

Ian Goodfellow ▪ Yoshua Bengio ▪ Aaron Courville

DEEP LEARNING

Systemy uczące się



DEEP LEARNING

Ian Goodfellow ▪ Yoshua Bengio ▪ Aaron Courville

DEEP LEARNING

Systemy uczące się



Dane oryginału

Copyright © 2016 Massachusetts Institute of Technology. Title of English-language original: **Deep learning** ISBN 978-1-59327-741-3, published by MIT Press. Polish-language edition copyright © 2018 by Polish Scientific Publishers PWN Wydawnictwo Naukowe PWN Spółka Akcyjna. All rights reserved.

Przekład **WITKOM Witold Sikorski**

Projekt okładki polskiego wydania **Hubert Zacharski**

Ilustracja na okładce: Central Park Azalea Walk Dreamscape, autor – Daniel Ambrosi (danielambrosi.com) Dreamscapes Daniela Ambrosiego są tworzone dzięki zastosowaniu otwartego oprogramowania DeepDream Google’a, zmodyfikowanego przez Josepha Smarra (Google) i Chrisa Lamba (NVIDIA) tak, aby z powodzeniem działało na panoramicznych obrazach złożonych z wielu setek pikseli tworzonych przez Ambrosiego.

Wydawca **Łukasz Łopuszański**

Redaktor prowadzący **Adam Kowalski**

Redaktor **Joanna Forysiak**

Koordynator produkcji **Anna Bączkowska**

Skład i łamanie **Fixpoint**

Zastrzeżonych nazw firm i produktów użyto w książce wyłącznie w celu identyfikacji.

Copyright © for the Polish edition by Wydawnictwo Naukowe PWN S.A.
Warszawa 2018

ISBN: 978-83-01-19583-0

Wydanie I
Warszawa 2018

Wydawnictwo Naukowe PWN SA
tel. 22 69 54 321; faks 22 69 54 288
infolinia 801 33 33 88
e-mail: pwn@pwn.com.pl, reklama@pwn.pl
www.pwn.pl

Druk i oprawa: OSDW Azymut Sp. z o.o.

Spis treści

1. Wprowadzenie	1
1.1. Kto powinien przeczytać tę książkę?	9
1.2. Historyczne trendy deep learningu	11
I Podstawy matematyki stosowanej i systemów uczących się	27
2. Algebra liniowa	29
2.1. Skalary, wektory, macierze i tensory	29
2.2. Mnożenie macierzy i wektorów	32
2.3. Macierze jednostkowe i odwrotne	34
2.4. Zależność liniowa i zakres	35
2.5. Normy	37
2.6. Macierze i wektory specjalne	38
2.7. Rozkład na wartości własne	40
2.8. Dekompozycja wartości osobliwej	42
2.9. Uogólniona macierz odwrotna (Moore’a–Penrose’a)	43
2.10. Operator śladowy	44
2.11. Wyznacznik	45
2.12. Przykład: analiza głównych składowych	45
3. Prawdopodobieństwo i teoria informacji	51
3.1. Dlaczego prawdopodobieństwo?	52
3.2. Zmienne losowe	54
3.3. Rozkłady prawdopodobieństwa	54
3.4. Prawdopodobieństwo brzegowe	56
3.5. Prawdopodobieństwo warunkowe	57
3.6. Reguła łańcuchowa w prawdopodobieństwie warunkowym	57
3.7. Niezależność oraz niezależność warunkowa	58
3.8. Wartość oczekiwana, wariancja i kowariancja	58

3.9.	Znane rozkłady prawdopodobieństwa	60
3.10.	Użyteczne cechy elementarnych funkcji	65
3.11.	Prawo Bayesa	68
3.12.	Techniczne szczegóły zmiennych ciągłych	68
3.13.	Teoria informacji	70
3.14.	Strukturalne modele probabilistyczne	73
4.	Obliczenia numeryczne	77
4.1.	Nadmiar i niedomiar	77
4.2.	Złe uwarunkowania	79
4.3.	Optymalizacja gradientowa	79
4.4.	Optymalizacja z ograniczeniami	89
4.5.	Przykład: liniowa metoda najmniejszych kwadratów	92
5.	Podstawy systemów uczących się	95
5.1.	Algorytmy uczenia się	96
5.2.	Pojemność, nadmierne dopasowanie i niedopasowanie	108
5.3.	Hiperparametry i zbiory walidacyjne	118
5.4.	Estymatory, obciążenie i wariancja	120
5.5.	Metoda maksymalnej wiarygodności	129
5.6.	Statystyki Bayesa	133
5.7.	Algorytmy nadzorowanego uczenia się	138
5.8.	Algorytmy nienadzorowanego uczenia się	143
5.9.	Metoda gradientu stochastycznego	150
5.10.	Tworzenie algorytmu dla systemu uczącego się	152
5.11.	Wyzwania motywujące deep learning	153
II	Głębokie sieci: nowoczesne praktyki	163
6.	Głębokie sieci jednokierunkowe	165
6.1.	Przykład: uczenie się funkcji XOR	168
6.2.	Uczenie się oparte na gradiencie	173
6.3.	Jednostki ukryte	188
6.4.	Projekt architektury	195
6.5.	Propagacja wsteczna i inne algorytmy różniczkowania	201
6.6.	Uwagi historyczne	221
7.	Regularyzacja w deep learningu	225
7.1.	Standardowe kary dla parametrów	227
7.2.	Standardowe kary jako optymalizacja z ograniczeniami	234

7.3.	Regularyzacja i problemy niedoograniczone	236
7.4.	Powiększanie zbioru danych	237
7.5.	Odporność na szum	239
7.6.	Uczenie się częściowo nadzorowane	241
7.7.	Uczenie się wielozadaniowe	242
7.8.	Wczesne zatrzymanie	243
7.9.	Wiązanie i współdzielenie parametrów	250
7.10.	Rzadko wypełnione reprezentacje	252
7.11.	Bagging i inne metody zespołowe	254
7.12.	Odrzucanie	256
7.13.	Szkolenie antagonistyczne	266
7.14.	Odległość styczna, propagacja stycznej oraz klasyfikator stycznej do różnorodności	268
8.	Optymalizacja w celu szkolenia głębokich modeli	273
8.1.	Czym uczenie się różni się od czystej optymalizacji	274
8.2.	Wyzwania związane z optymalizacją sieci neuronowej	281
8.3.	Podstawowe algorytmy	293
8.4.	Strategie nadawania parametrom wartości początkowych	299
8.5.	Algorytmy z adaptacyjną szybkością uczenia się	306
8.6.	Aproksymacyjne metody drugiego rzędu	310
8.7.	Strategie optymalizacji i meta-algorytmy	317
9.	Sieci splotowe	331
9.1.	Splot jako działanie	332
9.2.	Uzasadnienie	334
9.3.	Redukcja	340
9.4.	Splot i redukcja jako nieskończenie silny rozkład aprioryczny	346
9.5.	Warianty podstawowej funkcji splotowej	347
9.6.	Strukturalne wyjścia	358
9.7.	Typy danych	359
9.8.	Efektywne algorytmy splotu	361
9.9.	Cechy losowe lub nienadzorowane	362
9.10.	Neuronaukowe podstawy sieci splotowych	364
9.11.	Sieci splotowe a historia deep learningu	371
10.	Modelowanie sekwencyjne: sieci rekurencyjne i rekursywne	373
10.1.	Rozwijanie grafów obliczeniowych	375
10.2.	Rekurencyjne sieci neuronowe	378

10.3.	Dwukierunkowe rekurencyjne sieci neuronowe	393
10.4.	Architektury koder-dekoder i sekwencja do sekwencji	394
10.5.	Głębokie sieci rekurencyjne	397
10.6.	Rekursywne sieci neuronowe	399
10.7.	Problem z zależnościami długoterminowymi	400
10.8.	Sieci stanu echa	403
10.9.	Nieszczelne jednostki i inne strategie dla wielu skali czasowych	406
10.10.	Długa pamięć krótkoterminowa i inne bramkowane sieci RNN	408
10.11.	Optymalizacja zależności długoterminowych	412
10.12.	Pamięć jawna	416
11.	Metodologia praktyczna	421
11.1.	Metryki wydajności	422
11.2.	Modele domyślnej linii bazowej	425
11.3.	Decyzja, czy zbierać więcej danych	426
11.4.	Wybór hiperparametrów	428
11.5.	Strategie debugowania	437
11.6.	Przykład: rozpoznawanie liczb wielocyfrowych	441
12.	Zastosowania	445
12.1.	Deep learning wielkoskalowy	445
12.2.	Rozpoznawanie obrazów	455
12.3.	Rozpoznawanie mowy	461
12.4.	Przetwarzanie języka naturalnego	464
12.5.	Inne zastosowania	482
III	Badania na polu deep learningu	491
13.	Liniowe modele czynnikowe	495
13.1.	Probabilistyczna analiza PCA i analiza czynnikowa	496
13.2.	Analiza składowych niezależnych (ICA)	497
13.3.	Powolna analiza cech	500
13.4.	Rzadkie kodowanie	502
13.5.	Poznawanie różnorodności w analizie PCA	506
14.	Autokodery	509
14.1.	Autokodery niekompletne	510
14.2.	Autokodery z regularyzacją	511

14.3.	Reprezentacyjna potęga, rozmiar warstwy i głębokość	515
14.4.	Stochastyczne kodery i dekodery	516
14.5.	Autokodery z odszumianiem	517
14.6.	Poznanwanie różnorodności z użyciem autokoderów	522
14.7.	Autokodery kurczliwe	527
14.8.	Przybliżona rzadka dekompozycja	530
14.9.	Zastosowania autokoderów	531
15.	Poznanwanie reprezentacji	533
15.1.	Zachłanne nienadzorowane szkolenie wstępne warstwa po warstwie	535
15.2.	Transfer poznawania i adaptacja dziedziny	544
15.3.	Częściowo nadzorowane oswabianie czynników przyczynowych	548
15.4.	Reprezentacja rozproszona	554
15.5.	Wykładnicze zyski z głębokości	560
15.6.	Wskazówki do wykrywania przyczyn podstawowych	562
16.	Strukturalne modele probabilistyczne deep learningu	567
16.1.	Trudności w modelowaniu niestukturalnym	568
16.2.	Używanie grafów do opisu struktury modelu	572
16.3.	Próbkowanie z modeli graficznych	589
16.4.	Zalety modelowania strukturalnego	591
16.5.	Poznanwanie zależności	591
16.6.	Wnioskowanie i wnioskowanie przybliżone	592
16.7.	Strukturalne modele probabilistyczne w ujęciu deep learningu	594
17.	Metody Monte Carlo	599
17.1.	Próbkowanie i metody Monte Carlo	599
17.2.	Próbkowanie istotnościowe	601
17.3.	Metody Monte Carlo z łańcuchem Markowa	604
17.4.	Próbkowanie Gibbsa	608
17.5.	Problem mieszania między odseparowanymi trybami	609
18.	Zmagania z funkcją podziału	615
18.1.	Gradient wiarygodności logarytmicznej	616
18.2.	Stochastyczna maksymalna wiarygodność i kontrastowa dywergencja	617
18.3.	Pseudowiarygodność	625
18.4.	Dopasowywanie oceny i stosunku	628

18.5.	Dopasowywanie ocen z odszumianiem	630
18.6.	Estymacja kontrastywna szumu	630
18.7.	Szacowanie funkcji podziału	633
19.	Wnioskowanie przybliżone	641
19.1.	Wnioskowanie jako optymalizacja	642
19.2.	Maksymalizacja oczekiwania	644
19.3.	Wnioskowanie MAP i rzadkie kodowanie	645
19.4.	Wariacyjne wnioskowanie i uczenie się	648
19.5.	Poznawanie wnioskowania przybliżonego	661
20.	Głębokie modele generatywne	665
20.1.	Maszyny Boltzmanna	665
20.2.	Ograniczone maszyny Boltzmanna	667
20.3.	Głębokie sieci przekonań	671
20.4.	Głębokie maszyny Boltzmanna	674
20.5.	Maszyny Boltzmanna dla danych rzeczywistych	688
20.6.	Splotowe maszyny Boltzmanna	695
20.7.	Maszyny Boltzmanna dla strukturalnych lub sekwencyjnych wartości wynikowych	697
20.8.	Inne maszyny Boltzmanna	698
20.9.	Propagacja wsteczna przez losowe działania	700
20.10.	Skierowane sieci generatywne	704
20.11.	Pobieranie próbek z autokoderów	724
20.12.	Generatywne sieci stochastyczne	727
20.13.	Inne schematy generowania	729
20.14.	Szacowanie modeli generatywnych	730
20.15.	Konkluzja	733
Bibliografia		735
Skorowidz		800

1

Wprowadzenie

Wynalazcy od dawna marzyli o myślących maszynach. Te marzenia sięgają przynajmniej do czasów starożytnej Grecji. Mityczne postaci Pigmaliona, Dedala czy Hefajstosa można interpretować jako legendarnych wynalazców, a Galateę, Talosa i Pandorę rozpatrywać jako sztuczne życie (Ovid and Martin 2004, Sparkes 1996, Tandy 1997).

Już wtedy, gdy powstały pierwsze programowane komputery, ludzie zaczęli się zastanawiać, czy mogą one stać się inteligentne (Lovelace 1842) – ponad 100 lat wcześniej zanim zbudowano pierwsze inteligentne maszyny. Dziś **sztuczna inteligencja** (AI) to ekscytująca dziedzina z wieloma zastosowaniami praktycznymi i programami badawczymi. Wprowadza się inteligentne oprogramowanie, by zautomatyzować rutynowe prace, zrozumieć mowę lub obrazy, wykonywać diagnostykę medyczną i wspomagać podstawowe badania naukowe.

We wczesnych latach sztucznej inteligencji szybko rozwiązywane były problemy trudne do wykonania dla ludzi, ale dość proste dla komputerów – problemy, które można opisać jako listy sformalizowanych matematycznych reguł. Prawdziwym wyzwaniem dla sztucznej inteligencji okazało się rozwiązywanie zadań łatwych do wykonania dla ludzi, ale trudnych w opisie formalnym – problemów, które rozwiązujemy intuicyjnie, niejako automatycznie, jak rozpoznawanie słów lub twarzy na obrazach.

Książka ta poświęcona jest rozwiązywaniu właśnie tych intuicyjnych przeszkód. By je pokonać, trzeba pozwolić komputerom, aby uczyły się na swoich doświadczeniach i rozumiały świat w kontekście hierarchii pojęć, z których każde jest zdefiniowane przez jego związek z pojęciem prostszym. Dzięki gromadzeniu wiedzy na podstawie doświadczeń unika się potrzeby formalnego opisanego przez ludzi całej wiedzy potrzebnej komputerowi. Hierarchia pojęć

pozwala mu na uczenie się skomplikowanych pojęć przez budowanie ich z pojęć prostszych. Jeśli narysujemy schemat, w którym pojęcia te nakładają się na siebie, to powstanie wiele warstw. Dlatego to podejście do AI nazywamy *głębokim uczeniem się**.

Wiele początkowych sukcesów na polu AI miało miejsce we względnie sterylnych i sformalizowanych środowiskach i nie wymagało od komputerów specjalnej wiedzy o świecie. Na przykład system gry w szachy komputera Deep Blue firmy IBM pokonał w 1997 roku mistrza świata Garriego Kasparowa (Hsu 2002). Szachy to prosty świat, obejmujący tylko 64 pola i 32 elementy, które mogą poruszać się według ściśle określonych reguł. Opracowanie prowadzącej do sukcesu strategii gry w szachy to niezwykle osiągnięcie, ale nie wynika ono z trudności opisanego komputerowi figur szachowych i dopuszczalnych ruchów. Szachy można dokładnie opisać za pomocą krótkiej listy całkiem sformalizowanych reguł podanych przez programistę.

Jak na ironię, abstrakcyjne i sformalizowane zadania, należące do najtrudniejszych dla człowieka, są najłatwiejszymi dla komputera. Komputery od dawna mogą pokonać najlepszego szachistę, ale dopiero ostatnio zaczęły nabywać niektóre zwykłe ludzkie umiejętności, jak rozpoznawanie obiektów czy mowy. Życie codzienne wymaga sporej wiedzy o świecie. Jej duża część jest subiektywna i intuicyjna, dlatego trudno ją formalnie opisać. Aby działać inteligentnie, komputery muszą zdobyć tę samą wiedzę. Jednym z kluczowych wyzwań w dziedzinie sztucznej inteligencji jest zatem przekazanie tej nieformalnej wiedzy komputerowi.

W kilku projektach z dziedziny sztucznej inteligencji poszukiwano twardej wiedzy o świecie podanej językiem sformalizowanym – komputer może automatycznie tworzyć argumenty na podstawie tych instrukcji za pomocą logicznych reguł wnioskowania. Nazywa się to podejściem do sztucznej inteligencji **opartym na wiedzy**. Żaden z tych projektów nie zakończył się sukcesem. Najbardziej znany jest Cyc (Lenat and Guha 1989). Cyc to mechanizm wnioskowania i baza danych instrukcji w języku o nazwie CycL. Instrukcje te są wprowadzane przez personel złożony z ludzi nadzorujących projekt. Jest to uciążliwy proces. Ludzie męczą się, tworząc sformalizowane reguły o takiej złożoności, aby opisywały świat. Cyc nie potrafił zrozumieć opowieści o człowieku imieniem Fred, który rano się golił (Linde 1992). Jego mechanizm rozumowania wykrył niespójność w tej opowieści: wiedział,

*Termin *machine learning* jest tłumaczony na język polski na kilka sposobów. W tłumaczeniu przyjęto określenie *systemy uczące się*, aby podkreślić jednoznacznie, że systemy same się uczą, a nie uczą kogoś, co sugeruje sformułowanie *uczenie maszynowe*. Podobnie więc mamy *głębokie uczenie się*. W tekście zachowano jednak określenie *deep learning*, jako nowe pojęcie, którego przekład nie jest w Polsce utrwalony (wszystkie przypisy oznaczone * pochodzą od tłumacza).

że ludzie nie mają części elektrycznych, ale Fred miał w ręku elektryczną golarkę, więc uwierzył, że jednostka „FredPodczasGolenia” zawierała elementy elektryczne. Stąd pojawiło się pytanie, czy podczas golenia Fred nadal jest osobą.

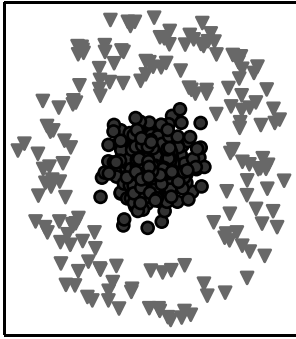
Trudności systemu opierającego się na twardo zapisanej wiedzy sugerują, że systemy AI potrzebują zdolności do przyswajania własnej wiedzy wyodrębnianej jako wzorce z surowych danych. Takie zdolności określa się mianem **systemów uczących się**. Ich wprowadzenie pozwoliło komputerom na zajęcie się problemami obejmującymi wiedzę o świecie i podejmowanie decyzji, które wydają się subiektywne. Prosty algorytm samouczenia się zwany regresją logistyczną pozwala określić, czy zalecać cesarskie cięcie (Mor-Yosef et al. 1990). Inny prosty algorytm samouczenia zwany „naiwny Baynes” potrafi oddzielić ważne e-maile od spamu.

Działanie tych algorytmów mocno zależy od reprezentatywności otrzymanych danych. Na przykład, gdy regresja logistyczna jest używana do wskazania cesarskiego cięcia, system AI nie bada pacjenta. To lekarz podaje systemowi kilka odpowiednich informacji, jak obecność blizn macicznych lub ich brak. Każdy fragment informacji ujęty w opisie pacjentki jest określany jako cecha. Regresja logistyczna uczy się, jaki związek ma każda z cech pacjentki z różnymi rezultatami. Nie może jednak w żaden sposób wpływać na sposób definiowania tych cech. Jeśli regresja logistyczna otrzyma skan rezonansu magnetycznego (MRI) pacjentki zamiast formalnego raportu lekarza, nie potrafi dać użytecznych wskazań. Pojedyncze piksele skanowania MRI mają pomijalną korelację z jakimikolwiek komplikacjami, jakie mogą pojawić się przy porodzie.

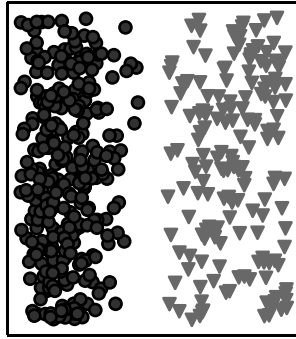
Ta zależność od sposobu reprezentacji jest ogólnym zjawiskiem w całej dziedzinie informatyki, a także w życiu codziennym. W naukach komputerowych działania takie, jak przeszukiwanie zbioru danych, są w wykładniczo szybsze, jeśli zbiór ma strukturę i jest inteligentnie indeksowany. Ludzie z łatwością wykonują działania arytmetyczne na cyfrach arabskich, ale działania na cyfrach rzymskich zajmują im znacznie więcej czasu. Nie jest więc zaskoczeniem, że wybrany sposób reprezentacji ma ogromny wpływ na wydajność algorytmów samouczących się. Prosty obrazowy przykład pokazano na rysunku 1.1.

Wiele zadań z zakresu sztucznej inteligencji można rozwiązać, projektując odpowiedni dla danego zadania zestaw cech, a następnie dostarczając te cechy prostemu algorytmowi uczącemu się. Na przykład użyteczną cechą identyfikacji mówcy na podstawie dźwięku jest ocena rozmiaru jego toru głosowego. Cecha ta daje nam możliwe przesłanki do określenia, czy mówiący jest kobietą, mężczyzną czy dzieckiem. Jednak dla wielu zadań trudno określić

Współrzędne kartezjańskie



Współrzędne



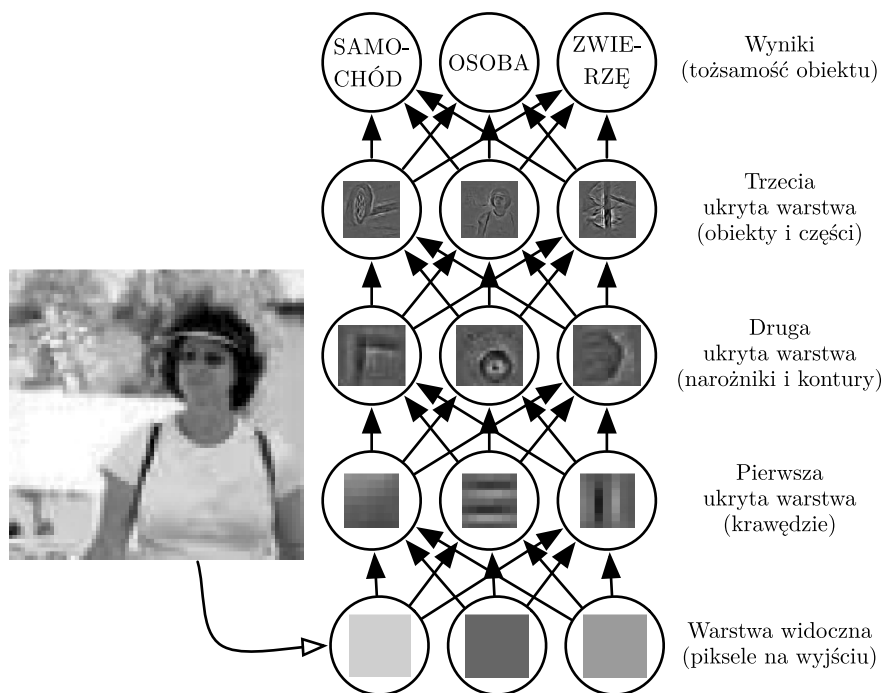
uczają się przechowywać możliwie wiele informacji, gdy dane przechodzą przez proces kodowania, a potem przez dekodery, ale także przypisują nowej reprezentacji różne właściwości. Różne rodzaje autokodowania mają na celu uzyskanie różnych rodzajów właściwości.

Gdy projektujemy cechy lub algorytmy do uczenia się cech, naszym celem jest zwykle oddzielenie czynników zmiennych, które tłumaczą obserwowane dane. W tym kontekście używamy słowa „czynniki”, aby odnieść się do oddzielnych źródeł wpływu; czynniki nie są zwykle powiązane z mnożeniem. Takie czynniki często nie są bezpośrednio obserwowanymi wielkościami. Mogą występować jako nieobserwowane obiekty lub siły w świecie rzeczywistym, które wpływają na obserwowalne wielkości. Mogą także występować jako konstrukcje w umyśle człowieka, które dostarczają upraszczających wyjaśnień lub wnioskowanych przyczyn obserwowanych danych. Można je traktować jak pojęcia lub abstrakcje pomagające nam zrozumieć bogatą zmienność danych. Podczas analizy nagrań mowy czynniki zmienne obejmują wiek mówcy, jego płeć i akcent oraz wypowiedane słowa. Podczas analizy obrazu samochodu: położenie samochodu, jego kolor oraz kąt padania promieni słonecznych.

Podstawowym źródłem trudności w aplikacjach sztucznej inteligencji jest fakt, że wiele czynników zmian wpływa na każdy element danych, które możemy obserwować. Pojedyncze piksele obrazu czerwonego samochodu mogą w nocy przypominać czerń. Sylwetka samochodu zależy od kąta, pod którym na nią patrzymy. Większość zastosowań wymaga od nas rozwikłania czynników zmian i odrzucenia tych, na których nam nie zależy. Oczywiście wybranie takich cech na wysokim poziomie abstrakcji z nieprzetworzonych danych może być trudne. Wiele czynników zmian, jak akcent mówcy, da się zidentyfikować tylko za pomocą wyrafinowanego rozumienia danych, na poziomie ludzkiego umysłu. Gdy uzyskanie reprezentatywnych informacji jest tak samo trudne, jak rozwiązanie problemu, wydaje się, że taka reprezentatywność na niewiele się zdaje.

Deep learning rozwiązuje ten podstawowy problem w reprezentatywnym uczeniu się dzięki wprowadzeniu reprezentatywności wyrażonej w formie innej, prostszej reprezentacji. Umożliwia komputerowi zbudowanie złożonych pojęć na podstawie prostszych. Na rysunku 1.2 pokazano, jak system deep learning może reprezentować pojęcie opisujące obraz osoby przez połączenie prostszych pojęć, jak narożniki i kontury, które z kolei są zdefiniowane jako krawędzie.

