bar{y}(x) = a^{(1)} + a^{(2)}x^{(1)} + a^{(3)}x^{(2)}$. Jak się okazuje, jest to wyraz ogólniejszej prawidłowości (por. [164]). Uogólnienie tego przykładu na regresję o większej liczbie zmiennych jest natychmiastowe, a punkty planu D-optymalnego rozmieszczone są we wszystkich wierzchołkach hiperkostki. Przypisane są im jednakowe wagi.



Rys. 6.3. Nośnik planów dla regresji kwadratowych względem zmiennych $x^{(1)}$ i $x^{(2)}$

Model kwadratowy z interakcjami

Rozważmy model będący wielomianem stopnia czwartego:

$$\bar{y}(x) = a^T \left\{ \begin{bmatrix} 1 \\ x^{(1)} \\ (x^{(1)})^2 \end{bmatrix} \otimes \begin{bmatrix} 1 \\ x^{(2)} \\ (x^{(2)})^2 \end{bmatrix} \right\}. \quad (6.4.48)$$

Łączna liczba estymowanych parametrów wynosi więc 9. Modele częściowe mają postać: $\alpha^{(1)} + \alpha^{(2)}x^{(i)} + \alpha^{(3)}(x^{(i)})^2$, a D-optymalnymi planami częściowymi na $[-1, 1]$ są $\hat{\xi}^{(i)} = \begin{bmatrix} -1 & 0 & 1 \\ 1/3 & 1/3 & 1/3 \end{bmatrix}$, $i = 1, 2$. D-optymalny plan dla (6.4.48) na $[-1, 1] \times [-1, 1]$ skupiony jest zatem we wszystkich dziewięciu punktach pokazanych na rysunku 6.3 z jednakowymi wagami $1/9$.

Planowanie w obszarze o kształcie walca

Rozważmy zdania planowania D- optymalnego dla funkcji regresji o postaci:

$$\begin{aligned} \bar{y}(x) = & a^{(0)} + a^{(1)}x^{(1)} + a^{(2)}x^{(2)} \\ & + a^{(3)}x^{(3)} + a^{(4)}x^{(1)}x^{(3)} + a^{(5)}x^{(2)}x^{(3)} \\ & + a^{(6)}(x^{(3)})^2 + a^{(7)}x^{(1)}(x^{(3)})^2 + a^{(8)}x^{(2)}(x^{(3)})^2. \end{aligned} \quad (6.4.49)$$

Obszarem planowania jest walec – uwzględniając brzeg i wnętrze. Podstawa walca ma promień 1 i umieszczona jest w płaszczyźnie zmiennych $x^{(1)}$ i $x^{(2)}$, podczas gdy tworząca umieszczona jest wzdłuż pionowej osi $x^{(3)} \in [-1, 1]$. Obszar planowania ma zatem postać $X_1 \times X_2$, gdzie $X_1 = \{(x^{(1)}, x^{(2)}) : (x^{(1)})^2 + (x^{(2)})^2 \leq 1\}$, $X_2 = \{x^{(3)} : x^{(3)} \in [-1, 1]\}$. Funkcja (6.4.49) może być zapisana w postaci

$$\bar{y}(x) = a^T \left(\mathbf{v}_1(\mathbf{x}^{(1)}) \otimes \mathbf{v}_2(\mathbf{x}^{(2)}) \right) \quad (6.4.50)$$

gdzie $\mathbf{v}_1(\mathbf{x}^{(1)}) = [1, x^{(1)}, x^{(2)}]$, $\mathbf{v}_2(\mathbf{x}^{(2)}) = [1, x^{(3)}, (x^{(3)})^2]^T$, a dla zachowania konsekwencji w notacji przyjęto $\mathbf{x}^{(1)} = (x^{(1)}, x^{(2)})$ oraz $\mathbf{x}^{(2)} = x^{(3)}$. Rozwiązanie zadania D- optymalnego planowania na X_1 , dla modelu częściowego rozpiętego przez $\mathbf{v}_1(\mathbf{x}^{(1)})$, przedstawiamy poniżej.

Rozważmy funkcję regresji $\bar{y}(x) = a^{(0)} + a^{(1)}x^{(1)} + a^{(2)}x^{(2)}$ określoną w kole jednostkowym $X = \{(x^{(1)}, x^{(2)}) : (x^{(1)})^2 + (x^{(2)})^2 \leq 1\}$. Gdy obliczymy macierz informacyjną planu pokazanego na rysunku 6.4 z wagami 1/4 i skorzystamy z twierdzenia Kiefera–Wolfowitza, to stwierdzimy, że plan ten jest D- optymalny.

Natomiast plan D- optymalny dla drugiego modelu częściowego $[1, x^{(3)}, (x^{(3)})^2]$, na $[-1, 1]$, skupiony jest w punktach ± 1 oraz 0 z wagami 1/3. Plan produktowy – zgodnie z Twierdzeniem 6.2 – D- optymalny dla planowania na walcu możemy zatem łatwo obliczyć. Jest on skupiony w 12 punktach, po cztery na dolnej i górnej podstawie walca oraz w jego środkowym przekroju kołowym. W każdym z tych przekrojów punkty rozmieszczone są tak jak na rysunku 6.4.

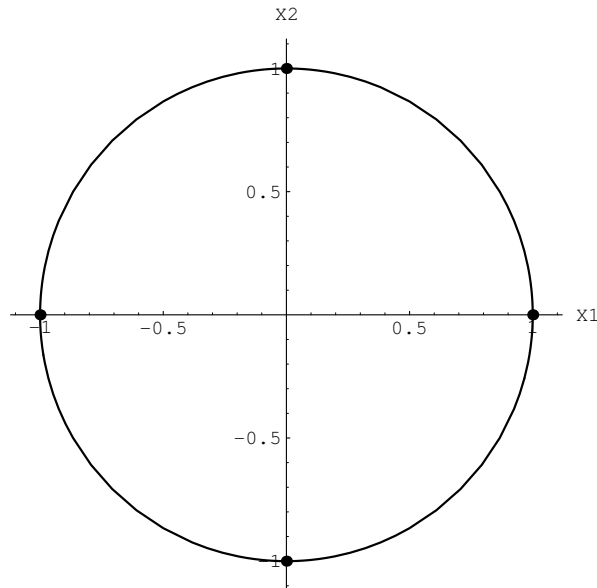
Plan A- optymalny dla wielomianu drugiego stopnia o dwóch zmiennych

Łatwo analitycznie sprawdzić, że plan

$$\begin{bmatrix} -1 & 0 & 1 \\ 1/4 & 1/2 & 1/4 \end{bmatrix} \quad (6.4.51)$$

jest A- optymalny w estymacji parametrów kwadratowej funkcji regresji rozpatrywanej na $X = [-1, 1]$. Rozważmy model

$$\bar{y}(x) = a^T \left\{ \begin{bmatrix} 1 \\ x^{(1)} \\ (x^{(1)})^2 \end{bmatrix} \otimes \begin{bmatrix} 1 \\ x^{(2)} \\ (x^{(2)})^2 \end{bmatrix} \right\} \quad (6.4.52)$$



Rys. 6.4. Plan D-optimalny dla regresji wielomianowej stopnia drugiego na kole

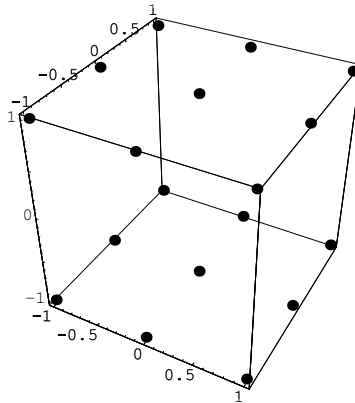
dla którego poszukiwany jest plan A-optimalny na $[-1, 1] \times [-1, 1]$. Nośnik planu optymalnego pokazano na rysunku 6.3. Wagi przypisywane punktom planu A-optimalnego przyjmują jedynie wartości $1/4$, $1/8$ i $1/16$, przy czym wagę $1/4$ ma punkt centralny planu, wagi $1/16$ punkty w narożach obszaru planowania, a wagi $1/8$ – punkty leżące w środkach brzegów kwadratu.

Plan dla ekstrapolacji

Rozważmy ponownie estymację parametrów modelu wielomianowego stopnia czwartego (6.4.52), lecz tym razem celem naszym jest takie zaplanowanie eksperymentu, by minimalizować wariancję predykcji w punkcie $(-2, 2)$. Leży on poza obszarem planowania i dlatego planowanie nazywa się także optymalnym dla ekstrapolacji w zadanym punkcie. Można pokazać (por. [35]), że dla kwadratowej regresji jednej zmiennej na $[-1, 1]$ plan optymalny dla ekstrapolacji w punkcie 2 skupiony jest w punktach $-1, 0, 1$, odpowiednio, z wagami $1/7, 3/7, 3/7$. Plan optymalny w naszym dwuwymiarowym problemie ma zatem postać pokazaną na rysunku 6.2.

Model liniowy o trzech zmiennych

Przypuśćmy, że w modelu liniowym o trzech zmiennych niezależnych x, u, v występują również interakcje między zmiennymi x i u oraz x i v , ale nie występuje



Rys. 6.5. Nośnik planu D- optymalnego

interakcja między u i v . Model zatem ma postać

$$\bar{y} = a_0 + a_1 x + a_2 u + a_3 v + a_4 xu + a_5 xv. \quad (6.4.53)$$

Zauważmy, że – z dokładnością do przenumerowania nieznanymi parametrów w wektorze a – model (6.4.53) zapisać można w postaci

$$\bar{y} = a^T \left\{ \left[\begin{array}{c} 1 \\ u \end{array} \right] \oplus [v] \right\} \otimes \left[\begin{array}{c} 1 \\ x \end{array} \right]. \quad (6.4.54)$$

Plan D- optymalny w obszarze $[-1, 1]^3$ skupiony jest w tych samych punktach co pełny plan czynnikowy typu 2^3 , czyli we wszystkich wierzchołkach kostki. Punktom tym przypisujemy wagi równe $1/8$.

Planowanie w czterowymiarowej przestrzeni o podstawie cylindrycznej

Niech obszarem planowania będzie $B \times [-1, 1] \times [-1, 1]$, gdzie $B = \{(x^{(1)})^2 + (x^{(2)})^2 \leq 1\}$ jest kołem jednostkowym. W obszarze tym określony jest model postaci

$$\bar{y} = a^T \left\{ \left[\begin{array}{c} 1 \\ x^{(1)} \\ x^{(2)} \end{array} \right] \oplus [u] \right\} \otimes \left[\begin{array}{c} 1 \\ v \end{array} \right] \quad (6.4.55)$$

Jak wiadomo, plan D- optymalny w obszarze B dla modelu rozpiętego przez funkcje $[1, x^{(1)}, x^{(2)}]^T$ jest postaci $(0, \pm 1), (\pm 1, 0)$ z wagami $1/4$. Plan D- optymalny dla

modelu (6.4.55) skupiony jest zatem w punktach $(0, \pm 1, \pm 1, \pm 1)$, $(\pm 1, 0, \pm 1, \pm 1)$ z wagami $1/16$. Estymowana funkcja regresji ma postać:

$$\begin{aligned}\bar{y}(x) = & a^{(1)} + a^{(2)}x^{(1)} + a^{(3)}x^{(2)} + a^{(4)}x^{(1)}x^{(2)} \\ & + a^{(5)}(x^{(1)})^2 + a^{(6)}(x^{(2)})^2 + a^{(7)}x^{(3)} + a^{(8)}x^{(1)}x^{(3)} + a^{(9)}x^{(2)}x^{(3)} \\ & + a^{(10)}x^{(1)}x^{(2)}x^{(3)} + a^{(11)}(x^{(1)})^2x^{(3)} + a^{(12)}(x^{(2)})^2x^{(3)}\end{aligned}$$

i określona jest na kostce $[-1, 1]^3$. Funkcji tej nie można przedstawić w postaci iloczynów Kroneckera względem wszystkich trzech zmiennych. Potrafimy natomiast przedstawić ją w postaci częściowo multiplikatywnej (6.3.27), jeśli przyjmiemy:

$$\begin{aligned}\mathbf{v}_2(\mathbf{x}^{(2)}) &= [1, x^{(3)}]^T, \text{ gdzie } \mathbf{x}^{(2)} = [x^{(3)}] \\ \mathbf{v}_1(\mathbf{x}^{(1)}) &= [1, x^{(1)}, x^{(2)}, x^{(1)}x^{(2)}, (x^{(1)})^2, (x^{(2)})^2]^T,\end{aligned}$$

a z powodu jednolitości zapisu utożsamimy $\mathbf{x}^{(1)} = [x^{(1)}, x^{(2)}]^T$. Model częściowy dla $\mathbf{x}^{(1)}$ ma postać:

$$\bar{y}(x) = a^{(1)} + a^{(2)}x^{(1)} + a^{(3)}x^{(2)} + a^{(4)}x^{(1)}x^{(2)} + a^{(5)}(x^{(1)})^2 + a^{(6)}(x^{(2)})^2, \quad (6.4.56)$$

a jego obszarem planowania jest $[-1, 1]^2$. Kiefer wykazał, że plan D-optimalny na kwadracie dla (6.4.56) ma postać pokazaną na rysunku 6.3, przy czym częstotliwości pomiarów w poszczególnych punktach mają postać:

- w wierzchołkach $p_1^{(1)} = p_2^{(1)} = p_3^{(1)} = p_4^{(1)} = 0.14805$,
- w centrum kwadratu $p_5^{(1)} = 0.0962$,
- w środkach boków $p_6^{(1)} = p_7^{(1)} = p_8^{(1)} = p_9^{(1)} = 0.08015$.

Numeracja punktów jest taka jak pokazano na rysunku 6.3. Natomiast model częściowy dla drugiej zmiennej jest funkcją liniową na $[-1, 1]$ i częściowy plan skupiony jest w punktach ± 1 z wagami $1/2$. D-optimalny plan produktowy dla (6.4.56), skupiony będzie zatem w pewnych (nie wszystkich) wierzchołkach, środkach krawędzi oraz dwóch środkach ścian prostopadłościanu. Aby plan ten opisać, przyjmijmy że zmienna $x^{(3)}$ umieszczona jest na pionowej osi układu współrzędnych. Wówczas siatkę przedstawioną na rysunku 6.3 należy umieścić na górnej i dolnej powierzchni kostki. Wagi przypisane tym punktom powinny być takie jak podano wcześniej, z tym że każdą wartość należy podzielić przez dwa. Nośnik tego planu pokazano na rysunku 6.5.

7. Zalety planów produktowych

W rozdziale tym zestawimy zalety planów produktowych zarówno z punktu widzenia numerycznych metod ich poszukiwania jak i ze względu na ich walory numeryczne w zadaniu estymacji parametrów modelu.

7.1. Poszukiwanie i realizacja planów produktowych

Celem tego podrozdziału jest pokazanie, że komponowanie optymalnych planów, będących iloczynem planów dla podmodeli, jest znacznie łatwiejsze numerycznie niż bezpośrednio poszukiwanie optymalnego planu wielowymiarowego.

Komponowanie planów

Zadanie to staje się jeszcze łatwiejsze, gdy posiadamy bazę danych optymalnych planów dla regresji jednej zmiennej lub o zmiennych zblokowanych. Wobec częściowej lub pełnej dekompozycji problemu, możliwe jest nie tylko sekwencyjne, ale także równoległe obliczanie brakujących w bazie planów.

Aby naszkicować źródła oszczędności obliczeniowych i sposób postępowania, posłużymy się modelem z pełnym zestawem interakcji wszystkich zmiennych. Podkreślić jednak należy, że dyskusja ta odnosi się również do modeli z pełnym zestawem interakcji zmiennych zgrupowanych w bloki. Przypomnijmy postać modelu

$$\bar{y}(x) = a^T \cdot \prod_{k=1}^s \mathbf{v}_k(x^{(k)}), \quad (7.1.1)$$

gdzie $\mathbf{v}_k(x^{(k)})$ są wektorami kolumnowymi funkcji zależnych jedynie od k -tej wielkości wejściowej $x^{(k)}$. Gdy $\dim(\mathbf{v}_k(x^{(k)})) = r_k$, to liczba wszystkich składników modelu, a zatem także $r = \dim(a)$, wynosi $\prod_{k=1}^s r_k$.

W celu zilustrowania konsekwencji tego faktu, przyjmijmy najprostsze modele częściowe $\mathbf{v}_k(x^{(k)}) = [1, x^{(k)}]^T$, $k = 1, 2, \dots, s$, odpowiadające liniowym funkcjom regresji. Jeśli obszarami planowania dla modeli częściowych jest $[-1, 1]$, a cały obszar planowania stanowi $X = [-1, 1]^s$, to plan D-optymalny jest oczywiście znany (por. Twierdzenie 4.6). W celach ilustracyjnych wyobraźmy sobie jednak, że plan ten chcemy obliczać numerycznie. Dla modeli częściowych $\mathbf{v}_k(x^{(k)}) = [1, x^{(k)}]^T$ wszystkie $r_k = 2$ i jeśli mamy pięć zmiennych wejściowych, to $r = 2^5 = 64$. Jak

wiemy, plan D-optimalny nie może mieć mniej niż $r = 64$ punktów. Poszukując go metodami numerycznymi opisanymi w poprzednim rozdziale, w każdym kroku znajdować musimy globalne maksimum funkcji $\phi(\xi_j, x)$ zależnej od 5 zmiennych. Sytuację utrudnia to, że wraz ze wzrostem jakości bieżącego planu ξ_j coraz bardziej zaznaczać się będą maksima lokalne funkcji $\phi(\xi_j, x)$. Dalsza poprawa planu związana będzie z ich coraz dokładniejszą lokalizacją. Punktów występowania lokalnego maksimum będzie co najmniej 64 i będzie im odpowiadać zbliżona wartość $\phi(\xi_j, x)$, gdyż dla każdego planu mamy $\sup_{x \in X} \phi(\xi_j, x) \geq r$, podczas gdy dla planu optymalnego $\sup_{x \in X} \phi(\xi_j, x) = r$ i supremum osiągnięte jest w punktach planu.

Zarysowane trudności wystąpiły przy rozpatrywaniu bardzo skromnego przykładu. W zadaniach o większej liczbie zmiennych i bardziej złożonych wektorach rozpinających częściowe funkcje regresji spodziewać się można wykładniczego narastania trudności obliczeniowych.

Jeśli struktura modelu na to pozwala, skorzystać możemy z Twierdzeń 6.1, 6.2 i prowadzić obliczenia w sposób¹ następujący:

Algorytm komponowania planów

Krok 0. Wybierz kryterium planowania. Ustaw licznik numeru zmiennej wejściowej $k = 1$ i dobierz plan optymalny ξ_1^* dla modelu częściowego rozpiętego przez $\mathbf{v}_1(x^{(1)})$. Jako plan bieżący ζ_1 przyjmij ξ_1^* .

Uwaga. Przez „dobór planu” rozumiemy tu i poniżej albo znalezienie planu w bazie planów optymalnych, albo obliczenie numerycznie jego przybliżenia, a w przypadku gdy plan optymalny nie jest jedyny, także wybór wariantu.

Krok 1. Zwiększ licznik $k := k + 1$ i jeśli $k = s + 1$, to zatrzymaj obliczenia, a jako wynik podaj plan ζ_k . W przeciwnym razie przejdź do kroku 2.

Krok 2. Dobierz plan optymalny ξ_k^* dla modelu częściowego rozpiętego przez $\mathbf{v}_1(x^{(k)})$. Oblicz plan $\zeta_k = \xi_k^* \otimes \zeta_{k-1}$. Wróć do kroku 1.

Zastosowanie tego algorytmu do omawianego poprzednio przykładu daje wielowymiarowy plan optymalny przy znikomym nakładzie obliczeń, gdyż plan D-optimalny dla regresji liniowej odczytujemy z bazy planów optymalnych (oczywiście, wcześniej musimy ją stworzyć). W nietypowych przypadkach, gdy plany optymalne dla modeli częściowych trzeba obliczać numerycznie, obliczenia także stają się dużo łatwiejsze niż w przypadku, gdy nie stosujemy Twierdzeń 6.2 lub 6.1. Dzieje się tak, gdyż w każdej iteracji algorytmu Wynna–Fedorova poszukujemy położenia kilku tylko maksimów globalnych i – co najważniejsze – poszukujemy ich na skończonym odcinku. W najprostszej implementacji poszukiwanie to można wykonać przeglądając wartości w równoodległych punktach, odległych

¹ Opis podajemy dla modelu (7.1.1), gdyż przypadek modelu addytywnego lub z interakcjami dla zablokowanych zmiennych rozpatruje się analogicznie.

od siebie o niezbędną dokładność nastawiania wielkości wejściowych w trakcie eksperymentu.

Opisany wyżej algorytm przeznaczony jest do realizacji na komputerze sekwencyjnym. W przypadku posiadania systemu z równoległymi procesorami działanie algorytmu może ulec radykalnemu przyspieszeniu, gdyż wyszukiwanie i/lub obliczanie planów dla modeli częściowych może odbywać się niezależnie, a po zakończeniu ich wyszukiwania lub obliczania następuje skomponowanie planu wielowymiarowego.

Uwagi o zaokrągłaniu planów produktowych

Wagi planu skomponowanego zgodnie z opisem w poprzednim podrozdziale pomnożyć można przez liczbę pomiarów i zaokrąglić do liczb naturalnych, tak jak każdy inny plan (por. s. 28 i następne), ignorując jego iloczynową strukturę. Omówimy tu jednak sposoby uzyskiwania realizowalnych planów, które zachowują strukturę planu produktowego. Jest to ważne, gdyż – jak pokażemy w następnym podrozdziale – prowadzi do znacznych ułatwień obliczeniowych.

Dla uproszczenia wzorów zakładamy, że odpowiedź badanego systemu zależy od $x^{(1)}$ i $x^{(2)}$, a zmienna wejściowa $x^{(1)}$ przyjmuje wartości ze skończonego niepustego zbioru

$$\mathbf{X}_1 = \{x_1^{(1)}, x_2^{(1)}, \dots, x_{K_1}^{(1)}\}.$$

Podobnie $x^{(2)}$ przyjmować może wartości ze zbioru

$$\mathbf{X}_2 = \{x_1^{(2)}, x_2^{(2)}, \dots, x_{K_2}^{(2)}\}.$$

Następujące dalej wzory będą miały sens także wówczas, gdy $x^{(1)}$ i $x^{(2)}$ zastąpimy przez wektory zblokowane zmiennych (wówczas także elementami zbiorów \mathbf{X}_1 i \mathbf{X}_2 są wektory).

Z punktami powyższymi skojarzone są częstości p_{ij} . Zgodnie z Definicją 6.1 częstości te spełniają warunek

$$p_{ij} = p_i^{(1)} p_j^{(2)}, \quad i = 1, 2, \dots, K_1, \quad j = 1, 2, \dots, K_2, \quad (7.1.2)$$

gdzie

$$p_i^{(1)} \geq 0, \quad \sum_{i=1}^{K_1} p_i^{(1)} = 1, \quad p_j^{(2)} \geq 0, \quad \sum_{j=1}^{K_2} p_j^{(2)} = 1. \quad (7.1.3)$$

Przypuśćmy, że powyższy plan, złożony z punktów i wag

$$((x_i^{(1)}, x_j^{(2)}), p_{ij}), \quad i = 1, 2, \dots, K_1, \quad j = 1, 2, \dots, K_2, \quad (7.1.4)$$

jest ciągłym planem optymalnym w sensie wybranego kryterium (lub planem, który z innych powodów decydujemy się realizować) i mamy do zrealizowania $N > 1$ pomiarów, które trzeba rozdysponować proporcjonalnie do p_{ij} .

Jak wspomniano, możemy postąpić następująco: obliczyć $\lfloor N p_{ij} \rfloor$, gdzie $\lfloor \cdot \rfloor$ oznacza największą liczbę całkowitą, nie większą niż liczba ujęta w $\lfloor \cdot \rfloor$. Pozostałe do rozdziału pomiary można rozmieścić losowo w punktach nośnika planu. W ten sposób zignorujemy jednak fakt, że p_{ij} spełniają (7.1.2).

W celu opisanego proponowanego sposobu przypisania liczby pomiarów poszczególnym punktom siatki $\mathbf{X}_1 \times \mathbf{X}_2$, potrzebne nam będą dwa pomocnicze zestawy liczb naturalnych $k_1^{(1)}, k_2^{(1)}, \dots, k_{K_1}^{(1)}$ oraz $k_1^{(2)}, k_2^{(2)}, \dots, k_{K_2}^{(2)}$. Zdefiniujemy liczbę pomiarów (powtórzeń eksperymentu) n_{ij} w punkcie $(x_i^{(1)}, x_j^{(2)})$ następująco:

$$n_{ij} = k_i^{(1)} k_j^{(2)}, \quad i = 1, 2, \dots, K_1, \quad j = 1, 2, \dots, K_2. \quad (7.1.5)$$

Całkowita liczba pomiarów wynosi zatem:

$$n = \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} n_{ij}. \quad (7.1.6)$$

Teraz interpretacja liczb $k_i^{(\cdot)}$ staje się jasna. Na przykład, $k_i^{(1)}$ oznacza liczbę powtórzeń ciągu par $(x_i^{(1)}, x_j^{(2)})$, gdzie $x_j^{(2)}$ przebiega cały zbiór \mathbf{X}_2 , natomiast $x_i^{(1)}$ pozostaje stałe.

Nasze zadanie polega teraz na dobraniu liczb $k_i^{(1)}$ oraz $k_j^{(2)}$, $i = 1, 2, \dots, K_1$, $j = 1, 2, \dots, K_2$ w taki sposób, by liczby n_{ij}/n były możliwie bliskie wagom planu p_{ij} , a liczba n była możliwie bliska N , lecz nie przekraczająca tej wartości. Wymagania te nie są precyzyjnym sformulowaniem zadania, które sformułować można na wiele sposobów, wybierając w różny sposób miary bliskości n_{ij}/n i p_{ij} .

Przykładowe sformułowanie może być następująco: znaleźć minimum względem $k_i^{(1)}$ oraz $k_j^{(2)}$, $i = 1, 2, \dots, K_1$, $j = 1, 2, \dots, K_2$

$$\min_{k_i^{(1)}, k_j^{(2)}} \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} \left(p_i^{(1)} p_j^{(2)} - \frac{k_i^{(1)} k_j^{(2)}}{n} \right)^2 \quad (7.1.7)$$

przy ograniczeniu $n \leq N$, gdzie

$$n = \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} k_i^{(1)} k_j^{(2)}. \quad (7.1.8)$$

W (7.1.7) odległości $\left(p_i^{(1)} p_j^{(2)} - \frac{k_i^{(1)} k_j^{(2)}}{n} \right)^2$ zamienić można na $\left| p_{ij} - \frac{k_i^{(1)} k_j^{(2)}}{n} \right|$, otrzymując inne zadanie optymalizacji. Rozpatrywać można także problem minimum

$$\min_{k_i^{(1)}, k_j^{(2)}} \max_{i,j} \left| p_i^{(1)} p_j^{(2)} - \frac{k_i^{(1)} k_j^{(2)}}{n} \right| \quad (7.1.9)$$

przy tych samych ograniczeniach. Niezależnie od tego, które z wymienionych sformułowań wybierzemy, otrzymamy zadanie optymalizacji dyskretnej. Jedynie dla małych wartości K_1 , K_2 i N do rozwiązywania tego typu problemów stosować można przegląd zupełny wszystkich wariantów. Trudności obliczeniowe znacznie wzrosną, jeśli wzrośnie liczba zmiennych niezależnych (lub bloków, na które zostały one podzielone).

Z tych względów warto rozważyć następujące przybliżone rozwiązanie problemu. Wybierzmy liczby naturalne $n^{(1)}$, $n^{(2)}$ takie, dla których zachodzi $n^{(1)} n^{(2)} \leq N$ i bliskie N . Obliczmy

$$k_i^{(1)} = \lfloor p_i^{(1)} n^{(1)} \rfloor, \quad i = 1, 2, \dots, K_1, \quad (7.1.10)$$

$$k_j^{(2)} = \lfloor p_j^{(2)} n^{(2)} \rfloor, \quad j = 1, 2, \dots, K_2. \quad (7.1.11)$$

Jako liczbę pomiarów n_{ij} przyjmujemy

$$n_{ij} = k_i^{(1)} k_j^{(2)}, \quad i = 1, 2, \dots, K_1, \quad j = 1, 2, \dots, K_2. \quad (7.1.12)$$

Łączna liczba rozdysponowanych pomiarów wyniesie zatem

$$\underline{n} = \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} k_i^{(1)} k_j^{(2)}. \quad (7.1.13)$$

Z konstrukcji wynika, że $\underline{n} \leq N$. Pozostałe pomiary, w liczbie $N - \underline{n}$, można rozdysponować losowo. Możemy teraz rozpatrywać zadanie optymalizacji

$$\min_{n^{(1)}, n^{(2)}} \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} \left(p_i^{(1)} p_j^{(2)} - \frac{k_i^{(1)} k_j^{(2)}}{\underline{n}} \right)^2 + \lambda (N - \underline{n})^2, \quad (7.1.14)$$

gdzie \underline{n} obliczane jest zgodnie z (7.1.13), natomiast $\lambda > 0$ jest wybieranym przez nas współczynnikiem wagowym, który ocenia wpływ spełnienia warunku wykorzystania dopuszczalnej liczby pomiarów, ale bez jej przekraczania. Zadanie to jest znacznie prostsze niż poprzednie, gdyż funkcja celu zależy tylko od dwóch zmiennych całkowitoliczbowych. W przypadku większej liczby zmiennych niezależnych rozwiązanie zadania (7.1.14) odbywać się będzie na siatce d -wymiarowej, podczas gdy w zadaniu (7.1.7) optymalizacja odbywa się w kostce $(\prod_{j=1}^d K_j)$ -wymiarowej.

Kolejną zaletą planów produktowych jest to, że realizując plan

$$n_{ij} = k_i^{(1)} k_j^{(2)}, \quad i = 1, 2, \dots, K_1, \quad j = 1, 2, \dots, K_2. \quad (7.1.15)$$

nie musimy nakładać ograniczeń na kolejność wykonywania pomiarów. Wręcz przeciwnie, istotna jest jedynie ich struktura. Zalecana jest raczej randomizacja niż systematyczne wykonywanie numerowanych obserwacji.

Z drugiej strony w sytuacjach, w których częścią wielkości wejściowych są na przykład parametry materiałowe próbek, produktowa struktura planu może ułatwić jego realizację lub zmniejszyć koszty przeprowadzenia eksperymentu. Systematycznie prowadzić można pomiary wzdłuż powiedzmy kolumn macierzy $\{(x_i^{(1)}, x_j^{(2)})\}$, kojarząc kolumny na przykład z parametrami materiałowymi próbek, a wiersze z poddającymi się łatwiejszej zmianie wielkościami wejściowymi. Taki sposób organizacji pomiarów nazywany jest w [155] eksperymentem wielostopniowym.

7.2. Zyski obliczeniowe w estymacji z zastosowaniem planów produktowych

W podrozdziale tym pokażemy, że realizacja planu uwzględniająca jego produktową strukturę prowadzi do:

- oszczędności obliczeniowych,
- znacznego zmniejszenia trudności obliczeniowych spowodowanych ewentualnym złym uwarunkowaniem układu równań normalnych MNK,
- ułatwień w prowadzeniu eksperymentu.

Algorytm obliczeniowy MNK dostosowany do planów produktowych

Rozważmy raz jeszcze model z pełnym zestawem interakcji względem zablokowanych zmiennych $\mathbf{x}^{(1)}$ i $\mathbf{x}^{(2)}$

$$\bar{y}(x) = a^T [\mathbf{v}_1(\mathbf{x}^{(1)}) \otimes \mathbf{v}_2(\mathbf{x}^{(2)})], \quad (7.2.16)$$

gdzie a jest kolumnowym wektorem nieznanych parametrów,

$$\dim(a) = \dim(\mathbf{v}_1(\mathbf{x}^{(1)})) \dim(\mathbf{v}_2(\mathbf{x}^{(2)})).$$

Modelami częściowymi dla (7.2.16) nazywać będziemy funkcje $\alpha^T \mathbf{v}_1(\mathbf{x}^{(1)})$ oraz $\beta^T \mathbf{v}_2(\mathbf{x}^{(2)})$. Funkcja $\alpha^T \mathbf{v}_1(\mathbf{x}^{(1)})$ to zależność wyjścia od zablokowanych wielkości wejściowych $\mathbf{x}^{(1)}$, przy ustalonych wartościach bloku wejść $\mathbf{x}^{(2)}$. Analogiczną interpretację można nadać drugiemu modelowi częściowemu. W modelach tych $\alpha \in \mathbf{R}^{\dim(\mathbf{v}_1(\mathbf{x}^{(1)}))}$ oraz $\beta \in \mathbf{R}^{\dim(\mathbf{v}_2(\mathbf{x}^{(2)}))}$ są kolumnowymi wektorami stałych parametrów.

Zadanie estymacji parametrów (7.2.16) rozpatrywane będzie na zbiorze punktów $(\mathbf{x}_i^{(1)}, \mathbf{x}_j^{(2)})$, przy czym w każdym z nich liczba powtórzeń pomiarów wynosi

$$n_{ij} = k_i^{(1)} k_j^{(2)}, \quad i = 1, 2, \dots, K_1, \quad j = 1, 2, \dots, K_2. \quad (7.2.17)$$

Niech \bar{y}_{ij} oznacza średnią arytmetyczną pomiarów odpowiedzi badanego obiektu wykonanych w punkcie $(\mathbf{x}_i^{(1)}, \mathbf{x}_j^{(2)})$, $i = 1, 2, \dots, K_1$, $j = 1, 2, \dots, K_2$.

W tych warunkach estymacja wektora a sprowadza się do minimalizacji

$$Q(a) = \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} k_i^{(1)} k_j^{(2)} \left(\bar{y}_{ij} - a^T [\mathbf{v}_1(\mathbf{x}_i^{(1)}) \otimes \mathbf{v}_2(\mathbf{x}_j^{(2)})] \right)^2. \quad (7.2.18)$$

Po obliczeniu gradientu $Q(a)$ względem a i przyrównaniu go do wektora zerowego dostaniemy następujący układ równań normalnych

$$M_{K_1 K_2} \cdot a = \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} k_i^{(1)} k_j^{(2)} \bar{y}_{ij} [\mathbf{v}_1(\mathbf{x}_i^{(1)}) \otimes \mathbf{v}_2(\mathbf{x}_j^{(2)})], \quad (7.2.19)$$

gdzie macierz $M_{K_1 K_2}$ określona jest następująco

$$M_{K_1 K_2} = \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} k_i^{(1)} k_j^{(2)} [\mathbf{v}_1(\mathbf{x}_i^{(1)}) \otimes \mathbf{v}_2(\mathbf{x}_j^{(2)})] \cdot [\mathbf{v}_1(\mathbf{x}_i^{(1)}) \otimes \mathbf{v}_2(\mathbf{x}_j^{(2)})]^T \quad (7.2.20)$$

Rozwiązanie (7.2.19) względem a daje oszacowania \hat{a} wektora nieznanych parametrów. Można je uzyskać przy mniejszym nakładzie obliczeniowym. W tym celu zauważmy, że własności (15.1.3) oraz (15.1.7) pozwalają zapisać (7.2.20) następująco

$$M_{K_1 K_2} = M_{K_1} \otimes M_{K_2}, \quad (7.2.21)$$

gdzie

$$M_{K_1} \stackrel{\text{def}}{=} \sum_{i=1}^{K_1} k_i^{(1)} \mathbf{v}_1(\mathbf{x}_i^{(1)}) \cdot \mathbf{v}_1(\mathbf{x}_i^{(1)})^T, \quad (7.2.22)$$

$$M_{K_2} \stackrel{\text{def}}{=} \sum_{j=1}^{K_2} k_j^{(2)} \mathbf{v}_2(\mathbf{x}_j^{(2)}) \cdot \mathbf{v}_2(\mathbf{x}_j^{(2)})^T. \quad (7.2.23)$$

Zakładając, że macierze M_{K_1} i M_{K_2} są nieosobliwe, stwierdzamy, że

$$M_{K_1 K_2}^{-1} = M_{K_1}^{-1} \otimes M_{K_2}^{-1} \quad (7.2.24)$$

(por. (15.1.9) w Dodatku). Korzystając z tego rezultatu, rozwiązanie układu równań (7.2.19) możemy zapisać w postaci

$$\hat{a} = \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} k_i^{(1)} k_j^{(2)} \bar{y}_{ij} \left[M_{K_1}^{-1} \cdot \mathbf{v}_1(\mathbf{x}_i^{(1)}) \otimes M_{K_2}^{-1} \cdot \mathbf{v}_2(\mathbf{x}_j^{(2)}) \right]. \quad (7.2.25)$$

Wygodnie będzie zdefiniować wektory

$$\alpha_i = M_{K_1}^{-1} \cdot \mathbf{v}_1(\mathbf{x}_i^{(1)}), \quad i = 1, 2, \dots, K_1. \quad (7.2.26)$$

oraz

$$\beta_j = M_{K_2}^{-1} \cdot \mathbf{v}_2(\mathbf{x}_j^{(2)}), \quad j = 1, 2, \dots, K_2. \quad (7.2.27)$$

W tej notacji (7.2.25) przepisać można następująco

$$\hat{a} = \sum_{i=1}^{K_1} \sum_{j=1}^{K_2} k_i^{(1)} k_j^{(2)} \bar{y}_{ij} [\alpha_i \otimes \beta_j]. \quad (7.2.28)$$

Powyższe ćwiczenia definicyjne pozwalają zapisać (7.2.26) i (7.2.27) w postaci, która jest matematycznie równoważna, lecz znacznie wygodniejsza do obliczeń numerycznych

$$M_{K_1} \cdot [\alpha_1, \alpha_2, \dots, \alpha_{K_1}] = [\mathbf{v}_1(\mathbf{x}_1^{(1)}), \mathbf{v}_1(\mathbf{x}_2^{(1)}), \dots, \mathbf{v}_1(\mathbf{x}_{K_1}^{(1)})], \quad (7.2.29)$$

$$M_{K_2} \cdot [\beta_1, \beta_2, \dots, \beta_{K_2}] = [\mathbf{v}_2(\mathbf{x}_1^{(2)}), \mathbf{v}_2(\mathbf{x}_2^{(2)}), \dots, \mathbf{v}_2(\mathbf{x}_{K_2}^{(2)})]. \quad (7.2.30)$$

We wzorach (7.2.29) i (7.2.30) uformowano macierze przez wpisanie ich kolumn w nawiasach [...]. Są to zatem równania macierzowe, których niewiadomymi są macierze $[\alpha_1, \alpha_2, \dots, \alpha_{K_1}]$ oraz $[\beta_1, \beta_2, \dots, \beta_{K_2}]$.

Rozważania te stanowią dowód następującego rezultatu i opartego na jego podstawie algorytmu.

Twierdzenie 7.1. *Załóżmy, że macierze M_{K_1} i M_{K_2} są nieosobliwe. Wówczas nieosobliwa jest również macierz układu równań normalnych (7.2.19). Wektor \hat{a} uzyskany w wyniku rozwiązania układów równań (7.2.29) oraz (7.2.30) i wstawieniu ich rozwiązań do (7.2.28) jest równocześnie rozwiązaniem układu równań normalnych (7.2.19) i tym samym minimalizuje sumę kwadratów błędów (7.2.18).*

Algorytm estymacji MNK dla planów produktowych i modeli z pełnym zestawem interakcji względem zablokowanych zmiennych.

Krok 1. Rozwiązać układ równań (7.2.29) względem $[\alpha_1, \alpha_2, \dots, \alpha_{K_1}]$.

Krok 2. Rozwiązać układ równań (7.2.30) względem $[\beta_1, \beta_2, \dots, \beta_{K_2}]$.

Krok 3. Obliczyć \hat{a} zgodnie ze wzorem (7.2.28).

Stosując metodę Housholdera (por. [71]) do rozwiązania równań (7.2.29) i (7.2.30) nie musimy formować macierzy M_{K_1} i M_{K_2} . Wystarczy obliczyć macierze

$$[\mathbf{v}_1(\mathbf{x}_1^{(1)}), \mathbf{v}_1(\mathbf{x}_2^{(1)}), \dots, \mathbf{v}_1(\mathbf{x}_{K_1}^{(1)})] \quad \text{oraz} \quad [\mathbf{v}_2(\mathbf{x}_1^{(2)}), \mathbf{v}_2(\mathbf{x}_2^{(2)}), \dots, \mathbf{v}_2(\mathbf{x}_{K_2}^{(2)})],$$

które i tak są potrzebne jako prawe strony tychże równań.

Zauważmy, że równania (7.2.29) i (7.2.30) rozwiązywać można niezależnie od siebie, jeśli tylko dysponujemy dwoma procesorami. W powyższych rozwiązaniach zastosować można także rekurencyjną metodę najmniejszych kwadratów (por. [153]).

Uogólnienie powyższych wyników na przypadek większej liczby zablokowanych zmiennych i planów będących produktem większej liczby planów częściowych jest zabiegiem czysto formalnym.

Walory numeryczne

Jak wiadomo (por. [71]), miarą kumulowania się błędów numerycznych, powstających przy rozwiązywaniu układu równań liniowych z symetryczną macierzą A , jest współczynnik

$$\kappa(A) \stackrel{\text{def}}{=} \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)},$$

gdzie $\lambda_{\max}(A)$ i $\lambda_{\min}(A)$ oznaczają, odpowiednio, maksymalną i minimalną wartość macierzy A . Przypomnijmy, że dla macierzy symetrycznych wartości te są liczbami rzeczywistymi.

Role $\kappa(A)$ można intuicyjnie wyjaśnić tak: wartość ta działa jak wzmacniacz błędów. Wiadomo też, że w źle uwarunkowanych układach równań normalnych wartość ta może przekraczać 10^4 .

Porównajmy teraz dwa sposoby rozwiązywania układu równań (7.2.19).

1. Sposób pierwszy polega na zastosowaniu opisanego wyżej algorytmu. Zgodnie z nim, rozwiązujemy dwa układy równań (7.2.29) i (7.2.30), odpowiednio, ze współczynnikami uwarunkowania $\kappa(M_{K_1})$ i $\kappa(M_{K_2})$.
2. Podejście drugie to bezpośrednie rozwiązanie układu równań (7.2.19) za współczynnikiem uwarunkowania $\kappa(M_{K_1 K_2})$.

Ponieważ $M_{K_1 K_2} = M_{K_1} \otimes M_{K_2}$, to (por. Dodatek)

$$\lambda_{\max}(M_{K_1 K_2}) = \lambda_{\max}(M_{K_1}) \lambda_{\max}(M_{K_2}) \quad (7.2.31)$$

oraz

$$\lambda_{\min}(M_{K_1 K_2}) = \lambda_{\min}(M_{K_1}) \lambda_{\min}(M_{K_2}). \quad (7.2.32)$$

Wzory te prowadzą do wniosku, że

$$\kappa(M_{K_1 K_2}) = \kappa(M_{K_1}) \kappa(M_{K_2}). \quad (7.2.33)$$

Ponieważ najczęściej $\kappa(M_{K_1}) > 1$ i $\kappa(M_{K_2}) > 1$, to – z punktu widzenia dokładności obliczeń – używanie proponowanego, tu algorytmu jest wielokrotnie korzystniejsze niż bezpośrednie rozwiązywanie układu równań normalnych (7.2.19).

