



MATEMATYKA W UCZENIU MASZYNOWYM

Marc Peter Deisenroth
A. Aldo Faisal
Cheng Soon Ong

Tytuł oryginału: Mathematics for Machine Learning

Tłumaczenie: Filip Kamiński

ISBN: 978-83-283-8459-0

© Marc Peter Deisenroth, A. Aldo Faisal, and Cheng Soon Ong 2020

This translation of *Mathematics for Machine Learning* is published by arrangement with Cambridge University Press.

Polish edition copyright © 2022 by Helion S.A.

All rights reserved.

All rights reserved. No part of this book may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording or by any information storage retrieval system, without permission from the Publisher.

Wszelkie prawa zastrzeżone. Nieautoryzowane rozpowszechnianie całości lub fragmentu niniejszej publikacji w jakiegokolwiek postaci jest zabronione. Wykonywanie kopii metodą kserograficzną, fotograficzną, a także kopiowanie książki na nośniku filmowym, magnetycznym lub innym powoduje naruszenie praw autorskich niniejszej publikacji.

Wszystkie znaki występujące w tekście są zastrzeżonymi znakami firmowymi bądź towarowymi ich właścicieli.

Autor oraz wydawca dołożyli wszelkich starań, by zawarte w tej książce informacje były kompletne i rzetelne. Nie biorą jednak żadnej odpowiedzialności ani za ich wykorzystanie, ani za związane z tym ewentualne naruszenie praw patentowych lub autorskich. Autor oraz wydawca nie ponoszą również żadnej odpowiedzialności za ewentualne szkody wynikłe z wykorzystania informacji zawartych w książce.

Drogi Czytelniku!

Jeżeli chcesz ocenić tę książkę, zajrzyj pod adres

<https://helion.pl/user/opinie/mawuma>

Możesz tam wpisać swoje uwagi, spostrzeżenia, recenzję.

Helion S.A.

ul. Kościuszki 1c, 44-100 Gliwice

tel. 32 231 22 19, 32 230 98 63

e-mail: helion@helion.pl

WWW: <https://helion.pl> (księgarnia internetowa, katalog książek)

Printed in Poland.

- Kup książkę
- Poleć książkę
- Oceń książkę

- Księgarnia internetowa
- Lubię to! » Nasza społeczność

Spis treści

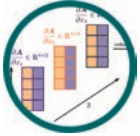
Lista symboli	9
Wstęp	13
Podziękowania	16
Część I. Podstawy matematyczne	19
1 Wprowadzenie i motywacje	21
1.1. Znajdowanie słów dla intuicji	21
1.2. Dwa sposoby na przeczytanie tej książki	23
1.3. Ćwiczenia i informacje zwrotne	25
2 Algebra liniowa	26
2.1. Układy równań liniowych	28
2.2. Macierze	31
2.3. Rozwiązywanie układów równań liniowych	36
2.4. Przestrzenie wektorowe	44
2.5. Niezależność liniowa	49
2.6. Baza i rząd	53
2.7. Przekształcenia liniowe	56
2.8. Przestrzenie afiniczne	69
2.9. Materiały dodatkowe	71
Ćwiczenia	71
3 Geometria analityczna	78
3.1. Normy	79
3.2. Iloczyny wewnętrzne	80
3.3. Długości i odległości	83

3.4. Kąty i ortogonalność	84
3.5. Baza ortonormalna	87
3.6. Dopełnienie ortogonalne	88
3.7. Iloczyn wewnętrzny funkcji	89
3.8. Rzuty ortogonalne	90
3.9. Obroty	99
3.10. Materiały dodatkowe	103
Ćwiczenia	103
4 Rozkłady macierzy	106
4.1. Wyznacznik i ślad	107
4.2. Wartości i wektory własne	113
4.3. Rozkład Choleskiego	122
4.4. Rozkład według wartości własnych i diagonalizacja	124
4.5. Rozkład według wartości osobliwych	127
4.6. Przybliżenie macierzy	139
4.7. Filogeneza macierzy	144
4.8. Materiały dodatkowe	145
Ćwiczenia	146
5 Rachunek wektorowy	149
5.1. Różniczkowanie funkcji jednowymiarowych	151
5.2. Pochodne cząstkowe i gradienty	156
5.3. Gradienty funkcji o wartościach wektorowych	160
5.4. Gradienty macierzy	165
5.5. Tożsamości przydatne w obliczeniach gradientów	169
5.6. Propagacja wsteczna i różniczkowanie automatyczne	170
5.7. Pochodne wyższych rzędów	175
5.8. Linearyzacja i wielowymiarowe szeregi Taylora	176
5.9. Materiały dodatkowe	181
Ćwiczenia	182
6 Prawdopodobieństwo i jego rozkłady	184
6.1. Struktura przestrzeni prawdopodobieństwa	184
6.2. Prawdopodobieństwo ciągłe i dyskretne	190
6.3. Reguły dodawania i mnożenia oraz twierdzenie Bayesa	195
6.4. Statystyki podsumowujące i niezależność	198
6.5. Rozkład Gaussa	208
6.6. Sprzężenie i rodzina wykładnicza	216
6.7. Zmiana zmiennych/przekształcenie odwrotne	225
6.8. Materiały dodatkowe	231
Ćwiczenia	232

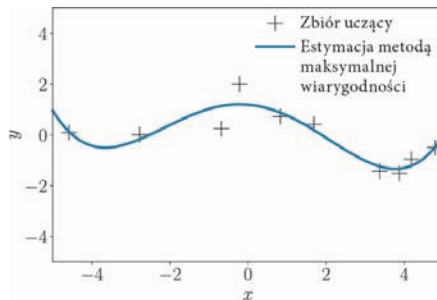
7	Optymalizacja ciągła	235
	7.1. Optymalizacja za pomocą metody gradientu prostego	238
	7.2. Optymalizacja z ograniczeniami i mnożniki Lagrange'a	243
	7.3. Optymalizacja wypukła	246
	7.4. Materiały dodatkowe	256
	Ćwiczenia	257
	Część II. Centralne problemy uczenia maszynowego	259
8	Gdy model spotyka dane	261
	8.1. Dane, modele i uczenie	261
	8.2. Minimalizacja ryzyka empirycznego	268
	8.3. Estymacja parametrów	274
	8.4. Modelowanie probabilistyczne i wnioskowanie	282
	8.5. Modele digrafowe	286
	8.6. Wybór modelu	291
9	Regresja liniowa	298
	9.1. Sformułowanie problemu	300
	9.2. Estymacja parametrów	301
	9.3. Bayesowska regresja liniowa	313
	9.4. Estymacja metodą maksymalnej wiarygodności jako rzut ortogonalny	323
	9.5. Materiały dodatkowe	325
10	Redukcja wymiarowości za pomocą analizy głównych składowych	327
	10.1. Sformułowanie problemu	328
	10.2. Perspektywa maksymalizacji wariancji	330
	10.3. Perspektywa rzutowania	335
	10.4. Znajdowanie wektora własnego i aproksymacja za pomocą macierzy niskiego rzędu	343
	10.5. PCA w dużej liczbie wymiarów	345
	10.6. Najważniejsze kroki algorytmu PCA z praktycznego punktu widzenia	346
	10.7. Perspektywa zmiennej ukrytej	350
	10.8. Materiały dodatkowe	353
11	Szacowanie gęstości za pomocą modeli mieszanin rozkładów Gaussa	358
	11.1. Model mieszaniny rozkładów Gaussa	359
	11.2. Uczenie parametrów za pomocą metody maksymalnej wiarygodności	360
	11.3. Algorytm EM	371
	11.4. Perspektywa zmiennej ukrytej	373
	11.5. Materiały dodatkowe	379

12	Klasyfikacja za pomocą maszyny wektorów nośnych	381
12.1.	Hiperpłaszczyzny rozdzielające	383
12.2.	Pierwotna maszyna wektorów nośnych	385
12.3.	Dualna maszyna wektorów nośnych	393
12.4.	Jądra	398
12.5.	Rozwiązanie numeryczne	401
12.6.	Materiały dodatkowe	402
	Bibliografia	405

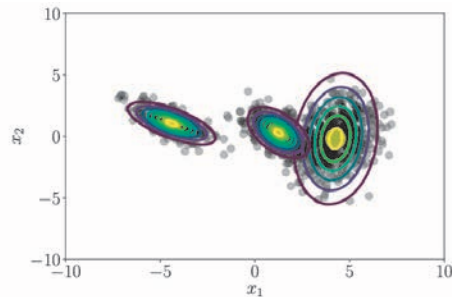
Rachunek wektorowy



W wielu algorytmach uczenia maszynowego wykorzystuje się optymalizację funkcji celu względem parametrów modelu. Parametry te kontrolują to, jak dobrze model wyjaśnia dane. Problem doboru odpowiednich parametrów można sformułować w postaci problemu optymalizacyjnego (patrz podrozdziały 8.2 i 8.3). Przykładami takich problemów są: (i) regresja liniowa (rozdział 9.), w której przyglądamy się problemom dopasowania krzywych i optymalizujemy parametry liniowo zależnych wag w celu maksymalizacji wiarygodności; (ii) autoenkodery w sieciach neuronowych wykorzystywane do redukcji liczby wymiarów i kompresji danych; ich parametrami są wagi i biasy każdej warstwy; w przypadku autoenkoderów minimalizujemy błąd rekonstrukcji poprzez wielokrotne wykorzystywanie reguły łańcuchowej; oraz (iii) modele mieszanin rozkładów Gaussa (patrz rozdział 11.) wykorzystywane do modelowania rozkładów danych, w których, aby zmaksymalizować wiarygodność, optymalizujemy parametry opisujące położenie i kształt każdego składnika mieszaniny. Na rysunku 5.1 pokazano wybrane problemy, które rozwiązuje się zazwyczaj za pomocą algorytmów optymalizacji wykorzystujących informacje o gradientach (podrozdział 7.1). Na rysunku 5.2 pokazano, jak pojęcia z tego rozdziału są powiązane ze sobą i z innymi rozdziałami tej książki.

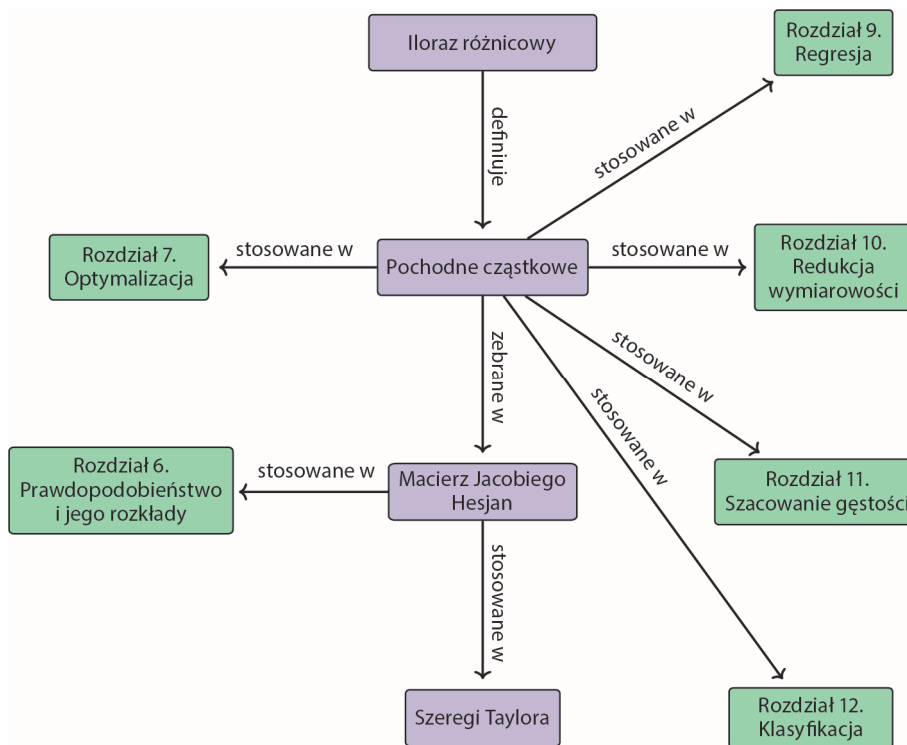


(a) Problem regresji: Dobierz parametry modelu tak, aby krzywa dobrze wyjaśniała obserwacje (zaznaczone krzyżykami)



(b) Szacowanie gęstości za pomocą modeli mieszanin rozkładów Gaussa: Znajdź średnie i kowariancje, które dobrze wyjaśniają dane (zaznaczone kropkami)

RYSUNEK 5.1. Rachunek wektorowy odgrywa kluczową rolę w (a) regresji (dopasowywanie krzywej) i (b) szacowaniu gęstości, czyli modelowaniu rozkładów danych



RYSUNEK 5.2. Mapa myśli prezentująca pojęcia przedstawione w tym rozdziale wraz z miejscami ich użycia w innych częściach książki

Centralnym punktem tego rozdziału jest pojęcie funkcji. Funkcja f to relacja, która wiąże ze sobą dwie wielkości. W tej książce wielkościami tymi są zazwyczaj dane wejściowe $\mathbf{x} \in \mathbb{R}^D$ i przewidywania (wartości funkcji) $f(\mathbf{x})$, o których zakładamy, że są wartościami rzeczywistymi (jeśli nie zaznaczono inaczej). Zbiór \mathbb{R}^D nazywamy **dziedzina** funkcji f . Wartości $f(\mathbf{x})$ nazywamy **zbiorem wartości**, **obrazem** lub **przeciwdziedzina** funkcji f . Funkcje liniowe zostały omówione bardziej szczegółowo w punkcie 2.7.3. Do zapisu funkcji często wykorzystujemy następującą notację:

$$f: \mathbb{R}^D \rightarrow \mathbb{R}, \quad (5.1a)$$

$$\mathbf{x} \mapsto f(\mathbf{x}). \quad (5.1b)$$

Równanie 5.1a mówi nam, że f jest odwzorowaniem prowadzącym z \mathbb{R}^D do \mathbb{R} , a równanie 5.1b jawnie określa, że wejściu \mathbf{x} przypisujemy dokładnie jedną wartość funkcji $f(\mathbf{x})$. Funkcja f przypisuje każdemu wejściu \mathbf{x} dokładnie jedną wartość funkcji $f(\mathbf{x})$.

Przykład 5.1

Przypomnijmy, że iloczyn skalarny jest szczególnym przypadkiem iloczynu wewnętrznego (podrozdział 3.2). W powyższej notacji funkcję $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{x}$, $\mathbf{x} \in \mathbb{R}^2$ zapisalibyśmy jako

$$f: \mathbb{R}^2 \rightarrow \mathbb{R} \quad (5.2a)$$

$$\mathbf{x} \mapsto x_1^2 + x_2^2. \quad (5.2b)$$

W tym rozdziale omówimy sposoby obliczania gradientów funkcji. Często gradienty są narzędziem niezbędnym do uproszczenia procesu uczenia modeli uczenia maszynowego, ponieważ wskazują kierunek najszybszego wzrostu funkcji. Z tego powodu rachunek wektorowy jest jednym z podstawowych narzędzi matematycznych wykorzystywanych w uczeniu maszynowym. W tej książce zakładamy, że badane przez nas funkcje są różniczkowalne. Po zastosowaniu pewnych nieomówionych w tej książce definicji wiele zaprezentowanych tu podejść można rozszerzyć również na funkcje subróżniczkowalne (funkcje, które są ciągłe, ale w pewnych punktach nie są różniczkowe). W rozdziale 7. znajdziesz przykład takiego rozszerzenia na funkcje z ograniczeniami.