Typowym przykładem modelu deep learningu jest głęboka sieć jednokierunkowa lub wielowarstwowy perceptron (MLP). Ten ostatni to funkcja matematyczna odwzorowująca pewien zestaw danych wejściowych na war-



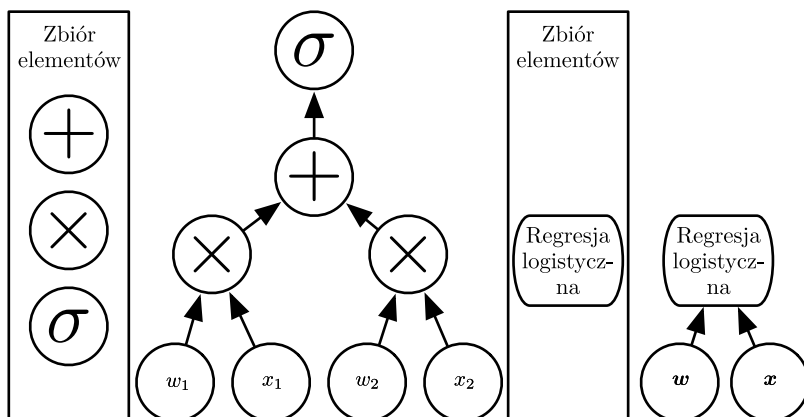
Rysunek 1.2. Ilustracja modelu deep learningu. Komputerowi trudno zrozumieć znaczenie surowych danych pochodzących ze zmysła, jak obraz reprezentowany przez zbiór pikseli. Funkcja odwzorująca zbiór pikseli na tożsamość obiektu jest bardzo skomplikowana. Nauczenie się lub wyznaczenie tego odwzorowania wydaje się nie do pokonania w bezpośrednim podejściu. Deep learning rozwiązuje tę trudność dzięki podziałowi skomplikowanego odwzorowania na serię zagnieżdżonych prostszych odwzorowań, z których każde jest opisane w innej warstwie modelu. Dane wejściowe są prezentowane w widocznej warstwie, którą nazywamy tak, gdyż zawiera zmienne, które można obserwować. Potem następuje ciąg warstw ukrytych, wyciągających z obrazu coraz bardziej abstrakcyjne cechy. Nazywamy je ukrytymi, gdyż ich wartości nie są podane w danych; model musi sam określić, które pojęcia są użyteczne do wyrozumienia związków w obserwowanych danych. Obrazy stanowią tu wizualizację pewnej cechy reprezentowanej przez każdą ukrytą jednostkę. Znajdź piksele, pierwsza warstwa może łatwo zidentyfikować krawędzie, porównując jasność sąsiednich pikseli. Mając opis brzegów zawarty w pierwszej warstwie, druga warstwa może łatwo odnaleźć narożniki i rozszerzone kontury. Mając opis obrazu z drugiej ukrytej warstwy w postaci narożników i konturów, trzecia ukryta warstwa może wykryć całe fragmenty określonych obiektów, dzięki znalezieniu specjalnego zbioru konturów i narożników. Wreszcie ten opis obrazu w kategorii części obiektu, które się na niego składają, może być wykorzystany do rozpoznania obiektów znajdujących się na obrazie (obrazy zamieszczono za Zeiler and Fergus 2014, za zgodą autorów)

tości wynikowe. Funkcja jest zbudowana z wielu mniejszych funkcji. Każde zastosowanie innej funkcji matematycznej możemy rozumieć jako nową reprezentację danych wejściowych. Cel poznania właściwej reprezentacji danych zapewnia jedną perspektywę deep learningu. Inną perspektywą jest fakt, że głębokość warstw pozwala komputerowi nauczyć się wielokrokowego programu komputerowego. Każda warstwa reprezentacji może być rozumiana jako stan pamięci komputera po równoległym wykonaniu kolejnego zbioru instrukcji. Sieci o większej głębokości mogą wykonywać więcej instrukcji sekwencyjnych. Mają one wielką moc, gdyż kolejne instrukcje mogą odwoływać się do wyników wcześniejszych instrukcji.

Zgodnie z tym spojrzeniem na deep learning nie wszystkie warstwy informacji muszą koniecznie kodować czynniki zmian, które objaśniają dane wejściowe. Reprezentacja przechowuje także informacje o stanie, co pomaga stworzyć program, który potrafi zrozumieć dane. Informacja o stanie może stanowić analogię dla licznika lub wskaźnika w tradycyjnym programie komputerowym. Nie ma nic wspólnego z zawartością konkretnych danych, ale pomaga modelowi w organizacji przetwarzania.

Są dwa sposoby pomiaru głębokości modelu. Pierwszy opiera się na liczbie kolejnych instrukcji, które należy wykonać, aby wyznaczyć architekturę problemu. Możemy to traktować jak najdłuższą ścieżkę w schemacie blokowym, który opisuje sposób obliczania każdego z wyników modelu na podstawie danych wejściowych. Podobnie dwa programy komputerowe mają różne długości, zależnie od tego, w jakim języku programowania napisano program. Funkcje można narysować w schemacie blokowym z różnym zagłębieniem, które zależy od tego, które funkcje mogą być używane w kolejnych krokach schematu. Na rysunku 1.3 widać, jak wybór języka może dać dwa różne pomiary tej samej architektury.

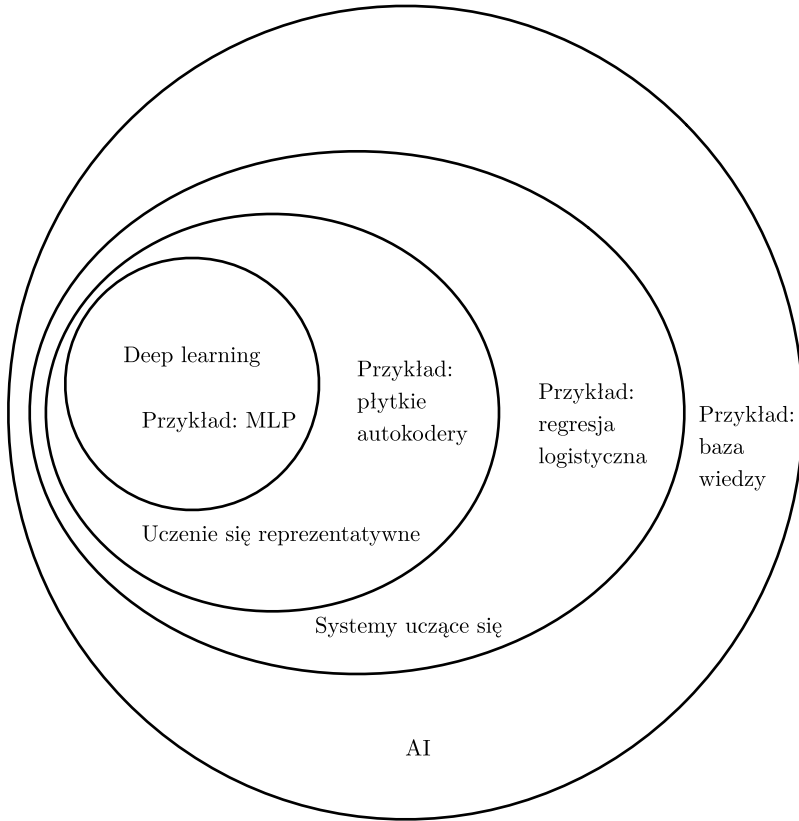
Inne podejście używane w modelach probabilistycznych dotyczy głębokości modelu nie jako głębokości grafu obliczeniowego, lecz jako głębokości grafu opisującego, jak pojęcia są ze sobą powiązane. W tym przypadku głębokość schematu blokowego obliczeń potrzebnych do wyznaczenia reprezentacji każdego z pojęć może być głębsza niż graf samych tych pojęć. Jest tak dlatego, że rozumienie przez system prostszych pojęć może być zdefiniowane przez informację o pojęciach złożonych. Na przykład system AI obserwujący obraz twarzy, gdzie jedno oko jest w cieniu, może z początku zobaczyć tylko jedno oko. Po wykryciu, że na obrazie jest twarz, system może wywnioskować, iż jest tam także drugie oko. W takim przypadku graf pojęć obejmuje tylko dwie warstwy – warstwę dla oczu i warstwę dla twarzy – ale graf obliczeniowy to $2n$ warstw, gdyż estymacja każdego pojęcia wykonywana jest n razy.



Rysunek 1.3. Grafy obliczeniowe odwzorujące wejście na wyjście, gdzie ka»dy węzeł oznacza wykonanie działania. Głębokość to długość najdłuższej ścieżki od wejścia do wyjścia, zależna od decyzji elementów składających się na możliwe kroki obliczeniowe. Obliczenia pokazane na tym grafie to wyniki modelu regresji logistycznej $\sigma(w^T x)$, gdzie σ to logistyczna funkcja sigmoidalna. Jeśli dodawanie, mnożenie i logistyczne sigmoidy będą elementami języka komputerowego, to model ten ma głębokość trzy. Jeśli elementem jest sama regresja logistyczna, to model ma głębokość jeden

Ponieważ nie jest zawsze jasne, które z tych dwóch podejść – głębokość grafu komputerowego czy głębokość grafu modelu probabilistycznego – jest odpowiednia, a także dlatego, że różni ludzie wybiorą inny zbiór najmniejszych elementów, z których zbudują graf, nie ma jednej poprawnej wartości opisującej długość programu komputerowego. Nie ma też zgody, jak dużej głębokości wymaga model, aby można go było nazwać głębokim. Jednak deep learning można bezpiecznie uznać za studium modeli obejmujących większą liczbę składników poznawanych funkcji czy pojęć niż w tradycyjnych systemach uczących się.

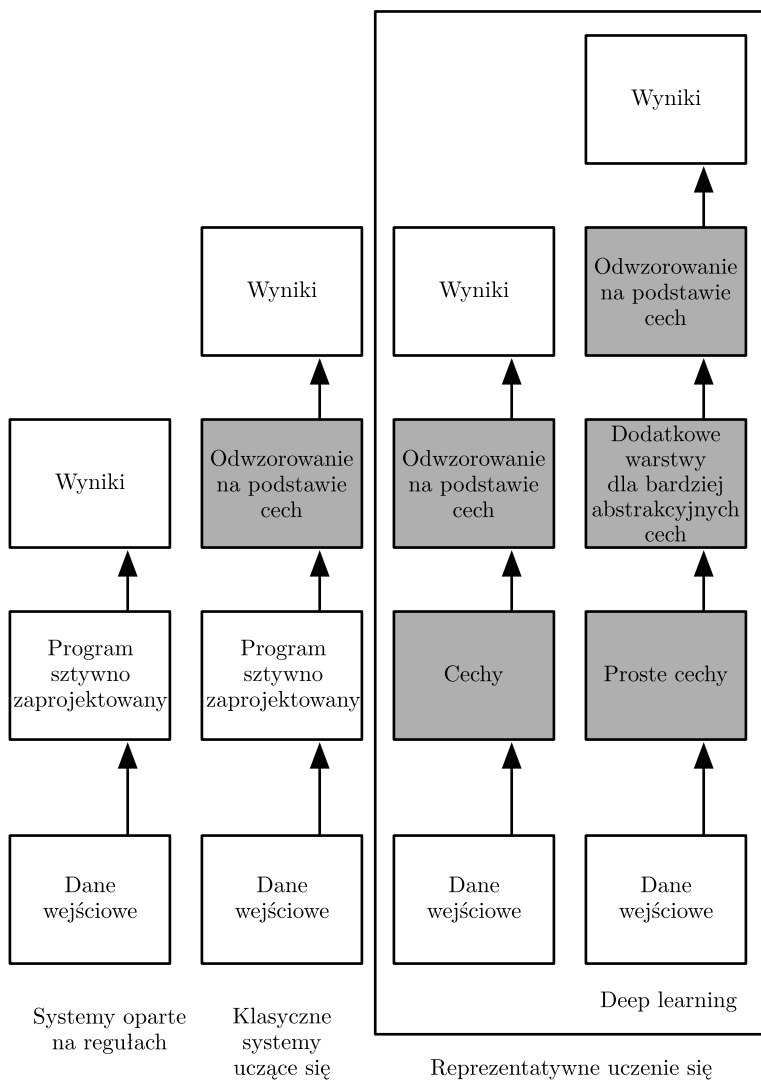
Podsumowując, deep learning, które jest tematem tej książki, stanowi podejście do AI. Konkretnie jest to rodzaj systemów uczących się, technika umożliwiająca systemom komputerowym ulepszenia oparte na doświadczeniu i danych. Systemy uczące się to jedyne sensowne podejście do budowy systemów AI, które mogą działać w skomplikowanych środowiskach świata rzeczywistego. Deep learning jest tym rodzajem systemów uczących się, który osiąga wielką moc i elastyczność, reprezentując świat jako zagnieżdżoną hierarchię pojęć, przy czym każde pojęcie jest zdefiniowane w powiązaniu z pojęciami prostszymi i bardziej abstrakcyjnymi reprezentacjami wyznaczanymi za pomocą mniej abstrakcyjnych. Rysunek 1.4 ilustruje związki między różnymi dyscyplinami AI. Rysunek 1.5 obrazuje schemat wysokiego poziomu pokazujący, jak to działa.



Rysunek 1.4. Diagram Venna pokazujący, dlaczego deep learning jest rodzajem uczenia sił reprezentacji, które z kolei jest typem systemu uczącego sił używanego w wielu, lecz nie wszystkich podejściach do AI. Diagram Venna stanowi przykład technologii AI

1.1. Kto powinien przeczytać tę książkę?

Książka ta przyda się wielu czytelnikom, ale została napisana dla dwóch rodzajów odbiorców. Jedną grupą są studenci (na studiach magisterskich i doktoranckich) zdobywający wiedzę o systemach uczących się, w tym także ci, którzy chcą specjalizować się w badaniach w dziedzinie deep learningu i sztucznej inteligencji. Drugą grupą docelową są inżynierowie oprogramowania, którzy nie mają za sobą nauki z dziedziny systemów uczących się i statystyki, ale chcą szybko nadrobić zaległości i zająć się wykorzystaniem deep learningu w swojej dziedzinie. Metody te sprawdziły się w wielu dziedzinach informatyki, jak wizualizacje komputerowe, przetwarzanie głosu i dźwięku, przetwarzanie języków naturalnych, robotyka, bioinformatyka i chemia, gry wideo, mechanizmy wyszukiwania czy ogłoszenia i finanse w trybie online.



Rysunek 1.5. Schematy blokowe pokazujące różne elementy systemu AI i ich wzajemne związki z różnymi dziedzinami AI. Zacieniowane ramki pokazują komponenty, których można nauczyć sił na podstawie danych

Książka składa się z trzech części, aby najlepiej dopasować ją do potrzeb rozmaitych czytelników. Część I obejmuje opis narzędzi matematycznych i pojęć z dziedziny systemów uczących się. Część II opisuje najważniejsze algorytmy deep learningu, które są poznanymi technikami. Część III zawiera nowe pomysły, które są powszechnie uważane za przyszłość deep learningu.

Czytelnicy mogą swobodnie opuszczać fragmenty, które są poza zakresem ich zainteresowań lub wiadomości. Osoby znające algebrę liniową, prawdopodobieństwo i podstawowe pojęcia systemów uczących się mogą na przykład pominąć część I, a ci, którzy chcą wdrożyć działający system, powinni skończyć na części II. Aby pomóc w wyborze czytanych rozdziałów, na rysunku 1.6 pokazano schemat blokowy ogólnej budowy książki.

Zakładamy, że wszyscy czytelnicy mają pewną wiedzę z dziedziny informatyki. Przyjmujemy, że znają zasady programowania, rozumieją podstawowe zagadnienia związane z wydajnością obliczeń, teorię złożoności, podstawy rachunku różniczkowego oraz pojęcia z teorii grafów.

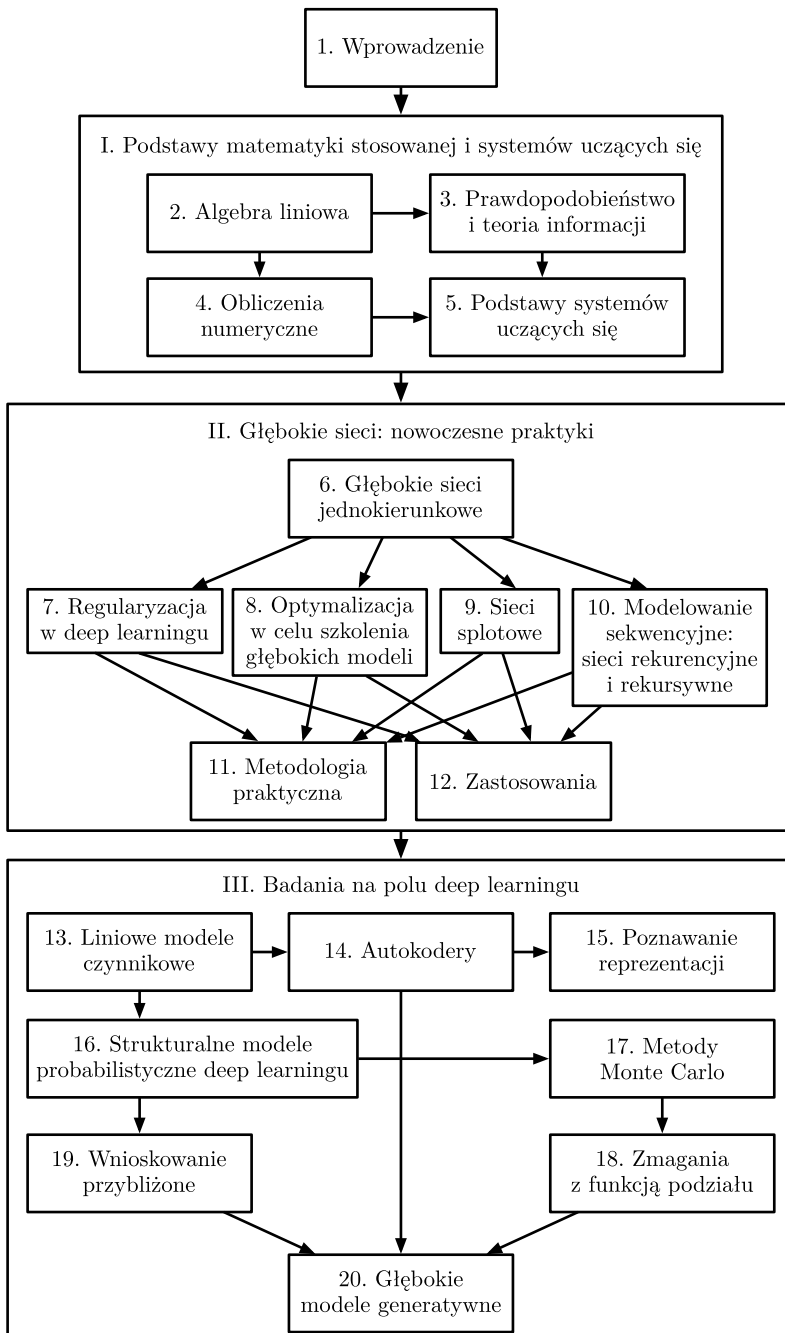
1.2. Historyczne trendy deep learningu

Najłatwiej zrozumieć deep learning w kontekście historycznym. Zamiast podawać jego szczegółową historię, poniżej wyszczególniono najważniejsze trendy:

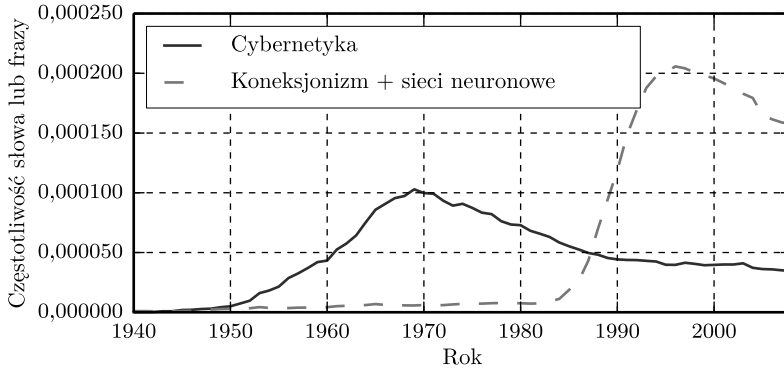
- deep learning ma długą i bogatą historię, ale znane jest pod różnymi nazwami, co odzwierciedla różne filozoficzne punkty widzenia; jego popularność pojawiała się i znikwała,
- deep learning stało się bardziej użyteczne, gdy wzrosła ilość dostępnych danych szkoleniowych,
- modele deep learningu z czasem rozrosły się, a dostępna infrastruktura komputerowa (sprzęt i oprogramowanie) uległa poprawie,
- deep learning pomaga w rozwiązywaniu coraz bardziej skomplikowanych aplikacji, a dokładność jego działania wzrasta.

1.2.1. Wiele nazw i zmienne powodzenie sieci neuronowych

Wielu czytelników tej książki z pewnością słyszało o deep learningu jako o nowej pasjonującej technice, więc są zdziwieni, że jest tu mowa o historii. W rzeczywistość deep learning ma korzenie w latach czterdziestych XX wieku. Wydaje się nowością, ale przez kilka lat poprzedzających obecną popularność był mało znany. Pojawiał się pod różnymi nazwami, a określenie „deep learning” pojawiło się niedawno. Dziedzina ta wiele razy zmieniała terminologię, co odzwierciedlało wpływy różnych naukowców i rozmaite perspektywy.



Rysunek 1.6. Ogólny sposób organizacji księgi. Strzałka z jednego rozdziału do drugiego wskazuje, że poprzedni rozdział zawiera materiał poprzedzający niezbędny do zrozumienia kolejnego



Rysunek 1.7. Dwie historyczne fale bada« nad sieciami neuronowymi wed«ug pomiaru cz«stotliwo«ci fraz cybernetyka i koneksjonizm lub sieci neuronowe wed«ug Google Botoks (trzecia fala jest zbyt nowa, »eby si« tu znale«ć). Pierwsza fala zacz«a si« od cybernetyki w latach 1940–1960, wraz z rozwojem teorii biologicznego uczenia si« (McCulloch and Pitts 1943, Hebb 1949) oraz implementacji pierwszych modeli, jak perceptron (Rosenblatt 1958), co pozwoli«o na szkolenie pojedynczego neuronu. Druga fala zacz«a si« od podej«cia koneksjonistycznego z okresu 1980–1995, z propagacj« wstecznej (Rumelhart et al. 1986), pozwalaj«cej szkoli«ć sieci neuronowe z jedn« lub dwiema ukrytymi warstwami. Obecna trzecia fala zacz«a si« oko«o roku 2006 (Hinton et al. 2006, Bengio et al. 2007, Ranzato et al. 2007), a jej opis w ksi«ce pojawi« si« w roku 2016. Poprzednie dwie fale uj«to w ksi«kach znacznie p«ó«niej ni« odpowiadaj«ce im badania naukowe

Wyczerpuj«ca historia deep learningu wykracza poza zakres tej ksi«żki, jednak dodatkowy kontekst pomaga w zrozumieniu tej dziedziny. Zaznaczy«y si« trzy fale jej rozwoju: deep learning pod nazw« cybernetyki w latach 1940–1960, deep learning znane jako koneksjonizm w latach 1980–1990 oraz obecne odrodzenie ju« jako deep learning pocz«wszy od roku 2006. Pokazano to na rysunku 1.7.

Jedne z pierwszych algorytm«w uczenia si«, kt«re znamy do dzi«s, mia«y na celu modelowanie komputerowe biologicznych proces«w uczenia si«, czyli modelowa«y, jak to si« dzieje, »e uczymy si« i co mo«e si« wtedy dzia« w m«zgu. Wynikiem by«a pierwsza nazwa deep learningu – sztuczne sieci neuronowe (ang. *artificial neural networks*, ANN). Odpowiadaj«ca temu perspektywa modeli deep learningu m«wi, »e s« to systemy in«ynierskie inspirowane przez m«zg biologiczny (ludzki lub zwier«cy). Cho« niekt«re rodzaje sieci neuronowych u«ywane w systemach ucz«cych si« by«y czasem wykorzystywane do zrozumienia dzia«ania m«zgu (Hinton and Shallice 1991), nie s« zwyk«le zaprojektowane jako realne modele funkcji biologicznych. Perspektywa neuronowa deep learningu opiera si« na dw«ch za«o«eniach. Jednym z nich jest twierdzenie, »e m«zg stanowi dow«d na to, i« mo«liwe jest inteligentne

zachowanie i stosunkowo prostą drogą do zbudowania inteligencji jest opracowanie komputerowych zasad działania mózgu i powielanie jego działania. Inną perspektywą jest przyjęcie, że ciekawe byłoby zrozumienie mózgu i zasad rządzących ludzką inteligencją, aby modele systemów uczących się rzucające światło na te naukowe pytania były użyteczne niezależnie od ich zdolności do rozwiązywania zastosowań inżynierskich.

Nowoczesny termin „deep learning” wychodzi poza neuronaukową perspektywę dotyczącą obecnej generacji modeli dla systemów uczących się. Odwołuje się do bardziej ogólnej zasady uczenia się na różnych poziomach złożoności, co można zastosować w schematach systemów uczących się, które nie muszą być inspirowane przez neurony.

Najwcześniejsze modele poprzedzające nowoczesne deep learningu stanowiły modele liniowe motywowane perspektywą neuronaukową. Modele te były zaprojektowane jako pobierające dane wejściowe w postaci zbioru n wartości x_1, \dots, x_n i wiążące je z wynikami w postaci wyników y . Modele te uczą się, jaki jest zbiór wag w_1, \dots, w_n , i obliczają wyniki według funkcji postaci $f(\mathbf{x}, \mathbf{w}) = x_1 w_1 + \dots + x_n w_n$. Ta pierwsza fala sieci neuronowych była znana jako cybernetyka, jak to pokazano na rysunku 1.7.

Model neuronów McCullocha-Pittsa (McCulloch and Pitts 1943) był pierwszym, liniowym modelem funkcji mózgu. Potrafił on rozróżnić dwa rodzaje danych wejściowych, sprawdzając czy wartość funkcji $f(\mathbf{x}, \mathbf{w})$ jest dodatnia, czy ujemna. Oczywiście, aby model odpowiadał pożądanym definicjom kategorii danych, wagi musiały być poprawnie ustawione. Mógł to robić człowiek. W lata pięćdziesiątych XX wieku pierwszym modelem, który potrafił uczyć się wag definiujących kategorie na podstawie przykładów danych w każdej z kategorii, był perceptron (Rosenblatt 1958, 1962). **Adaptacyjny liniowy element** (ADALINE) pochodzący z tego samego okresu po prostu zwracał wartość funkcji $f(\mathbf{x})$, aby przewidzieć liczbę rzeczywistą (Widrow and Hoff 1960), i potrafił nauczyć się przewidywać te wartości na podstawie danych.