Teraz pokażemy, że odnosimy także korzyści mierzone liczbą operacji zmienoprzecinkowych. Dla uproszczenia porównań załóżmy, że

$$p \stackrel{\text{def}}{=} \dim(\mathbf{v}_1(\mathbf{x}^{(1)})) = \dim(\mathbf{v}_2(\mathbf{x}^{(2)})),$$

co prowadzi do

$$\dim(a) = \dim(\mathbf{v}_1(\mathbf{x}^{(1)})) \dim(\mathbf{v}_2(\mathbf{x}^{(2)})) = p^2.$$

Przyjmijmy też, że zarówno bezpośrednie rozwiązanie układu równań normalnych (7.2.19), jak i rozwiązania układów równań (7.2.29) oraz (7.2.30) uzyskujemy klasyczną metodą eliminacji Gaussa o złożoności obliczeniowej proporcjonalnej do trzeciej potęgi wymiaru rozwiązywanego układu równań. Przyjmijmy też, że nakłady obliczeniowe potrzebne do obliczenia

$$[\mathbf{v}_1(\mathbf{x}_1^{(1)}), \mathbf{v}_1(\mathbf{x}_2^{(1)}), \dots, \mathbf{v}_1(\mathbf{x}_{K_1}^{(1)})], [\mathbf{v}_2(\mathbf{x}_1^{(2)}), \mathbf{v}_2(\mathbf{x}_2^{(2)}), \dots, \mathbf{v}_2(\mathbf{x}_{K_2}^{(2)})] \quad (7.2.34)$$

zostały już poniesione i zauważmy, że ten sam nakład obliczeniowy trzeba ponieść w obu podejściach.

N	Liczba parametrów p^2			
	9	16	25	36
100	10	18	30	47
600	10	17	26	37
1100	10	17	26	37

Tabela 7.1. Iloraz nakładów obliczeniowych (7.2.35) klasycznej MNK i proponowanego algorytmu, wykorzystującego strukturę planu i modelu

Licząc także wstawienie (7.2.34) do (7.2.28), otrzymamy następujące wyrażenie na iloraz nakładu obliczeniowego podejścia bezpośredniego do liczby operacji zmiennoprzecinkowych proponowanego algorytmu.

$$\zeta(N, p) = \frac{p^6 + N p^4}{p^3 + N p^2}, \quad (7.2.35)$$

gdzie $N = K_1 K_2$. Zauważmy, że

$$\lim_{N \rightarrow \infty} \zeta(N, p) = p^2$$

co oznacza, że dla dużej liczby pomiarów uzyskujemy zmniejszenie nakładów obliczeniowych, tym większe, im większa jest liczba estymowanych parametrów modelu. Co więcej, oszczędności obliczeniowe dla małych i średnich wartości N są nawet większe niż wskazywane przez powyższe wyrażenie asymptotyczne. Stwierdzenie to dokumentuje tabela 7.1. Jej elementy wskazują na krotność redukcji liczby operacji zmiennoprzecinkowych dla poszczególnych liczb pomiarów N

i rozmiarów modelu. Tabela ta wskazuje, że nawet w przypadku niedużej liczby pomiarów i prostych modeli zysk obliczeniowy jest dziesięciokrotny.

Nietrudno zauważyć, że w przypadku modeli o większej liczbie zmiennych (lub bloków zmiennych) z pełnym zestawem interakcji zyski obliczeniowe będą znacząco większe, jeśli tylko zastosować możemy plany produktowe. Jak pamiętamy z lektury poprzednich rozdziałów, plany o tej strukturze są często planami optymalnymi. Warto jednak odnotować, że powyższe oszczędności nakładów obliczeniowych uzyskamy także wówczas, gdy plan produktowy nie jest optymalny.

CZĘŚĆ III

Eksperyment w testowaniu jakości wyrobów

8. Diagnozowanie i poprawa odporności wyrobów

Rolę planowania eksperymentu we wstępnych etapach przygotowania wyrobu do produkcji oraz w aktywnym diagnozowaniu nadmiernej zmienności parametrów charakteryzujących poszczególne egzemplarze prześledzić można w wielu monografiach i artykułach [186], [187], [15], [206], [185], [144], [80], [24], [188].

8.1. Bezpośredni model odporności wyrobów na warunki eksploatacji

W rozdziale tym przedstawiono pewne aspekty racjonalizacji badań eksperymentalnych w celu oceny odporności wyrobu na zmiany warunków eksploatacji. Zaproponowano też nowy wskaźnik, który może służyć do oceny odporności produktu eksploatowanego w różnych warunkach. W odróżnieniu od podejścia Taguchi (por. rozdz. 8.2), które bazuje na eksperymentalnym poszukiwaniu parametrów wyrobu odpornych na warunki eksploatacji, tutaj proponujemy bezpośrednią budowę modelu dla zależności między odległością od zadanego celu a warunkami eksploatacji. Model ten budujemy dla istniejącego już wyrobu lub procesu i celem jest jedynie ocena jakości. Jeśli wyrób (proces) okaże się zbyt wrażliwy na warunki eksploatacji, to należy wrócić do fazy projektowej i wówczas posłużyć się można metodyką Taguchi.

Przypuśćmy, że testowaniu poddano pewien proces produkcyjny, którego wynikiem jest produkt scharakteryzowany (dla uproszczenia wzorów) jedną wielkością liczbową. Standardowych przykładów dostarczają procesy walcowania blach stalowych lub wyciągania drutu miedzianego. W procesach tych podstawowym parametrem wyjściowym jest grubość blachy lub średnica drutu, a utrzymanie zadanej wartości tego parametru jest głównym zadaniem technologa procesu. Tę zadaną wartość nazywać będziemy celem prowadzenia procesu (lub krótko – celem procesu). Podstawową ideę zilustrujemy na przykładzie badania odporności wartości wyjściowej od jednej zmiennej wejściowej procesu x , którą w wymienionych przykładach może być siła nacisku walców.

Przyjmijmy następujący prosty model zależności średniej odległości od celu $E(Y(x))$ od x :

$$E(Y(x)) = a_0 + a_1 x + a_2 x^2, \quad (8.1.1)$$

gdzie zmienna standaryzowana x przyjmuje wartość $x = -1$ w przypadku minimalnej dopuszczalnej wartości wejścia (w naszym przykładzie minimalny nacisk

walców, dający jeszcze produkt o wartości wyjścia zbliżonej do celu). Przyjmujemy $x = 1$, gdy proces prowadzony jest przy maksymalnej wartości wejścia. Wartość $x = 0$ odpowiada przeciętnym warunkom prowadzenia procesu. Naszym pierwszym zadaniem jest znalezienie takiej wartości x , dla której średnia odległość $E(Y(x))$ od wartości zadanej jest najmniejsza.

Model (8.1.1) jest najprostszą funkcją, która pozwala modelować charakterystykę posiadającą ekstremum, występujące dla $x^* = -a_1/(2a_2)$ i mające wartość $E(Y(x^*)) = a_0 - a_1^2/(4a_2)$. Wartości tej nie można obliczyć, gdyż nie są znane parametry a_0, a_1, a_2 . W celu ich oszacowania przeprowadzamy eksperyment A-optimalny o postaci pokazanej na rysunku 8.1. Następnie obliczamy oszacowania parametrów a_0, a_1, a_2 stosując klasyczną metodę najmniejszych kwadratów. Oznaczmy tak obliczone wartości przez $\hat{a}_0, \hat{a}_1, \hat{a}_2$ i wstawmy je do modelu (8.1.1), co daje:

$$\hat{Y}(x) = \hat{a}_0 + \hat{a}_1 x + \hat{a}_2 x^2, \quad (8.1.2)$$

gdzie $\hat{Y}(x)$ jest oszacowaniem średniej odległości od celu przy zastosowaniu wejścia x .

Oszacujmy teraz najkorzystniejszą wartość wejścia dla testowanego procesu: $\hat{x} = -\hat{a}_1/(2\hat{a}_2)$ oraz spodziewaną najmniejszą średnią odległość od celu: $\hat{Y}(\hat{x}) = \hat{a}_0 - \hat{a}_1^2/(4\hat{a}_2)$.

Mając oszacowane te wartości, wprowadzić możemy wskaźnik odporności, oznaczany dalej przez κ , który wskazuje względny wzrost średniej odległości od celu w niekorzystnych warunkach, odniesiony do najlepszej możliwej do uzyskania średniej odległości od celu badanego procesu. Dokładniej, teoretyczna wartość

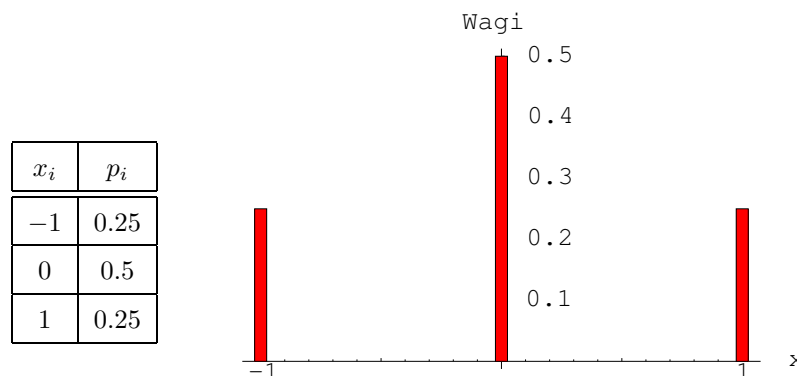
$$\kappa = \frac{0.5(E(Y(-1)) + E(Y(1)))}{E(Y(x^*))}, \quad (8.1.3)$$

natomiast oszacowanie wskaźnika κ , uzyskane na podstawie empirycznych oszacowań parametrów modelu, wynosi:

$$\hat{\kappa} = \frac{\hat{a}_0 + \hat{a}_2}{\hat{a}_0 - \hat{a}_1^2/(4\hat{a}_2)}. \quad (8.1.4)$$

Zauważmy, że z konstrukcji κ wynika, że zawsze $\kappa \geq 1$, przy czym wartość $\kappa = 1$ wystąpi jedynie wtedy, gdy średnia odległość od celu dla badanego procesu nie pogarsza się w skrajnych warunkach, odpowiadających $x = \pm 1$.

Wzrost wartości κ powyżej pewnej, wybieranej przez użytkownika, wartości wskazuje na nadmierną utratę średniej odległości od celu w skrajnych warunkach. Jeśli na przykład dopuścimy 50% pogorszenie średniej odległości od celu przy próbach w skrajnych warunkach, to $\kappa > 1.5$ wskazuje, że testowany proces jest nadmiernie wrażliwy na pracę w niekorzystnych sytuacjach. W praktyce musimy oczywiście korzystać z warunku na przewyższanie zadanego poziomu przez $\hat{\kappa}$ (w naszym przykładzie jest to warunek $\hat{\kappa} > 1.5$).



Rys. 8.1. Plan eksperymentu dla testowania średniej odległości od wartości zadanej w zależności od jednego czynnika (np. siły nacisku walców)

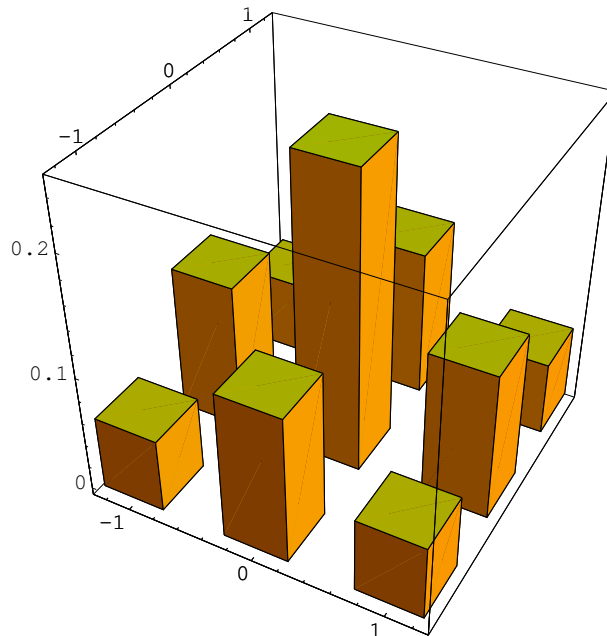
Powyższe rozważania można uogólnić na przypadek dwóch i większej liczby czynników wpływających na średnią odległość od celu, korzystając z podejścia do konstruowania planów produktowych, które opisano w rozdziale 6.

Plan A-optimalny dla pełnego modelu będącego funkcją kwadratową dwóch zmiennych (np. nacisku walców i temperatury walcowanego materiału) pokazano na rysunku 8.2. W tym przypadku wzory na κ i $\hat{\kappa}$ należy odpowiednio zmodyfikować.

8.2. Plany produktowe w modelu Taguchi

W podrozdziale tym przedstawiono krótko podejście Taguchi, używane w empirycznej optymalizacji procesów wytwórczych. Celem tego podejścia jest uzyskanie pomiarów do takiego projektowania wyrobów, by utrzymać cechy produktu na zadanym poziomie, minimalizując jednocześnie ich zmienność wokół wartości zadanej. Następnie proponujemy alternatywne podejście, w którym planowany eksperyment używany jest do budowy modelu matematycznego. Służy on zebraniu pomiarów do estymacji parametrów modelu. Następnie, model ten wykorzystać można w tych samych celach, które nakreślił Taguchi. Z dokładnością do błędów estymacji uzyskuje się zatem taki sam efekt finalny jak w podejściu Taguchi. W planowaniu wspomnianego wyżej eksperymentu ponownie przydatne okazują się plany produktowe. Plany, które oryginalnie proponował Taguchi, mają wiele cech wspólnych z planami produktowymi. Plany Taguchi były wielokrotnie krytykowane w literaturze (por. dyskusje w [98], [89], [209]) za to, że z góry zakłada się w nich strukturę planu, który składa się z dwóch tablic i dla każdego eksperymentu z jednej tablicy wykonać trzeba wszystkie eksperymenty z drugiej tablicy. Struktura ta jest podobna do konstrukcji planów produktowych.

$x_i^{(1)}$	$x_i^{(2)}$	p_i
-1	-1	1/16
+1	-1	1/16
-1	+1	1/16
+1	+1	1/16
0	0	1/4
+1	0	1/8
0	+1	1/8
-1	0	1/8
0	-1	1/8



Rys. 8.2. Plan eksperymentu do testowania średniej odległości od wartości zadanej wyjścia w zależności od dwóch czynników (np. siły nacisku walców i temperatury walcowanego materiału)

Postaramy się pokazać, że dla modeli z pełnym zestawem interakcji ograniczenie, które nakłada Taguchi, nie musi oznaczać stosowania planów nieoptymalnych. Trzeba jednak zastrzec, że omawiane tu plany produktowe są planami ciągłymi (aproxymacyjnymi) w sensie omawianym w części I i II tej książki, podczas gdy Taguchi propagował użycie ortogonalnych planów czynnikowych. Trzeba też podkreślić słuszność uwag zawartych w [89] na temat braku możliwości wpływania na wariancję wyjścia w przypadku, gdy przyjmuje się model addytywny względem grup zmiennych sterujących i zakłócających. Dlatego też rozpatrujemy model z pełnym zestawem interakcji. Opis oryginalnego podejścia Taguchi znaleźć można w [186], [187], [185].

Rozważmy następujący opis pewnego wyrobu:

$$y = H(u, z) + \epsilon, \quad (8.2.5)$$

gdzie wielkości występujące w (8.2.5) określone są następująco:

1. y oznacza jednowymiarowe wyjście badanego procesu.

2. ϵ to zmienna losowa reprezentująca niemierzalną część zakłóceń, które nie mogą być kontrolowane nawet w warunkach laboratoryjnych. Będziemy zakładać, że $\text{var}(\epsilon) = \sigma^2 < \infty$, gdzie wariancja σ^2 nie zależy od zmiennych x, u, z .
3. Przez z oznaczamy wektor wielkości wejściowych, których wartości można wybierać, mierzyć i kontrolować w warunkach laboratoryjnych. Natomiast w procesie normalnej eksploatacji wyrobu wartości te zmieniają się w sposób losowy i wpływają na wariancję zmiennej losowej y . Zwykle też wartości z nie są mierzone w trakcie eksploatacji wyrobu.
4. Przez u oznaczamy wektor czynników (wejść) o charakterze parametrów projektowych danego wyrobu. W trakcie prób laboratoryjnych na prototypach ich wartości mogą podlegać świadomym zmianom – są zatem czynnikami eksperymentu. W trakcie eksploatacji wartości tych nie zmienia się lub można je dostrajać w niewielkim zakresie.
5. Funkcja H zależy od obu rodzajów wymuszeń u i z , natomiast jej postać funkcyjna jest znana lub założona przez eksperymentatora jako aproksymacja pewnej nieznannej funkcji. Obie te funkcje zależą od nieznanymi współczynników.

To ostatnie założenie jawnie nie występuje w podejściu Taguchi, gdyż stara się on opierać swoją metodykę wyłącznie na wynikach eksperymentów. Założenie o postaci modelu wprowadzamy po to, by pokazać, że proponowane przez Taguchi plany mają cechy planów optymalnych, jeśli H ma postać modelu z pełnym zestawem interakcji między u i z . Ponadto wydaje się oczywiste, że mając wyniki eksperymentu zawsze warto podjąć próbę zbudowania modelu, gdyż może on pozwolić na bardziej precyzyjny dobór interesujących nas wartości. Założenia o tym, że przed eksperymentem jesteśmy w stanie wydzielić wektory u i z jest jednym z najczęściej krytykowanych w literaturze na temat podejścia Taguchi (por. [86]).

Zgodnie z metodologią Taguchi

celem eksperymentów jest zapewnienie, by wartość oczekiwana y osiągała zadany poziom y^ przy jednoczesnym dążeniu do minimalizacji rozrzutu y -ków y^* . Za miarę rozrzutu Taguchi proponuje przyjmować wariancję y .*

Przyjmijmy, że dla ustalonego z zależność $H(u, z)$ od wektora u przybliżyć można następująco:

$$H(u, z) = \alpha^T v_1(u) = v_1^T(u) \alpha, \quad (8.2.6)$$

gdzie $v_1(u)$ jest $r_1 \times 1$ wektorem wybranych przez nas liniowo niezależnych funkcji zmiennych projektowych. Wektor α nie jest znany i dla każdego ustalonego z może być estymowany na podstawie wyników eksperymentu wykonanych w punktach następującego zbioru

$$\mathbf{U} = \{u_j, j = 1, 2, \dots, K_1\}. \quad (8.2.7)$$

W ogólnym przypadku zarówno położenie punktów u_j , jak i ich liczba K_1 mogą być różne dla różnych z . Zgodnie ze wspomnianym wcześniej podejściem Taguchi, przyjmujemy jednak, że ten sam zestaw (8.2.7) stosowany jest dla wszystkich wartości z , dla których estymowany będzie wektor α . Oznacza to, że dla każdego elementu z zestawu

$$\mathbf{Z} = \{z_i, i = 1, 2, \dots, K_2\} \quad (8.2.8)$$

realizowany jest cały zestaw eksperymentów \mathbf{U} . Jeśli teraz zastosowalibyśmy metodę najmniejszych kwadratów do znalezienia oszacowań dla wektora α , to dla każdego $z_i \in \mathbf{Z}$ otrzymamy inny wektor oszacowań $\hat{\alpha}(z_i)$, $i = 1, 2, \dots, K_2$. W ten sposób każda składowa wektora α staje się funkcją wektora z , a zestaw $\hat{\alpha}(z_i)$, $i = 1, 2, \dots, K_2$ można traktować jako obserwacje wektora funkcji $\alpha(z)$.

Oznaczmy przez $\alpha^{(k)}(z)$, $k = 1, 2, \dots, r_1$ elementy wektora $\alpha(z)$. Wówczas (8.2.6) przyjmie postać

$$H(u, z) = \alpha^T(z) v_1(u). \quad (8.2.9)$$

Na podstawie $\hat{\alpha}(z_i)$, $i = 1, 2, \dots, K_2$ można podjąć próbę estymacji zależności $\alpha^{(k)}(z)$ od z . Każda z tych zależności może być aproksymowana za pomocą innego szeregu funkcji. W celu uproszczenia dalszych rozważań przyjmujemy jednak, że wszystkie funkcje $\alpha^{(k)}(z)$, $k = 1, 2, \dots, r_1$ można dostatecznie dokładnie aproksymować w tej samej, r_2 -wymiarowej, podprzestrzeni funkcji, która rozpięta jest przez wektor funkcji $v_2(z)$. Różne są natomiast wektory współczynników każdej z funkcji $\alpha^{(k)}(z)$. Wektory te oznaczamy będziemy przez $\beta^{(k)}$, $k = 1, 2, \dots, r_1$. Podsumowując, przyjmujemy, że dla $k = 1, 2, \dots, r_1$

$$\alpha^{(k)}(z) = (\beta^{(k)})^T v_2(z) \quad \text{lub} \quad \alpha(z) = B^T \cdot v_2(z), \quad (8.2.10)$$

gdzie macierz B o wymiarach $r_2 \times r_1$ utworzona jest z zestawienia kolumnowych wektorów $\beta^{(k)}$, $k = 1, 2, \dots, r_1$. Podstawiając w (8.2.9) prawą stronę (8.2.10) i biorąc pod uwagę, że $\alpha^T(z) = v_2^T(z) \cdot B$, otrzymamy

$$H(u, z) = v_2^T(z) B v_1(u) = \text{tr} \left[v_2^T(z) B v_1(u) \right] = \text{tr} \left[B v_1(u) v_2^T(z) \right]. \quad (8.2.11)$$

W powyższych przekształceniach wykorzystaliśmy to, że wyrażenie $v_2^T(z) B v_1(u)$ jest wielkością skalarną oraz następującą własność śladu macierzy: $\text{tr}[AB] = \text{tr}[BA]$, jeżeli wykonalne są mnożenia BA i AB .

Zauważmy jeszcze, że – z dokładnością co do numeracji – zbiór elementów macierzy $v_1(u) v_2^T(z)$ i zbiór składowych wektora $v_1(u) \otimes v_2(z)$ są identyczne. Zachowując odpowiedniość numeracji, możemy zatem tak uszeregować elementy macierzy B , by tworzyły one $r_1 r_2$ wymiarowy wektor kolumnowy, oznaczany dalej przez b , dla którego zachodzi

$$H(u, z) = b^T v_1(u) \otimes v_2(z). \quad (8.2.12)$$

Pokazaliśmy więc, że model H ma następującą własność:

Własność 9. *Jeśli spełnione są założenia (8.2.9) i (8.2.10), to funkcję $H(u, z)$ przedstawić można w postaci (8.2.12), czyli w postaci modelu z pełnym zestawem interakcji.*

W celach ilustracyjnych rozważmy następującą klasę modeli

$$y = b^T v_1(u) \otimes \begin{bmatrix} 1 \\ z \end{bmatrix} + \epsilon. \quad (8.2.13)$$

Na etapie eksploatacji wyrobu z traktujemy jako zmienną losową, o której zakładamy, że $E(z) = 0$, $E(z^2) = \sigma_z^2 < \infty$, a zakłócenia ϵ , o zerowej wartości oczekiwanej i skończonej wariancji σ_ϵ^2 są nieskorelowane z wymuszeniem z , a więc $E(z\epsilon) = 0$.

Wygodnie będzie zapisać wektor b w postaci $b^T = [b_1^T, b_2^T]^T$, gdzie b_1 i b_2 są kolumnowymi wektorami o takiej samej liczbie elementów, jaką ma wektor $v_1(u)$. Możemy wtedy przedstawić (8.2.13) w postaci

$$y = b_1^T v_1(u) + b_2^T v_1(u) z + \epsilon. \quad (8.2.14)$$

Łatwo sprawdzić, że

$$E(y) = b_1^T v_1(u), \quad \text{var}(y) = \left(b_2^T v_1(u) \right)^2 \sigma_z^2 + \sigma_\epsilon^2. \quad (8.2.15)$$

Ponadto $\text{var}(y) \geq \sigma_\epsilon^2$, przy czym równość osiągana jest, gdy $b_2^T v_1(u) = 0$. Jeśli zatem istnieje takie u^* , dla którego spełniony jest układ równań

$$b_1^T v_1(u^*) = y^*, \quad b_2^T v_1(u^*) = 0, \quad (8.2.16)$$

to spełnić można wymagania Taguchi – takiego doboru parametrów projektowych wyrobu, by $E(y) = y^*$ i $\text{var}(y)$ osiągała minimalną wartość. Zauważmy, że wymagań tych spełnić nie można, jeśli nie dysponujemy co najmniej dwoma parametrami projektowymi. W najprostszym przypadku, gdy $v_1(u) = [u^{(1)}, u^{(2)}]^T$, (8.2.16) jest układem równań liniowych. Jeśli $\dim(u) > 2$, to może istnieć nieskończenie wiele zestawów u , które spełniają (8.2.16). Można wówczas nałożyć dodatkowe wymagania na rozkład prawdopodobieństwa y -ków, na przykład, wymaganie jego symetrii względem y^* . Aby zrealizować ten przepis, musimy znać wartości wektorów b_1 i b_2 lub oszacować je na podstawie wyników eksperymentu.

Wróćmy do ogólniejszego zadania estymacji parametrów b modelu (8.2.12) na podstawie obserwacji (8.2.5), czyli

$$y = b^T v_1(u) \otimes v_2(z) + \epsilon. \quad (8.2.17)$$

Zgodnie z podejściem Taguchi, teraz – w warunkach laboratoryjnych – dopuszczamy nie tylko możliwość doboru wartości wymuszeń z , ale także możliwość

prowadzenia eksperymentów w następujący sposób: dla poszczególnych $z_i \in \mathbf{Z}$, $i = 1, 2, \dots, K_2$ podać należy wymuszenia

$$(u_j, z_i), \quad j = 1, 2, \dots, K_1. \quad (8.2.18)$$

Łatwo zauważyć, że (8.2.18) ma strukturę planu produktowego (opisaną w rozdziale 7), a (8.2.17) jest modelem z pełnym zestawem interakcji (por. rozdz. 6.3). O optymalności planu (8.2.18) nie potrafimy się wypowiedzieć. Możemy jednak potraktować (8.2.18) jako nośnik pewnego planu produktowego z klasy $\Xi(U \times Z)$, gdzie U i Z są zbiorami dopuszczalnych wartości, odpowiednio zmiennych u i z . Interpretując w taki szerszy sposób plany o strukturze proponowanej przez Taguchi, z Twierdzenia 6.3 otrzymamy.

Wniosek 8.1. *Jeśli spełnione są założenia Twierdzenia 6.3 i zrealizowana zostanie opisana tam procedura znajdowania Φ -optymalnych planów dla modeli częściowych w (8.2.17), to nośnik planu Φ -optymalnego będzie miał strukturę taką, jaką proponował Taguchi (por. (8.2.18)).*

W swoich postulatach Taguchi nakładał dodatkowe warunki, mianowicie, wymagał od planów dla zmiennych u i z , by były one planami ortogonalnymi. Zawężając odpowiednio Wniosek 8.1, można podać warunki dostateczne optymalności takiego postępowania.

Wniosek 8.2. *Założmy, że modele częściowe w (8.2.17) są liniowe, tzn. $v_1(u) = [1, u^T]^T$ oraz $v_2(z) = [1, z^T]^T$, gdzie u i z są wektorami, które mogą przyjmować wartości, odpowiednio, z kostek $[-1, 1]^{\dim(u)}$ i $[-1, 1]^{\dim(z)}$. Wówczas plan D -optymalny jest planem produktowym złożonym z planów ortogonalnych, będących pełnymi planami czynnikowymi na dwóch poziomach z jednakowymi wagami. Ponadto plan ten ma nośnik o postaci (8.2.18), przy czym wszystkie u_j oraz z_i przyjmują wartości ± 1 .*

9. Plany o minimalnym koszcie i zadanej jakości

Klasyczne plany ortogonalne (por. [62]) i plany o symetrii obrotowej (zwane też czasem planami rotatabilnymi) są nadal najczęściej wykorzystywane w praktycznych zastosowaniach, w tym również w zagadnieniach poprawy jakości wyrobów w fazie ich projektowania i produkcji (por. [89]). Zagadnienia uwzględniania kosztów takich planów są zwykle rozważane w tle, jak gdyby poza teorią, i najczęściej sprowadzają się do dążenia do uzyskania planu o pożądanych własnościach przy małej liczbie eksperymentów. Dążenia te tłumaczą popularność ortogonalnych planów ułamkowych.

Motywacji do rozważania kosztów eksperymentów dostarczają te procesy produkcyjne, w których dla oceny jakości produktu lub półproduktów ten sam eksperyment przeprowadzany jest na każdym produkowanym egzemplarzu. W produkcji masowej, nawet jeśli pojedynczy eksperyment jest prosty i względnie tani, to obniżka jego kosztów prowadzi do znacznych oszczędności.

W podrozdziale tym opisujemy takie postawienie problemu planowania, w którym dążymy do minimalizacji kosztów eksperymentu, zachowując jednocześnie cechy jakościowe planu, a dokładniej, generowanej przez ten plan macierzy informacyjnej. Korzystać będziemy z wyników zawartych w pracy autora [129].

Podobnie jak w poprzednich częściach tej książki, rozważać będziemy liniową względem parametrów $a \in R^r$ funkcję regresji rozpiętą przez liniowo niezależne funkcje $v_k(x)$, $k = 1, 2, \dots, r$, określone i ciągłe na zwartym (domkniętym oraz ograniczonym) zbiorze $X \subset R^s$. Z funkcji tych formujemy wektor $v(x) = [v_1(x), v_2(x), \dots, v_r(x)]^T$.

Odpowiedź na zestaw wejść x oznaczamy przez $Y(x)$ i zakładamy, że jest to zmienna losowa taka, że

- dla pewnego nieznanego $a \in R^r$ zachodzi $E(Y(x)) = a^T v(x)$,
- $\text{var}(Y(x)) = \sigma^2 < \infty$.

Jak wiemy, dokładność estymatora \hat{a} parametrów zależy od planu eksperymentu. W literaturze wyróżnić można trzy następujące główne nurty prac na temat planowania eksperymentów.

Nurt klasyczny. Celem planowania jest taki dobór eksperymentu, aby uzyskać pożądane własności macierzy kowariancji $\text{cov}(\hat{a})$, takie jak

- a) diagonalność tej macierzy,
- b) zależność funkcji $v^T(x) \text{cov}(\hat{a}) v(x)$ jedynie od $\|x\|$.

Planów zapewniających diagonalność macierzy $\text{cov}(\hat{a})$ szuka się zwykle wśród planów czynnikowych, czyli takich, w których na poszczególnych składowych

wektora x umieścić można wartości tylko ze skończonego (i z góry zadanego) zbioru. Jak wspomniano na wstępie tego rozdziału, koszt takich eksperymentów uwzględnia się jedynie pośrednio, poprzez dążenie do znalezienia planu, który zapewnia diagonalność $\text{cov}(\hat{a})$ przy nie nazbyt rozrzutnym gospodarowaniu liczbą eksperymentów.

Plany mające własność b), to plany o symetrii obrotowej, które zapewniają jednakową dokładność estymacji funkcji regresji $a^T v(x)$ za pomocą funkcji $\hat{a}^T v(x)$, w jednakowej odległości od centrum układu współrzędnych.

Plany optymalne. Planom tej grupy poświęcona jest większość rozdziałów tej książki. Jak wiemy, plany te są normalizowane tak, aby liczba eksperymentów, a zatem i zależny od niej koszt, jawnie nie występowała w zadaniu optymalizacji. O koszty tych planów zatroszczyć się możemy dopiero na etapie dyskretyzacji planu optymalnego, czyli wówczas, gdy dokonujemy zaokrąglenia do liczb całkowitych wartości $N p_i^*$, $i = 1, 2, \dots, m$.

Minimalizacja liczby pomiarów dla danej dokładności estymacji. W tej klasie problemów planowania jawnie stawia się zadanie minimalizacji liczby eksperymentów N przy ograniczeniu, że pewien funkcjonal macierzy $\text{cov}(\hat{a})$, mierzący dokładność estymacji \hat{a} , osiągnie zadaną z góry wartość. Za miarę jakości przyjmuje się zwykle te same funkcjonały, które rozważaliśmy w tej książce, a więc, np. wyznacznik czy ślad macierzy $\text{cov}(\hat{a})$. Jedną z pierwszych prac tego nurtu był rozdział w monografii [35].

9.1. Minimalizacja kosztu przy zadanej macierzy informacyjnej

Rozważany w tym rozdziale problem ma cechy zarówno pierwszej, jak i trzeciej z wymienionych klas zadań. Założmy, że wybraliśmy pewną symetryczną i dodatnio określoną macierz D o r wierszach i kolumnach. Macierz ta pełnić będzie dalej rolę macierzy informacyjnej, którą chcielibyśmy osiągnąć poprzez odpowiedni dobór planu eksperymentu, o ile przy zadanej strukturze funkcji regresji i zadanym zbiorze dopuszczalnych wymuszeń X jest to możliwe. Jedną z podstawowych własności macierzy informacyjnej jest jej symetria i nieujemna określoność. Dlatego też macierz D musi je spełniać. Dodatkowo założymy, że macierz ta jest nieosobliwa, aby zapewnić estymowalność wszystkich parametrów funkcji regresji. Innymi słowy, zadajemy macierz kowariancji D^{-1} oszacowań parametrów modelu i chcemy ją osiągnąć przez dobór planu. Jednocześnie chcielibyśmy minimalizować całkowity koszt eksperymentu, na który składają się koszty zastosowania poszczególnych zestawów wejść. Będziemy zakładać, że dana jest ciągła w X i nieujemna funkcja $w(x)$, która określa koszt pojedynczego eksperymentu wykonanego w punkcie $x \in X$.

Kolejnym elementem potrzebnym do sformułowania zadania jest klasa eksperymentów, które chcemy brać pod uwagę. W poprzednich rozdziałach ogra-

niczaliśmy się zwykle do klasy eksperymentów Ξ , będących dyskretnymi miarami probabilistycznymi określonymi na X . Tutaj będziemy potrzebować jeszcze szerszego pojęcia planu, po to, by uzyskać możliwie „duży” zbiór osiągalnych macierzy informacyjnych. W rozdziale tym jako plan eksperymentu traktować będziemy dowolną miarę skończoną, określoną na zbiorze $X \subset R^s$ dopuszczalnych zestawów wejść. Plan taki oznaczać będziemy przez F , a klasę wszystkich takich planów określonych na X oznaczymy przez \mathcal{F} . Nie wskazujemy jawnie zależności klasy planów \mathcal{F} od X , aby nie komplikować dalszych oznaczeń. Warto zwrócić uwagę na to, czym różnią się plany z klasy \mathcal{F} od tych, które rozważaliśmy dotychczas. W istocie zrezygnowaliśmy z warunku sumowania się wag planu do 1. W szczególności klasa \mathcal{F} zawiera plany o postaci

$$F_N = \begin{bmatrix} x_1, & x_2, & \dots, & x_m \\ N p_1, & N p_2, & \dots, & N p_m \end{bmatrix}, \quad (9.1.1)$$

gdzie x_i oraz p_i są takie jak w planach z klasy $\Xi(X)$, natomiast $N > 0$ pełni rolę liczby eksperymentów, przy czym w rozdziale tym nie wymagamy, by N było liczbą naturalną.

Nie muszą być liczbami naturalnymi także iloczyny $N p_j$, $j = 1, 2, \dots, m$. Dopiero na etapie realizacji będziemy zaokrąglać wartości $N p_j$, $j = 1, 2, \dots, m$ do najbliższych liczb naturalnych.

Dla podkreślenia różnic między nieunormowanymi planami z klasy \mathcal{F} a unormowanymi planami używanymi dotychczas w całym tym rozdziale stosować będziemy inny niż dotąd zapis całki Lebesgue'a, a mianowicie $\int_X w(x) dF(x)$. Podkreślamy, że rozumienie pojęcia całki pozostaje bez zmiany, tyle że liczymy ją względem innej klasy miar.

Sformułowanie problemu

Po tych objaśnieniach możemy zdefiniować koszt $I(F)$ eksperymentu $F \in \mathcal{F}$

$$I(F) = \int_X w(x) dF(x), \quad (9.1.2)$$

gdzie całka rozumiana jest w sensie Lebesgue'a. Dla planów o postaci (9.1.1) koszt ten obliczyć możemy następująco

$$I(F_N) = \sum_{j=1}^m N p_j w(x_j). \quad (9.1.3)$$

Jako koszt eksperymentu $w(x)$ możemy wybrać funkcję charakterystyczną zbioru X , to znaczy

$$w(x) = \begin{cases} 1, & x \in X, \\ 0, & x \notin X. \end{cases}$$

W tym przypadku, jeśli plan jest postaci (9.1.1), to całkowity koszt równy jest liczbie pomiarów $I(F_N) = N$. Innym przykładem wyboru funkcji kosztu jest $w(x) = \sum_{l=1}^s |x^{(l)}|$, którą możemy interpretować jako sumę amplitud wymuszeń zastosowanych na poszczególnych wejściach $x^{(l)}$, $l = 1, 2, \dots, s$.

Kolejnym elementem potrzebnym do sformułowania zadania planowania eksperymentu są ograniczenia. Zdefiniujmy podzbiór $\mathcal{F}(D)$ zbioru wszystkich planów \mathcal{F} poprzez nałożenie wymagania, by plan $F \in \mathcal{F}$ był elementem $\mathcal{F}(D)$ tylko wówczas, gdy zapewnia on spełnienie równości $D = \int_X v(x) \cdot v^T(x) dF(x)$. Badanie czy zbiór $\mathcal{F}(D)$ zawiera choćby jeden element jest dość trudne. Pewne komentarze na ten temat podamy w dalszej części tego rozdziału.

Możemy teraz sformułować problem rozważany w tym rozdziale. Zadanie polega na znalezieniu

$$I^* = \inf_{F \in \mathcal{F}(D)} \left[\int_X w(x) dF(x) \right], \quad (9.1.4)$$

Oczywiście interesuje nas nie tylko koszt „najtańszego” planu I^* , ale także plan $F^* \in \mathcal{F}(D)$, dla którego infimum w (9.1.4) jest osiągnięte, jeśli taki plan istnieje.

Warto odnotować, że jeśli

- zbiór X jest domknięty i ograniczony,
- funkcje $v(x)$ oraz $w(x)$ są na tym zbiorze ciągle,
- dla pewnego planu $F^* \in \mathcal{F}(D)$ zachodzi $I^* = I(F^*) < \infty$,
- zbiór macierzy

$$\left\{ \int_X v(x) v^T(x) dF(x) : F \in \mathcal{F}(D) \right\}$$

jest domknięty i ograniczony, to istnieje plan, powiedzmy $\hat{F} \in \mathcal{F}(D)$, który jest miarą dyskretną skupioną w skończonej liczbie punktów zbioru X i taką, że $I(F^*) = I(\hat{F})$. Dowód tego faktu jest bezpośrednim wnioskiem z twierdzenia Carathodory’ego.

W następnych podrozdziałach przedstawiamy dwie klasy problemów, które są szczególnymi przypadkami zadania (9.1.4), a jednocześnie wpisują się w klasyczne nurty teorii planowania eksperymentu.

Plany ortogonalne

Wybermy jako pożądaną macierz informacyjną $D = I_r$, czyli macierz jednostkową $r \times r$. Wówczas zadanie (9.1.4) interpretować można jako zadanie znalezienia ortogonalnego planu, ale nie dowolnego – jak w klasycznym sformułowaniu planowania ortogonalnego – lecz planu o minimalnym koszcie.

Jest jeszcze jedna różnica między naszym sformułowaniem zadania a klasycznym planowaniem ortogonalnym. Mianowicie, w tym ostatnim wymaga się, by macierz informacyjna była macierzą diagonalną, lecz – w przeciwieństwie do naszego zadania – nie precyzuje się dokładnie wartości diagonalnych elementów tej

macierzy. Z tego powodu rozważane przez nas zadanie nazywać będziemy zadaniem planowania ortogonalnego z ograniczeniami. Ograniczenia takie warto nałożyć, gdyż, w przypadku gdy macierz informacyjna jest diagonalna, odwrotności diagonalnych elementów macierzy informacyjnej są proporcjonalne do wariancji oszacowań parametrów. Uzyskanie oszacowań o wariancji mniejszej niż niezbędna może się zatem okazać nadmiernie kosztowne.

Plany o symetrii obrotowej

Wybermy macierz D w ten sposób, by spełniony był warunek

$$v^T(x) D^{-1} v(x) = h(\|x\|), \quad x \in X \quad (9.1.5)$$

gdzie $\|x\| = (x^T x)^{1/2}$ natomiast $h(\cdot)$ jest pewną nieujemną i niemalejącą funkcją, którą także my wybieramy. Warunek (9.1.5) oznacza, że pożądana macierz informacyjna ma zapewniać tę samą wariancję oszacowania wyjścia we wszystkich punktach, które są równoodległe od centrum¹ eksperymentu. Cechę tę posiadają klasyczne plany o symetrii obrotowej (por. [29]). W klasycznym sformułowaniu zadania planowania o symetrii obrotowej zwykle nie zadaje się funkcji $h(\cdot)$. My będziemy funkcję tę traktować jako zadaną i dlatego zadanie doboru planu o minimalnym koszcie i takiego, który zapewnia spełnienie (9.1.5) nazywać będziemy planowaniem o symetrii obrotowej z ograniczeniami.

9.2. Charakteryzacje planów optymalnych

Przypomnijmy, że aby sprecyzować zadanie rozważane w tym rozdziale, podać trzeba wektor funkcji $v(x)$, który rozpina badaną funkcję regresji, pożądaną macierz informacyjną D , funkcję kosztu w oraz obszar planowania X . W celu skrócenia zapisów zadanie takie oznaczamy będziemy jako problem (w, D, v, X) .

Zauważmy, że – na skutek braku unormowania planów F – zbiór osiągalnych macierzy informacyjnych jest teraz zbiorem nieograniczonym, a dokładniej, jest stożkiem o postaci

$$\mathcal{M} \stackrel{\text{def}}{=} \left\{ M : M = \int_X v(x)v^T(x) dF(x), \quad F \in \mathcal{F} \right\}. \quad (9.2.6)$$

Termin stożek w odniesieniu do \mathcal{M} jest uzasadniony w tym sensie, że jeśli pewna macierz D należy do \mathcal{M} to także $\gamma D \in \mathcal{M}$ dla dowolnego $\gamma > 0$. Własność ta wynika z następującego prostego rozumowania: jeśli $D \in \mathcal{M}$, to istnieje $\tilde{F} \in \mathcal{F}$

¹ Dla wygody przyjęto, że centrum położone jest w punkcie 0.

takie, że $D = \int_X v(x)v^T(x) d\tilde{F}(x)$. Do zrealizowania macierzy informacyjnej γD wystarczy użyć planu $\gamma \tilde{F}$, który także należy do \mathcal{F} , gdyż nie nałożyliśmy warunku unormowania.

Aby zadanie (w, D, v, X) było warte rozważania, założyć musimy, że $D \in \mathcal{M}$. Aby móc scharakteryzować plany optymalne, potrzebne nam będzie założenie, które jest nieco mocniejsze:

C_1) zadana macierz informacyjna D leży we wnętrzu stożka \mathcal{M} , co zapisujemy jako $D \in \text{Int}\mathcal{M}$. Dla klasycznych funkcji regresji, rozpiętych na

$$v(x) = [1, x, \dots, x^{r-1}]^T, \quad x \in [-1, 1]$$

oraz

$$v(x) = [1, \sin(x), \cos(x), \sin(2x), \cos(2x) \dots]^T, \quad x \in [0, 2\pi]$$

znane są warunki dostateczne dla zachodzenia warunku C_1) (por. [64]). Udowodnimy następujący lemat, który także podaje warunki dostateczne dla spełnienia C_1). W celu zapewnienia jednolitości odwołań, lemat ten nazywać będziemy dalej warunkiem (dostatecznym) C_2).