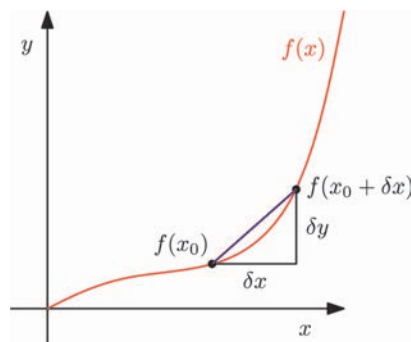
5.1. Różniczkowanie funkcji jednowymiarowych

Poniżej pokrótce omówimy różniczkowanie funkcji jednej zmiennej, które możesz znać z lekcji matematyki w szkole średniej. Zaczniemy od ilorazu różnicowego jednowymiarowej funkcji $y = f(x)$, $x, y \in \mathbb{R}$. Następnie na jego podstawie zdefiniujemy pochodną.

Definicja 5.1 (iloraz różnicowy). **Iloraz różnicowy**

$$\frac{\delta y}{\delta x} := \frac{f(x + \delta x) - f(x)}{\delta x} \quad (5.3)$$

opisuje nachylenie siecznej przechodzącej przez dwa punkty na wykresie funkcji f . Na rysunku 5.3 są nimi punkty, których współrzędne na osi x są równe x_0 i $x_0 + \delta x$.



RYСУNEK 5.3. Średnie nachylenie funkcji f pomiędzy x_0 i $x_0 + \delta x$ jest równe nachyleniu siecznej (kolor niebieski) przechodzącej przez $f(x_0)$ i $f(x_0 + \delta x)$. To nachylenie to $\delta y / \delta x$

Jeżeli przyjmiemy, że f jest funkcją liniową, to za iloraz różnicowy można również uznać średnie nachylenie f pomiędzy x_0 i $x_0 + \delta x$. Jeśli f jest różniczkowalna, to w granicy $\delta x \rightarrow 0$ otrzymamy styczną do f w punkcie x . W takim przypadku styczna ta jest pochodną f w punkcie x .

Definicja 5.2 (pochodna). Bardziej formalnie: dla $h > 0$ **pochodna** f względem x jest zdefiniowana jako granica

$$\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}. \quad (5.4)$$

Sieczna z rysunku 5.3 staje się w tej definicji styczną.

Pochodna wskazuje kierunek najszybszego wzrostu wartości f .

Przykład 5.2 (pochodna wielomianu)

Chcemy obliczyć pochodną funkcji $f(x) = x^n$, gdzie $n \in \mathbb{N}$. Być może wiesz już, że będzie ona równa nx^{n-1} , ale w tym przykładzie chcemy wyprowadzić tę wartość z definicji pochodnej jako granicy ilorazu różnicowego.

Z definicji pochodnej z równania 5.4 otrzymujemy

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad (5.5a)$$

$$= \lim_{h \rightarrow 0} \frac{(x+h)^n - x^n}{h} \quad (5.5b)$$

$$= \lim_{h \rightarrow 0} \frac{\sum_{i=0}^n \binom{n}{i} x^{n-i} h^i - x^n}{h}. \quad (5.5c)$$

Zauważ, że $x^n = \binom{n}{0} x^{n-0} h^0$. Z tego powodu x^n upraszcza się, a w liczniku pozostaje jedynie suma indeksowana od jeden:

$$\frac{df}{dx} = \lim_{h \rightarrow 0} \frac{\sum_{i=1}^n \binom{n}{i} x^{n-i} h^i}{h} \quad (5.6a)$$

$$= \lim_{h \rightarrow 0} \sum_{i=1}^n \binom{n}{i} x^{n-i} h^{i-1} \quad (5.6b)$$

$$= \lim_{h \rightarrow 0} \binom{n}{1} x^{n-1} + \underbrace{\sum_{i=2}^n \binom{n}{i} x^{n-i} h^{i-1}}_{\text{dąży do 0 gdy } h \rightarrow 0} \quad (5.6c)$$

$$= \frac{n!}{1!(n-1)!} x^{n-1} = nx^{n-1}. \quad (5.6d)$$

5.1.1. Szereg Taylora

Szereg Taylora jest reprezentacją funkcji f w postaci nieskończonej sumy. Składniki tej sumy są określane przy pomocy pochodnych f w punkcie x_0 .

Definicja 5.3 (wielomian Taylora). **Wielomian Taylora** n -tego stopnia funkcji $f: \mathbb{R} \rightarrow \mathbb{R}$ w x_0 jest zdefiniowany jako

$$T_n(x) := \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k, \quad (5.7)$$

gdzie $f^{(k)}(x_0)$ jest k -tą pochodną f w x_0 (zakładamy, że pochodna istnieje), a $\frac{f^{(k)}(x_0)}{k!}$ to współczynniki wielomianu.

Definicja 5.4 (szereg Taylora). Dla gładkiej funkcji $f \in C^\infty$, $f: \mathbb{R} \rightarrow \mathbb{R}$ **szeregiem Taylora** f w x_0 nazywamy

$$T_\infty(x) = \sum_{k=0}^{\infty} \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k. \quad (5.8)$$

Dla $x_0 = 0$ otrzymujemy specjalny przypadek szeregu Taylora, nazywany **szeregiem Maclaurina**. Jeśli $f(x) = T_\infty(x)$, to f nazywamy **funkcją analityczną**.

Uwaga. Wielomian Taylora n -tego stopnia jest przybliżeniem funkcji, która niekoniecznie jest wielomianem. Wielomian Taylora jest podobny do funkcji f w otoczeniu punktu x_0 . Wielomian Taylora n -tego stopnia jest dokładną reprezentacją funkcji wielomianowej f k -tego stopnia takiej, że $k \leq n$, ponieważ wszystkie pochodne $f^{(i)}$ dla $i > k$ są równe zero.

Przykład 5.3 (wielomian Taylora)

Rozważamy wielomian

$$f(x) = x^4. \quad (5.9)$$

Chcemy wyznaczyć wielomian Taylora T_6 w punkcie $x_0 = 1$. Zaczynamy od znalezienia współczynników $f^{(k)}(1)$ dla $k = 0, \dots, 6$:

$$f(1) = 1 \quad (5.10)$$

$$f'(1) = 4 \quad (5.11)$$

$$f''(1) = 12 \quad (5.12)$$

$$f^{(3)}(1) = 24 \quad (5.13)$$

$$f^{(4)}(1) = 24 \quad (5.14)$$

$$f^{(5)}(1) = 0 \quad (5.15)$$

$$f^{(6)}(1) = 0. \quad (5.16)$$

$t^0 := 1$
dla wszystkich $t \in \mathbb{R}$.

$f \in C^\infty$ oznacza, że f jest ciągła i różniczkowalna nieskończenie wiele razy.

Po podstawieniu otrzymujemy następujący wielomian Taylora

$$T_6(x) = \sum_{k=0}^6 \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k \quad (5.17a)$$

$$= 1 + 4(x - 1) + 6(x - 1)^2 + 4(x - 1)^3 + (x - 1)^4 + 0. \quad (5.17b)$$

Po wymnożeniu i zmianie kolejności składników daje to

$$T_6(x) = (1 - 4 + 6 - 4 + 1) + x(4 - 12 + 12 - 4) \quad (5.18a)$$

$$+ x^2(6 - 12 + 6) + x^3(4 - 4) + x^4$$

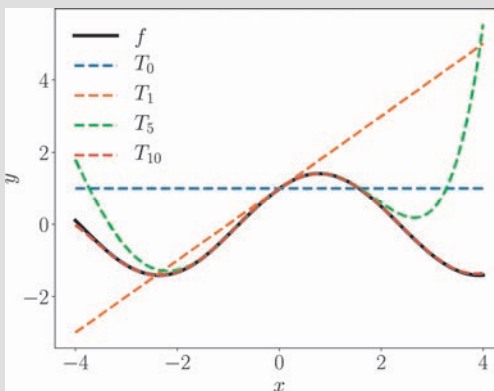
$$= x^4 = f(x). \quad (5.18b)$$

Otrzymaliśmy więc dokładną reprezentację naszej pierwotnej funkcji.

Przykład 5.4 (szereg Taylora)

Rozważmy funkcję z rysunku 5.4, która jest zdefiniowana w następujący sposób:

$$f(x) = \sin(x) + \cos(x) \in \mathcal{C}^\infty. \quad (5.19)$$



RYSUNEK 5.4. Wielomiany Taylora. Oryginalna funkcja $f(x) = \sin(x) + \cos(x)$ (czarna ciągła linia) jest aproksymowana wielomianami Taylora (kolorowe kreski) w otoczeniu punktu $x_0 = 0$. Wielomiany Taylora wyższego rzędu przybliżają funkcję f lepiej i bardziej globalnie. Wielomian T_{10} jest podobny do f w przedziale $[-4, 4]$

Chcemy znaleźć rozwinięcie f w szereg Taylora w punkcie $x_0 = 0$, czyli rozwinięcie f w szereg Maclaurina. Funkcja f ma następujące pochodne:

$$f(0) = \sin(0) + \cos(0) = 1 \quad (5.20)$$

$$f'(0) = \cos(0) - \sin(0) = 1 \quad (5.21)$$

$$f''(0) = -\sin(0) - \cos(0) = -1 \quad (5.22)$$

$$f^{(3)}(0) = -\cos(0) + \sin(0) = -1 \quad (5.23)$$

$$f^{(4)}(0) = \sin(0) + \cos(0) = f(0) = 1 \quad (5.24)$$

⋮

Na powyższej liście pochodnych możemy zauważyć pewną prawidłowość. Współczynniki w szeregu Taylora będą równe jedynie ± 1 (ponieważ $\sin(0) = 0$). Każda z tych wartości występuje w dwóch kolejnych wyrazach, po których następuje zmiana wartości na przeciwną. Ponadto $f^{(k+4)}(0) = f^{(k)}(0)$.

Dlatego pełnym rozwinięciem f w szereg Taylora w punkcie $x_0 = 0$ jest

$$T_\infty(x) = \sum_{k=0}^{\infty} \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k \quad (5.25a)$$

$$= 1 + x - \frac{1}{2!}x^2 - \frac{1}{3!}x^3 + \frac{1}{4!}x^4 + \frac{1}{5!}x^5 - \dots \quad (5.25b)$$

$$= 1 - \frac{1}{2!}x^2 + \frac{1}{4!}x^4 - \dots + x - \frac{1}{3!}x^3 + \frac{1}{5!}x^5 - \dots \quad (5.25c)$$

$$= \sum_{k=0}^{\infty} (-1)^k \frac{1}{(2k)!} x^{2k} + \sum_{k=0}^{\infty} (-1)^k \frac{1}{(2k+1)!} x^{2k+1} \quad (5.25d)$$

$$= \cos(x) + \sin(x). \quad (5.25e)$$

W powyższych przekształceniach wykorzystaliśmy reprezentację funkcji trygonometrycznych w postaci **szeregów potęgowych**:

$$\cos(x) = \sum_{k=0}^{\infty} (-1)^k \frac{1}{(2k)!} x^{2k} \quad (5.26)$$

$$\sin(x) = \sum_{k=0}^{\infty} (-1)^k \frac{1}{(2k+1)!} x^{2k+1}. \quad (5.27)$$

Na rysunku 5.4 pokazano kilka pierwszych wielomianów T_n dla $n = 0, 1, 5, 10$.

Uwaga. Szereg Taylora to szczególny przypadek szeregu potęgowego

$$f(x) = \sum_{k=0}^{\infty} a_k (x - c)^k, \quad (5.28)$$

w którym a_k są współczynnikami, a c jest stałą, która w szeregu Taylora przybiera formę określoną w definicji 5.4.

5.1.2. Zasady różniczkowania

Poniżej zamieściliśmy krótką listę podstawowych reguł różniczkowania, w których pochodną f oznaczyliśmy symbolem f' .

$$\text{Pochodna iloczynu: } (f(x)g(x))' = f'(x)g(x) + f(x)g'(x) \quad (5.29)$$

$$\text{Pochodna ilorazu: } \left(\frac{f(x)}{g(x)}\right)' = \frac{f'(x)g(x) - f(x)g'(x)}{(g(x))^2} \quad (5.30)$$

$$\text{Pochodna sumy: } (f(x) + g(x))' = f'(x) + g'(x) \quad (5.31)$$

$$\text{Reguła łańcuchowa: } (g(f(x)))' = (g \circ f)'(x) = g'(f(x))f'(x) \quad (5.32)$$

Zapis $g \circ f$ oznacza złożenie funkcji $x \mapsto f(x) \mapsto g(f(x))$.

Przykład 5.5 (reguła łańcuchowa)

Wykorzystajmy regułę łańcuchową do znalezienia pochodnej $h(x) = (2x + 1)^4$:

$$h(x) = (2x + 1)^4 = g(f(x)), \quad (5.33)$$

$$f(x) = 2x + 1, \quad (5.34)$$

$$g(f) = f^4. \quad (5.35)$$

Obliczamy pochodne f i g

$$f'(x) = 2, \quad (5.36)$$

$$g'(f) = 4f^3. \quad (5.37)$$

Pochodna funkcji h jest więc równa

$$h'(x) = g'(f)f'(x) = (4f^3) \cdot 2 \stackrel{(5.34)}{=} 4(2x + 1)^3 \cdot 2 = 8(2x + 1)^3. \quad (5.38)$$

Do wyznaczenia pochodnej wykorzystaliśmy regułę łańcuchową (równanie 5.32). W pochodnej $g'(f)$ wykorzystaliśmy definicję funkcji f .

5.2. Pochodne cząstkowe i gradienty

W podrozdziale 5.1 przedstawiliśmy różniczkowanie funkcji f zmiennej skalarnej $x \in \mathbb{R}$. Poniżej rozważymy ogólny przypadek, w którym funkcja f zależy od jednej lub więcej zmiennych $\mathbf{x} \in \mathbb{R}^n$, np. $f(\mathbf{x}) = f(x_1, x_2)$. Uogólnienie pochodnej na funkcje kilku zmiennych nazywamy **gradientem**.

Gradient funkcji f względem wektora \mathbf{x} obliczamy, *analizując jedną zmienną naraz*. Pozostałe zmienne traktujemy jako stałe. Gradient jest zatem zbiorem **pochodnych cząstkowych**.

Definicja 5.5 (pochodna cząstkowa). **Pochodne cząstkowe** funkcji $f: \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto f(\mathbf{x})$, $\mathbf{x} \in \mathbb{R}^n$ n -zmiennych x_1, \dots, x_n są zdefiniowane w następujący sposób:

$$\begin{aligned} \frac{\partial f}{\partial x_1} &= \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(\mathbf{x})}{h} \\ &\quad \vdots \\ \frac{\partial f}{\partial x_n} &= \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{n-1}, x_n + h) - f(\mathbf{x})}{h}. \end{aligned} \quad (5.39)$$

Pochodne te zapisuje się w postaci wektora wierszowego:

$$\nabla_{\mathbf{x}} f = \text{grad} f = \frac{df}{d\mathbf{x}} = \left[\frac{\partial f(\mathbf{x})}{\partial x_1} \quad \frac{\partial f(\mathbf{x})}{\partial x_2} \quad \dots \quad \frac{\partial f(\mathbf{x})}{\partial x_n} \right] \in \mathbb{R}^{1 \times n}, \quad (5.40)$$

gdzie n to liczba zmiennych, a 1 to wymiar obrazu/przeciwdziedziny/zbioru wartości funkcji f . W powyższych równaniach \mathbf{x} jest wektorem kolumnowym: $\mathbf{x} = [x_1, \dots, x_n]^T \in \mathbb{R}^n$. Wektor wierszowy z równania 5.40 nazywany jest **gradientem** funkcji f lub jej *macierzą Jacobiego*. Wektor ten jest uogólnieniem pochodnej z podrozdziału 5.1.

