

Uczenie maszynowe z użyciem Scikit-Learn i TensorFlow

POJĘCIA, TECHNIKI I NARZĘDZIA SŁUŻĄCE
DO TWORZENIA INTELIGENTNYCH SYSTEMÓW



Tytuł oryginału: Hands-On Machine Learning with Scikit-Learn and TensorFlow

Tłumaczenie: Krzysztof Sawka

ISBN: 978-83-283-4373-3

© 2018 Helion S.A.

Authorized Polish translation of the English edition of Hands-On Machine Learning with Scikit-Learn and TensorFlow ISBN 9781491962299 © 2017 Aurélien Géron

This translation is published and sold by permission of O'Reilly Media, Inc., which owns or controls all rights to publish and sell the same.

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

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

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

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

Wydawnictwo HELION

ul. Kościuszki 1c, 44-100 GLIWICE

tel. 32 231 22 19, 32 230 98 63

e-mail: helion@helion.pl

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

Drogi Czytelniku!

Jeżeli chcesz ocenić tę książkę, zajrzyj pod adres

<http://helion.pl/user/opinie/uczema>

Możesz tam wpisać swoje uwagi, spostrzeżenia, recenzję.

Printed in Poland.

- [Kup książkę](#)
- [Poleć książkę](#)
- [Oceń książkę](#)

- [Księgarnia internetowa](#)
- [Lubię to! » Nasza społeczność](#)

Przedmowa	13
-----------------	----

Część I. Podstawy uczenia maszynowego 21

1. Krajobraz uczenia maszynowego	23
Czym jest uczenie maszynowe?	24
Dlaczego warto korzystać z uczenia maszynowego?	24
Rodzaje systemów uczenia maszynowego	27
Uczenie nadzorowane/nienadzorowane	28
Uczenie wsadowe/przyrostowe	34
Uczenie z przykładów/z modelu	37
Główne problemy uczenia maszynowego	42
Niedobór danych uczących	42
Niereprezentatywne dane uczące	43
Dane kiepskiej jakości	45
Nieistotne cechy	45
Przetrenowanie danych uczących	45
Niedotrenowanie danych uczących	47
Podsumowanie	48
Testowanie i ocenianie	48
Ćwiczenia	50
2. Nasz pierwszy projekt uczenia maszynowego	53
Praca z rzeczywistymi danymi	53
Przeanalizuj całokształt projektu	55
Określ zakres problemu	55
Wybierz metrykę wydajności	57
Sprawdź założenia	59

Zdobądź dane	59
Stwórz przestrzeń roboczą	59
Pobierz dane	62
Rzut oka na strukturę danych	63
Stwórz zbiór testowy	67
Odkrywaj i wizualizuj dane, aby zdobywać nowe informacje	71
Wizualizowanie danych geograficznych	71
Poszukiwanie korelacji	73
Eksperymentowanie z kombinacjami atrybutów	76
Przygotuj dane pod algorytmy uczenia maszynowego	77
Oczyszczanie danych	78
Obsługa tekstu i atrybutów kategorialnych	80
Niestandardowe transformatory	81
Skalowanie cech	82
Potoki transformujące	83
Wybór i uczenie modelu	85
Trenowanie i ocena modelu za pomocą zbioru uczącego	85
Dokładniejsze ocenianie za pomocą sprawdzianu krzyżowego	86
Wyreguluj swój model	88
Metoda przeszukiwania siatki	88
Metoda losowego przeszukiwania	90
Metody zespołowe	91
Analizuj najlepsze modele i ich błędy	91
Oceń system za pomocą zbioru testowego	92
Uruchom, monitoruj i utrzymuj swój system	92
Teraz Twoja kolej!	93
Ćwiczenia	94
3. Klasyfikacja	95
Zbiór danych MNIST	95
Uczenie klasyfikatora binarnego	97
Miary wydajności	98
Pomiar dokładności za pomocą sprawdzianu krzyżowego	98
Macierz pomyłek	100
Precyzja i pełność	101
Kompromis pomiędzy precyzją a pełnością	102
Wykres krzywej ROC	106
Klasyfikacja wieloklasowa	108
Analiza błędów	111
Klasyfikacja wieloetykietowa	114
Klasyfikacja wielowyjściowa	115
Ćwiczenia	117

4. Uczenie modeli	119
Regresja liniowa	120
Równanie normalne	121
Złożoność obliczeniowa	123
Gradient prosty	124
Wsadowy gradient prosty	127
Stochastyczny spadek wzdłuż gradientu	130
Schodzenie po gradiencie z minigrupami	132
Regresja wielomianowa	133
Krzywe uczenia	135
Regularyzowane modele liniowe	139
Regresja grzbietowa	139
Regresja metodą LASSO	141
Metoda elastycznej siatki	143
Wczesne zatrzymywanie	144
Regresja logistyczna	145
Szacowanie prawdopodobieństwa	146
Funkcje ucząca i kosztu	146
Granice decyzyjne	148
Regresja softmax	150
Ćwiczenia	153
5. Maszyny wektorów nośnych	155
Liniowa klasyfikacja SVM	155
Klasyfikacja miękkiego marginesu	156
Nieliniowa klasyfikacja SVM	158
Jądro wielomianowe	159
Dodawanie cech podobieństwa	160
Gaussowskie jądro RBF	161
Złożoność obliczeniowa	163
Regresja SVM	163
Mechanizm działania	165
Funkcja decyzyjna i prognozy	165
Cel uczenia	166
Programowanie kwadratowe	168
Problem dualny	169
Kernelizowane maszyny SVM	170
Przyrostowe maszyny SVM	172
Ćwiczenia	173