Te proste algorytmy uczące się wywarły duży wpływ na dzisiejszy krajobraz systemów uczących się. Algorytm szkoleniowy używany do adaptacji wag w ADALINE stanowił szczególny przypadek algorytmu określanego jako algorytm **metodą gradientu stochastycznego**. Nieco zmodyfikowana wersja tego algorytmu dominuje w dzisiejszych modelach szkoleniowych w deep learningu.

Modele oparte na $f(\mathbf{x}, \mathbf{w})$, używane przez perceptron i ADALINE, są określane jako **modele liniowe**. Nadal są najszerzej używanymi modelami w systemach uczących się, choć w wielu przypadkach szkolenie przebiega w inny sposób, niż miało to miejsce w oryginalnych wersjach modelu.

Modele liniowe mają wiele ograniczeń. Najbardziej znanym jest fakt, że nie mogą nauczyć się funkcji XOR, gdzie $f([0, 1], \mathbf{w}) = 1$ i $f([1, 0], \mathbf{w}) = 1$, ale $f([1, 1], \mathbf{w}) = 0$ i $f([0, 0], \mathbf{w}) = 0$. Krytycy, którzy widzą te wady modeli liniowych, spowodowali ogólne odejście od uczenia się inspirowanego biologicznie (Minsky and Papert 1969). Było to pierwsze tąpnięcie popularności sieci neuronowych. Dziś neuronauka jest uznawana za ważne źródło inspiracji dla badaczy deep learningu, ale nie stanowi podstawowej wskazówki w tej dziedzinie.

Podstawową przyczyną mniejszej roli neuronauki w obecnych badaniach nad deep learningiem jest po prostu fakt, że nie mamy dość informacji o mózgu, które mogą stanowić wskazówkę. Aby dojść do głębokiego zrozumienia prawdziwych algorytmów używanych w mózgu, powinniśmy mieć możliwość jednoczesnego monitorowania aktywności tysięcy (przynajmniej jako minimum) połączonych ze sobą neuronów. Ponieważ nie możemy tego zrobić, jesteśmy nadal daleko od zrozumienia nawet najprostszych i najlepiej poznanych części mózgu (Olshausen and Field 2005).

Neuronauka dała nadzieję na to, że jeden algorytm deep learningu umie rozwiązać wiele różnych zadań. Neuronaukowcy odkryli, że fretki potrafią nauczyć się „widzieć” za pomocą obszaru mózgu przetwarzającego dźwięki, jeśli ich mózgi zostaną zmienione tak, aby sygnały wzrokowe były przesyłane do tego obszaru (Von Melchner et al. 2000). Sugeruje to, że duża część mózgu ssaków może korzystać z jednego algorytmu do rozwiązywania w mózgu różnych zadań. Zanim powstała ta hipoteza, badania nad systemami uczącymi się były podzielone, a różne społeczności naukowców zajmowały się przetwarzaniem języka naturalnego, wizją, planowaniem ruchu i rozpoznaniem mowy. Dziś te społeczności nadal pracują oddzielnie, ale wszystkie grupy badające deep learningu badają jednocześnie wiele lub wszystkie te dziedziny zastosowań.

Możemy nakreślić ogólne wskazówki na bazie neuronauki. Podstawowa zasada wielu jednostek obliczeniowych, które stają się inteligentne tylko dzięki wzajemnym połączeniom, jest inspirowana działaniem mózgu. *Neocognitron* (Fukushima 1980) wprowadza potężną architekturę modelu przetwarzania obrazów zainspirowaną strukturą systemu wzrokowego ssaków i stał się podstawą nowoczesnych sieci spłotowych (LeCun et al. 1998b), jak to pokazano w punkcie 9.10. Większość dzisiejszych sieci neuronowych opiera się na modelu znanym jako **poprawiona jednostka liniowa** (ang. *rectified linear unit*, ReLU). Oryginalny *cognitron* (Fukushima 1975) wprowadził bardziej skomplikowaną wersję w dużym stopniu inspirowaną naszą wiedzą o działaniu mózgu. Uproszczona współczesna wersja została opracowana dzięki połączeniu pomysłów z różnych punktów widzenia, przy czym Nair i Hinton

(2010) oraz Glorot et al. (2011a) podają neuronaukę jako źródło wpływu, a Jarrett et al. (2009) powołują się na wpływy bardziej inspirowane wiedzą inżynierską. Wprawdzie neuronauka jest ważnym źródłem inspiracji, jednak nie powinna być traktowana jako ściśle wskazówki. Wiemy, że pojedyncze neurony obliczają całkiem inne funkcje niż nowoczesne jednostki ReLU, lecz większy realizm w tym zakresie wcale nie doprowadza do lepszej wydajności systemów uczących się. Poza tym, choć neuronauka z powodzeniem zainspirowała szereg architektur sieci neuronowych, nadal nie mamy dość wiedzy o biologicznym uczeniu się, aby neuronauka dawała wiele wskazówek dla algorytmów uczących się używanych przez nas do szkolenia.

Media często podkreślają podobieństwo deep learningu do pracy mózgu. Wprawdzie prawdą jest, że pracujący nad deep learningiem podkreślają wpływ sposobu działania mózgu bardziej niż badacze w innych dziedzinach systemów uczących się, jak w przypadku jądra maszyn lub statystyk bayesowskich, ale nie należy widzieć deep learningu jako próby symulacji mózgu. Nowoczesny deep learning czerpie inspirację z wielu dziedzin, jak matematyka stosowana, w tym algebra liniowa, prawdopodobieństwo, a także teoria informacji i optymalizacja numeryczna. Niektórzy badacze deep learningu przytaczają neuronaukę jako istotne źródło inspiracji, inni wcale nie poruszają tego tematu.

Warto zauważyć, że wysiłek wkładany w zrozumienie, jak mózg pracuje na poziomie algorytmicznym, jest nadal żywy. Te starania są znane jako „komputerowa neuronauka” i stanowią oddzielne pole badań w stosunku do deep learningu. Badacze często przechodzą między tymi dwiema dziedzinami. Deep learning koncentruje się przede wszystkim na tym, jak zbudować systemy komputerowe, które potrafią z powodzeniem rozwiązywać zadania wymagające inteligencji, a komputerowa neuronauka zajmuje się budowaniem dokładniejszych modeli pracy mózgu.

W latach osiemdziesiątych XX wieku pojawiła się druga fala badań nad sieciami neuronowymi, głównie dzięki ruchowi zwanemu **koneksjonizmem** lub **równoległym przetwarzaniem rozproszonym** (Rumelhart et al. 1986c, McClelland et al. 1995). Koneksjonizm powstał w kontekście wiedzy kognitywnej. Jest to interdyscyplinarne podejście do zrozumienia umysłu, łączące wiele różnych poziomów analizy. We wczesnych latach osiemdziesiątych większość kognitywistów analizowała modele symbolicznego rozumowania. Mimo ich popularności modele symboliczne były trudne do wyjaśnienia w sensie zrozumienia, jak mózg może je wdrażać, używając neuronów. Koneksjonisci zaczęli analizować modele poznawcze, które mogły mieć źródło w implementacjach neuronowych (Touretzky and Minton 1985), ożywiając wiele pomysłów pochodzących z prac psychologa Donalda Hebba z lat czterdziestych (Hebb 1949).

Podstawowa idea koneksjonizmu stanowi, że duża liczba prostych elementów obliczeniowych może razem w sieci osiągnąć inteligentne zachowanie. To spostrzeżenie odnosi się zarówno do neuronów w systemach biologicznych, jak i do ukrytych jednostek w modelach obliczeniowych.

W czasach rozwoju ruchu koneksjonalistów z lat osiemdziesiątych pojawiło się kilka kluczowych pojęć, które do dziś pozostają w centrum deep learningu. Jednym z nich jest **rozproszona reprezentacja** (Hinton et al. 1986). Mówi ona, że każde dane wejściowe do systemu powinny być reprezentowane przez wiele cech, a każda z nich powinna być częścią reprezentacji wielu możliwych danych wejściowych. Przypuśćmy na przykład, że mamy system wizyjny, który potrafi rozróżniać samochody, ciężarówki i ptaki, a obiekty te mogą być czerwone, zielone lub niebieskie. Jednym ze sposobów prezentacji danych wejściowych może być oddzielny neuron lub ukryta jednostka aktywizująca każdą z dziewięciu możliwych kombinacji: czerwona ciężarówka, czerwony ptak, czerwony samochód, zielona ciężarówka i tak dalej. Wymaga to dziewięciu różnych neuronów, a każdy z nich musi niezależnie nauczyć się pojęcia koloru i tożsamości obiektu. Sposobem poprawy tej sytuacji jest skorzystanie z rozproszonej reprezentacji, gdzie trzy neurony opisują kolor, a trzy tożsamość obiektu. Wymaga to tylko sześciu neuronów zamiast dziewięciu, a neuron opisujący czerwień może dowiedzieć się o tym kolorze na podstawie obrazów samochodów, ciężarówek i ptaków, nie zaś z obrazów jednego typu obiektów. Pojęcie rozproszonej reprezentacji stanowi istotę tej książki i jest w szczegółach opisane w rozdziale 15.

Innym istotnym osiągnięciem ruchu koneksjonistycznego jest udane zastosowanie wstecznej propagacji w celu szkolenia głębokich sieci neuronowych za pomocą wewnętrznych reprezentacji i popularyzacja algorytmu propagacji wstecznej (Rumelhart et al. 1986a, LeCun 1987). Ten algorytm zyskiwał i tracił na popularności, ale w czasie pisania tej książki stanowi dominujące podejście do modeli deep learningu.

W latach dziewięćdziesiątych XX wieku badacze poczynili znaczące postępy w modelowaniu sekwencji z wykorzystaniem sieci neuronowych. Hochreiter (1991) i Bengio et al. (1994) rozwiązali niektóre podstawowe trudności matematyczne przy modelowaniu długich sekwencji, opisane w punkcie 10.7. Hochreiter i Schmidhuber (1997) wprowadzili sieć z długą pamięcią krótkoterminową (ang. *long short-term memory*, LSTM) do rozwiązywania niektórych z tych trudności. Dziś LSTM jest szeroko wykorzystywana w wielu sekwencyjnych zadaniach modelowania, w tym przy przetwarzaniu języka naturalnego w Google.

Druga fala badań nad sieciami neuronowymi trwała do połowy lat dziewięćdziesiątych. Działania opierające się na sieciach neuronowych i innych

technologiach AI zaczęły brnąć w nierealistyczne obietnice, szukając inwestorów. Gdy nie spełniły tych nierealnych oczekiwań, inwestorzy byli zawiedzeni. Jednocześnie w innych dziedzinach systemów uczących się dokonywano postępów. Maszyny oparte na jądrze (ang. *kernel machines*; Boser et al. 1992, Cortes and Vapnik 1995, Schölkopf et al. 1999) oraz modele graficzne (Jordan 1998) osiągnęły dobre wyniki w wielu ważnych zadaniach. Te dwa czynniki doprowadziły do spadku popularności sieci neuronowych, co trwało do roku 2007.

W tym czasie sieci neuronowe nadal pozwalały uzyskać świetną wydajność przy wykonywaniu niektórych zadań (LeCun et al. 1998b, Bengio et al. 2001). Kanadyjski instytut CIFAR (Canadian Institute for Advanced Research) pomógł w utrzymaniu badań nad sieciami neuronowymi poprzez swoją inicjatywę badawczą NCAP (Neural Computation and Adaptive Perception). Program ten zjednoczył zespoły badawcze pracujące nad systemami uczącymi się pod kierunkiem Geoffreya Hintona z Uniwersytetu w Toronto, Yoshuy Bengio z Uniwersytetu w Montrealu oraz Yanna LeCun z Uniwersytetu w Nowym Jorku. Multidyscyplinarne badania z inicjatywy CIFAR NCAP obejmowały neuronaukowców oraz ekspertów z zakresu widzenia ludzkiego i komputerowego.

Obecne głębokie sieci miały wówczas opinię trudnych do nauczenia. Teraz wiemy, że algorytmy dostępne od lat dziewięćdziesiątych działają dość dobrze, ale około roku 2006 nie było to oczywiste. Problem polega zapewne na tym, że algorytmy te są zbyt kosztowne obliczeniowo, aby można było pozwolić sobie na eksperymenty na dostępnym wówczas sprzęcie.

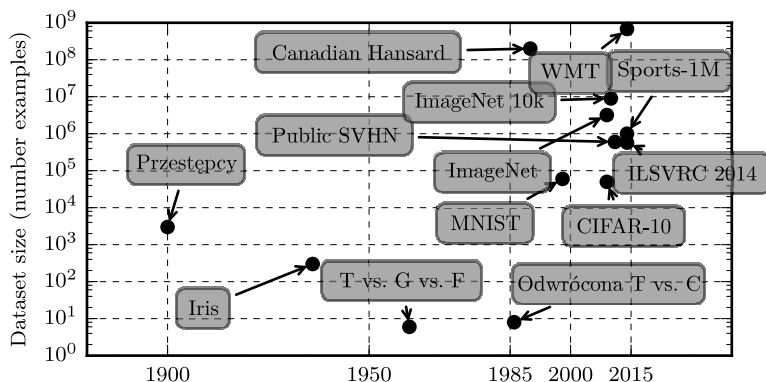
Trzecia fala badań nad sieciami neuronowymi zaczęła się od przełomu z roku 2006. Geoffrey Hinton pokazał rodzaj sieci neuronowej o nazwie *głęboka sieć przekonañ*, która mogła być szkolona za pomocą strategii zachłannego wstępnego szkolenia opartego na warstwach (Hinton et al. 2006), co szczegółowo jest omawiane w punkcie 15.1. Inna grupa badawcza powiązana z CIFAR szybko pokazała, że ta sama strategia może być wykorzystana do wielu innych rodzajów sieci (Bengio et al. 2007, Ranzato et al. 2007a) i systematycznie pomagała w poprawie uogólnienia jej na przykładach testowych. Ta fala badań nas sieciami neuronowymi spopularyzowała określenie „deep learning”, aby podkreślić, że badacze potrafią teraz głębiej szkolić sieci neuronowe, niż to było możliwe wcześniej, oraz aby skoncentrować uwagę na tym, jak teoretycznie ważna jest głębokość (Bengio and LeCun 2007, Delalleau i Bengio 2011, Pascanu et al. 2014a, Montufar et al. 2014). W tamtym czasie głębokie sieci neuronowe przewyższyły inne systemy AI, które opierają się na innych technikach uczenia się, a także opracowane ręcznie funkcje. Ta trzecia fala popularności sieci neuronowych trwa do dziś, choć główny cel badań nad deep learningiem znacznie zmieniły się z upływem czasu. Trzecia

fala rozpoczęła się od koncentracji na nowych nienadzorowanych technikach uczenia się i na zdolności głębokich modeli do generalizowania na podstawie małych zbiorów danych, lecz dziś nie ma już zainteresowania starszymi nadzorowanymi algorytmami uczenia się i zdolnością głębokich modeli do stosowania dużych etykietowanych zbiorów danych.

1.2.2. Rosnące rozmiary zbiorów danych

Można się zastanawiać, dlaczego deep learning dopiero ostatnio uznano za ważną technikę, choć pierwsze eksperymenty ze sztucznymi sieciami neuronowymi prowadzono już w latach pięćdziesiątych XX wieku. Deep learning był z powodzeniem używany w zastosowaniach komercyjnych od lat dziewięćdziesiątych, ale często uważany za sztukę, a nie technikę; twierdzono też, że do obsługi wymaga eksperta. To prawda, że dobra wydajność algorytmu wymaga posiadania umiejętności. Na szczęście ich zakres maleje wraz ze wzrostem ilości danych szkoleniowych. Algorytmy uczenia się, które z trudem rozwiązywały banalnie proste problemy w latach osiemdziesiątych, dzięki modelom szkoleniowym przeszły znaczące zmiany upraszczające szkolenia w przypadku bardzo głębokich architektur. Najważniejszą nowością jest fakt, że dziś istnieją algorytmy o zasobach wystarczających do odniesienia sukcesu. Na rysunku 1.8 pokazano, jak wielkość wzorcowych zbiorów danych rozrosła się z upływem czasu. Ten trend jest wynikiem rosnącej digitalizacji społeczeństwa. W miarę jak coraz więcej naszych działań odbywa się na komputerach, coraz więcej naszych działań jest rejestrowanych.

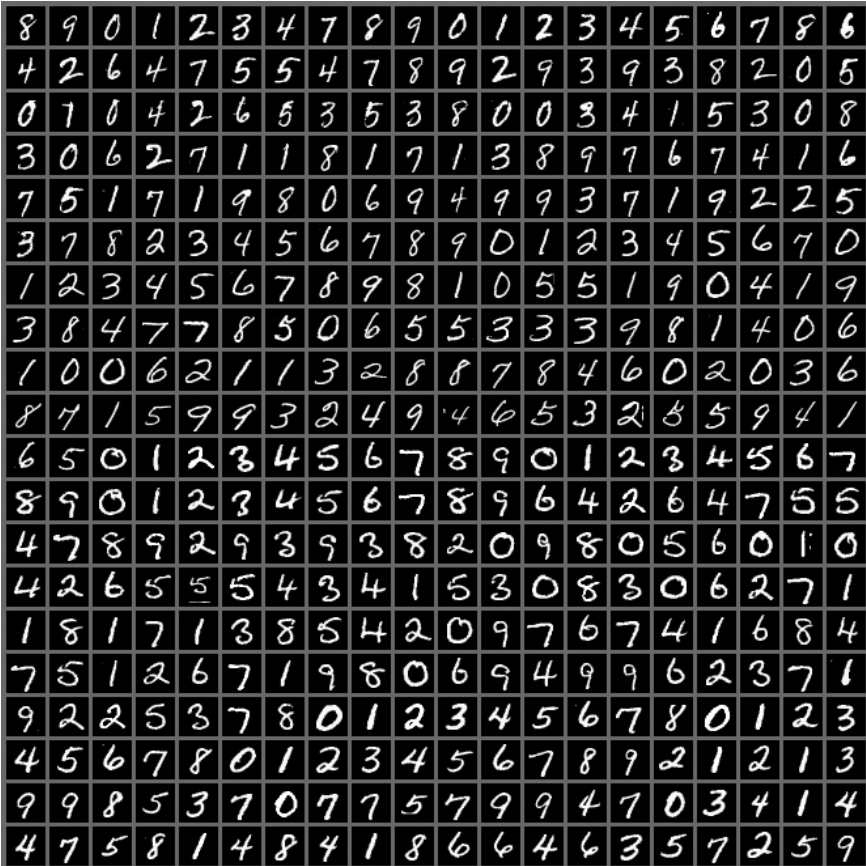
Nasze komputery są w coraz większym stopniu połączone z siecią, więc łatwiej scentralizować zapisy i dołączyć je do zbioru danych odpowiedniego dla aplikacji systemów uczących się. Czasy Big Data sprawiły, że systemy uczące się są prostsze, gdyż podstawowy ciężar estymacji statystycznej – uogólnianie nowych danych na podstawie obserwacji ich niewielkich ilości – stał się znacząco mniejszy. W roku 2016 zasadą jest, że nadzorowany algorytm deep learning osiąga akceptowalną wydajność przy około 5000 etykietowanych przykładach w danej kategorii i potrafi doścignąć lub przewyższyć wydajność ludzi, jeśli jest szkolony na podstawie nie mniej niż 10 milionów etykietowanych przykładów. Powodzenie w działaniu przy zbiorach danych mniejszych niż te wielkości stanowi ważny obszar badawczy, skupiając się w szczególności na możliwości korzystania z dużej liczby nieetykietowanych przykładów przy uczeniu się nienadzorowanym lub częściowo nadzorowanym.



Rysunek 1.8. Rosnące z upływem czasu zbiory danych. Na początku XX wieku statystycy analizowali zbiory danych za pomocą setek lub tysięcy ręcznie utworzonych miar (Garson 1900, Gosset 1908, Anderson 1935, Fisher 1936). Od 1950 roku do lat osiemdziesiątych XX wieku pionierzy systemów uczących się inspirowanych biologicznie często pracowali na małych syntetycznych zbiorach danych, jak mapy bitowe liter o niewielkiej rozdzielczości, które były projektowane tak, aby zmniejszyć koszty obliczeniowe i pokazać, że sieci neuronowe mogą nauczyć się określonych funkcji (Widrow and Ho 1960, Rumelhart et al. 1986b). W latach osiemdziesiątych i dziewięćdziesiątych XX wieku systemy uczące się zaczęły obejmować więcej statystyki i korzystać z większych zbiorów danych zawierających tysiące przykładów, jak na przykład zbiór danych MNIST (pokazany na rysunku 1.9) zawierający skany ręcznie pisanych liczb (LeCun et al. 1998b). W pierwszej dekadzie XXI wieku zaczęły być tworzone bardziej skomplikowane zbiory danych tych samych rozmiarów, jak CIFAR-10 (Krizhevsky i Hinton 2009). Pod koniec dekady i w pierwszej połowie drugiej dekady XXI wieku znacznie większe zbiory danych, zawierające setki tysięcy do dziesiątek milionów przykładów, całkowicie zmieniły możliwości deep learningu. Te zbiory danych obejmowały zbiory danych Street View House Numbers (Netzer et al. 2011), różne wersje zbioru danych ImageNet (Deng et al. 2009, 2010a; Russakovsky et al. 2014a) oraz zbiorów danych Sports-1M (Karpathy et al. 2014). Na górze wykresu widzimy, że zbiory danych przetłumaczonych zdań, jak zbiór IBM zbudowany na podstawie Canadian Hansard (Brown et al. 1990) oraz zbiory danych WMT 2014 English to French (Schwenk 2014), znacznie wyprzedzają wielkość innych zbiorów danych.

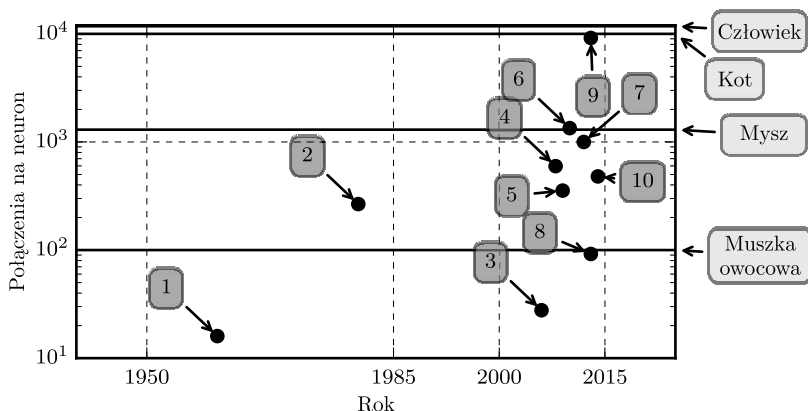
1.2.3. Rosnące rozmiary modeli

Innym ważnym powodem dzisiejszego sukcesu sieci neuronowych w porównaniu z umiarkowanym powodzeniem w latach dziewięćdziesiątych jest dostępność zasobów pozwalających na wykonywanie znacznie większych modeli. Jedną z głównych tez koneksjonizmu jest twierdzenie, że zwierzęta stają się inteligentne, gdy wiele z ich neuronów współpracuje ze sobą. Pojedynczy neuron lub niewielka ich grupa nie są szczególnie przydatne.



Rysunek 1.9. Przykład danych wejściowych ze zbioru danych MNIST. Skrót NIST oznacza National Institute of Standards and Technology (Narodowy Instytut Standardów i Technologii), instytucji, która zgromadziła te dane. Litera M pochodzi od słowa *machine* (maszyna), gdy dane zostały wstępnie przetworzone, aby łatwiej było je wykorzystać w algorytmach systemów uczących się. Zbiór danych MNIST składa się ze skanów ręcznie pisanych cyfr oraz związanych z nimi etykiet opisujących, która z cyfr od 0 do 9 znajduje się na obrazku. Ten prosty sposób klasyfikacji jest najprostszym i najczęściej używanym testem stosowanym w badaniach nad deep learningiem. Pozostaje popularny, mimo że jak na dzisiejsze techniki jest stosunkowo łatwy do rozwiązania. Geoffrey Hinton opisał go jako owocówkę systemów uczących się, co oznacza, że pozwala badaczom systemów uczących się na analizie ich algorytmów w kontrolowanych warunkach laboratoryjnych, podobnie jak biologzy badają muszki owocowe.

Neurony biologiczne nie mają zbyt gęstej sieci połączeń. Jak pokazano na rysunku 1.10, nasze modele systemów uczących się mają liczbę połączeń na neuron od długiego czasu o rząd wielkości większą niż mózgi ssaków.

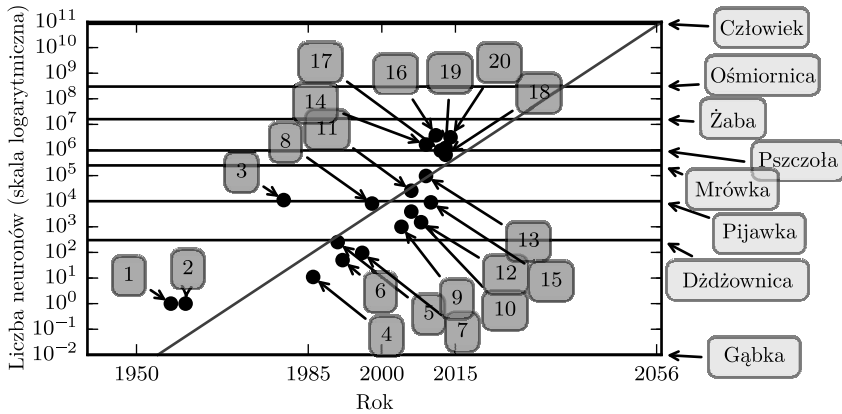


Rysunek 1.10. Liczba połączeń na neuron względem czasu. Dawniej liczby połączeń między neuronami w sztucznych sieciach neuronowych ograniczały możliwości sprżtów. Dziś liczba połączeń między neuronami wynika głównie z zaawansowanych projektowych. Niektóre sztuczne sieci neuronowe mają niemal tyle połączeń na neuron, ile mniejsze ssaki, np. myszy. Nawet ludzki mózg nie ma tak wielkiej liczby połączeń na neuron. Wielkości sieci biologicznych podano za Wikipedię (2015):

1. Adaptacyjny element liniowy (*linear element*, Widrow and Hoff 1960)
2. Neokognitron (Fukushima 1980)
3. Sieć spłotowa przyspieszana przez GPU (Chellapilla et al. 2006)
4. Głęboka maszyna Boltzmanna (Salakhutdinov and Hinton 2009a)
5. Nienadzorowana sieć spłotowa (Jarrett et al. 2009)
6. Wielowarstwowy perceptron przyspieszany przez GPU (Ciresan et al. 2010)
7. Rozproszony autokoder (Le et al. 2012)
8. Sieć spłotowa z wieloma GPU (Krizhevsky et al. 2012)
9. Nienadzorowana sieć spłotowa COTS HPC (Coates et al. 2013)
10. GoogLeNet (Szegedy et al. 2014a)

W kontekście liczby neuronów sieci neuronowe do niedawna były niezwykle małe. Zmieniło się to dopiero ostatnio, co pokazuje rysunek 1.11, gdzie widać rozrastanie się sieci neuronowych w funkcji czasu.