Uwaga. Powyższa definicja macierzy Jacobiego jest szczególnym przypadkiem ogólnej definicji macierzy Jacobiego (macierzy pochodnych cząstkowych) dla funkcji o wartościach wektorowych. Powróćmy do niej w podrozdziale 5.3.

Przykład 5.6 (obliczanie pochodnych cząstkowych za pomocą reguły łańcuchowej)

Funkcja $f(x, y) = (x + 2y^3)^2$ ma następujące pochodne cząstkowe

$$\frac{\partial f(x, y)}{\partial x} = 2(x + 2y^3) \frac{\partial}{\partial x} (x + 2y^3) = 2(x + 2y^3), \quad (5.41)$$

$$\frac{\partial f(x, y)}{\partial y} = 2(x + 2y^3) \frac{\partial}{\partial y} (x + 2y^3) = 12(x + 2y^3)y^2. \quad (5.42)$$

Do ich obliczenia wykorzystaliśmy regułę łańcuchową (równanie 5.32).

Uwaga (gradient jako wektor wierszowy). W literaturze nierzadko definiuje się gradient jako wektor kolumnowy. Wynika to z ogólnie przyjętej konwencji, która zakłada, że wektory są zazwyczaj wektorami kolumnowymi. Istnieją dwa powody, dla których zdefiniowaliśmy wektor gradientu jako wektor wierszowy. Po pierwsze, pozwoli nam to w konsekwentny sposób uogólnić gradient na funkcje o wartościach wektorowych $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ (gradient stanie się macierzą). Po drugie, taka reprezentacja pozwala nam zastosować regułę łańcuchową dla funkcji wielowymiarowych bez zwracania uwagi na wymiary gradientu. Oba zagadnienia omówimy w podrozdziale 5.3.

W obliczeniach możemy wykorzystać wyniki z różniczkowania funkcji jednej zmiennej. Każda pochodna cząstkowa jest pochodną względem jednej zmiennej.

Przykład 5.7 (gradient)

Pochodne cząstkowe (pochodne f względem x_1 i x_2) funkcji $f(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3 \in \mathbb{R}$ to

$$\frac{\partial f(x_1, x_2)}{\partial x_1} = 2x_1 x_2 + x_2^3 \quad (5.43)$$

$$\frac{\partial f(x_1, x_2)}{\partial x_2} = x_1^2 + 3x_1 x_2^2. \quad (5.44)$$

Gradientem f jest

$$\frac{df}{dx} = \left[\frac{\partial f(x_1, x_2)}{\partial x_1} \quad \frac{\partial f(x_1, x_2)}{\partial x_2} \right] = [2x_1 x_2 + x_2^3 \quad x_1^2 + 3x_1 x_2^2] \in \mathbb{R}^{1 \times 2}. \quad (5.45)$$

5.2.1. Podstawowe zasady obliczania pochodnych cząstkowych

W przypadku wielowymiarowym, w którym $\mathbf{x} \in \mathbb{R}^n$, nadal obowiązują podstawowe reguły różniczkowania, które znamy ze szkoły (np. reguła sumy, reguła iloczynu, reguła łańcuchowa; patrz też punkt 5.1.2). Podczas obliczania pochodnych względem wektorów $\mathbf{x} \in \mathbb{R}^n$ musimy zachować ostrożność, ponieważ w gradientach będą teraz występowały wektory i macierze, których mnożenie nie jest przemienne (tj. kolejność ma znaczenie, patrz punkt 2.2.1).

Pochodna iloczynu: $(fg)' = f'g + fg'$

Pochodna sumy: $(f + g)' = f' + g'$

Reguła łańcuchowa: $(g(f))' = g'(f)f'$

Oto ogólne reguły pozwalające obliczyć pochodną iloczynu i sumy oraz uogólniona reguła łańcuchowa:

$$\text{Pochodna iloczynu:} \quad \frac{\partial}{\partial \mathbf{x}} (f(\mathbf{x})g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} g(\mathbf{x}) + f(\mathbf{x}) \frac{\partial g}{\partial \mathbf{x}} \quad (5.46)$$

$$\text{Pochodna sumy:} \quad \frac{\partial}{\partial \mathbf{x}} (f(\mathbf{x}) + g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} + \frac{\partial g}{\partial \mathbf{x}} \quad (5.47)$$

$$\text{Reguła łańcuchowa:} \quad \frac{\partial}{\partial \mathbf{x}} (g \circ f)(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (g(f(\mathbf{x}))) = \frac{\partial g}{\partial f} \frac{\partial f}{\partial \mathbf{x}} \quad (5.48)$$

Przyjrzyjmy się bliżej regule łańcuchowej. Przypomina ona do pewnego stopnia reguły mnożenia macierzy, które wymagają, aby „sąsiadujące” ze sobą wymiary macierzy były ze sobą zgodne (patrz punkt 2.2.1). Spojrzenie na równanie 5.48 od lewej do prawej pozwala zauważyć, że reguła łańcuchowa wykazuje podobne właściwości. ∂f pojawia się w „mianowniku” pierwszego czynnika i w „liczniku” drugiego. Jeśli pomnożymy te dwa „ułamki”, to zauważymy, że mnożenie jest możliwe, ponieważ „wymiar” obu ∂f są „zgodne”. W efekcie czynnik ∂f ulega „skróceniu” i w wyniku pozostanie jedynie $\partial g / \partial \mathbf{x}$.

To tylko matematycznie niepoprawna intuicja. Pochodne cząstkowe nie są ułamkami.

5.2.2. Reguła łańcuchowa

Rozważmy funkcję $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ dwóch zmiennych x_1 i x_2 . Niech $x_1(t)$ i $x_2(t)$ będą funkcjami zmiennej t . Aby obliczyć gradient f względem t , musimy zastosować regułę łańcuchową (równanie 5.48) dla funkcji wielowymiarowych:

$$\frac{df}{dt} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1(t)}{\partial t} \\ \frac{\partial x_2(t)}{\partial t} \end{bmatrix} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}, \quad (5.49)$$

gdzie d oznacza gradient, a ∂ pochodne cząstkowe.

Przykład 5.8

Jeżeli $f(x_1, x_2) = x_1^2 + 2x_2$ oraz $x_1 = \sin t$ i $x_2 = \cos t$, to

$$\frac{df}{dt} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t} \quad (5.50a)$$

$$= 2 \sin t \frac{\partial \sin t}{\partial t} + 2 \frac{\partial \cos t}{\partial t} \quad (5.50b)$$

$$= 2 \sin t \cos t - 2 \sin t = 2 \sin t (\cos t - 1) \quad (5.50c)$$

jest pochodną f względem t .

Jeśli $f(x_1, x_2)$ jest funkcją x_1 i x_2 , a $x_1(s, t)$ i $x_2(s, t)$ są funkcjami dwóch zmiennych s i t , to za pomocą reguły łańcuchowej możemy obliczyć ich pochodne cząstkowe w następujący sposób:

$$\frac{\partial f}{\partial s} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial s} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial s}, \quad (5.51)$$

$$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial t} + \frac{\partial f}{\partial x_2} \frac{\partial x_2}{\partial t}. \quad (5.52)$$

Gradient otrzymamy z mnożenia macierzy:

$$\frac{df}{d(s, t)} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial (s, t)} = \underbrace{\begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix}}_{=\frac{\partial f}{\partial \mathbf{x}}} \underbrace{\begin{bmatrix} \frac{\partial x_1}{\partial s} & \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial s} & \frac{\partial x_2}{\partial t} \end{bmatrix}}_{=\frac{\partial \mathbf{x}}{\partial (s, t)}}. \quad (5.53)$$

Regułę łańcuchową możemy zapisać w postaci mnożenia macierzy.

Ten zwięzły sposób zapisu reguły łańcuchowej w postaci mnożenia macierzy ma sens tylko wtedy, gdy gradient jest zdefiniowany jako wektor wierszowy. W przeciwnym razie musielibyśmy transponować gradienty, tak aby zgadzały się wymiary mnożonych macierzy. Nie jest to skomplikowane, o ile gradient jest wektorem lub macierzą. Jednak gdy gradient stanie się tensorem (co omówimy w dalszej części rozdziału), znalezienie transpozycji nie będzie już trywialnym zadaniem.

Uwaga (weryfikacja poprawności implementacji gradientu). Do numerycznego sprawdzenia wyników obliczeń gradientów w programach komputerowych możemy wykorzystać definicję pochodnych cząstkowych w postaci granicy ilorazu różnicowego (patrz równanie 5.39). Do sprawdzenia poprawności implementacji możemy wykorzystać różnice skończone. Wybieramy małą wartość h (np. $h = 10^{-4}$) i porównujemy przybliżenie różnicy skończonej z równania 5.39 z naszą (analityczną) implementacją gradientu. Jeśli błąd jest mały, to prawdopodobnie nasza implementacja gradientu jest prawidłowa. „Mały” może oznaczać, że $\sqrt{\frac{\sum_i (dh_i - df_i)^2}{\sum_i (dh_i + df_i)^2}} < 10^{-6}$, gdzie dh_i jest przybliżeniem za pomocą różnic skończonych, a df_i jest różniczką funkcji względem i -tej zmiennej x_i .

5.3. Gradienty funkcji o wartościach wektorowych

Dotychczas omawialiśmy pochodne cząstkowe i gradienty funkcji $f: \mathbb{R}^n \rightarrow \mathbb{R}$ przekształcających wektory na liczby rzeczywiste. Poniżej uogólnimy pojęcie gradientu na funkcje o wartościach wektorowych (pola wektorowe) $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$, gdzie $n \geq 1$ oraz $m > 1$.

Wektorem wartości funkcji $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$, która przyjmuje wektor $\mathbf{x} = [x_1, \dots, x_n]^T \in \mathbb{R}^n$, jest

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^m. \quad (5.54)$$

Powyższy zapis pozwala nam traktować funkcję o wartościach wektorowych $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ jako wektor funkcji $[f_1, \dots, f_m]^T$ postaci $f_i: \mathbb{R}^n \rightarrow \mathbb{R}$, których zbiorem wartości jest \mathbb{R} . Różniczkowanie funkcji f_i odbywa się za pomocą reguł różniczkowania omówionych w podrozdziale 5.2.

A zatem pochodną cząstkową funkcji o wartościach wektorowych $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ względem $x_i \in \mathbb{R}$ ($i = 1, \dots, n$) jest wektor

$$\frac{\partial \mathbf{f}}{\partial x_i} = \begin{bmatrix} \frac{\partial f_1}{\partial x_i} \\ \vdots \\ \frac{\partial f_m}{\partial x_i} \end{bmatrix} = \begin{bmatrix} \lim_{h \rightarrow 0} \frac{f_1(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_n) - f_1(\mathbf{x})}{h} \\ \vdots \\ \lim_{h \rightarrow 0} \frac{f_m(x_1, \dots, x_{i-1}, x_i + h, x_{i+1}, \dots, x_n) - f_m(\mathbf{x})}{h} \end{bmatrix} \in \mathbb{R}^m. \quad (5.55)$$

Z równania 5.40 wiemy, że gradientem \mathbf{f} względem wektora jest wektor wierszowy pochodnych cząstkowych. W równaniu 5.55 każda pochodna cząstkowa $\partial \mathbf{f} / \partial x_i$ jest wektorem kolumnowym. Gradient $\mathbf{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$ względem $\mathbf{x} \in \mathbb{R}^n$ to zbiór tych wszystkich pochodnych cząstkowych:

$$\frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \left[\boxed{\frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_1}} \dots \boxed{\frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_n}} \right] \quad (5.56a)$$

$$= \begin{bmatrix} \boxed{\frac{\partial f_1(\mathbf{x})}{\partial x_1}} & \dots & \boxed{\frac{\partial f_1(\mathbf{x})}{\partial x_n}} \\ \vdots & & \vdots \\ \boxed{\frac{\partial f_m(\mathbf{x})}{\partial x_1}} & \dots & \boxed{\frac{\partial f_m(\mathbf{x})}{\partial x_n}} \end{bmatrix} \in \mathbb{R}^{m \times n} \quad (5.56b)$$

Gradient funkcji $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ jest macierzą o wymiarach $m \times n$.

Definicja 5.6 (macierz Jacobiego). Zbiór wszystkich pochodnych cząstkowych pierwszego rzędu funkcji o wartościach wektorowych $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ względem wektora $\mathbf{x} \in \mathbb{R}^n$ nazywamy **macierzą Jacobiego** (ang. *Jacobian*). Macierz Jacobiego J jest macierzą o wymiarach $m \times n$, którą definiujemy w następujący sposób:

$$J = \nabla_{\mathbf{x}} f = \frac{d\mathbf{f}(\mathbf{x})}{d\mathbf{x}} = \begin{bmatrix} \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial \mathbf{f}(\mathbf{x})}{\partial x_n} \end{bmatrix} \quad (5.57)$$

$$= \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_1(\mathbf{x})}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix}, \quad (5.58)$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \quad J(i, j) = \frac{\partial f_i}{\partial x_j}. \quad (5.59)$$

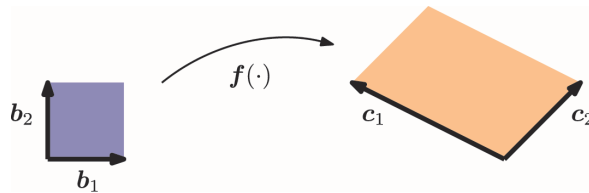
Szczególnym przypadkiem równania 5.58 jest funkcja $f: \mathbb{R}^n \rightarrow \mathbb{R}^1$, która odwzorowuje wektor $\mathbf{x} \in \mathbb{R}^n$ na skalar (np. $f(\mathbf{x}) = \sum_{i=1}^n x_i$). Macierz Jacobiego tej funkcji jest wektorem wierszowym (macierzą o wymiarach $1 \times n$, patrz równanie 5.40).

Uwaga. W tej książce używamy macierzy Jacobiego zapisanej w postaci licznikowej, tj. układzie, w którym pochodna $d\mathbf{f}/d\mathbf{x}$ funkcji $\mathbf{f} \in \mathbb{R}^m$ względem $\mathbf{x} \in \mathbb{R}^n$ jest macierzą o wymiarach $m \times n$, w której wiersze odpowiadają elementom \mathbf{f} , a kolumny elementom \mathbf{x} (patrz równanie 5.58). Istnieje również postać mianownikowa, która jest transpozycją postaci licznikowej. W tej książce stosujemy postać licznikową.

W podrozdziale 6.7 zobaczysz przykład wykorzystania macierzy Jacobiego do zmiany zmiennych w rozkładach prawdopodobieństwa. Stopień skalowania spowodowany zmianą zmiennych jest określony przez wyznacznik macierzy Jacobiego (nazywany *jakobianem*).

W podrozdziale 4.1 wspomnieliśmy, że wyznacznik można wykorzystać do obliczenia pola równoległoboku. Jeśli mamy dwa wektory $\mathbf{b}_1 = [1, 0]^T$ i $\mathbf{b}_2 = [0, 1]^T$ reprezentujące boki kwadratu jednostkowego (zaznaczone kolorem niebieskim, patrz rysunek 5.5), to pole tego kwadratu jest równe

$$\left| \det \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right| = 1. \quad (5.60)$$



RYСУNEK 5.5. Jakobian funkcji f może być wykorzystany do ustalenia współczynnika skalowania niebieskiego obszaru w pomarańczowy

Wartość bezwzględna poniższego wyznacznika (patrz podrozdział 4.1) odpowiada polu powierzchni równoległoboku o bokach $\mathbf{c}_1 = [-2, 1]^T$, $\mathbf{c}_2 = [1, 1]^T$ (na rysunku 5.5 zaznaczony kolorem pomarańczowym)

$$\left| \det \begin{pmatrix} -2 & 1 \\ 1 & 1 \end{pmatrix} \right| = |-3| = 3. \quad (5.61)$$

Pole tego równoległoboku jest trzykrotnością pola kwadratu jednostkowego. Ten współczynnik skalowania możemy wyznaczyć, znajdując odwzorowanie, które przekształca jeden kształt w drugi. W języku algebry liniowej oznacza to przeprowadzenie zmiany zmiennych z $(\mathbf{b}_1, \mathbf{b}_2)$ na $(\mathbf{c}_1, \mathbf{c}_2)$. W naszym przypadku odwzorowanie jest liniowe i wartość bezwzględna wyznacznika tego odwzorowania daje nam poszukiwany współczynnik skalowania.