6. Drzewa decyzyjne	175
Uczenie i wizualizowanie drzewa decyzyjnego	175
Wyliczanie prognoz	176
Szacowanie prawdopodobieństw przynależności do klas	178
Algorytm uczący CART	179
Złożoność obliczeniowa	180
Wskaźnik Giniego czy entropia?	180
Hiperparametry regularyzacyjne	181
Regresja	182
Niestabilność	184
Ćwiczenia	185
7. Uczenie zespołowe i losowe lasy	187
Klasyfikatory głosujące	187
Agregacja i wklejanie	190
Agregacja i wklejanie w module Scikit-Learn	191
Ocena OOB	192
Rejony losowe i podprzestrzenie losowe	194
Losowe lasy	194
Zespół Extra-Trees	195
Istotność cech	195
Wzmacnianie	197
AdaBoost	197
Wzmacnianie gradientowe	200
Kontaminacja	204
Ćwiczenia	207
8. Redukcja wymiarowości	209
Kłątwa wymiarowości	210
Główne strategie redukcji wymiarowości	211
Rzutowanie	211
Uczenie różnorodnościowe	213
Analiza PCA	215
Zachowanie wariancji	215
Główne składowe	215
Rzutowanie na d wymiarów	217
Implementacja w module Scikit-Learn	217
Współczynnik wariancji wyjaśnionej	217
Wybór właściwej liczby wymiarów	218
Algorytm PCA w zastosowaniach kompresji	219
Przyrostowa analiza PCA	220
Losowa analiza PCA	220

Jądrowa analiza PCA	221
Wybór jądra i strojenie hiperparametrów	221
Algorytm LLE	224
Inne techniki redukcji wymiarowości	225
Ćwiczenia	226

Część II. Sieci neuronowe i uczenie głębokie

229

9. Instalacja i używanie modułu TensorFlow	231
Instalacja	234
Tworzenie pierwszego grafu i uruchamianie go w sesji	234
Zarządzanie grafami	236
Cykl życia wartości w węźle	236
Regresja liniowa przy użyciu modułu TensorFlow	237
Implementacja metody gradientu prostego	238
Ręczne obliczanie gradientów	238
Automatyczne różniczkowanie	239
Korzystanie z optymalizatora	240
Dostarczanie danych algorytmowi uczącemu	241
Zapisywanie i wczytywanie modeli	242
Wizualizowanie grafu i krzywych uczenia za pomocą modułu TensorBoard	243
Zakresy nazw	246
Modułowość	247
Udostępnianie zmiennych	250
Ćwiczenia	253
10. Wprowadzenie do sztucznych sieci neuronowych	255
Od biologicznych do sztucznych neuronów	255
Neurony biologiczne	256
Operacje logiczne przy użyciu neuronów	257
Perceptron	259
Perceptron wielowarstwowy i propagacja wsteczna	262
Uczenie sieci MLP za pomocą zaawansowanego interfejsu API modułu TensorFlow	265
Uczenie głębokiej sieci neuronowej za pomocą standardowego interfejsu TensorFlow	266
Faza konstrukcyjna	266
Faza wykonawcza	270
Korzystanie z sieci neuronowej	271
Strojenie hiperparametrów sieci neuronowej	271
Liczba ukrytych warstw	271
Liczba neuronów tworzących warstwę ukrytą	272
Funkcje aktywacji	273
Ćwiczenia	273

11. Uczenie głębokich sieci neuronowych	275
Problemy zanikających/eksplodujących gradientów	275
Inicjacje wag Xavier'a i He	276
Nienasycające funkcje aktywacji	278
Normalizacja wsadowa	281
Obcinanie gradientu	285
Wielokrotne stosowanie gotowych warstw	286
Wielokrotne stosowanie modelu TensorFlow	287
Wykorzystywanie modeli utworzonych w innych środowiskach	288
Zamrażanie niższych warstw	289
Zapamiętywanie warstw ukrytych	290
Modyfikowanie, usuwanie lub zastępowanie górnych warstw	291
Repozytoria modeli	291
Nienadzorowane uczenie wstępne	291
Uczenie wstępne za pomocą dodatkowego zadania	293
Szybsze optymalizatory	293
Optymalizacja momentum	294
Przyśpieszony spadek wzdłuż gradientu (algorytm Nesterova)	295
AdaGrad	296
RMSProp	298
Optymalizacja Adam	298
Harmonogramowanie współczynnika uczenia	300
Regularyzacja jako sposób unikania przetrenowania	302
Wczesne zatrzymywanie	302
Regularyzacja ℓ_1 i ℓ_2	303
Porzucanie	304
Regularyzacja typu max-norm	306
Dogenerowanie danych	308
Praktyczne wskazówki	309
Ćwiczenia	310
12. Rozdzielanie operacji TensorFlow pomiędzy urządzenia i serwery	313
Wiele urządzeń na jednym komputerze	314
Instalacja	314
Zarządzanie pamięcią operacyjną karty graficznej	316
Umieszczanie operacji na urządzeniach	318
Przetwarzanie równoległe	321
Zależności sterujące	322
Wiele urządzeń na wielu serwerach	323
Otwieranie sesji	324
Usługi nadrzędna i robocza	325

Przypinanie operacji w wielu zadaniach	325
Rozdzielanie zmiennych pomiędzy wiele serwerów parametrów	326
Udostępnianie stanu rozproszonych sesji za pomocą kontenerów zasobów	327
Komunikacja asynchroniczna za pomocą kolejek	328
Wczytywanie danych bezpośrednio z grafu	333
Przetwarzanie równoległe sieci neuronowych w klastrze TensorFlow	339
Jedna sieć neuronowa na każde urządzenie	339
Replikacja wewnątrzgrafowa i międzygrafowa	341
Zrównoleglanie modelu	343
Zrównoleglanie danych	345
Ćwiczenia	349
13. Splotowe sieci neuronowe	351
Architektura kory wzrokowej	352
Warstwa splotowa	353
Filtry	354
Stosy map cech	356
Implementacja w module TensorFlow	357
Zużycie pamięci operacyjnej	359
Warstwa łącząca	360
Architektury splotowych sieci neuronowych	362
LeNet-5	362
AlexNet	364
GoogLeNet	366
ResNet	369
Ćwiczenia	373
14. Rekurencyjne sieci neuronowe	377
Neurony rekurencyjne	378
Komórki pamięci	380
Sekwencje wejść i wyjść	380
Podstawowe sieci RSN w module TensorFlow	381
Statyczne rozwijanie w czasie	382
Dynamiczne rozwijanie w czasie	384
Obsługa sekwencji wejściowych o zmiennej długości	385
Obsługa sekwencji wyjściowych o zmiennej długości	386
Uczenie rekurencyjnych sieci neuronowych	386
Uczenie klasyfikatora sekwencji	387
Uczenie w celu przewidywania szeregów czasowych	389
Twórcza sieć rekurencyjna	392