Od czasu wprowadzenia ukrytych jednostek sztuczne sieci neuronowe podwajają swoje rozmiary średnio co 2,4 roku. Ten wzrost wynika z użycia szybszych komputerów z większą pamięcią oraz dostępu do większych zbiorów danych. Większe sieci mogą osiągać większą dokładność bardziej złożonych zadań. Wygląda na to, że ten trend utrzymał się na lata. Jeśli nowe techniki nie pozwolą na szybsze skalowanie, sztuczne sieci neuronowe nie będą miały tylu neuronów, jak ludzki mózg, co najmniej do lat pięćdziesiątych XXI wieku. Neurony biologiczne mogą reprezentować bardziej skomplikowane funkcje niż obecne neurony sztuczne, więc biologiczne sieci neuronowe mogą być nawet większe, niż to wynika z rysunku.



Rysunek 1.11. Wzrost rozmiaru sieci neuronowych względem czasu. Od czasu wprowadzenia ukrytych jednostek sztuczne sieci neuronowe podwajają swoje rozmiary mniej więcej co 2,4 roku. Oto rozmiary sieci neuronowych wg Wikipedii (2015):

1. Perceptron (Rosenblatt 1958, 1962)
2. Adaptacyjny element liniowy (*linear element*, Widrow and Hoff 1960)
3. Neocognitron (Fukushima 1980)
4. Wczesne sieci propagacji wstecznej (Rumelhart et al. 1986b)
5. Rekurencyjna sieć neuronowa rozpoznawania mowy (Robinson and Fallside 1991)
6. Wielowarstwowy perceptron rozpoznawania mowy (Bengio et al. 1991)
7. Sigmoidalna sieć średniego pola przekonań (Saul et al. 1996)
8. LeNet-5 (LeCun et al. 1998b)
9. Sieć stanu echa (Jaeger and Haas 2004)
10. Sieć głębokich przekonań (Hinton et al. 2006)
11. Sieć spłotowa przyspieszana przez GPU (Chellapilla et al. 2006)
12. Głęboka maszyna Boltzmanna (Salakhutdinov and Hinton 2009a)
13. Głęboka sieć przekonań przyspieszana przez GPU (Raina et al. 2009)
14. Nienadzorowana sieć spłotowa (Jarrett et al. 2009)
15. Wielowarstwowy perceptron przyspieszany przez GPU (Ciresan et al. 2010)
16. Sieć OMP-1 (Coates and Ng 2011)
17. Rozproszony autokoder (Le et al. 2012)
18. Sieć spłotowa z wieloma GPU (Krizhevsky et al. 2012)
19. Nienadzorowana sieć spłotowa COTS HPC (Coates et al. 2013)
20. GoogLeNet (Szegedy et al. 2014a).

Patrząc retrospektywnie, możemy uznać, że nie jest szczególnie zaskakujące, iż sieci neuronowe, mające mniej neuronów niż pijawki, nie potrafiły rozwiązywać skomplikowanych problemów sztucznej inteligencji. Nawet dzisiejsze sieci, dość duże z obliczeniowego punktu widzenia, mają system nerwowy mniejszy niż względnie prymitywne zwierzęta, takie jak żaba.

Powiększanie się modeli z upływem czasu, dostępność szybszych procesorów, nadejście GPU ogólnego zastosowania (opisane w punkcie 12.1.2), szybsze połączenia sieciowe i lepsza infrastruktura programowa do celów

rozproszonych obliczeń – to najważniejsze trendy w historii deep learningu. Spodziewamy się, że taki trend utrzyma się w przyszłości.

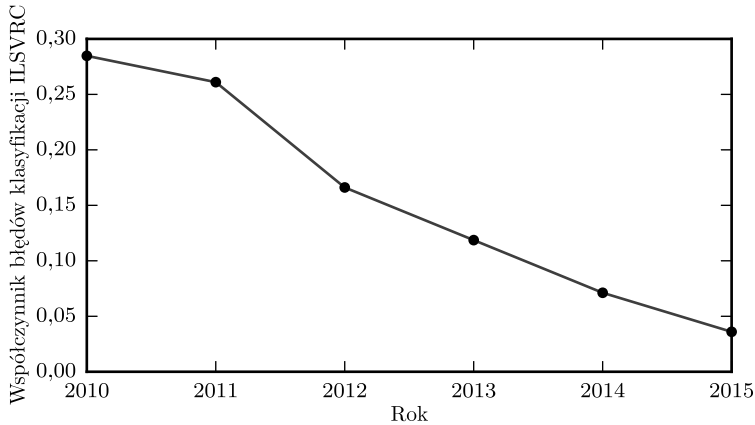
1.2.4. Rosnąca dokładność, złożoność i wpływ świata rzeczywistego

Od lat osiemdziesiątych XX wieku deep learning wciąż poprawia swoją zdolność do dokładnego rozpoznawania i przewidywania. Co więcej, jest on stale w powodzeniem stosowany do coraz szerszego spektrum zastosowań.

Najwcześniejsze modele deep learningu były używane do rozpoznawania pojedynczych obiektów na dobrze przyciętych bardzo małych obrazach (Rumelhart et al. 1986a). Od tego czasu następuje stały wzrost wielkości obrazów, które mogą być przetwarzane przez sieci neuronowe. Nowoczesne sieci rozpoznawania obiektów przetwarzają pełne fotografie w wysokiej rozdzielczości i nie wymagają, aby zdjęcie było przycięte wokół rozpoznawanego obiektu (Krizhevsky et al. 2012). Pierwsze sieci mogły rozpoznawać tylko dwa rodzaje obiektów (lub w niektórych przypadkach brak jakiegoś typu obiektu), a nowoczesne sieci zwykle rozpoznają co najmniej 1000 różnych rodzajów obiektów. Największy konkurs rozpoznawania obiektów to odbywający się co roku ImageNet Large Scale Visual Recognition Challenge (ILSVRC). Momentem błyskawicznego wzrostu deep learningu był dzień, gdy sieć splotowa po raz pierwszy wyraźnie wygrała konkurs, zmniejszając współczynnik 5 górnych błędów z 26,1% na 15,3% (Krizhevsky et al. 2012), co oznaczało, że ta sieć splotowa tworzy listę rankingową możliwych kategorii dla każdego obrazu, a poprawna kategoria pojawia się wśród pierwszych pięciu pozycji we wszystkich testach, poza 15,3%. Od tego czasu te konkursy są stale wygrywane przez głębokie sieci splotowe, a w momencie pisania tej książki postępy w rozwoju deep learningu sprowadziły procent błędnych odpowiedzi w górnych pięciu do 3,5%, jak to pokazano na rysunku 1.12.

Deep learning ma również ogromny wpływ na rozpoznawanie mowy. Po znacznej poprawie w latach dziewięćdziesiątych współczynnik błędów w rozpoznawaniu mowy zatrzymał się około roku 2000. Wprowadzenie deep learningu (Dahl et al. 2010, Deng et al. 2010b, Seide et al. 2011, Hinton et al. 2012a) do procesu rozpoznawania mowy sprawiło, że niektóre współczynniki błędów spadły o połowę. Ten temat jest bardziej szczegółowo omówiony w punkcie 12.3.

Głębokie sieci odniosły też spektakularny sukces przy wykrywaniu pieszych oraz segmentacji obrazów (Sermanet et al. 2013, Farabet et al. 2013, Couprie et al. 2013) i osiągnęły nadludzką wydajność w rozpoznawaniu znaków drogowych.



Rysunek 1.12. Zmniejszanie współczynnika błędów względem czasu. Od czasu, gdy głębokie sieci osiągnęły skalę potrzebną do konkurencji w ImageNet Large Scale Visual Recognition Challenge, co roku wygrywają ten konkurs, z każdym razem osiągając niższy współczynnik błędów. Dane za Russakovsky et al. 2014b i He et al. 2015

W czasie gdy wzrosła skala i dokładność głębokich sieci, wzrosła także złożoność zadań, które potrafią one rozwiązać. Goodfellow et al. (2014d) pokazali, że sieci neuronowe mogą nauczyć się wyprowadzać całe sekwencje znaków pobranych z obrazu, zamiast po prostu identyfikować poszczególne obiekty. Poprzednio powszechnie sądzono, że tego rodzaju uczenie się wymaga etykietowania poszczególnych elementów sekwencji (Gülçehre and Bengio 2013). Rekurencyjne sieci neuronowe, jak wspomniany wcześniej model sekwencyjny LSTM, są teraz używane do modelowania związków pomiędzy sekwencjami oraz innymi sekwencjami, o nieustalonych danych wejściowych. Uczenie się wedle zasady od sekwencji do sekwencji zapewne zrewolucjonizuje kolejne zastosowanie: tłumaczenia maszynowe (Bahdanau et al. 2015).

Trend rosnącej złożoności doprowadził w konsekwencji do wprowadzenia neuronowych maszyn Turinga (Graves et al. 2014a), które uczą się czytać z komórek pamięci i zapisują arbitralną zawartość do tych komórek. Takie sieci neuronowe mogą nauczyć się prostych programów, na przykładach pożądanego zachowania. Mogą nauczyć się sortować listy liczb, mając jako dane sekwencje danych pomieszanych i uporządkowanych. Ta technika samoprogramowania dopiero się rodzi, ale w przyszłości może być w zasadzie zastosowana do niemal każdego zadania.

Innym osiągnięciem deep learningu jest jego rozwinięcie na dziedzinę **wspomagane uczenia się**. W kontekście wzmocnianego uczenia się autonomiczny agent musi nauczyć się wykonywać zadanie metodą prób i błędów, bez żadnej pomocy ze strony człowieka. DeepMind pokazało, że system

wzmacnianego uczenia się oparty na deep learningu potrafi nauczyć się grać w gry wideo na Atari, osiągając w wielu zadaniach wyniki na poziomie człowieka (Mnih et al. 2015). Deep learning znacząco poprawił też wydajność wzmacnianego uczenia się w dziedzinie robotyki (Finn et al. 2015).

Wiele z tych zastosowań deep learningu przynosi wysokie zyski. Jest on obecnie wykorzystywany przez czołowe firmy, w tym Google, Microsoft, Facebook, IBM, Baidu, Apple, Adobe, Netflix, NVIDIA i NEC.

Postępy w dziedzinie deep learningu zależą w dużym stopniu od postępów w dziedzinie infrastruktury programistycznej. Biblioteki programów, jak Theano (Bergstra et al. 2010, Bastien et al. 2012), PyLearn2 (Goodfellow et al. 2013c), Torch (Collobert et al. 2011b), DistBelief (Dean et al. 2012), Caffe (Jia 2013), MXNet (Chen et al. 2015) i TensorFlow (Abadi et al. 2015), zawsze wspomagały ważne projekty badawcze, a także produkty komercyjne.

Deep learning ma wkład także w inne dziedziny nauki. Wiele splotowych sieci rozpoznawania obrazów zapewnia modelowanie procesów widzenia, które mogą badać neuronaukowcy (DiCarlo 2013). Deep learning zapewnia także użyteczne narzędzia do przetwarzania ogromnej ilości danych, pozwalając na użyteczne przewidywania w wielu dziedzinach nauki. Z powodzeniem wykorzystano je do przewidywania, jaki wpływ będą mieć na siebie molekuly, co pomaga firmom farmaceutycznym na tworzenie nowych leków (Dahl et al. 2014), poszukiwanie cząstek atomowych (Baldi et al. 2014) i automatyzowanie analizy obrazów mikroskopowych w celu utworzenia trójwymiarowej mapy ludzkiego mózgu (Knowles-Barley et al. 2014). Spodziewamy się, że deep learning pojawi się w przyszłości w wielu kolejnych dziedzinach nauki.

Podsumowując, deep learning, rozwijany w ciągu kilku ostatnich dekad, stanowi podejście do systemów uczących się, które wpłynęło mocno na naszą wiedzę o mózgu człowieka, statystykę i matematykę stosowaną. W ostatnich latach deep learning znacznie zyskał na popularności i zastosowaniach, głównie ze względu na komputery o większej mocy, większe zbiory danych oraz techniki szkolenia głębszych sieci. Nadchodzące lata są pełne wyzwań i możliwości, aby dalej poprawiać deep learning, przekraczając nowe granice.

I

Podstawy matematyki
stosowanej i systemów
uczących się

Ta część książki stanowi wprowadzenie do podstawowych pojęć matematycznych potrzebnych do zrozumienia deep learningu. Zaczyna się od ogólnych pojęć matematyki stosowanej, które pozwolą zdefiniować funkcje wielu zmiennych, znaleźć maksima i minima tych funkcji oraz określić stopień przekonania. Następnie opisywane są podstawowe cele systemów uczących się. W tej części wskazuje się, jak osiągnąć te cele, opisując model reprezentujący określone przekonania, projektując funkcje kosztów, które sprawdzają, jak te przekonania korespondują z rzeczywistością, oraz używając algorytmu szkoleniowego pozwalającego zminimalizować funkcję kosztów. Ta podstawowa struktura stanowi podstawę wielu algorytmów stosowanych w systemach uczących się, w tym podejścia, które nie jest głębokie. W kolejnych częściach książki w ramach tej struktury rozwiniemy algorytmy deep learningu.

2

Algebra liniowa

Algebra liniowa to dziedzina matematyki, która jest szeroko stosowana w nauce i inżynierii. Jest to matematyka ciągła, a nie dyskretna, więc wielu naukowców zajmujących się informatyką niewiele o niej wie. Dobre zrozumienie algebry liniowej jest konieczne do zrozumienia i wykorzystania wielu algorytmów systemów uczących się, w tym algorytmów deep learningu. Dlatego wprowadzenie do deep learningu zostało poprzedzone skoncentrowaną prezentacją kluczowych pojęć algebry liniowej.

Osoby znające algebrę liniową mogą pominąć ten rozdział. Czytelnikom, którzy mają doświadczenie i znają pojęcia, ale potrzebują przypomnienia kluczowych wzorów, warto polecić podręcznik *The Matrix Cookbook* (Petersen and Pedersen 2006). Tym, którzy wcale nie znają algebry liniowej, ten rozdział wystarczy, aby przeczytać resztę książki, ale gorąco zaleca się sięgnięcie również po inne źródła z zakresu nauki algebry liniowej, jak Shilov (1977). W tym rozdziale pominięto wiele ważnych tematów algebry liniowej, które nie są konieczne do zrozumienia deep learningu.

2.1. Skalary, wektory, macierze i tensory

Nauka algebry liniowej obejmuje kilka typów obiektów matematycznych.

- **Skalary.** Skalar to jedna liczba, co odróżnia go od innych obiektów algebry liniowej, które są zwykle tablicami złożonymi z wielu liczb. Skalary zapisujemy tu kursywą. Zwykle nadaje się im nazwy rozpoczynające się małą literą. Gdy są wprowadzane, podajemy, jakiego typu są liczbą. Na przykład można powiedzieć: „Niech $\mathbf{s} \in \mathbb{R}$ stanowi nachylenie prostej”, definiując skalar rzeczywisty, lub: „Niech $n \in \mathbb{N}$ będzie liczbą jednostek” – gdy definiujemy liczbę naturalną.

- **Wektory.** Wektor to tablica liczb. Liczby są ustawione po kolei. Można zidentyfikować kolejne liczby za pomocą indeksu wskazującego położenie. Zwykle nadaje się wektorom nazwy pisane małą literą, pisząc je pogrubioną czcionką, jak \mathbf{x} . Elementy wektora identyfikujemy przez nazwę pisaną kursywą z indeksem. Pierwszy element wektora \mathbf{x} to x_1 , drugi to x_2 i tak dalej. Musimy też określić, jakiego typu liczby będą przechowywane w wektorze. Jeśli każdy element należy do \mathbb{R} , a wektor ma n elementów, to leży on w zbiorze powstałym z iloczynu kartezjańskiego $\mathbb{R}n$ razy, oznaczonym jako \mathbb{R}^n . Gdy musimy szybko określić elementy wektora, piszemy je w postaci kolumny ujętej w nawiasy kwadratowe:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}. \quad (2.1)$$

Traktujemy wektory jako identyfikację punktów w przestrzeni, przy czym każdy podaje współrzędną na innej osi.

Czasami trzeba indeksować zbiór elementów wektora. W takim przypadku definiuje się zbiór zawierający indeksy i zapisuje zbiór jako indeks dolny. Na przykład, aby uzyskać dostęp do x_1, x_3 i x_6 , definiuje się zbiór $S = \{1, 3, 6\}$ i pisze \mathbf{x}_S . Używa się znaku $-$, aby indeksować dopełnienie tego zbioru. Na przykład \mathbf{x}_{-1} to wektor zawierający wszystkie elementy \mathbf{x} , poza elementem x_1 , natomiast \mathbf{x}_{-S} to wektor zawierający elementy \mathbf{x} poza x_1, x_3 i x_6 .

- **Macierze.** Macierz to dwuwymiarowa tablica liczb, więc każdy element jest identyfikowany przez dwa indeksy. Zwykle nazwy macierzy pisze się wielką literą z pogrubieniem, jak na przykład \mathbf{A} . Jeśli macierz rzeczywista \mathbf{A} ma wysokość m i szerokość n , to mówimy, że $\mathbf{A} \in \mathbb{R}^{n \times m}$. Zwykle identyfikuje się elementy macierzy, pisząc jej nazwę kursywą, a nie czcionką pogrubioną, a indeksy oddziela przecinkami. Na przykład $A_{1,1}$ to lewy górny element \mathbf{A} , zaś $A_{m,n}$ to prawy dolny element. Wszystkie liczby o współrzędnej pionowej i zapisuje się jako “:” dla współrzędnej poziomej. Na przykład $A_{i,:}$ oznacza poprzeczny przekrój \mathbf{A} ze współrzędną pionową i . Jest to określane jako i -ty wiersz macierzy \mathbf{A} . Podobnie $A_{:,i}$ to i -ta kolumna macierzy \mathbf{A} . Gdy chce się podać wszystkie elementy macierzy, zapisuje się je jako tablicę ujętą w nawiasy kwadratowe:

$$\begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}. \quad (2.2)$$

Czasami trzeba indeksować wyrażenia macierzowe, gdzie nie pisze się pojedynczej litery. W takim przypadku używane są indeksy dolne po wyrażeniu, ale nie zmienia się liter na małe. Na przykład $f(\mathbf{A})_{i,j}$ daje element (i, j) macierzy obliczony po zastosowaniu na niej funkcji f .

- **Tensory.** W niektórych przypadkach potrzebne są tablice o większej liczbie wymiarów. W ogólnym przypadku liczby tablicy tworzą regularną siatkę ze zmienną liczbą osi, którą nazywamy tensorem. Oznaczamy tensor o nazwie „A” za pomocą innej czcionki: \mathbf{A} . Element \mathbf{A} o współrzędnych (i, j, k) zapisujemy jako $A_{i,j,k}$.

Istotnym działaniem macierzowym jest **transpozycja**. Jest to lustrzane odbicie macierzy względem jej przekątnej, określanej jako **główna przekątna**, przebiegającej w dół i na prawo począwszy od lewego górnego narożnika. Na rysunku 2.1 widać graficzny obraz tego działania.

$$\mathbf{A} = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \\ A_{3,1} & A_{3,2} \end{bmatrix} \Rightarrow \mathbf{A}^\top = \begin{bmatrix} A_{1,1} & A_{2,1} & A_{3,1} \\ A_{1,2} & A_{2,2} & A_{3,2} \end{bmatrix}$$

Rysunek 2.1. Transpozycja macierzy to jej lustrzane odbicie według jej głównej przekątnej

Transpozycję macierzy \mathbf{A} oznaczamy jako \mathbf{A}^\top i zgodnie z definicją ma ona postać:

$$(\mathbf{A}^\top)_{i,j} = A_{j,i}. \quad (2.3)$$

Wektory można traktować jak macierze zawierające tylko jedną kolumnę. Transpozycja wektora przekształca go więc w macierz o jednym wierszu. Czasami definiuje się wektor, wypisując jego elementy w wierszu, jako macierz wierszową, a wtedy transpozycja przekształca ją na standardowy wektor kolumnowy, na przykład $\mathbf{x} = [x_1, x_2, x_3]^\top$.

Skalar można traktować jako macierz o jednym elemencie. Z tego wynika, że skalar stanowi swoją własną transpozycję: $a = a^\top$.

Macierze można do siebie dodawać, jeśli mają ten sam kształt, dodając ich odpowiadające sobie elementy: $\mathbf{C} = \mathbf{A} + \mathbf{B}$, gdzie $C_{i,j} = A_{i,j} + B_{i,j}$.

Można też dodać do macierzy skalar lub pomnożyć macierz przez skalar, wykonując odpowiednie działanie na każdym elemencie macierzy: $\mathbf{D} = a \cdot \mathbf{B} + c$, gdzie $D_{i,j} = aB_{i,j} + c$.

W kontekście deep learningu używa się także mniej konwencjonalnej notacji. Pozwalamy na dodanie macierzy i wektora, otrzymując nową macierz:

$\mathbf{C} = \mathbf{A} + \mathbf{b}$, gdzie $C_{i,j} = A_{i,j} + b_j$. Innymi słowy, wektor \mathbf{b} zostaje dodany do każdego wiersza macierzy. Taki skrót eliminuje potrzebę definiowania macierzy, w której przed dodawaniem \mathbf{b} jest kopiowany do każdego jej wiersza. Takie domyślne kopiowanie \mathbf{b} do wielu lokalizacji nazywa się **rozgłaszaniem**.

2.2. Mnożenie macierzy i wektorów

Jednym z najważniejszych działań związanych z macierzami jest mnożenie dwóch macierzy. **Iloczynem macierzy \mathbf{A} i \mathbf{B}** jest trzecia macierz \mathbf{C} . Aby można było uzyskać iloczyn, macierz \mathbf{A} musi mieć tyle samo kolumn, ile wierszy ma macierz \mathbf{B} . Jeśli \mathbf{A} ma wymiary m na n , a \mathbf{B} n na p , to \mathbf{C} będzie miała rozmiar m na p . Iloczyn macierzy możemy zapisać, umieszczając razem dwie macierze lub więcej, na przykład:

$$\mathbf{C} = \mathbf{AB}. \quad (2.4)$$

Iloczyn definiowany jest jako:

$$C_{i,j} = \sum_k A_{i,k} B_{k,j}. \quad (2.5)$$

Zauważmy, że standardowy iloczyn dwóch macierzy to nie jest po prostu macierz zawierająca iloczyny poszczególnych elementów. Takie działanie też istnieje i nosi nazwę **iloczynu po współrzędnych** lub **iloczynu Hadamarda** i jest oznaczane jako $\mathbf{A} \odot \mathbf{B}$.

Iloczyn skalarny (ang. *dot product*) między dwoma wektorami \mathbf{x} i \mathbf{y} o tych samych wymiarach jest iloczynem macierzy $\mathbf{x}^\top \mathbf{y}$. Możemy traktować iloczyn macierzy w postaci $\mathbf{C} = \mathbf{AB}$ jako wyznaczenie $C_{i,j}$ jako iloczynu skalarnego między wierszem macierzy \mathbf{A} i kolumną macierzy \mathbf{B} .

Iloczyny macierzy mają wiele użytecznych własności, które sprawiają, że analiza matematyczna macierzy jest wygodniejsza. Na przykład mnożenie macierzy jest rozdzielne:

$$\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}. \quad (2.6)$$

Jest także łączne:

$$\mathbf{A}(\mathbf{BC}) = (\mathbf{AB})\mathbf{C}. \quad (2.7)$$

W przeciwieństwie do mnożenia skalarnego mnożenie macierzy nie jest przemienne (warunek $\mathbf{AB} = \mathbf{BA}$ nie zawsze jest spełniony). Jednak iloczyn z kropką dwóch wektorów jest przemienny:

$$\mathbf{x}^\top \mathbf{y} = \mathbf{y}^\top \mathbf{x}. \quad (2.8)$$

Transpozycja iloczynu macierzy ma prostą postać:

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top. \quad (2.9)$$

To pozwala nam na pokazanie równania 2.8, dzięki wykorzystaniu faktu, że wartość takiego iloczynu jest skalar, a więc jest równa swojej transpozycji:

$$\mathbf{x}^\top \mathbf{y} = (\mathbf{x}^\top \mathbf{y})^\top = \mathbf{y}^\top \mathbf{x}. \quad (2.10)$$

Ponieważ przedmiotem tej książki nie jest algebra liniowa, nie próbujemy tu podać wyczerpującej listy użytecznych własności iloczynu macierzy. Jednak czytelnicy powinni być świadomi, że istnieje o wiele więcej własności.

Wiemy już dość na temat zapisu algebry liniowej, aby zapisać układ równań liniowych:

$$\mathbf{Ax} = \mathbf{b}, \quad (2.11)$$

gdzie $\mathbf{A} \in \mathbb{R}^{m \times n}$ to znana macierz, $\mathbf{b} \in \mathbb{R}^m$ to znany wektor, a $\mathbf{x} \in \mathbb{R}^n$ jest wektorem nieznanymi zmiennymi, które chcemy znaleźć. Każdy element x_i wektora \mathbf{x} jest jedną z tych nieznanymi zmiennymi. Równanie 2.11 można zapisać jako:

$$\mathbf{A}_{1,:} \mathbf{x} = b_1 \quad (2.12)$$

$$\mathbf{A}_{2,:} \mathbf{x} = b_2 \quad (2.13)$$

$$\dots \quad (2.14)$$

$$\mathbf{A}_{m,:} \mathbf{x} = b_m \quad (2.15)$$

albo w bardziej jawnej postaci:

$$\mathbf{A}_{1,1}x_1 + \mathbf{A}_{1,2}x_2 + \dots + \mathbf{A}_{1,n}x_n = b_1 \quad (2.16)$$

$$\mathbf{A}_{2,1}x_1 + \mathbf{A}_{2,2}x_2 + \dots + \mathbf{A}_{2,n}x_n = b_2 \quad (2.17)$$

$$\dots \quad (2.18)$$

$$\mathbf{A}_{m,1}x_1 + \mathbf{A}_{m,2}x_2 + \dots + \mathbf{A}_{m,n}x_n = b_m. \quad (2.19)$$

Notacja iloczynu macierzowo-wektorowa pozwala na bardziej zwięzłą reprezentację równań.