Poniżej przedstawimy dwa podejścia pozwalające znaleźć to odwzorowanie. W pierwszym wykorzystamy fakt, że odwzorowanie jest liniowe, co pozwoli nam wykorzystać narzędzia z rozdziału 2. do zidentyfikowania odwzorowania. W drugim znajdziemy odwzorowanie za pomocą pochodnych cząstkowych i narzędzi, które omówiliśmy w tym rozdziale.

Podejście 1. W podejściu algebraicznym zauważamy, że $\{\mathbf{b}_1, \mathbf{b}_2\}$ i $\{\mathbf{c}_1, \mathbf{c}_2\}$ są bazami przestrzeni \mathbb{R}^2 (patrz punkt 2.6.1). Interesującym nas odwzorowaniem jest zmiana bazy z $(\mathbf{b}_1, \mathbf{b}_2)$ na $(\mathbf{c}_1, \mathbf{c}_2)$. W tym podejściu będziemy więc szukać macierzy transformacji, która realizuje tę zmianę bazy. W oparciu o informacje z punktu 2.7.2 otrzymujemy następującą macierz zmiany bazy

$$\mathbf{J} = \begin{bmatrix} -2 & 1 \\ 1 & 1 \end{bmatrix}. \quad (5.62)$$

$\mathbf{J}\mathbf{b}_1 = \mathbf{c}_1$ a $\mathbf{J}\mathbf{b}_2 = \mathbf{c}_2$. Wartość bezwzględna wyznacznika \mathbf{J} jest poszukiwanym współczynnikiem skalowania. $|\det(\mathbf{J})| = 3$, co oznacza, że powierzchnia równoległoboku rozpinanego przez $(\mathbf{c}_1, \mathbf{c}_2)$ jest trzykrotnie większa niż powierzchnia kwadratu rozpinanego przez $(\mathbf{b}_1, \mathbf{b}_2)$.

Podejście 2. Podejście oparte na algebrze liniowej sprawdza się jedynie w przypadku przekształceń liniowych. W przypadku przekształceń nieliniowych (które odegrają ważną rolę w podrozdziale 6.7) stosujemy bardziej ogólne podejście, wykorzystujące pochodne cząstkowe.

W tym podejściu rozważamy funkcję $\mathbf{f}: \mathbb{R}^2 \rightarrow \mathbb{R}^2$, która dokonuje zmiany zmiennych. W naszym przykładzie \mathbf{f} odwzorowuje współrzędne dowolnego wektora $\mathbf{x} \in \mathbb{R}^2$ w bazie $(\mathbf{b}_1, \mathbf{b}_2)$ na współrzędne $\mathbf{y} \in \mathbb{R}^2$ w bazie $(\mathbf{c}_1, \mathbf{c}_2)$. Chcemy zidentyfikować odwzorowanie, aby ustalić, jak w wyniku przekształcenia \mathbf{f} zmienia się powierzchnia (lub objętość). W tym celu musimy ustalić, jak zmienia się $\mathbf{f}(\mathbf{x})$, jeśli odrobinę zmodyfikujemy \mathbf{x} . Na to pytanie odpowiada dokładnie macierz Jacobiego $\frac{d\mathbf{f}}{d\mathbf{x}} \in \mathbb{R}^{2 \times 2}$. Możemy więc zapisać, że

$$y_1 = -2x_1 + x_2 \quad (5.63)$$

$$y_2 = x_1 + x_2. \quad (5.64)$$

Otrzymaliśmy zależność funkcyjną pomiędzy \mathbf{x} i \mathbf{y} , która pozwala nam obliczyć pochodne cząstkowe

$$\frac{\partial y_1}{\partial x_1} = -2, \quad \frac{\partial y_1}{\partial x_2} = 1, \quad \frac{\partial y_2}{\partial x_1} = 1, \quad \frac{\partial y_2}{\partial x_2} = 1. \quad (5.65)$$

Z geometrycznego punktu widzenia jacobian określa współczynnik powiększenia/skalowania przekształcanego obszaru lub objętości.

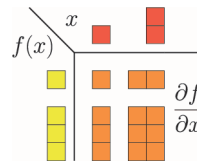
Po ich zapisaniu w macierzy otrzymujemy następującą macierz Jacobiego:

$$J = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} \end{bmatrix} = \begin{bmatrix} -2 & 1 \\ 1 & 1 \end{bmatrix}. \quad (5.66)$$

Macierz Jacobiego reprezentuje poszukiwaną przez nas transformację współrzędnych. Reprezentacja ta jest dokładna, jeżeli badane przekształcenie jest liniowe (tak jest w tym przypadku). Równanie 5.66 odpowiada dokładnie macierzy zmiany bazy z równania 5.62. Jeśli transformacja współrzędnych jest nieliniowa, to macierz Jacobiego lokalnie aproksymuje tę nieliniową transformację za pomocą odwzorowania liniowego. Wartość bezwzględna wyznacznika macierzy Jacobiego $|\det(J)|$ jest współczynnikiem skalowania powierzchni lub objętości podczas przekształcania współrzędnych. W tym przykładzie $|\det(J)| = 3$.

Jakobiany i przekształcenia zmiennych odegrają ważną rolę w podrozdziale 6.7, w którym będziemy przekształcać zmienne losowe i rozkłady prawdopodobieństwa. Przekształcenia te są niezwykle istotne w uczeniu głębokich sieci neuronowych z użyciem tzw. *reparametryzacji*, zwanej również **analizą nieskończonych perturbacji**.

W tym rozdziale omówiliśmy pochodne funkcji. Podsumowanie ich wymiarów pokazano na rysunku 5.6. Jeśli $f: \mathbb{R} \rightarrow \mathbb{R}$, to gradient jest po prostu skalar (lewy górny róg rysunku). Dla $f: \mathbb{R}^D \rightarrow \mathbb{R}$ gradient jest wektorem wierszowym o wymiarach $1 \times D$ (prawy górny róg rysunku). Dla $f: \mathbb{R} \rightarrow \mathbb{R}^E$ gradient jest wektorem kolumnowym o wymiarach $E \times 1$, a dla $f: \mathbb{R}^D \rightarrow \mathbb{R}^E$ otrzymujemy macierz gradientu o rozmiarze $E \times D$.



RYСУNEK 5.6. Wymiary pochodnych (cząstkowych)

Przykład 5.9 (gradient funkcji o wartościach wektorowych)

Dane są:

$$f(\mathbf{x}) = \mathbf{A}\mathbf{x}, \quad f(\mathbf{x}) \in \mathbb{R}^M, \quad \mathbf{A} \in \mathbb{R}^{M \times N}, \quad \mathbf{x} \in \mathbb{R}^N$$

Aby obliczyć gradient $d\mathbf{f}/d\mathbf{x}$, najpierw określamy jego wymiary. Ponieważ $f: \mathbb{R}^N \rightarrow \mathbb{R}^M$, gradient $d\mathbf{f}/d\mathbf{x} \in \mathbb{R}^{M \times N}$. Następnie wyznaczamy pochodne cząstkowe f względem każdego x_j :

$$f_i(\mathbf{x}) = \sum_{j=1}^N A_{ij}x_j \Rightarrow \frac{\partial f_i}{\partial x_j} = A_{ij} \quad (5.67)$$

i zapisujemy je w macierzy Jacobiego

$$\frac{df}{dx} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_N} \\ \vdots & & \vdots \\ \frac{\partial f_M}{\partial x_1} & \cdots & \frac{\partial f_M}{\partial x_N} \end{bmatrix} = \begin{bmatrix} A_{11} & \cdots & A_{1N} \\ \vdots & & \vdots \\ A_{M1} & \cdots & A_{MN} \end{bmatrix} = \mathbf{A} \in \mathbb{R}^{M \times N}. \quad (5.68)$$

Przykład 5.10 (reguła łańcucha)

Rozważmy funkcję $h: \mathbb{R} \rightarrow \mathbb{R}$, $h(t) = (f \circ g)(t)$, gdzie

$$f: \mathbb{R}^2 \rightarrow \mathbb{R} \quad (5.69)$$

$$g: \mathbb{R} \rightarrow \mathbb{R}^2 \quad (5.70)$$

$$f(\mathbf{x}) = \exp(x_1 x_2^2) \quad (5.71)$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = g(t) = \begin{bmatrix} t \cos t \\ t \sin t \end{bmatrix}. \quad (5.72)$$

Wyznamy gradient h względem t . Ponieważ $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ i $g: \mathbb{R} \rightarrow \mathbb{R}^2$, możemy zauważyć, że

$$\frac{\partial f}{\partial \mathbf{x}} \in \mathbb{R}^{1 \times 2}, \quad \frac{\partial g}{\partial t} \in \mathbb{R}^{2 \times 1}. \quad (5.73)$$

Interesującą nas gradient obliczamy za pomocą reguły łańcuchowej:

$$\frac{dh}{dt} = \frac{\partial f}{\partial \mathbf{x}} \frac{\partial \mathbf{x}}{\partial t} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} \begin{bmatrix} \frac{\partial x_1}{\partial t} \\ \frac{\partial x_2}{\partial t} \end{bmatrix} \quad (5.74a)$$

$$= [\exp(x_1 x_2^2) x_2^2 \quad 2 \exp(x_1 x_2^2) x_1 x_2] \begin{bmatrix} \cos t - t \sin t \\ \sin t + t \cos t \end{bmatrix} \quad (5.74b)$$

$$= \exp(x_1 x_2^2) (x_2^2 (\cos t - t \sin t) + 2 x_1 x_2 (\sin t + t \cos t)), \quad (5.74c)$$

gdzie $x_1 = t \cos t$ i $x_2 = t \sin t$; patrz równanie 5.72.

Model ten omówimy znacznie bardziej szczegółowo w rozdziale 9., w którym przedstawimy go w kontekście regresji liniowej. W regresji liniowej wykorzystuje się pochodne minimalnkwadratowej funkcji straty L względem parametrów $\boldsymbol{\theta}$.

Przykład 5.11 (gradient minimalnkwadratowej straty w modelu liniowym)

Rozważmy model liniowy

$$\mathbf{y} = \boldsymbol{\Phi} \boldsymbol{\theta}, \quad (5.75)$$

gdzie $\boldsymbol{\theta} \in \mathbb{R}^D$ jest wektorem parametrów, $\boldsymbol{\Phi} \in \mathbb{R}^{N \times D}$ to cechy wejściowe, a $\mathbf{y} \in \mathbb{R}^N$ to obserwacje. Zdefiniujmy funkcje

$$L(\mathbf{e}) := \|\mathbf{e}\|^2, \quad (5.76)$$

$$\mathbf{e}(\boldsymbol{\theta}) := \mathbf{y} - \boldsymbol{\Phi} \boldsymbol{\theta}. \quad (5.77)$$

Interesuje nas gradient $\frac{\partial L}{\partial \theta}$. Do jego wyznaczenia wykorzystamy regułę łańcuchową. L jest nazywana **minimalnokwadratową funkcją straty**.

Przed przystąpieniem do obliczeń określamy wymiary gradientu:

$$\frac{\partial L}{\partial \theta} \in \mathbb{R}^{1 \times D}. \quad (5.78)$$

Reguła łańcuchowa pozwala nam zapisać gradient jako

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial e} \frac{\partial e}{\partial \theta}. \quad (5.79)$$

d -ty element gradientu to

$$\frac{\partial L}{\partial \theta} [1, d] = \sum_{n=1}^N \frac{\partial L}{\partial e} [n] \frac{\partial e}{\partial \theta} [n, d]. \quad (5.80)$$

Korzystamy z faktu, że $\|e\|^2 = e^T e$ (patrz podrozdział 3.2), który pozwala zapisać, że

$$\frac{\partial L}{\partial e} = 2e^T \in \mathbb{R}^{1 \times N}. \quad (5.81)$$

Ponadto

$$\frac{\partial e}{\partial \theta} = -\Phi \in \mathbb{R}^{N \times D}. \quad (5.82)$$

Naszą pożądaną pochodną jest zatem

$$\frac{\partial L}{\partial \theta} = -2e^T \Phi \stackrel{(5.77)}{=} - \underbrace{2(\mathbf{y}^T - \boldsymbol{\theta}^T \Phi^T)}_{1 \times N} \underbrace{\Phi}_{N \times D} \in \mathbb{R}^{1 \times D}. \quad (5.83)$$

Uwaga. Ten sam wynik możemy otrzymać bez stosowania reguły łańcuchowej. Wystarczy spojrzeć na funkcję

$$L_2(\theta) := \|\mathbf{y} - \Phi \theta\|^2 = (\mathbf{y} - \Phi \theta)^T (\mathbf{y} - \Phi \theta). \quad (5.84)$$

To podejście sprawdza się w przypadku prostych funkcji, takich jak L_2 , ale staje się niepraktyczne w bardziej złożonych przypadkach.

5.4. Gradienty macierzy

Tensory możemy traktować jako wielowymiarowe tablice.

W niektórych sytuacjach musimy obliczyć gradienty macierzy względem wektorów (lub innych macierzy). Obliczenia te dają w wyniku wielowymiarowe tensory. Tensory możemy traktować jako wielowymiarowe tablice, w których zapisano pochodne cząstkowe. Na przykład gradient macierzy \mathbf{A} o wymiarach $m \times n$ względem macierzy \mathbf{B} o wymiarach $p \times q$ to obiekt o wymiarach $(m \times n) \times (p \times q)$, tj. czterowymiarowy tensor \mathbf{J} , którego elementy są zdefiniowane jako $J_{ijkl} = \partial A_{ij} / \partial B_{kl}$.

Ponieważ macierze reprezentują odwzorowania liniowe, możemy wykorzystać fakt istnienia izomorfizmu (liniowego, odwracalnego odwzorowania) przestrzeni wektorowej $\mathbb{R}^{m \times n}$ macierzy $m \times n$ w przestrzeń \mathbb{R}^{mn} mn -elementowych wektorów. Pozwala to przekształcić nasze macierze w wektory

Macierze można przekształcić w wektory, układając kolumny macierzy jedna za drugą (tzw. spłaszczanie).

o długościach mn i pq . Gradienty względem mn -elementowych wektorów to macierze Jacobiego o wymiarach $mn \times pq$. Oba podejścia pokazano na rysunku 5.7. W zastosowaniach praktycznych często chcemy przekształcić macierze w wektory, które pozwolą nam kontynuować pracę z macierzami Jacobiego. W przypadku macierzy Jacobiego reguła łańcuchowa (równanie 5.48) sprowadza się do prostego mnożenia macierzy, podczas gdy w przypadku tensorów Jacobiego musielibyśmy pilnować zgodności wymiarów, względem których sumujemy.