Głębokie sieci rekurencyjne	393
Rozmieszczanie głębokiej sieci rekurencyjnej pomiędzy wiele kart graficznych	394
Wprowadzanie metody porzucania	395
Problem uczenia w sieciach wielotaktowych	396
Komórka LSTM	397
Połączenia przezierne	400
Komórka GRU	400
Przetwarzanie języka naturalnego	401
Reprezentacje wektorowe słów	402
Sieć typu koder-dekoder służąca do tłumaczenia maszynowego	403
Ćwiczenia	406
15. Autokodery	407
Efektywne reprezentacje danych	407
Analiza PCA za pomocą niedopełnionego autokodera liniowego	409
Autokodery stosowe	410
Implementacja w module TensorFlow	411
Wiązanie wag	412
Uczenie autokoderów pojedynczo	413
Wizualizacja rekonstrukcji	415
Wizualizowanie cech	416
Nienadzorowane uczenie wstępne za pomocą autokoderów stosowych	417
Autokodery odszumiające	419
Implementacja w module TensorFlow	420
Autokodery rzadkie	420
Implementacja w module TensorFlow	422
Autokodery wariacyjne	423
Generowanie cyfr	426
Inne autokodery	427
Ćwiczenia	427
16. Uczenie przez wzmacnianie	431
Uczenie się optymalizowania nagród	432
Wyszukiwanie polityki	433
Wprowadzenie do narzędzia OpenAI gym	435
Sieci neuronowe jako polityki	438
Ocenianie czynności — problem przypisania zasługi	440
Gradienty polityk	441
Procesy decyzyjne Markowa	446

Uczenie metodą różnic czasowych i algorytm Q-uczenia	449
Polityki poszukiwania	451
Przybliżający algorytm Q-uczenia	451
Nauka gry w Ms. Pac-Man za pomocą głębokiego Q-uczenia	452
Ćwiczenia	459
Dziękuję!	460

Dodatki	461
A Rozwiązania ćwiczeń	463
B Lista kontrolna projektu uczenia maszynowego	487
C Problem dualny w maszynach wektorów nośnych	493
D Różniczkowanie automatyczne	497
Skorowidz	513

Krajobraz uczenia maszynowego

Większość osób, słysząc o uczeniu maszynowym, wyobraża sobie robota: jedni myślą o posłusznym kamerdynerze, inni o zabójczym Terminatorze. Zjawisko to jednak nie stanowi futurystycznej fantazji, lecz jest już elementem współczesnego świata. W rzeczywistości istnieje już od dziesiątek lat w pewnych wyspecjalizowanych zastosowaniach, takich jak techniki **optycznego rozpoznawania znaków** (ang. *Optical Character Recognition* — OCR). Jednak techniki uczenia maszynowego trafiły pod strzechy, poprawiając komfort życia milionów osób, dopiero w latach dziewięćdziesiątych — w postaci **filtrów spamu**. Niekoniecznie przypominają one Skynet¹, ale pod względem technicznym są klasyfikowane jako aplikacje wykorzystujące uczenie maszynowe (istotnie, filtry te uczą się tak skutecznie, że bardzo rzadko musisz własnoręcznie oznaczać wiadomości jako spam). W następnych latach nastąpił wysyp różnych innych zastosowań uczenia maszynowego, stanowiących trzon wielu aplikacji i usług, z których korzystamy na co dzień, począwszy od polecanych produktów aż do wyszukiwania głosowego.

Gdzie się zaczyna i kończy uczenie maszynowe? Co to właściwie oznacza, że dana maszyna **uczy się** określonego zagadnienia? Czy jeśli pobiorę na dysk artykuł z Wikipedii, to mój komputer „nauczył się” jego treści? Czy stał się nagle mądrzejszy? W tym rozdziale wyjaśnię, co uznajemy za uczenie maszynowe oraz dlaczego ta dziedzina może Cię zainteresować.

Następnie, zanim wyruszymy na wyprawę w głąb kontynentu Uczenia Maszynowego, zerkniemy na mapę i przyjrzymy się jego najważniejszym regionom oraz punktom orientacyjnym: porównamy uczenie nadzorowane z nienadzorowanym, przyrostowe ze wsadowym, a także metody uczenia z przykładów z metodami uczenia z modelu. Następnie przeanalizujemy cykl tworzenia typowego projektu uczenia maszynowego, zastanowimy się nad głównymi wyzwaniami, jakim należy stawić czoła, a także przeanalizujemy sposoby oceny i strojenia szybkości modelu.

W niniejszym rozdziale wprowadzam wiele podstawowych pojęć (i terminów), które musi znać każdy szanujący się analityk danych. Informacje tu zawarte są bardzo ogólne (jest to jedyny rozdział pozbawiony listingów kodu) i raczej nieskomplikowane, zanim jednak przejdiesz do następnego rozdziału, upewnij się, że doskonale je rozumiesz. Zatem kawa w dłoń i do dzieła!

¹ Samoświadoma sztuczna inteligencja dążąca do zniszczenia ludzkości, istniejąca w serii filmów *Terminator* — *przyp. tłum.*



Jeśli znasz już podstawy uczenia maszynowego, możesz przejść od razu do rozdziału 2. Jeżeli jednak nie masz pewności, spróbuj najpierw odpowiedzieć na wszystkie pytania umieszczone na końcu rozdziału.

Czym jest uczenie maszynowe?

Uczeniem maszynowym nazywamy dziedzinę nauki (i sztukę) programowania komputerów w sposób umożliwiający im **uczenie się z danych**.

Poniżej przedstawiam nieco ogólniejszą definicję:

[Uczenie maszynowe to] dziedzina nauki dająca komputerom możliwość uczenia się bez konieczności ich jawnego programowania.

— Arthur Samuel, 1959

A tu bardziej techniczna:

Mówimy, że program komputerowy uczy się na podstawie doświadczenia E w odniesieniu do jakiegoś zadania T i pewnej miary wydajności P, jeśli jego wydajność (mierzona przez P) wobec zadania T wzrasta wraz z nabywaniem doświadczenia E.