2.3. Macierze jednostkowe i odwrotne

Algebra liniowa oferuje potężne narzędzie nazywane **odwracaniem macierzy**, które pozwala analitycznie rozwiązać równanie 2.11 dla wielu wartości \mathbf{A} . Aby opisać odwracanie macierzy, musimy najpierw zdefiniować **macierze jednostkowe**. Macierz jednostkowa to macierz, która nie zmienia żadnego wektora, gdy mnożymy go przez tę macierz. Macierz tę, która zachowuje wartość wektorów n -wymiarowych, oznaczamy jako \mathbf{I}_n . Formalnie $\mathbf{I}_n \in \mathbb{R}$ oraz:

$$\forall \mathbf{x} \in \mathbb{R}^n, \mathbf{I}_n \mathbf{x} = \mathbf{x}. \quad (2.20)$$

Struktura macierzy jednostkowej jest prosta: wszystkie elementy na głównej przekątnej mają wartość 1, a pozostałe mają wartość 0. Przykład pokazano na rysunku 2.2.

Macierz odwrotną do macierzy \mathbf{A} oznaczamy jako \mathbf{A}^{-1} i definiujemy jako:

$$\mathbf{A}^{-1} \mathbf{A} = \mathbf{I}_n. \quad (2.21)$$

Równanie 2.11 możemy rozwiązać w pokazanych niżej krokach:

$$\mathbf{A} \mathbf{x} = \mathbf{b} \quad (2.22)$$

$$\mathbf{A}^{-1} \mathbf{A} \mathbf{x} = \mathbf{A}^{-1} \mathbf{b} \quad (2.23)$$

$$\mathbf{I}_n \mathbf{x} = \mathbf{A}^{-1} \mathbf{b} \quad (2.24)$$

$$\mathbf{x} = \mathbf{A}^{-1} \mathbf{b}. \quad (2.25)$$

Oczywiście proces ten zależy od możliwości znalezienia \mathbf{A}^{-1} . Warunki istnienia macierzy odwrotnej opisujemy w kolejnym punkcie.

Gdy istnieje macierz odwrotna \mathbf{A}^{-1} , możemy ją wyznaczyć, stosując jeden z kilku istniejących algorytmów. W teorii ta sama macierz odwrotna może zostać użyta do wielokrotnego rozwiązania równania dla różnych wartości \mathbf{b} . Jednak \mathbf{A}^{-1} jest użyteczna jako narzędzie teoretyczne i w praktyce nie powinna być używana w większości aplikacji programistycznych. Ponieważ macierz odwrotna może być wyznaczona na komputerach tylko z określoną dokładnością, algorytmy wykorzystujące wartość \mathbf{b} pozwalają otrzymać bardziej dokładne estymację \mathbf{x} .

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Rysunek 2.2. Przykład macierzy jednostkowej \mathbf{I}_3

2.4. Zależność liniowa i zakres

Aby istniała macierz \mathbf{A}^{-1} , równanie 2.11 musi mieć dokładnie jedno rozwiązanie dla każdej wartości \mathbf{b} . Możliwe jest, że układ równań nie ma żadnego rozwiązania lub ma nieskończenie wiele rozwiązań dla określonej wartości \mathbf{b} . Nie jest jednak możliwe, aby dla konkretnej wartości \mathbf{b} było więcej niż jedno rozwiązanie, a mniej niż nieskończenie wiele. Jeśli zarówno \mathbf{x} , jak i \mathbf{y} są rozwiązaniami, to:

$$\mathbf{z} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \quad (2.26)$$

jest także rozwiązaniem dla dowolnej wartości rzeczywistej α .

Aby stwierdzić, ile rozwiązań ma równanie, można potraktować kolumny \mathbf{A} jako różne kierunki podróży z **początku** (punktu określonego przez wektor zawierający same zera), a potem określić, ile jest dróg prowadzących do \mathbf{b} . Z tego punktu widzenia każdy element \mathbf{x} określa, jak daleko mamy podróżować w każdym z kierunków, a x_i określa, jak daleko mamy się poruszać w kierunku kolumny i :

$$\mathbf{Ax} = \sum_i x_i \mathbf{A}_{:,i}. \quad (2.27)$$

Taki rodzaj działania nosi nazwę **kombinacji liniowej**. Formalnie kombinacja liniowa pewnego zbioru wektorów $\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}\}$ jest otrzymywana przez pomnożenie każdego wektora $\mathbf{v}^{(i)}$ przez odpowiadający mu współczynnik skalarny i dodanie do siebie wyników:

$$\sum_i c_i \mathbf{v}^{(i)}. \quad (2.28)$$

Zakres (ang. *span*) zbioru wektorów to zbiór wszystkich punktów, które można otrzymać na podstawie kombinacji liniowej oryginalnych wektorów.

Określenie, czy $\mathbf{Ax} = \mathbf{b}$ ma rozwiązanie, sprowadza się więc do sprawdzenia, czy \mathbf{b} znajduje się w zakresie kolumn \mathbf{A} . Zakres ten określamy jako **przestrzeń kolumny** lub **zasięg \mathbf{A}** .

Aby układ $\mathbf{Ax} = \mathbf{b}$ miał rozwiązanie dla wszystkich wartości $\mathbf{b} \in \mathbb{R}^m$, musimy wymagać, aby przestrzeń kolumn \mathbf{A} w całości należała do \mathbb{R}^m . Jeśli jakiś punkt w \mathbb{R}^m jest poza przestrzenią kolumn, stanowi on potencjalnie wartość \mathbf{b} , dla której nie ma rozwiązania. Z wymagania, aby przestrzeń kolumn \mathbf{A} należała w całości do \mathbb{R}^m , wynika, że \mathbf{A} musi mieć co najmniej m kolumn, czyli $n \geq m$. W przeciwnym przypadku wymiarowość przestrzeni kolumn będzie mniejszy od m . Jako przykład rozpatrzmy przypadek macierzy 3 na 2. Docelowe \mathbf{b} jest trójwymiarowe, lecz \mathbf{x} jest tylko dwuwymiarowe, więc

modyfikowanie wartości \mathbf{x} co najwyżej pozwala na określenie płaszczyzny dwuwymiarowej w obrębie \mathbb{R}^3 . Równanie ma rozwiązanie wtedy i tylko wtedy, gdy \mathbf{b} leży na tej płaszczyźnie.

Spełnienie $n \geq m$ jest jedynym koniecznym warunkiem istnienia rozwiązania. Nie jest to jednak warunek wystarczający, gdyż niektóre kolumny mogą być nadmiarowe. Rozważmy macierz 2 na 2, gdzie obie kolumny są identyczne. Jest to ta sama przestrzeń kolumn, jak dla macierzy 2 na 1, zawierającej jedną kopię powielonej kolumny. Innymi słowy, przestrzeń jest linią prostą i nie wskazuje wszystkich \mathbb{R}^2 , mimo że ma dwie kolumny.

Formalnie ten rodzaj nadmiarowości jest znany jako **zależność liniowa**. Zbiór wektorów jest **liniowo niezależny**, jeśli żaden wektor ze zbioru nie jest kombinacją liniową pozostałych wektorów. Jeśli dodamy do zbioru wektor stanowiący kombinację liniową innych wektorów ze zbioru, to nowy wektor nie dodaje żadnych punktów do zakresu zbioru. Oznacza to, że aby przestrzeń kolumn macierzy obejmowała wszystkie \mathbb{R}^m , to macierz musi zawierać co najmniej jeden zbiór m liniowo niezależnych kolumn. Ten warunek jest konieczny i wystarczający do tego, aby równanie 2.11 miało rozwiązanie dla każdej wartości \mathbf{b} . Warto zauważyć, że zgodnie z wymaganiami zbiór musi mieć dokładnie m niezależnych kolumn, a nie co najmniej m . Żaden zbiór m -wymiarowych wektorów nie może mieć więcej niż m wzajemnie niezależnych liniowo kolumn, ale macierz, która ma więcej niż m kolumn, może mieć więcej niż jeden taki zbiór.

Aby macierz miała macierz odwrotną, musimy dodatkowo zapewnić, że równanie 2.11 ma nie więcej niż jedno rozwiązanie dla każdej wartości \mathbf{b} . Aby tego dokonać, musimy mieć pewność, że macierz ma nie więcej niż m kolumn. W przeciwnym przypadku jest więcej niż jeden sposób parametryzacji każdego rozwiązania. W sumie oznacza to, że macierz musi być **kwadratowa**, czyli $m = n$, a wszystkie kolumny muszą być liniowo niezależne. Macierz kwadratowa zawierająca liniowo zależne kolumny jest określana jako **osobliwa**. Jeśli \mathbf{A} nie jest kwadratowa lub jest, ale jest osobliwa, to rozwiązanie równania jest nadal możliwe, ale nie można używać w tym celu macierzy odwrotnej.

Omówiliśmy dotąd odwrotność macierzy jako mnożoną lewostronnie. Można też zdefiniować odwrotność, która jest mnożona z prawej:

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{I}. \quad (2.29)$$

Dla macierzy kwadratowych odwrotności lewa i prawa są równe.

2.5. Normy

Czasami trzeba zmierzyć wielkość wektora. W systemach uczących się zwykle stosujemy w tym celu funkcję określaną jako **norma**. Formalnie norma L^p jest opisana wzorem:

$$\|\mathbf{x}\|_p = \left(\sum_i |x_i|^p \right)^{\frac{1}{p}} \quad (2.30)$$

dla $p \in \mathbb{R}, p \geq 1$.

Normy, w tym norma L^p , to funkcje odwzorowujące wektory na nieujemne wartości. Intuicyjnie norma wektora \mathbf{x} mierzy odległość od początku układu współrzędnych do punktu \mathbf{x} . Ścisłej normy to dowolna funkcja, która spełnia następujące cechy:

- $f(\mathbf{x}) = 0 \Rightarrow \mathbf{x} = 0$
- $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y})$ (**nierówność trójkąta**)
- $\forall \alpha \in \mathbb{R}, f(\alpha \mathbf{x}) = |\alpha|f(\mathbf{x})$

Norma L^2 przy $p = 2$ znana jest jako **norma euklidesowa**, co oznacza odległość euklidesową od początku układu współrzędnych do punktu wskazanego przez \mathbf{x} . Norma L^2 jest często używana w systemach uczących się i często jest oznaczana $\|\mathbf{x}\|$, a indeks dolny 2 pomijamy. Często też używa się jej jako miary wielkości wektora za pomocą normy L^2 do kwadratu, obliczanej jako $\mathbf{x}^\top \mathbf{x}$.

Kwadratowa norma L^2 jest wygodniejsza matematycznie i obliczeniowo niż sama norma L^2 . Na przykład każda pochodna kwadratowej normy L^2 względem każdego elementu \mathbf{x} zależy tylko od odpowiedniego elementu \mathbf{x} , podczas gdy pochodne zwykłej normy L^2 zależą od całego wektora. W wielu kontekstach kwadratowa norma L^2 może być niepożądana, gdyż w pobliżu początku układu współrzędnych rośnie bardzo powoli. W niektórych zastosowaniach systemów uczących się ważne jest rozróżnienie elementów, które są dokładnie równe 0, od elementów małych, ale różnych od 0. W tych przypadkach sięga się po funkcję, która rośnie w takim samym stopniu w każdym miejscu, ale pozostaje prosta matematycznie: normy L^1 . Norma L^1 w prostej postaci wygląda tak:

$$\|\mathbf{x}\|_1 = \sum_i |x_i|. \quad (2.31)$$

Norma L^1 jest powszechnie używana w systemach uczących się, gdy różnica między elementami o wartości zero i niezerowymi jest bardzo istotna. Za każdym razem, gdy element \mathbf{x} oddala się od 0 o wartość ϵ , norma L^1 wzrasta o ϵ .

Czasami mierzymy wielkość wektora, licząc, ile jest w nim elementów różnych od 0. Niektórzy autorzy odwołują się do tej funkcji jako do „normy L^0 ”, ale nie jest to poprawna terminologia. Liczba niezerowych elementów w wektorze nie jest normą, gdyż skalowanie wektora przez α nie zmienia liczby tych elementów. Norma L^1 jest często używana jako zastępstwo dla liczby elementów niezerowych.

Inną normą, która często pojawia się w systemach uczących się, jest L^∞ , znana też jako **norma maksymalna**. Jest ona uproszczeniem i odwołuje się do wartości bezwzględnej największego co do modułu elementu wektora:

$$\|\mathbf{x}\|_\infty = \max_i |x_i|. \quad (2.32)$$

Czasami chcemy zmierzyć rozmiar macierzy. W kontekście deep learningu najpopularniejsze jest zastosowanie normy Frobeniusa, uważanej zwykle za przestarzałą:

$$\|A\|_F = \sqrt{\sum_{i,j} A_{i,j}^2}, \quad (2.33)$$

Stanowi ona analogię normy L^2 dla wektora.

Iloczyn z kropką dwóch wektorów można zapisać, uwzględniając normy:

$$\mathbf{x}^\top \mathbf{y} = \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \cos \theta, \quad (2.34)$$

gdzie θ jest kątem między \mathbf{x} a \mathbf{y} .

2.6. Macierze i wektory specjalne

Niektóre specjalne rodzaje macierzy i wektorów są szczególnie przydatne. Macierze **diagonalne** zawierają głównie zera, a elementy niezerowe mają tylko na głównej przekątnej. Formalnie macierz \mathbf{D} jest macierzą diagonalną wtedy i tylko wtedy, jeśli $D_{i,j} = 0$ dla każdego $i \neq j$. Widzieliśmy już przykład takiej macierzy – macierz jednostkową, gdzie na głównej przekątnej są same jedynki. Kwadratową macierz diagonalną oznaczamy jako $\text{diag}(\mathbf{v})$. Wartości na przekątnej podane są jako elementy wektora \mathbf{v} . Macierze diagonalne są w kręgu naszego zainteresowania, gdyż mnożenie przez nie jest wydajne obliczeniowo. Aby obliczyć $\text{diag}(\mathbf{v})\mathbf{x}$, trzeba tylko przeskalować każdy element przez v_i . Innymi słowy, $\text{diag}(\mathbf{v})\mathbf{x} = \mathbf{v} \odot \mathbf{x}$. Odwroćenie kwadratowej macierzy diagonalnej jest także wydajne. Macierz odwrotna istnieje tylko wtedy, gdy wszystkie elementy na głównej przekątnej są różne od 0. Wtedy $\text{diag}(\mathbf{v})^{-1} = \text{diag}([1/v_1, \dots, 1/v_n]^\top)$. W wielu przypadkach możemy wyprowadzić ogólny

algorytm systemów uczących się w postaci dowolnych macierzy, ale uzyskać mniej kosztowny (i prostszy w opisie) algorytm dzięki ograniczeniu części macierzy do macierzy diagonalnych.

Nie wszystkie macierze diagonalne muszą być kwadratowe. Możliwe jest zbudowanie prostokątnej macierzy diagonalnej. Nie można wprawdzie takiej macierzy odwrócić, ale można przez nią tanio mnożyć. Dla niekwadratowej macierzy D iloczyn $D\mathbf{x}$ będzie wymagał skalowania każdego elementu \mathbf{x} i skalania niektórych zer w wyniku, jeśli D jest wyższa niż szersza, albo odrzucenia niektórych końcowych elementów wektora, jeśli D jest szersza niż wyższa.

Macierz **symetryczna** to każda macierz, która równa się swojej transpozycji:

$$\mathbf{A} = \mathbf{A}^\top. \quad (2.35)$$

Macierze symetryczne pojawiają się, gdy elementy są wynikiem działania funkcji dwóch argumentów, przy czym wartość nie zależy od ich kolejności. Jeśli na przykład \mathbf{A} jest macierzą pomiarów odległości, gdzie $\mathbf{A}_{i,j}$ daje odległość od punktu i do punktu j , to $\mathbf{A}_{i,j} = \mathbf{A}_{j,i}$, gdyż funkcje odległości są symetryczne.

Wektor jednostkowy to wektor o **normie jednostkowej**:

$$\|\mathbf{x}\|_2 = 1. \quad (2.36)$$

Wektory \mathbf{x} i \mathbf{y} są względem siebie **ortogonalne**, jeśli $\mathbf{x}^\top \mathbf{y} = 0$. Jeśli oba wektory mają normy różne od 0, to oznacza, że są względem siebie pod kątem 90 stopni. W \mathbb{R}^n nie więcej niż n wektorów może być wzajemnie ortogonalnych z niezerową normą. Jeśli wektory są nie tylko ortogonalne, lecz także mają jednostkową normę, to nazywamy je **ortonormalnymi**.

Macierz ortogonalna to macierz kwadratowa, której wiersze są wzajemnie ortonormalne i kolumny są wzajemnie ortonormalne:

$$\mathbf{A}^\top \mathbf{A} = \mathbf{A} \mathbf{A}^\top = \mathbf{I}. \quad (2.37)$$

Stąd wynika, że:

$$\mathbf{A}^{-1} = \mathbf{A}^\top, \quad (2.38)$$

więc macierze ortogonalne są w zakresie naszego zainteresowania, gdyż ich odwrotność jest bardzo prosta do obliczenia. Warto zwrócić uwagę na definicję macierzy ortogonalnych. Wbrew intuicji ich wiersze są nie tylko ortogonalne, lecz także w pełni ortonormalne. Nie ma oddzielnej nazwy na macierz, której wiersze lub kolumny są ortogonalne, ale nie są ortonormalne.

2.7. Rozkład na wartości własne

Wiele obiektów matematycznych można lepiej zrozumieć, dzieląc je na składniki lub znajdując ich właściwości, które są uniwersalne, niezależne od tego, jaki wybierzemy sposób prezentacji. Na przykład liczby całkowite można rozłożyć na czynniki pierwsze. Sposób reprezentacji liczby 12 będzie się różnić w zależności od tego, czy zapiszemy ją w systemie dziesiętnym czy binarnym, ale zawsze prawdziwa będzie zależność $12 = 2 \times 2 \times 3$. Z tego sposobu prezentacji możemy wyprowadzić cenne właściwości, takie jak to, że 12 nie jest podzielne przez 5 oraz że każda liczba całkowita będąca wielokrotnością 12 będzie podzielna przez 3.

Podobnie jak dowiadujemy się czegoś o liczbie całkowitej, rozkładając ją na czynniki pierwsze, możemy rozłożyć macierze w taki sposób, żeby uzyskać informacje o ich właściwościach, które nie są oczywiste na podstawie samej prezentacji ich elementów. Jednym z szeroko używanych sposobów dekompozycji macierzy jest **rozkład na wartości własne**, gdzie rozkładamy macierz na zbiór wektorów własnych i wartości własnych. **Wektor własny** (ang. *eigenvector*) macierzy kwadratowej \mathbf{A} to niezerowy wektor \mathbf{v} , którego pomnożenie przez λ zmienia tylko skalę \mathbf{v} :

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}. \quad (2.39)$$

Skalar λ to **wartość własna** odpowiadająca wektorowi własnemu (można też znaleźć **lewostronny wektor własny**, taki, że $\mathbf{v}^\top \mathbf{A} = \lambda \mathbf{v}^\top$, ale zwykle zajmujemy się prawostronnymi wektorami własnymi). Jeśli \mathbf{v} jest wektorem własnym \mathbf{A} , to jest nim także każdy przeskalowany wektor $s\mathbf{v}$, dla każdego $s \in \mathbb{R}, s \neq 0$. Co więcej, $s\mathbf{v}$ ma tę samą wartość własną, dlatego zwykle szukamy tylko jednostkowych wektorów własnych.

Przypuśćmy, że macierz \mathbf{A} ma n niezależnych liniowo wektorów własnych $\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}\}$ z odpowiadającymi im wartościami własnymi $\{\lambda_1, \dots, \lambda_n\}$. Możemy połączyć wszystkie wektory własne w macierz \mathbf{V} z jednym wektorem własnym w każdej kolumnie: $\mathbf{V} = [\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}]$. Podobnie możemy scalić wartości własne w postaci wektora $\boldsymbol{\lambda} = [\lambda_1, \dots, \lambda_n]^\top$. **Dekompozycja własna** macierzy \mathbf{A} ma postać:

$$\mathbf{A} = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}. \quad (2.40)$$

Widzimy, że budowa macierzy z określonymi wartościami własnymi i wektorami własnymi pozwala nam na rozszerzenie przestrzeni w pożądanym kierunku. Czasem chcemy **rozłożyć** macierz na jej wartości i wektory własne. Pomaga to w analizie określonych właściwości macierzy, podobnie jak rozkład

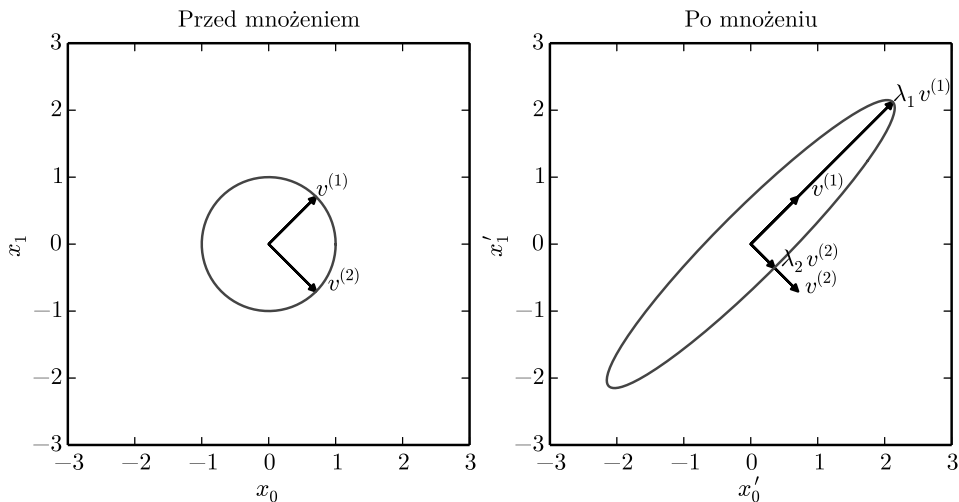
liczby całkowitej na czynniki pierwsze pozwala zrozumieć zachowanie się liczby całkowitej.

Nie każdą macierz można rozłożyć na wartości i wektory własne. W niektórych przypadkach dekompozycja jest możliwa, ale obejmuje liczby zespolone, a nie rzeczywiste. Na szczęście w tej książce musimy dokonywać rozkładu tylko określonego typu macierzy, które mają prosty rozkład. Konkretnie, każdą rzeczywistą macierz symetryczną można rozłożyć na wyrażenia, wykorzystując tylko rzeczywiste wartości i wektory własne:

$$\mathbf{A} = \mathbf{Q}\mathbf{\Lambda}\mathbf{Q}^T, \tag{2.41}$$

gdzie \mathbf{Q} jest macierzą ortogonalną złożoną z wektorów własnych \mathbf{A} , natomiast $\mathbf{\Lambda}$ jest macierzą diagonalną. Wartość własna $\Lambda_{i,i}$ jest powiązana z wektorem z kolumny i macierzy \mathbf{Q} i oznaczona jako $Q_{:,i}$. Ponieważ macierz \mathbf{Q} jest ortogonalna, można traktować \mathbf{A} jako przestrzeń skalowania względem λ_i w kierunku $v^{(i)}$. Przykład pokazano na rysunku 2.3.

Podczas gdy każdą rzeczywistą symetryczną macierz \mathbf{A} na pewno można rozłożyć na elementy własne, dekompozycja ta może nie być unikatowa. Jeśli jeden lub więcej wektorów własnych ma tę samą wartość własną, to każdy zbiór ortogonalnych wektorów leżący w ich zakresie stanowi wektory własne



Rysunek 2.3. Efekt rozkładu na wartości i wektory własne. Macierz \mathbf{A} ma dwa ortogonalne wektory własne, gdzie $v^{(1)}$ ma wartość własną λ_1 , a $v^{(2)}$ wartość własną λ_2 . Po lewej: kreślimy zbiór wszystkich wektorów jednostkowych $u \in \mathbb{R}^2$ jako okrąg jednostkowy. Po prawej: kreślimy zbiór wszystkich punktów $\mathbf{A}u$. Obserwując sposób, w jaki \mathbf{A} deformuje okrąg jednostkowy, możemy zobaczyć, że skaluje on przestrzeń w kierunku $v^{(i)}$ przez λ_i

z tą samą wartością własną i można wybrać macierz \mathbf{Q} , wykorzystującą te wektory własne zamiast innych. Przyjęto, że sortujemy elementy macierzy Λ w kolejności malejącej. Zgodnie z tą konwencją dekompozycja na elementy własne jest unikatowa tylko wtedy, gdy wartości własne są unikatowe.

Dekompozycja własna macierzy podaje wiele faktów o tej macierzy. Macierz jest osobliwa wtedy i tylko wtedy, gdy jedna z jej wartości własnych równa się 0. Dekompozycja własna rzeczywistej macierzy symetrycznej może być wykorzystana do optymalizacji wyrażeń kwadratowych w postaci $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$ względem $\|\mathbf{x}\|_2 = 1$. Gdy \mathbf{x} jest równe wektorowi własnemu \mathbf{A} , to f przyjmuje wartość odpowiedniej wartości własnej. Maksymalna wartość f w ograniczonym obszarze stanowi maksymalną wartość własną, a wartość minimalna – minimalną wartość własną.

Macierz, której wszystkie wartości własne są dodatnie, jest nazywana **dodatnio określona**. Jeśli wartości są większe lub równe 0, to macierz nazywamy **półokreślona dodatnio**. Podobnie przy wartościach ujemnych macierz jest **ujemnie określona**, a przy wartościach własnych mniejszych lub równych 0 jest **półokreślona ujemnie**. Macierze półokreślone dodatnio są interesujące, gdyż gwarantują, że $\forall \mathbf{x}, \mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$. Macierze określone dodatnio dodatkowo zapewniają, że $\mathbf{x}^\top \mathbf{A} \mathbf{x} = 0 \Rightarrow \mathbf{x} = \mathbf{0}$.

2.8. Dekompozycja wartości osobliwej

W punkcie 2.7 opisano, jak dokonać dekompozycji macierzy na wartości i wektory własne. **Dekompozycja wartości osobliwej** (ang. *singular value decomposition*, SVD) daje jeszcze jeden sposób rozkładu macierzy na czynniki na **wektory osobliwe** i **wartości osobliwe**. SVD pozwala nam otrzymać te same informacje, które daje dekompozycja na czynniki własne; jednak SVD ma szersze zastosowanie. Każda macierz rzeczywista ma wartość osobliwą dekompozycji, co nie jest zawsze prawdą w odniesieniu do wartości własnych. Jeśli na przykład macierz nie jest kwadratowa, rozkład na czynniki własne nie jest zdefiniowany i musimy stosować dekompozycję na wartości osobliwe.