Przykład 5.12 (gradient wektorów względem macierzy)

Rozważmy następujący przykład:

$$\mathbf{f} = \mathbf{A}\mathbf{x}, \quad \mathbf{f} \in \mathbb{R}^M, \quad \mathbf{A} \in \mathbb{R}^{M \times N}, \quad \mathbf{x} \in \mathbb{R}^N. \quad (5.85)$$

Chcemy wyznaczyć gradient $d\mathbf{f}/d\mathbf{A}$. Ponownie rozpoczynamy od ustalenia wymiarów:

$$\frac{d\mathbf{f}}{d\mathbf{A}} \in \mathbb{R}^{M \times (M \times N)}. \quad (5.86)$$

Z definicji gradient jest zbiorem pochodnych cząstkowych:

$$\frac{d\mathbf{f}}{d\mathbf{A}} = \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{A}} \\ \vdots \\ \frac{\partial f_M}{\partial \mathbf{A}} \end{bmatrix}, \quad \frac{\partial f_i}{\partial \mathbf{A}} \in \mathbb{R}^{1 \times (M \times N)}. \quad (5.87)$$

W obliczaniu pochodnych cząstkowych pomoże nam jawne zapisanie formuły iloczynu macierzy z wektorem

$$f_i = \sum_{j=1}^N A_{ij}x_j, \quad i = 1, \dots, M. \quad (5.88)$$

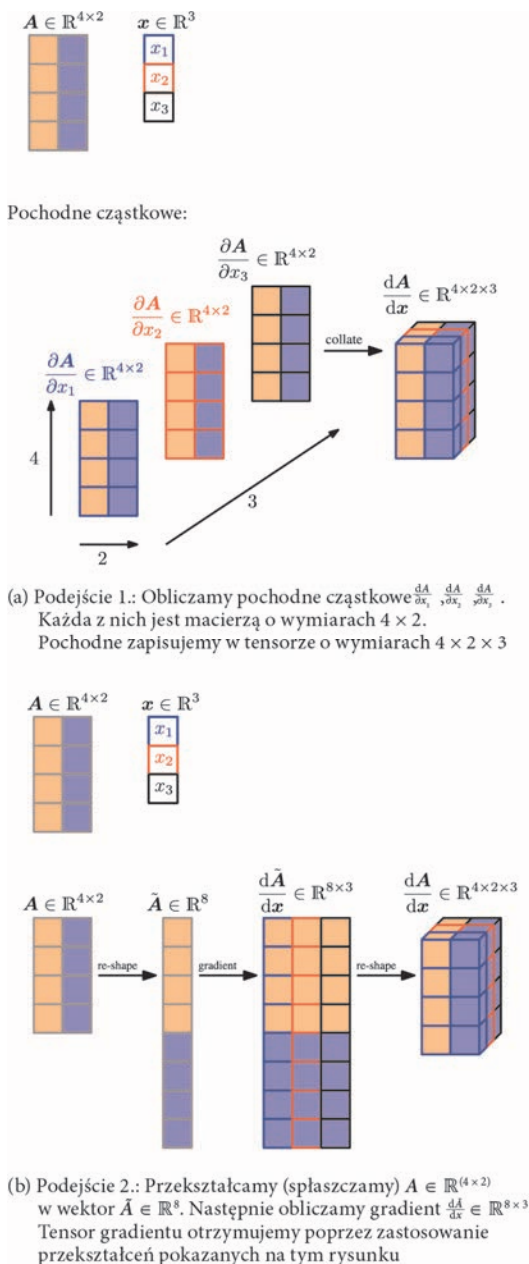
Pochodne cząstkowe mają więc następującą postać:

$$\frac{\partial f_i}{\partial A_{iq}} = x_q. \quad (5.89)$$

To pozwala nam obliczyć pochodne cząstkowe f_i względem wiersza macierzy \mathbf{A} :

$$\frac{\partial f_i}{\partial A_{i,:}} = \mathbf{x}^T \in \mathbb{R}^{1 \times 1 \times N}, \quad (5.90)$$

$$\frac{\partial f_i}{\partial A_{k \neq i, :}} = \mathbf{0}^T \in \mathbb{R}^{1 \times 1 \times N}. \quad (5.91)$$



- (a) Podejście 1.: Obliczamy pochodne cząstkowe $\frac{dA}{dx_1}$, $\frac{dA}{dx_2}$, $\frac{dA}{dx_3}$. Każda z nich jest macierzą o wymiarach 4×2 . Pochodne zapisujemy w tensorze o wymiarach $4 \times 2 \times 3$

- (b) Podejście 2.: Przekształcamy (spłaszczamy) $A \in \mathbb{R}^{(4 \times 2)}$ w wektor $\tilde{A} \in \mathbb{R}^8$. Następnie obliczamy gradient $\frac{d\tilde{A}}{dx} \in \mathbb{R}^{8 \times 3}$. Tensor gradientu otrzymujemy poprzez zastosowanie przekształceń pokazanych na tym rysunku

RYSUNEK 5.7. Wizualizacja procesu obliczania gradientu macierzy względem wektora. Chcemy wyznaczyć gradient $A \in \mathbb{R}^{4 \times 2}$ względem wektora $x \in \mathbb{R}^3$. Wiemy, że $\frac{dA}{dx} \in \mathbb{R}^{4 \times 2 \times 3}$. Do wyznaczenia gradientu stosujemy dwa podejścia: (a) zapisujemy pochodne cząstkowe w tensorze Jacobiego; (b) spłaszczamy macierz do wektora i obliczamy macierz Jacobiego, którą następnie przekształcamy w tensor Jacobiego

Musimy zwrócić uwagę na wymiary. Pochodną cząstkową f_i względem wiersza \mathbf{A} jest tensor o rozmiarze $1 \times 1 \times N$, ponieważ przeciwdziedzina f_i jest \mathbb{R} , a każdy wiersz \mathbf{A} ma rozmiar $1 \times N$. Pochodne cząstkowe z równania 5.91 układamy w stos i otrzymujemy interesujący nas gradient (równanie 5.87):

$$\frac{\partial f_i}{\partial \mathbf{A}} = \begin{bmatrix} \mathbf{0}^\top \\ \vdots \\ \mathbf{0}^\top \\ \mathbf{x}^\top \\ \mathbf{0}^\top \\ \vdots \\ \mathbf{0}^\top \end{bmatrix} \in \mathbb{R}^{1 \times (M \times N)}. \quad (5.92)$$

Przykład 5.13 (gradient macierzy względem macierzy)

Rozważmy macierz $\mathbf{R} \in \mathbb{R}^{M \times N}$ oraz funkcję $f: \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{N \times N}$ zdefiniowaną jako

$$\mathbf{f}(\mathbf{R}) = \mathbf{R}^\top \mathbf{R} =: \mathbf{K} \in \mathbb{R}^{N \times N}. \quad (5.93)$$

Chcemy znaleźć gradient $d\mathbf{K}/d\mathbf{R}$.

Aby rozwiązać ten trudny problem, najpierw zapiszemy to, co już wiemy:

Gradient ma wymiary

$$\frac{d\mathbf{K}}{d\mathbf{R}} \in \mathbb{R}^{(N \times N) \times (M \times N)}, \quad (5.94)$$

a więc jest tensorem. Ponadto dla $p, q = 1, \dots, N$

$$\frac{dK_{pq}}{d\mathbf{R}} \in \mathbb{R}^{1 \times M \times N}, \quad (5.95)$$

gdzie K_{pq} jest (p, q) -tym elementem $\mathbf{K} = \mathbf{f}(\mathbf{R})$. Jeżeli i -tą kolumnę \mathbf{R} oznaczymy przez \mathbf{r}_i , to każdy element \mathbf{K} będzie określony przez iloczyn skalarny dwóch kolumn macierzy \mathbf{R} :

$$K_{pq} = \mathbf{r}_p^\top \mathbf{r}_q = \sum_{m=1}^M R_{mp} R_{mq}. \quad (5.96)$$

Obliczymy teraz pochodną cząstkową $\frac{\partial K_{pq}}{\partial R_{ij}}$:

$$\frac{\partial K_{pq}}{\partial R_{ij}} = \sum_{m=1}^M \frac{\partial}{\partial R_{ij}} R_{mp} R_{mq} = \partial_{pqij}, \quad (5.97)$$

gdzie

$$\partial_{pqij} = \begin{cases} R_{iq} & \text{jeżeli } j = p, p \neq q \\ R_{ip} & \text{jeżeli } j = q, p \neq q \\ 2R_{iq} & \text{jeżeli } j = p, p = q \\ 0 & \text{w innych przypadkach} \end{cases} \quad (5.98)$$

Z równania 5.94 wiemy, że interesujący nas gradient ma wymiar $(N \times N) \times (M \times N)$. Elementami tego tensora są ∂_{pqij} z równania 5.98 ($p, q, j = 1, \dots, N$, a $i = 1, \dots, M$).

5.5. Tożsamości przydatne w obliczeniach gradientów

Poniżej prezentujemy kilka przydatnych tożsamości, które są często wykorzystywane w uczeniu maszynowym (Petersen i Pedersen, 2012). W poniższych równaniach $\text{tr}(\cdot)$ oznacza ślad macierzy (patrz definicja 4.4), $\det(\cdot)$ wyznacznik (patrz podrozdział 4.1), a $\mathbf{f}(\mathbf{X})^{-1}$ odwrotność funkcji $\mathbf{f}(\mathbf{X})$ (przy założeniu, że taka odwrotność istnieje).

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{f}(\mathbf{X})^\top = \left(\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right)^\top \quad (5.99)$$

$$\frac{\partial}{\partial \mathbf{X}} \text{tr}(\mathbf{f}(\mathbf{X})) = \text{tr} \left(\frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right) \quad (5.100)$$

$$\frac{\partial}{\partial \mathbf{X}} \det(\mathbf{f}(\mathbf{X})) = \det(\mathbf{f}(\mathbf{X})) \text{tr} \left(\mathbf{f}(\mathbf{X})^{-1} \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \right) \quad (5.101)$$

$$\frac{\partial}{\partial \mathbf{X}} \mathbf{f}(\mathbf{X})^{-1} = -\mathbf{f}(\mathbf{X})^{-1} \frac{\partial \mathbf{f}(\mathbf{X})}{\partial \mathbf{X}} \mathbf{f}(\mathbf{X})^{-1} \quad (5.102)$$

$$\frac{\partial \mathbf{a}^\top \mathbf{X}^{-1} \mathbf{b}}{\partial \mathbf{X}} = -(\mathbf{X}^{-1})^\top \mathbf{a} \mathbf{b}^\top (\mathbf{X}^{-1})^\top \quad (5.103)$$

$$\frac{\partial \mathbf{x}^\top \mathbf{a}}{\partial \mathbf{x}} = \mathbf{a}^\top \quad (5.104)$$

$$\frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}^\top \quad (5.105)$$

$$\frac{\partial \mathbf{a}^\top \mathbf{X} \mathbf{b}}{\partial \mathbf{X}} = \mathbf{a} \mathbf{b}^\top \quad (5.106)$$

$$\frac{\partial \mathbf{x}^\top \mathbf{B} \mathbf{x}}{\partial \mathbf{x}} = \mathbf{x}^\top (\mathbf{B} + \mathbf{B}^\top) \quad (5.107)$$

$$\text{Dla symetrycznych macierzy } \mathbf{W}: \frac{\partial}{\partial \mathbf{s}} (\mathbf{x} - \mathbf{A} \mathbf{s})^\top \mathbf{W} (\mathbf{x} - \mathbf{A} \mathbf{s}) = -2(\mathbf{x} - \mathbf{A} \mathbf{s})^\top \mathbf{W} \mathbf{A} \quad (5.108)$$

Uwaga. W tej książce zajmujemy się tylko śladami i transpozycjami macierzy. Wiesz jednak, że pochodne mogą być tensorami o większej liczbie wymiarów. W tym przypadku zwykły ślad i transpozycja nie są zdefiniowane. Ślad tensora o wymiarach $D \times D \times E \times F$ byłby macierzą o wymiarach $E \times F$. Jest to szczególny przypadek kontrakcji/zwężania tensorów. W podobny sposób, gdy „transponujemy” tensor, mamy na myśli zamianę dwóch jego pierwszych wymiarów. Jeżeli nie zamienimy macierzy na wektory, tak jak omówiono to w podrozdziale 5.4, to w przypadku obliczeń pochodnych funkcji wielu zmiennych $f(\cdot)$ względem macierzy w równaniach od 5.99 do 5.102 pojawią się obliczenia na tensorach.

5.6. Propagacja wsteczna i różniczkowanie automatyczne

W wielu rozwiązaniach z obszaru uczenia maszynowego do znalezienia odpowiednich parametrów modelu wykorzystujemy metodę gradientu prostego (podrozdział 7.1), która pozwala obliczyć gradient funkcji celu w odniesieniu do parametrów modelu. Gradient ten możemy wyznaczyć za pomocą reguł rachunku różniczkowego oraz reguły łańcuchowej (patrz punkt 5.2.2). Przykład takiego gradientu pokazaliśmy w podrozdziale 5.3, w którym przyjrzeliliśmy się gradientowi minimalnokwadratowej straty w odniesieniu do parametrów modelu regresji liniowej.

Rozważmy funkcję

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2)). \quad (5.109)$$

Skorzystamy z faktu, że różniczkowanie jest operacją liniową. Po zastosowaniu reguły łańcuchowej otrzymujemy następujący gradient:

$$\begin{aligned} \frac{df}{dx} = & \frac{2x + 2x\exp(x^2)}{2\sqrt{x^2 + \exp(x^2)}} - \sin(x^2 + \exp(x^2))(2x + 2x\exp(x^2)) \\ & - 2x \left(\frac{1}{2\sqrt{x^2 + \exp(x^2)}} - \sin(x^2 + \exp(x^2)) \right) (1 + \exp(x^2)). \end{aligned} \quad (5.110)$$

Zapisywanie gradientów w powyższej jawnej postaci nie jest zbyt praktyczne i często prowadzi do bardzo długich wyrażeń opisujących pochodne. Z praktycznego punktu widzenia oznacza to, że przy nieuważnej implementacji koszt obliczenia gradientu może być znacznie wyższy niż koszt ewaluacji różniczkowanej funkcji. Stosowanie takiego podejścia może więc niepotrzebnie spowolnić nasze obliczenia. W przypadku uczenia głębokich sieci neuronowych wydajnym sposobem obliczania gradientu funkcji błędu w odniesieniu do parametrów modelu jest algorytm propagacji wstecznej (Kelley, 1960; Bryson, 1961; Dreyfusa, 1962; Rumelhart i in., 1986).

5.6.1. Gradienty w uczeniu głębokim

Obszarem, w którym silnie eksploatuje się regułę łańcuchową, jest uczenie głębokie, w którym funkcja y jest wielopoziomowym złożeniem funkcji

$$y = (f_K \circ f_{K-1} \circ \dots \circ f_1)(x) = f_K(f_{K-1}(\dots(f_1(x))\dots)), \quad (5.111)$$

Dobre omówienie propagacji wstecznej i reguły łańcuchowej znajdziesz na blogu Tima Viery pod adresem <https://tinyurl.com/yefm2yrw>.