— Tom Mitchell, 1997

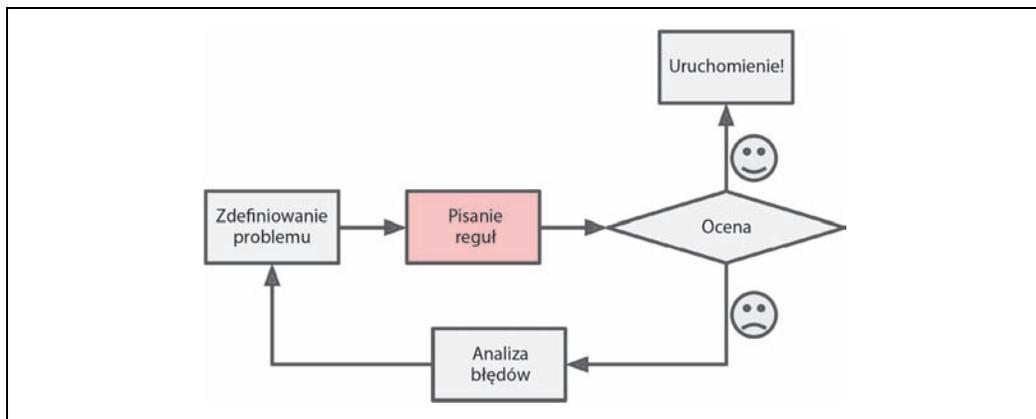
Na przykład filtr spamu jest programem wykorzystującym algorytmy uczenia maszynowego do rozpoznawania spamu na podstawie przykładowych wiadomości e-mail (tj. spamu oznaczonego przez użytkowników oraz zwykłych wiadomości niebędących spamem, znanych jako „ham”²). Przykładowe dane używane do trenowania systemu noszą nazwę **zbioru/zestawu uczącego** (ang. *training set*). Każdy taki element uczący jest nazywany **przykładem uczącym (próbką uczącą)**. Zgodnie z powyższą definicją naszym zadaniem T jest tu oznaczanie spamu, doświadczeniem E — **dane uczące**; pozostaje nam wyznaczyć miarę wydajności P. Może być nią, na przykład, stosunek prawidłowo oznaczonych wiadomości do przykładów nieprawidłowo sklasyfikowanych. Ta konkretna miara wydajności jest nazywana **dokładnością** (ang. *accuracy*) i znajduje często zastosowanie w zadaniach klasyfikujących.

Nawet jeśli pobierzesz całą Wikipedię na dysk, Twój komputer będzie przechowywał mnóstwo informacji, ale nie wykorzysta ich w żaden sposób. Dlatego nie uznajemy tego za uczenie maszynowe.

Dlaczego warto korzystać z uczenia maszynowego?

Żałujemy, że chcemy napisać filtr spamu za pomocą tradycyjnych technik programistycznych (rysunek 1.1):

² Angielska nazwa niechcianych wiadomości pocztowych — „spam” (mielonka) — została zapożyczona ze słynnego skeczu grupy Monty Python pt. „Spam”. Nazwa standardowych wiadomości („ham”) oznacza w dosłownym tłumaczeniu szynkę — *przyp. tłum.*



Rysunek 1.1. Tradycyjne podejście

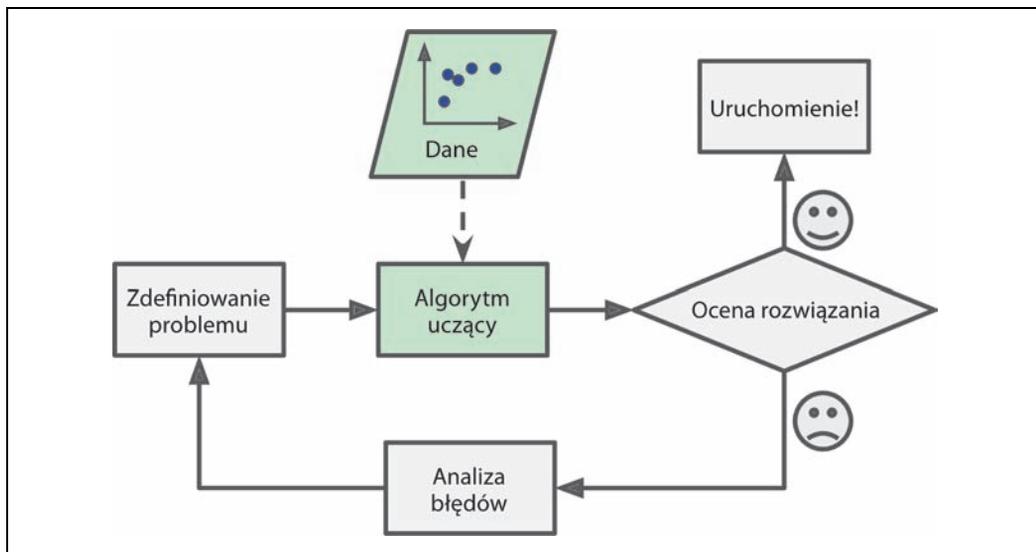
1. Najpierw musielibyśmy zastanowić się, jak wygląda klasyczny spam. Pewnie zauważysz, że niektóre wyrazy lub zwroty (np. „karta kredytowa”, „darmowe”, „zdumiewające”, „bez limitów”) bardzo często występują w temacie wiadomości. Być może odkryjesz również kilka innych szablonów w nazwie nadawcy, ciele wiadomości itd.
2. Następnym etapem byłoby napisanie algorytmu wykrywającego każdy z zaobserwowanych szablonów; program oznaczałby wiadomości jako spam, jeśli wykrywałby w nich kilka spośród zdefiniowanych wzorców.
3. Teraz należałoby przetestować ten program i usprawniać go poprzez ciągłe powtarzanie kroków 1. i 2.

Problem ten nie jest trywialny i omawiany program będzie się rozrastał o coraz większą liczbę skomplikowanych reguł — jego własnoręczne utrzymanie stałoby się w końcu bardzo uciążliwe.