Przypomnijmy, że dekompozycja na wartości własne obejmuje analizę macierzy \mathbf{A} w celu znalezienia macierzy \mathbf{V} wektorów własnych i wektora wartości własnych $\boldsymbol{\lambda}$, dzięki którym możemy zapisać macierz \mathbf{A} jako:

$$\mathbf{A} = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^{-1}. \quad (2.42)$$

Dekompozycja na wartości osobliwe jest podobna, ale tym razem \mathbf{A} zapisuje się jako iloczyn trzech macierzy:

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^\top. \quad (2.43)$$

Przypuśćmy, że \mathbf{A} jest macierzą $m \times n$. Wtedy \mathbf{U} jest macierzą $m \times m$, \mathbf{D} macierzą $m \times n$, a \mathbf{V} macierzą $n \times n$. Każda z tych macierzy ma zdefiniowaną specjalną strukturę. Macierze \mathbf{U} i \mathbf{V} są ortogonalne. Macierz \mathbf{D} jest macierzą diagonalną. Zauważmy, że \mathbf{D} nie musi być kwadratowa. Elementy na przekątnej macierzy \mathbf{D} są **wartościami osobliwymi** macierzy \mathbf{A} . Kolumny \mathbf{U} są **wektorami lewostronnie osobliwymi**, natomiast kolumny \mathbf{V} są **wektorami prawostronnie osobliwymi**.

Dekompozycję macierzy \mathbf{A} na wartości osobliwe interpretujemy jako dekompozycję na wartości własne funkcji \mathbf{A} . Wektory lewostronnie osobliwe macierzy \mathbf{A} są wektorami osobliwymi $\mathbf{A}\mathbf{A}^\top$. Wektory prawostronnie osobliwe to wektory osobliwe $\mathbf{A}^\top\mathbf{A}$. Niezerowe wartości osobliwe macierzy \mathbf{A} to pierwiastki kwadratowe wartości osobliwych $\mathbf{A}^\top\mathbf{A}$. To samo odnosi się do $\mathbf{A}\mathbf{A}^\top$.

Zapewne najbardziej użyteczną funkcją SVD jest to, że można jej użyć do częściowego uogólnienia odwracania niekwadratowych macierzy, co zostanie zaprezentowane w kolejnym punkcie.

2.9. Uogólniona macierz odwrotna (Moore'a–Penrose'a)

Odwracanie macierzy nie jest zdefiniowane dla macierzy niekwadratowych. Przypuśćmy, że chcemy dokonać lewostronnej odwrotności \mathbf{B} macierzy \mathbf{A} , aby rozwiązać równanie liniowe:

$$\mathbf{A}\mathbf{x} = \mathbf{y}, \quad (2.44)$$

mnożąc lewostronnie obie strony, aby uzyskać:

$$\mathbf{x} = \mathbf{B}\mathbf{y}. \quad (2.45)$$

Zależnie od struktury problemu jednoznaczne odwzorowanie z \mathbf{A} na \mathbf{B} może nie być możliwe. Jeśli macierz \mathbf{A} jest wyższa niż szersza, równanie może nie mieć rozwiązania. Jeśli \mathbf{A} jest szersza niż wyższa, to może być wiele możliwych rozwiązań.

Uogólniona macierz odwrotna (ang. *Moore-Penrose pseudoinverse*) pozwala na pewien postępek w tych przypadkach. Uogólniona odwrotność macierzy \mathbf{A} zdefiniowana jest jako macierz:

$$\mathbf{A}^+ = \lim_{\alpha \searrow 0} (\mathbf{A}^\top\mathbf{A} + \alpha\mathbf{I})^{-1}\mathbf{A}^\top. \quad (2.46)$$

Praktycznie algorytm wyznaczenia uogólnionej macierzy odwrotnej nie opiera się na definicji, a na wzorze:

$$\mathbf{A}^+ = \mathbf{V}\mathbf{D}^+\mathbf{U}^\top, \quad (2.47)$$

gdzie \mathbf{U} , \mathbf{D} i \mathbf{V} to wartości osobliwe dekompozycji \mathbf{A} , natomiast uogólnioną odwrotność \mathbf{D}^+ diagonalnej macierzy \mathbf{D} otrzymujemy, biorąc odwrotność jej niezerowych elementów, a następnie transponując macierz wynikową.

Gdy \mathbf{A} ma więcej kolumn niż wierszy, rozwiązanie równania liniowego za pomocą uogólnionej macierzy odwrotnej daje jedno z wielu możliwych rozwiązań. Konkretnie daje to spośród wielu rozwiązań $\mathbf{x} = \mathbf{A}^+\mathbf{y}$, które ma najmniejszą ze wszystkich normę euklidesową $\|\mathbf{x}\|_2$. Gdy \mathbf{A} ma więcej wierszy niż kolumn, możliwe, że nie ma żadnego rozwiązania. W takim przypadku zastosowanie uogólnionej macierzy odwrotnej daje wartość \mathbf{x} , dla której $\mathbf{A}\mathbf{x}$ jest możliwie najbliższe \mathbf{y} w sensie normy euklidesowej $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2$.

2.10. Operator śladowy

Operator śladowy daje sumę wszystkich elementów na przekątnej macierzy:

$$\text{Tr}(\mathbf{A}) = \sum_i \mathbf{A}_{i,i}. \quad (2.48)$$

Operator śladowy jest przydatny z wielu powodów. Niektóre działania, które są trudne do określenia bez sprowadzenia do operacji dodawania, można określić za pomocą iloczynu macierzy i operatora śladowego. Na przykład operator śladowy daje alternatywny sposób zapisu normy Frobeniusa dla macierzy:

$$\|\mathbf{A}\|_F = \sqrt{\text{Tr}(\mathbf{A}\mathbf{A}^\top)}. \quad (2.49)$$

Pisanie wyrażenia w kontekście operatora śladowego otwiera możliwości manipulowania wyrażeniem za pomocą wielu użytecznych tożsamości. Na przykład operator śladowy jest niezmienny względem operatora transpozycji:

$$\text{Tr}(\mathbf{A}) = \text{Tr}(\mathbf{A}^\top). \quad (2.50)$$

Ślad macierzy kwadratowej składa się z wielu czynników i też jest niezmienny względem przeniesienia ostatniego czynnika na pierwszą pozycję, jeśli kształty odpowiednich macierzy pozwalają na zdefiniowanie powstałego iloczynu:

$$\text{Tr}(\mathbf{ABC}) = \text{Tr}(\mathbf{CAB}) = \text{Tr}(\mathbf{BCA}) \quad (2.51)$$

lub bardziej ogólnie:

$$\mathrm{Tr}\left(\prod_{i=1}^n \mathbf{F}^{(i)}\right) = \mathrm{Tr}\left(\mathbf{F}^{(n)} \prod_{i=1}^{n-1} \mathbf{F}^{(i)}\right). \quad (2.52)$$

Ta niezmiennosc względem cyklicznych permutacji utrzymuje się nawet wtedy, gdy wynikowy iloczyn ma inny kształt. Na przykład dla $\mathbf{A} \in \mathbb{R}^{m \times n}$ i $\mathbf{B} \in \mathbb{R}^{n \times m}$ mamy:

$$\mathrm{Tr}(\mathbf{AB}) = \mathrm{Tr}(\mathbf{BA}), \quad (2.53)$$

mimo że $\mathbf{AB} \in \mathbb{R}^{m \times m}$, a $\mathbf{BA} \in \mathbb{R}^{n \times n}$.

Innym pożytecznym faktem jest to, że skalar jest swoim własnym śladem: $\mathbf{a} = \mathrm{Tr}(\mathbf{a})$.

2.11. Wyznacznik

Wyznacznik macierzy kwadratowej, oznaczany jako $\det(\mathbf{A})$, to funkcja odwzorowująca macierz na skalary rzeczywiste. Wyznacznik jest równy iloczynowi wszystkich wartości własnych macierzy. Wartość bezwzględna wyznacznika można traktować jako miarę, na ile mnożenie przez macierz rozszerza lub zawęża przestrzeń. Jeśli wyznacznik jest równy 0, to przestrzeń zostaje całkowicie zmniejszona wzdłuż przynajmniej jednego wymiaru, co powoduje utratę objętości. Jeśli wyznacznik jest równy 1, to przekształcenie zachowuje objętość.

2.12. Przykład: analiza głównych składowych

Na podstawie samej znajomości algebry liniowej można wyprowadzić jeden prosty algorytm dla systemów uczących się. Jest to **analiza głównych składowych** (ang. *principal components analysis*, PCA).

Przypuśćmy, że mamy zbiór m punktów $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}$ w \mathbb{R}^n i chcemy zastosować ich stratną kompresję. Kompresja stratna oznacza zapisanie punktów w sposób, który wymaga mniej pamięci, ale powoduje utratę dokładności. Chcemy stracić jak najmniej na dokładności. Jednym ze sposobów na zakodowanie tych punktów jest reprezentowanie ich wersji w mniejszej liczbie wymiarów. Dla każdego punktu $\mathbf{x}^{(i)} \in \mathbb{R}^n$ znajdziemy odpowiadający mu wektor kodowania $\mathbf{c}^{(i)} \in \mathbb{R}^l$. Jeśli l jest mniejsze niż n , przechowywanie punktów kodu zajmie mniej pamięci, niż zapisanie przechowywanie oryginalnych danych. Będziemy chcieli znaleźć jakąś funkcję kodowania, która utworzy kod

dla danych wejściowych $f(\mathbf{x}) = \mathbf{c}$, oraz funkcję dekodowania, która utworzy zrekonstruowane dane wejściowe na podstawie ich kodu $\mathbf{x} \approx g(f(\mathbf{x}))$.

Wybieramy PCA jako naszą funkcję dekodowania. Aby dekodery był bardzo prosty, wybieramy mnożenie macierzy do odwzorowania kodu z powrotem na \mathbb{R}^n . Niech $g(\mathbf{c}) = \mathbf{D}\mathbf{c}$, gdzie $\mathbf{D} \in \mathbb{R}^{n \times l}$ jest macierzą definiującą dekodowanie. Wyznaczenie optymalnego kodu dla tego dekodera może stanowić problem. Aby zagadnienie kodowania było proste, PCA ogranicza kolumny macierzy \mathbf{D} do takich, które są względem siebie ortogonalne (zauważmy, że macierz \mathbf{D} nie jest formalnie macierzą ortogonalną, jeśli nie jest spełniony warunek $l = n$). Przy tak zdefiniowanym problemie możliwe jest wiele rozwiązań, gdyż możemy zwiększyć skalę macierzy $D_{:,i}$, jeśli proporcjonalnie zmniejszymy c_i dla wszystkich punktów. Aby problem miał jedno rozwiązanie, ograniczamy kolumny \mathbf{D} , tak aby miały jednostkową normę.

Aby przemienić ten pomysł w algorytm, który można wdrożyć, trzeba określić, jak wygenerować optymalny punkt kodu \mathbf{c}^* dla każdego punktu wejściowego \mathbf{x} . Jednym ze sposobów jest zminimalizowanie odległości między punktem wejściowym \mathbf{x} a jego rekonstrukcją $g(\mathbf{c}^*)$. Można zmierzyć tę odległość za pomocą normy. W algorytmie głównych składowych używamy normy L^2 :

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2. \quad (2.54)$$

Można przejść na normę kwadratową L^2 zamiast samej normy L^2 , gdyż obie są minimalizowane przez tę samą wartość \mathbf{c} . Dzieje się tak dlatego, że norma L^2 jest nieujemna i podniesienie jej do kwadratu monotonicznie zwiększa nieujemne argumenty:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} \|\mathbf{x} - g(\mathbf{c})\|_2^2. \quad (2.55)$$

Minimalizowana funkcja sprowadza się do:

$$(\mathbf{x} - g(\mathbf{c}))^\top (\mathbf{x} - g(\mathbf{c})) \quad (2.56)$$

(według definicji normy L^2 z równania 2.30)

$$= \mathbf{x}^\top \mathbf{x} - \mathbf{x}^\top g(\mathbf{c}) - g(\mathbf{c})^\top \mathbf{x} + g(\mathbf{c})^\top g(\mathbf{c}) \quad (2.57)$$

(zgodnie z własnością rozdzielności)

$$= \mathbf{x}^\top \mathbf{x} - 2\mathbf{x}^\top g(\mathbf{c}) + g(\mathbf{c})^\top g(\mathbf{c}) \quad (2.58)$$

(gdyż skalar $g(\mathbf{c})^\top \mathbf{x}$ jest równy swojej własnej transpozycji).

Teraz można ponownie zmienić minimalizowaną funkcję, aby ominąć pierwszy składnik, gdyż nie jest on zależny od \mathbf{c} :

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} -2\mathbf{x}^\top g(\mathbf{c}) + g(\mathbf{c})^\top g(\mathbf{c}). \quad (2.59)$$

Dalej musimy zrobić podstawienie w definicji $g(\mathbf{c})$:

$$\mathbf{c}^* = \arg \min_{\mathbf{c}} -2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{D}^\top \mathbf{D}\mathbf{c} \quad (2.60)$$

$$= \arg \min_{\mathbf{c}} -2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{I}_l \mathbf{c} \quad (2.61)$$

(przez ortogonalność i ograniczenia normy jednostkowej \mathbf{D})

$$= \arg \min_{\mathbf{c}} -2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{c}. \quad (2.62)$$

Możemy rozwiązać ten problem optymalizacyjny za pomocą rachunku wektorowego (patrz punkt 4.3, jeśli nie wiesz, jak to się robi):

$$\nabla_{\mathbf{c}}(-2\mathbf{x}^\top \mathbf{D}\mathbf{c} + \mathbf{c}^\top \mathbf{c}) = \mathbf{0} \quad (2.63)$$

$$-2\mathbf{D}^\top \mathbf{x} + 2\mathbf{c} = \mathbf{0} \quad (2.64)$$

$$\mathbf{c} = \mathbf{D}^\top \mathbf{x}. \quad (2.65)$$

To sprawia, że algorytm jest wydajny: możemy optymalnie zakodować \mathbf{x} za pomocą działania wektorowo macierzowego. Aby zakodować wektor, stosuje się funkcję kodera:

$$f(\mathbf{x}) = \mathbf{D}^\top \mathbf{x}. \quad (2.66)$$

Korzystając dalej z mnożenia macierzy, można także zdefiniować działanie rekonstrukcji PCA:

$$r(\mathbf{x}) = g(f(\mathbf{x})) = \mathbf{D}\mathbf{D}^\top \mathbf{x}. \quad (2.67)$$

Dalej należy wybrać macierz kodowania \mathbf{D} . Aby to zrobić, wracamy do minimalizacji odległości L^2 między wejściami a rekonstrukcjami. Ponieważ używamy tej samej macierzy \mathbf{D} do odkodowania wszystkich punktów, nie możemy dalej rozpatrywać tych punktów oddzielnie. Przeciwnie, musimy zminimalizować normę Frobeniusa macierzy błędów względem wszystkich wymiarów i punktów:

$$\mathbf{D}^* = \arg \min_{\mathbf{D}} \sqrt{\sum_{i,j} \left(x_j^{(i)} - r(\mathbf{x}^{(i)})_j\right)^2} \text{ dla } \mathbf{D}^\top \mathbf{D} = \mathbf{I}_l. \quad (2.68)$$

Aby wyprowadzić algorytm znajdowania \mathbf{D}^+ , rozpoczynamy od rozważenia przypadku, gdzie $l = 1$. Wtedy \mathbf{D} jest po prostu pojedynczym wektorem \mathbf{d} . Po podstawieniu równania 2.67 do równania 2.68 i uproszczeniu \mathbf{D} na \mathbf{d} problem sprowadza się do:

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{d}\mathbf{d}^\top \mathbf{x}^{(i)}\|_2^2 \text{ dla } \|\mathbf{d}\|_2 = 1. \quad (2.69)$$

Powyższe sformułowanie jest najbardziej bezpośrednim sposobem wykonania podstawienia, ale nie jest to najlepszy sposób zapisu równania. Umieszcza się wartość skalarną $d_{>x^{(i)}}$ na prawo od wektora \mathbf{d} . Skalarne współczynniki są zwykle zapisywane po lewej stronie wektora, względem którego są stosowane. Dlatego zwykle zapisuje się ten wzór w postaci:

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{d}^\top \mathbf{x}^{(i)} \mathbf{d}\|_2^2 \text{ dla } \|\mathbf{d}\|_2 = 1, \quad (2.70)$$

lub, wykorzystując fakt, że skalar jest swoją własną transpozycją, jako:

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{x}^{(i)\top} \mathbf{d}\mathbf{d}\|_2^2 \text{ dla } \|\mathbf{d}\|_2 = 1. \quad (2.71)$$

Czytelnicy powinni przyzwyczać się do takich kosmetycznych przekształceń.

W tym miejscu pomocne może być przepisanie problemu według przykładów pojedynczego projektu macierzy, zamiast stosowania sumy oddzielnych przykładowych wektorów. To pozwoli na bardziej zwięzły zapis. Niech $\mathbf{X} \in \mathbb{R}^{m \times n}$ będzie macierzą zdefiniowaną przez ułożenie w stos wszystkich wektorów opisujących punkty $X_{i,:} = \mathbf{x}^{(i)\top}$. Można teraz przepisać problem jako:

$$\mathbf{d}^* = \arg \min_{\mathbf{d}} \|\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top\|_F^2 \text{ dla } \mathbf{d}^\top \mathbf{d} = 1. \quad (2.72)$$

Pomijając na chwilę ograniczenie, można uprościć część normy Frobeniusa w następujący sposób:

$$\arg \min_{\mathbf{d}} \|\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top\|_F^2 \quad (2.73)$$

$$= \arg \min_{\mathbf{d}} \text{Tr} \left(\left(\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top \right)^\top \left(\mathbf{X} - \mathbf{X}\mathbf{d}\mathbf{d}^\top \right) \right) \quad (2.74)$$

(zgodnie z równaniem 2.49)

$$= \arg \min_{\mathbf{d}} \text{Tr}(\mathbf{X}^\top \mathbf{X} - \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top - \mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X} + \mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) \quad (2.75)$$

$$= \arg \min_{\mathbf{d}} \text{Tr}(\mathbf{X}^\top \mathbf{X}) - \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) - \text{Tr}(\mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}) + \text{Tr}(\mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) \quad (2.76)$$

$$= \arg \min_{\mathbf{d}} - \text{Tr}(\mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) - \text{Tr}(\mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}) + \text{Tr}(\mathbf{d}\mathbf{d}^\top \mathbf{X}^\top \mathbf{X}\mathbf{d}\mathbf{d}^\top) \quad (2.77)$$

(ponieważ składniki nie dotyczące \mathbf{d} nie wpływają na $\arg \min$)

$$= \arg \min_{\mathbf{d}} -2 \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{d} \mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \quad (2.78)$$

(ponieważ można zmieniać cyklicznie kolejność macierzy w śladzie, zgodnie z równaniem 2.52)

$$= \arg \min_{\mathbf{d}} -2 \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top \mathbf{d} \mathbf{d}^\top) \quad (2.79)$$

(używając ponownie tej samej własności).

W tym miejscu ponownie wprowadzamy ograniczenie: $\arg \min$:

$$\arg \min_{\mathbf{d}} -2 \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top \mathbf{d} \mathbf{d}^\top) \text{ dla } \mathbf{d}^\top \mathbf{d} = 1 \quad (2.80)$$

$$= \arg \min_{\mathbf{d}} -2 \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) + \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ dla } \mathbf{d}^\top \mathbf{d} = 1 \quad (2.81)$$

(w związku z ograniczeniem)

$$= \arg \min_{\mathbf{d}} -\operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ dla } \mathbf{d}^\top \mathbf{d} = 1 \quad (2.82)$$

$$= \arg \max_{\mathbf{d}} \operatorname{Tr}(\mathbf{X}^\top \mathbf{X} \mathbf{d} \mathbf{d}^\top) \text{ dla } \mathbf{d}^\top \mathbf{d} = 1 \quad (2.83)$$

$$= \arg \max_{\mathbf{d}} \operatorname{Tr}(\mathbf{d}^\top \mathbf{X}^\top \mathbf{X} \mathbf{d}) \text{ dla } \mathbf{d}^\top \mathbf{d} = 1. \quad (2.84)$$

Ten problem optymalizacyjny może zostać rozwiązany za pomocą dekompozycji własnej. W szczególności optymalna wartość \mathbf{d} jest podana przez wektor własny $\mathbf{X}^\top \mathbf{X}$ odpowiadający największej wartości własnej.

To wyprowadzenie jest właściwa dla przypadku $l = 1$ i podaje tylko pierwszą główną składową. Bardziej ogólnie, jeśli chcemy uzyskać bazę głównych składowych, macierz \mathbf{D} jest określona przez l wektorów własnych odpowiadających największej wartości własnej. To można udowodnić przez indukcję. Zalecamy jako ćwiczenie przeprowadzenie takiego dowodu.

Algebra liniowa jest jedną z podstawowych dziedzin matematyki potrzebnych do zrozumienia deep learningu. Innym kluczowym obszarem matematyki niezbędnym w systemach uczących się jest teoria prawdopodobieństwa, prezentowana w kolejnym rozdziale.

3

Prawdopodobieństwo i teoria informacji

W tym rozdziale opisujemy teorię prawdopodobieństwa oraz teorię informacji. Teoria prawdopodobieństwa to matematyczny system reprezentacji niepewności. Daje środki pomiaru niepewności, a także aksjomaty wyprowadzania nowych stwierdzeń związanych z niepewnością. W zastosowaniach sztucznej inteligencji używa się teorii prawdopodobieństwa na dwa podstawowe sposoby. Po pierwsze prawa prawdopodobieństwa mówią, jak systemy AI powinny rozumować, aby można projektować algorytmy do wyznaczania lub aproksymacji wyrażeń wyprowadzonych za pomocą teorii prawdopodobieństwa. Po drugie możemy używać prawdopodobieństwa i statystyki do teoretycznej analizy zachowania proponowanych systemów AI.

Teoria prawdopodobieństwa to podstawowe narzędzie w wielu dziedzinach nauki i inżynierii. Ten rozdział służy temu, aby czytelnicy, którzy mają wykształcenie w dziedzinie programowania i mało poznali teorię prawdopodobieństwa, mogli zrozumieć zawartość tej książki.

Podczas gdy teoria prawdopodobieństwa pozwala tworzyć niepewne stwierdzenia i rozważać zagadnienia z niepewnością, teoria informacji pozwala nam na oszacowanie niepewności ilościowo jako rozkład prawdopodobieństwa.

Każdy, kto zna teorię prawdopodobieństwa oraz teorię informacji, może pominąć ten rozdział, poza punktem 3.14, w którym opisano grafy stosowane do opisu strukturalnych modeli probabilistycznych w systemach uczących się. Jeśli nie macie żadnego doświadczenia w tych dziedzinach, ten rozdział wystarczy, aby prowadzić projekty deep learningu, ale warto rozszerzyć swoją wiedzę z innych źródeł, jak na przykład Jaynes (2003).

3.1. Dlaczego prawdopodobieństwo?

Wiele gałęzi nauk komputerowych ma do czynienia z elementami, które są całkowicie deterministyczne i pewne. Programista może zwykle bezpiecznie założyć, że procesor wykona każdą instrukcję bez błędów. Błędy sprzętowe się zdarzają, ale są na tyle rzadkie, że większość aplikacji programistycznych nie musi ich brać pod uwagę. W związku z tym, że wielu informatyków i programistów działa we względnie czystym i pewnym środowisku, zaskoczeniem może być dla nich szerokie zastosowanie teorii prawdopodobieństwa w deep learningu.

Systemy uczące się mają do czynienia z niepewnymi ilościami, a czasami też z elementami stochastycznymi (niedeterministycznymi). Niepewność i stochastyka mogą pochodzić z wielu źródeł. Począwszy od lat osiemdziesiątych XX wieku naukowcy przedstawili istotne argumenty na ocenę niepewności za pomocą prawdopodobieństwa. Wiele przedstawianych tu argumentów stanowi podsumowanie lub jest zainspirowanych pracami Pearl'a (1988).

Prawie wszystkie działania wymagają wnioskowania w warunkach niepewności. W zasadzie poza twierdzeniami matematycznymi, które są z definicji prawdziwe, trudno znaleźć twierdzenia, które są z pewnością prawdziwe, lub zdarzenia, które na pewno będą miały miejsce. Oto trzy możliwe źródła niepewności:

1. Stochastyczna natura modelowanego systemu. Na przykład większość interpretacji w mechanice kwantowej opisuje dynamikę cząsteczek subatomowych jako probabilistyczne. Możemy też stworzyć teoretyczne scenariusze, które mają dynamikę losową, jak hipotetyczna gra w karty, gdzie zakładamy, że karty są naprawdę potasowane w kolejności losowej.
2. Niepełne możliwości obserwacji. Nawet systemy deterministyczne mogą okazać się stochastyczne, gdy nie można obserwować wszystkich zmiennych, które wpływają na zachowanie systemu. Na przykład w paradoksie Monty'ego Halla* uczestnicy gry mogą wybrać jedno z trojga drzwi, a wygrana znajduje się za zamkniętymi drzwiami. Dwoje drzwi prowadzi do kozy, a trzecie do samochodu. Wynik wyboru uczestnika gry jest deterministyczny, ale z jego punktu widzenia jest on losowy.
3. Niekompletne modelowanie. Gdy korzysta się z modelu, w którym trzeba odrzucić część zaobserwowanych informacji, to odrzucone informacje wprowadzają niepewność do przewidywań modelu. Przypuśćmy, że budujemy robota, który potrafi dokładnie określić położenie każdego

*Od nazwiska prezentera teleturnieju popularnego w USA i Kanadzie.

otaczającego go obiektu. Jeśli robot podczas przewidywania przyszłego położenia przetwarza przestrzeń na wartości dyskretne, to zabieg ten sprawia, że nie potrafi on precyzyjnie przewidzieć położenia obiektów: każdy obiekt może być w dowolnym miejscu dyskretnej komórki, w której ma się znajdować.

W wielu przypadkach bardziej praktyczne jest skorzystanie z prostej, lecz niepewnej zasady niż z pewnej, ale skomplikowanej, jeśli nawet zasada jest deterministyczna i system modelowania ma możliwość spełnienia skomplikowanej zasady. Na przykład prosta zasada: „Większość ptaków lata” jest prostsza do wprowadzenia i szerzej użyteczna niż zasada w postaci: „Ptaki latają, poza bardzo młodymi, które się jeszcze nie nauczyły latać, chorymi lub skaleczonymi, które utraciły możliwość latania, niektórymi gatunkami ptaków nielotnych, w tym kazuarów, strusi, kiwi...”, która jest trudna do zastosowania, utrzymania i rozpowszechniania, a po pokonaniu trudności jest nadal niepewna i narażona na niepowodzenie.