Lemat 9.1. (C_2) *Jeśli macierz D jest dodatnio określona i dla każdego $y \in R^r$, $y \neq 0$ istnieje $x \in X \subset R^s$ oraz $\gamma \neq 0$ takie, że $y = \gamma v(x)$, to $D \in \text{Int}\mathcal{M}$, czyli spełniony jest warunek C_1).*

Dowód. Wskażemy F , które realizuje macierz D . W tym celu zauważmy, że macierz D – jako, że jest macierzą symetryczną i dodatnio określoną – może być przedstawiona w postaci

$$D = \sum_{j=1}^r \lambda_j d_j d_j^T, \quad (9.2.7)$$

gdzie d_j są wektorami własnymi macierzy D , a $\lambda_j > 0$ są odpowiadającymi im wartościami własnymi tej macierzy, a zatem spełnione są

$$D d_j = \lambda_j d_j, \quad j = 1, 2, \dots, r. \quad (9.2.8)$$

Założenie poczynione w lemacie pozwala nam stwierdzić, że wektory d_j potrafimy zrealizować wybierając $x_j \in X$ oraz mnożniki γ_j , dla których zachodzi $d_j = \gamma_j v(x_j)$, $j = 1, 2, \dots, r$. Jako plan F realizujący D możemy teraz wybrać dyskretną miarę skupioną w punktach x_j , którym przypisujemy wagi $\lambda_j \gamma_j^2$, $j = 1, 2, \dots, r$. W ten sposób pokazaliśmy, że $\mathcal{F}(D) \neq \emptyset$. Rezultat ten i to, że D jest macierzą dodatnio określoną pozwala stwierdzić, że $D \in \text{Int}\mathcal{M}$ (por. [64], rozdział XIII, § 1). •

Przytoczymy rezultat Karlina i Issi [64], Chapter XII, § 2 w notacji używanej w tejże monografii. Rezultat ten pozwoli nam scharakteryzować rozwiązania na-

szego problemu. Niech \mathcal{T} będzie odcinkiem na prostej i niech będzie dany wektor $c^0 \in R^{n+1}$ oraz wektor funkcji $u : \mathcal{T} \rightarrow R^{n+1}$. Zdefiniujemy zbiór

$$V(c^0) = \left\{ \sigma \in \mathcal{F} : \int_{\mathcal{T}} u(t) d\sigma(t) = c^0 \right\}.$$

Niech dana będzie także funkcja $\Omega : \mathcal{T} \rightarrow R$. Rozważać będziemy następujący problem

$$I_{\min} \stackrel{\text{def}}{=} \inf_{\sigma \in V(c^0)} \int_{\mathcal{T}} \Omega(t) d\sigma(t).$$

A oto zapowiadane twierdzenie Karlina i Issi.

Twierdzenie 9.1. *Załóżmy, że*

$$c^0 \in \text{Int} \left\{ c : c = \int_{\mathcal{T}} u(t) d\sigma(t), \sigma \in \mathcal{F} \right\}.$$

Niech funkcja $\Omega(t)$ będzie taka, że następujący zbiór funkcji

$$P_- \stackrel{\text{def}}{=} \left\{ v(\cdot) = a^T u(\cdot) : v(t) \leq \Omega(t), t \in \mathcal{T} \right\}$$

nie jest zbiorem pustym.

Zdefiniujemy jeszcze zbiór \mathcal{A} wszystkich takich wektorów $a \in R^{n+1}$, dla których $a^T u(\cdot) \in P_-$. Wówczas

$$I_{\min} = \sup_{a \in \mathcal{A}} (a^T c^0).$$

Zastosujmy ten rezultat do naszego problemu (w, D, v, X) .

Wniosek 9.1. *Niech \mathcal{B} będzie zbiorem wszystkich symetrycznych $r \times r$ macierzy takich, że*

$$\forall x \in X \quad v^T(x) B v(x) \leq w(x). \quad (9.2.9)$$

Jeśli D jest macierzą dodatnio określoną, i spełnione jest założenie C_1), to istnieje taka symetryczna $r \times r$ macierz $B^ \in \mathcal{B}$, dla której*

$$I^* = \text{tr}[B^* \cdot D] = \sup_{B \in \mathcal{B}} \text{tr}[B \cdot D], \quad (9.2.10)$$

gdzie I^ jest zdefiniowane przez (9.1.4).*

Jako instruktażowy przykład zastosowania Wniosku 9.1 rozpatrzmy następujące zadanie. Niech kula o promieniu jeden $S = \{x : \|x\| \leq 1\}$ będzie naszym obszarem planowania (zbiorem X). Na zbiorze tym estymujemy regresję

liniową bez wyrazu wolnego, tzn. $v(x) = x$. Jako funkcję kosztu przyjmujemy liczbę eksperymentów, czyli $w(x) = \chi_S(x)$ jest funkcją charakterystyczną zbioru S ($\chi_S(x) = 1$, gdy $x \in S$ i zero poza S).

Rozważamy zatem problem (χ_S, D, x, S) , przy założeniu że macierz D jest zadana i dodatnio określona. Oznaczmy przez d_j i λ_j , $j = 1, 2, \dots, r$ unormowane (tzn. $\|d_j\| = 1$) wektory i wartości własne macierzy D .

Wykażemy, korzystając z Wniosku 9.1, że plan F^* skupiony w punktach $x_j = d_j$ z wagami λ_j , $j = 1, 2, \dots, r$ minimalizuje liczbę eksperymentów. Zbiór \mathcal{B} składa się z takich macierzy B , dla których $x^T B x \leq 1$. Zgodnie z Wnioskiem 9.1 i równaniem (9.2.7),

$$I^* = \sup_{B \in \mathcal{B}} \text{tr}[B \cdot D] = \sup_{B \in \mathcal{B}} \sum_{j=1}^r \lambda_j d_j^T B d_j \leq \sum_{j=1}^r \lambda_j, \quad (9.2.11)$$

przy czym równość w ostatnim fragmencie wzoru (9.2.11) osiągnięta jest dla macierzy $B^* = I_r$. Wnioskujemy zatem, że $I^* = \sum_{j=1}^r \lambda_j$, a co więcej, opisany wyżej plan zapewnia, że ten minimalny koszt zostaje osiągnięty, gdyż

$$I(F^*) = \int_S \chi_S(x) dF^*(x) = \sum_{j=1}^r \lambda_j$$

(ostatnia równość wynika z faktu, że wagi punktów d_j wynoszą λ_j a $\chi_S(d_j) = 1$).

Reinterpretując rezultaty z ([64], Chapter XII, § 2), można podać także następującą charakteryzację planu optymalnego. Jeśli plan $F^* \in \mathcal{F}(D)$ jest kosztowo optymalny, to jego nośnik zawarty jest w następującym zbiorze

$$X^* = \left\{ x \in X : v^T(x) B^* v(x) = w(x) \right\},$$

gdzie B^* jest takie jak we Wniosku 9.1.

W celu zwięzłego sformułowania następnego rezultatu warto przypomnieć pojęcie ortonormalnego układu funkcji. Zestaw funkcji $v_j(x)$, $j = 1, 2, \dots, r$ nazwiemy ortonormalnym na zbiorze X z funkcją wagową $p(x) > 0$, jeżeli spełnione są następujące warunki

$$\int_X v_j(x) v_k(x) p(x) dx = 0, \quad k \neq j, \quad k, j = 1, 2, \dots, r, \quad (9.2.12)$$

$$\int_X v_j^2(x) p(x) dx = 1, \quad j = 1, 2, \dots, r. \quad (9.2.13)$$

Poniższe nierówności dla I^* pozwalają na znalezienie planów optymalnych w tych przypadkach, gdy potrafimy wykazać, że dolna granica jest osiągnięta przez pewien plan dopuszczalny.

Twierdzenie 9.2. *Jeśli spełnione są założenia poczynione we Wniosku 9.1, to:*

1. *koszt minimalny spełnia nierówność*

$$I^* \geq \mu r, \quad \text{gdzie} \quad \mu \stackrel{\text{def}}{=} \inf_{x \in X} \left[\frac{w(x)}{v^T(x) D^{-1} v(x)} \right], \quad (9.2.14)$$

2. *a jeśli ponadto:*

- *D jest macierzą diagonalną o jednakowych elementach na diagonalu o wartości, powiedzmy, d ,*
- *składowe wektora $v(x)$ są ortonormalne na X z funkcją wagową $p(x) > 0$, $x \in X$,*

to koszt minimalny

$$I^* \leq d \int_X w(x) p(x) dx. \quad (9.2.15)$$

Dowód. Zauważmy, że dla każdego $0 < \alpha \leq \mu$ zachodzi $\alpha D^{-1} \in \mathcal{B}$. Teraz nierówność (9.2.14) wynika z następującego ciągu zależności

$$I^* = \sup_{B \in \mathcal{B}} \text{tr}[B D] \geq \sup_{\alpha \leq \mu} \text{tr}[\alpha D^{-1} \cdot D] = \mu r.$$

W celu udowodnienia nierówności (9.2.15) zdefiniujemy zbiór \mathcal{B}_{ex} wszystkich symetrycznych $r \times r$ macierzy H takich, że

$$\text{tr}[H] \leq \int_X w(x) p(x) dx. \quad (9.2.16)$$

Wykażemy teraz, że $\mathcal{B} \subset \mathcal{B}_{ex}$. W tym celu wybierzmy dowolną macierz $B \in \mathcal{B}$, pomnóżmy obie strony (9.2.9) przez $p(x)$, a następnie scałkujmy po zbiorze X . W wyniku stwierdzamy, że dla B zachodzi nierówność $\text{tr}[B] \leq \int_X w(x) p(x) dx$, która pozwala stwierdzić, że $B \in \mathcal{B}_{ex}$ (por. (9.2.16)). Kładąc $D = d I_r$ i korzystając z Wniosku 9.1, otrzymamy

$$I^* = \sup_{B \in \mathcal{B}} \text{tr}[B D] \leq \sup_{B \in \mathcal{B}_{ex}} \text{tr}[B D] = d \int_X p(x) w(x) dx,$$

co kończy dowód drugiej części twierdzenia. •

W następnym podrozdziale wykażemy, że znak równości w oszacowaniach (9.2.14) i (9.2.15) jest w pewnych przypadkach osiągalny, a więc nierówności tych nie można w ogólnym przypadku uściślić.

9.3. Czy klasyczne eksperymenty są planami o minimalnym koszcie?

W podrozdziale tym spojrzymy na klasyczne plany ortogonalne i plany o symetrii obrotowej z punktu widzenia kosztu.

Planowanie ortogonalne jako problem kanoniczny

Rozpocznijmy ten podrozdział od kilku, łatwych do udowodnienia, rezultatów, które pokazują, że zadanie planowania eksperymentów ortogonalnych o minimalnej liczbie obserwacji jest, w pewnym sensie, zadaniem kanonicznym dla całej klasy rozważanych w tym rozdziale problemów (w, D, v, X) . Kanoniczność rozumiana jest tu w ten sposób, że większość „niezdegenerowanych” problemów (w, D, v, X) sprowadzić można do odpowiednio dobranego problemu planowania ortogonalnego o minimalnej liczbie obserwacji i to ostatnie zadanie rozważać można w obszarze planowania o prostej geometrii, a następnie przetransformować plan do oryginalnego zbioru X .

Lemat 9.2. *Rozważmy problem (w, D, v, X) . Oznaczmy przez Λ diagonalną macierz o r wierszach i kolumnach, o diagonalnych elementach równych wartościom własnym macierzy D . Przez P oznaczmy $r \times r$ macierz, której kolumnami są unormowane wektory własne macierzy D . Jeśli $F^* \in \mathcal{F}(D)$ i I^* są optymalnym rozwiązaniem problemu (w, D, v, X) , to są one jednocześnie rozwiązaniem problemu $(w, \Lambda, P^T \cdot v, X)$, to znaczy problemu planowania dla estymacji regresji rozpiętej przez zestaw funkcji $P^T \cdot v(x)$ i zadaną macierzą informacyjną Λ .*

Dowód. Macierz D możemy przedstawić w postaci $D = P \Lambda P^T$, a macierz P spełnia $P \cdot P^T = I_r$. Plan F^* jest optymalny dla problemu (w, D, v, X) i zapewnia, że

$$D = P \Lambda P^T = \int_X v(x) v^T(x) dF^*(x). \quad (9.3.17)$$

Mnożąc obie strony tych równości lewostronnie przez $P^{-1} = P^T$ i prawostronnie przez $(P^T)^{-1} = P$, stwierdzamy, że ten sam plan realizuje również macierz Λ dla regresji rozpiętej przez $P^T \cdot v(x)$. Pozostaje jeszcze sprawdzić, że plan F^* jest również optymalnym planem dla problemu $(w, \Lambda, P \cdot v, X)$. Przypuszczenie, że dla tego zadania istnieje plan o koszcie mniejszym niż I^* prowadzi natychmiast do sprzeczności, bo gdyby taki plan istniał, to moglibyśmy wykorzystać go w problemie (w, D, v, X) . •

Dowody dwóch następujących lematów są równie elementarne jak poprzedni i dlatego je pominiemy.

Lemat 9.3. *Niech plan $F^* \in \mathcal{F}(D)$ o koszcie I^* będzie optymalnym rozwiązaniem problemu (w, D, v, X) , przy czym zakładamy, że $w(x) > 0$, $x \in X$.*

- Oznaczmy przez χ_X funkcję charakterystyczną zbioru X . Dla dowolnych zbiorów borelowskich $A \subset X$ zdefiniujmy miarę

$$C^*(A) = \int_A w(x) dF^*(x).$$

Miara ta jest optymalnym rozwiązaniem problemu $(\chi_X, D, w^{-1/2}v, X)$, a koszt tego rozwiązania wynosi I^* . Jest to zadanie minimalizacji liczby eksperymentów zamiast ogólnej funkcji kosztu, lecz ze zmodyfikowanymi funkcjami rozpinającymi regresję (w miejsce wektora $v(x)$ wstawić należy $w^{-1/2}(x)v(x)$).

- Odwrotnie, jeśli G^* jest planem optymalnym w zadaniu $(\chi_X, D, w^{-1/2}v, X)$, to miara

$$F^*(A) = \int_A w^{-1}(x) dG^*(x)$$

jest optymalnym rozwiązaniem problemu (w, D, v, X) .

Lemat 9.4. Niech $X \subset R^k$ oraz $X' \subset R^k$ będą obszarami planowania (zbiorami domkniętymi i ograniczonymi). Oznaczmy przez $\varphi : X \rightarrow X'$ mierzalną i różnowartościową funkcję, która odwzorowuje X na X' , a przez $\psi : X' \rightarrow X$ oznaczmy funkcję odwrotną do φ . Jeśli F^* jest optymalnym rozwiązaniem problemu (w, D, v, X) , to miara zdefiniowana na zbiorach borelowskich $C \subset X'$ następująco

$$H^*(C) = F^*({x \in X : x = \psi(x'), x' \in C})$$

jest optymalnym rozwiązaniem zadania $(w \circ \varphi, D, v \circ \varphi, X')$, gdzie przez \circ oznaczono złożenie odwzorowań.

Plany o symetrii obrotowej dla regresji liniowej na kuli

Z Twierdzenia 9.2 otrzymujemy następujący wniosek dla planów o symetrii obrotowej.

Wniosek 9.2. Niech $v(x)$ oraz D będą takie, że dla pewnej niemalejącej i ciągłej w X funkcji $h(x) > 0$, $x \in X$ spełniony jest warunek

$$v^T(x) D^{-1} v(x) = h(\|x\|), \quad x \in X. \quad (9.3.18)$$

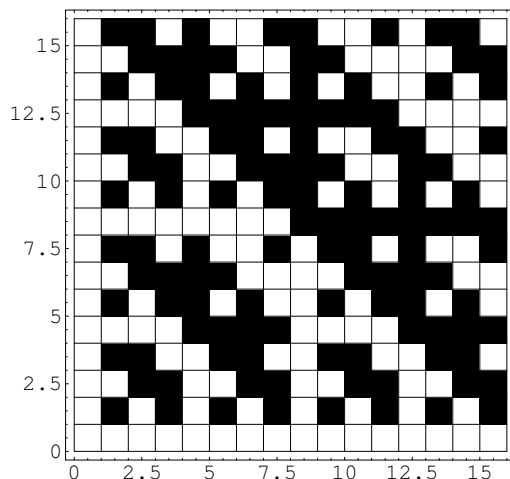
Zdefiniujmy $h^* = \max_{x \in X} h(\|x\|)$. Wówczas dla funkcji kosztu $w(x) = \chi_X(x)$ całkowity koszt minimalny ograniczony jest przez $I^* \geq r/h^*$.

Zastosujemy ten wniosek w celu znalezienia planu o minimalnej liczbie pomiarów ($w(x) = \chi_X(x)$) na kuli jednostkowej ($X \stackrel{\text{def}}{=} \{x : \|x\| \leq 1\}$) w problemie estymacji parametrów regresji liniowej $v(x) = x \in R^r$ (bez wyrazu wolnego).

Jako zadaną macierz informacyjną przyjmijmy macierz jednostkową $D = I_r$. Mamy wówczas $h(\|x\|) = \|x\|^2$ i $h^* = 1$. Na podstawie Wniosku 9.2 otrzymujemy zatem $I^* \geq r$. Wystarczy teraz wskazać plan, który tę dolną granicę osiąga i realizuje macierz $D = I_r$.

Bezpośrednim rachunkiem można sprawdzić, że planem takim jest F^* , który przypisuje wagi jednostkowe wszystkim punktom

$$x_i = [0, \dots, 0, 1, 0, \dots, 0]^T \quad (1 \text{ na } i\text{-tej pozycji}), \quad i = 1, 2, \dots, r.$$



Rys. 9.1. Rozkład wartości ± 1 w transponowanej macierzy Hadamarda H_{16}

Nie jest to jedyny plan optymalny. Tę samą macierz informacyjną i ten sam poziom kosztu uzyskamy, stosując plan skupiony w $2r$ punktach o postaci:

$$[0, \dots, 0, \pm 1, 0, \dots, 0]^T,$$

przy czym znaki \pm występują na i -tej pozycji, $i = 1, 2, \dots, r$. Każdemu z takich punktów przypisujemy wagi $1/2$.

Zauważmy, że plany te są także planami ortogonalnymi. Omówieniu planów ortogonalnych poświęcony jest następny podrozdział.

Plany ortogonalne dla regresji liniowej na kostce

W zadaniu estymacji liniowej regresji $v(x) = x$ jako obszar planowania wybierzmy kostkę jednostkową $X \stackrel{\text{def}}{=} \{x : |x^{(i)}| \leq 1, i = 1, 2, \dots, r\}$, jako zadaną macierz informacyjną przyjmijmy $D = r I_r$, natomiast $w(x) = \chi_X(x)$.

Na podstawie Twierdzenia 9.2 otrzymujemy $I^* \geq r$. Plany, które zapewniają osiągnięcie dolnej granicy można opisać następująco:

- wszystkim opisanym niżej punktom planu przypisujemy wagi równe 1,
- jako punkty planu przyjmujemy kolejne kolumny macierzy Hadamarda o wymiarach $r \times r$.

Oznaczmy taką macierz przez H_r . Jej elementami są liczby ± 1 dobrane w taki sposób, że $H_r \cdot H_r^T = r \cdot I_r$. Jako źródło informacji na temat istnienia i konstrukcji macierzy Hadamarda dla różnych wartości r można wskazać [82]. Tutaj opiszemy najprostsza z takich konstrukcji, która pozwala uzyskiwać macierze Hadamarda

$$H_{16} = \begin{bmatrix} + & + & + & + & + & + & + & + & + & + & + & + & + & + & + & + \\ + & - & + & - & + & - & + & - & + & - & + & - & + & - & + & - \\ + & + & - & - & + & + & - & - & + & + & - & - & + & + & - & - \\ + & - & - & + & + & - & - & + & + & - & - & + & + & - & - & + \\ + & + & + & + & - & - & - & - & + & + & + & + & - & - & - & - \\ + & - & + & - & - & + & - & + & + & - & + & - & - & + & - & + \\ + & + & - & - & - & - & + & + & + & + & - & - & - & - & + & + \\ + & - & - & + & - & + & + & - & + & - & - & + & - & + & + & - \\ + & + & + & + & + & + & + & + & - & - & - & - & - & - & - & - \\ + & - & + & - & + & - & + & - & - & + & - & + & - & + & - & + \\ + & + & - & - & + & + & - & - & - & - & + & + & - & - & + & + \\ + & - & - & + & + & - & - & + & - & + & + & - & - & + & + & - \\ + & + & + & + & - & - & - & - & - & - & - & - & + & + & + & + \\ + & - & + & - & - & + & - & + & - & + & - & + & + & - & + & - \\ + & + & - & - & - & - & + & + & - & - & + & + & + & + & - & - \\ + & - & - & + & - & + & + & - & - & + & + & - & + & - & - & + \end{bmatrix}.$$

Tabela 9.1. Struktura macierzy Hadamarda H_{16}

o wymiarach $2^k \times 2^k$, $k = 1, 2, \dots$. Punktem startowym jest macierz H_2 o postaci:

$$H_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}, \quad H_4 = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix}. \quad (9.3.19)$$

Macierz Hadamarda o wymiarach $2^k \times 2^k$ otrzymujemy w wyniku zastosowania iloczynu Kroneckera do macierzy o wymiarach $2^{k-1} \times 2^{k-1}$, $k = 2, 3, \dots$. Na przykład, macierz H_4 , pokazaną w prawej części wzoru (9.3.19), obliczamy następująco: $H_4 = H_2 \otimes H_2$. Podobnie, $H_{16} = H_4 \otimes H_4$. Macierz tę pokazano w tabeli 9.1, a jej transpozycję przedstawiono na rysunku 9.1.

10. Sekwencje planów nadążających za zmianami otoczenia

Celem tego podrozdziału jest przedstawienie rezultatów, które pozwalają planować sekwencje eksperymentów nadążających za zmianami otoczenia. Zmiany te muszą być mierzalne (obserwowalne), a do czynników wywołujących zmienność zaliczyć można

- upływ czasu,
- zmienne w czasie, mierzalne zakłócenia,
- starzenie się materiałów.

W dwóch ostatnich przypadkach mamy do czynienia z pośrednimi skutkami upływu czasu. Celem nadążania planu za zmianami otoczenia jest – jak zwykle – możliwie dokładna estymacja parametrów modelu. Celem użytkowym może być, na przykład zastosowanie w ocenie jakości wytwarzania, gdyż stwierdzenie zmienności estymowanych parametrów wraz ze zmianami otoczenia może być sygnałem, że nastąpiło rozregulowanie procesu lub zmiana jakości surowca.

Załóżmy, na przykład, że nasze obserwacje są pomiarami temperatury wzdłuż wlewka miedzianego, który jest podgrzewany przed walcowaniem. Jeśli w wyniku estymacji parametrów stwierdzimy, że współczynnik przewodnictwa ciepła nie jest stały na całej długości wlewka, to w wyniku równomiernego nagrzewu wlewka ten nie zostanie równomiernie nagrzany. Gdyby w takiej sytuacji wlewek ten poddany został operacji walcowania i wyciągania, to powstały w wyniku drut miedziany będzie miał nierównomierną strukturę i może się rwać podczas dalszej przeróbki.

Omawiane tu rezultaty są uogólnieniem wyników z pracy autora [127], w której badano zagadnienie sterowania ruchomymi czujnikami.

10.1. Sformułowanie problemu

Niech $q(x, t; a^0)$ oznacza odpowiedź obiektu na zestaw wejść $x \in X$ w sytuacji, gdy wartość zewnętrznego czynnika wynosi t . q zależy także od wektora $a \in R^r$ nieznanych parametrów. Oznaczmy przez $x_i(t) \in X$, wektor wartości i -tego wejścia zastosowanego wówczas, gdy zewnętrzny czynnik ma wartość t . Zakładamy przy tym, że zewnętrzny czynnik przyjmuje wartości:

- albo z domkniętego przedziału $\mathcal{T} = [0, \tau]$, $\tau > 0$,
- albo z dyskretnego zbioru $\mathcal{T} = \{t_1, t_2, \dots, t_p\}$, przy czym t_1, t_2, \dots traktować możemy jak „etykiety” o dość dowolnej naturze matematycznej.

W przypadku a) przez $|\mathcal{T}|$ oznacza długość przedziału $[0, \tau]$ lub licznosc zbioru $|\mathcal{T}|$, gdy jest on zbiorem skończonym (przypadek b)).

Cała teoria prezentowana w tym rozdziale przenosi się na przypadek dowolnego innego niż $[0, \tau]$ przedziału oraz na przypadek, gdy t jest wektorem przyjmującym wartości z domkniętego i ograniczonego zbioru.

Obserwacje wyjścia badanego procesu mają następującą postać:

$$y_i(t) = q(x_i(t), t; a^0) + \varepsilon(x_i(t), t), \quad t \in \mathcal{T}, \quad (10.1.1)$$

gdzie $i = 1, 2, \dots, I$ przy czym $I \geq 1$ oznacza liczbę wektorów wejść zastosowanych, gdy zewnętrzny czynnik ma wartość t . Dla uproszczenia przyjęliśmy, że dla różnych t liczba stosowanych wejść jest taka sama. Uogólnienie na przypadek, gdy I zależy od t jest łatwe.

W odniesieniu do zakłóceń pomiarowych i modelu q przyjmujemy następujące założenia:

A₁ $\varepsilon(x, t)$, $x \in X$, $t \in \mathcal{T}$ jest realizacją gausowskiego pola losowego o wartości średniej równej zeru.

A₂ Pole $\varepsilon(x, t)$ jest nieskorelowane w tym sensie, że funkcja korelacji R_ε spełnia $R_\varepsilon(x, y, t, \tau) = \delta(x - y) \cdot \delta(t - \tau)$, gdzie δ jest funkcją delta Diraca.

A₃ Dla każdego $x \in X$ i $t \in \mathcal{T}$ funkcja $q(x, t; a)$ jest różniczkowalna względem a w punkcie $a = a^0$. (W praktyce nie znamy a^0 , więc założenie to oznacza wymaganie różniczkowalności w pewnym obszarze).

Dzięki temu założeniu mamy zagwarantowane, że funkcja

$$g(x, t) \stackrel{\text{def}}{=} \text{grad}_a q(x, t; a)|_{a=a^0} \quad (10.1.2)$$

jest poprawnie zdefiniowana. Zauważmy, że $g(x, t)$ jest r -wymiarowym wektorem kolumnowym.

A₄ Niech $\| \cdot \|$ oznacza normę euklidesową w R^r . Zakładamy, że

$$\sup_{x \in X} \int_{\mathcal{T}} \|g(x, t)\|^2 dt < \infty. \quad (10.1.3)$$

W przypadku, gdy zbiór \mathcal{T} jest dyskretny, całkę $\int_{\mathcal{T}}$ interpretować będziemy tu i dalej jako sumę.

A₅ Dla każdego $t \in \mathcal{T}$ składowe wektora $g(x, t)$ są liniowo niezależnymi i ciągłymi funkcjami w $x \in X$.

Wektor parametrów $a^0 \in R^r$ estymujemy na podstawie pomiarów (10.1.1). W rozważanym przypadku nie możemy podać dokładnej zależności pomiędzy dokładnością estymacji a użytym planem eksperymentu, ponieważ zależność wyjścia od parametrów nie jest liniowa. Narzędzi do tworzenia przybliżonych zależności w tak zwanych regularnych przypadkach dostarczają: nierówność Rao–Craméra oraz teoria estymatorów asymptotycznie efektywnych. Warunki regularności, które są warunkami dostatecznymi dla prawdziwości nierówności Rao–Craméra i asymptotycznej efektywności estymatorów są dość skomplikowane i, co ważniejsze, trudne do weryfikacji. Z tego powodu nie będziemy ich tu szczegółowo omawiać. Ograniczymy się do stwierdzenia, że jeżeli

1. w modelu (10.1.1) zakłócenia mają rozkład normalny,
2. spełnione są podane wyżej warunki A₁–A₅,
3. następująca granica istnieje

$$\lim_{|\mathcal{T}| \rightarrow \infty} |\mathcal{T}|^{-1} \sum_{i=1}^I \int_{\mathcal{T}} g(x_i(t), t) \cdot g^T(x_i(t), t) dt = \Upsilon \quad (10.1.4)$$

gdzie macierz Υ jest nieosobliwa (macierz ta zależy od planu eksperymentu, lecz zależność tę pomijamy w notacji). Symbol $|\mathcal{T}|$ oznacza licznosc zbioru, \mathcal{T} w przypadku gdy jest on skończony lub długość przedziału $(0, T)$ – w przypadku gdy czynnik zewnętrzny t przyjmuje wartości w skończonym przedziale,

4. parametry modelu $q(x, t; a)$ są identyfikowalne na podstawie pomiarów w punktach $x_i(t)$, $i = 1, 2, \dots, I$, $t \in \mathcal{T}$,

to spełnione są warunki wystarczające dla poprawności nierówności Rao–Cramera i asymptotycznej efektywności estymatora parametrów a , uzyskanego metodą najmniejszych kwadratów (por. [173]).

Pojęcie identyfikowalności oznacza, że dwa dowolne, lecz różne zestawy parametrów, powiedzmy θ' i θ'' , prowadzą do różnych zestawów odpowiedzi systemu, co oznacza, że $\theta' \neq \theta''$ implikuje

$$\{q(x_i(t), t; \theta') : i = 1, 2, \dots, I, t \in \mathcal{T}\} \neq \{q(x_i(t), t; \theta'') : i = 1, 2, \dots, I, t \in \mathcal{T}\}.$$

Oznaczmy przez $\hat{\theta}$ estymator uzyskany w wyniku minimalizacji błędu średniokwadratowego

$$\sum_{i=1}^I \int_{\mathcal{T}} (y_i(t) - q(x_i(t), t; \theta))^2 dt \quad (10.1.5)$$

względem θ . Wybierzmy $|\mathcal{T}|$ na tyle duże, że można pominąć obciążenie estymatora $\hat{\theta}$. Wówczas nierówność Rao–Craméra ma postać

$$\text{cov}(\hat{\theta}) \geq M_{\mathcal{T}}^{-1}, \quad (10.1.6)$$

gdzie $M_{\mathcal{T}}$ jest macierzą informacyjną zdefiniowaną wzorem (10.1.9), natomiast nierówność pomiędzy macierzami rozumiana jest w ten sposób, że $\text{cov}(\hat{\theta}) - M_{\mathcal{T}}^{-1}$ jest macierzą nieujemnie określoną. Statystyczne uzasadnienie posługiwania się macierzą informacyjną do oceny dokładności w nieliniowych problemach estymacji znaleźć można w [60] i [173].

W celu określenia $M_{\mathcal{T}}$ oceńmy wpływ i -tego wejścia na informację o parametrach. Oznaczmy przez $M_{(i)}$ składnik macierzy informacyjnej odpowiadający temu wymuszeniu. Wówczas

$$M^{(i)} = \int_{\mathcal{T}} g(x_i(t), t) g^T(x_i(t), t) dt \quad (10.1.7)$$

lub równoważnie

$$M^{(i)} = \int_{\mathcal{T}} \int_X g(x, t) g^T(x, t) \delta(x - x_i(t)) dx dt, \quad (10.1.8)$$

gdzie δ oznacza deltę Diraca. Biorąc pod uwagę założoną niezależność zakłóceń pomiarowych, stwierdzamy, że

$$M_{\mathcal{T}} = \sum_{i=1}^I M^{(i)}. \quad (10.1.9)$$

Zakładając, że $I \geq 1$ oraz $|\mathcal{T}| > 0$, wygodnie będzie posługiwać się unormowaną macierzą informacyjną, zdefiniowaną jako $M = M_{\mathcal{T}} / (I \cdot |\mathcal{T}|)$. Macierz tę wyrazić można w postaci

$$M = \frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \int_X g(x, t) g^T(x, t) \xi_I(x, t) dx dt, \quad (10.1.10)$$

gdzie

$$\xi_I(x, t) \stackrel{\text{def}}{=} \sum_{i=1}^{\bar{I}(t)} p_i(t) \cdot \delta(x - x_i(t)). \quad (10.1.11)$$

We wzorze (10.1.11) założono, że – dla ustalonego t – wejścia zostały ponumerowane w ten sposób, by

$$x_i(t) \neq x_j(t) \text{ jeśli } i \neq j; \quad i, j = 1, 2, \dots, \bar{I}(t), \quad (10.1.12)$$

gdzie $\bar{I}(t) \leq I$ jest liczbą istotnie różnych wejść stosowanych dla tego ustalonego $t \in \mathcal{T}$. Oznaczmy przez $n_i(t)$ krotność stosowania wejścia $x_i(t)$, $i = 1, 2, \dots, \bar{I}(t)$ w sytuacji, gdy $t \in \mathcal{T}$ ma ustaloną wartość. Teraz możemy objaśnić znaczenie

symbolu $p_i(t)$ w wyrażeniu (10.1.11). Definiujemy $p_i(t) \stackrel{\text{def}}{=} n_i(t)/I$ – czyli jest to częstość użycia i -tego wymuszenia, przy ustalonym $t \in \mathcal{T}$. Z określenia tego wynika, że każdego $t \in \mathcal{T}$

$$\sum_{i=1}^{\bar{I}(t)} p_i(t) = 1, \quad p_i(t) \geq 0, \quad i = 1, 2, \dots, \bar{I}(t). \quad (10.1.13)$$

$\xi_I(x, t)$, $x \in X$ ma zatem wszystkie cechy dyskretnego rozkładu prawdopodobieństwa.

W naszych rozważaniach ξ_I jest prototypem zmiennej decyzyjnej. W celu uproszczenia zadania optymalizacji zrezygnujemy z założenia, że $I p_i(t)$ jest liczbą całkowitą. Jest to uproszczenie analogiczne do tego, które stosowaliśmy w poprzednich rozdziałach.

Oznaczmy przez $\mu_t(x)$ miarę odpowiadającą (niewłaściwej) gęstości rozkładu $\xi_I(x, t)$, lecz bez wymagania, by $I p_i(t)$ było liczbą naturalną, to znaczy dla $t \in \mathcal{T}$

$$\mu_t(dx) = \sum_{i=1}^{\bar{I}(t)} p_i(t) \delta(x - x_i(t)) dx. \quad (10.1.14)$$

Wobec założenia A_4 i skończoności miar $\mu_t(x)$ możemy zagwarantować, że poniższa całka

$$\frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \int_X \|g(x, t)\|^2 \mu_t(dx) dt \quad (10.1.15)$$

ma skończoną wartość. W wyrażeniu tym oraz w dalszych wzorach stosujemy następującą konwencję: dla ustalonej funkcji, powiedzmy $\psi(x)$, ciągłej na X całka iloczynu tej funkcji z μ_t rozumiana jest zgodnie ze wzorem

$$\int_X \psi(x) \mu_t(dx) = \sum_{i=1}^{\bar{I}(t)} p_i(t) \psi(x_i(t)), \quad t \in \mathcal{T}. \quad (10.1.16)$$

Zbiór wszystkich dyskretnych rozkładów prawdopodobieństwa postaci (10.1.14) oznaczamy będziemy przez \mathcal{M}_d . Oznaczmy przez μ całą sekwencję planów, to znaczy

$$\mu = \{\mu_t \in \mathcal{M}_d : t \in \mathcal{T}\}. \quad (10.1.17)$$

Niech dalej \mathcal{M} oznacza klasę wszystkich sekwencji postaci (10.1.17).

Analizując (10.1.7–10.1.11) stwierdzamy, że dla $\mu \in \mathcal{M}$ unormowana macierz informacyjna M ma postać

$$M(\mu) = \frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \int_X g(x, t) \cdot g^T(x, t) \mu_t(dx) dt. \quad (10.1.18)$$

Możemy teraz sformułować zadanie D- optymalnego doboru sekwencji planów nadszających za zmianami otoczenia.

Definicja 10.1. *Sekwencję planów $\hat{\mu} \in \mathcal{M}$ nazwiemy D- optymalną, jeśli*

$$\det M(\hat{\mu}) = \sup_{\mu \in \mathcal{M}} \det M(\mu). \quad (10.1.19)$$

W następnym podrozdziale przedstawiamy charakteryzację optymalnego rozwiązania.

10.2. Warunki optymalności sekwencji planów

Podamy tu pełną charakteryzację optymalnej sekwencji planów w postaci warunku, który jest jednocześnie konieczny i dostateczny dla jej optymalności. Warunek ten można sprawdzać analitycznie, a jego przybliżone – numeryczne – sprawdzanie pozwala znajdować sekwencje bliskie optymalnym. Jednocześnie, w pewnych przypadkach wystarczy odwołać się do warunków, które są tylko dostateczne, lecz prostsze do sprawdzenia. Podamy dwa zestawy takich warunków i ich zastosowania do znajdowania optymalnych trajektorii ruchu czujników pomiarowych.

Charakteryzacja optymalnej sekwencji planów

Wykażemy najpierw następujący lemat.

Lemat 10.1. *Macierz informacyjna $M(\hat{\mu})$, odpowiadająca D- optymalnej sekwencji planów $\hat{\mu}$, jest macierzą nieosobliwą.*

Dowód. Istnienie sekwencji $\mu \in \mathcal{M}$, dla której $M(\mu)$ jest nieosobliwa, wynika bezpośrednio z (A₅). •

Dla sekwencji $\mu \in \mathcal{M}$ o nieosobliwej macierzy $M(\mu)$ zdefiniujmy funkcję $\varphi(x, t; \mu)$ następująco

$$\varphi(x, t; \mu) = g^T(x, t) M^{-1}(\mu) g(x, t); \quad x \in X, \quad t \in \mathcal{T}. \quad (10.2.20)$$

Dla danego ustalonego $t \in \mathcal{T}$ funkcja ta pełni taka samą rolę, jak funkcja wariancji odpowiedzi modelu, której używaliśmy w poprzednich rozdziałach. Jak pokazuje następnny lemat i poniższe twierdzenie, podobna jest też jej rola w formułowaniu warunków optymalności.

Lemat 10.2. *Dla każdej sekwencji planów $\mu \in \mathcal{M}$ o nieosobliwej macierzy $M(\mu)$ zachodzi następująca nierówność*

$$\frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \left[\max_{x \in X} \varphi(x, t; \mu) \right] dt \geq r. \quad (10.2.21)$$

Dowód. Rozważmy ciąg elementarnych równości

$$r = \operatorname{tr} [M^{-1}(\mu) M(\mu)] = \frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \int_X \varphi(x, t; \mu) \mu_t(dx) dt. \quad (10.2.22)$$

Ponieważ funkcja $\varphi(x, t; \mu)$ jest nieujemna oraz ciągła na zwartym zbiorze i dla każdego $t \in \mathcal{T}$ zachodzi $\int_X \mu_t(dx) = 1$, to nierówność (10.2.21) otrzymujemy bezpośrednio z (10.2.22). •

Twierdzenie 10.1. *Sekwencja planów $\hat{\mu} \in \mathcal{M}$ jest sekwencją D- optymalną wtedy i tylko wtedy, gdy spełniony jest następujący warunek:*

$$\frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \left[\max_{x \in X} \varphi(x, t; \hat{\mu}) \right] dt = r, \quad (10.2.23)$$

gdzie $r \geq 1$ jest liczbą estymowanych parametrów.

Dowód. Zdefiniujmy szczególną sekwencję planów $\mu^0 \in \mathcal{M}$, która lokuje wszystkie pomiary w jednym punkcie, którego położenie może zależeć od $t \in \mathcal{T}$. Oznaczmy ten punkt przez $x_0(t) \in X$ i wybierzmy jego współrzędne następująco:

$$x_0(t) = \arg \max_{x \in X} \varphi(x, t; \hat{\mu}). \quad (10.2.24)$$

Ponieważ $\varphi(x, t; \hat{\mu})$, traktowana jako funkcja x , jest ciągłą na domkniętym i ograniczonym zbiorze X , więc supremum w (10.2.24) jest osiągnięte w punktach zbioru X . Możemy zatem powiedzieć, że $x_0(t)$ to punkt, w którym osiągnięte jest maksimum $\varphi(x, t; \hat{\mu})$ względem $x \in X$. Punktów, w których to maksimum jest osiągnięte może być oczywiście więcej niż jeden. Nadużywając nieco notacji, zbiór takich punktów nadal oznaczać będziemy przez $x_0(t)$.

Wypukłość zbioru \mathcal{M} gwarantuje nam, że dla dowolnego $\alpha \in [0, 1]$ zachodzi

$$\mu^\alpha \stackrel{\text{def}}{=} [(1 - \alpha)\hat{\mu} + \alpha\mu^0] \in \mathcal{M}. \quad (10.2.25)$$

Załóżmy, że $\hat{\mu}$ jest optymalną sekwencją planów. Wówczas,

$$\left. \frac{\partial \ln \det M(\mu^\alpha)}{\partial \alpha} \right|_{\alpha=0} \leq 0, \quad (10.2.26)$$

ponieważ odstrojenie od sekwencji optymalnej nie może dawać sekwencji o większej wartości kryterium. Po wykonaniu różniczkowania w (10.2.26), otrzymamy nierówność

$$\frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \varphi(x_0(t), t; \hat{\mu}) dt \leq r, \quad (10.2.27)$$

która dowodzi, że warunek (10.2.23) jest warunkiem koniecznym optymalności, gdyż zachodzi także (10.2.21) i (10.2.24).

W celu udowodnienia dostateczności warunku (10.2.23) założmy, że jest on spełniony, lecz sekwencja planów $\hat{\mu}$ nie jest D-optymalna. Istnieje wówczas inna sekwencja planów, oznaczmy ją przez $\mu^* \in \mathcal{M}$, taka, że $\det M(\mu^*) > \det M(\hat{\mu})$. Nierówność ta i ścisła wypukłość funkcji $\ln \det[\cdot]$ na zbiorze \mathcal{M} prowadzą do nierówności

$$\frac{\partial}{\partial \alpha} \ln \det M[(1 - \alpha)\hat{\mu} + \alpha\mu^*] \Big|_{\alpha=0} > 0. \quad (10.2.28)$$

Różniczkowanie w (10.2.28) prowadzi do

$$r < \frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \int_X \varphi(x, t; \mu) \mu_t^*(dx) dt \leq \frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \left[\max_{x \in X} \varphi(x, t; \hat{\mu}) \right] dt. \quad (10.2.29)$$

Otrzymana sprzeczność między tą nierównością a warunkiem (10.2.23) dowodzi jego dostateczności. •

Wniosek 10.1. *Warunkiem dostatecznym optymalności sekwencji planów $\tilde{\mu} \in \mathcal{M}$ jest spełnienie poniższego warunku*

$$\forall t \in \mathcal{T} \quad \max_{x \in X} \varphi(x, t; \tilde{\mu}) = r. \quad (10.2.30)$$

Dowód. Dostateczność (10.2.30) wynika ze wstawienia go do (10.2.23). •

I warunek dostateczny optymalności sekwencji planów

W podrozdziale tym podamy zestawy warunków dostatecznych, które – w pewnych przypadkach – pozwalają uprościć zadanie znajdowania sekwencji optymalnych.

Twierdzenie 10.2. *Niech $\mu^* \in \mathcal{M}$ będzie sekwencją planów o nieosobliwej macierzy informacyjnej $M(\mu^*)$. Oznaczmy przez $x_i^*(t)$, $i = 1, 2, \dots, I^*(t)$, $t \in \mathcal{T}$ trajektorie zmian wejść, które odpowiadają sekwencji planów μ^* . Jeśli spełniony jest warunek*

$$\forall t \in \mathcal{T} \quad \max_{x \in X} \varphi(x, t; \mu^*) = \varphi(x_i^*(t), t; \mu^*) \quad (10.2.31)$$

to μ^ jest D-optymalną sekwencją planów. Ponadto, jeśli spełniony jest warunek (10.2.31), to zachodzi także*

$$\max_{t \in \mathcal{T}} \varphi(x_i^*(t), t; \mu^*) = r. \quad (10.2.32)$$

Dowód. Pomnóżmy obie strony (10.2.31) przez wagi $p_i^*(t)$, $i = 1, 2, \dots, I^*(t)$, odpowiadające sekwencji planów μ^* . Zsumujmy powstałe w ten sposób równości stronami i scałkujmy obie strony tych sum względem t . W rezultacie otrzymamy

$$\frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \left[\sup_{x \in X} \varphi(x, t; \mu^*) \right] dt = \frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \int_X \varphi(x, t; \mu^*) \mu_t^*(dx) dt = r, \quad (10.2.33)$$

gdzie ostatnia równość w poprzednim ciągu przekształceń wynika z (10.2.22). Teraz D-optymalność sekwencji μ^* wynika bezpośrednio z Twierdzenia 10.1.