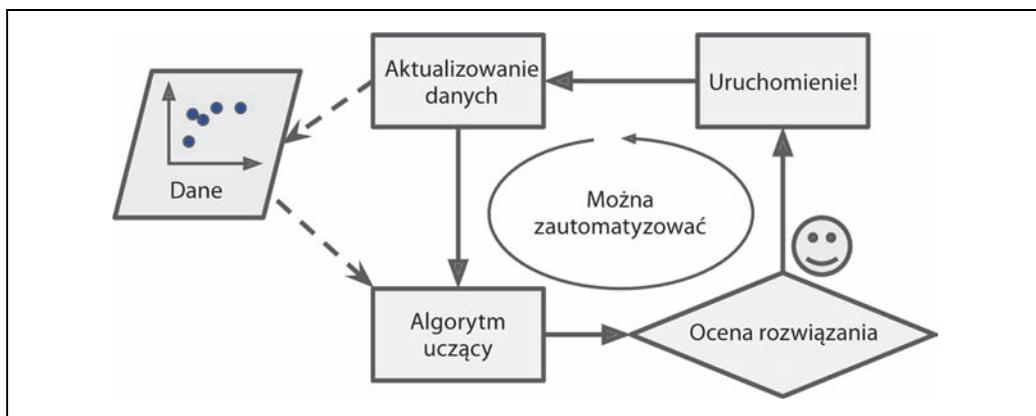
Z drugiej strony filtr spamu bazujący na algorytmach uczenia maszynowego automatycznie rozpoznaje, które słowa i wyrażenia stanowią dobre czynniki prognostyczne, poprzez wykrywanie podejrzanie często powtarzających się wzorców wyrazów w przykładowym spamie (w porównaniu do przykładów standardowych wiadomości); patrz rysunek 1.2. Kod takiej aplikacji jest znacznie krótszy, łatwiejszy do utrzymywania i prawdopodobnie dokładniejszy.

Poza tym jeśli spamerzy zauważą, że blokowane są wszystkie wiadomości zawierające zwrot „bez limitu”, mogą go zmodyfikować na „bez ograniczeń”. Filtr spamu stworzony za pomocą tradycyjnych technik programistycznych musiałby wtedy zostać zaktualizowany tak, aby oznaczał wiadomości zawierające nowy zwrot. Gdyby spamerzy próbowali cały czas w jakiś sposób oszukiwać Twój filtr spamu, musiałabyś/musiałbyś bez przerwy dopisywać nowe reguły.

Natomiast filtr spamu stworzony za pomocą technik uczenia maszynowego automatycznie „zauważy” nagły wzrost częstotliwości pojawiania się zwrotu „bez ograniczeń” w spamie oznaczonym przez użytkowników i zacznie go blokować bez Twojego udziału (rysunek 1.3).



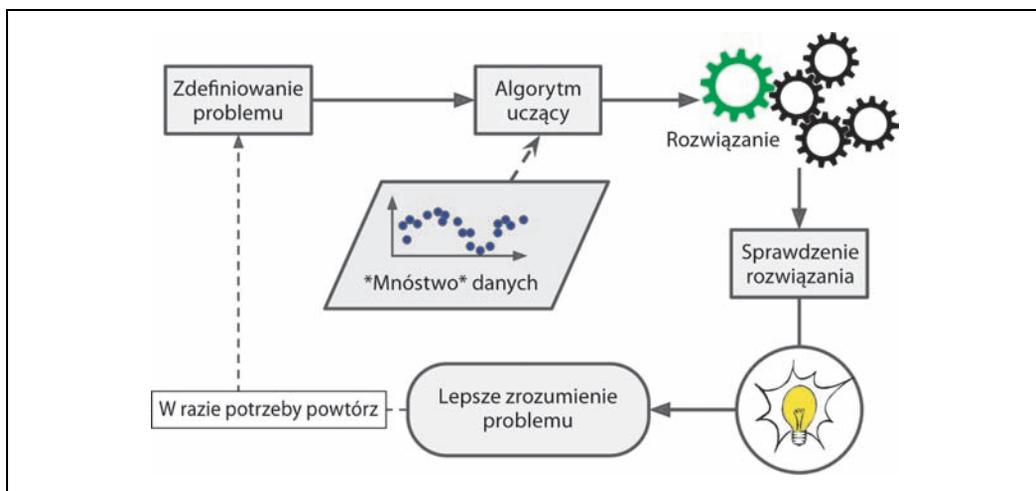
Rysunek 1.2. Podejście wykorzystujące uczenie maszynowe



Rysunek 1.3. Automatyczne dostosowywanie się do zmian

Kolejnym obszarem, w którym wyróżnia się mechanizm uczenia maszynowego, są problemy, które są zbyt złożone dla tradycyjnych metod lub dla których nie istnieją gotowe algorytmy. Przyjrzyjmy się, na przykład, rozpoznawaniu mowy: założmy, że chcesz stworzyć prostą aplikację zdolną do rozróżniania wyrazów „raz” i „trzy”. Zwróć uwagę, że słowo „trzy” zaczyna się od dźwięku o wysokim tonie („T”), dlatego moglibyśmy stworzyć algorytm wykrywający wysokość dźwięku i użyć go do rozróżniania obydwu wyrazów. Oczywiście, technika ta zupełnie nie nadaje się do rozpoznawania tysięcy słów wypowiedzianych w hałaśliwych miejscach przez miliony różnych ludzi, porozumiewających się w dziesiątkach języków. Najlepszym rozwiązaniem (przynajmniej przy obecnym stanie naszej wiedzy) jest napisanie samouczącego się algorytmu na podstawie przykładowych nagrań danego słowa.

Uczenie maszynowe pomaga również w trenowaniu ludzi (rysunek 1.4): możemy przeglądać algorytmy, aby sprawdzić, czego się nauczyły (choć nie zawsze jest to łatwo sprawdzić). Przykładowo, po wytrenowaniu filtra spamu do rozpoznawania odpowiedniej ilości spamu, można z łatwością sprawdzić listę wyuczonych słów i kombinacji wyrazów uznawanych za najlepsze czynniki prognozy. Czasem w ten sposób możemy wykryć niespodziewane korelacje lub nowe trendy, dzięki czemu jesteśmy w stanie lepiej pojąć dany problem.



Rysunek 1.4. Uczenie maszynowe pomaga w trenowaniu ludzi

Wykorzystanie technik uczenia maszynowego do analizowania olbrzymich ilości danych może pomóc w wykrywaniu nieoczywistych wzorców. Proces ten nazywamy **wydobywaniem danych** (ang. *data mining*).