Podczas gdy powinno być jasne, że musimy mieć środki do prezentowania i rozważania niepewności, nie jest oczywiste, że teoria prawdopodobieństwa może dać wszystkie narzędzia potrzebne w zastosowaniach sztucznej inteligencji. Teoria prawdopodobieństwa powstała w celu analizy częstości występowania zdarzeń. Łatwo zobaczyć, jak stosować ją do analizy zdarzeń, takich jak wylosowanie określonych kart w rozdaniu podczas gry w pokera. Tego rodzaju zdarzenia są często powtarzalne. Gdy mówimy, że wynik uzyskamy z prawdopodobieństwem p , to znaczy, że jeśli powtórzymy eksperyment (np. rozdawanie kart) nieskończenie wiele razy, to w wynikach nastąpi p powtórzeń. Ten sposób rozumowania nie wydaje się możliwy do natychmiastowego zastosowania do propozycji, które się nie powtarzają. Jeśli lekarz zbada pacjenta i powie, że ma on na 40% grype, to znaczy coś całkiem innego – nie możemy mieć nieskończenie wielu kopii pacjenta ani nie ma powodu wierzyć, że kopie pacjenta miałyby te same objawy przy różnych warunkach początkowych. W przypadku diagnozy lekarskiej pacjenta korzystamy z prawdopodobieństwa, aby pokazać **stopień przekonania**, gdzie 1 oznacza absolutną pewność, że pacjent ma grype, a 0 to absolutna pewność, że grypy nie ma. Pierwszy rodzaj prawdopodobieństwa, powiązany bezpośrednio z częstością występowania zdarzeń, znany jest jako **prawdopodobieństwo częstotliwościowe**, natomiast ten drugi, który wiąże się z jakościowymi poziomami pewności, znany jest jako **prawdopodobieństwo bayesowskie**.

Jeśli wymienimy listę kilku właściwości zdroworozsądkowego rozumowania o niepewności, to jedynym sposobem, aby je spełnić, jest traktowanie prawdopodobieństwa bayesowskiego jako zachowującego się identycznie jak prawdopodobieństwo częstotliwościowe. Jeśli na przykład chcemy wyznaczyć

prawdopodobieństwo wygranej w pokera, zakładając, że gracz ma określony zestaw kart, używamy dokładnie tych samych wzorów, jak w przypadku wyznaczania prawdopodobieństwa, iż pacjent ma daną chorobę, pod warunkiem że ma określone objawy. Więcej przykładów tego, że mały zbiór zdroworozsądkowych założeń pozwala stosować te same aksjomaty do obu rodzajów prawdopodobieństwa, można znaleźć w książce Ramsey'a (1926).

Można patrzeć na prawdopodobieństwo jak na logiczny układ z niepewnością. Logika podaje zbiór formalnych reguł do określenia, które propozycje mogą być prawdziwe, a które fałszywe. Teoria prawdopodobieństwa daje zbiór formalnych reguł do określenia wiarygodności tezy, że propozycja jest prawdziwa, z uwagi na prawdopodobieństwo innych propozycji.

3.2. Zmienne losowe

Zmienna losowa to zmienna, która może losowo przyjmować różne wartości. Zwykle oznaczana jest małą literą pisaną zwykłą czcionką, a przyjmowane przez nią wartości przyjmują postać indeksów dolnych. Na przykład x_1 i x_2 to dwie możliwe wartości zmiennej losowej x . Dla zmiennych wektorowych zapiszemy zmienną losową jako \mathbf{x} , a jedną z jego wartości jako \mathbf{x} . Pojedyncza zmienna losowa opisuje po prostu możliwe stany; musi być powiązana z rozkładem prawdopodobieństwa, który określa, jaka jest możliwość wystąpienia każdego z tych stanów.

Zmienne losowe mogą być dyskretne lub ciągłe. Dyskretna zmienna losowa to taka, która ma skończoną lub też nieskończoną, ale policzalną liczbę stanów. Warto zwrócić uwagę, że stany te nie muszą być wartościami całkowitymi; mogą też być nazwane, więc wcale nie muszą mieć wartości liczbowej. Ciągła zmienna losowa jest powiązana z wartością rzeczywistą.

3.3. Rozkłady prawdopodobieństwa

Rozkład prawdopodobieństwa to opis, jaka jest możliwość ustawienia zmiennej losowej lub wielu zmiennych losowych na każdą możliwą wartość. Sposób opisu rozkładu zależy od tego, czy zmienne są dyskretne, czy ciągłe.

3.3.1. Zmienne dyskretne i funkcje masy prawdopodobieństwa

Rozkład prawdopodobieństwa względem zmiennych dyskretnych można opisać za pomocą **funkcji masy prawdopodobieństwa** (PMF). Zwykle oznaczamy funkcje masy prawdopodobieństwa wielką literą P . Często łączymy

każdą zmienną losową z różną funkcją masy prawdopodobieństwa i czytelnik musi wywnioskować, której PMF użyć na podstawie tożsamości zmiennej losowej, a nie na podstawie nazwy funkcji; $P(x)$ jest zwykle różna od $P(y)$.

Funkcja masy prawdopodobieństwa odwzorowuje stan zmiennej losowej na prawdopodobieństwo przyjęcia przez nią danego stanu. Prawdopodobieństwo, że $x = x$, jest oznaczone jako $P(x)$, gdzie prawdopodobieństwo 1 wskazuje, że równość $x = x$ jest pewna, natomiast 0 wskazuje, że równość $x = x$ jest niemożliwa*. Czasami, aby ujednoznaczyć, którą PMF zastosować, piszemy nazwę zmiennej losowej jawnie $P(x = x)$. Czasami najpierw definiujemy zmienną, a następnie używamy zapisu \sim , aby określić, jaki ma ona rozkład: $x \sim P(x)$.

Funkcja masy prawdopodobieństwa może w tym samym czasie działać na wielu zmiennych. Taki rozkład prawdopodobieństwa znany jest jako **wspólny rozkład prawdopodobieństwa**. $P(x = x, y = y)$ wskazuje prawdopodobieństwo, że jednocześnie $x = x$ i $y = y$. Zapis możemy skrócić do postaci $P(x, y)$.

Aby funkcja P była PMF względem zmiennej x , musi mieć następujące właściwości:

- dziedzina P musi być zbiorem wszystkich możliwych stanów x
- $\forall x \in x, 0 \leq P(x) \leq 1$ – zdarzenie niemożliwe ma prawdopodobieństwo równe 0 i żaden stan nie może mieć mniejszego prawdopodobieństwa; podobnie zdarzenie z pewnością się wydarzy, jeśli ma prawdopodobieństwo równe 1, a żaden stan nie może mieć większej szansy zdarzenia się
- $\sum_{x \in x} P(x) = 1$ – tę właściwość traktujemy jako **znormalizowaną**; bez niej można by otrzymać prawdopodobieństwa większe od 1, jeśli obliczamy prawdopodobieństwo zajścia jednego z wielu zdarzeń.

Rozważmy na przykład pojedynczą dyskretną zmienną losową x , która może mieć k stanów. Można wprowadzić **rozkład jednostajny** x – czyli każdy stan będzie jednakowo prawdopodobny – nadając jego funkcji PMF wartość:

$$P(x = x_i) = \frac{1}{k} \tag{3.1}$$

dla każdego i . Możemy zobaczyć, że wszystkie wymagania funkcji masy prawdopodobieństwa są spełnione. Wartość $\frac{1}{k}$ jest dodatnia, gdyż k jest liczbą całkowitą dodatnią. Widzimy także, że:

$$\sum_i P(x = x_i) = \sum_i \frac{1}{k} = \frac{k}{k} = 1, \tag{3.2}$$

więc rozkład jest poprawnie znormalizowany.

*Zwykle mówimy „bardzo mało prawdopodobne”, gdyż najczęściej wartość 0 jest wartością graniczną.

3.3.2. Zmienne ciągłe i funkcje gęstości prawdopodobieństwa

Podczas pracy z ciągłymi zmiennymi losowymi opisujemy rozkład prawdopodobieństwa za pomocą **funkcji gęstości prawdopodobieństwa** (PDF), a nie za pomocą funkcji masy prawdopodobieństwa. Aby funkcja p była funkcją gęstości prawdopodobieństwa, musi mieć następujące właściwości:

- dziedziną p musi być zbiór wszystkich możliwych stanów x
- $\forall x \in x, p(x) \geq 0$ – zauważmy, że nie ma wymogu, aby $p(x) \leq 1$
- $\int p(x)dx = 1$.

Funkcja gęstości prawdopodobieństwa $p(x)$ nie daje bezpośrednio prawdopodobieństwa określonego stanu. Zamiast tego podaje prawdopodobieństwo znalezienia się wewnątrz nieskończenie małego obszaru, którego objętość δx podana jest jako $p(x)\delta x$.

Można scałkować funkcję gęstości, aby znaleźć masę prawdopodobieństwa dla zbioru punktów. W szczególności prawdopodobieństwo, że x leży w pewnym zbiorze \mathbb{S} , jest podane jako całka z $p(x)$ względem tego zbioru. W przykładzie jednowymiarowym prawdopodobieństwo, że x należy do przedziału $[a, b]$ jest określone jako $\int_{[a,b]} p(x)dx$.

Dla funkcji PDF odpowiadającej określonej gęstości prawdopodobieństwa względem ciągłej zmiennej losowej rozpatrzmy rozkład jednostajny przedziału liczb rzeczywistych. Można to zrobić za pomocą funkcji $u(x; a, b)$, gdzie a i b to końce przedziału, przy czym $b > a$. Notacja z przecinkiem oznacza „względem parametru”; rozpatrujemy x jako argument funkcji, podczas gdy a i b są parametrami definiującymi funkcję. Aby sprawić, że żadna masa prawdopodobieństwa nie znajduje się poza przedziałem, mówimy, że $u(x; a, b) = 0$ dla każdego $x \notin [a, b]$. Wewnątrz $[a, b]$, $u(x; a, b) = \frac{1}{b-a}$. Jest ona wszędzie nie mniejsza od 0, ponadto całka wynosi 1. Często oznaczamy, że x ma rozkład jednostajny w przedziale $[a, b]$, zapisując to jako $x \sim U(a, b)$.

3.4. Prawdopodobieństwo brzegowe

Czasami znamy rozkład prawdopodobieństwa względem zbioru zmiennych i chcemy poznać rozkład prawdopodobieństwa dla ich podzbioru. Rozkład prawdopodobieństwa dla podzbioru określany jest jako **brzegowy rozkład prawdopodobieństwa**. Przypuśćmy na przykład, że mamy dyskretne zmienne losowe x i y oraz znamy $P(x, y)$. Możemy znaleźć $P(x)$ wg **reguły sum**:

$$\forall x \in x, P(x = x) = \sum_y P(x = x, y = y). \quad (3.3)$$

Nazwa „prawdopodobieństwa brzegowe” pochodzi od procesu jego obliczana na papierze. Gdy wartości $P(x, y)$ zapisujemy na siatce z różnymi wartościami x w wierszach i y w kolumnach, w naturalny sposób sumujemy je wzdłuż wiersza siatki, a potem zapisujemy $P(x)$ na brzegu kartki papieru na prawo od wiersza. Dla zmiennych ciągłych musimy korzystać z całkowania w miejsce sumowania:

$$p(x) = \int p(x, y) dy. \quad (3.4)$$

3.5. Prawdopodobieństwo warunkowe

W wielu przypadkach interesuje nas prawdopodobieństwo jakiegoś zdarzenia, pod warunkiem że zaszło już inne zdarzenie. Określamy je jako **prawdopodobieństwo warunkowe**. Oznaczamy je tak, że $y = y$, wiedząc, że $x = x$ jako $P(y = y \mid x = x)$. Takie prawdopodobieństwo warunkowe możemy wyznaczyć wg wzoru:

$$P(y = y \mid x = x) = \frac{P(y = y, x = x)}{P(x = x)}. \quad (3.5)$$

Prawdopodobieństwo warunkowe jest zdefiniowane tylko wtedy, gdy $P(x = x) > 0$; nie możemy obliczyć prawdopodobieństwa warunkowego uwarunkowanego zdarzeniem, które nigdy nie zajdzie.

Ważne jest, aby nie pomylić prawdopodobieństwa warunkowego z wyznaczeniem tego, co się zdarzy, jeśli podejmiemy jakieś działanie. Prawdopodobieństwo warunkowe, że osoba pochodzi z Niemiec, jeżeli mówi po niemiecku, jest dość wysokie, ale jeśli losowo wybrana osoba zostanie nauczona mówić po niemiecku, jej kraj pochodzenia się nie zmieni. Wyznaczenie konsekwencji działania określa się jako wykonanie **interwencyjnego zapytania**. Zapytania takie stanowią domenę **modelowania przyczynowego** i nie są przedmiotem tej książki.

3.6. Reguła łańcuchowa w prawdopodobieństwie warunkowym

Każdy wspólny rozkład prawdopodobieństwa dla wielu zmiennych losowych może zostać rozłożony na rozkłady warunkowe względem tylko jednej zmiennej:

$$P(x^{(1)}, \dots, x^{(n)}) = P(x^{(1)}) \prod_{i=2}^n P(x^{(i)} \mid x^{(1)}, \dots, x^{(i-1)}). \quad (3.6)$$

Nosi to nazwę **reguły łańcuchowej** lub **reguły iloczynowej prawdopodobieństwa** i wynika bezpośrednio z prawdopodobieństwa warunkowego zawartego we wzorze 3.5. Na przykład dwukrotne zastosowanie tej definicji da poniższy wynik:

$$\begin{aligned} P(a, b, c) &= P(a \mid b, c)P(b, c) \\ P(b, c) &= P(b \mid c)P(c) \\ P(a, b, c) &= P(a \mid b, c)P(b \mid c)P(c) \end{aligned}$$

3.7. Niezależność oraz niezależność warunkowa

Dwie zmienne losowe x i y są **niezależne**, jeśli ich rozkład prawdopodobieństwa można wyrazić jako iloczyn dwóch czynników, gdzie jeden dotyczy tylko x , a drugi tylko y :

$$\forall x \in X, y \in Y, p(x = x, y = y) = p(x = x)p(y = y). \quad (3.7)$$

Dwie zmienne losowe są **warunkowo niezależne** względem zmiennej losowej z , jeśli prawdopodobieństwo warunkowe względem x i y rozkłada się na czynniki dla każdej wartości z wg zależności:

$$\forall x \in X, y \in Y, z \in Z, p(x = x, y = y \mid z = z) = p(x = x \mid z = z)p(y = y \mid z = z). \quad (3.8)$$

Możemy oznaczyć niezależność i warunkową niezależność za pomocą zwartego zapisu $x \perp y$, co oznacza, że x i y są niezależne, natomiast $x \perp y \mid z$ oznacza, że x i y są warunkowo niezależne dla z .

3.8. Wartość oczekiwana, wariancja i kowariancja

Wartość oczekiwana pewnej funkcji $f(x)$ względem rozkładu prawdopodobieństwa $P(x)$ to wartość średnia, którą f przyjmuje, gdy x zostanie wylosowane z P . Dla zmiennych dyskretnych można ją wyznaczyć jako sumę:

$$\mathbb{E}_{x \sim P}[f(x)] = \sum_x P(x)f(x), \quad (3.9)$$

natomiast dla zmiennych ciągłych obliczamy ją za pomocą całki:

$$\mathbb{E}_{x \sim P}[f(x)] = \int p(x)f(x)dx. \quad (3.10)$$

Gdy tożsamość rozkładu wynika jasno z kontekstu, można po prostu napisać nazwę zmiennej losowej, względem której liczymy wartość oczekiwaną, jako $\mathbb{E}_x[f(x)]$. Jeśli jest jasne, której zmiennej losowej dotyczy wartość oczekiwana, można całkowicie pominąć indeks dolny jak w wyrażeniu $\mathbb{E}[f(x)]$. Domyślnie przyjmuje się, że $\mathbb{E}[\cdot]$ jest średnią względem wartości wszystkich zmiennych losowych wewnątrz nawiasu. Podobnie gdy nie ma niejednoznaczności, można pominąć nawiasy kwadratowe.

Wartości oczekiwane są liniowe, jak na przykład:

$$\mathbb{E}_x[\alpha f(x) + \beta g(x)] = \alpha \mathbb{E}_x[f(x)] + \beta \mathbb{E}_x[g(x)], \quad (3.11)$$

gdy α i β nie są zależne od x .

Wariancja określa miarę tego, jak wartości zmiennej losowej x będą się różnić przy próbkowaniu różnych wartości x z jej rozkładu prawdopodobieństwa:

$$\text{Var}(f(x)) = \mathbb{E} \left[(f(x) - \mathbb{E}[f(x)])^2 \right]. \quad (3.12)$$

Gdy wariancja jest niska, wartości $f(x)$ skupiają się w pobliżu wartości oczekiwanej. Pierwiastek kwadratowy z wariancji określany jest jako **odchylenie standardowe**.

Kowariancja daje nam miarę tego, na ile dwie wartości są liniowo ze sobą powiązane, a także skalę tych zmiennych:

$$\text{Cov}(f(x), g(y)) = \mathbb{E} [(f(x) - \mathbb{E}[f(x)]) (g(y) - \mathbb{E}[g(y)])]. \quad (3.13)$$

Duże wartości bezwzględne kowariancji oznaczają, że wartości ulegają dużym zmianom i są jednocześnie daleko od swoich odpowiednich średnich. Jeśli znak kowariancji jest dodatni, to obie zmienne będą miały jednocześnie względnie duże wartości. Jeśli znak kowariancji jest ujemny, to jedna zmienna ma relatywnie duże wartości, a jednocześnie druga zmienna relatywnie małe. Inne miary, jak **korelacja**, normalizują wkład każdej zmiennej, pozwalając tylko stwierdzić, w jakim stopniu zmienne są ze sobą powiązane, nie ma natomiast na nie wpływu wielkość pojedynczych zmiennych.

Pojęcia kowariancji i zależności są ze sobą związane, ale należy je traktować odrębnie. Są powiązane, gdyż dwie zmienne niezależne mają kowariancję równą 0, a zmienne o niezerowej kowariancji są zależne od siebie. Jednak niezależność to co innego niż kowariancja. Aby dwie zmienne miały zerową kowariancję, nie może być między nimi zależności liniowej. Niezależność jest mocniejszym wymaganiem niż zerowa wartość kowariancji, gdyż niezależność wyklucza także związki nieliniowe. Możliwe jest, aby dwie zmienne były zależne, ale miały zerową kowariancję. Przypuśćmy na przykład, że pobieramy pierwszą próbkę liczby rzeczywistej x z rozkładem jednostajnym w przedziale

$[-1, 1]$. Następna próbka to zmienna losowa s . Z prawdopodobieństwem $\frac{1}{2}$ określamy, że wartość s wynosi 1. W przeciwnym przypadku określamy, że wartością s jest -1 . Możemy wygenerować zmienną losową y , przypisując jej $y = sx$. Jasno widać, że x i y nie są niezależne, gdyż x w pełni określa wielkość y . Jednak $\text{Cov}(x, y) = 0$.

Macierz kowariancji losowego wektora $x \in \mathbb{R}^n$ jest macierzą $n \times n$, jak na przykład:

$$\text{Cov}(\mathbf{x})_{i,j} = \text{Cov}(x_i, x_j). \quad (3.14)$$

Elementy diagonalne kowariancji dają wariancję:

$$\text{Cov}(x_i, x_i) = \text{Var}(x_i). \quad (3.15)$$

3.9. Znane rozkłady prawdopodobieństwa

Poniżej kilka prostych rozkładów prawdopodobieństwa, przydatnych w kontekście systemów uczących się.

3.9.1. Rozkład Bernoulliego

Rozkład Bernoulliego to rozkład jednej binarnej zmiennej losowej. Jest zależny od jednego parametru $\phi \in [0, 1]$, który daje prawdopodobieństwo, że zmienna losowa przyjmie wartość 1. Rozkład ma następujące cechy:

$$P(x = 1) = \phi \quad (3.16)$$

$$P(x = 0) = 1 - \phi \quad (3.17)$$

$$P(x = x) = \phi^x (1 - \phi)^{1-x} \quad (3.18)$$

$$\mathbb{E}_x[x] = \phi \quad (3.19)$$

$$\text{Var}_x(x) = \phi(1 - \phi) \quad (3.20)$$

3.9.2. Rozkład wielopunktowy

Rozkład wielopunktowy (inaczej kategoryczny) jest rozkładem jednej zmiennej dyskretnej o k różnych stanach, gdzie k ma wartość skończoną¹. Rozkład

¹„Rozkład wielopunktowy” to termin ukuty przez Gustava Lacerdo i spopularyzowany przez Murphy’ego (2012). Rozkład wielopunktowy jest specjalnym przypadkiem rozkładu wielomianowego. Rozkład wielomianowy to rozkład względem wektorów w $\{0, \dots, n\}^k$, reprezentującej, ile razy każda z kategorii k jest odwiedzana, gdy pobierzemy n próbek z rozkładu wielopunktowego. W wielu tekstach termin „wielomianowy” oznacza rozkład wielopunktowy bez określenia, że odnosi się tylko do przypadku $n = 1$.

wielopunktowy jest parametryzowany przez wektor $p \in [0, 1]^{k-1}$, gdzie p_i podaje prawdopodobieństwo i -tego stanu. Końcowe prawdopodobieństwo stanu k -tego jest określone jako $1 - \mathbf{1}^\top p$. Zwróćmy uwagę, że musimy ograniczyć $\mathbf{1}^\top p \leq 1$. Rozkłady wielopunktowe są często używane w celu odwołania się do rozkładu względem kategorii obiektów, więc zwykle nie zakładamy, że stan 1 ma wartość liczbową 1 i tak dalej. Dlatego nie musimy zwykle wyznaczać wartości oczekiwanej ani wariancji zmiennych losowych o rozkładzie wielopunktowym.

Rozkłady Bernoulliego i wielopunktowy wystarczają do opisanie dowolnego rozkładu w ich dziedzinie. Mogą opisać każdy rozkład, ale nie dlatego, że mają szczególnie duże możliwości, lecz dlatego, że mamy do czynienia z prostą dziedziną: modelują zmienne dyskretne, w których można policzyć wszystkie możliwe stany. Gdy mamy do czynienia ze zmiennymi ciągłymi, istnieje niepoliczalna liczba stanów, więc każdy rozkład opisywany przez niewielką liczbę parametrów musi nakładać ściśle ograniczenia na rozkład.

3.9.3. Rozkład Gaussa

Najczęściej używanym rozkładem względem liczb rzeczywistych jest **rozkład normalny**, znany jako **rozkład Gaussa**:

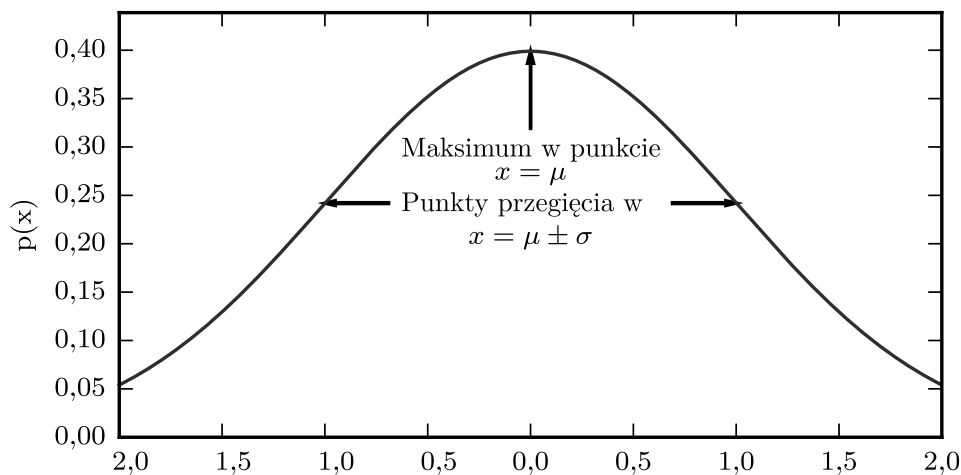
$$\mathcal{N}(x; \mu, \sigma^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right). \quad (3.21)$$

Na rysunku 3.1 pokazano wykres funkcji gęstości rozkładu normalnego.

Dwa parametry $\mu \in \mathbb{R}$ i $\sigma \in (0, \infty)$ sterują rozkładem normalnym. Parametr μ daje współrzędną punktu szczytowego. Jest to również średnia arytmetyczna rozkładu: $\mathbb{E}[x] = \mu$. Odchylenie standardowe rozkładu jest oznaczone jako σ , a wariancja jako σ^2 . Gdy wyznaczamy wartość PDF, musimy podnieść do kwadratu i odwrócić σ . Gdy musimy często wyznaczać wartość PDF dla różnych wartości parametru, bardziej efektywnym sposobem parametryzacji rozkładu jest użycie parametru $\beta \in (0, \infty)$ do kontroli **precyzji** lub odwrotnej wariancji rozkładu:

$$\mathcal{N}(x; \mu, \beta^{-1}) = \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{1}{2}\beta(x - \mu)^2\right). \quad (3.22)$$

Rozkłady normalne są dobrym wyborem w wielu zastosowaniach. Gdy nie znamy formy rozkładu względem liczb rzeczywistych, rozkład normalny to dobry wybór domyślny z dwóch podstawowych powodów. Po pierwsze wiele rozkładów, które chcemy modelować, ma rozkład naprawdę bliski rozkładu



Rysunek 3.1. Rozkład normalny $N(x; \mu, \sigma^2)$ ma postać klasycznej krzywej dzwonekowej, gdzie współrzędna x szczytu jest określona przez μ , a szerokość krzywej jest dana przez σ . W tym przykładzie pokazujemy standardowy rozkład normalny z wartościami $\mu = 0$ i $\sigma = 1$

normalnego. **Centralne twierdzenie graniczne** pokazuje, że suma wielu niezależnych zmiennych losowych ma średni rozkład normalny. W praktyce oznacza to, że wiele skomplikowanych systemów można z powodzeniem modelować jako szum o rozkładzie normalnym, jeśli nawet można je dekomponować na części o bardziej strukturalnym zachowaniu. Po drugie spośród wszystkich możliwych rozkładów prawdopodobieństwa o tej samej wariancji rozkład normalny koduje najwięcej niepewności względem liczb rzeczywistych. Możemy więc traktować rozkład normalny jako ten, który wykorzystuje w modelu najmniej wstępnej wiedzy. Pełne rozwinięcie i uzasadnienie tej tezy wymaga więcej narzędzi matematycznych i wrócimy do niego w punkcie 19.4.2.