Aby uprościć notację, omawiamy przypadek, w którym w każdej warstwie mamy identyczne funkcje aktywacji.

gdzie \mathbf{x} to dane wejściowe (np. obrazy), a \mathbf{y} obserwacje (np. etykiety klas). Każda funkcja $f_i, i = 1, \dots, K$ posiada własne parametry. W wielowarstwowej sieci neuronowej z i -tą warstwą sieci związana jest funkcja $f_i(\mathbf{x}_{i-1}) = \sigma(\mathbf{A}_{i-1}\mathbf{f}_{i-1} + \mathbf{b}_{i-1})$, w której \mathbf{x}_{i-1} jest wyjściem z $i - 1$ warstwy, a σ jest funkcją aktywacji, na przykład funkcją sigmoidalną $\frac{1}{1+e^{-x}}$, tangensem hiperbolicznym lub funkcją ReLU. Proces uczenia takich sieci wymaga obliczenia gradientu funkcji straty L względem wszystkich parametrów $\mathbf{A}_j, \mathbf{b}_j (j = 1, \dots, K)$ modelu. Konieczne jest również znalezienie gradientu L względem danych wejściowych każdej warstwy. Rozważmy przykład, w którym \mathbf{x} to dane wejściowe, a \mathbf{y} to obserwacje. Struktura sieci jest określona przez

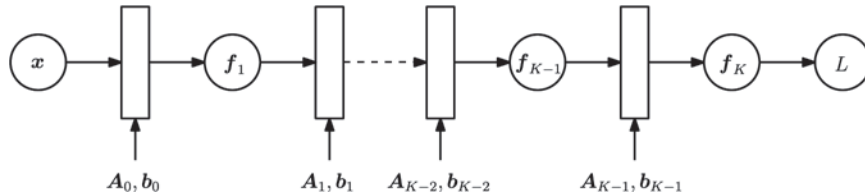
$$\mathbf{f}_0 := \mathbf{x} \tag{5.112}$$

$$\mathbf{f}_i := \sigma_i(\mathbf{A}_{i-1}\mathbf{f}_{i-1} + \mathbf{b}_{i-1}), i = 1, \dots, K. \tag{5.113}$$

Siec tę pokazano też na rysunku 5.8. Chcemy znaleźć parametry $\mathbf{A}_j, \mathbf{b}_j (j = 0, \dots, K - 1)$ minimalizujące kwadratową funkcję straty

$$L(\boldsymbol{\theta}) = \|\mathbf{y} - \mathbf{f}_K(\boldsymbol{\theta}, \mathbf{x})\|^2, \tag{5.114}$$

gdzie $\boldsymbol{\theta} = \{\mathbf{A}_0, \mathbf{b}_0, \dots, \mathbf{A}_{K-1}, \mathbf{b}_{K-1}\}$.



RYСУNEK 5.8. Przejście do przodu przez wielowarstwową sieć neuronową w celu obliczenia straty L w funkcji wejść \mathbf{x} i parametrów $\mathbf{A}_i, \mathbf{b}_i$

Do wyznaczenia gradientów względem zbioru parametrów $\boldsymbol{\theta}$ potrzebne nam będą pochodne cząstkowe L względem parametrów $\boldsymbol{\theta}_j = \{\mathbf{A}_j, \mathbf{b}_j\}$ każdej warstwy $j = 0, \dots, K - 1$. Reguła łańcuchowa pozwala nam zapisać pochodne cząstkowe jako

$$\frac{\partial L}{\partial \boldsymbol{\theta}_{K-1}} = \frac{\partial L}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \boldsymbol{\theta}_{K-1}} \tag{5.115}$$

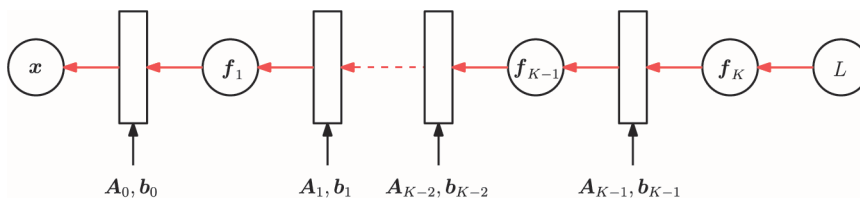
$$\frac{\partial L}{\partial \boldsymbol{\theta}_{K-2}} = \frac{\partial L}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \mathbf{f}_{K-1}} \frac{\partial \mathbf{f}_{K-1}}{\partial \boldsymbol{\theta}_{K-2}} \tag{5.116}$$

$$\frac{\partial L}{\partial \boldsymbol{\theta}_{K-3}} = \frac{\partial L}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \mathbf{f}_{K-1}} \frac{\partial \mathbf{f}_{K-1}}{\partial \mathbf{f}_{K-2}} \frac{\partial \mathbf{f}_{K-2}}{\partial \boldsymbol{\theta}_{K-3}} \tag{5.117}$$

$$\frac{\partial L}{\partial \boldsymbol{\theta}_i} = \frac{\partial L}{\partial \mathbf{f}_K} \frac{\partial \mathbf{f}_K}{\partial \mathbf{f}_{K-1}} \dots \frac{\partial \mathbf{f}_{i+2}}{\partial \mathbf{f}_{i+1}} \frac{\partial \mathbf{f}_{i+1}}{\partial \boldsymbol{\theta}_i} \tag{5.118}$$

Bardziej szczegółowe omówienie gradientów w sieciach neuronowych znajdziesz w notatkach Justina Domkego pod adresem <https://tinyurl.com/yalcxgtv>.

Składniki oznaczone kolorem pomarańczowym to pochodne cząstkowe wyjścia warstwy względem jej wejścia. Natomiast niebieskie składniki to pochodne cząstkowe wyjścia warstwy względem jej parametrów. Jeżeli obliczyliśmy już pochodne cząstkowe $\partial L / \partial \theta_{i+1}$, to większość obliczonych wartości możemy wykorzystać ponownie do wyznaczenia $\partial L / \partial \theta_i$. Dodatkowe składniki, które musimy obliczyć, zaznaczono ramkami. Na rysunku 5.9 pokazano, w jaki sposób gradienty są przekazywane wstecz sieci.



RYSUNEK 5.9. Przejście wstecz przez wielowarstwową sieć neuronową w celu znalezienia gradientu funkcji straty

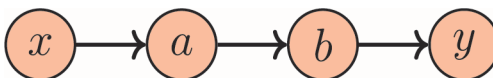
5.6.2. Różniczkowanie automatyczne

Różniczkowanie automatyczne różni się od różniczkowania symbolicznego i numerycznych metod aproksymacji gradientu np. przy użyciu różnic skończonych.

Okazuje się, że propagacja wsteczna jest szczególnym przypadkiem ogólnej techniki analizy numerycznej znanej jako **różniczkowanie automatyczne**. Różniczkowanie automatyczne możesz traktować jako zbiór metod służących do numerycznego (w odróżnieniu od symbolicznego) wyznaczania gradientów z dokładnością maszynową. Techniki różniczkowania automatycznego wykorzystują zmienne pośrednie i regułę łańcuchową. Różniczkowanie automatyczne wykorzystuje szereg elementarnych operacji arytmetycznych, takich jak na przykład dodawanie i mnożenie, oraz funkcje elementarne, takie jak sinus, cosinus, funkcja wykładnicza i logarytm. Zastosowanie reguły łańcuchowej pozwala wyznaczać gradient w sposób automatyczny. W ogólności różniczkowanie automatyczne dotyczy programów komputerowych i występuje przede wszystkim w dwóch wariantach: w przód i wstecz. Świetny przegląd różniczkowania automatycznego w kontekście uczenia maszynowego znajdziesz w pracy Baydina i in. (2018).

Na rysunku 5.10 pokazano prosty graf przedstawiający przepływ danych z wejść x do wyjść y przez zmienne pośrednie a i b . Gdybyśmy chcieli wyznaczyć pochodną dy/dx , to korzystając z reguły łańcuchowej, moglibyśmy zapisać, że:

$$\frac{dy}{dx} = \frac{dy}{db} \frac{db}{da} \frac{da}{dx} \quad (5.119)$$



RYSUNEK 5.10. Prostý graf ilustrujący przepływ danych od x do y przez zmienne pośrednie a i b

Z intuicyjnego punktu widzenia warianty w przód i wstecz różnią się kolejnością wykonywania mnożenia. Ponieważ mnożenie macierzy jest łączne, możemy wybierać pomiędzy

W ogólnym przypadku pracujemy ze strukturami, którymi mogą być gradienty oraz macierze i tensory Jacobiego.

$$\frac{dy}{dx} = \left(\frac{dy}{db} \frac{db}{da} \right) \frac{da}{dx} \quad (5.120)$$

i

$$\frac{dy}{dx} = \frac{dy}{db} \left(\frac{db}{da} \frac{da}{dx} \right). \quad (5.121)$$

Równanie 5.120 to przykład metody różniczkowania automatycznego wstecz, ponieważ gradienty są w nim przekazywane od końca do początku grafu, tj. odwrotnie do przepływu danych. Równanie 5.121 jest przykładem metody w przód, w której informacje o gradientach są przekazywane zgodnie z kierunkiem danych, tj. od lewej do prawej.

Poniżej skupimy się na różniczkowaniu wstecz, którego przykładem jest algorytm propagacji wstecznej. W przypadku sieci neuronowych, w których rozmiar wejścia jest często znacznie większy niż rozmiar wyjścia, różniczkowanie wstecz jest bardziej opłacalne obliczeniowo niż różniczkowanie metodą w przód. Zacznijmy od przykładu.

Przykład 5.14

Rozważmy funkcję

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2)) \quad (5.122)$$

z równania 5.109. Podczas implementacji tej funkcji na komputerze moglibyśmy zaoszczędzić trochę czasu, stosując tzw. **zmiennie pośrednie**:

$$a = x^2 \quad (5.123)$$

$$b = \exp(a) \quad (5.124)$$

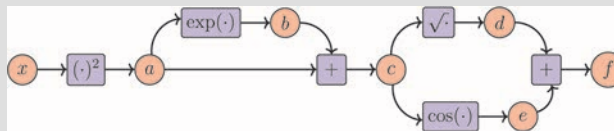
$$c = a + b \quad (5.125)$$

$$d = \sqrt{c} \quad (5.126)$$

$$e = \cos(c) \quad (5.127)$$

$$f = d + e. \quad (5.128)$$

W ten sam sposób rozumiemy również podczas stosowania reguły łańcuchowej. Zauważ, że powyższy zestaw równań wymaga wykonania mniejszej liczby operacji niż bezpośrednia implementacja funkcji $f(x)$ w postaci pokazanej w równaniu 5.109. Graf obliczeniowy tej funkcji, pokazujący przepływ danych i obliczenia wymagane do wyznaczenia wartości funkcji f , pokazano na rysunku 5.11.



RYСУNEK 5.11. Graf obliczeniowy z danymi wejściowymi x , wartościami funkcji f i zmiennymi pośrednimi a, b, c, d, e

Zbiór równań, które zawierają zmienne pośrednie, możemy traktować jako graf obliczeniowy. Ten sposób przedstawiania obliczeń jest powszechnie stosowany w implementacjach bibliotek przeznaczonych do tworzenia sieci neuronowych. Pochodne zmiennych pośrednich względem odpowiadających im danych wejściowych możemy obliczyć za pomocą definicji pochodnych funkcji elementarnych. Otrzymujemy następujące pochodne:

$$\frac{\partial a}{\partial x} = 2x \quad (5.129)$$

$$\frac{\partial b}{\partial a} = \exp(a) \quad (5.130)$$

$$\frac{\partial c}{\partial a} = 1 = \frac{\partial c}{\partial b} \quad (5.131)$$

$$\frac{\partial d}{\partial c} = \frac{1}{2\sqrt{c}} \quad (5.132)$$

$$\frac{\partial e}{\partial c} = -\sin(c) \quad (5.133)$$

$$\frac{\partial f}{\partial d} = 1 = \frac{\partial f}{\partial e} \quad (5.134)$$

Analiza grafu z rysunku 5.11 pozwala zauważyć, że pochodną $\partial f / \partial x$ możemy obliczyć metodą wstecz, rozpoczynając obliczenia w wierzchołku wyjściowym:

$$\frac{\partial f}{\partial c} = \frac{\partial f}{\partial d} \frac{\partial d}{\partial c} + \frac{\partial f}{\partial e} \frac{\partial e}{\partial c} \quad (5.135)$$

$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial c} \frac{\partial c}{\partial b} \quad (5.136)$$

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial b} \frac{\partial b}{\partial a} + \frac{\partial f}{\partial c} \frac{\partial c}{\partial a} \quad (5.137)$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial x} \quad (5.138)$$

Zauważ, że do wyznaczenia $\partial f / \partial x$ zastosowaliśmy regułę łańcuchową. Po podstawieniu pochodnych funkcji elementarnych do powyższych równań otrzymujemy

$$\frac{\partial f}{\partial c} = 1 \cdot \frac{1}{2\sqrt{c}} + 1 \cdot (-\sin(c)) \quad (5.139)$$

$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial c} \cdot 1 \quad (5.140)$$

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial b} \exp(a) + \frac{\partial f}{\partial c} \cdot 1 \quad (5.141)$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial a} \cdot 2x. \quad (5.142)$$

Jeżeli każdą z powyższych pochodnych potraktujemy jako zmienną, to proces wyznaczania pochodnej będzie miał podobną złożoność obliczeniową co proces obliczania wartości funkcji. Jest to sprzeczne z naszą intuicją, ponieważ wyrażenie matematyczne opisujące pochodną $\frac{\partial f}{\partial x}$ (równanie 5.110) jest znacznie bardziej skomplikowane niż wyrażenie opisujące funkcję $f(x)$ (równanie 5.109).

Różniczkowanie automatyczne jest formalizacją obliczeń z przykładu 5.14. Niech x_1, \dots, x_d będą zmiennymi wejściowymi, x_{d+1}, \dots, x_{D-1} zmiennymi pośrednimi, a x_D zmienną wyjściową. Graf obliczeniowy możemy wtedy zapisać jako

$$x_i = g_i(x_{\text{Pa}(x_i)}) \text{ gdzie } i = d + 1, \dots, D. \quad (5.143)$$

W powyższym równaniu $g_i(\cdot)$ to funkcje elementarne, a $x_{\text{Pa}(x_i)}$ to węzły rodzice zmiennej x_i w grafie. Pochodną tak zdefiniowanej funkcji możemy obliczyć za pomocą reguły łańcuchowej. Przypomnijmy, że z definicji $f = x_D$. Stąd

$$\frac{\partial f}{\partial x_D} = 1. \quad (5.144)$$

Dla pozostałych zmiennych x_i stosujemy regułę łańcuchową

$$\frac{\partial f}{\partial x_i} = \sum_{x_j: x_i \in \text{Pa}(x_j)} \frac{\partial f}{\partial x_j} \frac{\partial x_j}{\partial x_i} = \sum_{x_j: x_i \in \text{Pa}(x_j)} \frac{\partial f}{\partial x_j} \frac{\partial g_j}{\partial x_i}, \quad (5.145)$$

gdzie $\text{Pa}(x_j)$ jest zbiorem węzłów rodziców wierzchołka x_j w grafie. Równanie 5.143 to przykład propagacji w przód, a równanie 5.145 reprezentuje propagację wsteczną gradientu w grafie obliczeniowym. W uczeniu sieci neuronowych propagacji wstecznej podlega błąd predykcji w odniesieniu do etykiety.

Różniczkowanie automatyczne działa zawsze, gdy różniczkujemy funkcję dającą się przedstawić w postaci grafu obliczeniowego zawierającego różniczkowalne funkcje elementarne. Różniczkowana funkcja nie musi być nawet funkcją w sensie matematycznym, a jedynie programem komputerowym. Niestety, nie wszystkie programy komputerowe mogą być różniczkowane automatycznie. Nie jest to możliwe np. wtedy, gdy nie możemy wyznaczyć pochodnych funkcji elementarnych. Struktury występujące w kodzie, takie jak pętle for i instrukcje if, również wymagają większej uwagi.

5.7. Pochodne wyższych rzędów

Do tej pory omówiliśmy jedynie gradienty, czyli pochodne pierwszego rzędu. Czasami interesują nas też pochodne wyższych rzędów. Przydają się one np. w metodzie optymalizacji Newtona, która wymaga pochodnych drugiego rzędu (Nocedal i Wright, 2006). W punkcie 5.1.1 omówiliśmy szereg Taylora, pozwalający przybliżyć funkcje za pomocą wielomianów. Przybliżenie takie możemy również wyznaczyć dla funkcji wielu zmiennych. W dalszej części tego rozdziału pokażemy, jak to zrobić. Zaczniemy jednak od wprowadzenia notacji.