Podsumowując, uczenie maszynowe nadaje się znakomicie do:

- problemów, których rozwiązanie wymaga mnóstwo ręcznego dostrajania algorytmu lub korzystania z długich list reguł; często jeden algorytm uczenia maszynowego upraszcza aplikację i poprawia jej szybkość,
- złożonych problemów, których nie można rozwiązać tradycyjnymi metodami; najlepsze algorytmy uczenia maszynowego są w stanie znaleźć rozwiązanie,
- zmiennych środowisk; mechanizm uczenia maszynowego potrafi dostosować się do nowych danych,
- pomagania człowiekowi w analizowaniu skomplikowanych zagadnień i olbrzymich ilości danych.

Rodzaje systemów uczenia maszynowego

Istnieje tak wiele typów uczenia maszynowego, że warto podzielić je na ogólne kategorie na podstawie takich czynników, jak:

- nadzór człowieka w procesie trenowania (uczenie nadzorowane, nienadzorowane, półnadzorowane i uczenie przez wzmacnianie),

- możliwość uczenia się w czasie rzeczywistym (uczenie przyrostowe i wsadowe),
- sposób pracy: proste porównywanie nowych punktów danych ze znanymi punktami lub, podobnie do naukowców, wykrywanie wzorców w danych uczących i tworzenie modelu predykcyjnego (uczenie z przykładów i uczenie z modelu).

Kryteria te nie wykluczają się wzajemnie; możesz łączyć je w dowolny sposób. Na przykład nowoczesny filtr spamu może uczyć się na bieżąco przy użyciu modelu głębokiej sieci neuronowej, korzystając z przykładowych wiadomości e-mail; w ten sposób zaprzęgamy do pracy przyrostowy, bazujący na modelu i nadzorowany system uczenia maszynowego.

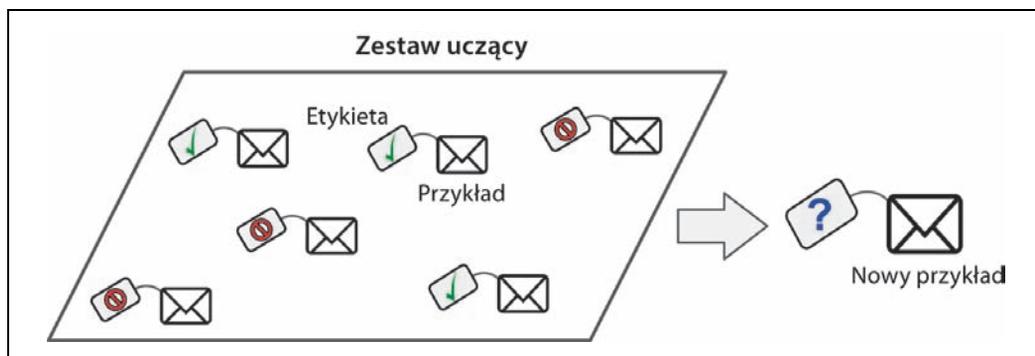
Przyjrzyjmy się nieco uważniej każdemu z wymienionych kryteriów.

Uczenie nadzorowane/nienadzorowane

Systemy uczenia maszynowego możemy podzielić na podstawie stopnia i rodzaju nadzorowania procesu uczenia. Pod tym względem uczenie maszynowe dzielimy na cztery zasadnicze rodzaje: uczenie nadzorowane, nienadzorowane, półnadzorowane i uczenie przez wzmacnianie).

Uczenie nadzorowane

W **uczeniu nadzorowanym** (ang. *supervised learning*) dane uczące przekazywane algorytmowi zawierają dołączone rozwiązania problemu, tzw. **etykiety** (ang. *labels*); patrz rysunek 1.5.



Rysunek 1.5. Zbiór danych uczących zaopatrzonych w etykiety, stosowany w uczeniu nadzorowanym (np. do klasyfikowania spamu)

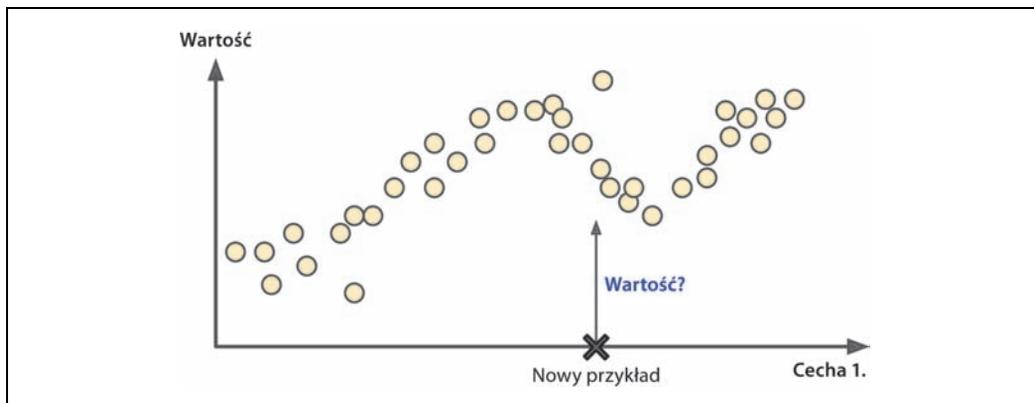
Klasycznym zadaniem uczenia nadzorowanego jest **klasyfikacja** (ang. *classification*). Filtr spamu stanowi tu dobry przykład: jest on trenowany za pomocą dużej liczby przykładowych wiadomości e-mail należących do danej **klasy** (spamu lub hamu), dzięki którym musi być w stanie klasyfikować nowe wiadomości.

Innym typowym zadaniem uczenia nadzorowanego jest przewidywanie **docelowej** wartości numerycznej, takiej jak cena samochodu, przy użyciu określonego zbioru cech (przebieg, wiek, marka itd.), zwanych **czynnikiemami prognostycznymi/predykcyjnymi** (ang. *predictors*). Ten typ zadania nosi nazwę

regresji (rysunek 1.6)³. Aby wyuczyć dany system, należy mu podać wiele przykładów samochodów, w tym zarówno etykiety (np. ceny), jak i czynniki prognostyczne.