Rozważmy drugie z wymienionych stwierdzeń. Z Lematu 10.2 otrzymujemy, że dla dowolnej sekwencji planów o nieosobliwej macierzy informacyjnej, a więc także dla μ^* , zachodzi

$$\sup_{t \in \mathcal{T}} \max_{x \in X} \varphi(x, t; \mu^*) \geq r. \quad (10.2.34)$$

Teraz drugie ze stwierdzeń udowodnimy przez sprowadzenie do sprzeczności. Przy założeniu, że jednocześnie warunek (10.2.31) jest spełniony, ale (10.2.32) nie zachodzi, prowadzi w świetle (10.2.34) do nierówności

$$\sup_{t \in \mathcal{T}} \left[\max_{x \in X} \varphi(x, t; \mu^*) \right] < r. \quad (10.2.35)$$

Skoro jednak supremum względem t funkcji w nawiasach kwadratowych nie przekracza r , to i sama ta funkcja nie przekracza r we wszystkich punktach \mathcal{T} . Innymi słowy

$$\forall t \in \mathcal{T} \quad \max_{x \in X} \varphi(x, t; \mu^*) < r. \quad (10.2.36)$$

Po scałkowaniu obu stron tej nierówności otrzymamy

$$\frac{1}{|\mathcal{T}|} \int_{\mathcal{T}} \left[\max_{x \in X} \varphi(x, t; \hat{\mu}) \right] dt < r, \quad (10.2.37)$$

co, w świetle Twierdzenia 10.1, przeczy udowodnionej wcześniej optymalności μ^* . Udowodniliśmy zatem, że

$$\sup_{t \in \mathcal{T}} \left[\max_{x \in X} \varphi(x, t; \mu^*) \right] = r. \quad (10.2.38)$$

Wystarczy teraz odwołać się do warunku (10.2.31), aby z (10.2.38) otrzymać (10.2.32). •

Można zauważyć, że warunek (10.2.31) ma postać podobną do zasady maksimum Pontriagina, znanej w teorii sterowania optymalnego. Warto wskazać jednak na dwie podstawowe różnice.

- W ogólnym przypadku zasada maksimum jest warunkiem koniecznym optymalności, podczas gdy warunek (10.2.31) jest warunkiem dostatecznym.
- W warunku (10.2.31) obie strony zależą od μ^* , jeśli tylko $r > 1$.

Przydatność Twierdzenia 10.2 do znajdowania sekwencji optymalnych pokażemy w następnym podrozdziale.

II warunek dostateczny optymalności sekwencji planów

Założmy, że gradient q względem parametrów ma postać (por. A₃)

$$g(x, t) = h(w(x, t)), \quad x \in X, \quad t \in \mathcal{T}, \quad (10.2.39)$$

gdzie w jest funkcją $w : X \times \mathcal{T} \rightarrow R^l$, $1 \leq l \leq r$, natomiast h jest r -wymiarowym wektorem funkcji $h : R^l \rightarrow R^r$. Założmy także, że zarówno funkcja w , jak i składowe wektora $h(\cdot)$ są funkcjami ciągłymi w swoich obszarach określoności. Przykłady problemów, dla których spełniony jest warunek (10.2.39) podamy w następnym podrozdziale.

Zdefiniujmy sekwencję zbiorów pomocniczych

$$Z_t \stackrel{\text{def}}{=} \{z : z = w(x, t), x \in X\}, \quad t \in \mathcal{T}. \quad (10.2.40)$$

Założenie o ciągłości funkcji $w(\cdot, \cdot)$ gwarantuje nam, że zbiory Z_t są zbiorami domkniętymi i ograniczonymi. Założmy, że zbiory Z_t nie zależą od $t \in \mathcal{T}$ i oznaczymy odpowiedni wspólny zbiór przez Z .

Rozważmy problem pomocniczy: znaleźć wektory $Z_i^* \in Z$, $i = 1, 2, \dots, L$ oraz liczby

$$\pi_i^* \geq 0, \quad i = 1, 2, \dots, L, \quad \sum_{i=1}^L \pi_i^* = 1$$

takie, że

$$\max \det \left[\sum_{i=1}^L \pi_i h(Z_i) h^T(Z_i) \right] = \det \left[\sum_{i=1}^L \pi_i^* h(Z_i^*) h^T(Z_i^*) \right], \quad (10.2.41)$$

gdzie supremum dotyczy wszystkich

$$Z_i \in Z \quad \text{oraz} \quad \pi_i \geq 0, \quad \sum_{i=1}^L \pi_i = 1.$$

Liczba wektorów $L \geq 1$ jest tu również zmienną decyzyjną i nie jest z góry ustalona.

Problem (10.2.41) jest już znanym nam zadaniem D-optymalnego planowania eksperymentu dla pomocniczo skonstruowanej funkcji regresji liniowej, rozpiętej

przez wektor funkcji $h(z)$ z obszarem planowania $Z \subset R^l$. Możemy przyjąć, że zadanie to potrafimy rozwiązać analitycznie lub numerycznie, co daje nam punkty planu Z_i^* i wagi π_i^* , $i = 1, 2, \dots, L$.

Mając to rozwiązanie, możemy znaleźć sekwencję wejść $x_i^*(t)$, $i = 1, 2, \dots, L$, rozwiązując względem x_i następujące równania dla każdego $t \in \mathcal{T}$ z osobna

$$Z_i^* = w(x_i, t), \quad i = 1, 2, \dots, L. \quad (10.2.42)$$

Możemy teraz utworzyć sekwencję planów μ^* w ten sposób, że dla każdego $t \in \mathcal{T}$ kojarzymy wejścia $x_i^*(t)$ ze stałymi (niezależnymi od t) wagami π_i^* , $i = 1, 2, \dots, L$.

Wniosek 10.2. *Jeśli spełnione są założenia (10.2.39) i (10.2.41), to tak utworzona sekwencja planów $\mu^* \in \mathcal{M}$ jest sekwencją D- optymalną.*

Dowód tego wniosku podajemy na końcu podrozdziału.

Uwaga 10.1. *W tym przypadku do skonstruowania D- optymalnej sekwencji planów wystarczy raz rozwiązać pomocniczy problem planowania (10.2.41), a następnie rozwiązać równania (10.2.42) dla poszczególnych $t \in \mathcal{T}$.*

Uwaga 10.2. *Istnienie rozwiązań równań (10.2.42) mamy zagwarantowane, gdyż każde $z \in Z_i^* \in Z$. Rozwiązania te nie muszą być jednoznaczne, co można wykorzystać następująco:*

- *gdy \mathcal{T} jest odcinkiem prostej, można szukać rozwiązań $x_i^*(t)$, które są ciągle jako funkcje zmiennej t ,*
- *gdy mamy dodatkowe kryterium wyboru wartości wejść, to możemy wybrać te spośród niejednoznacznie wyznaczonych $x_i^*(t)$, które zapewniają minimum kryterium pomocniczego.*

W celu zilustrowania sposobu posługiwania się Wnioskiem 10.2, rozważmy następującą klasę modeli określonych dla $x \in [0, 2\pi]$:

$$q(x, t; a) = \alpha_0 + \sum_{k=1}^K \alpha_k \cos[k(x - vt)] + \beta_k \sin[k(x - vt)] \quad (10.2.43)$$

gdzie

$$a \stackrel{\text{def}}{=} [\alpha_0, \alpha_1, \dots, \alpha_K, \beta_1, \dots, \beta_K]^T$$

jest wektorem nieznanych parametrów, natomiast $v > 0$ jest znane. Modele tej klasy można interpretować jako trygonometryczną aproksymację pewnej funkcji argumentu x , której położenie liniowo zmienia się w czasie z prędkością v .

Po obliczeniu gradientu (10.2.43) względem a , stwierdzamy, że wektor wrażliwości na zmiany parametrów a określony jest wzorem $g(x, t) = \tilde{h}(x - vt)$, gdzie

$$\tilde{h}(z) \stackrel{\text{def}}{=} [1, \cos z, \dots, \cos(Kz), \sin z, \dots, \sin(Kz)]^T.$$

Aby móc skorzystać z Wniosku 10.2, powinniśmy rozwiązać zadanie (10.2.41) z $h(z) = \tilde{h}(z)$, gdzie supremum obliczane jest w przedziale $[0, 2\pi]$. Problem taki już napotkaliśmy, rozwiązując zadanie D- optymalnego planowania dla regresji trygonometrycznej. Przypomnijmy, że jego rozwiązanie ma postać

$$\hat{z}_i = 2\pi(i-1)/r, \quad \hat{\pi}_i = 1/r, \quad i = 1, 2, \dots, r,$$

gdzie $r = 2K + 1$. Wybierzmy miarę $\hat{\mu} \in \mathcal{M}$ w ten sposób, by przypisać masy $\hat{\pi}_i$, $i = 1, 2, \dots, r$ do następujących trajektorii

$$\hat{x}_i(t) = \hat{z}_i + vt, \quad i = 1, 2, \dots, r \quad (10.2.44)$$

w każdym punkcie $t \in \mathcal{T}$. Wówczas, dla każdego $t \in \mathcal{T}$ spełnione będzie

$$\max_{x \in X} \tilde{h}^T(x - vt) M^{-1}(\hat{\mu}) \tilde{h}(x - vt) = r. \quad (10.2.45)$$

Zgodnie z Wnioskiem 10.1 $\hat{\mu}$ jest rozwiązaniem optymalnym.

Dowód wniosku 10.2.

Z twierdzenia Kiefiera–Wolfowitza o równoważności planów D- i G- optymalnych wynika, że zadanie (10.2.41) jest równoważne ze spełnieniem warunku

$$\max_{z \in Z} h^T(z) (H^*)^{-1} h(z) = r, \quad (10.2.46)$$

gdzie

$$H^* \stackrel{\text{def}}{=} \sum_{i=1}^L \pi_i^* h(Z_i^*) h^T(Z_i^*).$$

Ponadto supremum w (10.2.46) osiągnęte jest dla wszystkich Z_i^* , $i = 1, 2, \dots, L$. Bezpośrednim rachunkiem można sprawdzić, że macierz informacyjna dla sekwencji planów μ^* ma postać $M(\mu^*) = H^*$. Dla każdego $t \in \mathcal{T}$ zachodzi zatem

$$\begin{aligned} \max_{x \in X} \varphi(x, t; \mu^*) &= \max_{z \in Z_t} h^T(z) (H^*)^{-1} h(z) = \max_{z \in Z} h^T(z) (H^*)^{-1} h(z) \\ &= \varphi(x_i^*(t), t; \mu^*) = r. \end{aligned} \quad (10.2.47)$$

Zestawienie tych równości z Twierdzeniem 10.2 kończy dowód. •

10.3. Zastosowanie – sterowanie ruchomymi czujnikami

Sytuację opisaną w (10.2.39) spotykamy często w systemach opisywanych równaniami o pochodnych cząstkowych. Rozwiązania tych równań zależą zwykle od tak zwanych kompleksów bezwymiarowych. Kompleksy te zawierają na ogół

zależności między zmiennymi przestrzennymi i czasem. Uwagi te zilustrowano w przykładach opisanych w tym podrozdziale.

W przykładach tych funkcja $q(x, t, a)$, która występuje w modelu obserwacji (10.1.1), jest rozwiązaniem równania o pochodnych cząstkowych, które zależy od zmiennych przestrzennych x , czasu t i wektora nieznanymi parametrów a . Parametry te opisują zwykle stałe materiałowe, a odpowiedź obiektu $q(x, t, a)$ zależy od nich nieliniowo. Zauważmy, że przy takim podejściu zmienne przestrzenne x odgrywają rolę wielkości wejściowych naszego modelu. Innymi słowy wartości $x_i(t)$ oznaczają położenie i -tego czujnika pomiarowego w przestrzeni, a zależność jego położenia od czasu t oznacza, że czujniki mogą zmieniać swoje położenie w poszczególnych chwilach. Takie ruchome czujniki są w wielu dziedzinach nauki dość łatwo realizowalne w dzisiejszym stanie technik pomiarowych. Traktowanie położenia czujników jako wielkości decyzyjnych (wejściowych) pozwala nam zinterpretować zadanie (10.1.19) jako zadanie szukania trajektorii ruchu czujników, które są D- optymalne ze względu na dokładność estymacji nieznanymi parametrów w równaniu o pochodnych cząstkowych, które opisuje badany proces.

Wnioski 10.1 i 10.2 oraz Twierdzenie 10.2 pozwalają znaleźć w sposób analityczny sekwencję planów D- optymalnych, interpretowanych tutaj jako sekwencja położenia czujników pomiarowych.

Problem 10.1. *Rozważmy uproszczony model wymiennika ciepła. Uproszczenia polegają na przyjęciu, że zewnętrzny płaszcz wymiennika jest na tyle duży, iż jego pojemność cieplną można uznać za nieskończenie dużą. Zakładamy też dla uproszczenia, że temperatura medium chłodzącego w zewnętrznym płaszczu wymiennika wynosi zero stopni. Oznaczmy przez q temperaturę chłodzonego medium w chwili t i w punkcie położonym o x jednostek od początku wymiennika. Zmiany tej temperatury opisać można (w przybliżeniu) następującym równaniem*

$$\frac{\partial q(x, t; a)}{\partial t} + v \frac{\partial q(x, t; a)}{\partial x} + a q(x, t; a) = 0 \quad \text{dla } x > 0. \quad (10.3.48)$$

W równaniu (10.3.48), $v > 0$ oznacza prędkość medium wewnątrz wymiennika. Wielkość tę traktować będziemy jako znaną i stałą. Współczynnik a opisuje intensywność wymiany ciepła między medium chłodzonym i chłodzącym. Będziemy zakładać, że jest on stały wzdłuż długości wymiennika i niezależny od temperatury. Wartość parametru a jest nieznaną i podlegać ma estymacji na podstawie pomiarów pochodzących z ruchomego czujnika. We wzorze (10.3.48) użyto oznaczenia $q(x, t; a)$ dla podkreślenia zależności rozwiązania od nieznanego parametru a .

Oznaczmy przez $f(t)$ temperaturę medium chłodzonego w punkcie $x = 0$ (czyli na wejściu do wymiennika) w chwili t . Warunek na lewym brzegu dla naszego równania (10.3.48) ma postać $q(0, t) = f(t)$. Przyjmijmy także, że w chwili uznanej za początkową ($t = 0$), temperatura medium wewnętrznego równa jest temperaturze medium chłodzącego i wynosi zero stopni. Prowadzi to do następującego warunku

początkowego $q(x, 0) = 0$, $x > 0$ dla równania (10.3.48). Dla zgodności powyższych warunków założyć musimy także, że $f(0) = 0$ oraz $f(s) = 0$ dla $s < 0$ (zamiast s można użyć dowolnej innej litery, gdyż s jest tutaj niemym argumentem). Dla tych warunków początkowych i brzegowych rozwiązanie równania (10.3.48) ma postać

$$q(x, t; a) = \exp(-a x/v) f\left(t - \frac{x}{v}\right), \quad x > 0, \quad t > 0. \quad (10.3.49)$$

Funkcja wrażliwości odpowiedzi systemu q względem a ma zatem postać

$$g(x, t) = -\frac{x}{v} \exp(-a x/v) f(t - x/v) \quad (10.3.50)$$

Ponieważ $r = 1$, mamy tu do czynienia z prostym przypadkiem i wystarczy jedna trajektoria $x^*(t)$, aby zapewnić spełnienie warunku optymalności. Twierdzenie 10.2 pozwala znaleźć tę trajektorię w wyniku maksymalizacji $g^2(x, t)$ względem x dla każdego $t \in [0, T]$ z osobna. Powodem łatwości znalezienia rozwiązania w tym przypadku jest to, że dla $r = 1$ macierz \dot{M} redukuje się do wielkości skalarnej.

Przebieg $x^*(t)$ będzie oczywiście zależał od zmienności w czasie temperatury chłodzonego medium $f(\cdot)$. W ogólnym przypadku do spełnienia warunku optymalności potrzebować będziemy ruchomego czujnika. Będzie tak wówczas, gdy maksimum funkcji $g^2(x, t)$ względem $x \in [0, L]$ dla poszczególnych $t \in [0, T]$ osiągane jest wewnątrz przedziału $[0, L]$. Wówczas $x^*(t)$ otrzymujemy w wyniku rozwiązania równania

$$\frac{\partial}{\partial x} [g^2(x, t)] = 0$$

względem x dla każdego $t \in [0, T]$. Z analizy (10.3.50) wynika, że spełnienie tego warunku jest możliwe tylko wtedy, gdy $x^*(t) - v t = \text{const}$ dla każdego $t \in [0, T]$, ponieważ tylko wtedy jesteśmy w stanie zapewnić stałą wartość argumentu $t - x/v$.

Wyposażeni w warunki optymalności możemy łatwo rozwiązać dwa następujące problemy doboru trajektorii ruchu czujników pomiarowych.

Problem 10.2. Zadanie polega na doborze trajektorii ruchomego czujnika temperatury tak, by zmaksymalizować dokładność estymacji współczynnika przewodnictwa ciepła a w dwuwymiarowej płycie nieskończonej, opisanej równaniem

$$\frac{\partial q(x, t; a)}{\partial t} = a \Delta q(x, t; a) + \delta(t) \delta(x), \quad x \in R^2, \quad (10.3.51)$$

gdzie $q(x, t; a)$ oznacza temperaturę w punkcie $x \in R^2$ w chwili $t > 0$. Zakładamy zerową temperaturę początkową i impulsowe wymuszenie $\delta(t) \delta(x)$, gdzie δ oznacza deltę Diraca, natomiast Δ jest operatorem Laplace'a w R^2 .

Jedynym rozwiązaniem równania (10.3.51), które dąży do zera, gdy $\|x\| \rightarrow \infty$, jest funkcja

$$q(x, t; a) = \frac{1}{4\pi a t} \exp[-\|x\|^2/4at], \quad t > 0 \quad (10.3.52)$$

Wrażliwość tego rozwiązania na zmiany parametru a , czyli pochodna q względem a , ma postać

$$g(x, t) = \frac{1}{4\pi a^2 t} (1 - z) \exp(-z), \quad \text{gdzie } z \stackrel{\text{def}}{=} \|x\|^2/4at. \quad (10.3.53)$$

Niech obszarem X , po którym może poruszać się czujnik, będzie koło

$$X = \{x \in R^2 : \|x\|^2 \leq R^2\}$$

o zadanym promieniu $R > 0$ i takim, że $R^2 > 8a|\mathcal{T}|$, gdzie przez $|\mathcal{T}|$ oznaczono długość przedziału czasu, w którym dokonujemy obserwacji. Zauważmy, że warunek $R^2 > 8a|\mathcal{T}|$ powinniśmy spełnić, mimo że a nie jest znane. Wystarczy jednak wybrać R dostatecznie duże, gdyż – jak poniżej stwierdzimy – warunek ten potrzebny jest po to, by poruszający się czujnik pozostał bądź wewnątrz, bądź na brzegu koła X .

Funkcja $h^2(z) \stackrel{\text{def}}{=} (1 - z)^2 \exp[-2z]$ osiąga maksimum dla $z^* = 2$. Każda trajektoria $x^*(t)$, która dla każdego $t \in \mathcal{T}$ spełnia równanie

$$\|x^*(t)\|^2 = 8at, \quad (10.3.54)$$

jest zatem optymalną ścieżką ruchu czujnika.

Uwaga 10.3. W powyższym problemie położenie czujnika powinno oddalać się od punktu $(0, 0)$ wzdłuż promienia koła X z prędkością zależną od estymowanego parametru a , a więc dokładnie nieznaną. Z trudnością tą spotykaliśmy się już niejednokrotnie w nieliniowych problemach estymacji. W omawianej sytuacji można oczywiście zastosować wstępne oszacowanie a dla określenia prędkości oddalania się czujnika od centrum koła. Pojawia się też druga możliwość, polegająca na estymowaniu a w miarę napływania pomiarów i odpowiednim dostosowywaniu prędkości ruchu czujnika. Jeśli przez $\hat{a}(t)$ oznaczymy oszacowanie parametru a , uzyskane na podstawie pomiarów dokonanych do chwili t , to prędkość oddalania się czujnika od centrum koła powinna w chwili t wynosić $8\hat{a}(t)$.

Rozwiązanie następnego zadania opiszemy w dużym skrócie, gdyż sposób postępowania jest analogiczny do tego, którego użyliśmy w problemie 10.2.

Problem 10.3. Oznaczmy przez $q(x, t)$ odkształcenie cienkiej nieskończonej płyty w punkcie x w chwili $t > 0$, którą poddano punktowemu, impulsowemu wymuszeniu, które umieszczone było w punkcie $(0, 0)$. Równanie dla q ma postać

$$\frac{\partial^2 q(x, t; a)}{\partial t^2} = a \Delta^2 q + \delta(x) \delta(t), \quad x \in R^2. \quad (10.3.55)$$

W powyższym równaniu przez Δ^2 oznaczono operator różniczkowy czwartego rzędu uzyskany przez formalne podniesienie do kwadratu operatora Laplace'a. Zakładając, że w chwili $t = 0$ zarówno położenie, jak i prędkość odkształceń płyty na całej jej powierzchni były zerowe, jako rozwiązanie (10.3.55) otrzymamy

$$q(x, t; a) = \frac{1}{4\pi a t} \sin^2[\|x\|^2/4at]. \quad (10.3.56)$$

Wybierając jako X koło o promieniu nie mniejszym niż $2\sqrt{a|T|}$ i posługując się taką samą techniką jak w problemie 10.2, stwierdzamy, że każda trajektoria ruchu czujnika $x^*(t)$, która spełnia równanie $z^* 4at = \|x^*(t)\|^2$, jest trajektorią optymalną, jeśli jako z^* przyjmiemy najmniejszy z dodatnich pierwiastków równania $\operatorname{tg}(z) = -(1+z)$, który w przybliżeniu wynosi 1.9.

Prostota rozwiązań dwóch przedstawionych problemów uzyskana została dzięki założeniu o impulsowym i punktowym charakterze wymuszenia. Przy zastosowaniu szerszej klasy wymuszeń wystąpi jawna zależność optymalnej trajektorii czujnika od wymuszenia.

Warunki optymalności trajektorii czujników przedstawione w tym podrozdziale można uogólnić na inne niż D-optymalność kryteria planowania eksperymentu. Warunki takie podano w monografii [199], (s. 125–129).

11. Pokrewne zadania planowania

W rozdziale tym postaramy się zebrać uwagi bibliograficzne na temat tych zagadnień planowania eksperymentu, które nie znalazły się w tej monografii, a zostały w pewnym stopniu zbadane. Literatura na te i pokrewne tematy jest już na tyle bogata, że – zdaniem autora – nie sposób dokonać jej pełnego przeglądu, tym bardziej, że dziedziny te nadal intensywnie się rozwijają. Przegląd ograniczymy do przedstawienia:

- innych niż omawiane w tej monografii aspektów planowania, których celem jest estymacja funkcji,
- wybranych zagadnień planowania, które ukierunkowane są na estymację parametrów systemów dynamicznych.

11.1. Eksperyment w zadaniach estymacji – wybrane aspekty

Poniższa lista zawiera wybrane zadania planowania eksperymentu i związane z nimi problemy estymacji skończonej i z góry znanej liczby parametrów funkcji. **EKSPERYMENTY W TESTOWANIU HIPOTEZ.** Celem pierwszych prac Fishera było rozstrzygnięcie hipotez o wpływie wybranych oddziaływań na pożądane cechy roślin i zwierząt, z jednoczesnym dążeniem do eliminacji wpływu czynników ubocznych. Ten ważny nurt planowania eksperymentu nadal się rozwija (por. [39], [96], [146], [201]).

PLANOWANIE BADAŃ SYMULACYJNYCH. W wielu dziedzinach nauki i techniki badania symulacyjne zyskują status pełnoprawnych narzędzi badawczych. Jednocześnie zaobserwować można wiele przypadków marnowania czasu ludzi i komputerów na prowadzenie źle zaprojektowanych eksperymentów obliczeniowych. Wśród tych eksperymentów wyróżnić można dwie klasy:

- badania o charakterze wielokrotnie przeprowadzanych eksperymentów losowych (np. eksperymenty typu Monte Carlo),
- badania deterministycznych, lecz złożonych modeli matematycznych, nie poddających się analitycznemu rozwiązaniu (zwykle mają one postać układów nieliniowych równań o pochodnych cząstkowych).

W obliczeniach typu Monte Carlo stosuje się inne niż omawiane tutaj środki planowania badań, natomiast w symulacjach drugiej grupy stosować można odpowiednio dobrane metody planowania eksperymentu (por. [156], [203], [160]). Planowanie eksperymentów symulacyjnych ma jednak swoją specyfikę,

a mianowicie, dla tych samych danych wejściowych programu wynik powtórných obliczeń jest taki sam i dlatego nie jest celowe powtarzanie eksperymentów dla tego samego zestawu wejść.

DOBÓR SKŁADU MIESZANIN. Problem polega na właściwym dobraniu procentowego udziału składników wyrobu (betonu, stopu itp.) tak, by uzyskać jego pożądane cechy (np. twardości czy wytrzymałości). Plany specjalizowane dla mieszanin omawiane są, między innymi, w [62], [5].

ODPORNOŚĆ NA DUŻE ZAKŁÓCENIA. Problem odporności estymatorów na tzw. grube błędy pomiarów jest od dawna intensywnie badany (por. [57], [150]). Tematyce właściwego zaplanowania eksperymentu, tak by minimalizować skutki ewentualnych grubych błędów, poświęcono bardzo mało prac (por. [91]).

ODPORNOŚĆ NA NIEPRAWIDŁOWĄ SPECYFIKACJĘ MODELU. Od dość dawna zdawano sobie sprawę z ograniczającej roli założenia o poprawności modelu (por. [78], [95]), lecz pierwsze istotne wyniki odnośnie odporności planów produkcyjnych na ten czynnik uzyskano niedawno [164], [167], [165].

ROZSZERZALNOŚĆ PLANU. Przez pojęcie to rozumiemy taką konstrukcję planu, która pozwala dodawać do niego dodatkowe punkty w taki sposób, by nowy plan, eksplorujący szerszy obszar, był nadal planem optymalnym. Konstrukcję takich planów zaproponowano w [140].

Ponadto, wiele innych aspektów planowania eksperymentu omówiono w następujących monografiach, pracach przeglądowych i artykułach [39], [7], [3], [4], [23], [63], [146], [214], [87] [88].

Wszystkie omawiane dotąd problemy planowania eksperymentu dotyczyły modeli funkcji regresji, które dało się (z założenia) opisać za pomocą skończonej liczby nieznaných parametrów i skończonej liczby znanych funkcji. Liczba tych funkcji była albo z góry ustalona, albo – jak w przypadku zadania doboru struktury regresji – była ograniczona od góry. Tego typu problemy nazywa się zadaniami estymacji parametrycznej, w odróżnieniu od nieparametrycznych zagadnień estymacji, które zamierzamy teraz krótko omówić. Warto zaznaczyć, że przymiotnikiem *nieparametryczny* określa się w statystyce kilka różnych zagadnień. Tutaj skupimy się na zagadnieniach planowania eksperymentu dla nieparametrycznej estymacji funkcji regresji. Użyty w tym kontekście termin *nieparametryczna estymacja* oznacza, że nieznaną funkcję regresji jest elementem pewnej nieskończonej wymiarowej przestrzeni funkcji, na przykład przestrzeni funkcji, których druga potęga jest całkowalna na pewnym zbiorze X z przestrzeni R^s . Przestrzeń taką oznaczamy będziemy przez $L_2(X)$. W tak bogatej klasie funkcji nie da się zbudować metody estymacji, która „zrekonstruuje” nieznaną funkcję na podstawie skończonej liczby pomiarów. Dlatego metody nieparametrycznej estymacji są w istocie sekwencją estymatorów, które konstruowane są tak, by coraz dokładniej przybliżać nieznaną funkcję, gdy liczba pomiarów rośnie do nieskończoności. Stosuje się różne miary dokładności owego przybliżenia. Przykładowo może nią być tzw.

błąd średniokwadratowy, czyli wartość oczekiwana kwadratu normy (w $L_2(X)$) różnicy między nieznaną funkcją a jej estymatorem.

W głównym nurcie badań nieparametrycznej estymacji funkcji regresji zakłada się, że zmienne niezależne (wejścia) są realizacjami pewnych zmiennych losowych (aktualną i obszerną bibliografię na ten temat zawiera monografia [49]). Innymi słowy, statystyk jest biernym obserwatorem wejść i wyjść badanego procesu, co nie pozostawia miejsca na planowanie eksperymentu. Dlatego w dalszym krótkim przeglądzie koncentrować się będziemy na tych pracach, w których dopuszcza się wpływ eksperymentatora na dobór wartości zmiennych niezależnych. Bibliografię i omówienie publikacji tego nurtu zawiera monografia [31]. W większości tych prac zakłada się, że zmienne niezależne są rozmieszczone w równo-odległych punktach odcinka lub kostki wielowymiarowej. O ile autorowi wiadomo [90] jest pierwszą pracą, w której optymalizowano dobór wartości zmiennych niezależnych. W pracy tej badano przypadek jednej zmiennej niezależnej, a jako estymatora funkcji regresji używano estymatora z jądrem.

Bliższy tematyce tej książki jest nurt prac, w których badano wpływ planu eksperymentu na tempo zbieżności błędu średniokwadratowego nieparametrycznych estymatorów funkcji regresji opartych na rozwinięciach ortogonalnych. Istotą tej grupy metod jest estymacja współczynników rozwinięcia nieznannej funkcji w wybrany szereg ortogonalny z jednoczesnym doбором tempa wzrostu liczby współczynników szeregu w zależności od liczby pomiarów. Różne aspekty tej klasy metod badano, między innymi, w [43], [151], [44], [32], [152], a w pracach [114] i [115] wykazano, że metoda najmniejszych kwadratów ma własności algorytmów nieparametrycznych, jeśli rozmiar podprzestrzeni funkcji rozpinających funkcję regresji¹ rośnie odpowiednio wolno wraz z narastaniem liczby pomiarów.

Liczba prac poświęconych zagadnieniom planowania eksperymentu, ukierunkowanego na późniejsze stosowanie nieparametrycznych metod estymacji nie jest zbyt duża. W [114] wykazano, że ze względu na szybkość zbieżności metody rozwinięć ortogonalnych korzystne jest rozmieszczanie pomiarów w węzłach kwadratur, odpowiednio dobranych do problemu. W pracy [141] zaproponowano, by punkty planu eksperymentu dla nieparametrycznej estymacji addytywnej funkcji regresji o s zmiennych wejściowych rozmieszczać zgodnie ze wzorem

$$x_n = [\text{frac}(n\theta_1), \text{frac}(n\theta_2), \dots, \text{frac}(n\theta_s)], \quad n = 1, 2, \dots, \quad (11.1.1)$$

gdzie $\text{frac}(\cdot)$ oznacza część ułamkową liczby podanej w nawiasach, natomiast θ_j , $j = 1, 2, \dots, s$ są – odpowiednio dobranymi – liczbami niewymiernymi. Ciąg (11.1.1) jest ciągiem deterministycznym, ale jednocześnie, ma on wiele cech rozkładu równomiernego w kostce s -wymiarowej. W pracach [141] i [142] wykazano, że – przy dodatkowych założeniach dotyczących gładkości estymowanej funkcji

¹ Podprzestrzenie te wybiera się z systemu funkcji ortogonalnych i zupełnych w $L_2(X)$.

– estymator rozwinięć ortogonalnych uzyskuje najszybszą z możliwych do osiągnięcia szybkość malenia do zera błędu średniokwadratowego. Istnienie dolnej granicy szybkości zbieżności wykazano w [184].

11.2. Składniki eksperymentu w estymacji systemów dynamicznych

Liczba elementów, które można wybrać przed eksperymentem, który ma na celu estymację parametrów systemów dynamicznych jest znacznie większa niż w przypadku estymacji funkcji i obejmuje następujące aspekty.

Próbkowanie i rekonstrukcja sygnałów. Podstawowym rezultatem w zakresie doboru częstotliwości próbkowania jest twierdzenie Shannona. Rekonstrukcja sygnałów na podstawie odpowiednio często pobieranych próbek należy do klasycznych zagadnień przetwarzania sygnałów i dlatego nie będziemy go tu szerzej omawiać. Ostatnio uzyskano wyniki dotyczące rekonstrukcji sygnałów o ograniczonym widmie, próbkowanych w obecności zakłóceń o charakterze szumu białego [99], [100] oraz filtracji takich pomiarów [101].

Dobór wymuszeń. W przypadku systemów dynamicznych wymuszenie (sygnał wejściowy) może być funkcją czasu i/lub zmiennych przestrzennych.

Rozmieszczenie czujników pomiarowych. Jeśli opis matematyczny systemu dynamicznego zależy od zmiennych przestrzennych, to celowe jest rozważenie, gdzie rozmieścić czujniki pomiarowe.

Odnosniki do literatury na dwa ostatnie tematy przytaczamy w następnych podrozdziałach.

Systemy dynamiczne o skończonej liczbie stopni swobody

Przez pojęcie systemów dynamicznych o skończonej liczbie stopni swobody rozumiemy tutaj takie przybliżone opisy procesów, które modelować można za pomocą układów równań różniczkowych zwyczajnych.

Dobór sterowań w estymacji systemów dynamicznych doceniany był od lat siedemdziesiątych XX wieku (por. [84]). Propagowano wówczas szum biały jako sygnał pobudzający dla systemów liniowych.

Kilka lat później sformułowano zadania optymalnego doboru sygnałów wymuszających w dziedzinie częstotliwości. Kryteria optymalności formułowano na podstawie różnych miar dokładności ocen parametrów. Zwykle nakładano ograniczenia na średnią moc sygnału. Wyniki na ten temat zebrano w [42], [213] i [189]. W [123] pokazano, że podobne wyniki uzyskać można, nakładając słabsze ograniczenia na tempo narastania sygnałów.

Rezultaty zawarte w cytowanych pracach w istotny sposób wykorzystywały liniowość systemu dynamicznego. W ostatnich kilkunastu latach rośnie liczba

prac na temat parametrycznej i nieparametrycznej identyfikacji nieliniowych systemów dynamicznych (por. [46], [47], [48], [53], [54]). W cytowanych pracach uzyskano teoretyczne własności estymatorów, które dają podstawy do oceny ich dokładności. Mogą być one punktem wyjścia do podjęcia prób doboru sterowań optymalnych z punktu widzenia dokładności estymacji.

Systemy o nieskończonej liczbie stopni swobody

Dokładniejszy opis procesów otrzymać można, jeśli uwzględnimy przebieg zjawisk nie tylko w czasie, ale i w przestrzeni. Systemy dynamiczne tego typu opisuje się za pomocą układów równań o pochodnych cząstkowych. Równania takie uzyskuje się zwykle w wyniku stosowania praw uznanych za obowiązujące w danej dziedzinie nauki. Dlatego przyjmować będziemy, że nieznanymi wielkościami są jedynie parametry tych równań.

Różne aspekty doboru wymuszeń w identyfikacji wartości własnych lub parametrów tej klasy systemów badano w pracach [117], [121], [124], [134], [126], [135], [125], [128], [131], [132].

Praca [85] zawiera przegląd wczesnych podejść do problemu rozmieszczenia czujników pomiarowych. Charakteryzowały się one dążeniem do znajdowania takich rozmieszczeń czujników, które ułatwiałyby obliczeniowe aspekty identyfikacji parametrów.

W pracach [118], [122], [193], [73], [195] dobierano rozmieszczenie czujników tak, by maksymalizować dokładność estymacji wartości własnych operatora opisującego system, jego parametrów lub stanu systemu.

Kolejnym etapem badań nad zwiększaniem dokładności estymacji parametrów systemów o nieskończonej liczbie stopni swobody było opracowanie algorytmów doboru trajektorii ruchomych czujników (por. [127], [198], [196], [197], [194]). Wyniki tego nurtu badań zebrano i znacznie poszerzono w wydanej ostatnio monografii [199], która zawiera także kompletną bibliografię prac o tej tematyce.

CZĘŚĆ IV

Eksperyment w diagnostyce procesów

12. Próbkowanie funkcyjnych charakterystyk wyrobów

W wielu dziedzinach produkcji poziom wymagań technologicznych i/lub wymagań klienta jest na tyle wysoki, że dla każdego egzemplarza wyrobu mierzy się jego indywidualną charakterystykę, która następnie porównywana jest z nominalnym przebiegiem charakterystyki idealnego wyrobu.

Powszechnie znanym przykładem takiego postępowania są indywidualne częstotliwościowe charakterystyki słuchawek i zestawów głośnikowych sprzętu hi-fi. Podobne procedury stosują producenci silników elektrycznych wyższych mocy.

Spodziewać się można, że tego rodzaju postępowanie będzie się upowszechniać wśród producentów innych wyrobów. Dlatego też, ważne jest opracowanie metod próbkowania, która pozwoli na osiągnięcie maksymalnej dokładności porównania charakterystyki wyrobu z charakterystyką nominalną, przy zachowaniu zadanej z góry liczby pomiarów.

W rozdziale tym zaproponowana zostanie metoda planowania eksperymentów z zakreślonym wyżej celem. Posługiwać się przy tym będziemy ciągłymi planami skupionymi w skończonej liczbie punktów, a rozdzielanie całkowitej liczby pomiarów pomiędzy poszczególne punkty pomiarowe następować będzie w sposób opisany w części I.

12.1. Ocena jakości funkcyjnej charakterystyki wyrobu

Pod pojęciem charakterystyki wyrobu rozumiemy funkcję $\eta : R^s \rightarrow R$, która dla danego zestawu wejść (pobudzeń) $x \in R^s$ podaje liczbową reakcję (wyjście, odpowiedź) badanego wyrobu.

W przykładzie ze wstępu do tego rozdziału charakterystyka głośnika to funkcja, której argumentem jest częstotliwość sygnału sinusoidalnego podanego na głośnik, a wartość charakterystyki to tłumienie sygnału o tej częstotliwości używane na wyjściu głośnika (zwyczajowo mierzone w decybelach).

Będziemy zakładać, że nominalna (idealna) charakterystyka wyrobu jest znana i z góry zadana. Oznaczmy charakterystykę nominalną przez $\eta_{id}(x)$ i załóżmy, że potrafimy ją dostatecznie dokładnie aproksymować skończonym szeregiem znanych (wybranych przez nas) funkcji $v(x) = [v_1(x), v_2(x), \dots, v_r(x)]^T$. Założenie to oznacza, że istnieje zestaw współczynników

$$a_{id} \stackrel{\text{def}}{=} [a_{id}^{(1)}, a_{id}^{(2)}, \dots, a_{id}^{(r)}]^T$$

zapewniających spełnienie równości

$$\eta_{id}(x) = \sum_{k=1}^r a_{id}^{(k)} v_k(x) = a_{id}^T v(x). \quad (12.1.1)$$

W odniesieniu do charakterystyki badanego egzemplarza wyrobu zakładać będziemy, że ma ona taką samą postać funkcyjną jak charakterystyka nominalna, lecz współczynniki rozwinięcia jej w szereg $v(x) = [v_1(x), v_2(x), \dots, v_r(x)]^T$ mogą być inne niż te w charakterystyce nominalnej.

Oznaczmy przez $a_0 = [a_0^{(1)}, a_0^{(2)}, \dots, a_0^{(r)}]^T$ wartości współczynników charakterystyki badanego egzemplarza wyrobu. Współczynniki te są nam nieznane i podlegać będą szacowaniu na podstawie pomiarów opisanych w dalszej części rozdziału.

Zgodnie z wymienionymi założeniami, charakterystyka badanego egzemplarza ma postać

$$\eta_0(x) = \sum_{k=1}^r a_0^{(k)} v_k(x) = a_0^T v(x). \quad (12.1.2)$$

Za miarę odległości między charakterystyką bieżącego egzemplarza a nominalną przyjąć można dowolną z metryk w przestrzeniach funkcyjnych, o ile jej interpretacja odpowiada intuicyjnym wymaganiom porównywania charakterystyk danego wyrobu. Tutaj posłużymy się błędem średniokwadratowym, gdyż pozwala on na stosunkowo łatwe przeliczenie odległości współczynników a_{id} i a_0 na odległość między samymi charakterystykami. Przyjmujemy zatem, że odległość ta dana jest wzorem

$$\rho(\eta_{id}, \eta_0) \stackrel{\text{def}}{=} \int_Z w(x) (\eta_{id}(x) - \eta_0(x))^2 dx, \quad (12.1.3)$$

gdzie $Z \subset R^s$ jest obszarem w przestrzeni zmiennych wejściowych, na którym porównujemy charakterystyki, natomiast $w(x) \geq 0$ jest zadaną funkcją wagową, która wskazuje, gdzie w obszarze Z bliskość charakterystyk powinna być większa. Przyjęcie $w(x) \equiv 1$ oznacza, że przywiązujemy jednakowe znaczenie do bliskości charakterystyk w całym obszarze Z .

Lemat 12.1. *Jeśli porównywane charakterystyki mają postać (12.1.1) i (12.1.2), to*

$$\rho(\eta_{id}, \eta_0) = \text{tr} \left[(a_{id} - a_0) \cdot (a_{id} - a_0)^T W \right], \quad (12.1.4)$$

gdzie macierz W zdefiniowana jest następująco

$$W = \int_Z w(x) v(x) \cdot v^T(x) dx. \quad (12.1.5)$$

Jeśli dodatkowo funkcje rozpinające charakterystykę są ortonormalne z wagą w , to znaczy,

$$\int_Z w(x) v_k(x) v_j(x) dx = \begin{cases} 1, & \text{gdy } k = j; \\ 0, & \text{gdy } k \neq j, \end{cases} \quad (12.1.6)$$

to wówczas

$$\rho(\eta_{id}, \eta_0) = \text{tr} \left[(a_{id} - a_0) \cdot (a_{id} - a_0)^T \right]. \quad (12.1.7)$$

Dowód. Dowód lematu jest czysto algebraiczny, dlatego przedstawimy tylko jego ważniejsze kroki.

$$\begin{aligned} \rho(\eta_{id}, \eta_0) &= \int_Z w(x) v^T(x) (a_{id} - a_0) \cdot (a_{id} - a_0)^T v(x) dx \\ &= \int_Z \text{tr} [w(x) v^T(x) (a_{id} - a_0) \cdot (a_{id} - a_0)^T v(x)] dx \\ &= \int_Z \text{tr} [(a_{id} - a_0) \cdot (a_{id} - a_0)^T v(x) \cdot v^T(x) w(x)] dx \\ &= \text{tr} \left\{ [(a_{id} - a_0) \cdot (a_{id} - a_0)^T] \int_Z [v(x) \cdot v^T(x) w(x)] dx \right\}. \end{aligned}$$

Powyżej skorzystaliśmy z tego, że $\text{tr}[AB] = \text{tr}[BA]$ oraz liniowości śladu macierzy. Druga teza lematu wynika natychmiast ze spostrzeżenia, że gdy spełnione jest (12.1.6), to macierz W jest macierzą jednostkową. •

12.2. Dobór planu i pomiary charakterystyki

Zakładamy, że pomiary charakterystyki badanego wyrobu dokonywane są zgodnie z klasycznym modelem addytywnych, nieskorelowanych zakłóceń losowych

$$y_i = \eta_0(x_i) + \epsilon_i, \quad i = 1, 2, \dots, n, \quad (12.2.8)$$

gdzie x_i , to wartości wymuszeń podawanych na badany wyrób, ϵ_i zakłócenia pomiarowe, $E(\epsilon_i) = 0$, $E(\epsilon_i^2) = \sigma^2 < \infty$.