Rozkład normalny jest uogólniony jako \mathbb{R}^n , co w tym przypadku określamy jako **rozkład normalny jednoczynnikowy**. Można go sparametryzować za pomocą symetrycznej macierzy dodatnio określonej Σ :

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^n \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (3.23)$$

Parametr $\boldsymbol{\mu}$ nadal podaje średnią arytmetyczną rozkładu, ale teraz ma wartość wektora. Parametr $\boldsymbol{\Sigma}$ podaje macierzy kowariancję rozkładu. Podobnie jak w przypadku jednoczynnikowym, gdy chcemy wyznaczyć PDF kilka razy dla różnych wartości parametrów, kowariancja nie jest efektywnym obliczeniowo sposobem parametryzacji rozkładu, gdyż w celu obliczenia PDF

musimy odwrócić Σ . Zamiast tego możemy użyć **macierzy precyzji** β :

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\beta}^{-1}) = \sqrt{\frac{\det(\boldsymbol{\beta})}{(2\pi)^n}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\beta}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (3.24)$$

Często ustala się macierz kowariancji jako macierz diagonalną. Jeszcze prostszy jest **równokierunkowy** rozkład Gaussa, którego kowariancja jest skalarem pomnożonym przez macierz jednostkową.

3.9.4. Rozkłady wykładniczy i Laplace’a

W kontekście deep learningu często chcemy uzyskać rozkład prawdopodobieństwa z ostrym punktem dla $x = 0$. Aby to osiągnąć, możemy wykorzystać **rozkład wykładniczy**:

$$p(x; \lambda) = \lambda \mathbf{1}_{x \geq 0} \exp(-\lambda x). \quad (3.25)$$

Rozkład wykładniczy wykorzystuje funkcję wskaźnikową $\mathbf{1}_{x \geq 0}$, aby przypisać zerowe prawdopodobieństwo wszystkim ujemnym wartościom x .

Blisko spokrewniony z nim rozkład prawdopodobieństwa, który pozwala umieścić ostry szczyt masy prawdopodobieństwa w arbitralnie wybranym punkcie μ , to **rozkład Laplace’a**:

$$\text{Laplace}(x; \mu, \gamma) = \frac{1}{2\gamma} \exp\left(-\frac{|x - \mu|}{\gamma}\right). \quad (3.26)$$

3.9.5. Rozkład Diraca i rozkład empiryczny

W niektórych przypadkach chcemy określić, że cała masa rozkładu prawdopodobieństwa skupia się wokół jednego punktu. Można to osiągnąć, definiując PDF za pomocą **delty Diraca** $\delta(x)$:

$$p(x) = \delta(x - \mu). \quad (3.27)$$

Funkcja delta Diraca jest zdefiniowana w ten sposób, że ma wartość 0 wszędzie poza zerem, a jej całka wynosi 1. Funkcja ta nie jest zwykłą funkcją, która łączy każdą wartość x z rzeczywistym wynikiem. Jest innym rodzajem obiektu matematycznego określanym jako **funkcja uogólniona** zdefiniowana w sensie jej właściwości podczas całkowania. Możemy traktować deltę Diraca jako punkt graniczny ciągu funkcji, które przypisują coraz mniej masy wszystkim punktom poza 0.

Przez zdefiniowanie $p(x)$ jako wartości δ przesuniętej o $-\mu$ otrzymujemy nieskończenie wąski szczyt masy prawdopodobieństwa, gdzie $x = \mu$.

Powszechnym zastosowaniem rozkładu delta Diraca jest użycie go jako składnika **rozkładu empirycznego**:

$$\hat{p}(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m \delta(\mathbf{x} - \mathbf{x}^{(i)}). \quad (3.28)$$

Przypisuje on masę prawdopodobieństwa $\frac{1}{m}$ do każdego z m punktów $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}$ na podstawie zbioru danych lub kolekcji próbek. Rozkład delta Diraca jest konieczny jedynie do zdefiniowania rozkładu empirycznego względem zmiennych ciągłych. Dla zmiennych dyskretnych sytuacja jest prostsza: rozkład empiryczny można określić jako rozkład wielopunktowy z prawdopodobieństwem związanym z każdą możliwą wartością wejściową, która jest równa **empirycznej częstotliwości** tej wartości w zbiorze szkoleniowym.

Można traktować rozkład empiryczny utworzony ze zbioru danych przykładów testowych jako określenie rozkładu próbkowanego podczas szkolenia modelu na tym zbiorze danych. Inna ważną perspektywą dla rozkładu empirycznego jest fakt, że jest to gęstość prawdopodobieństwa maksymalizująca wiarygodność danych szkoleniowych (patrz punkt 5.5).

3.9.6. Rozkłady mieszane

Często definiuje się rozkłady prawdopodobieństwa, łącząc ze sobą prostsze rozkłady. Jednym ze sposobów łączenia rozkładów jest budowa **rozkładu mieszanego**. Składa się on z kilku rozkładów składowych. Przy każdej próbie wybór, którego składnika użyć do wygenerowania próbki, jest określony przez próbkowanie tożsamości składnika z rozkładu wielopunktowego:

$$P(\mathbf{x}) = \sum_i P(c = i)P(\mathbf{x} | c = i), \quad (3.29)$$

gdzie $P(c)$ to rozkład wielopunktowy względem tożsamości składników.

Widzieliśmy już jeden przykład rozkładu mieszanego – rozkład empiryczny względem wartości rzeczywistych to rozkład mieszany z jednym składnikiem Diraca dla każdego przykładu szkoleniowego.

Model mieszany jest prostą strategią łączenia rozkładów prawdopodobieństwa, aby utworzyć rozkład bogatszy. W rozdziale 16 bardziej szczegółowo analizujemy sztukę budowania złożonych rozkładów prawdopodobieństwa na podstawie tych prostych.

Model mieszany pozwala pobieżnie zapoznać się z pojęciem, które będzie później miało ogromne znaczenie – **zmiennej utajonej**. Zmienna utajona to zmienna losowa, której nie można bezpośrednio obserwować. Składowa

zmienna tożsamościowa c w modelu mieszanym posłuży jako przykład. Utajone zmienne mogą odnosić się do x za pośrednictwem wspólnego rozkładu, w tym przypadku $P(x, c) = P(x|c)P(c)$. Rozkład $P(c)$ względem zmiennej utajonej oraz rozkład $P(x|c)$ wiążący zmienne utajone ze zmiennymi widocznymi określa kształt rozkładu $P(x)$, mimo że można opisać $P(x)$ bez odwoływania się do zmiennej utajonej. Zmienne utajone są szerzej omawiane w punkcie 16.5.

Często stosowanym modelem mieszanym o dużych możliwościach jest **model mieszaný Gaussa**, w którym składniki $p(x | c = i)$ są gaussowskie. Każdy składnik ma oddzielnie parametryzowaną średnią arytmetyczną $\mu^{(i)}$ oraz kowariancję $\Sigma^{(i)}$. Niektóre mieszanki mogą mieć więcej ograniczeń, na przykład określa się, że kowariancje powinny być wspólne dla składników przez ograniczenie $\Sigma^{(i)} = \Sigma, \forall i$. Podobnie jak w pojedynczym rozkładzie Gaussa, mieszanka rozkładów gaussowskich może ograniczać macierz kowariancji, wymuszając diagonalność lub izotropowość dla każdego ze swoich składników.

Poza średnimi arytmetycznymi i kowariancjami parametry mieszanki gaussowskiej określają **rozkład aprioryczny** (*a priori*) $\alpha_i = P(c = i)$ dla każdego składnika i . Określenie *a priori* wskazuje, że wyrażone są oczekiwania modelu dotyczące c zanim nastąpi obserwacja \mathbf{x} . Analogicznie $P(c | \mathbf{x})$ to prawdopodobieństw *a posteriori*, gdyż jest wyznaczane *po* obserwacji \mathbf{x} . Model mieszaný Gaussa to **metoda uniwersalnej aproksymacji** gęstości w tym sensie, że każda gładka gęstość może zostać aproksymowana z określonym niezerowym błędem przez mieszaný model Gaussa z dostateczną liczbą składników. Na rysunku 3.2 pokazano próbki z mieszanego modelu Gaussa.

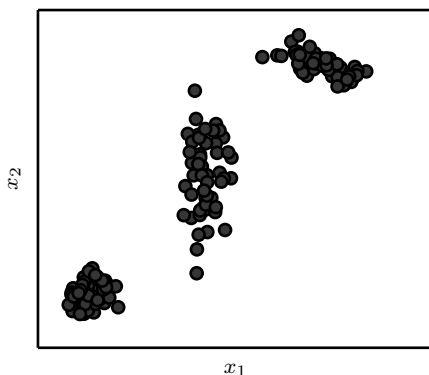
3.10. Użyteczne cechy elementarnych funkcji

Często podczas pracy z rozkładami prawdopodobieństwa, zwłaszcza rozkładami prawdopodobieństwa używanymi w modelach deep learning, pojawiają się pewne funkcje. Jedną z nich jest funkcja logistyczno-sigmoidalna:

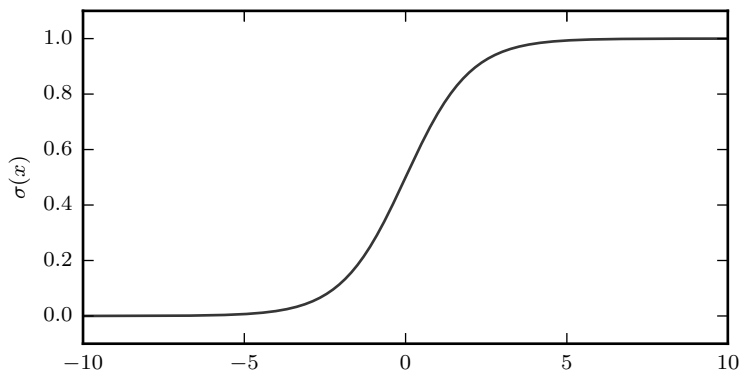
$$\sigma(x) = \frac{1}{1 + \exp(-x)}. \quad (3.30)$$

Sigmoida logistyczna jest powszechnie używana do pokazania parametru ϕ rozkładu Bernoulliego, gdyż jej zakresem jest $(0, 1)$, co odpowiada dopuszczalnemu zakresowi wartości parametru ϕ . Funkcję tę pokazano na rysunku 3.3.

Funkcja ulega nasyceniu, gdy jej argument jest mocno dodatni lub mocno ujemny, co oznacza, że funkcja staje się bardzo płaska i nieczuła na niewielkie zmiany argumentu wejściowego.



Rysunek 3.2. Przykład mieszanego modelu Gaussa. W tym przykładzie mamy trzy składowe. Od lewej do prawej: pierwszy składnik ma izotropową macierz kowariancji, co oznacza, że ma tę samą wartość wariancji w każdym kierunku; drugi ma diagonalną macierz kowariancji, co oznacza, że można sterować wariancjami oddzielnie wzdłuż każdej osi w danym kierunku (w tym przykładzie mamy większą wariancję wzdłuż osi x_2 niż wzdłuż osi x_1); trzeci składnik ma pełną macierz kowariancji, co umożliwia oddzielne sterowanie wariancją wzdłuż dowolnie wybranych podstawowych kierunków



Rysunek 3.3. Sigmoida logistyczna

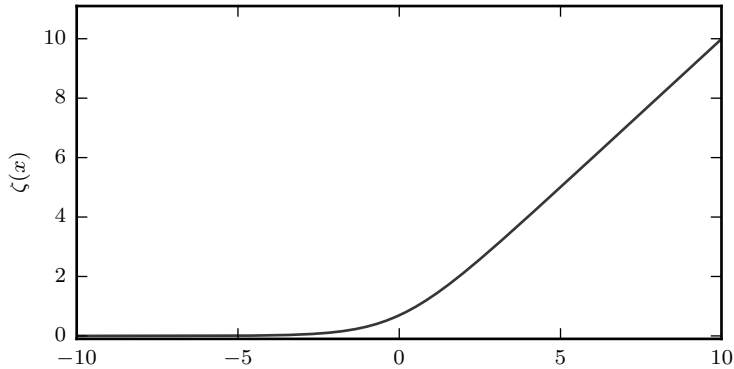
Inną często spotykaną funkcją jest **funkcja softplus** (Dugas et al. 2001):

$$\zeta(x) = \log(1 + \exp(x)). \quad (3.31)$$

Funkcja softplus może być wykorzystana do tworzenia parametru β lub σ rozkładu normalnego, gdyż jej zakresem jest przedział $(0, \infty)$. Pojawia się też często podczas przekształcania wyrażeń związanych z sigmoidami. Nazwa „softplus” wynika stąd, że jest wygładzoną, czyli „zmiękczoną” wersją funkcji:

$$x^+ = \max(0, x). \quad (3.32)$$

Jej wykres widać na rysunku 3.4.



Rysunek 3.4. Funkcja softplus

Poniższe właściwości są także użyteczne, więc warto je zapamiętać:

$$\sigma(x) = \frac{\exp(x)}{\exp(x) + \exp(0)} \quad (3.33)$$

$$\frac{d}{dx}\sigma(x) = \sigma(x)(1 - \sigma(x)) \quad (3.34)$$

$$1 - \sigma(x) = \sigma(-x) \quad (3.35)$$

$$\log \sigma(x) = -\zeta(-x) \quad (3.36)$$

$$\frac{d}{dx}\zeta(x) = \sigma(x) \quad (3.37)$$

$$\forall x \in (0, 1), \sigma^{-1}(x) = \log\left(\frac{x}{1-x}\right) \quad (3.38)$$

$$\forall x > 0, \zeta^{-1}(x) = \log(\exp(x) - 1) \quad (3.39)$$

$$\zeta(x) = \int_{-\infty}^x \sigma(y)dy \quad (3.40)$$

$$\zeta(x) - \zeta(-x) = x \quad (3.41)$$

Funkcja $\sigma^{-1}(x)$ jest w statystyce nazywana **logit**, ale ten termin rzadko występuje w systemach uczących się.

W równaniu 3.41 podano dodatkowe uzasadnienie używania nazwy soft-plus. Jest to bowiem wygładzona wersja **dodatniej części funkcji** $x^+ = \max\{0, x\}$. Dodatnia część funkcji jest odpowiednikiem **ujemnej części funkcji** $x^- = \max\{0, -x\}$. Aby otrzymać gładką funkcję, która jest analogiczna do części ujemnej, można użyć $\zeta(-x)$. Podobnie jak x można otrzymać z jej dodatniej części i jej ujemnej części przez tożsamość $x^+ - x^- = x$, można też otrzymać x za pomocą tego samego związku między $\zeta(x)$ a $\zeta(-x)$, co pokazano w równaniu 3.41.

3.11. Prawo Bayesa

Często znamy $P(y | x)$, a chcemy poznać $P(x | y)$. Na szczęście, jeśli znamy także $P(x)$, możemy wyznaczyć pożądaną ilość za pomocą **prawa Bayesa**:

$$P(x | y) = \frac{P(x)P(y | x)}{P(y)}. \quad (3.42)$$

Zauważmy, że choć $P(y)$ pojawia się we wzorze, można zwykle obliczyć $P(y) = \sum_x P(y | x)P(x)$, więc nie trzeba wychodzić od znajomości $P(y)$.

Prawo Bayesa można bezpośrednio wyprowadzić z definicji prawdopodobieństwa warunkowego, ale warto znać nazwę tego wzoru, która pojawia się w wielu tekstach. Funkcję nazwano od nazwiska wielkiego Thomasa Bayesa, który pierwszy odkrył szczególny przypadek wzoru. Ogólna wersja, którą tu prezentujemy, została niezależnie opracowana przez Pierre-Simona Laplace'a.

3.12. Techniczne szczegóły zmiennych ciągłych

Poprawne formalne zrozumienie ciągłych zmiennych losowych i funkcji gęstości prawdopodobieństwa wymaga rozwinięcia teorii prawdopodobieństwa jako gałęzi matematyki znanej jako **teoria miary**. Teoria miary wykracza poza zakres tej książki, ale możemy krótko naszkicować zagadnienia, jakie teoria miary pomaga rozwiązywać.

W punkcie 3.3.2 pokazano, że prawdopodobieństwo ciągłych wartości wektorowych \mathbf{x} leżących w pewnym zbiorze \mathbb{S} jest dane całką funkcji $p(\mathbf{x})$ na zbiorze \mathbb{S} . Niektóre z wybranych zbiorów \mathbb{S} mogą prowadzić do paradoksów. Na przykład można zbudować dwa zbiory \mathbb{S}_1 i \mathbb{S}_2 takie, że $p(\mathbf{x} \in \mathbb{S}_1) + p(\mathbf{x} \in \mathbb{S}_2) > 1$, ale $\mathbb{S}_1 \cap \mathbb{S}_2 = \emptyset$. Te zbiory są zbudowane z mocnym wykorzystaniem nieskończonej dokładności liczb rzeczywistych, na przykład przez tworzenie zbiorów o kształcie fraktalnym lub zbiorów, które są zdefiniowane przez transformację zbioru liczb wymiernych². Jednym z kluczowych wartości teorii miary jest zapewnienie charakterystyki zbioru zbiorów, dla których możemy wyznaczyć prawdopodobieństwo bez napotykania na paradoksy. W tej książce całkujemy tylko zbiory o względnie prostym opisie, więc aspekty teorii miary nie są nigdy przedmiotem większej troski.

Dla naszych celów teoria miary jest bardziej użyteczna do opisu twierdzeń, które odnoszą się do większości punktów w \mathbb{R}^n , ale nie mają zastosowania do niektórych przypadków brzegowych. Teoria miary podaje ścisły sposób opisu, kiedy zbiór punktów jest pomijalnie mały. Taki zbiór określa się jako zbiór

²Paradoks Banacha-Tarskiego daje dobry przykład tego rodzaju zbiorów.

z **miarą zero**. W tej książce nie przedstawiamy tego pojęcia formalnie. Dla naszych celów wystarczy, aby zrozumieć intuicyjnie, że zbiór z miarą zero nie zajmuje objętości w przestrzeni, która jest przedmiotem naszej analizy. Na przykład w ramach \mathbb{R}^2 prosta ma miarę zero, a wypełniony wielokąt ma miarę dodatnią. Podobnie pojedynczy punkt ma miarę zero. Każda suma policzalnej liczby zbiorów, z których każdy ma miarę zero, także ma miarę zero (więc zbiór wszystkich liczb wymiernych ma przykładowo miarę zero).

Inne użyteczne terminy z teorii miary można znaleźć **niemal wszędzie**. Właściwość, która ma sens niemal wszędzie, utrzymuje się w całej przestrzeni poza zbiorem mającym miarę 0. Ponieważ wyjątki zajmują pomijalną ilość przestrzeni, mogą w wielu zastosowaniach zostać bezpiecznie zignorowane. Niektóre ważne wyniki w teorii prawdopodobieństwa są prawdziwe dla wszystkich wartości dyskretnych, ale są określone tylko dla wartości ciągłych.

Inny techniczny szczegół zmiennych ciągłych wiąże się z obsługą ciągłych zmiennych losowych, które są funkcjami wzajemnie deterministycznymi. Przypuśćmy, że mamy dwie zmienne losowe \mathbf{x} i \mathbf{y} , takie, że $\mathbf{y} = g(\mathbf{x})$, gdzie g jest przekształceniem nieodwracalnym, ciągłym i różniczkowalnym. Można by się spodziewać, że $p_y(\mathbf{y}) = p_x(g^{-1}(\mathbf{y}))$. Ale okazuje się, że tak nie jest. Prosty przykład: przypuśćmy, że mamy skalarnie zmienne losowe x i y . Przypuśćmy dalej, że $y = \frac{x}{2}$ i $x \sim U(0, 1)$. Jeśli skorzystamy z zasady, że $p_y(y) = p_x(2y)$, to p_y będzie miało wartość 0 wszędzie poza przedziałem $[0, \frac{1}{2}]$, a w tym przedziale będzie miało wartość 1. Oznacza to:

$$\int p_y(y)dy = \frac{1}{2}, \tag{3.43}$$

co narusza definicję rozkładu prawdopodobieństwa. To częsty błąd. Problem w tym podejściu stanowi fakt, że nie bierze pod uwagę zakłócenia przestrzeni wprowadzonego przez funkcję g . Przypomnijmy, że prawdopodobieństwo, że \mathbf{x} leży w nieskończenie małym obszarze o objętości $\delta\mathbf{x}$, jest opisane przez $p(\mathbf{x})\delta\mathbf{x}$. Ponieważ g może rozszerzać lub zawężać przestrzeń, nieskończenie mała objętość otaczająca \mathbf{x} w przestrzeni \mathbf{x} może mieć inną objętość w przestrzeni \mathbf{y} .

Aby zobaczyć, jak naprawić ten problem, wróćmy do przypadku skalarnego. Musimy zachować właściwość:

$$|p_y(g(x))dy| = |p_x(x)dx|. \tag{3.44}$$

Rozwiązując ją, otrzymujemy:

$$p_y(y) = p_x(g^{-1}(y)) \left| \frac{\partial x}{\partial y} \right| \tag{3.45}$$

lub ekwiwalentnie:

$$p_x(x) = p_y(g(x)) \left| \frac{\partial g(x)}{\partial x} \right|. \quad (3.46)$$

Przy większej wymiarowości pochodnia uogólnia wyznacznik **macierzy Jacobiego** – macierzy przy $J_{i,j} = \frac{\partial x_i}{\partial y_j}$. Stąd dla wektorów \mathbf{x} i \mathbf{y} o wartościach rzeczywistych:

$$p_x(\mathbf{x}) = p_y(g(\mathbf{x})) \left| \det \left(\frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right) \right|. \quad (3.47)$$

3.13. Teoria informacji

Teoria informacji to gałąź matematyki stosowanej, która obraca się wokół określenia, ile informacji znajduje się w sygnale. Powstała w celu analizy wysyłania komunikatów w dyskretnych alfabetach przez kanał pełen szumów, jak transmisje radiowe. W tym kontekście teoria informacji mówi, jak zaprojektować optymalne kody i wyznaczyć oczekiwaną długość komunikatów próbkowanych według specjalnych rozkładów prawdopodobieństwa za pomocą różnych schematów kodowania. W kontekście systemów uczących się możemy zastosować teorię informacji do zmiennych ciągłych, w których nie mają zastosowania niektóre z tych interpretacji długości komunikatu. Ta dziedzina jest podstawą w wielu obszarach elektryki i nauk komputerowych. W książce skorzystamy z kilku kluczowych pojęć z teorii informacji, aby scharakteryzować rozkłady prawdopodobieństwa lub ocenić podobieństwa między rozkładami prawdopodobieństwa. Więcej szczegółów na temat teorii informacji można znaleźć w Cover and Thomas (2006) lub MacKay (2003).

Intuicyjnie rozumiemy teorię informacji tak, że wiedza o tym, iż zaszło mało prawdopodobne zdarzenie, niesie więcej informacji niż to, że zaszło zdarzenie prawdopodobne. Komunikat mówiący „dzisiaj weszło słońce” niesie tak mało informacji, że nie warto go wysyłać, natomiast komunikat „dzisiaj rano było zaćmienie słońca” ma dużą zawartość informacyjną.

Warto określić ilościowo informacje tak, aby sformalizować to intuicyjne podejście.

- Zdarzenia prawdopodobne powinny mieć małą zawartość informacyjną, a w skrajnych przypadkach zdarzenia, które na pewno zajdą, nie powinny mieć jakiegokolwiek zawartości informacyjnej.
- Zdarzenia mniej prawdopodobne powinny mieć wyższą zawartość informacyjną.
- Zdarzenia niezależne powinny być informacyjnie addytywne. Na przykład stwierdzenie, że rzut monetą dał dwa razy orła, powinno zawierać dwa razy więcej informacji niż stwierdzenie, że orzeł wypadł tylko raz.

Aby spełnić te trzy właściwości/własności, definiujemy **własną informację** zdarzenia $x = x$ jako:

$$I(x) = -\log P(x). \quad (3.48)$$

W tej książce \log zawsze oznacza logarytm naturalny, czyli o podstawie e . Nasza definicja $I(x)$ jest więc zapisana w jednostkach nat. Jeden nat to ilość informacji uzyskana przez zaobserwowanie zdarzenia o prawdopodobieństwie $\frac{1}{e}$. W innych tekstach używane są logarytmy o podstawie 2 oraz jednostki nazywane **bitami** (w literaturze anglosaskiej czasem występują jako **shannon**); informacja mierzona w bitach stanowi przeskalowanie informacji mierzonej w natach.

Gdy x jest ciągle, to przez analogię używamy tej samej definicji informacji, ale niektóre właściwości w przypadku dyskretnym zostają zagubione. Na przykład zdarzenie o jednostkowej gęstości ma nadal zerową informację, choć nie jest zdarzeniem, które zdarzy się na pewno.

Informacje własne dotyczą tylko pojedynczego wyniku. Możemy ocenić ilościowo poziom niepewności całego rozkładu prawdopodobieństwa za pomocą **entropii Shannona**:

$$H(x) = \mathbb{E}_{x \sim P}[I(x)] = -\mathbb{E}_{x \sim P}[\log P(x)], \quad (3.49)$$

oznaczanej także jako $H(P)$. Innymi słowy, entropia Shannona rozkładu to oczekiwana ilość informacji w zdarzeniu losowanym z tego rozkładu. Daje dolną granicę liczby bitów (jeśli logarytm ma podstawę 2, w innym przypadku jednostki są inne) potrzebną średnio do zakodowania symboli wylosowanych z rozkładu P . Rozkłady, które są niemal deterministyczne (gdzie wynik jest niemal pewny), mają niską entropię; rozkłady bliższe rozkładowi jednostajnego mają wysoką entropię. Pokazuje to rysunek 3.5. Gdy x jest ciągle, entropia Shannona określana jest jako **entropia różniczkowa**.

Jeśli mamy dwa oddzielne rozkłady prawdopodobieństwa $P(x)$ i $Q(x)$ względem tej samej zmiennej losowej x , możemy zmierzyć, na ile się one różnią, za pomocą **dywergencji Kullbacka-Leiblera (KL)**:

$$D_{\text{KL}}(P\|Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)]. \quad (3.50)$$

W przypadku zmiennych dyskretnych mamy dodatkową ilość informacji (mierzoną w bitach, jeśli stosujemy logarytm przy podstawie 2, choć w systemach uczących się stosuje się zwykle naty i logarytmy naturalne), potrzebną do wysłania komunikatu zawierającego symbole wyprowadzone z rozkładu prawdopodobieństwa P , gdy korzystamy z kodu opracowanego na potrzeby minimalizacji długości komunikatu na podstawie rozkładu prawdopodobieństwa Q .