W terminologii uczenia maszynowego **atrybutem** (ang. *attribute*) nazywamy typ danych (np. „Przebieg”), natomiast **cecha** (ang. *feature*) w zależności od kontekstu ma kilka różnych znaczeń, zazwyczaj jednak mówiąc o cesze, mamy na myśli atrybut wraz z jego wartością (np. „Przebieg = 15 000”). Wiele osób korzysta jednak naprzemiennie z wyrazów **atrybut** i **cecha**.



Rysunek 1.6. Regresja

Zwróć uwagę, że niektóre algorytmy regresyjne mogą być również używane do klasyfikowania danych i odwrotnie. Na przykład algorytm **regresji logistycznej** jest powszechnie stosowany w zadaniach klasyfikacyjnych, ponieważ może podawać wynikową wartość odpowiadającą prawdopodobieństwu przynależności do danej klasy (np. 20% szans na to, że dana wiadomość jest spamem).

Poniżej wymieniam niektóre z najważniejszych algorytmów nadzorowanego uczenia maszynowego (opisywanych w tej książce):

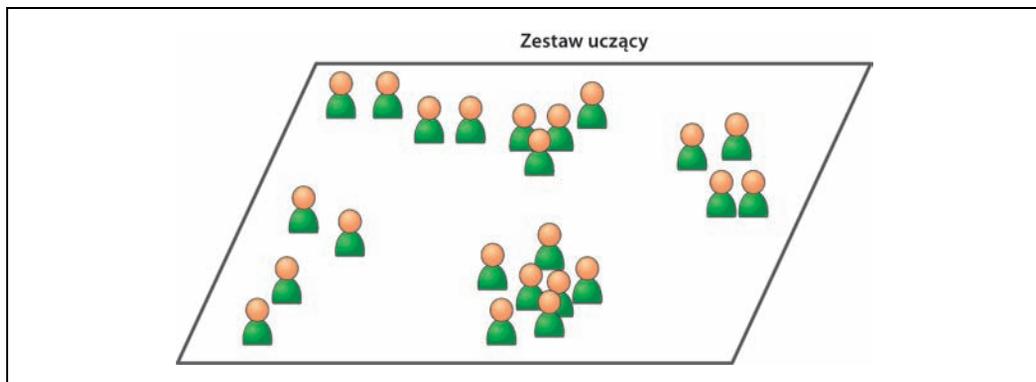
- metoda k-najbliższych sąsiadów,
- regresja liniowa,
- regresja logistyczna,
- maszyny wektorów nośnych,
- drzewa decyzyjne i losowe lasy,
- sieci neuronowe⁴.

³ Ciekawostka: ta dziwnie brzmiąca nazwa wywodzi się z teorii statystyki i została zaproponowana przez Francisca Galtona, gdy analizował regułę, zgodnie z którą dzieci wysokich rodziców często bywają od nich niższe. Zjawisko redukcji wzrostu w kolejnym pokoleniu badacz ten nazwał **regresją do średniej**. Nazwą tą zostały następnie określone użyte przez niego metody służące do analizowania korelacji pomiędzy zmiennymi.

⁴ Architektury niektórych sieci neuronowych mogą być nienadzorowane (np. autokodery czy ograniczone maszyny Boltzmanna). Istnieją także architektury półnadzorowane, takie jak głębokie sieci przekonań i nienadzorowane uczenie wstępne.

Uczenie nienadzorowane

Jak można się domyślać, w **uczeniu nienadzorowanym** (ang. *unsupervised learning*) dane uczące są nieoznakowane (rysunek 1.7). System próbuje uczyć się bez nauczyciela.

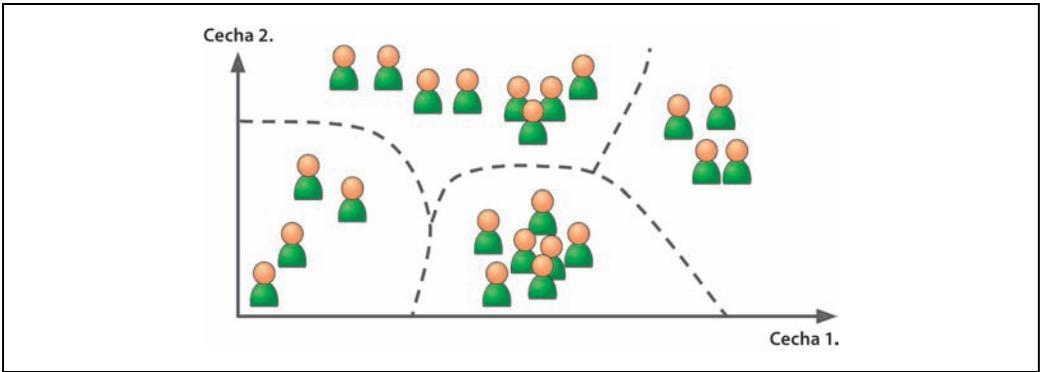


Rysunek 1.7. Zbiór nieoznakowanych danych uczących używanych w uczeniu nienadzorowanym

Poniżej wymieniamy kilka najważniejszych algorytmów uczenia nienadzorowanego (w rozdziale 8. zajmiemy się zagadnieniem redukcji wymiarowości):

- analiza skupień:
 - metoda k-średnich,
 - hierarchiczna analiza skupień (HCA),
 - algorytm oczekiwanie-maksymalizacja,
- wizualizacja i redukcja wymiarowości:
 - analiza głównych składowych (PCA),
 - jądrowa analiza głównych składowych,
 - lokalnie liniowe zanurzenie (ang. *locally linear embedding* — LLE),
 - stochastyczne zanurzenie sąsiadów przy użyciu rozkładu t (ang. *t-Distributed Stochastic Neighbor Embedding* — t-SNE),
- uczenie przy użyciu reguł asocjacyjnych:
 - algorytm Apriori,
 - algorytm Eclat.