Do estymacji wektora a_0 używamy metody najmniejszych kwadratów. Oznaczmy tak uzyskany estymator przez \hat{a} . Jak wiemy, \hat{a} jest nieobciążonym estymatorem a_0 , a zatem $E \hat{a} = a_0$. Wówczas nieobciążonym estymatorem dla $\eta_0(x)$ jest funkcja

$$\hat{\eta}(x) \stackrel{\text{def}}{=} \hat{a}^T v(x). \quad (12.2.9)$$

Postawmy hipotezę, że charakterystyka aktualnie badanego wyrobu $\eta_0(x)$ jest dokładnie równa charakterystyce idealnej $\eta_{id}(x)$. Zakładając poprawność tej hipotezy, do oceny jakości aktualnej charakterystyki użyć możemy $\eta_0(x)$ zamiast $\eta_{id}(x)$. Innymi słowy, interesuje nas odległość $\rho(\eta_0, \hat{\eta})$ powtarzając kroki dowodowe Lematu 12.1 otrzymamy

$$\rho(\eta_0, \hat{\eta}) = \text{tr}[(\hat{a} - a_0) \cdot (\hat{a} - a_0)^T W] = (\hat{a} - a_{id})^T W (\hat{a} - a_{id}). \quad (12.2.10)$$

Uzasadnienie ostatniej równości podamy we Wniosku 12.1. Wyrażenie (12.2.10) może służyć jako statystyka testowa, gdyż jeśli nasza hipoteza jest spełniona, to z większym prawdopodobieństwem spodziewać się można, że wartość $\rho(\eta_0, \hat{\eta})$ będzie mała. Stwierdzeniu temu można nadać precyzyjny sens, jeśli założymy, że rozkład zakłóceń jest rozkładem normalnym, co pozwala na określenie rozkładu zmiennej losowej $\rho(\eta_0, \hat{\eta})$, przy założeniu poprawności hipotezy zerowej. Jeśli wariancja zakłóceń jest znana, to wartości krytyczne dla $\rho(\eta_0, \hat{\eta})$ odczytujemy z tablic rozkładu χ^2 . Szczegółowego omówienia tego testu nie podajemy, gdyż naszym celem jest zaplanowanie eksperymentu, który zapewni statystyczne zmniejszenie wartości $\rho(\eta_0, \hat{\eta})$, wówczas gdy rzeczywiście $\eta_{id}(x) = \eta_0(x)$.

Liniowa niezależność składowych $v(x)$ prowadzi do następującego wniosku.

Wniosek 12.1. *Jeśli każdego $x \in X$ $\eta_{id}(x) = \eta_0(x)$, to także $a_{id} = a_0$, a więc w przestrzeni parametrów posługiwać się możemy odległością wektora \hat{a} od a_0 .*

Zauważmy, że zamiast a_0 możemy użyć $E(\hat{a})$, gdyż zachodzi $E(\hat{a}) = a_0$. Korzystając z tych spostrzeżeń i powtarzając kroki dowodowe Lematu 12.1 otrzymamy

$$\rho(\eta_0, \hat{\eta}) = \text{tr}[(\hat{a} - E(\hat{a})) \cdot (\hat{a} - E(\hat{a}))^T W]. \quad (12.2.11)$$

Wyrażenie to jeszcze nie może służyć jako kryterium planowania eksperymentu, gdyż \hat{a} jest wektorem losowym. Możemy natomiast posłużyć się wartością oczekiwaną $\rho(\eta_0, \hat{\eta})$, co prowadzi do następującego wskaźnika jakości planu

$$E[\rho(\eta_0, \hat{\eta})] = \text{tr} \left[E \left[(\hat{a} - E(\hat{a})) \cdot (\hat{a} - E(\hat{a}))^T \right] W \right] = \sigma^2 \text{tr}[M_n^{-1} W]. \quad (12.2.12)$$

M_n oznacza macierz informacyjną dla planu x_1, x_2, \dots, x_n . Ostatnia równość wynika ze spostrzeżenia, że $E \left[(\hat{a} - E(\hat{a})) \cdot (\hat{a} - E(\hat{a}))^T \right]$ jest macierzą kowariancji estymatora MNK, która jest równa $\sigma^2 M_n^{-1}$ (por. część I).

Alternatywnym wobec kryterium (12.2.12) może być maksyminowe kryterium T- optymalności rozważane w pracy [200].

Wniosek 12.2. *Przy założeniach poczynionych w tym rozdziale, problem planowania eksperymentu dla testowania zgodności charakterystyki wyrobu z charakterystką nominalną sprowadzić można do zadania planowania L- optymalnego z kryterium $\text{tr}[M_n^{-1} W]$. W przypadku, gdy $v(x)$ spełnia warunek (12.1.6), zadanie to sprowadza się do problemu planowania A- optymalnego.*

By móc efektywnie korzystać z tego wniosku, warto przejść do minimalizacji kryterium L-optymalności w klasie planów ciągłych, skupionych w skończonej liczbie punktów.

Załóżmy, że naszym celem jest zbieranie danych o jakości poszczególnych egzemplarzy badanego wyrobu, wówczas powyższe rezultaty i uwagi prowadzą do następującego algorytmu.

ALGORYTM AKWIZYCJI DANYCH O JAKOŚCI CHARAKTERYSTYK FUNKCYJNYCH

- Krok 1.** Ustalić pożądaną (idealną) charakterystykę wyrobu $\eta_{id}(x)$, $x \in X$.
- Krok 2.** Wybrać wektor funkcji $v(x)$, który pozwala dostatecznie dokładnie aproksymować $\eta_{id}(x)$ i obliczyć macierz W oraz wektor a_{id} współczynników rozwinięcia $\eta_{id}(x)$ w bazie $v(x)$.
- Krok 3.** Znaleźć (analitycznie lub numerycznie) plan L-optymalny $\epsilon^* \in \Xi(X)$, który minimalizuje $\text{tr}[M^{-1}(\epsilon)W]$. Wybrać liczbę pomiarów charakterystyki każdego egzemplarza wyrobu N i przybliżyć plan ϵ^* odpowiednim planem dyskretnym (zarys algorytmu podano na s. 28). Kroki te dają w wyniku pary (x_i^*, n_i^*) , $i = 1, 2, \dots, m$, gdzie x_i^* wskazują punkty próbkowania charakterystyki wyrobu, a n_i^* – liczbę próbek pobranych w każdym z nich, $\sum_{i=1}^m n_i^* = N$.
- Krok 4.** Dla danego egzemplarza wyrobu, o numerze, powiedzmy j , zmierzyć wartości charakterystyki w punktach x_i^* z krotnością powtórzeń n_i^* i stosując MNK, obliczyć oszacowanie parametrów \hat{a} .
- Krok 5.** Korzystając z (12.2.11), obliczyć oszacowanie ϱ_j odległości charakterystyki badanego egzemplarza od charakterystyki idealnej.
- Krok 6.** Powtarzać krok 4 i 5 dla każdego egzemplarza, uzyskując ciąg wartości ϱ_j , $j = 1, 2, \dots$

W następnych rozdziałach ciąg wartości ϱ_j , $j = 1, 2, \dots$ stanowić będzie podstawę do oceny, czy nie nastąpiło rozregulowanie procesu wytwórczego.

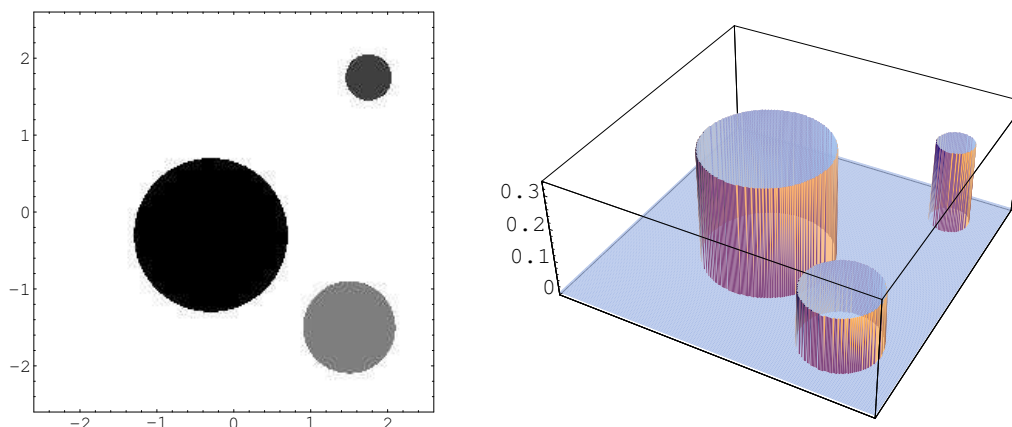
13. Próbkowanie obrazów do celów diagnostycznych

Sekwencje obrazów dostarczane przez kamery monitorujące przebieg procesów wytwórczych coraz częściej stają się źródłem informacji o jakości aktualnej produkcji. Nadal jednak zastosowania kamer w przemyśle są zbyt rzadkie jak na możliwości techniczne i diagnostyczne, które można uzyskać w wyniku ich stosowania. Można się jednak spodziewać wzrostu zainteresowania tymi technikami, gdyż w ostatnich latach ceny kamer maleją, a ich zdolności rozdzielcze znacznie się zwiększyły.

Wydaje się, że jednym z powodów wciąż zbyt małego korzystania z kamer jest paradoksalnie nadmiar informacji płynącej z kamery w krótkim czasie. Aby informacja taka mogła być użyteczna, musi być przetworzona „na bieżąco” w dane użyteczne do celów diagnostycznych. Tempo bieżącego przetwarzania zależy oczywiście od prędkości procesu, który chcemy monitorować. W wielu przypadkach (np. w procesie ciągłego odlewania miedzi) produkcja przebiega tak szybko, że na przetwarzanie poszczególnych klatek obrazu zostają ułamki sekundy. Uzyskanie takiej prędkości przetwarzania obrazów wymaga szybkich algorytmów ekstrakcji informacji istotnej w ocenie jakości danego procesu, gdyż nawet współczesne procesory o prędkości taktowania rzędu 2–3 GHz nie są w stanie w pełni kontrolować strumienia danych, napływającego z prędkością kilkudziesięciu klatek na sekundę, gdy objętość każdej takiej klatki jest rzędu kilku milionów bajtów.

Z wymienionych względów w rozdziale tym proponujemy dwa różne podejścia do próbkowania obrazów, które są dostosowane do zadań diagnozowania powierzchni powstających w trakcie różnych procesów wytwórczych. W odróżnieniu od klasycznych podejść do próbkowania obrazów, które starają się zachować informacje o całym obrazie, tutaj świadomie redukujemy ją tylko do cech istotnych do oceny szeroko rozumianej równomierności powierzchni. Proponowane podejścia różnią się sposobem traktowania powierzchni i nakładami obliczeniowymi.

Podobnie jak w poprzednich częściach tej monografii, oba proponowane podejścia bazują na technikach planowania eksperymentu i statystycznego przetwarzania jego wyników, lecz używać będziemy terminu *próbkowanie* obrazów, gdyż ma on w tej dziedzinie swoje tradycje. Warto też zaznaczyć, że diagnostyka stanu procesu jest tutaj rozumiana w dość wąskim zakresie i sprowadza się do stwierdzenia, czy proces przebiega poprawnie, czy też nie, bez próby tworzenia szczegółowego opisu matematycznego całego zjawiska. Odsyłamy Czytelnika do wydanych ostatnio monografii [74], [75], w których zagadnienia diagnostyki traktowane są znacznie szerzej.



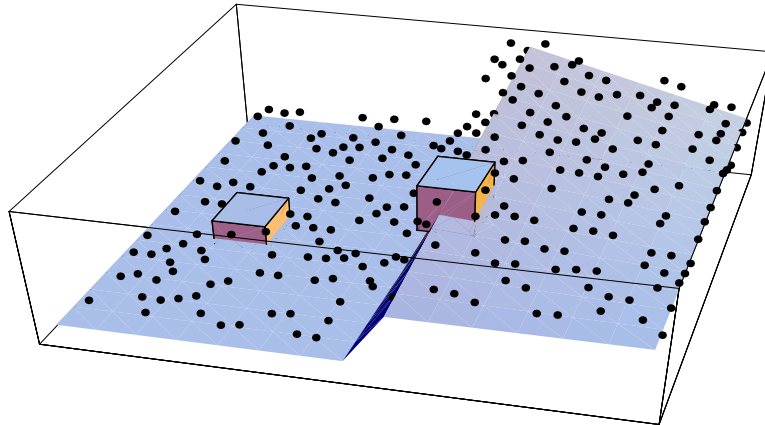
Rys. 13.1. Przykładowy obraz (po lewej) i jego reprezentacja w postaci funkcji dwóch zmiennych (po prawej)

13.1. Próbkowanie zmienności obrazu. Szybki algorytm okonturowywania

W podrozdziale tym zaproponowano algorytm próbkowania do oceny równomierności powierzchni przedstawionej na obrazie cyfrowym. Przyjmujemy, że obraz reprezentowany jest w ustalonej skali szarości. W trakcie próbkowania bierze się pod uwagę zmienność poziomów szarości – można powiedzieć, że jest ono sterowane zawartością obrazu. Dokładniej, algorytm bazuje na okonturowaniu fragmentów obrazu, które odróżniają się od równomiernego tła.

Powodem wyboru okonturowywania jest konieczność, chociaż częściowego, uniezależnienia się od zmian oświetlenia obrazu, które mogą być wywołane czynnikami zewnętrznymi, nie wpływającymi na proces produkcyjny, na przykład, zmiennością pór dnia lub zachmurzeniem. Takiego uniezależnienia nie gwarantują metody oceny równomierności powierzchni korzystające tylko z wartości poziomów szarości poszczególnych pikseli. Powody, dla których proponujemy modyfikację gradientowego algorytmu okonturowywania, są następujące:

1. Znane metody okonturowywania (por. [59], [41]) są funkcjami dyskretnych wersji gradientu lub laplasjanu. Jeśli nie zastosuje się odpowiednich metod redukcji szumów, to zastosowanie ich do obrazu ze znacznymi błędami prowadzi do ich wzmocnienia i wykrywania „fałszywych” krawędzi.
2. Jeśli, tak jak w przypadku metody LoG (*Laplacian of Gaussian*), przed okonturowaniem zastosuje się filtr, to wydłuża się znacznie czas przetwarzania obrazu i metodę trudno stosować do przetwarzania sekwencji obrazów na bieżąco.



Rys. 13.2. Idea algorytmu próbkowania krawędzi – trójwymiarowa reprezentacja oryginalnego obrazu

W tym kontekście, proponowana metoda ma dwie zalety.

- Może być zrealizowana w arytmetyce całkowitoliczbowej i wymaga jedynie operacji odejmowania, porównywania i zliczania.
- Ma „wbudowaną” odporność na zakłócenia, gdyż decyzję o tym czy dany piksel zaliczyć do brzegu obszaru decyduje tylko zliczanie i porównywanie z progami.

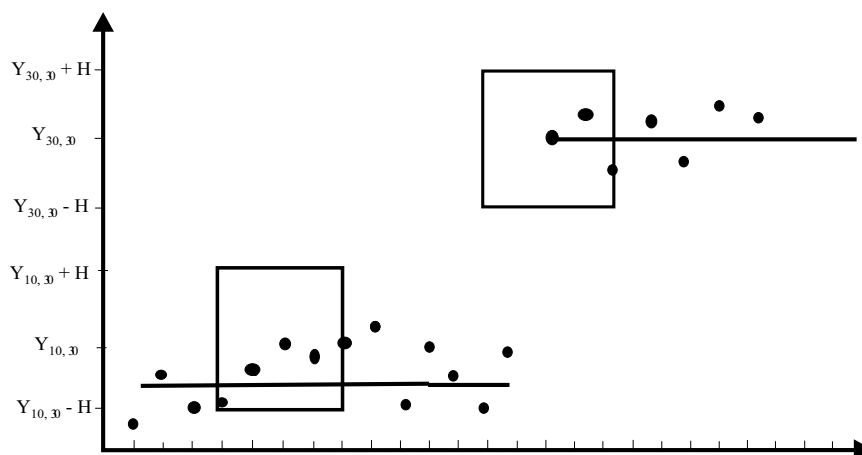
W pracach [21] i [149] znaleźć można dyskusje na temat zaawansowanych metod wykrywania krawędzi i odnośniki do wcześniejszych publikacji. Proponowane podejście inspirowane jest przez ideę wertykalnie ważonej regresji (por. [103]), chociaż w opisie metody pojęcia tego nie będziemy używać.

Opis algorytmu

Rozważmy obraz o rozmiarach $n_x \times n_y$. Niech y_{ij} , $i = 1, 2, \dots, n_x$, $j = 1, 2, \dots, n_y$ oznacza poziom szarości piksela (i, j) . Będziemy zakładać, że y_{ij} zostały przeskalowane do przedziału $[0, 1]$, chociaż w praktyce implementacja algorytmu jest szybsza, jeśli realizowana jest w arytmetyce całkowitoliczbowej o zakresie $[0, 255]$.

Uwaga 13.1. *Poziom szarości y_{ij} może być odzwierciedleniem „czystego” obrazu lub $y_{ij} = f_{ij} + \epsilon_{ij}$, gdzie f_{ij} oznacza poziom szarości nieobserwowanego, idealnego obrazu, na który oddziałują losowe zakłócenia ϵ_{ij} .*

Istotnym elementem proponowanej metody jest prostopadłościan o rozmiarach $(2h_x + 1) \times (2h_y + 1) \times 2H$, gdzie $0 < h_x < n_x$ i $0 < h_y < n_y$ oznaczają połowy rozmiarów okna położonego w płaszczyźnie obrazu, natomiast $0 < H < 1$ jest wysokością prostopadłościanu. Jego centrum znajduje się na wysokości poziomu



Rys. 13.3. Idea algorytmu próbkowania krawędzi — przekrój

szerości y_{ij} piksela, o którym mamy zdecydować, czy ma stać się próbką konturu. Zawartość $\mathcal{B}_{ij}(h_x, h_y, H)$ tego prostopadłościanu definiujemy następująco

$$\mathcal{B}_{ij}(h_x, h_y, H) \stackrel{\text{def}}{=} \left\{ (i+k, j+l, z) : |z - y_{i,j}| \leq H, k \in \overline{[-h_x, h_x]}, \right. \quad (13.1.1)$$

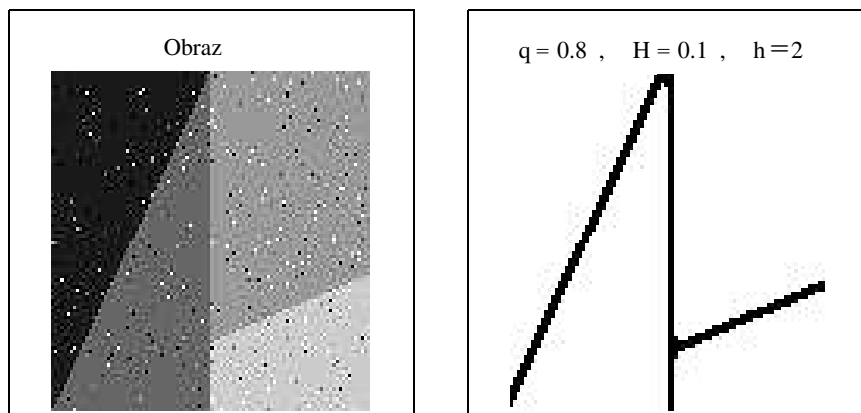
$$\left. l \in \overline{[-h_y, h_y]}, z \in [0, 1] \right\}.$$

W celu skrócenia wzorów, użyto następującej notacji: dla k będącego liczbą całkowitą, zapis $k \in \overline{[-h_x, h_x]}$ oznacza, że k przyjmuje wszystkie wartości całkowite z przedziału $[-h_x, -h_x + 1, \dots, 0, 1, \dots, h_x]$.

Uwaga 13.2. Prostopadłościan $\mathcal{B}_{ij}(h_x, h_y, H)$ ma centrum w punkcie o współrzędnych (i, j, y_{ij}) . Jego elementami są współrzędne i i poziomy szerokości tych pikseli, które

1. odległe są od piksela (i, j) o mniej niż h_x i h_y pikseli,
2. a ich poziom szerokości różni się od poziomu szerokości piksela centralnego y_{ij} o mniej niż H .

Zawartość tego prostopadłościanu zinterpretować można także w terminach różnicowej aproksymacji pochodnych w kierunku w następujący sposób. Z piksela (i, j) szacujemy pochodne $y_{i+k, j+l} - y_{i,j}$ w kierunku wszystkich sąsiednich pikseli $|k| < h_x$, $|l| < h_y$. Do „pudełka” $\mathcal{B}_{ij}(h_x, h_y, H)$ zaliczamy tylko te trójki $(i+k, j+l, y_{i+k, j+l})$, w kierunku, których oszacowania pochodnych $|y_{i+k, j+l} - y_{i,j}|$ nie przekraczają H . Innymi słowy, prostopadłościan ten zawiera tylko te piksele sąsiednie, których poziom szerokości mało różni się od wartości piksela centralnego.



Rys. 13.4. Próbkowanie zmian w obrazie w obecności zakłóceń typu „sól z pieprzem”

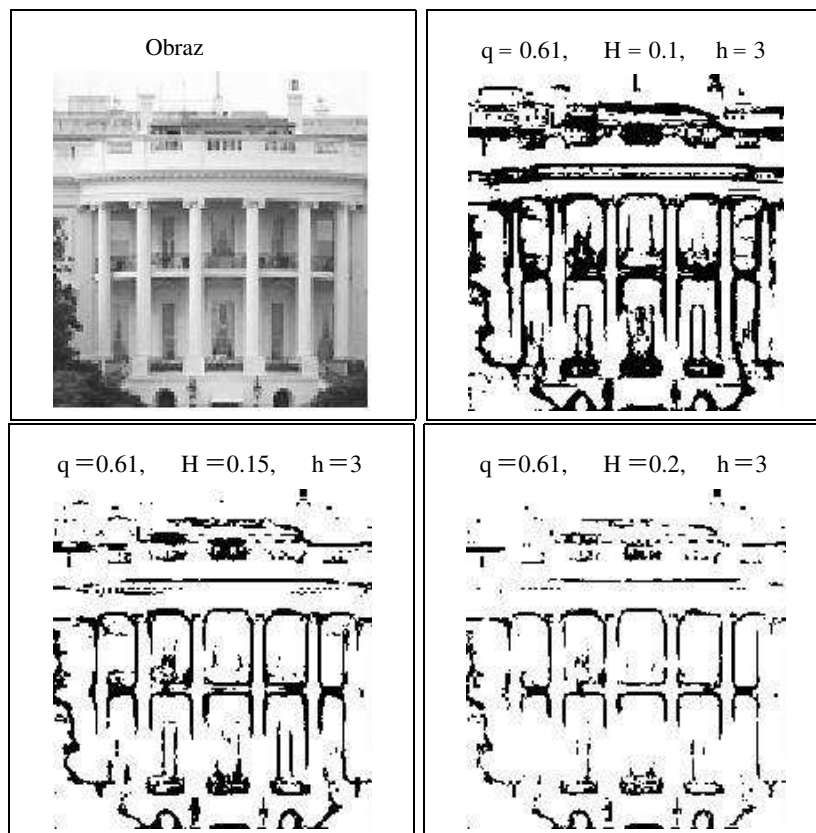
Oznaczmy przez $b_{ij}(h_x, h_y, H)$ liczbę elementów zawartych w prostopadłościanie $\mathcal{B}_{ij}(h_x, h_y, H)$. Zauważmy, że zachodzą następujące nierówności

$$1 \leq b_{ij}(h_x, h_y, H) \leq L \stackrel{\text{def}}{=} (2h_x + 1)(2h_y + 1), \quad (13.1.2)$$

ponieważ $\mathcal{B}_{ij}(h_x, h_y, H)$ zawsze zawiera (i, j, y_{ij}) i nie może zawierać więcej pikseli, niż jest ich w oknie wokół elementu centralnego. Poniższe rozważania stanowią intuicyjne przesłanki proponowanej metody. Załóżmy, że wybraliśmy h_x i h_y rzędu kilku (3 – 7) pikseli i H rzędu 0.1 – 0.2. Rozważmy zawartość prostopadłościanu $\mathcal{B}_{ij}(h_x, h_y, H)$ w dwóch sytuacjach.

1. Przyjmijmy, że piksel (i, j) znajduje się w części obrazu, w której zmienność poziomów szarości nie jest zbyt duża, tzn., $y_{i+k, j+l}$, $k \in [\overline{h_x}, \overline{h_x}]$, $l \in [\overline{-h_y}, \overline{h_y}]$ różnią się między sobą od kilku do kilkunastu procent. Wówczas, $\mathcal{B}_{ij}(h_x, h_y, H)$ zawiera prawie wszystkie $(i+k, j+l, y_{i+k, j+l})$ (por. „pudełko” po lewej stronie na rys. 13.2).
2. Jeśli (i, j) znajdzie się w pobliżu krawędzi, wówczas $\mathcal{B}_{ij}(h_x, h_y, H)$ zawiera tylko pewną część pikseli $(i+k, j+l, y_{i+k, j+l})$ (por. „pudełko” po prawej stronie na rysunku 13.2. Sytuację tę zilustrowano także w przekroju pokazanym na rysunku 13.3). Widać na nim, że jeśli prostopadłościan znajdzie się w części o wyższych wartościach y -ków, to nie będzie zawierał pikseli po lewej stronie od uskoku. Podobnie, jeśli znajdzie się na lewo od krawędzi, to nie będzie zawierał pikseli o wyższych wartościach y -ków. Aby stwierdzenia te były prawdziwe, musimy odpowiednio dobrać wartość H , gdyż przy zbyt dużej – w stosunku do wielkości skoku – wartości H , do „pudełka” trafią również piksele znajdujące się po przeciwnej stronie krawędzi.

Patrząc na rysunek 13.2, stwierdzamy, że wraz ze zmianami (i, j) , „pudełko” $\mathcal{B}_{ij}(h_x, h_y, H)$ porusza się gładko po obszarach, gdzie intensywność poziomów



Rys. 13.5. Próbkowanie zmian w rzeczywistym obrazie – wpływ doboru H

szarości jest w miarę jednolita i skacze (w górę lub w dół) w tych obszarach, gdzie znajdują się krawędzie. Wraz z tymi zmianami wysokości położenia, zmienia się także liczba pikseli zawartych w $\mathcal{B}_{ij}(h_x, h_y, H)$, a śledzenie tej liczby pozwala na wypróbkowanie i zaznaczenie tych fragmentów obrazu, gdzie następują duże (większe niż H) zmiany poziomów szarości. Ideę tę zrealizowano w opisanym następującym algorytmie, który nazywamy algorytmem próbkowania¹ krawędzi.

ALGORYTM PRÓBKOWANIA KRAWĘDZI

Krok 0. Wybrać następujące parametry algorytmu.

1. Wybrać rozmiary okna poruszającego się w płaszczyźnie obrazu, które stanowi podstawę prostopadłościanu $\mathcal{B}_{ij}(h_x, h_y, H)$. Parametry $1 \leq h_x \leq n_x$,

¹ W odróżnieniu od algorytmów śledzenia konturów, które mają zwykle na celu aproksymację konturów, naszym celem jest jedynie ich zaznaczenie, co wyjaśnia termin *próbkowanie krawędzi*.

$1 \leq h_y \leq n_y$ wyznaczają połowę długości i szerokości tego okna. Dokładniej, ma ono wymiary $(2h_x + 1) \times (2h_y + 1)$. W badaniach symulacyjnych, prowadzonych przez autora, przyjmowano $h_x = h_y$ i wybierano je z przedziału $[1, 7]$.

2. Wybrać wysokość prostopadłościanu $0 < H < 1$, biorąc pod uwagę naturalną zmienność obrazu w obszarach o małych zmianach poziomów szarości oraz wysokość zmian, które chcemy wypróbować. H nie powinno być zbyt małe, gdyż wówczas algorytm staje się nadmiernie czuły na drobne, naturalne zmiany w obrazie. Zbyt duże wartości H prowadzą z kolei do „znieczulenia” algorytmu na istotne zmiany poziomów szarości, dlatego H powinno być tylko nieco większe niż najmniejsze zmiany, które chcemy wykrywać. W badaniach symulacyjnych przyjmowano H w przedziale $[0.05, 0.2]$.
3. Wybrać $0 < q < 1$, które oznacza ułamek całkowitej liczby obserwacji $L = (2h_x + 1) \cdot (2h_y + 1)$ w $\mathcal{B}_{ij}(h_x, h_y, H)$, jeśli prostopadłościan znajduje się w obszarze o małej zmienności poziomów szarości. Gdy obraz składa się z płaskich obszarów, oddzielonych uskokami, to można wybrać q bliskie 1. Jeśli natomiast występują zakłócenia lub dopuszczamy, że małe zmiany poziomów szarości nie zostaną zaznaczone, to możemy wybrać q rzędu 0.6–0.7.
4. Przygotować tablicę binarną Υ o rozmiarach $(n_x - 2h_x) \times (n_y - 2h_y)$, w której przechowywane będą wyniki próbkowania, traktowane jako rezultaty pośrednie. Oznaczmy wartości elementów tej tablicy przez v_{ij} . Tablicę Υ wypełniamy zerami ($v_{ij} = 0$ oznacza, że odpowiedni piksel (i, j) obrazu nie został zaklasyfikowany jako fragment obrazu o dużej zmienności).

Krok 1. Dla bieżących wartości (i, j) ustal zawartość $\mathcal{B}_{ij}(h_x, h_y, H)$ i oblicz liczbę $b_{ij}(h_x, h_y, H)$ elementów w tym prostopadłościanie.

Krok 2. Jeśli liczba ta spełnia warunek

$$b_{ij}(h_x, h_y, H) < qL, \quad \text{to} \quad v_{ij} = 1, \quad (13.1.3)$$

w przeciwnym razie pozostaw $v_{ij} = 0$.

Krok 3. Powtarzaj kroki 1 i 2 zmieniając $i = h_x, h_x + 1, \dots, (n_x - h_x)$ oraz $j = h_y, h_y + 1, \dots, (n_y - h_y)$.

W odróżnieniu od części znanych algorytmów okonturowywania, proponowany algorytm nie preferuje żadnego z kierunków. Kosztem wzrostu nakładów obliczeniowych, można wybrać $\mathcal{B}_{ij}(h_x, h_y, H)$ o innym kształcie niż prostopadłościan. Kandydatem może być walec, którego podstawa leży w płaszczyźnie obrazu.

Rezultaty badań empirycznych

Zanim przejdziemy do zastosowań opisanego algorytmu w celach diagnostycznych, przedstawimy wyniki jego testowania na wybranych obrazach. We wszyst-

kich przypadkach przyjęto, że rozmiary okna w płaszczyźnie obrazu są jednakowe i oznaczają będziemy je przez $h \stackrel{\text{def}}{=} h_x = h_y$.

Analiza rysunków 13.4, 13.5, 13.6 i innych obrazów, których tu nie pokazano, prowadzi do następujących wniosków:

1. Jeśli q jest zbyt duże, to również piksele w otoczeniu krawędzi są zaznaczane. Zbyt małe wartości tego parametru skutkują tym, że nie wszystkie krawędzie zostają wykryte.
2. Jeśli rozmiary okna w płaszczyźnie obrazu są zbyt małe ($h = 1$), to kontury zaznaczane są cienką linią, lecz dobór właściwych wartości q staje się trudniejszy.
3. Zbyt małe wartości H czynią metodę nadmiernie czułą na zmiany w obrazie i wówczas również drobne (punktowe) zakłócenia oznaczane są jako istotne zmiany. Wybór zbyt dużych wartości H powoduje, że nie wszystkie krawędzie zostają wykryte.

Podsumowując, jeśli właściwie dobierzemy parametry proponowanego algorytmu, to uzyskujemy narzędzie do szybkiego i odpornego na zakłócenia próbkowania tych fragmentów obrazu, w których następują duże zmiany poziomów szarości. Sterując parametrami, możemy dość precyzyjnie określać rodzaj wykrywanych zmian.

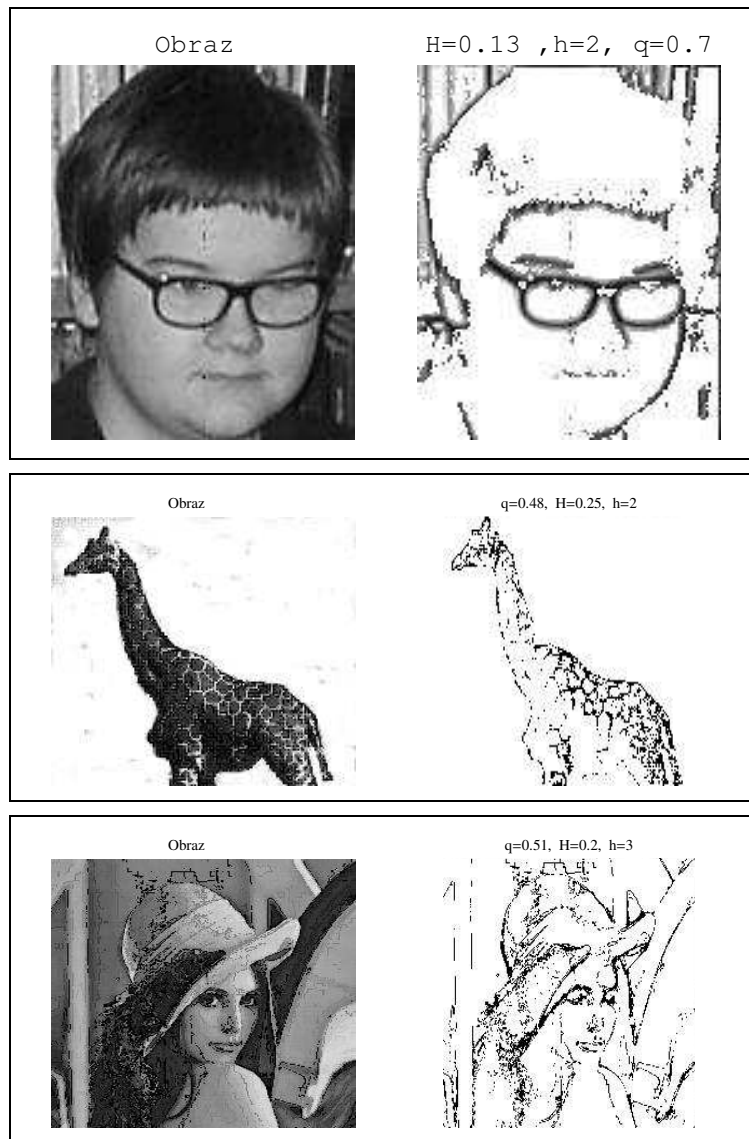
Przygotowanie danych do wykrywania zmian w sekwencji obrazów

Naszym głównym celem jest przystosowanie zaproponowanego algorytmu do zadań związanych z zastosowaniami przemysłowymi. Z tego punktu widzenia algorytm ten wymaga pewnej modyfikacji, gdyż okazuje się on nadmiernie wrażliwy na najdrobniejsze nawet nierówności powierzchni, które są trudne do uniknięcia w trakcie produkcji. Chodzi więc o taką modyfikację, która pozostawi algorytm wrażliwym na większe niejednorodności powierzchni i jednocześnie nie będzie wskazywać niejednorodności o powierzchni kilku pikseli.

Proponowana modyfikacja polega na przeglądnięciu pikseli obrazu Υ i zbadaniu ich sąsiadów, dla których $v_{ij} = 1$. Jeśli liczba sąsiadów zaznaczonych również jako „czarne” przekroczy zadany poziom (np. 50%), to (i, j) pozostaje zaznaczony jako piksel o dużej zmienności obrazu. W przeciwnym razie zmieniamy jego kwalifikację i kładziemy $v_{ij} = 0$.

Jako dane wyjściowe do wykrywania zmian w sekwencji obrazów przyjmujemy procent pikseli uznanych za te, w których występują istotne zmiany w obrazie (procent „czarnych” pikseli), bez względu na ich położenie. Ten sposób przygotowania danych jest szybki, prosty i daje w wyniku informację o procencie powierzchni, gdzie występują niejednorodności.

Algorytm próbkowania krawędzi wraz z opisaną tu modyfikacją zastosowano do sekwencji obrazów próbek miedzi podgrzewanych przed walcowaniem. Oryginalne obrazy i wyniki ich przetworzenia pokazano na rysunku 13.7. Na obrazie



Rys. 13.6. Dalsze przykłady działania algorytmu próbkowania zmian w obrazach

nr 7 wykryty został większy obszar o temperaturze różniącej się od temperatury otoczenia. Obraz ten pokazano w powiększeniu na rysunku 13.8 po to, by wykryty obszar można było porównać z oryginałem. Na rysunku 13.9 pokazano procentowy udział obszarów o większej zmienności w poszczególnych obrazach. Obraz nr 7 wyraźnie się wyróżnia, ale uwagę zwraca również próbka nr 5, w której rozrzucone i niewielkie obszary nierównomierności zajmują łącznie 1.84% powierzchni,

co zauważalnie przekracza średnią dla wszystkich siedmiu obrazów, która wynosi 1.54%.

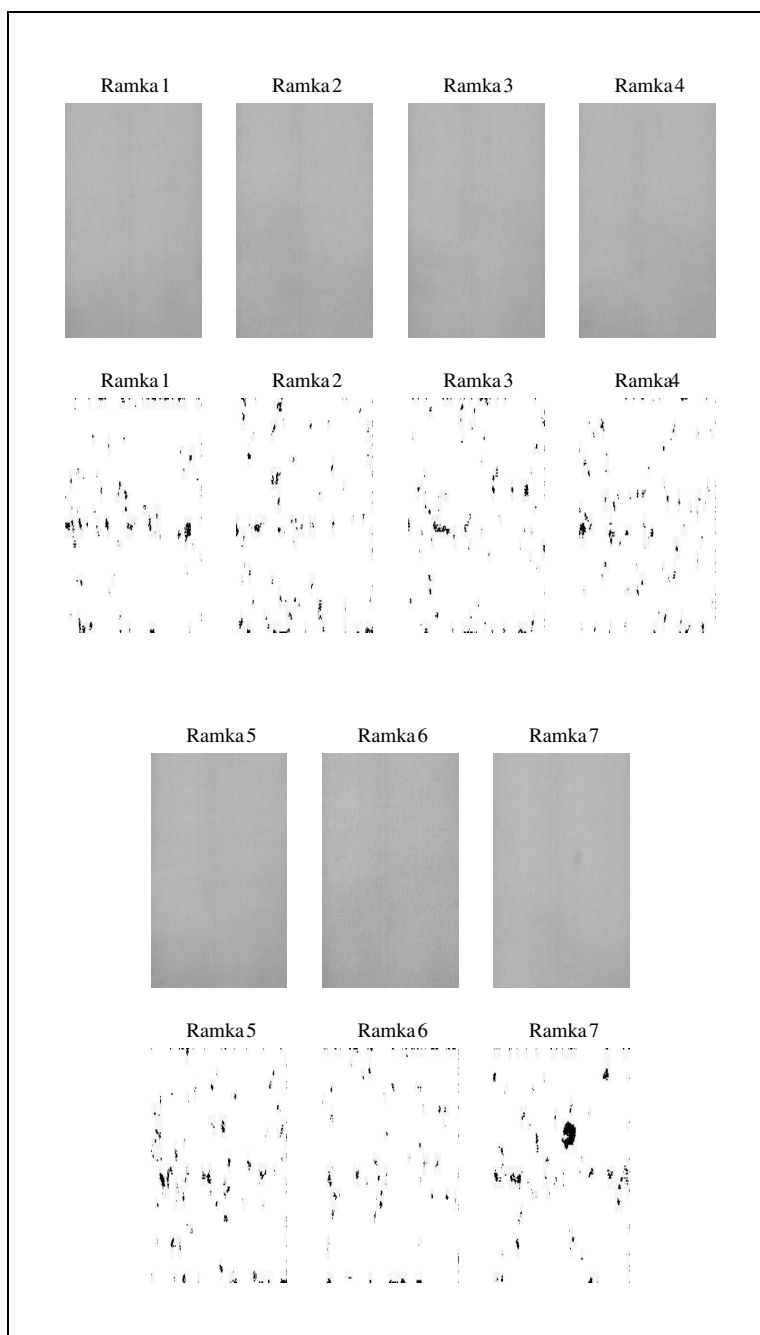
13.2. Próbkowanie wymiaru fraktalnego obrazu jako wskaźnika diagnostycznego

Na rysunku 13.10 pokazano przebiegi temperatury w wybranych przekrojach miedzianego wlewka, którego obraz poznaliśmy już wcześniej (na rys. 13.7 oznaczono go jako ramka 1). Funkcje te sprawiają wrażenie losowych i charakteryzuje je duża zmienność, która widoczna jest na przedstawionych powiększeniach. Jednocześnie, same wartości poziomów układają się względnie blisko wartości średnich, obliczonych dla każdego z przekrojów osobno. Powoduje to, że jesteśmy skłonni traktować powierzchnię na rysunku 13.7 ramka 1 jako jednolitą w skali makroskopowej. Nasuwa się więc pomysł wykorzystania wymiaru fraktalnego do liczbowej oceny stopnia (nie-)równomierności powierzchni, bazujący na następujących przesłankach.

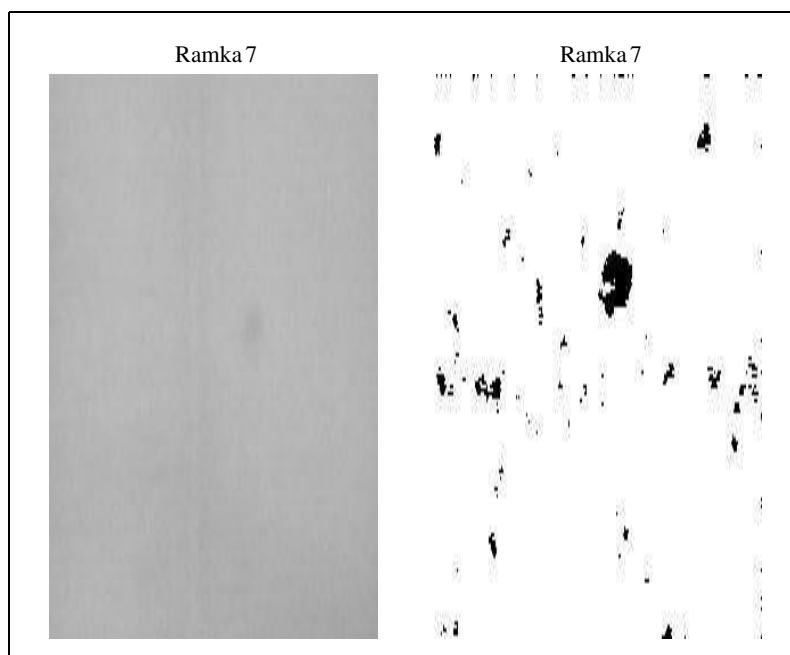
1. Wymiar fraktalny procesów stochastycznych można traktować jako miarę szybkości ich zmienności. Jak wiadomo (por. [162], [34]), funkcje różniczkowalne mają wymiar fraktalny równy 1, a bardzo szybko zmieniające się procesy, takie jak klasyczny ruch Browna, mają wymiar fraktalny równy 2.
2. Jeśli ocenimy wymiar fraktalny przebiegu poziomów szarości w danym przekroju i stwierdzimy, że ma on znaczną wartość (powiedzmy, powyżej 1.75) i jednocześnie obserwowane poziomy szarości nieznacznie różnią się od wartości średniej, to można uznać, że przebieg jest równomierny. Odwrotnie, małe wartości wymiaru fraktalnego mogą świadczyć o tym, że na obrazie znajduje się obszar, w którym zmienność jest zbyt mała jak na równomierną powierzchnię. Można się tutaj odwołać do przykładu dobrze naostrzonej żyłki – ma ona ostrze, które w powiększeniu wygląda jak krajobraz Tatr, natomiast zużyta żyłka ma znaczne obszary równej powierzchni. Podobnie wyglądają obrazy narzędzi używanych do skrawania, gdzie dodatkowym czynnikiem wywołującym mniejsze wartości wymiaru fraktalnego jest złamanie się fragmentu ostrza.