Różniczkowanie automatyczne metodą wstecz wymaga stworzenia drzewa składniowego.

Rozważmy funkcję $f: \mathbb{R}^2 \rightarrow \mathbb{R}$ dwóch zmiennych x, y . W przypadku pochodnych cząstkowych wyższego rzędu (oraz gradientów) stosujemy następujący zapis:

- $\frac{\partial^2 f}{\partial x^2}$ to druga pochodna cząstkowa f względem x .
- $\frac{\partial^n f}{\partial x^n}$ to pochodna cząstkowa n -tego rzędu f względem x .
- $\frac{\partial^2 f}{\partial y \partial x} = \frac{\partial}{\partial y} \left(\frac{\partial f}{\partial x} \right)$ to pochodna cząstkowa otrzymana w wyniku różniczkowania względem x , a następnie względem y .
- $\frac{\partial^2 f}{\partial x \partial y}$ to pochodna cząstkowa otrzymana w wyniku różniczkowania względem y , a następnie względem x .

Hesjan (macierz Hessego) to zbiór wszystkich pochodnych cząstkowych drugiego rzędu. Jeśli $f(x, y)$ jest funkcją dwukrotnie różniczkowalną (a jej druga pochodna jest ciągła), to

$$\frac{\partial^2 f}{\partial x \partial y} = \frac{\partial^2 f}{\partial y \partial x} \quad (5.146)$$

tj. kolejność różniczkowania nie ma znaczenia. Hesjanem tej funkcji jest symetryczna macierz

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial x \partial y} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix}. \quad (5.147)$$

Hesjan oznaczamy również symbolem $\nabla_{x,y}^2 f(x, y)$. Dla $\mathbf{x} \in \mathbb{R}^n$ i $f: \mathbb{R}^n \rightarrow \mathbb{R}$ hesjan jest macierzą o wymiarach $n \times n$. Hesjan mierzy lokalną krzywiznę funkcji wokół punktu (x, y) .

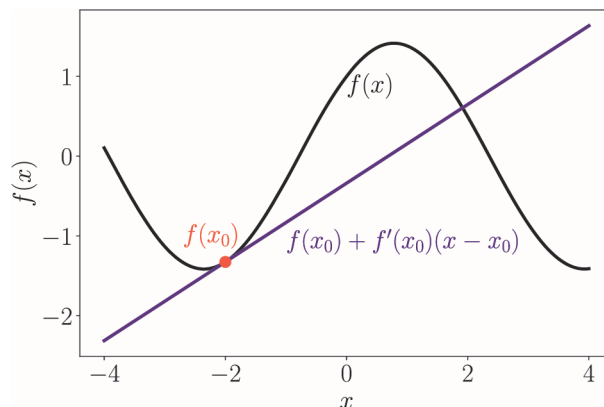
Uwaga (hesjan pola wektorowego). Jeśli $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ jest polem wektorowym, to hesjan jest tensorem o wymiarach $m \times n \times n$.

5.8. Linearyzacja i wielowymiarowe szeregi Taylora

Gradient ∇f funkcji f często wykorzystuje się do lokalnego liniowego przybliżania funkcji f w otoczeniu punktu \mathbf{x}_0 :

$$f(\mathbf{x}) \approx f(\mathbf{x}_0) + (\nabla_{\mathbf{x}} f)(\mathbf{x}_0)(\mathbf{x} - \mathbf{x}_0). \quad (5.148)$$

$(\nabla_{\mathbf{x}} f)(\mathbf{x}_0)$ jest gradientem f względem \mathbf{x} obliczonym w \mathbf{x}_0 . Liniową aproksymację funkcji f w otoczeniu \mathbf{x}_0 pokazano na rysunku 5.12. Oryginalna funkcja jest przybliżana za pomocą prostej. To przybliżenie jest lokalnie dokładne, ale im dalej od \mathbf{x}_0 , tym prosta gorzej przybliża funkcję f . Równanie 5.148 jest szczególnym przypadkiem rozwinięcia f w wielowymiarowym szeregu Taylora w otoczeniu punktu \mathbf{x}_0 , w którym korzystamy jedynie z dwóch pierwszych wyrazów. Poniżej omówimy bardziej ogólny przypadek tego szeregu, który pozwoli uzyskać lepsze przybliżenia.



RYSUNEK 5.12. Liniowa aproksymacja funkcji. Oryginalna funkcja f jest linearyzowana w otoczeniu $x_0 = -2$ za pomocą rozwinięcia w szereg Taylora pierwszego rzędu

Definicja 5.7 (szereg Taylora funkcji wielu zmiennych). Rozważamy funkcję

$$f: \mathbb{R}^D \rightarrow \mathbb{R} \quad (5.149)$$

$$\mathbf{x} \mapsto f(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^D, \quad (5.150)$$

która jest gładka w punkcie \mathbf{x}_0 . Jeżeli zdefiniujemy wektor różnic jako $\boldsymbol{\delta} := \mathbf{x} - \mathbf{x}_0$, to **szereg Taylora funkcji wielu zmiennych** f w punkcie \mathbf{x}_0 będzie zdefiniowany jako

$$f(\mathbf{x}) = \sum_{k=0}^{\infty} \frac{D_{\mathbf{x}}^k f(\mathbf{x}_0)}{k!} \boldsymbol{\delta}^k, \quad (5.151)$$

gdzie $D_{\mathbf{x}}^k f(\mathbf{x}_0)$ jest k -tą (zupelną) pochodną f względem \mathbf{x} obliczoną w \mathbf{x}_0 .

Definicja 5.8 (wielomian Taylora). **Wielomian Taylora** n -tego stopnia funkcji f w \mathbf{x}_0 zawiera pierwsze $n + 1$ elementów szeregu z równania 5.151 i jest zdefiniowany jako

$$T_n(\mathbf{x}) = \sum_{k=0}^n \frac{D_{\mathbf{x}}^k f(\mathbf{x}_0)}{k!} \boldsymbol{\delta}^k. \quad (5.152)$$

W równaniach 5.151 i 5.152 zastosowaliśmy nieco nieprecyzyjny zapis z użyciem $\boldsymbol{\delta}^k$, które jest niezdefiniowane dla wektorów $\mathbf{x} \in \mathbb{R}^D$ takich, że $D > 1$ i $k > 1$. Zauważ, że $D_{\mathbf{x}}^k f$ i $\boldsymbol{\delta}^k$ są tensorami k -tego rzędu, tj. k -wymiarowymi tablicami. Tensor k -tego rzędu $\boldsymbol{\delta}^k \in \mathbb{R}^{\overbrace{D \times D \times \dots \times D}^{k \text{ razy}}}$ to wynik k -krotnego iloczynu zewnętrznego (oznaczanego symbolem \otimes) wektora $\boldsymbol{\delta} \in \mathbb{R}^D$. Na przykład

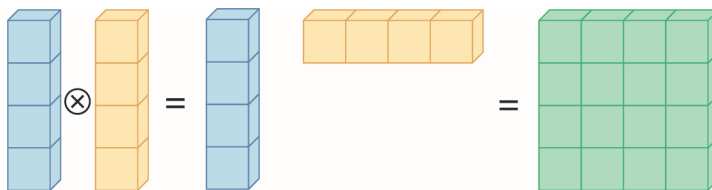
$$\boldsymbol{\delta}^2 := \boldsymbol{\delta} \otimes \boldsymbol{\delta} = \boldsymbol{\delta} \boldsymbol{\delta}^T, \quad \boldsymbol{\delta}^2[i, j] = \delta[i] \delta[j] \quad (5.153)$$

$$\boldsymbol{\delta}^3 := \boldsymbol{\delta} \otimes \boldsymbol{\delta} \otimes \boldsymbol{\delta}, \quad \boldsymbol{\delta}^3[i, j, k] = \delta[i] \delta[j] \delta[k]. \quad (5.154)$$

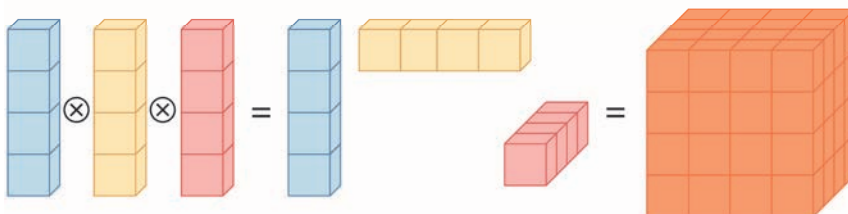
Na rysunku 5.13 pokazano dwa takie iloczyny zewnętrzne. W ogólnym przypadku w szeregu Taylora pojawiają się następujące składniki

Wektor można zaimplementować jako tablicę jednowymiarową, a macierz w postaci tablicy dwuwymiarowej.

$$D_x^k f(\mathbf{x}_0) \boldsymbol{\delta}^k = \sum_{i_1=1}^D \cdots \sum_{i_k=1}^D D_x^k f(\mathbf{x}_0)[i_1, \dots, i_k] \delta[i_1] \cdots \delta[i_k]. \quad (5.155)$$



(a) Iloczyn zewnętrzny wektora $\boldsymbol{\delta} \in \mathbb{R}^4$ daje $\boldsymbol{\delta}^2 := \boldsymbol{\delta} \otimes \boldsymbol{\delta} = \boldsymbol{\delta} \boldsymbol{\delta}^T \in \mathbb{R}^{4 \times 4}$, czyli macierz



(b) Iloczyn zewnętrzny $\boldsymbol{\delta}^3 := \boldsymbol{\delta} \otimes \boldsymbol{\delta} \otimes \boldsymbol{\delta} \in \mathbb{R}^{4 \times 4 \times 4}$ daje tensor trzeciego rzędu (macierz trójwymiarową), tj. macierz z trzema indeksami (osiąmi)

RYSUNEK 5.13. Wizualizacja iloczynów zewnętrznych. Każdy iloczyn zewnętrzny wektorów zwiększa o jeden liczbę wymiarów wynikowej tablicy. (a) Iloczyn zewnętrzny dwóch wektorów daje macierz; (b) iloczyn zewnętrzny trzech wektorów daje tensor trzeciego rzędu

Element $D_x^k f(\mathbf{x}_0) \boldsymbol{\delta}^k$ zawiera wielomiany k -tego rzędu.

Teraz, gdy zdefiniowaliśmy już szereg Taylora dla pól wektorowych, wyznaczmy pierwsze wyrazy $D_x^k f(\mathbf{x}_0) \boldsymbol{\delta}^k$ rozwinięcia w szereg Taylora dla $k = 0, \dots, 3$ i $\boldsymbol{\delta} := \mathbf{x} - \mathbf{x}_0$:

$$k = 0: D_x^0 f(\mathbf{x}_0) \boldsymbol{\delta}^0 = f(\mathbf{x}_0) \in \mathbb{R} \quad (5.156)$$

$$k = 1: D_x^1 f(\mathbf{x}_0) \boldsymbol{\delta}^1 = \underbrace{\nabla_x f(\mathbf{x}_0)}_{1 \times D} \underbrace{\boldsymbol{\delta}}_{D \times 1} = \sum_{i=1}^D \nabla_x f(\mathbf{x}_0)[i] \delta[i] \in \mathbb{R} \quad (5.157)$$

$$k = 2: D_x^2 f(\mathbf{x}_0) \boldsymbol{\delta}^2 = \text{tr} \left(\underbrace{\mathbf{H}(\mathbf{x}_0)}_{D \times D} \underbrace{\boldsymbol{\delta}}_{D \times 1} \underbrace{\boldsymbol{\delta}^T}_{1 \times D} \right) = \boldsymbol{\delta}^T \mathbf{H}(\mathbf{x}_0) \boldsymbol{\delta} \quad (5.158)$$

$$= \sum_{i=1}^D \sum_{j=1}^D H[i, j] \delta[i] \delta[j] \in \mathbb{R} \quad (5.159)$$

$$k = 3: D_x^3 f(\mathbf{x}_0) \boldsymbol{\delta}^3 = \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D D_x^3 f(\mathbf{x}_0)[i, j, k] \delta[i] \delta[j] \delta[k] \in \mathbb{R}. \quad (5.160)$$

W powyższych równaniach $\mathbf{H}(\mathbf{x}_0)$ jest hesjanem f w punkcie \mathbf{x}_0 .

Przykład 5.15 (rozwińnięcie funkcji dwóch zmiennych w szereg Taylora)

Rozważmy funkcję

$$f(x, y) = x^2 + 2xy + y^3. \quad (5.161)$$

Chcemy znaleźć jej rozwinięcie w szereg Taylora w punkcie $(x_0, y_0) = (1, 2)$. Zanim zaczniemy, omówmy, czego możemy się spodziewać. Funkcja z równania 5.161 jest wielomianem trzeciego stopnia. Szukamy rozwinięcia w szereg Taylora, które samo w sobie jest kombinacją liniową wielomianów. Spodziewamy się więc, że w przypadku wielomianu trzeciego stopnia rozwinięcie nie będzie zawierało wyrazów czwartego lub wyższego stopnia. Oznacza to, że pierwsze cztery wyrazy z sumy z równania 5.151 powinny wystarczyć do uzyskania dokładnej reprezentacji funkcji z równania 5.161.

Proces wyznaczania rozwinięcia zaczynamy od wyrazu stałego i pochodnych pierwszego rzędu, które dane są wzorem

$$f(1, 2) = 13 \quad (5.162)$$

$$\frac{\partial f}{\partial x} = 2x + 2y \Rightarrow \frac{\partial f}{\partial x}(1, 2) = 6 \quad (5.163)$$

$$\frac{\partial f}{\partial y} = 2x + 3y^2 \Rightarrow \frac{\partial f}{\partial y}(1, 2) = 14. \quad (5.164)$$

Stąd

$$D_{x,y}^1 f(1, 2) = \nabla_{x,y} f(1, 2) = \left[\frac{\partial f}{\partial x}(1, 2) \quad \frac{\partial f}{\partial y}(1, 2) \right] = [6 \quad 14] \in \mathbb{R}^{1 \times 2}, \quad (5.165)$$

co daje

$$\frac{D_{x,y}^1 f(1, 2)}{1!} \boldsymbol{\delta} = [6 \quad 14] \begin{bmatrix} x - 1 \\ y - 2 \end{bmatrix} = 6(x - 1) + 14(y - 2). \quad (5.166)$$

Zauważ, że $D_{x,y}^1 f(1, 2) \boldsymbol{\delta}$ zawiera tylko składniki liniowe, tj. wielomiany pierwszego rzędu.

Pochodne cząstkowe drugiego rzędu to

$$\frac{\partial^2 f}{\partial x^2} = 2 \Rightarrow \frac{\partial^2 f}{\partial x^2}(1, 2) = 2 \quad (5.167)$$

$$\frac{\partial^2 f}{\partial y^2} = 6y \Rightarrow \frac{\partial^2 f}{\partial y^2}(1, 2) = 12 \quad (5.168)$$

$$\frac{\partial^2 f}{\partial y \partial x} = 2 \Rightarrow \frac{\partial^2 f}{\partial y \partial x}(1, 2) = 2 \quad (5.169)$$

$$\frac{\partial^2 f}{\partial x \partial y} = 2 \Rightarrow \frac{\partial^2 f}{\partial x \partial y}(1, 2) = 2. \quad (5.170)$$

Po zapisaniu ich w macierzy otrzymujemy hesjan

$$\mathbf{H} = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix} = \begin{bmatrix} 2 & 2 \\ 2 & 6y \end{bmatrix}. \quad (5.171)$$

A zatem

$$\mathbf{H}(1,2) = \begin{bmatrix} 2 & 2 \\ 2 & 12 \end{bmatrix} \in \mathbb{R}^{2 \times 2}. \quad (5.172)$$

A więc kolejny wyraz rozwinięcia w szereg Taylora to

$$\frac{D_{x,y}^2 f(1,2)}{2!} \boldsymbol{\delta}^2 = \frac{1}{2} \boldsymbol{\delta}^\top \mathbf{H}(1,2) \boldsymbol{\delta} \quad (5.173a)$$

$$= \frac{1}{2} [x-1 \quad y-2] \begin{bmatrix} 2 & 2 \\ 2 & 12 \end{bmatrix} \begin{bmatrix} x-1 \\ y-2 \end{bmatrix} \quad (5.173b)$$

$$= (x-1)^2 + 2(x-1)(y-2) + 6(y-2)^2. \quad (5.173c)$$

W tym przypadku $D_{x,y}^2 f(1,2) \boldsymbol{\delta}^2$ zawiera tylko wyrazy kwadratowe, tj. wielomiany drugiego rzędu.