Obszerny przegląd literatury na temat zastosowań i technik przetwarzania obrazów w diagnostyce procesów wytwórczych przedstawiono w [81].

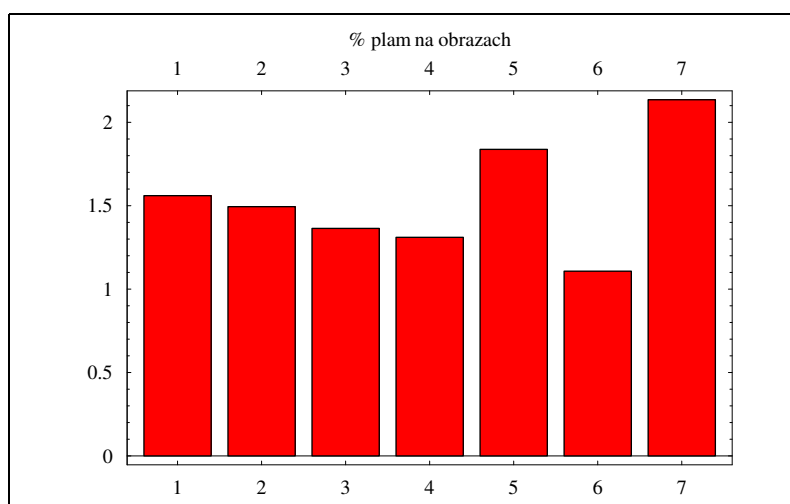
W pracy [22] zastosowano szacowanie tzw. wymiaru pudełkowego do wykrywania defektów w produkcji tkanin. Proponowane tu podejście opiera się na wymiarze korelacyjnym, jako bardziej odpowiednim do opisu procesów stochastycznych i pól losowych, a obszar potencjalnych zastosowań obejmuje badanie stopnia (nie-)równomierności powierzchni wykazujących jedynie statystyczne samopodobieństwo, które jest znacznie częściej spotykane niż samopodobieństwo



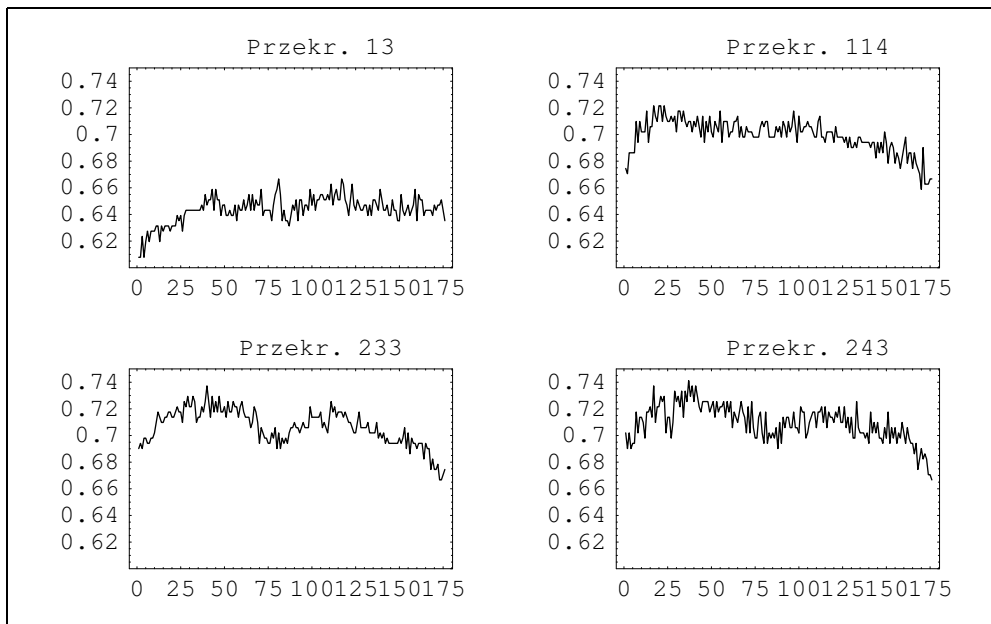
Rys. 13.7. Wyniki próbkowania obszarów o większej zmienności temperatury w próbkach miedzi



Rys. 13.8. Wynik próbkowania obszarów o większej zmienności temperatury w próbkach miedzi – powiększony obraz nr 7



Rys. 13.9. Procentowy udział obszarów o większej nierównomierności w obrazach 1 – 7, pokazanych na rys. 13.7



Rys. 13.10. Przekroje poziomów szarości na obrazie nr 1

geometryczne, mające charakter deterministyczny. Definicje powyższych pojęć i podstawowe własności wymiarów fraktalnych zebrano między innymi w [162], [34], [8].

Estymacja korelacyjnego wymiaru fraktalnego

Przedstawimy metodę estymacji korelacyjnego wymiaru fraktalnego, oznaczonego dalej przez F_{dim} , w wersji opisanej w [162]. Mimo że proponowane podejście nie jest w pełni zależne od sposobu estymacji wymiaru fraktalnego fragmentów obrazu, to wybór metody korelacyjnej ma, oprócz opisanych wyżej zalet, także inne walory, a mianowicie,

- względną prostotę implementacji,
- mały nakład obliczeniowy,
- znaczną dokładność w porównaniu z innymi metodami, jeśli rozważa się cały zakres wymiarów $[1, 2]$ fraktalnych sygnałów, częściej spodziewając się sygnałów o wymiarach powyżej 1.25.

Uwaga 13.3. *Jeśli spodziewamy się, że $1 \leq F_{dim} \leq 1.25$, to wówczas wymiar korelacyjny można estymować dokładniej (z mniejszym błędem średniokwadratowym) metodą wariacji kwadratowej (MWK). Asymptotykę błędów zbadano w [58], [10], a przedstawione w [143] badania symulacyjne potwierdziły większą dokładność MWK dla próbek o średnich rozmiarach.*

Popularna jest też wersja Grassberga–Procaccia szacowania wymiaru korelacyjnego (por., np. [97]).

Zastosować można także wymiar pudełkowy (ang. box-counting dimension). W pracy [180] zaproponowano algorytm szacowania tego wymiaru z użyciem krzywych wypełniających.

Niech $s(t)$ będzie stacjonarnym procesem stochastycznym o skończonej wariancji. Funkcję kowariancji $s(t)$ oznaczamy będziemy przez $\gamma(t) = \text{cov}(s(t), s(0))$.

Założymy dodatkowo, że funkcja kowariancji ma (dla małych opóźnień t) następującą postać:

$$\gamma(t) = \gamma(0) - c|t|^{2H} + o(|t|), \quad \text{gdy } t \rightarrow 0, \quad (13.2.4)$$

gdzie $c > 0$ i $0 < H \leq 1$ są pewnymi stałymi.

Parametr H charakteryzuje gładkość trajektorii procesu $s(t)$. Można wykazać, że jeśli $s(\cdot)$ jest procesem gaussowskim i $s(\cdot)$ ma różniczkowalne trajektorie, to $H = 1$.

Dla szerokiej klasy procesów drugiego rzędu zachodzi bezpośredni związek między H i wymiarem fraktalnym

$$F_{dim}(s) = 2 - H, \quad (13.2.5)$$

ale zależność ta nie obejmuje wszystkich procesów tej klasy (w monografiach [1], [192] i pracach [10], [11] znaleźć można warunki dostateczne dla stosowalności (13.2.5)).

Niech $s_i = s(i\tau)$, $i = 1, 2, \dots, n$ oznaczają próbki wartości procesu $s(\cdot)$ wykonane z okresem próbkowania $\tau > 0$. Należy wybrać maksymalną liczbę etapów opóźnień $1 < M < n$, dla których obliczany będzie wariogram. M powinno być ułamkiem całkowitej liczby próbek n . Wariogram obliczany jest następująco

$$g_j = (n - j)^{-1} \sum_{i=1}^{n-j} (s_{i+j} - s_i)^2, \quad j = 1, 2, \dots, M. \quad (13.2.6)$$

Wartości wariogramu g_j są oszacowaniami dla $2(\gamma(0) - \gamma(j\tau))$, $j = 1, 2, \dots, M$.

Na podstawie (13.2.4), dla $|t|$ dostatecznie małych, otrzymujemy

$$\log(\gamma(0) - \gamma(t)) = 2H \log(|t|) + \text{const}. \quad (13.2.7)$$

Możemy potraktować $\log(g_j)$ jako dostępne „obserwacje” $\gamma(0) - \gamma(j\tau)$, a (13.2.7) interpretujemy jako regresję liniową o współczynniku nachylenia $2H$.

Oszacowanie \hat{H} dla H możemy zatem obliczyć jako połowę współczynnika regresji liniowej, w której $\log(g_j)$ i $\log(j)$, $j = 1, 2, \dots, M$ traktujemy jak obserwacje zmiennej zależnej i niezależnej.

Warunki dostateczne zbieżności, wraz ze wzrostem n , opisanego wyżej estymatora \hat{H} do H znaleźć można w pracy [51].

Opis algorytmu metody korelacyjnej

Krok 1. Pobierz próbki $s_i, i = 1, 2, \dots, n$ sygnału $s(\cdot)$. Dla $j = 1, 2, \dots, M$ oblicz g_j zgodnie ze wzorem (13.2.6).

Krok 2. Znajdź \hat{a} i \hat{b} , które minimalizują względem a i b

$$\sum_{j=1}^M (\log(g_j) - (a \log(j) + b))^2. \quad (13.2.8)$$

Krok 3. Oblicz oszacowanie wymiaru fraktalnego sygnału $s(\cdot)$ według wzoru $\hat{F}_{dim}(s) = 4 - \hat{a}/2$.

Można podać jawne wzory dla a i b , które minimalizują (13.2.8).

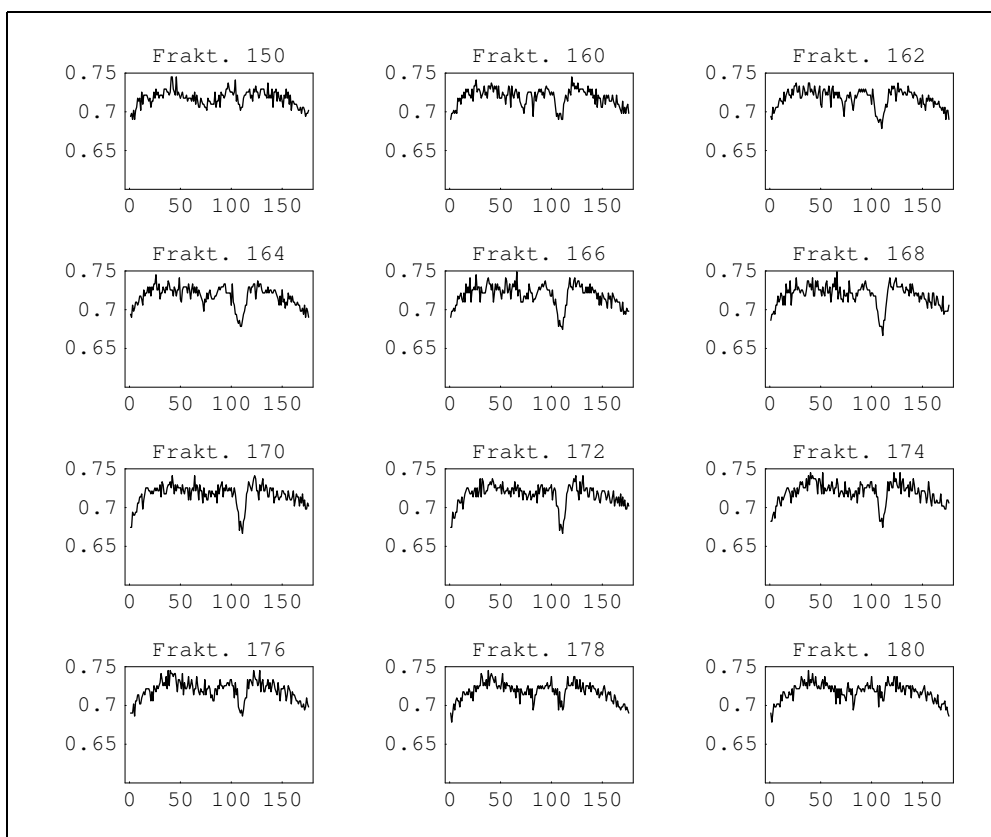
Z badań symulacyjnych przeprowadzonych przez autora wynika, że dla rozkład prawdopodobieństwa błędów $\hat{F}_{dim}(s) - F_{dim}(s)$ można traktować jako rozkład normalny o średniej 0 i dyspersji 0.05.

Przygotowanie danych do wykrywania zmian w obrazach

Pojęcie wymiaru fraktalnego jest na tyle szerokie, że można zastosować je zarówno do całych obrazów, jak i do – różnie dobranych – fragmentów obrazu.

W zastosowaniach tego pojęcia do obrazów pochodzących z przebiegających w sposób ciągły procesów produkcyjnych (wytop metali, produkcja papieru czy tkanin) celowe wydaje się szacowanie wymiaru fraktalnego każdej bieżącej linii pikseli obrazu. Zakładamy, że linie te są prostopadłe do kierunku przesuwania się wytwarzanego materiału. Taki wybór sposobu próbkowania obrazu pozwala na bieżącą kontrolę każdego pojawiającego się fragmentu wyrobu. Jednocześnie opisany wyżej algorytm szacowania wymiaru korelacyjnego można zastosować bez żadnych zmian, gdyż każda linia pikseli z poziomami szarości obrazu może być traktowana jako próbki jednowymiarowego sygnału. Jeśli przez $s_i^{(k)}$ oznaczymy poziom szarości i -tego piksela w k -tej linii obrazu, $i = 1, 2, \dots, n$, $k = 1, 2, \dots$, to jako dane wejściowe do algorytmu szacowania wymiaru fraktalnego podajemy każdorazowo zestaw $s_i^{(k)}$, $i = 1, 2, \dots, n$, a w wyniku dostajemy jedną liczbę $F_{dim}(s^{(k)}(\cdot))$ dla każdej z linii $k = 1, 2, \dots$

Na rysunku 13.7 linie te przebiegają poziomo, a przykładowe przekroje wartości poziomów szarości $s_i^{(k)}$, $i = 1, 2, \dots, n$ (dla $n = 175$) pokazano na rysunkach 13.10 i 13.11. Mimo że na rysunkach tych widoczna jest pewna zmienność poziomów szarości w przekrojach, to ich wartości średnie – liczone dla całych obrazów – różnią się między sobą zbyt mało, by mogły stanowić wskaźnik diagnostyczny (por. tab. 13.1). Jeśli natomiast do każdej linii poziomej w każdym z tych



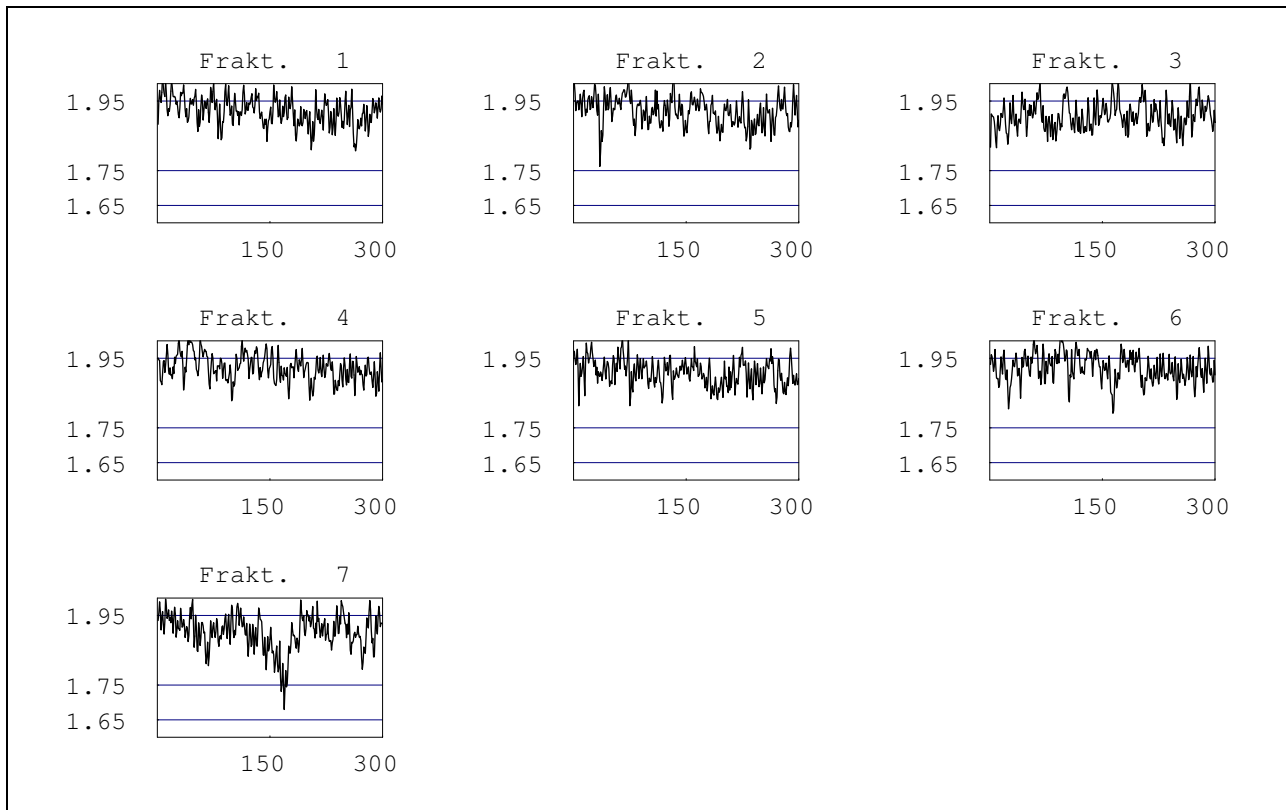
Rys. 13.11. Przekroje poziomów szarości na obrazie nr 7

Ramka Nr	1	2	3	4	5	6	7
Średnia	0.690	0.689	0.691	0.687	0.689	0.687	0.708

Tabela 13.1. Wartości średnie poziomów szarości dla obrazów pokazanych na rys. 13.7

obrazów zastosujemy procedurę szacowania wymiaru fraktalnego, to otrzymamy wyniki pokazane na rysunku 13.12. W następnym rozdziale wyniki te posłużą jako dane wejściowe do karty kontrolnej, której zadaniem jest wykrywanie skokowych zmian w szeregu czasowym. W rozważanym przykładzie wymiar fraktalny okazał się na tyle wrażliwym wskaźnikiem diagnostycznym, że znaczny spadek wymiaru widoczny jest na ostatnim wykresie (ramka 7) rysunku 13.12. Zauważmy, że spadek poniżej wymiaru 1.75 nastąpił w tych kilkunastu liniach, które swym położeniem odpowiadają położeniu „plamy” widocznej na rysunku 13.8.

Rys. 13.12. Oszacowania wymiarów fraktalnych przekrojów obrazów 1-7



14. Wykrywanie zmian jakości w sekwencjach obrazów

W dwóch poprzednich rozdziałach przedstawiono różne podejścia do uzyskania danych diagnostycznych o charakterystykach wyrobów. Charakterystyki te były albo funkcjami, albo danymi dwuwymiarowymi w postaci obrazów. Wspólnym efektem końcowym różnych podejść omawianych w tych rozdziałach była sekwencja liczb, charakteryzująca jakość poszczególnych egzemplarzy wyrobów lub ich fragmentów, jak np. w produkcji ciągłej.

Pozostaje jedynie wykrywać na bieżąco znaczące zmiany w takiej sekwencji liczb, by móc monitorować jakość badanego procesu. Zagadnienia tego rodzaju były intensywnie badane od połowy XX wieku i są nadal obiektem zainteresowania, zarówno w wielu ośrodkach badawczych, jak i w praktyce przemysłowej. Tradycyjnie metody monitorowania i wykrywania zmian w jakości produkcji nazywane są kartami kontrolnymi, chociaż dzisiaj znacznie częściej używane są one jako algorytmy komputerowe i wykresy na monitorze.

14.1. Nieparametryczna karta kontrolna

Opisom i badaniu własności statystycznych kart kontrolnych poświęcona jest bardzo duża liczba prac i monografii (por. bibliografie w [89], [188], [56]). Mimo istnienia wielu dobrze zbadanych kart o uznanych zaletach (karta Shewharta, EWMA, CUSUM itp.), w rozdziale tym proponujemy modyfikację kart opartych na teście znaków Wilcoxona. Proponowana karta ma cechy wspólne z kartą CUSCORE zaproponowaną w [92]. Zasadnicza różnica polega na tym, że rozważana tutaj karta ma skończony bufor pamięci, co powoduje jej szybszą reakcję na bardzo małe (rzędu 0.1 zmienności procesu) zmiany jakości. Dzięki tej własności, proponowana karta ma, w zamysle autora, być dobrze dostosowana do potrzeb wykrywania zmian w sekwencjach obrazów. Bardziej szczegółowe uzasadnienie podajemy w tym podrozdziale.

Uzasadnienie kierunku badań. Analiza prac cytowanych w wymienionych monografiach wskazuje, że większość popularnych kart kontrolnych w istotny sposób korzysta z założenia o normalności rozkładu zakłóceń mierzonych wielkości. W zastosowaniu do monitorowania sekwencji obrazów takiego założenia a priori przyjmować nie możemy, gdyż ciąg liczb, którego zmiany chcemy badać nie składa się z wielkości mierzonych wprost, lecz powstaje w wyniku dość skomplikowanych i nieliniowych operacji na obrazach. Z tego powodu potrzebna

jest nam karta nieparametryczna. Przymiotnik *nieparametryczna(-y)* używany jest w statystyce w kilku, nieco różnych, znaczeniach. Tutaj, używać będziemy go dla zaznaczenia, że nie zakładamy żadnej konkretnej klasy rozkładów prawdopodobieństwa odchyłeń obserwowanych wielkości od wartości nominalnej. Jako jedyne założenie przyjmować będziemy symetrię tego rozkładu wokół zera. Jest to minimalne założenie niezbędne do tego, by zakłócenia nie wprowadzały systematycznego błędu. Z drugiej strony, jest ono tak mało restykcyjne, że dopuszczamy rozkłady zakłóceń, które mają nieskończoną wariancję (np. rozkład Cauchy'ego), co pozwala modelować zakłócenia o większych niż w rozkładzie normalnym prawdopodobieństwach pojawienia się dużych błędów. Dzięki temu karta jest na nie bardziej odporna, niż karty bazujące na rozkładzie normalnym. Zauważmy też, że potrzebna jest nam karta operująca na pojedynczych obserwacjach, a nie na średnich z podgrup (jak, przykładowo, w karcie Shewharta), co wyklucza odwołanie się do centralnego twierdzenia granicznego. Przegląd literatury na temat nieparametrycznych kart kontrolnych zawiera praca [19]. Inne podejścia do konstruowania nieparametrycznych kart kontrolnych zaproponowano w [102] i [104], gdzie nieparametryczny jest model zmienności procesu, który może być funkcją z dość szerokiej klasy.

Drugim powodem rozważania proponowanej modyfikacji kart opartych na badaniu znaków odchyłeń jest, wspomniana już, własność wykrywania małych zmian jakości w czasie krótszym – w sensie średnim – niż wykrywają je znane karty. Odbywa się to kosztem wydłużonej reakcji na duże zmiany jakości. Analiza znanych kart wskazuje jednak, że nie są znane karty „uniwersalne”, o krótszym niż inne średnim czasie detekcji zmian w szerokim zakresie ich amplitud. Ponadto w obecnym stanie techniki komputerowej można prowadzić monitorowanie za pomocą kilku kart równocześnie, dobierając je tak, by miały krótkie średnie czasy reakcji w poszczególnych przedziałach potencjalnej zmienności procesu.

Model zmian jakości i opis karty

Będziemy przyjmować klasyczny model obserwacji procesu

$$Y_n = Y + m \mathbf{1}(n - q) + \varepsilon_n, \quad n = 1, 2, \dots, \quad (14.1.1)$$

gdzie Y oznacza znany, pożądany poziom wielkości charakteryzującej jakość procesu, na przykład temperatury metalu przed walcowaniem. Dalej, dla wygody zakładać będziemy $Y = 0$, gdyż w celu uwzględnienia niezerowej wartości Y wystarczy odejmować ją od bieżących, obserwowanych wartości Y_n charakterystyki jakości tegoż procesu. Przez ε_n , $n = 1, 2, \dots$ oznaczamy wartości losowych, niemie-rzalnych zakłóceń w jakości i obserwacjach badanego procesu. Dalsze założenia na ich temat przedstawimy nieco później.

Przyjmujemy najprostszy, skokowy model możliwej trwałej zmiany (zwykle pogorszenia) charakterystyki procesu. Skokowa zmiana o nieznaną wartość pa-

parametru $m \neq 0$ może pojawić się w nieznannej chwili czasu $q > 0$. Dla uproszczenia przyjmujemy, że chwila ta jest jedną z tych, w których dokonywane są obserwacje $n = 1, 2, \dots$. W modelu (14.1.1) $\mathbf{1}(t)$ oznacza skok jednostkowy w chwili 0, tzn.

$$\mathbf{1}(t) = \begin{cases} 0, & \text{gdy } t < 0, \\ 1, & \text{gdy } t \geq 0. \end{cases} \quad (14.1.2)$$

Wyrażenie $m\mathbf{1}(n - q)$ reprezentuje zatem skok o wartości m , który wystąpił w chwili $q > 0$.

O zakłóceniach ε_n , $n = 1, 2, \dots$ zakładamy, że są one niezależnymi zmiennymi losowymi o jednakowych rozkładach prawdopodobieństw.

Zwykle przyjmuje się, że ε_n mają wartość oczekiwaną zero, skończoną wariancję i rozkład normalny. Będziemy unikać tych założeń, dążąc do tego, by karta działała nie tylko bez założenia o normalności rozkładu zakłóceń, ale także bez zakładania jakiegokolwiek konkretnej postaci ich rozkładu. Ceną za to jest brak możliwości podania rozkładu statystyki testowej, a – co za tym idzie – brak możliwości podania średniego czasu do wykrycia zmiany jakości. Konsekwencją jest konieczność posilkowania się albo rozkładami asymptotycznymi, albo badaniami symulacyjnymi w celu ustalenia granic dopuszczalnych odchyłeń statystyki testowej.

Co więcej, nie wymagamy, by rozkład ε_n posiadał wartość oczekiwaną, gdyż chcemy dopuścić do rozważań takie rozkłady, jak rozkład Cauchy'ego. Chcemy jednak zachować podstawową intuicję modelu addytywnych zakłóceń, a mianowicie, wymaganie, by zakłócenia nie wносиły systematycznego błędu. W tym celu zakładamy, że ε_n , $n = 1, 2, \dots$ mają rozkład prawdopodobieństwa, który jest symetryczny względem zera, tzn. jego dystrybuanta $F(x)$ spełnia warunek

$$F(x) = 1 - F(-x), \quad x \in R. \quad (14.1.3)$$

Zauważmy, że nie wymagamy nawet istnienia gęstości rozkładu prawdopodobieństwa. ε_n , a więc dopuszczamy także dyskretne rozkłady prawdopodobieństwa, o ile tylko są one symetryczne względem zera. Jeśli natomiast gęstość istnieje, to warunek (14.1.3) przyjmuje postać $f(x) = f(-x)$, gdzie f jest pochodną F . Warunek ten spełnia nie tylko rozkład Cauchy'ego, ale także na przykład rozkład Laplace'a.

Pożądane cechy karty. Celem działania karty kontrolnej jest wykrycie zmiany o nieznannej wartości m , która wystąpiła w nieznannej chwili czasu $q > 0$. Zakładamy przy tym, że nieznaną dystrybuanta zakłóceń F spełnia warunek (14.1.3). Pożądane jest, by wartość oczekiwana czasu do wykrycia zmiany jakości była możliwie mała, ale jednocześnie karta powinna mieć parametry, które pozwolą ustawić dostatecznie długi średni czas reakcji karty, jeśli zmiana jakości nie nastąpiła. Innymi słowy użytkownik powinien móc wpływać na wartość oczekiwaną czasu, w którym następuje fałszywy alarm.

Pożądaną cechą karty powinna być też prostota jej implementacji programowej i możliwość wykorzystania doświadczeń zaczerpniętych z działania innych kart kontrolnych lub testów statystycznych. Opisywana w tym podrozdziale karta posiada te cechy, a zanim ją opiszemy, przedstawimy motywacje, które prowadzą do jej konstrukcji

Uwagi na temat karty do badania liczby defektów. Punktem wyjścia jest proste spostrzeżenie, że jeśli przyjmiemy model obserwacji (14.1.1) z warunkiem symetrii zakłóceń (14.1.3) i zmiana jakości procesu nie zaszła ($m = 0$), to – z dokładnością do statystycznych fluktuacji – można oczekiwać, że około połowa obserwacji Y_n będzie mieć wartości dodatnie, a pozostałe będą ujemne.

Wprowadźmy zmienną losową Z_n , która wskazywać będzie nam znak danej obserwacji

$$Z_n \stackrel{\text{def}}{=} \text{sign}(Y_n) = \begin{cases} 0, & \text{jeżeli } Y_n < 0, \\ 1, & \text{jeżeli } Y_n \geq 0, \end{cases} \quad n = 1, 2, \dots, N, \quad (14.1.4)$$

gdzie N jest liczbą danych aktualnie dostępnych. Zdefiniujmy zliczającą zmienną losową

$$I_N \stackrel{\text{def}}{=} \text{card}\{Z_i = 1, i = 1, 2, \dots, N\} = \sum_{i=1}^N Z_i. \quad (14.1.5)$$

Przy poczynionych założeniach wartość oczekiwana tej zmiennej losowej wynosi $E(I_N) = N/2$, gdyż I_N ma dwumianowy rozkład prawdopodobieństwa z prawdopodobieństwem sukcesu w jednej próbie $p_0 = 1/2$.

Jeśli wystąpiła zmiana jakości procesu o wartość $m \neq 0$, to rozkład Y_n nie jest już rozkładem symetrycznym względem zera dla $n > q$, a prawdopodobieństwo zdarzenia $Z_n = 1$ wynosi

$$p_1 = 1 - F(-m), \quad (14.1.6)$$

przy czym p_1 może być zarówno większe, jak i mniejsze od $1/2$, w zależności od tego czy m jest dodatnie czy ujemne. Możemy zatem wykryć zmianę jakości, testując hipotezę $H_0 : p_0 = 1/2$ przeciw alternatywom, że prawdopodobieństwo sukcesu w pojedynczej próbie jest różne od $1/2$.

Jeśli zmiana jakości nie wystąpiła ($m = 0$), to dyspersja zmiennej losowej I_N jest równa $\sqrt{N p_0 (1 - p_0)}$. I_N/N ma zatem wartość oczekiwaną $1/2$ i dyspersję $\sqrt{p_0 (1 - p_0)}/N$. Jeśli N jest dostatecznie duże, to rozkład dwumianowy aproksymować można rozkładem normalnym i na jego podstawie obliczać prawdopodobieństwa zdarzeń. W szczególności, możemy obliczyć prawdopodobieństwo zdarzenia polegającego na przekroczeniu zadanych granic przez zmienną losową I_N/N , zakładając, że $m = 0$. Oznaczmy dolną granicę kontrolną dla I_N/N przez

LCL, a górną przez UCL. Nazwy te są skrótami od *lower (upper) chart limit*. Przyjmować będziemy, że granice kontrolne mają postać

$$\text{UCL} = p_0 + k \sqrt{p_0(1-p_0)/N}, \quad (14.1.7)$$

$$\text{LCL} = p_0 - k \sqrt{p_0(1-p_0)/N}, \quad (14.1.8)$$

gdzie k jest stałą, która podlega wyborowi. Wybieramy ją z tablic rozkładu normalnego tak, by prawdopodobieństwo, że $I_N/N < \text{LCL}$ lub $I_N/N > \text{UCL}$ nie przekraczało wybranej przez nas wartości, na przykład, 0.001 lub 0.01. Innym, choć w tym przypadku równoważnym, sposobem doboru k jest ustalenie z góry wartości oczekiwanej czasu do zasygnalizowania fałszywego alarmu, czyli przekroczenia powyższych granic, mimo że $m = 0$. Często wybiera się czas rzędu 465, co prowadzi do przyjęcia $k = 3$. Po wybraniu k , karta ta działa w ten sposób, że jeśli I_N/N przekroczy jedną z granic UCL lub LCL, to sygnalizowane jest wyjście jakości procesu poza dopuszczalny obszar.

Przedstawiona karta jest bardzo dobrze znana w literaturze jako karta do oceny częstości zdarzeń, na przykład napotkania wadliwych egzemplarzy produktu. Tutaj służy nam ona do objaśnienia pomysłu modyfikacji jej tak, by stała się kartą dla zmiennej mierzalnej, ale odporną na zakłócenia i wrażliwą na małe odchylenia wartości m .

Powtarzając przedstawione rozumowanie, możemy otrzymać analogiczną kartę dla liczby obserwacji, których spodziewamy się w zadanych granicach, gdy $m = 0$. Linie graniczne takiej karty, kontrolującej wartości I_N , wyznaczamy następująco:

$$N p_0 \pm k \sqrt{N p_0 (1 - p_0)}, \quad (14.1.9)$$

gdzie k wybieramy według tych samych zasad, które już opisano. Karta ta służy zwykle do kontrolowania liczby egzemplarzy mających defekt.

Opis karty zmodyfikowanej Jeśli opisana wyżej karta ma być przydatna do wykrywania małych zmian jakości, to wymaga ona modyfikacji. Modyfikacja ta jest potrzebna, gdyż szerokość przedziału między granicami kontrolnymi (14.1.9) wynosi $k\sqrt{N}$ i rośnie wraz ze wzrostem liczby obserwacji N . Jeśli skok o wartości m wystąpi po długim okresie stanów, które mieszczą się wewnątrz granic jakości, to odstęp między tymi granicami narodzi się do znacznych wartości i potrzebny będzie długi czas oczekiwania na skumulowanie się skutków skoku m aż do chwili kiedy nastąpi zdarzenie

$$|N/2 - I_N| > k \sqrt{N}/2. \quad (14.1.10)$$

Z drugiej strony, jeśli skok wystąpi w czasie niezbyt odległym od chwili rozpoczęcia działania karty, to N jest względnie małe i wzrośnie prawdopodobieństwo, że zdarzenie (14.1.10) nastąpi w krótkim czasie.

Spostrzeżenia te sugerują następujący sposób modyfikacji karty (14.1.9). Zamiast pamiętać i zliczać wszystkie obserwacje Z_1, Z_2, \dots, Z_N od startu karty do chwili N -tej, utworzymy bufor (listę obserwacji) o długości $M > 1$, w którym zapamiętane będzie M ostatnio zaobserwowanych wartości Z_i . W momencie, gdy następna obserwacja – powiedzmy $(n+1)$ -sza – stanie się dostępna, wówczas Z_n zastępowane jest przez Z_{n+1} , a Z_n zajmuje miejsce zajmowane dotąd przez Z_{n-1} itd. W każdej kolejnej chwili n zawartość bufora jest badana, aby sprawdzić, czy proces nadal znajduje się w normalnej pracy. Jeśli tak jest, to pobierana jest następna obserwacja, jeśli zaś nie – to deklarowane jest przekroczenie dopuszczalnych poziomów jakości i podejmowane są działania korygujące.

Dokładniej, oznaczmy przez J_n liczbę dodatnich obserwacji w buforze w chwili, gdy n -ta obserwacja została już dokonana, tzn.

$$J_n = \text{card}\{Z_i = 1, i = n, (n-1), \dots, n - (M-1)\}. \quad (14.1.11)$$

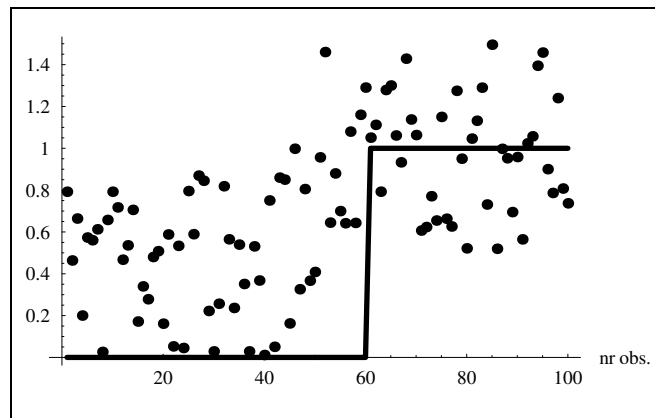
Powyższy wzór nie wymaga komentarza, gdy $n = M, (M+1), \dots$. Jeśli natomiast $n < M$, wystąpi efekt brzegowy, polegający na formalnym braku obserwacji, które mamy wstawić do (14.1.11). Dalej będziemy zakładać, że dostępne są „dobre” dane z wcześniejszych etapów pracy tego procesu, kiedy to znajdował się on w stanie poprawnej pracy. Te właśnie dane mogą być użyte do załadowania bufora w fazie rozruchu karty. Jeśli danych takich nie posiadamy, na przykład w trakcie symulacyjnych badań karty, to rozsądną alternatywą jest wstępne zapełnienie bufora zerami i jedynkami w liczbie i kolejności, jakiej dostarcza ich generator liczb pseudolosowych, który obu tym wartościom przypisuje prawdopodobieństwo wystąpienia równe $1/2$.

Dalej zakładamy, że w chwili $n = -1$ bufor zawiera „historyczne” lub sztucznie wygenerowane obserwacje o numerach Z_{-1}, \dots, Z_{-M} . Formalnie startujemy nanoszenie nowych danych na kartę w chwili $n = 0$, gdy pojawi się obserwacja Z_0 . Dolną i górną linię kontrolną definiujemy przez analogię z (14.1.9), kładąc $p_0 = 1/2$ i wstawiając M na miejsce N , gdyż zawsze tylko M ostatnich obserwacji bierzemy pod uwagę. W ten sposób otrzymujemy

$$UCL = M p_0 + k \sqrt{M p_0 (1 - p_0)} = M/2 + k \sqrt{M}/2, \quad (14.1.12)$$

$$LCL = M p_0 - k \sqrt{M p_0 (1 - p_0)} = M/2 - k \sqrt{M}/2. \quad (14.1.13)$$

Gdy $J_n > UCL$ lub $J_n < LCL$, deklarujemy wyjście poza stan poprawnego funkcjonowania procesu, w przeciwnym razie obserwacje są kontynuowane, a zawartość bufora podlega aktualizacji. Zasygnalizowanie wyjścia procesu poza stan poprawnego funkcjonowania, zwane też alarmem, powinno wywołać czynności zmierzające do przywrócenia stanu normalnej pracy. Za każdym razem po usunięciu przyczyn alarmu, karta powinna być tworzona od początku, tzn. od zapełnienia bufora obserwacjami pochodzącymi z etapów poprawnej pracy procesu.



Rys. 14.1. Przykład działania karty — objaśnienia w tekście

W celu wstępnego zilustrowania działania karty rozważmy następujący przykład symulacyjny. Obserwacje poprawnych stanów funkcjonowania procesu pochodzą z rozkładu równomiernego na odcinku $[0, 1]$, a jako poziom odniesienia przyjęto $Y = 0.5$. Po $q = 50$ obserwacjach tego procesu symulowano jego przejście w stan nieprawidłowej pracy. Trwała zmiana wartości średniej wynosiła $m = 0.5$ (por. rys. 14.1). W celu jej wykrycia prowadzono kartę kontrolną z granicami (14.1.12), (14.1.13) przy $M = 12$ i $k = 2.31$. Wybór tej ostatniej wartości poddyktowany został wynikami badań, które przytaczamy w dalszej części tego rozdziału. W wyniku działania karty zmiana została wykryta w sześćdziesiątej obserwacji, a więc opóźnienie wyniosło dziesięć jednostek czasu. Podkreślić należy, że jest to tylko jeden z możliwych scenariuszy sekwencji obserwacji i działania tej karty. Jej uśrednione zachowanie omawiamy nieco później. Różni się ona od znanych kart kontrolnych pod następującymi względami:

- W odróżnieniu od karty (14.1.9) dla liczby defektów, proponowana modyfikacja charakteryzuje się stałą, niezależną od czasu pracy karty, odległością między dolną i górną granicą kontrolną. W wyniku tej własności, wrażliwość karty na małe zmiany jakości procesu nie maleje wraz ze wzrostem liczby obserwacji.
- Zmodyfikowana karta ma dwa parametry M i k , co pozwala na dokładniejsze jej dostrojenie do konkretnych wymagań.
- Wzór pozwala porównać kartę zmodyfikowaną z kartą CUSCORE (por. [92]). Karta CUSCORE przypisuje „punkty” o wartościach ± 1 tylko za dostatecznie duże (większe od pewnego $A > 0$) odchylenia, w górę lub w dół, od wartości zadanej. Sumę punktów porównuje się z wartością progową. Odpowiedni dobór A wprawdzie stabilizuje wariancję czasu, który upływa do fałszywego alarmu, ale jednocześnie karta ta – z powodu swej konstrukcji – albo nie wykryje trwałej zmiany jakości o wartości poniżej A , albo wykryje ją z bar-

dzo dużym opóźnieniem. Proponowana karta zmodyfikowana sumuje każdą, nawet najmniejszą odchyłkę od wartości nominalnej, natomiast stabilizację wariancji uzyskuje się, dobierając odpowiednio długość bufora. Jednakże zbyt duża długość bufora także niekorzystnie wpływa na wykrywanie małych zmian jakości.

Dobór parametrów karty zmodyfikowanej. Proponowana karta ma dwa parametry M i k , które należy starannie wybrać, by dopasować algorytm detekcji do warunków konkretnego procesu. Zauważmy, że nawet najlepsza i najlepiej zestrojona karta kontrolna, pracująca w warunkach oddziaływania losowych zakłóceń, będzie popełniać dwa rodzaje błędów. Mianowicie, będzie czasem sygnalizować wyjście wskaźnika jakości procesu poza granice kontrolne mimo, że pogorszenie jakości nie nastąpiło. Przyczyną takiego zachowania karty są pojawiające się duże wartości zakłóceń. Prawdopodobieństwo ich wystąpienia jest zwykle małe, lecz niezerowe. Sytuację taką nazywać będziemy fałszywym alarmem. Drugi rodzaj błędów, to opóźnienia w zasygnalizowaniu pogorszenia wskaźnika jakości w stosunku do chwili, w której wystąpiły. Zarówno dla teorii, jak i dla zastosowań, ważna jest średnia wartość opóźnienia w zasygnalizowaniu powyższego pogorszenia. Wartość tę oznaczamy będziemy przez Out-C-ARL (od angielskiego terminu *out of control average run length*). Podobnie, ważny jest średni czas przebiegu do zasygnalizowania przez kartę fałszywego alarmu. Czas ten oznaczamy będziemy jako In-C-ARL (od angielskiego terminu *in-control average run length*).

Jak już wspominaliśmy, nie jest możliwa jednoczesna minimalizacja obu tych średnich czasów i dlatego zwykle nakłada się ograniczenie, że In-C-ARL powinien mieć wartość nie mniejszą niż pewna dolna granica, natomiast dąży się do skrócenia Out-C-ARL. Dla proponowanej karty trudno jest znaleźć analityczne zależności In-C-ARL i Out-C-ARL od M i k , obowiązujące dla skończonej liczby obserwacji. Wyprowadzenie asymptotycznych odpowiedników tych zależności wymaga zastosowania równań Chapmana–Kołmogorowa wyprowadzanych dla łańcuchów Markowa. Wyprowadzenia takie leżą poza zakresem tej książki i dlatego ograniczymy się do opisu procedury strojenia karty w oparciu o wielokrotne symulacje, które są wprawdzie czasochłonne, ale wystarczy wykonać je raz, a wyniki zestawić w formie tabelarycznej. Ponadto, wynikające z tych badań wskazania nie mają asymptotycznego charakteru.

Poniżej przytaczamy uwagi i podpowiedzi dotyczące strojenia omawianej karty kontrolnej. Zostały one zebrane przez autora w trakcie badań symulacyjnych.

Uwaga 1. Standardowy wybór $k = 3$, prowadzący do granic kontrolnych $\pm 3\sigma$ nie jest – na ogół – zalecany, gdyż prowadzi do nadmiernie długich In-C-ARL, a przez to również do wydłużenia średniego czasu do wykrycia pogorszenia jakości. Oczywiście In-C-ARL zależy także od długości bufora M , ale nawet dla niezbyt dużych M wartość In-C-ARL przekracza 1000. Skrajny przypadek osiąga się dla $M = 9$ i $k = 3$. Wówczas $LCL = 0$ a $UCL = 9$, co prowadzi

do nieskończenie długiego In-C-ARL, gdyż warunek $J_n \leq 9$ jest zawsze spełniony, a więc granice kontrolne nigdy nie będą naruszone. Jednakże, jeśli dla tego samego $M = 9$ wybierzemy $k = 2.34$, to otrzymamy często używaną w praktyce wartość In-C-ARL, wynoszącą około 500.

Uwaga 2. Dla ustalonej pojemności bufora M ta sama wartość In-C-ARL osiągnięta jest dla różnych wartości k z pewnego dość szerokiego przedziału. Powodem takiego szczególnego zachowania się tej karty jest fakt, że J_n przyjmować może tylko wartości całkowite i dlatego zmiany k , a zatem także UCL i LCL, nie zawsze prowadzą do zmiany In-C-ARL. Opisane zachowanie się karty zilustrowano na rysunku 14.2, na którym pokazano logarytm In-C-ARL jako funkcję k dla różnych wartości M .

Uwaga 3. Analiza rysunku 14.2 wskazuje, że nie jesteśmy w stanie osiągnąć dowolnego zadanego poziomu In-C-ARL. Jednakże odstęp między osiągalnymi poziomami są na tyle małe, że fakt ten nie powinien nastęrczać trudności w zastosowaniach.

Uwaga 4. Ten sam wykres sugeruje, że korzystne jest wybieranie wartości parametru k blisko lewego krańca przedziału, w którym osiągany jest wybrany przez nas poziom In-C-ARL. Uzasadnienie tej sugestii wynika z faktu, że – zachowując ten sam poziom In-C-ARL – redukujemy jednocześnie szerokość przedziału między LCL i UCL, co prowadzi do skrócenia średniego czasu do wykrycia pogorszenia jakości.

Uwaga 5. W świetle uwag 2–4 zaproponować można następującą procedurę dostrajania karty do konkretnych warunków.

1. Wybrać pożądaną wartość średnią czasu do fałszywego alarmu. Często In-C-ARL wybiera się w przedziale 400–500.
2. Wybrać pojemność bufora $M > 1$, biorąc pod uwagę, że zbyt duże wartości M redukują wrażliwość karty na duże zmiany jakości. Z drugiej strony, większe wartości M zmniejszają średni czas reakcji karty na małe (rzędu 0.1σ – 0.25σ) zmiany jakości. Szerszą dyskusję na temat doboru M przedstawiamy w dalszej części rozdziału, po prezentacji wyników badań symulacyjnych.
3. Przeprowadzić badania symulacyjne do oceny In-C-ARL w zależności od k , zmieniając k w zakresie $[1, 3]$ i dobrać takie k , by zapewnić pożądaną wartość In-C-ARL dla ustalonego M . Dalsze sugestie na temat metodyki symulacji i podpowiedzi dotyczące doboru k przedstawiamy dalej.

14.2. Dostrajanie karty i wyniki porównań

Znalezienie przybliżonej wartości k w opisany wyżej sposób nie jest trudne. Zadanie staje się trudniejsze, gdy chcemy precyzyjnie zlokalizować najmniejsze k , które zapewnia pożądaną wartość In-C-ARL. Lokalizacja k z dokładnością do

$M =$	12	23	28	71	90	150	212	441
$k =$	2.31	2.30	2.27	2.02	2.0	1.8	1.65	1.39
In-C-ARL =	395	415	423	411	450	452	440	456

Tabela 14.1. Zestawy parametrów (M, k) karty kontrolnej, które zapewniają poziom średniego czasu do fałszywego alarmu rzędu 435

0.01 wymaga kilkunastu minut obliczeń na komputerze z procesorem Pentium IV 2.4 GHz, jeśli dla każdej wartości k z pewnego przedziału uśredniamy przebiegi 10 000–30 000 razy. Tak duże krotności uśrednień są niezbędne wobec znacznej zmienności czasów, które wpływają do wystąpienia fałszywego alarmu w poszczególnych przebiegach. Warto dodać, że cecha ta nie dotyczy wyłącznie badanej tutaj karty.

Szczegółowe uwagi na temat dostrajania karty

Omówiony w poprzednim podrozdziale przepis na dobór parametrów karty został zastosowany wielokrotnie zarówno w celu dostrojenia do średnich (rzędu 400), jak i dużych (rzędu 800) wartości In-C-ARL.

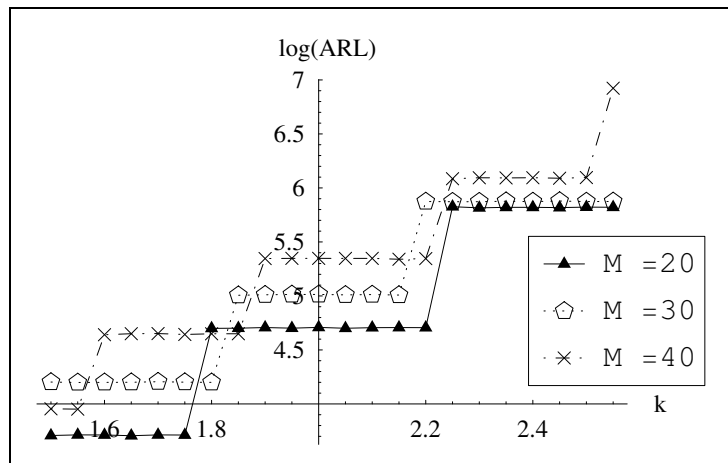
Z myślą o potencjalnych zastosowaniach, w tabelicy 14.1 zebrano pary (M, k) – wielkości bufora i progę karty, które zapewniają średni czas do wystąpienia fałszywego alarmu na poziomie 435. Dobór wartości k prowadzono z dokładnością do 0.01, natomiast wartości In-C-ARL, które się uzyskuje, podano w ostatnim wierszu tej tabeli. Nie można oczekiwać większych dokładności, gdyż J_n przyjmuje wartości całkowite.

Gdyby informacje podane w tabelicy 14.1 okazały się niewystarczające, należy przeprowadzić opisane obliczenia dla wybranej wartości i znaleźć odpowiednią wartość progę k . Poszukiwania takie można przyspieszyć, posługując się empirycznym wykresem pokazanym na rysunku 14.4. Dla wybranej pojemności bufora M wykres ten pozwala odczytać taką przybliżoną wartość k , która zapewnia In-C-ARL rzędu 435.

Dalszym ułatwieniem może być empiryczny wzór

$$\log(\log(k)) = -0.137 - 0.0026 M, \quad (14.2.14)$$

który jest liniowym przybliżeniem zależności z rysunku 14.4. Wzór ten otrzymano jako regresję liniową (estymowaną metodą minimalizacji sumy kwadratów błędów) między wartościami $\log(\log(k))$ oraz M w punktach oznaczonych grubymi kropkami na rysunku 14.4. Wykresy na rysunku 14.3 pozwalają odczytać dalsze sugestie dotyczące wyboru pojemności bufora M . Na wykresach tych pokazano zależność średniego opóźnienia od momentu wystąpienia skokowej zmiany jakości do chwili jego wykrycia przez kartę kontrolną (Out-C-ARL) w funkcji M .



Rys. 14.2. Zależność logarytmu In-C-ARL od progu k dla różnych wartości M (wyniki otrzymano dla zakłóceń o rozkładzie $N(0, 1)$, w wyniku uśrednienia 10^4 przebiegów)

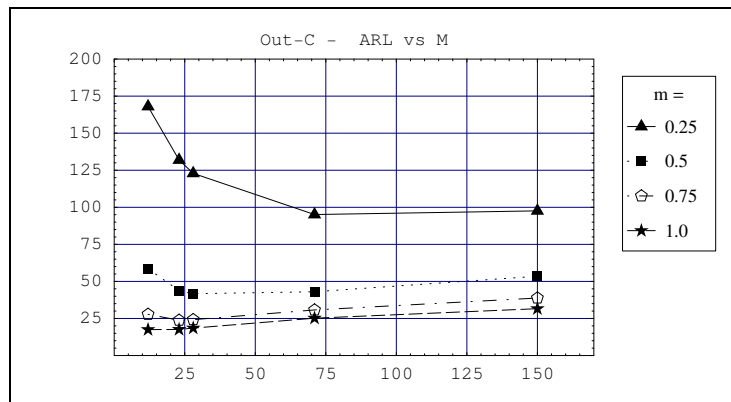
Poszczególne wykresy ilustrują jak przebiega ta funkcja w zależności od wysokości skokowej zmiany jakości (parametr m). Można zauważyć, że dla $m = 0.25$, $m = 0.5$ i $m = 0.75$ istnieje optymalna wartość M , dla której Out-C-ARL osiąga minimum. Mimo że przebieg zależności Out-C-ARL od M jest dość płaski w otoczeniu minimum, to odczytać można, że dla $m = 0.25$, $m = 0.5$ i $m = 0.75$ optymalne wartości M wynoszą, odpowiednio, $M = 71$, $M = 28$ i $M = 23$.

Zauważmy, że dla $m = 1$ wykres jest funkcją rosnącą i można się spodziewać, że optymalna wielkość bufora byłaby mniejsza niż 12, ale – jak stwierdziliśmy wcześniej – dla $M < 12$ nie da się dobrać takiej wartości k , by In-C-ARL był rzędu 400.

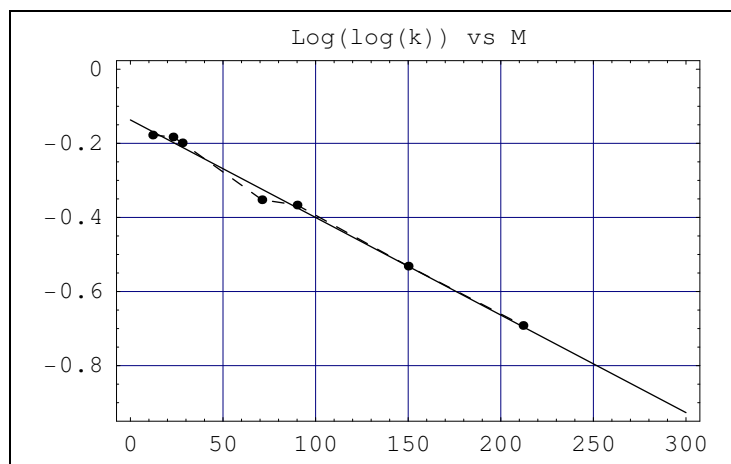
Na drugim krańcu pozostaje przypadek $m = 0.1$, który nie został przedstawiony na omawianych wykresach. Przebieg zależności Out-C-ARL od M w tym przypadku nie ma minimum, ale maleje i dla $M = 436$ otrzymuje się dużą wartość Out-C-ARL równą 219.8. Z tego powodu stosowanie buforów o pojemności M przekraczającej 100 wydaje się niecelowe. Bezpośrednie korzystanie z optymalnych długości buforów też nie jest możliwe, gdyż zwykle nie wiemy jakiej wysokości zmiana jakości może się pojawić. Możemy jednak „uczulić” kartę na skoki o mniejszych wartościach, wybierając M rzędu 70. Jeśli natomiast wybierzemy M rzędu 25, to karta będzie szybciej reagować na skoki o większych amplitudach.

Wyniki badań symulacyjnych

Badania symulacyjne zmodyfikowanej karty kontrolnej miały kilka celów:



Rys. 14.3. Średnie opóźnienie wykrycia skoku o wartości m w funkcji M (zakłócenia o rozkładzie normalnym)



Rys. 14.4. Wykres ułatwiający dobór progu karty k w zależności od wybranej wielkości bufora M tak, by uzyskać In-C-ARL rzędu 435

- Znalezienie par (pojemność bufora, próg karty), które zapewniają pożądany poziom In-C-ARL.
- Zebranie obserwacji opóźnień karty w wykrywaniu zmian jakości o różnym poziomie przy założeniu, że zakłócenia mają rozkład normalny.
- Przetworzenie wyżej wymienionych obserwacji w celu porównania Out-C-ARL karty zmodyfikowanej ze znanymi kartami opracowanymi dla zakłóceń o rozkładzie normalnym.

- Zebranie danych i gromadzenie doświadczeń o zachowaniu się Out-C-ARL karty zmodyfikowanej, gdy zakłócenia mają rozkłady inne niż rozkład normalny.

Działanie zmodyfikowanej karty. Zakłócenia gaussowskie – In-C-ARL \approx 435

W tabelach 14.2 i 14.3 zestawiono średnie czasy do wykrycia skoku (kolumna oznaczona ARL), uzyskane za pomocą zmodyfikowanej karty kontrolnej dla małych, średnich i dużych pojemności bufora M . Dla każdej z wartości M podano też stosowaną wartość progu k , która zapewnia In-C-ARL rzędu 435. Kontrolnie w pierwszym wierszu podano faktycznie uzyskany In-C-ARL, który odpowiada wartości skoku $m = 0$. Kolejne wiersze zawierają wartości średnich czasów do wykrycia skoku (Out-C-ARL), które uzyskano dla skokowych zmian jakości o wartościach m z przedziału $[0.1, 3]$, z krokiem 0.25. Tabele te zawierają też kolumnę oznaczoną symbolem Dysp. W kolumnie tej zebrano empiryczne oszacowania dyspersji czasu, który upływa od chwili wystąpienia skoku do momentu jego wykrycia przez kartę.

Wszystkie badania symulacyjne zawarte w tych tabelach prowadzono przy symulowanych zakłóceniach o rozkładzie normalnym $N(0, 1)$, a więc wartości skokowych zmian interpretować można w jednostkach $m\sigma$, gdzie σ oznacza dyspersję zakłóceń.

Analiza tych tabel sugeruje następujące wnioski:

1. Zaobserwować można spodziewane skrócenie Out-C-ARL wraz ze wzrostem wysokości skoku m . Jednakże po przekroczeniu wartości $m = 1.5$ dalsza redukcja Out-C-ARL jest nieznaczna.
2. Wraz ze wzrostem skoku m znacząco maleje też dyspersja czasu do wykrycia skoku. Jest to bardzo korzystna własność tej karty, gdyż faktyczne czasy do wykrycia skoku „skupione” są bliżej wartości Out-C-ARL. Jednakże tempo redukcji dyspersji maleje dla $m > 1.5$.
3. Tabele potwierdzają wnioski dotyczące wyboru M , które przedstawiono na stronie 183.

Porównamy teraz działanie proponowanej karty z klasycznymi kartami EWMA i CUSUM. Dane do porównań zaczerpnięto z niedawno wydanej pracy [52], która zawiera między innymi rezultaty starannie opisanych badań symulacyjnych wymienionych kart. Badania autora wykonane zostały z zachowaniem tej samej metodologii dla In-C-ARL rzędu 435 z uwzględnieniem zakłóceń o rozkładzie normalnym. A oto wnioski z tych porównań:

1. Dla $m = 0.1\sigma$ w pracy [52] otrzymano Out-C-ARL = 297 dla karty EWMA oraz Out-C-ARL = 295 dla karty CUSUM. Z tabeli 14.3 w odpowiednim wierszu możemy odczytać, że dla karty zmodyfikowanej Out-C-ARL wynosi 254.9, 243.5 i 234.3, w zależności od pojemności bufora M (równego odpowiednio 71, 150, 212). Czasy te są więc znacząco krótsze.

2. Dla $m = 0.25\sigma$ w cytowanej pracy otrzymano $\text{Out-C-ARL} = 110$ dla karty EWMA oraz $\text{Out-C-ARL} = 132$ dla karty CUSUM. Natomiast proponowana karta pozwala uzyskać czasy od 95.1 do 101.2 (por. trzeci wiersz tab. 14.3). Są one zauważalnie krótsze, a średnie zmniejszenie czasu do wykrycia skoku wynosi od 10 do 30 procent.
3. W całym zakresie małych zmian jakości $m \in [0.1\sigma, 0.25\sigma]$ również dyspersja czasów do wykrycia skoku przemawia na korzyść karty zmodyfikowanej. Gdy $m = 0.1\sigma$, to dyspersja ta dla kart EWMA i CUSUM wynosi, odpowiednio, 288 i 323 (por. [52]). Przy tej samej wartości skoku karta zmodyfikowana zapewnia dyspersję od 167 do 183, w zależności od wybranego bufora (por. drugi wiersz tabeli 14.3). Analogicznie, gdy $m = 0.25\sigma$, to dyspersja dla kart EWMA i CUSUM wynosi, odpowiednio, 102 i 123, podczas, gdy dla proponowanej tu karty otrzymujemy dyspersję 61.
4. W zakresie średnich zmian jakości $m \in [0.5\sigma, 1.0\sigma]$ karty klasyczne zapewniają krótsze średnie czasy do wykrycia skoku niż karta zmodyfikowana. Przykładowo, gdy $m = 0.5\sigma$, to Out-C-ARL dla karty EWMA wynosi 32.4, a dla karty CUSUM w [52] znajdujemy 37.2. W tych samych warunkach, najlepszy czas odczytany z tabeli 14.3 wynosi 43 (dla $M = 71$). Przewaga kart klasycznych rośnie ze wzrostem m i gdy $m = 1.0$, to zapewniają one Out-C-ARL około 10σ (dyspersja około 5), podczas gdy dla karty zmodyfikowanej otrzymujemy wartość 25 (dyspersja 13).
5. W zakresie dużych zmian jakości ($m \in [1.25\sigma, 3.0\sigma]$) przewaga kart klasycznych jest już bardzo znacząca. Przykładowo dla $m = 2.0\sigma$ karty te zapewniają Out-C-ARL około 4 (dyspersja 1.3). a karta zmodyfikowana pozwala uzyskać wartość 18.6 (dyspersja 9.5) dla $M = 71$, a w najlepszym przypadku Out-C-ARL wynosi 9.46 (dyspersja 2.9) dla $M = 12$.

Działanie zmodyfikowanej karty. Zakłócenia gaussowskie — $\text{In-C-ARL} \approx 840$

Porównamy teraz działanie karty zmodyfikowanej z kartami EWMA i CUSUM w przypadku, gdy założony średni czas do fałszywego alarmu (In-C-ARL) jest długi i wynosi ponad 800 obserwacji (w trakcie badań symulacyjnych stosowano wartości 840–860). Wyniki badań symulacyjnych zebrano w tabeli 14.4. Z jakościowego punktu widzenia obraz jest podobny do omówionego wyżej przypadku średnich wartości $\text{In-C-ARL} \approx 435$. Dlatego przejdziemy od razu do krótkiego porównania ilościowego omawianej karty i kart klasycznych. Źródłem danych o średnich czasach reakcji tych ostatnich jest cytowana już praca [52].

1. Podobnie jak poprzednio, w zakresie zmian jakości o małej amplitudzie, karta zmodyfikowana wykrywa je szybciej niż karty klasyczne. Na przykład, gdy $m = 0.1\sigma$, to
 - Out-C-ARL dla karty EWMA wynosi 524,

$M=12, k=2.31$			$M=23, k=2.3$			$M=28, k=2.27$		
Skok	ARL	Dysp.	Skok	ARL	Dysp.	Skok	ARL	Dysp.
0	395.27	171.09	0	415.66	181.42	0	423.12	185.75
0.1	328.33	144.18	0.1	305.80	133.14	0.1	303.43	133.17
0.25	168.09	72.47	0.25	131.89	56.00	0.25	122.90	51.72
0.5	58.65	24.52	0.5	43.78	17.08	0.5	41.66	15.73
0.75	27.84	10.91	0.75	23.76	8.35	0.75	24.18	8.27
1	17.51	6.35	1	17.60	5.76	1	18.53	6.00
1.25	12.98	4.41	1.25	14.99	4.76	1.25	16.10	5.12
1.5	10.96	3.54	1.5	13.66	4.31	1.5	14.76	4.68
1.75	10.00	3.14	1.75	12.88	4.05	1.75	13.99	4.42
2	9.46	2.94	2	12.45	3.91	2	13.54	4.28
2.25	9.19	2.84	2.25	12.26	3.84	2.25	13.25	4.18
2.5	9.09	2.80	2.5	12.12	3.80	2.5	13.14	4.14
2.75	9.05	2.79	2.75	12.04	3.77	2.75	12.96	4.09
3	9.01	2.77	3	11.96	3.75	3	13.09	4.12

Tabela 14.2. Średnie czasy do wykrycia skoku dla zmodyfikowanej karty, nastrojonej na In-C-ARL ≈ 435 (bufory o małej pojemności M , zakłócenia o rozkładzie normalnym)

- karta CUSUM reaguje w czasie 592,
 - karta zmodyfikowana reaguje szybciej i zapewnia Out-C-ARL nieco poniżej 400, gdy $M = 111$ lub $M = 131$ (por. drugi wiersz tabeli 14.4).
2. Powyżej $m = 0.5\sigma$ szybsze czasy reakcji zapewniają karty klasyczne. Przykładowo, gdy $m = 1.0\sigma$, to Out-C-ARL wynosi 34, ale dla kart EWMA i CUSUM otrzymujemy czasy rzędu 11.
 3. Jakościowo podobnie zachowuje się dyspersja czasów do wykrycia skoku, tzn. w obszarze $m < 0.5\sigma$ jest ona znacząco mniejsza dla karty zmodyfikowanej. Gdy $m > 0.5\sigma$, dyspersja kart klasycznych jest mniejsza niż karty zmodyfikowanej.

Karta zmodyfikowana przy zakłóceniach niegaussowskich, In-C-ARL ≈ 435

Celem tego podrozdziału jest podsumownie badań symulacyjnych, działania karty zmodyfikowanej, gdy skokowe zmiany jakości obserwowane są w obecności zakłóceń o rozkładach innych niż normalny. Do badań wybrano rozkład Laplace'a (podwójnie wykładniczy) o wartości oczekiwanej zero i dyspersji 1 oraz rozkład Cauchy'ego. Ponieważ wartość oczekiwana i wariancja tego rozkładu nie istnieją,

$M=71, k=2.02$			$M=150, k=1.8$			$M=212, k=1.65$		
Skok	ARL	Dysp.	Skok	ARL	Dysp.	Skok	ARL	Dysp.
0	411.23	301.39	0	452.05	337.19	0	440.32	334.70
0.1	254.91	182.71	0.1	243.54	172.58	0.1	234.27	166.92
0.25	95.12	60.68	0.25	97.58	58.68	0.25	101.26	60.62
0.5	43.03	23.33	0.5	53.50	29.52	0.5	56.87	32.41
0.75	30.75	16.08	0.75	38.80	21.12	0.75	41.30	23.18
1	25.23	13.06	1	31.60	17.07	1	33.77	18.76
1.25	22.11	11.39	1.25	27.71	14.87	1.25	29.38	16.24
1.5	20.28	10.41	1.5	25.20	13.50	1.5	26.92	14.81
1.75	19.22	9.83	1.75	23.82	12.74	1.75	25.38	13.93
2	18.61	9.50	2	23.10	12.30	2	24.57	13.48
2.25	18.13	9.25	2.25	22.64	12.05	2.25	24.17	13.19
2.5	17.91	9.14	2.5	22.31	11.86	2.5	23.76	13.00
2.75	17.83	9.08	2.75	22.17	11.80	2.75	23.70	12.95
3	17.69	9.03	3	22.15	11.77	3	23.70	12.94

Tabela 14.3. Średnie czasy do wykrycia skoku przez zmodyfikowaną kartę, nastrojoną na In-C-ARL ≈ 435 (bufory o średniej i dużej pojemności M , zakłócenia gaussowskie)

więc – w celu zachowania podobnych cech tego rozkładu – wybrano go tak, by był symetryczny względem zera i by różnica między kwantylami rzędu $3/4$ i rzędu $1/4$ była taka jak dla rozkładu $N(0, 1)$.

Wyniki badań karty przy takich zakłóceniach zestawiono w tabeli 14.5. Pozwala ona stwierdzić, że zmodyfikowana karta działa poprawnie w zakresie małych i średnich zmian jakości. Fakt, że rozkład Cauchy’ego ma nieskończoną wariancję, zauważalnie wydłuża średni czas reakcji karty na wystąpienie skoku, ale karta nadal działa poprawnie w tak trudnych warunkach, a spowolnienie reakcji w porównaniu z sytuacją, gdy występują zakłócenia gaussowskie jest mniejsze niż można byłoby oczekiwać.

Proponowana karta może skutecznie konkurować ze znanymi kartami w wykrywaniu małych zmian jakości (w zakresie do 0.25σ) nawet wówczas, gdy zakłócenia mają rozkład normalny. Jest ona skuteczna w wykrywaniu takich zmian nawet wówczas, gdy rozkład zakłóceń nie jest gaussowski i ma nieskończoną wariancję. Dodatkową zaletą jest to, że rozkład ten może nie być znany.

$M= 111 , k =2,19$			$M= 131 , k =1,84$			$M= 453 , k =1,35$		
Skok	ARL	Dysp.	Skok	ARL	Dysp.	Skok	ARL	Dysp.
0	836,64	370,04	0	841,83	370,73	0	840,02	650,13
0,1	398,04	170,64	0,1	399,22	171,86	0,1	337,79	226,64
0,25	122,34	46,80	0,25	124,06	46,98	0,25	149,54	87,46
0,5	57,56	19,21	0,5	57,63	19,18	0,5	82,94	47,32
0,75	41,71	13,77	0,75	42,10	13,83	0,75	59,05	33,34
1	34,13	11,17	1	34,18	11,20	1	48,22	27,03
1,25	29,78	9,72	1,25	29,60	9,66	1,25	42,12	23,43
1,5	27,12	8,84	1,5	27,18	8,84	1,5	38,14	21,23
1,75	25,66	8,35	1,75	25,61	8,32	1,75	36,05	20,03
2	24,76	8,04	2	24,66	8,02	2	34,93	19,38
2,25	24,24	7,88	2,25	24,21	7,87	2,25	34,33	18,98
2,5	24,00	7,78	2,5	23,96	7,77	2,5	33,79	18,71
2,75	23,97	7,78	2,75	23,78	7,74	2,75	33,36	18,48
3	23,79	7,73	3	23,69	7,70	3	33,55	18,57

Tabela 14.4. Średnie czasy do wykrycia skoku przez zmodyfikowaną kartę, In-C-ARL ≈ 840 (bufory o średniej i dużej pojemności, zakłócenia o rozkładzie normalnym)

14.3. Metodyka wykrywania zmian w sekwencjach obrazów

Przedstawiany w tym podrozdziale tok postępowania przy wykrywaniu zmian w sekwencjach obrazów ukierunkowany jest na potencjalne zastosowania w wykrywaniu zmian jakości z użyciem kamer przemysłowych. Przedstawimy teraz proponowane etapy postępowania, które łącznie prowadzą do opracowania metody wykrywania zmian jakości na podstawie sekwencji obrazów lub sygnałów.

Akwizycja i wstępna obróbka obrazów. Obrazy pozyskiwane w warunkach przemysłowych mogą charakteryzować się znacznymi zniekształceniami i błędami. Ich źródłem mogą być na przykład:

- drgania powietrza wywołane wysoką temperaturą odlewów,
- drobiny pyłu i kurzu unoszące się w powietrzu,
- zmienność warunków oświetlenia (światło dzienne o różnym natężeniu, światło sztuczne).

Laplace (DbExp)			Cauchy		
$M=40, k=2.22$			$M=28, k=2.28$		
Skok	ARL	Dysp.	Skok	ARL	Dysp.
0	437.69	315.91	0	420.79	300.12
0.1	191.35	133.31	0.1	334.82	240.04
0.25	59.51	37.02	0.25	167.28	116.28
0.5	28.51	15.02	0.5	64.17	41.76
0.75	22.04	11.13	0.75	37.52	22.47
1	19.33	9.69	1	27.27	15.21
1.25	17.70	8.84	1.25	22.70	12.03
1.5	16.85	8.37	1.5	20.51	10.58
1.75	16.15	8.02	1.75	18.86	9.55
2	15.78	7.83	2	17.93	9.00
2.25	15.55	7.69	2.25	17.23	8.58
2.5	15.34	7.59	2.5	16.67	8.28
2.75	15.19	7.52	2.75	16.29	8.07
3	15.07	7.46	3	15.98	7.90

Tabela 14.5. Średnie czasy do wykrycia skoku (Out-C-ARL) uzyskane dla zmodyfikowanej karty kontrolnej, nastrojonej na In-C-ARL ≈ 435 . Stosowano bufor o średniej pojemności M , a zakłócenia miały rozkłady niegaussowskie

Z tych względów każdy napływający obraz powinien być poddany wstępnej obróbce, której celem jest poprawa jakości obrazu. Szczegółowe omówienie tego etapu pozostaje poza zakresem tematycznym tej książki. Do dyspozycji mamy całą gamę technik opracowanych przez teorię przetwarzania obrazów (por. [41]), ze szczególnym uwzględnieniem metod wyspecjalizowanych i dostosowanych do obrazów przemysłowych (por. [81]). Wybierając na tym etapie techniki filtracji i poprawy jakości obrazów, pamiętać trzeba, że nie powinny one być zbyt czasochłonne, gdyż przetwarzanie każdego obrazu musi przebiegać w tempie ich napływania. Ponadto wybrane techniki nie powinny zbyt mocno redukować tych cech obrazu, które w następnych etapach służyć będą do obliczenia syntetycznego wskaźnika jakości. Na przykład, jeśli wskaźnik ten obliczany będzie na podstawie oceny wymiaru fraktalnego, to nie powinniśmy używać filtrów cyfrowych typu ruchomej średniej, gdyż redukują one wymiar fraktalny, a przez to wpływają na wartość wskaźnika.

Wybór cech i próbkowanie obrazów. Jeśli naszym celem jest jedynie wykrywanie zmian jakości, a nie na przykład archiwizowanie pełnych obrazów z przebiegu procesu produkcyjnego, to takie nastawienie pozwala na znaczną redukcję zawartości poszczególnych obrazów i wybieranie z nich tylko tych cech, które są istotne dla wykrycia zmian jakości tego konkretnego produktu, który jest badany. Można to nazwać dedykowanym wyborem cech z obrazów. Na przykład, kamera dozoru na parkingu samochodowym powinna z obserwowanej sekwencji wybierać tylko te cechy, które są istotne dla stwierdzenia, czy któryś z samochodów nie został skradziony. Nie powinna ona traktować jako istotnej zmiany w sekwencji obrazów czynników związanych ze zmianą oświetlenia. W innym przypadku, to właśnie zmiany jasności obrazu, na przykład wlewka miedzianego, są tą właśnie cechą, która nie może zostać zredukowana.

Powyższe szkice przykładów wskazują na to, że sposób wyboru cech i redukcji pozostałej zawartości informacyjnej obrazów nie będzie łatwo poddawał się algorytmizacji i w dalszych rozważaniach musimy przyjąć, że etap ten został już wcześniej zrealizowany. Pomocne w jego realizacji okazać się mogą techniki uczenia i sieci neuronowe. Wydaje się także, że proponowane w poprzednich rozdziałach techniki okonturowywania i empirycznej oceny wymiaru fraktalnego mogą mieć zastosowanie jako sposób wyboru cech dla całych klas zagadnień wykrywania zmian jakości.

Przetwarzanie sygnałów i obrazów we wskaźniki jakości. Problemom takim poświęcone były rozdziały 12 i 13. W rozdziałach tych oceniano procentowy udział zanieczyszczeń i wymiar fraktalny liniowych fragmentów sekwencji obrazów i używano tych wielkości jako wskaźników syntetycznych. Nie są to oczywiście jedyne możliwe wskaźniki, które można stosować do oceny jakości. Teoria przetwarzania sygnałów i obrazów dostarcza wielu innych narzędzi, które mogą być użyteczne. Naturalnymi kandydatami do formowania syntetycznego wskaźnika jakości są współczynniki rozwinięcia sygnału lub obrazu w ortogonalne układy funkcji trygonometrycznych Haara, Walsh, itp. Przy wyborze takiego układu warto brać pod uwagę możliwość zastosowania szybkich algorytmów do obliczenia współczynników rozwinięcia w szereg ortogonalny (np. szybkiej transformaty Fouriera). Istotne są także własności wybranego układu ortogonalnego w zestawieniu z cechami jakości, które mają zostać wykryte. Na przykład, wiadomo, że układ funkcji ortogonalnych na kole, znany jako wielomiany Zernike'a, ma współczynniki rozwinięcia niezmiennicze względem obrotów wokół środka koła. Dzięki niezmienniczości można oceniać jakość wytworzenia przedmiotów położonych przed kamerą bez dbałości o ich precyzyjne ułożenie. Z drugiej jednak strony, jeśli ocenie jakości podlegają na przykład kątowe położenia otworów wywierconych na obrzeżu metalowej

tarczy, to układ Zernike'a nie jest właściwym kandydatem do konstruowania wskaźnika jakości.

Wybór karty kontrolnej i strojenie jej parametrów. Na tym etapie możemy założyć, że posiadamy stale uzupełnianą sekwencję liczb, które charakteryzują jakość. Jak już wskazywano wcześniej, te wskaźniki jakości powstawać mogą w wyniku:

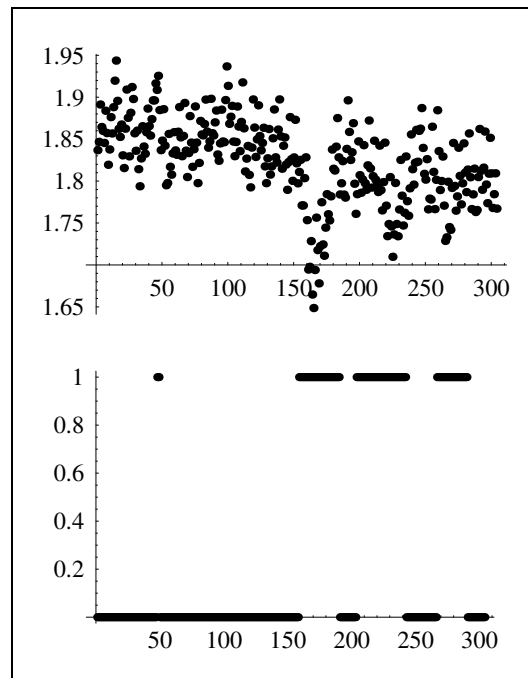
- bezpośredniego pomiaru,
- przetworzenia charakterystyki funkcyjnej wyrobu,
- obróbki sygnału, obrazu lub jego fragmentów.

Pozostaje zatem wykrycie ewentualnej zmiany wskaźnika jakości w sytuacji, gdy wykazuje on także naturalną zmienność, wywołaną czynnikami losowymi. W najprostszych przypadkach, gdy trwała zmiana wskaźnika jakości jest bardzo duża, do wykrycia jej wystarczy posłużyć się zwykłym porównywaniem wskaźnika z wartością progową. Jednakże w sytuacjach trwałych, ale niezbyt dużych zmian wskaźnika, korzystniej jest posłużyć się jedną z kart kontrolnych, na przykład kartą Shewharta, CUSUM lub EWMA. Przedstawione w poprzednim podrozdziale wyniki pozwalają sugerować także użycie proponowanej karty nieparametrycznej, która zapewnia krótsze średnie opóźnienia wykrycia małych zmian jakości w porównaniu z kartami klasycznymi. W przypadku, gdy spodziewać się można zarówno małych, jak i znacznych zmian wskaźnika, korzystne jest równoczesne prowadzenie omawianej karty i jednej z kart klasycznych.

Zastosowanie do oceny jakości powierzchni miedzi

Jako przykładu ilustrującego powyżej zarysowaną metodologię użyjemy danych pochodzących z obrazów opisanych w rozdziale 13. W rozdziale tym opisano próbkowanie obrazu i wykorzystanie lokalnego (liczonego z przekrojów) wymiaru fraktalnego jako wskaźnika jakości. Wystarczy teraz dostarczać wartości syntetycznego wskaźnika jakości jako danych wejściowych dla karty kontrolnej. Ponieważ w dobrze zestrojonym procesie produkcji miedzi spodziewać się można jedynie niewielkich zmian jakości, to jako kartę kontrolną wybrano zmodyfikowaną kartę, która została opisana w poprzednim podrozdziale.

Na rysunku 14.5 pokazano przykład wykrywania nierównomierności powierzchni. W górnej jego części kropki oznaczają oszacowania wymiarów fraktalnych poszczególnych przekrojów obrazu 7, a w jego części dolnej odpowiedzi karty kontrolnej, które uzyskano dla bufora o $M = 12$ komórkach pamięci i progu $k = 2.31$. Oś poziome obu wykresów na rysunku 14.5 są zsynchronizowane, a poziom 1 oznacza, że karta zasygnalizowała wykrycie zmiany jakości. Porównanie obu części tego rysunku wskazuje, że estymacja lokalnego wymiaru fraktalnego w połączeniu z kartą kontrolną bardzo szybko, bo w czasie równym przebiegowi kilkunastu linii obrazu, potrafi wykryć nierównomierność powierzchni.



Rys. 14.5. Wykres górny – lokalne wymiary fraktalne z obrazu 7, wykres dolny – działanie zmodyfikowanej karty kontrolnej

Zauważmy, że jako danych wejściowych dla karty kontrolnej moglibyśmy także użyć innego syntetycznego wskaźnika, a mianowicie, procentowej zawartości obszarów o większej nierównomierności. Przykład takich danych pokazano na rysunku 13.9. Jednakże wówczas ewentualne zasygnalizowanie przez kartę wyjścia ze strefy dopuszczalnych zmian jakości dotyczyłoby całych klatek obrazów, podczas gdy w wersji pokazanej na rysunku 14.5 informację taką uzyskujemy już w trakcie obróbki pojedynczej klatki z sekwencji obrazów.

CZĘŚĆ V

Dodatek

15. Iloczyn Kroneckera i jego własności

Iloczyn Kroneckera jest klasycznym narzędziem algebry (por. [82], [77]). W dodatku tym używać będziemy oznaczeń używanych w algebrze liniowej. x oznacza dowolny wektor kolumnowy, a duże litery A, B, C oznaczać będą macierze kwadratowe. Przedstawiając własności iloczynu Kroneckera, będziemy zakładać, że występujące we wzorach macierze mają rozmiary dobrane tak, by operacje na nich miały sens.

15.1. Definicja i podstawowe własności

Niech $A = [a_{ij}]$ i $B = [b_{ij}]$ będą macierzami o wymiarach $\dim(A) = n_A \times m_A$, $\dim(B) = n_B \times m_B$, odpowiednio.

Definicja 15.1. *Iloczynem Kroneckera macierzy A i B , oznaczanym dalej przez $A \otimes B$, nazywa się macierz $C = A \otimes B$ o postaci*

$$C = \begin{bmatrix} a_{11} B, & a_{12} B, & \dots & a_{1m_A} B \\ a_{21} B, & a_{22} B, & \dots & a_{2m_A} B \\ & & \dots & \\ a_{n_A 1} B, & a_{n_A 2} B, & \dots & a_{n_A m_A} B \end{bmatrix}. \quad (15.1.1)$$

Iloczyn ten jest lewostronnym mnożeniem macierzy A, B w tym sensie, że macierz C budowana jest z kłatek macierzy B , które mnożone są przez poszczególne elementy macierzy A .

Zauważmy, że $\dim(C) = n_C \times m_C$, gdzie $n_C = n_A \cdot n_B$, $m_C = m_A \cdot m_B$. Macierz $A \otimes B$ zbudowana jest następująco: w miejsce każdego elementu macierzy A wstawiana jest klatka o rozmiarach macierzy B , przy czym zawartość każdej z kłatek jest postaci $a_{ij} B$.

Przykładem iloczynu Kroneckera macierzy jest iloczyn wektorów kolumnowych $\mathbf{a} = [a_1, a_2, \dots, a_{n_a}]^T$ i $\mathbf{b} = [b_1, b_2, \dots, b_{n_b}]^T$:

$$\mathbf{c} = \mathbf{a} \otimes \mathbf{b} = [a_1 \mathbf{b}^T, a_2 \mathbf{b}^T, \dots, a_{n_a} \mathbf{b}^T]^T. \quad (15.1.2)$$

Dowody zebranych tu własności iloczynu Kroneckera znaleźć można w [82], [77] i [9]:

$$(A + B) \otimes C = A \otimes C + B \otimes C, \quad (15.1.3)$$

$$A \otimes (B + C) = A \otimes B + A \otimes C, \quad (15.1.4)$$

$$(A \otimes B)^T = A^T \otimes B^T, \quad (15.1.5)$$

$$\text{tr}[A \otimes B] = \text{tr}[A] \cdot \text{tr}[B]. \quad (15.1.6)$$

Jeśli macierze A, B, C, D są takie, że mają sens iloczynu (klasyczne) macierzy $A \cdot C$ oraz $B \cdot D$, to

$$(A \otimes B) \cdot (C \otimes D) = (A \cdot C) \otimes (B \cdot D), \quad (15.1.7)$$

gdzie \cdot oznacza klasyczny iloczyn macierzy. Powyższą własność można uogólnić na ciągi macierzy $A_i, B_i, i = 1, 2, \dots, p$, a mianowicie

$$(A_1 \otimes B_1) \cdot \dots \cdot (A_p \otimes B_p) = \left[\prod_{i=1}^p A_i \right] \otimes \left[\prod_{i=1}^p B_i \right] \quad (15.1.8)$$

Jeśli kwadratowe macierze A, B są nieosobliwe, to również macierz $A \otimes B$ jest nieosobliwa i zachodzi

$$(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}. \quad (15.1.9)$$

Warto zwrócić uwagę, że – w odróżnieniu od odwrotności klasycznego iloczynu macierzy – po prawej stronie (15.1.9) nie następuje zmiana kolejności macierzy A^{-1} oraz B^{-1} .

Iloczyn Kroneckera nie jest przemienny. Jest to operacja łączna, co pozwala zdefiniować iloczyn ciągu macierzy $A_i, i = 1, 2, \dots, p$ następująco:

$$\prod_{i=1}^p A_i = \left[\prod_{i=1}^{p-1} A_i \right] \otimes A_p. \quad (15.1.10)$$

15.2. Wartości własne iloczynu Kroneckera macierzy

Niech $\lambda(A)$ oznacza jedną z wartości własnych kwadratowej macierzy A . Przez $u(A)$ oznaczmy odpowiadający jej (prawostronny) wektor własny, to znaczy, że zachodzi

$$A \cdot u(A) = \lambda(A) u(A). \quad (15.2.11)$$

Jeśli będziemy chcieli wskazać numery wektorów i wartości własnych, to oznaczymy przez $\lambda_i(A), i = 1, 2, \dots, r_A$ wartości własne $r_A \times r_A$ macierzy A , a przez $u_i(A)$,

$i = 1, 2, \dots, r_A$ odpowiadające im (prawostronne) wektory własne. Analogiczne oznaczenia stosować będziemy do innych macierzy kwadratowych, B, C, \dots

Niech A i B będą macierzami kwadratowymi. Każda z wartości własnych macierzy $A \otimes B$ jest postaci

$$\lambda(A \otimes B) = \lambda_i(A) \lambda_j(B). \quad (15.2.12)$$

Wektor własny odpowiadający tej wartości własnej można przedstawić następująco:

$$u(A \otimes B) = u_i(A) \otimes u_j(B), \quad i = 1, 2, \dots, r_A, j = 1, 2, \dots, r_B. \quad (15.2.13)$$

Innych wartości i wektorów własnych, macierz $A \otimes B$ nie ma.