Następnie wyznaczamy pochodne trzeciego rzędu

$$D_{x,y}^3 f = \begin{bmatrix} \frac{\partial \mathbf{H}}{\partial x} & \frac{\partial \mathbf{H}}{\partial y} \end{bmatrix} \in \mathbb{R}^{2 \times 2 \times 2} \quad (5.174)$$

$$D_{x,y}^3 f[:, :, 1] = \frac{\partial \mathbf{H}}{\partial x} = \begin{bmatrix} \frac{\partial^3 f}{\partial x^3} & \frac{\partial^3 f}{\partial x^2 \partial y} \\ \frac{\partial^3 f}{\partial x \partial y \partial x} & \frac{\partial^3 f}{\partial x \partial y^2} \end{bmatrix} \quad (5.175)$$

$$D_{x,y}^3 f[:, :, 2] = \frac{\partial \mathbf{H}}{\partial y} = \begin{bmatrix} \frac{\partial^3 f}{\partial y \partial x^2} & \frac{\partial^3 f}{\partial y \partial x \partial y} \\ \frac{\partial^3 f}{\partial y^2 \partial x} & \frac{\partial^3 f}{\partial y^3} \end{bmatrix}. \quad (5.176)$$

Ponieważ większość pochodnych cząstkowych drugiego rzędu w hesjanie z równania 5.171 to stałe, jedyną niezerową pochodną cząstkową trzeciego rzędu jest

$$\frac{\partial^3 f}{\partial y^3} = 6 \implies \frac{\partial^3 f}{\partial y^3}(1,2) = 6. \quad (5.177)$$

Pochodne wyższego rzędu i pochodne mieszane trzeciego rzędu (np. $\frac{\partial^3 f}{\partial x^2 \partial y}$) zanikają:

$$D_{x,y}^3 f[:, :, 1] = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad D_{x,y}^3 f[:, :, 2] = \begin{bmatrix} 0 & 0 \\ 0 & 6 \end{bmatrix}. \quad (5.178)$$

A zatem

$$\frac{D_{x,y}^3 f(1,2)}{3!} \delta^3 = (y-2)^3. \quad (5.179)$$

Składnik ten zawiera w sobie wszystkie sześciennne wyrazy szeregu. (Dokładne) rozwinięcie f w szereg Taylora w punkcie $(x_0, y_0) = (1, 2)$ jest zatem równe

$$f(x) = f(1,2) + D_{x,y}^1 f(1,2) \delta + \frac{D_{x,y}^2 f(1,2)}{2!} \delta^2 + \frac{D_{x,y}^3 f(1,2)}{3!} \delta^3 \quad (5.180a)$$

$$\begin{aligned} &= f(1,2) + \frac{\partial f(1,2)}{\partial x} (x-1) + \frac{\partial f(1,2)}{\partial y} (y-2) \\ &+ \frac{1}{2!} \left(\frac{\partial^2 f(1,2)}{\partial x^2} (x-1)^2 + \frac{\partial^2 f(1,2)}{\partial y^2} (y-2)^2 \right) \end{aligned} \quad (5.180b)$$

$$\begin{aligned} &+ 2 \frac{\partial^2 f(1,2)}{\partial x \partial y} (x-1)(y-2) + \frac{1}{6} \frac{\partial^3 f(1,2)}{\partial y^3} (y-2)^3 \\ &= 13 + 6(x-1) + 14(y-2) \\ &+ (x-1)^2 + 6(y-2)^2 + 2(x-1)(y-2) + (y-2)^3 \end{aligned} \quad (5.180c)$$

W tym przypadku otrzymaliśmy dokładne rozwinięcie wielomianu z równania 5.161 w szereg Taylora, tzn. wielomian z równania 5.180c jest identyczny z pierwotnym wielomianem z równania 5.161. W tym konkretnym przykładzie wynik ten nie jest zaskakujący, ponieważ pierwotną funkcją był wielomian trzeciego rzędu, który wyraziliśmy poprzez kombinację liniową stałych oraz wielomianów pierwszego, drugiego i trzeciego stopnia (równanie 5.180c).

5.9. Materiały dodatkowe

Więcej informacji o różniczkowaniu macierzy oraz krótki przegląd związanej z nim algebry liniowej można znaleźć u Magnusa i Neudeckera (2007). Różniczkowanie automatyczne ma długą historię. Więcej szczegółów znajdziesz w pracach Griewanka i Walthera (2003), Griewanka i Walthera (2008) oraz Elliotta (2009), a także w pracach wymienionych w bibliografiach dołączonych do tych pozycji.

W uczeniu maszynowym (i w innych dziedzinach) często interesują nas wartości oczekiwane, czyli całki postaci

$$\mathbb{E}_x[f(\mathbf{x})] = \int f(\mathbf{x})p(\mathbf{x})d\mathbf{x}. \quad (5.181)$$

Nawet jeśli $p(\mathbf{x})$ ma dogodną formę (np. jest funkcją gęstości rozkładu normalnego), to w ogólnym przypadku całka ta nie ma rozwiązania analitycznego. Rozwinięcie f w szereg Taylora jest jednym ze sposobów na znalezienie przybliżonego rozwiązania. Jeżeli przyjmiemy, że $p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, czyli, że $p(\mathbf{x})$ jest funkcją gęstości prawdopodobieństwa rozkładu normalnego, to rozwinięcie w szereg Taylora pierwszego rzędu w otoczeniu $\boldsymbol{\mu}$ będzie lokalną linearyzacją nieliniowej funkcji f . Jeśli $p(\mathbf{x})$ jest funkcją gęstości rozkładu normalnego (patrz podrozdział 6.5), to w przypadku liniowych funkcji f możemy dokładnie wyznaczyć średnią (i kowariancję). Ta własność jest intensywnie

wykorzystywana przez **rozszerzony filtr Kalmana** (ang. *extended Kalman filter*; Maybeck, 1979) do bieżącej estymacji stanu w nieliniowych układach dynamicznych (zwanymi również *modelami przestrzeni stanów*, ang. *state-space models*). Inne deterministyczne sposoby przybliżania całki z równania 5.181 to **transformacja bezśladowa** (ang. *unscented transform*, Julier i Uhlmann, 1997), która nie wymaga znajomości gradientów, oraz **aproksymacja Laplace’a** (MacKay, 2003; Bishop, 2006; Murphy, 2012), która wykorzystuje rozwinięcie lokalnego przybliżenia funkcji gęstości rozkładu normalnego w szereg Taylora drugiego rzędu (wymagana jest znajomość hesjanu) wokół jego mody.

Ćwiczenia

- 5.1. Oblicz pochodną $f'(x)$ funkcji

$$f(x) = \log(x^4) \sin(x^3).$$

- 5.2. Wyznacz pochodną $f'(x)$ sigmoidalnej funkcji logistycznej

$$f(x) = \frac{1}{1 + \exp(-x)}.$$

- 5.3. Oblicz pochodną $f'(x)$ funkcji

$$f(x) = \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right),$$

gdzie $\mu, \sigma \in \mathbb{R}$ są stałymi.

- 5.4. Dla $n = 0, \dots, 5$ znajdź wielomiany Taylora T_n funkcji $f(x) = \sin(x) + \cos(x)$ w punkcie $x_0 = 0$.

- 5.5. Rozważ następujące funkcje:

$$\begin{aligned} f_1(\mathbf{x}) &= \sin(x_1) \cos(x_2), & \mathbf{x} &\in \mathbb{R}^2 \\ f_2(\mathbf{x}, \mathbf{y}) &= \mathbf{x}^\top \mathbf{y}, & \mathbf{x}, \mathbf{y} &\in \mathbb{R}^n \\ f_3(\mathbf{x}) &= \mathbf{x} \mathbf{x}^\top, & \mathbf{x} &\in \mathbb{R}^n \end{aligned}$$

a) Jakie są wymiary $\frac{\partial f_i}{\partial \mathbf{x}}$?

b) Wyznacz gradienty/macierze Jacobiego powyższych funkcji.

- 5.6. Znajdź pochodną f względem \mathbf{t} i pochodną g względem \mathbf{X}

$$\begin{aligned} f(\mathbf{t}) &= \sin(\log(\mathbf{t}^\top \mathbf{t})), & \mathbf{t} &\in \mathbb{R}^D \\ g(\mathbf{X}) &= \text{tr}(\mathbf{A} \mathbf{X} \mathbf{B}), & \mathbf{A} &\in \mathbb{R}^{D \times E}, \quad \mathbf{X} \in \mathbb{R}^{E \times F}, \quad \mathbf{B} \in \mathbb{R}^{F \times D}, \end{aligned}$$

gdzie $\text{tr}(\cdot)$ oznacza ślad.

- 5.7. Za pomocą reguły łańcuchowej oblicz pochodne $df/d\mathbf{x}$ następujących funkcji. Podaj wymiary każdej pochodnej cząstkowej. Szczegółowo opisz swoje działania.

a)

$$f(z) = \log(1 + z), \quad z = \mathbf{x}^\top \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^D$$

b)

$$f(\mathbf{z}) = \sin(\mathbf{z}), \quad \mathbf{z} = \mathbf{A} \mathbf{x} + \mathbf{b}, \quad \mathbf{A} \in \mathbb{R}^{E \times D}, \quad \mathbf{x} \in \mathbb{R}^D, \quad \mathbf{b} \in \mathbb{R}^E,$$

gdzie $\sin(\cdot)$ jest obliczany dla każdego elementu wektora \mathbf{z} .

5.8. Oblicz pochodne $df/d\mathbf{x}$ następujących funkcji. Opisz szczegółowo wykonywane kroki.

a) Wykorzystaj regułę łańcuchową. Podaj wymiary każdej pochodnej cząstkowej.

$$f(\mathbf{z}) = \exp\left(-\frac{1}{2}\mathbf{z}\right)$$

$$\mathbf{z} = g(\mathbf{y}) = \mathbf{y}^\top \mathbf{S}^{-1} \mathbf{y}$$

$$\mathbf{y} = h(\mathbf{x}) = \mathbf{x} - \boldsymbol{\mu},$$

gdzie $\mathbf{x}, \boldsymbol{\mu} \in \mathbb{R}^D, \mathbf{S} \in \mathbb{R}^{D \times D}$.

b)

$$f(\mathbf{x}) = \text{tr}(\mathbf{x}\mathbf{x}^\top + \sigma^2 \mathbf{I}), \mathbf{x} \in \mathbb{R}^D,$$

gdzie $\text{tr}(\mathbf{A})$ jest śladem \mathbf{A} , czyli sumą elementów A_{ii} leżących na przekątnej. *Wskazówka: zapisz iloczyn zewnętrzny w jawnej postaci.*

c) Wykorzystaj regułę łańcuchową. Podaj wymiary każdej pochodnej cząstkowej. Nie musisz jawnie obliczać iloczynu pochodnych cząstkowych.

$$\mathbf{f} = \text{tgh}(\mathbf{z}) \in \mathbb{R}^M$$

$$\mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{x} \in \mathbb{R}^N, \mathbf{A} \in \mathbb{R}^{M \times N}, \mathbf{b} \in \mathbb{R}^M,$$

gdzie $\text{tgh}(\cdot)$ jest obliczany dla każdego elementu wektora \mathbf{z} .

5.9. Niech

$$g(\mathbf{z}, \mathbf{v}) := \log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}, \mathbf{v})$$

$$\mathbf{z} := t(\boldsymbol{\epsilon}, \mathbf{v}),$$

gdzie p, q i t są funkcjami różniczkowalnymi oraz $\mathbf{x} \in \mathbb{R}^D, \mathbf{z} \in \mathbb{R}^E, \mathbf{v} \in \mathbb{R}^F, \boldsymbol{\epsilon} \in \mathbb{R}^G$. Za pomocą reguły łańcuchowej wyznacz gradient

$$\frac{d}{d\mathbf{v}} g(\mathbf{z}, \mathbf{v}).$$

PROGRAM PARTNERSKI

— GRUPY HELION —

1. ZAREJESTRUJ SIĘ
2. PREZENTUJ KSIĄŻKI
3. ZBIERAJ PROWIZJĘ

Zmień swoją stronę WWW w działający bankomat!

Dowiedz się więcej i dołącz już dzisiaj!

<http://program-partnerski.helion.pl>

GRUPA
Helion 

Uczenie maszynowe staje się wszechobecne. Dzięki coraz lepszym narzędziom służącym do tworzenia aplikacji szczegóły techniczne związane z obliczeniami i modelami matematycznymi są często pomijane przez projektantów. Owszem, to wygodne podejście, ale wiąże się z ryzykiem braku świadomości co do wszystkich konsekwencji wybranych rozwiązań projektowych, szczególnie ich mocnych i słabych stron. A zatem bez ugruntowanych podstaw matematyki nie można mówić o profesjonalnym podejściu do uczenia maszynowego.

Ten podręcznik jest przeznaczony dla osób, które chcą dobrze zrozumieć matematyczne podstawy uczenia maszynowego i nabrać praktycznego doświadczenia w używaniu pojęć matematycznych. Wyjaśniono tutaj stosowanie szeregu technik matematycznych, takich jak algebra liniowa, geometria analityczna, rozkłady macierzy, rachunek wektorowy, optymalizacja, probabilistyka i statystyka. Następnie zaprezentowano matematyczne aspekty czterech podstawowych metod uczenia maszynowego: regresji liniowej, analizy głównych składowych, modeli mieszanin rozkładów Gaussa i maszyn wektorów nośnych. W każdym rozdziale znalazły się przykłady i ćwiczenia ułatwiające przyswojenie materiału.

W książce między innymi:

- podstawy algebry: układy równań, macierze, przestrzenie afiniczne
- rachunek prawdopodobieństwa, sprzężenia, optymalizacja
- wnioskowanie z wykorzystaniem różnego rodzaju modeli
- regresja liniowa i redukcja wymiarowości
- maszyna wektorów nośnych i rozwiązania numeryczne

**Matematyka: koniecznie,
jeśli chcesz zrozumieć istotę
sztucznej inteligencji!**

Marc Peter Deisenroth

kieruje zakładem sztucznej inteligencji w University College London. W swojej pracy badawczej zajmuje się efektywnym uczeniem, modelowaniem probabilistycznym i autonomicznym podejmowaniem decyzji.

A. Aldo Faisal

kieruje laboratorium Brain & Behavior w Imperial College London, gdzie jest również wykładowcą i członkiem Data Science Institute. W swoich badaniach zajmuje się zagadnieniami na styku neuronauki i uczenia maszynowego.

Cheng Soon Ong

jest głównym badaczem w Machine Learning Research Group i adiunktem na Australian National University. Koncentruje się na rozwijaniu statystycznych metod uczenia maszynowego.

Helion



helion.pl



HELION SA
ul. Kościuszki 1c
44-100 Gliwice
tel.: 32 230 98 63
helion@helion.pl

KOD KORZYŚCI
Sięgnij po więcej! ▶



ISBN 978-83-283-8459-0



9 788328 384590

Cena: 99,00 zł