Z (15.2.12) otrzymujemy wzór dla wyznacznika iloczynu Kroneckera macierzy

$$\det(A \otimes B) = [\det(A)]^{r_B} \cdot [\det(B)]^{r_A}. \quad (15.2.14)$$

Literatura

- [1] Adler J.R., *The Geometry of Random Fields*. Chichester, Wiley, 1981.
- [2] Atkinson A.C., *Plots, Transformations and Regression. An Introduction to Graphical Methods of Diagnostic Regression Analysis*. Clarendon Press, Oxford 1985, 1985.
- [3] Atkinson A.C., Developments in the design of experiments. *Int. Statist. Rev.*, 50, s. 161–177, 1982.
- [4] Atkinson A.C., Recent developments in the methods of optimum and related experimental designs. *Int. Statist. Rev.*, 56, s. 99–115, 1988.
- [5] Atkinson A.C., Donev A.N., *Optimum Experimental Designs*. Clarendon Press, Oxford, 1992.
- [6] Atwood C.L., Convergent design sequences for sufficiently regular optimality criteria. *Ann. Stat. Math.*, 4(6), s. 1124–1138, 1976.
- [7] Bandemer H. et al., *Theorie und Anwendung der optimalen Versuchsplanung I. Handbuch zur Theorie*. Akademie-Verlag, Berlin, 1977.
- [8] Barnsley M., *Fractals Everywhere*. New York: Academic Press, 1988.
- [9] Bellman, R., *Introduction to Matrix Analysis*. McGraw-Hill, 1960.
- [10] Benassi A., Cohen S., Istas J., Identification and properties of Real Harmonizable Fractional Levy Motions, *Bernoulli*, Vol. 8, s. 97–115, 2002.
- [11] Benassi A., Cohen S., Istas J., Local self-similarity and Hausdorff dimension, *C.R. Acad. Sci. Paris*, T. 336, s. 267–272, 2003.
- [12] Billingsley P., *Prawdopodobieństwo i miara*. PWN, Warszawa, 1987.
- [13] Boltze L., Näther W., Some remarks on experimental design for estimating the expectation of a stationary random process. *Math. Operationsforsch. Statist., Ser. Statistics*, 11, s. 475–481, 1980.
- [14] Boltze L., Näther W., On effective observation methods in regression models with correlated errors. *Math. Operationsforsch. Statist., Ser. Statistics*, 13, s. 507–519, 1982.
- [15] Box G.E.P., Meyer R.D., Dispersion effects from factorial designs. *Technometrics*, 28, 19–27, 1986.
- [16] Boyd S., Vandenberghe L., *Convex Optimization*. Cambridge: Cambridge University Press, 2004.
- [17] Caines P.E., *Linear Stochastic Systems*. Wiley, New York, 1988.
- [18] Chan D.S., Wong C.S., A general approach to optimal control of a regression experiment. *J. Multiv. Anal.*, 11, s. 85–101, 1981.
- [19] Chakraborti S., Van Der Laan P., Bakir S.T., Nonparametric control charts: An overview and some results, *Journal of Quality Technology*, Vol. 33, s. 304–315, 2001.
- [20] Champ Ch.W., Rigdon S.E., An analysis of the run sum control chart *Journal of Quality Technology*, Vol. 29, s. 407–417, 1997.

- [21] Chen Y.W., Zeng X.Y., LU H., Edge detection and texture segmentation based on independent component analysis. *Proc. 16th Int. Conf. Pattern Recognition ICPR2002, Quebec City August 2002, Session II.9*, s. 351–354, 2002.
- [22] Concia A., Proenca C. B., A fractal image analysis system for fabric inspection based on a box-counting method. *Computer Networks and ISDN Systems*, 30, s. 1887–1895, 1998.
- [23] Daniel C., *Applications of Statistics to Industrial Experiments*. Wiley, New York, 1976.
- [24] Dehnad K., editor. *Quality Control, Robust Design, and the Taguchi Method*. Wadsworth and Brooks/Cole, Pacific Grove, 1989.
- [25] Dette H., A generalization of d - and d_1 -optimal designs in polynomial regression. *Ann. Statist.*, 18:1784–1804, 1990.
- [26] Devroye L., *Non-Uniform Random Variate Generation*. Springer-Verlag New York, Berlin, Heidelberg, Tokyo, 1986.
- [27] Draper M.R., John, R.C., D-optimality for regression. A review. *Technometrics*, 17, s. 15–23, 1975.
- [28] Ermakov S.M., Zhiglyavsky A.A., *Mathematical Theory of Optimal Experiments*. Nauka, Moscow, 1987 (in Russian).
- [29] Ermakov S.M. et al., *Mathematical Theory of Experimental Design*. Nauka, Moscow, 1983 (in Russian).
- [30] Eubank P.L., Smith R.L., Smith P.W., Uniqueness and eventual uniqueness of optimal designs in some time series. *Ann. Stat.* 9, s. 486–493, 1981.
- [31] Eubank R.L., *Spline Smoothing and Nonparametric Regression*. Marcell Dekker, Inc. New York, Basel, 1988.
- [32] Eubank R.L., Speckman P., Convergence rates for trigonometric and polynomial-trigonometric regression estimators. *Statist. and Probab. Letters*, 11, s. 119–124, 1991.
- [33] Eubank, R.L., Smith P.W., On the computation of optimal designs for certain time series models with application to optimal quantile selection for location and scale parameter estimation. *SIAM J. Sci. Statist. Comput.* 3, 2, s. 238–249, 1982.
- [34] Falconer K., *Fractal Geometry*, New York: Wiley, 1990.
- [35] Fedorov V.V., *Theory of optimal experiments*. Academic Press, New York, 1972.
- [36] Fedorov V.V., Hackl P., *Model-oriented Design of Experiments*. Springer, New York, 1997.
- [37] Fellman J., *On the allocation of linear observations*. Commentationes Physico-Mathematicae, Vol. 44, No 2–3, 1974.
- [38] Fellman J., An empirical study of a class of iterative searches for optimal designs. *J. Statist. Planning and Inference*, 21, s. 85–92, 1989.
- [39] Fisher R., *The Design of Experiments*. Oliver Boyd, London, 1935.
- [40] Golub G.H., Heath M.T., Wahba G., Generalized cross-validation as a method for choosing a good ridge parameter. *Technometrics*, 21, s. 215–223, 1979.
- [41] Gonzales, R.C., Woods, R.E., *Digital Image Processing*, Addison- Wesley Publishing Company, Inc., 1992.
- [42] Goodwin, Payne L.C., *Dynamic System Identification Experiment Design and Data Analysis*. Academic Press, London, 1977.

- [43] Greblicki W., *Asymptotycznie optymalne algorytmy rozpoznawania i identyfikacji w warunkach probabilistycznych*. Zeszyty Naukowe Instytutu Cybernetyki Technicznej, 1974.
- [44] Greblicki W., Pawlak M., Fourier and Hermite series estimates of regression functions. *Ann. Inst. Stat. Math.*, 37, s. 443–454, 1985.
- [45] Greblicki W., Pawlak M., Krzyżak A., Distribution – free pointwise consistency of kernel regression estimate. *Ann. Statist.*, 12, s. 1570–1575, 1984.
- [46] Greblicki W., Nonparametric identification of Wiener systems by orthogonal series. *IEEE Trans. Autom. Control*, Vol. 39, No 10, s. 2077–2086, 1994.
- [47] Greblicki W., Nonlinearity estimation in Hammerstein systems based on ordered observations. *IEEE Trans. Signal Process*, Vol. 44, No 5 s. 1224–1233, 1996.
- [48] Greblicki W., Continuous-time Hammerstein system identification. *IEEE Trans. Autom. Control*, Vol. 45 nr 6 s. 1232–1236, 2000.
- [49] Györfi L., Kohler M., Krzyżak A., Walk H. *A Distribution-Free Theory of Nonparametric Regression*. Springer-Verlag, New York, 2003.
- [50] Hajek J., Kimeldorf G., Regression designs in autoregressive stochastic processes. *Ann. Stat.*, 2, s. 520–527, 1974.
- [51] Hall P., Davies S., Fractal analysis of surface roughness by using spatial data, *J. Roy. Statist. Soc., Ser. B*, Vol. 61, s. 3–37, 1999.
- [52] Han D., Tsung F., A generalized EWMA control chart and its comparison with the optimal EWMA, CUSUM and GLR schemes. *Ann. Statist.*, 32, s. 316–340, 2004.
- [53] Hasiwicz Z., Mzyk G., Combined parametric-nonparametric identification of Hammerstein systems. *IEEE Trans. Autom. Control*, Vol. 49, s. 1370–1375, 2004.
- [54] Hasiwicz Z., Pawlak M., Śliwiński P., Nonparametric identification of nonlinearities in block-oriented systems by orthogonal wavelets with compact support. *IEEE Trans. Circuits Systems*, Vol. 52, s. 427–442, 2005.
- [55] Hoel P.G., Minimax design in two dimensional regression. *Ann. Math. Statist.*, Vol. 36, s. 1097–1106, 1965.
- [56] Hryniewicz O., *Współczesne metody statystyczne w sterowaniu jakością*. IBS PAN, Warszawa 1996.
- [57] Huber P.J., *Robust Statistical Procedures*. CBMS-NSF Regional Conference Series in Appl. Math. SIAM, Philadelphia, 1977.
- [58] Istaş J., Lang G., Quadratic variations and estimation of the local Hölder index of a Gaussian process, *Ann. Inst. Poincaré*, Vol. 33, No 4, s. 407–436, 1997.
- [59] Jain A.K., *Fundamentals of Digital Image Processing*. Prentice Hall, NY, 1996.
- [60] Jennrich R.J., Asymptotic properties of nonlinear least squares estimators. *Ann. Math. Stat.*, 40, s. 633–643, 1969.
- [61] Kang L., Albin S.L., Shea G., An X and EWMA chart for individual observations *Journal of Quality Technology*, 1997, Vol. 29, s. 41–47.
- [62] Mańczak K., *Technika planowania eksperymentu*. WNT, Warszawa, 1976.
- [63] Kacprzyński B., *Planowanie eksperymentów. Podstawy matematyczne*. WNT, Warszawa, 1974.
- [64] Karlin, S., Studden W.J., *Chebyshev Systems, with Applications in Analysis and Statistics*. Interscience, New York, 1966.

- [65] Kiefer J., Optimum experimental designs. *J. R. Statist. Soc., Ser. B*, 21, s. 272–304, 1959.
- [66] Kiefer J., Optimum designs in regression problems, ii. *Ann. Math. Statist.*, 32, s. 298–325, 1961.
- [67] Kiefer J., General equivalence theory for optimum designs approximate theory. *Ann. Statist.*, 2, s. 849–879, 1971.
- [68] Kiefer J., Wolfowitz J., Optimum designs in regression problems. *Ann. Math. Statist.*, 30, s. 271–294, 1959.
- [69] Kiefer J., Wolfowitz J., The equivalence of two extremum problems. *Can. J. Math.*, 12, s. 363–366, 1960.
- [70] Kiefer J., Studden W.J., Optimal designs for large degree polynomial regression. *Ann. Statist.*, 4, s. 1113–1123, 1976.
- [71] Kielbasiński A., Schwetlick H., *Numerische Lineare Algebra*. VEB Deutscher Verlag der Wissenschaften, Berlin, 1988. (książkę wydano także w języku polskim)
- [72] Györfi L., Kohler M., Krzyżak A., Walk H., *A Distribution-Free Theory of Non-parametric Regression*. Springer-Verlag, New York, 2003.
- [73] Korbicz J., Uciński D., Sensors allocation for state and parameter estimation of distributed systems. In Gutkowski W. Bauer J., editors, *Proc. IUTAM Int. Symp. Discrete Structural Optimization*, s. 178–198. Int. Symp. *Discrete Structural Optimization*, Zakopane, Poland, August 31 – September 3, Springer-Verlag, 1993.
- [74] Korbicz J., Kościelny J.M., Kowalczyk Z., Cholewa W., (Eds.) *Diagnostyka procesów. Modele. Metody sztucznej inteligencji. Zastosowania*. WNT, Warszawa, 2002.
- [75] Korbicz J., Kościelny J.M., Kowalczyk Z., Cholewa W., (Eds.) *Fault Diagnosis. Models, Artificial Intelligence, Applications*, Berlin Heidelberg, Springer-Verlag, 2004.
- [76] Kurotschka G., A general approach to optimum design of experiments with qualitative and quantitative factors. In *Proceedings of the Indian Statistical Institute, Golden Jubilee International Conference on Statistics: Applications and New Directions, Calcuta 16-19 December 1981*, s. 353–368, Calcuta, 1981.
- [77] Lankaster P., *Theory of Matrices*. Academic Press, London, 1969.
- [78] Läuter E., Experimental design in a class of models. *Math. Operationsforsch. Statist.*, 5, s. 379–398, 1974.
- [79] Lim Y.B., Studden W.J., Efficient d_s -optimal designs for multivariate polynomial regression on the q -cube. *Ann. Statist.*, 16, s. 1225–1240, 1988.
- [80] Logothetis N., Wynn H.P., *Quality Through Design. Experimental Design, Off-line Quality Control and Taguchi's Contributions*. Clarendon Press, Oxford, 1989.
- [81] Malamas E.N. et al, A survey on industrial vision systems, applications and tools, *Image and Vision Computing* 21, s. 171–188, 2003.
- [82] Marcus M., Minc H., *A Survey of Matrix Theory and Matrix Inequalities*. Allyn and Bacon, Boston, 1964.
- [83] Margavio T.M., Alarm rates for quality control charts, *Stat. and Probab. Lett.*, Vol. 24, s. 219–224, 1995.
- [84] Mehra R.K., Optimal input signals for parameter estimation in dynamic systems - survey and new results. *IEEE Trans. Auto. Contr.*, AC-19, s. 753–768, 1974.

- [85] Mehra R.K., Optimization of measurement schedules and sensor design for linear dynamic systems. *IEEE Trans. Auto. Contr.*, 21, s. 55–64, 1976.
- [86] Miller A., Analysis of Parameter Design Experiments for Signal-Response Systems. *J. Quality Tech.*, Vol. 32, No 2, 139–151, 2002.
- [87] Miller A.J., Nam Ky N., A Fedorov exchange algorithm for D-optimal design. *Appl. Statist.*, 43, s. 669–678, 1994.
- [88] Mitchell T.J., An algorithm for the construction of D-optimal experimental designs. *Technometrics*, 16, s. 203–210, 1974.
- [89] Montgomery D.C., *Introduction to Statistical Quality Control*. Wiley, 4th Edition, 2001.
- [90] Müller H.G., Optimal designs for nonparametric kernel regression, *Statistics and Probability Letters*, 2, s. 285–290, 1984.
- [91] Müller C.H., Optimal breakdown point maximizing designs. *Technical Report A-11-95, Fachbereich Mathematik und Informatik, Freie Universität Berlin-West*, 1995.
- [92] Munford A.G., A Control Chart based on Cumulative Scores, *Applied Statistics*, Vol. 29, s. 252–258, 1980.
- [93] Myszka W. (Ed.). *Komputerowy system obsługi eksperymentu*. WNT, Warszawa, 1991.
- [94] Näther W., Optimal observing of stochastic processes. *Math. Operationsforsch. Statist., Ser. Statistics*, 15, s. 239–247, 1984.
- [95] O'Hagan A., Curve fitting and optimal design for prediction. *J. Roy. Statist. Soc., Ser B.*, 40, s. 1–42, 1978.
- [96] Oktaba W. *Metody statystyki matematycznej i metodyka doświadczalnictwa*. PWN, Warszawa, 1971.
- [97] Ott E., *Chaos in Dynamical Systems*. Cambridge University Press: Cambridge, 1993.
- [98] Phadke M. S., *Quality Engineering Using Robust Design*, Prentice-Hall, Englewood Cliffs, NJ, 1989.
- [99] Pawlak M., Rafajłowicz E. On restoration of band-limited signals. *IEEE Trans. Information Theory*, 40, s. 1490–1503, 1994.
- [100] Pawlak M., Rafajłowicz E., Krzyżak A. Exponential weighting algorithms for reconstruction of band-limited signals. *IEEE Trans. Signal Processing*, Vol. 44, s. 538–545, 1996.
- [101] Pawlak M., Rafajłowicz E., Krzyżak A., Postfiltering Versus Prefiltering for Signal Recovery From Noisy Samples. *IEEE Trans. Information Theory*, Vol. IT-49, s. 3195–3212, 2003.
- [102] Pawlak, M., Rafajłowicz, E., Vertically weighted regression – a tool for nonlinear data analysis and constructing control charts. *Journal of the German Statistical Association*, Vol. 84, s. 367–388, 1999.
- [103] Pawlak, M., Rafajłowicz, E. Non-Linear Local Harmonic Filters for Edge-Preserving Image Denoising. *Proc. 16th Int. Conf. Pattern Recognition ICPR2002, Quebec City August 2002*, Session III.11, s. 895–898, 2002.
- [104] Pawlak M., Rafajłowicz E., Steland A., On detecting jumps in time series: nonparametric setting. *J. Nonparametric Statistics*, Vol. 16, nr 3/4, s. 329–347, 2004.

- [105] Pazman A. *Foundations of Optimum experimental design*. D. Reidel, Dordrecht, 1986.
- [106] Pilz J., *Bayesian estimation and experimental design in linear regression models*. Teubner-Texte zur Mathematik Band 55, Teubner-Verlag, Leipzig, 1983.
- [107] Pilz J., *Bayesian Estimation and Experimental Design in Linear Regression Models*. Wiley, Chichester, 2nd ed. 1991.
- [108] Polański Z. *Planowanie doświadczeń w technice*. PWN, Warszawa, 1984.
- [109] Pronzato L., Removing non-optimal support points in D-optimum design algorithms, *Statistics and Probability Letters*, Vol. 63, s. 223–228, 2003.
- [110] Pukelsheim F., *Optimal Experimental Design*, Wiley, New York, 1993.
- [111] Rafajłowicz E., Optimum experiment design for parameter identification in distributed systems. brief survey and new results. In *9th World Congress of the International Federation of Automatic Control*, s. 134–138, Budapest, Hungary, July, 2-6 1984. Preprints. Vol. X, colloquia 14–1, 11–1, 1984.
- [112] Rafajłowicz E., Information equivalence of sensors-controllers configurations in identification of homogeneous static distributed systems. In *Modelling and Simulation of Distributed Parameter Systems. Proceedings of the IMACS/IFAC International Symposium*, s. 553–557, Hiroshima 1987. Hiroshima Inst. Technol., Japan, October 6th–9th, 1987.
- [113] Rafajłowicz E., *Dobór sterowań optymalnych w identyfikacji systemów liniowych o parametrach rozłożonych*. Prace Naukowe Instytutu Cybernetyki Technicznej Politechniki Wrocławskiej, Monografie 13, 1986.
- [114] Rafajłowicz E., Nonparametric orthogonal series estimators of regression: a class attaining the optimal convergence rate in L_2 . *Statist. and Probab. Letters*, s. 219–224, 1987.
- [115] Rafajłowicz E., Nonparametric least squares estimation of a regression function. *Statistics*, 19, s. 349–358, 1988.
- [116] Rafajłowicz E., Myszka W., Optimum experimental design for a regression on a hypercube – generalization of Hoel’s result. *Ann. Inst. Statist. Math.*, 40, s. 821–827, 1988.
- [117] Rafajłowicz E., Optimal experiment design for identification of linear distributed-parameter systems (frequency domain approach). *IEEE Trans. Automatic Control*, Vol. AC-28, Nr 7, s. 806–808, 1983.
- [118] Rafajłowicz E., Design of experiments for eigenvalue identification in distributed-parameter systems. *Int. J. Control*, 34, s. 1079–1094, 1981.
- [119] Rafajłowicz E., Optimal input signals for parameter estimation in linear distributed-parameter systems. *Int. J. Syst. Sci.*, 13, s. 799–808, 1982.
- [120] Rafajłowicz E., Nonparametric algorithm for input signals identification in static distributed-parameter systems. *IEEE Trans. Automatic Control*, 29, s. 631–633, 1984.
- [121] Rafajłowicz E., Optimization of actuators for distributed parameter systems identification. *Problems of Control and Information Theory*, 13, s. 39–51, 1984.
- [122] Rafajłowicz E., Optimization of measurements for state estimation in parabolic distributed systems. *Kybernetika*, 20, s. 44–48, 1984.

- [123] Rafajłowicz E., Unbounded power input signals in optimum experiment design for parameter estimation in linear systems. *Int. J. Control*, 40, s. 383–391, 1984.
- [124] Rafajłowicz E., Spectral synthesis of optimal input signals for linear distributed-parameter systems identification. *Int. J. Syst. Sci.*, 16, s. 667–675, 1985.
- [125] Rafajłowicz E., Adaptive input sequence design for linear distributed-parameter systems identification. *Large Scale Systems*, 11, s. 43–58, 1986.
- [126] Rafajłowicz E., L-optimum input signals for distributed-parameter systems identification. *Problems of Control and Information Theory*, 15, s. 79–89, 1986.
- [127] Rafajłowicz E., Optimum choice of moving sensor trajectories for distributed-parameter system identification. *Int. J. Contr.*, 43, s. 1441–1451, 1986.
- [128] Rafajłowicz E., Sequential identification algorithm and controller choice for a certain class of distributed systems. *Kybernetika*, 22, s. 471–486, 1986.
- [129] Rafajłowicz E., Minimum cost experimental design with a prescribed information matrix. *Theory of Probability and its Applications*, t. 34 s. 412–416, 1989.
- [130] Rafajłowicz E., Reduction of distributed systems identification complexity using intelligent sensors. *Int. J. Contr.*, 50, s. 1571–1576, 1989.
- [131] Rafajłowicz E. Time-domain optimization of input signals for distributed - parameter systems identification. *Journal of Optimization Theory and Applications*, 60, s. 67–79, 1989.
- [132] Rafajłowicz E., Optimum input signals for parameter estimation in systems described by linear integral equations. *Computational Statistics and Data Analysis*, 9, s. 11–19, 1990.
- [133] Rafajłowicz E., Nonparametric identification with errors in independent variables. *Int. J. Syst. Sci.*, 1994.
- [134] Rafajłowicz E., Myszka W., Computational algorithm for input-signal optimization in distributed-parameter systems identification. *Int. J. Syst. Sci.*, 17, s. 911–924, 1986.
- [135] Rafajłowicz E., Myszka W., L-optimum input signals for distributed systems identification. part 2. Computational aspects. *Problems of Control and Information Theory*, 16, s. 193–209, 1987.
- [136] Rafajłowicz E., Myszka W., When product type experimental design is optimal? brief survey and new results. *Metrika*, 39, s. 321–333, 1992.
- [137] Rafajłowicz E., Myszka W., Computational algorithm for generating optimum experimental design on a hypercube. In *Stochastic Methods in Experimental Sciences. Proceedings of the 1989 Cosmex Meeting.*, s. 394–406, Singapore, 1989. Tech. Univ. Wrocław, Szklarska Poręba, 8–14 September 1989, World Scientific 1990.
- [138] Rafajłowicz E., *Algorytmy planowania eksperymentu. Z implementacjami w środowisku Mathematica.* Akademicka Oficyna Wydawnicza PLJ, Warszawa 1996.
- [139] Rafajłowicz E., Selective random search for optimal experiment designs. MODA 5, Advances in Model-Oriented Data Analysis and Experimental Design. *Proceedings of the 5th International Workshop.* Anthony C. Atkinson, Luc Pronzato, Henry P. Wynn (Eds.), Marseilles, France, June 22–26, 1998. Heidelberg, New York: Physica-Verlag, s. 75–83, 1998.
- [140] Rafajłowicz E., Schwabe R., Experimental design for (semi-)local regression. *Communications in Statistics - Theory and Methods*, 32, s. 1035–1055, 2003.

- [141] Rafajłowicz E., Schwabe R., Equidistributed designs in nonparametric regression. *Statistica Sinica* 13, s. 129–142, 2003.
- [142] Rafajłowicz E., Schwabe R., Halton and Hammersley sequences in multivariate nonparametric regression. *Otto-von-Guericke-Universität Magdeburg, Fakultät für Mathematik, Preprint Nr. 20*, 2004.
- [143] Rafajłowicz E., Testing (Non-)Existence of Input–Output Relationships by Estimating Fractal Dimensions. *IEEE Trans. Signal Processing*, Vol. 52, No. 11, s. 3151–3159, 2004.
- [144] Ranky P.G., *Total Quality Control and Management in CIM*. CIMware Ltd, Guilford, 1989.
- [145] Rao C.R., *Linear Statistical Inference and Its Applications*. Wiley, New York, 1973.
- [146] Rasch D., Herrendorfer A., *Statystyczne planowanie doświadczeń*. PWN, Warszawa, 1991.
- [147] Regalia P.A., Mitra S.K., Kronecker products, unitary matrices and signal processing applications, *SIAM Review*, Vol. 31, No 4, s. 586–613, 1989.
- [148] Reynolds M.R. Jr, Stoumbos Z.G., A CUSUM chart for monitoring a proportion when inspecting continuously *Journal of Quality Technology*, Vol. 31, s.87–108, 1999.
- [149] Robinson J.A., Optimal Detection of Blurred Edges, *Proc. 16th Int. Conf. Pattern Recognition ICPR 2002, Quebec City August 2002*, Session III.11, s. 831–834, 2002.
- [150] Rousseeuw P.J., Leroy A.M. *Robust regression and Outlier detection*. Wiley, New York, 1987.
- [151] Rutkowski L., On system identification by nonparametric function fitting. *IEEE Trans. Automat. Contr.*, AC 27, s. 225 – 227, 1982.
- [152] Rutkowski L., Rafajłowicz E., On optimal global rate of convergence of some nonparametric identification procedures. *IEEE Trans. Automatic Control*, AC-34, s. 1089–1092, 1989.
- [153] Rutkowski L. *Filtry adaptacyjne i adaptacyjne przetwarzanie sygnałów*. WNT, Warszawa, 1994.
- [154] Ryan T.P., Schwertman N. C., Optimal limits for attributes control charts, *Journal of Quality Technology*, Vol. 29, s.86–98, 1997.
- [155] Rybaczuk M., Kasprzak W., Lysik B., *Dimensional Analysis in the Identification of Mathematical Models*. World Scientific, Singapore, New Jersey, 1990.
- [156] Sacks J. et al, Design and analysis of computer experiments. *Statist. Science*, 4, s. 409–435, 1989.
- [157] Sacks J., Ylvisaker D., Design for regression problems with correlated errors. *Ann. Math. Stat.*, 37, s. 66–89, 1966.
- [158] Sacks J., Ylvisaker D., Design for regression problems with correlated errors: many parameters. *Ann. Math. Stat.*, 39, s. 49–69, 1968.
- [159] Sacks J., Ylvisaker D., Designs for regression problems with correlated errors, iii. *Ann. Math. Statist.*, 41, s. 2057–2074, 1970.
- [160] Santner T.J., Williams B.J., Notz W.I., *The Design and Analysis of Computer Experiments*. Springer-Verlag, New York, 2003.

- [161] Schoenfelder C., Cambanis S., Random designs for estimating integrals of stochastic processes. *Ann. Stat.*, 10, s. 526–538, 1982.
- [162] Schuster H.G., *Deterministic Chaos*. Weinheim: VGH Verlagsgesellschaft, 1988.
- [163] Schwabe R., *Optimum Designs for Multi-Factor Models*. Freie Universität Berlin West, 1992.
- [164] Schwabe R., Model robust experimental design in the presence of interactions. the orthogonal case. In *PROBSTAT '94 Proc. Int. Conf. Math. Ststist., Smolenice 1994, Tatra Mountains Math. Publ. 7*, 1994.
- [165] Schwabe R., Experimental design for linear models with higher interaction terms. In Mammschitz and Schneeweiss, editors, *Symposia Gaussiana*. Gruyter & Co, Berlin, New York, 1995.
- [166] Schwabe R., Uncertain resources and designing for additional information. Technical Report A-17-95, Fachbereich Mathematik und Informatik, Freie Universität Berlin-West, 1995.
- [167] Schwabe R., Wong W.K., Some relationships between efficiencies and marginal efficiencies of product designs. *Technical Report A-18-95, Fachbereich Mathematik und Informatik, Freie Universität Berlin-West*, 1995.
- [168] Schwabe R., Wong W.K., A relationship between efficiencies of marginal designs and the product design. *Metrika* 45, s. 253–257, 1997.
- [169] Schwabe R., Optimal designs for hierarchical interaction structures. *Journal of Statistical Planning and Inference* 70, s. 181–190, 1998.
- [170] Schwabe R., Wong W.K., Efficient product designs for quadratic models on the hypercube. *Sankhya* 65, s. 649–659, 2003.
- [171] Schwabe R., On an adaptive design in regression. *Statistics*, 18, s. 521–525, 1987.
- [172] Seber G.A., *Linear regression Analysis*. Wiley, New York, 1977.
- [173] Seber G.A., Wild C. J., *Nonlinear Regression.*, Wiley, New York, 1989.
- [174] Silvey S.D., Titterington D.M., Torsney, B., An algorithm for optimal designs on a infinite design space. *Comm. Stat. Th. Meth. A*, 7, No 14, s. 1379–1389, 1978.
- [175] Silvey S.D., *Optimal Design*. Chapman and Hall, London, 1980.
- [176] Skubalska-Rafajłowicz E., Rafajłowicz E., Searching for optimal experimental designs using space-filling curves. *Appl. Math. Comput. Sci.*, Vol. 8, s. 647–656, 1998.
- [177] Skubalska-Rafajłowicz E., *Krzywe wypełniające w problemach decyzyjnych*, Wyd. Politechniki Wrocławskiej, Seria Monografie, 2001.
- [178] Skubalska-Rafajłowicz E., Neural networks with orthogonal activation function approximating space-filling curves. *Proceedings of the 9th IEEE International Conference on Methods and Models in Automation and Robotics. MMAR 2003. Ed. R. Kaszyński*, s. 927–934, 2003.
- [179] Skubalska-Rafajłowicz E., Recurrent network structure for computing quasi- inverses of the Sierpiński space-filling curves. *Artificial intelligence and soft computing – ICAISC 2004. 7th International conference. Proceedings*, Leszek Rutkowski (ed.). Zakopane, June 7–11, 2004. Springer, 2004, s. 272–277, (Lecture Notes in Computer Science, Lecture Notes in Artificial Intelligence Vol. 3070), 2004.
- [180] Skubalska-Rafajłowicz E., A new method of estimation the box-counting dimension of the multivariate objects using space-filling curves. *Nonlinear Analysis* (w druku).

- [181] Spruill M.C., Optimal designs for second order processes with general linear means. *Ann. Stat.*, 8, s. 652–663, 1980.
- [182] Spruill M.C., Studden W.J., Optimum designs when the observations are second order processes. *J. Multivar. Anal.*, 8, s. 153–172, 1978.
- [183] Spruill M.C., Studden W.J., A Kiefer–Wolfowitz theorem in a stochastic process setting. *Ann. Stat.*, 7, s. 1329–1332, 1979.
- [184] Stone C.J., Optimal global rate of convergence for nonparametric regression. *Ann. Statist.*, 10, s. 1040–1053, 1982.
- [185] Taguchi G., *System of Experimental Design.*, Volume 1 and 2. UNIPUB/Kraus International, White Plains, 1987.
- [186] Taguchi G. Performance analysis design. *Int. J. Proc. Res.*, 16, s. 521–530, 1978.
- [187] Taguchi G. *Introduction to Quality Engineering. Designing Quality into Products and Processes.* Asian Productivity Organization, Tokyo, 1986.
- [188] Thompson J. R., Koronacki J., *Statystyczne sterowanie procesem. Metoda Deminga etapowej optymalizacji jakości.* Akademicka Oficyna Wydawnicza PLJ, Warszawa, 1994.
- [189] Titterton D.M., Aspects of optimal designs in dynamic systems. *Technometrics*, 22, s. 287–297, 1980.
- [190] Torsney B., Computing optimizing distributions with applications in designs, estimation and image processing. In Fedorov V.V. Dodge, Y., H. Wynn, Ed., *Optimal Design and Analysis of Experiments*, s. 361–370. North-Holland, Amsterdam, New York, 1988.
- [191] Torsney B., Mandal S., Multiplicative Algorithms for Constructing Optimizing Distributions: Further Developments. In Di Bucchianico, Lauter H., H. Wynn, Ed., *MODA 7 – Advances in Model-Oriented Design and Analysis*, s. 361–370. Physica-Verlag, 2004.
- [192] Tricot C., *Curves and Fractal Dimension.* New York: Springer, 1995.
- [193] Uciński D., Optimal sensors location for parameter identification of distributed systems. *Appl. Math. and Comp. Sci.*, 2, s. 119–134, 1992.
- [194] Uciński D., Optimal design of moving sensors trajectories for identification of distributed parameter systems. *Proc. 1st Int. Symp. Mathematical Models in Automation and Robotics*, Międzyzdroje, Poland, September 1–3, 1994.
- [195] Uciński D., Optimal selection of measurement locations for identification of parameters in distributed systems. *Proc. 2nd Int. Symp. Methods and Models in Automation and Robotics*, Międzyzdroje, Poland, 30 August – 2 September, 1995.
- [196] Uciński D., Korbicz J., An algorithm for the computation of optimal measurement trajectories in distributed parameter systems identification. *Proc. 3rd European Control Conference*, Rome, Italy, September 5–8, 1995.
- [197] Uciński D., Korbicz J., Moving sensors in state and parameter estimation of dps – a survey of some techniques. – *Proc. IFIP Conf. Modelling and Optimization of Distributed Parameter Systems with Applications to Engineering*, Warsaw, 1995.
- [198] Uciński D., Korbicz J., Zaremba M., On optimization of sensors motions in parameter identification of two-dimensional distributed systems. In *Proc. 2nd European Control Conference*, Groningen, The Netherlands, June/July, s. 1359–1364. European Control Conference, 1993.

- [199] Uciński D., *Optimal Measurement Methods for Distributed Parameter System Identification*. CRC Press, London, New York, 2005.
- [200] Uciński D., Bogacka B., T-optimum designs for discrimination between two multivariate dynamic models, *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, Vol. 67, s. 3–18, 2005.
- [201] Wallis W.D., *Combinatorial Designs*. Marcel Dekker, New York, Basel, 1988.
- [202] Walter E., Pronzato L., *Identification of Parametric Models from Experimental Data*. Springer, London, 1997.
- [203] Welch W.J. et al, Screening, predicting, and computer experiments. *Technometrics*, 34, s. 15–25, 1992.
- [204] Wierich W., On optimal designs and complete class theorems for experiments with continuous and discrete factors of influence. *J. Statist. Plann. Inference*, 15, s. 19–27, 1986.
- [205] Wierich W., A-optimal design measures for one-way layouts with additive regression. *J. Statist. Plann. Inference*, 18, s. 57–68, 1988.
- [206] Winterbottom A., Basso L., Wynn H.P., A review of the Taguchi methods for off-line quality control. *Quality Reliability Eng.*, 2, s. 71–79, 1986.
- [207] Woodall W. H., Control charts based on attribute data: Bibliography and review *Journal of Quality Technology*, Vol. 29, s. 172–183, 1997.
- [208] Wu Z., Spedding T.A., A synthetic control chart for detecting small shifts in the process mean *Journal of Quality Technology*, Vol. 32, s. 32–38, 2000.
- [209] Wu C.F.J., Hamada M., *Experiments: Planning, Analysis and Parameter Design Optimization*. John Wiley and Sons, Inc., New York, 2000.
- [210] Wynn H.P., The sequential generation of d-optimum experimental designs. *Ann. Math. Stat.*, 1, s. 1655–1664, 1970.
- [211] Wynn H.P., Chien Fu W., The convergence of general step length algorithms for regular optimum design criteria. *Ann. Stat.*, 6(6), s. 1273–1285, 1976.
- [212] Yang Z., On the Performance of Geometric Charts with Estimated Control Limits, *Journal of Quality Technology*, Vol. 34, No. 4, s. 448–458, 2002.
- [213] Zarrop M.B., *Optimal experimental design for dynamic system identification*, Lecture Notes in Control and Information Science 21. Springer-Verlag, Berlin, Heidelberg, New York, 1979.
- [214] Zieliński R., *Wybrane zagadnienia optymalizacji statystycznej*. PWN, Warszawa, 1974.

Skorowidz

- Φ -optymalność
 - definicja, 36
 - warunek optymalności, 42
- A-optymalność
 - przykład, 44, 89
 - warunek optymalności, 43
- algorytm Wyna–Fedorova
 - idea, 52
 - opis, 52
 - zbieżność, 53
- algorytm Wyna–Fedorova
 - modyfikacje, 55
- ciąg reszt, 17
- D-optymalność
 - metody numeryczne
 - I rzędu, 53
 - II rzędu, 56
 - modyfikacje, 55
 - optymalizacja wag, 53
 - szybkość zbieżności, 53
 - wielowymiarowe, 53
 - nierówność, 41
 - plan równomierny, 46
 - regresja liniowa, 44
 - obszar – kostka, 46
 - regresja trygonometryczna, 45
- estymacja
 - dokładność, 16–18
 - macierz kowariancji, 17
 - obciążenie, 16
 - ocena wariancji, 18
 - oceny regresji, 17
 - nieobciążona, 16
 - systemy rozłożone
 - dobór wymuszeń, 148
- estymator
 - grzebieniowy, 19
 - Jamesa–Steina, 19
 - liniowy, 18
 - o minimalnej wariancji, 18
 - wariancji regresji, 18
- G-optymalność
 - definicja, 38
 - twierdzenie o równoważności, 40
 - własności kryterium, 38
- iloczyn Kroneckera, 127, 197
 - definicja, 197
 - macierzy, 197
 - własności, 198
 - wektorów, 197
 - wielokrotny, 198
- karta kontrolna
 - badanie liczby defektów, 176
 - dobór parametrów, 180
 - dostrajanie, 181
 - model zmian jakości, 174
 - nieparametryczna, 173
 - pożądane cechy, 175
 - porównanie z EWMA i CUSUM, 186
 - symulacje, 183
 - w sekwencjach obrazów, 189
 - zakłócenia niegaussowskie, 187
 - zmodyfikowana, 177
- kombinacja wypukła planów, 31
- kontury obrazu, 162
- kosztowa optymalność
 - klasyczne plany, 123
 - optymalne plany, 122
- regresja liniowa
 - obszar – kostka, 127
 - obszar – kula, 125
 - twierdzenie Karlina i Issi, 121
- kryteria planowania, 35, 43
 - częściowe, 38
 - D-optymalność, 36, 43

- E- optymalność, 37
- G- optymalność, 38, 40
- L_p - optymalność, 37
- L- optymalność, 36, 37
- Q- optymalność, 38
- w przestrzeni funkcji, 38
- kryterium D- optymalności
 - pochođna w kierunku, 39
 - twierdzenie Kiefera i Wolfowitza, 40
 - własności, 38, 39
 - warunki optymalności, 38, 40
- krzywe wypelniające, 71
- macierz Hadamarda, 126, 127
- macierz informacyjna
 - definicja, 30
 - nieujemna określoność, 31
 - osobliwa, 31
 - planu ciągłego, 30
 - symetria, 31
 - unormowanego planu, 30
 - własności, 31
 - zbiór realizowalnych macierzy, 31
- maksimum globalne, 57
- metoda
 - najmniejszych kwadratów
 - klasyczna, 14
 - z wagami, 15
 - największej wiarygodności, 15
- metoda najmniejszych kwadratów, 14
 - błąd średniokwadratowy, 14
 - estymator, 14
 - istnienie rozwiązania, 16
 - jednoznaczność rozwiązania, 16
 - oszacowanie funkcji regresji, 16
 - oszacowanie parametrów, 16
 - równania normalne, 15
- MNK, 14
 - błąd średniokwadratowy, 14
 - estymator, 14
 - istnienie rozwiązania, 16
 - jednoznaczność rozwiązania, 16
 - równania normalne, 15
 - rozkład normalny, 15
 - z wagami, 15
- model
 - addytywny
 - blokowo, 74
 - plany optymalne, 78
 - w pełni, 75
 - multiplikatywny, 98
 - częściowo, 79, 81
 - w pełni, 80, 93
 - z pełnym zestawem interakcji, 113
- model liniowy
 - funkcje rozpinające, 13
 - liniowa niezależność, 11
 - przykłady, 13
 - opis, 11
 - rozdzielność parametrów, 13
 - założenie o x -ach, 12
 - założenie o x -ach, 12
 - zakłócenia, 11
 - addytywność, 11
 - wariancja, 12
- nośnik planu
 - D- optymalnego, 43
- odporność procesów
 - model, 107
 - planowanie, 107
- plan
 - ciągły
 - korzystanie, 28
 - macierz informacyjna, 30
 - kombinacja wypukła, 31
 - symetryczny, 77
 - unormowany, 26
 - macierz informacyjna, 30
 - plan czynnikowy, 46, 85
 - plan D- optymalny
 - heterogeniczna wariancja, 45
 - liczba punktów, 43
 - liczba punktów nośnika, 43
 - regresja liniowa, 44
 - regresja wielomianowa, 44
 - własności, 38, 39
- planowanie eksperymentu

- badania symulacyjne, 145
- badanie wyrobów, 151, 153, 155
- D- optymalne, 36, 43
- do testowania odporności, 106
- kryteria, 35, 43
- model charakterystyki, 152
- nadążające za zmianami, 128
 - model, 129
 - sformułowanie zadania, 128
- o symetrii obrotowej, 125
- ocena odporności, 107
- ortogonalne, 124
 - kanoniczność, 124
 - na kostce, 126
- problemy
 - czynniki uboczne, 144
 - dobór składu mieszanin, 145
 - testowanie hipotez, 144
 - w symulacji, 145
- sekwencje planów, 128
- wymagania
 - czynniki ilościowe, 21
 - czynniki jakościowe, 21
 - informacja aprioryczna, 21
 - odporność, 145
 - optymalność, 21
 - ortogonalność, 20
 - sekwencyjność, 21
 - symetria obrotowa, 20
 - walory numeryczne, 20
- założenia
 - informacja aprioryczna, 21
 - korelacja zakłóceń, 21
- plany produktowe
 - macierz informacyjna, 76
 - optymalność
 - kryteria planowania, 84
 - w modelach multiplikatywnych, 84
 - warunki, 77, 85
 - symetria, 76
- plany wielowymiarowe
 - algorytm losowy, 57
 - D- optymalność, 53
 - komponowanie, 93
 - metody numeryczne, 53
 - optymalizacja wag, 53
 - porównanie algorytmów, 69
 - produktowe, 93
 - przekleństwo wymiarowości, 56
 - przykłady numeryczne, 67
 - wieloekstremalność, 56
- regresja
 - addytywna, 74
 - funkcje rozpinające, 13
 - liniowa, 9, 46
 - model, 9
 - modele częściowe, 75
 - plany produktowe, 82
 - przykłady, 81
 - przykłady interakcji, 86
 - przykłady planów, 85
 - testowanie hipotez, 18
 - trygonometryczna, 45
 - wielomianowa, 44
 - wielowymiarowa, 81
 - z interakcjami, 78
- rozkład normalny, 15
- rozmieszczenie czujników
 - dobór trajektorii, 148
 - maksymalizacja dokładności, 148
- suma prosta wektorów, 74
- twierdzenie
 - o równoważności, 40
 - Gaussa-Markova, 18
 - Karlina-Issi, 121
 - Kiefera i Wolfowitza, 40
 - Kiefera-Wolfowitza uogólnione, 42
 - o planach czynnikowych, 46
 - o planach produktowych, 77
 - o zwartości, 33
- wielomiany Legendre'a, 45
- założenia
 - o x -ach, 12
 - o modelu liniowym, 11
 - o zakłóceniach, 11
 - w modelu Taguchi, 111