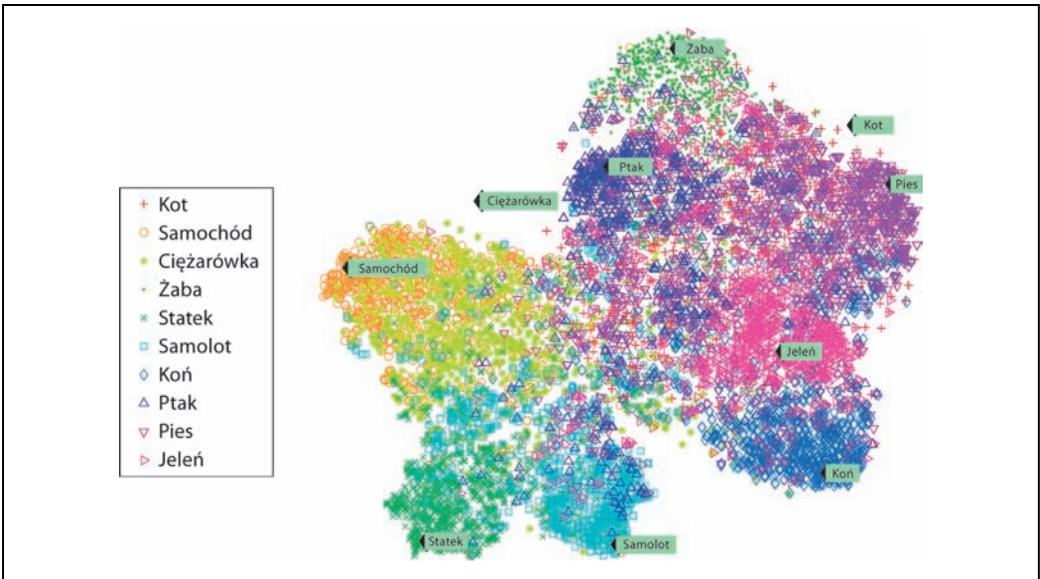
Załóżmy, na przykład, że masz do dyspozycji mnóstwo danych dotyczących osób odwiedzających Twój blog. Możesz chcieć skorzystać z **analizy skupień** (ang. *clustering*), aby spróbować określić grupy podobnych użytkowników (rysunek 1.8). W dowolnym momencie możesz sprawdzić, do której grupy zalicza się dana odwiedzająca osoba: powiązania pomiędzy poszczególnymi użytkownikami są określane bez Twojej pomocy. Algorytm ten może, przykładowo, zauważyć, że 40% odwiedzających osób to mężczyźni miłośnicy komiksów przeglądający Twoją stronę głównie wieczorami, podczas gdy grupa 20% innych użytkowników to młodociani czytelnicy fantastyki naukowej, którzy zaglądają na Twój blog wyłącznie w weekendy itd. Jeśli użyjesz algorytmu **hierarchicznej analizy skupień**



Rysunek 1.8. Analiza skupień

(ang. *hierarchical clustering*), może on dodatkowo rozdzielić grupy na mniejsze podjednostki. W ten sposób możesz łatwiej określić, jakie wpisy umieszczają dla poszczególnych grup.

Dobrym przykładem algorytmów uczenia nienadzorowanego są również algorytmy **wizualizujące**: wprowadzasz mnóstwo złożonych, nieoznakowanych danych, które są następnie wyświetlane w postaci punktów na dwu- lub trójwymiarowym wykresie (rysunek 1.9). Algorytmy te starają się w maksymalnym stopniu zachować pierwotną strukturę danych (np. próbując rozdzielać poszczególne skupienia w przestrzeni wejściowej, redukując zjawiska nakładania się poszczególnych klastrów na siebie), dzięki czemu możesz łatwiej przeanalizować te dane i być może odkryć jakieś nieprzewidziane wzorce.



Rysunek 1.9. Przykład wizualizacji t-SNE wskazującej semantyczny podział skupień⁵

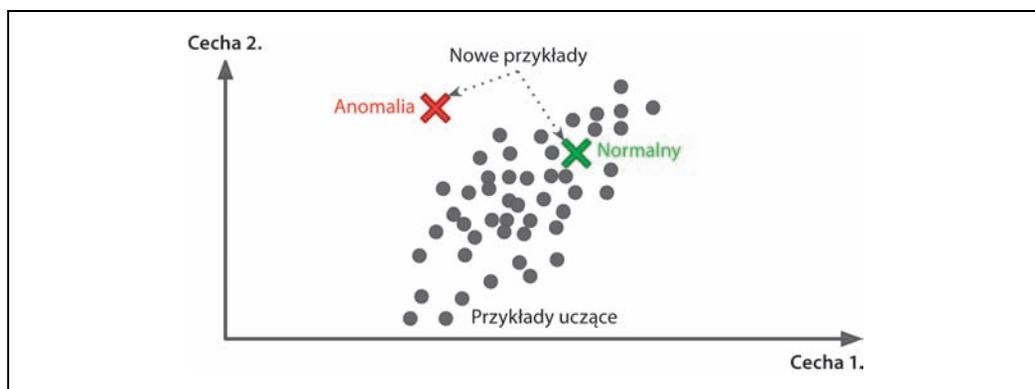
⁵ Zwróć uwagę, jak ładnie są rozdzielone zwierzęta od pojazdów, konie znajdują się blisko jeleni, ale daleko od ptaków itd. Rysunek wykorzystany za pozwoleniem Sochera, Ganjoo, Manninga i Ng (2013), *T-SNE visualization of the semantic word space*.

Podobnym zadaniem jest **redukcja wymiarowości** (ang. *dimensionality reduction*), której celem jest uproszczenie danych bez utraty nadmiernej ilości informacji. Możemy to osiągnąć między innymi poprzez scalenie kilku skorelowanych cech w jedną. Przykładowo, możemy skorelować przebieg samochodu z jego wiekiem, dzięki czemu algorytm będzie w stanie połączyć je w jedną cechę reprezentującą zużycie pojazdu. Proces ten jest nazywany **wydobyciem cech** (ang. *feature extraction*).



Zawsze warto spróbować zredukować wymiarowość danych uczących przed dostarczeniem ich do algorytmu uczenia maszynowego (np. algorytmu uczenia nadzorowanego). Algorytm będzie działał wtedy szybciej, dane będą zabierały mniej miejsca na dysku i w pamięci operacyjnej, a w niektórych przypadkach wzrośnie także wydajność uczenia maszynowego.

Kolejnym ważnym nienadzorowanym zadaniem jest **wykrywanie anomalii** (ang. *anomaly detection*), takich jak np. nietypowe transakcje wykorzystujące kartę kredytową, co pozwala zapobiegać nielegalnym operacjom. Służy ono także do wyłapywania usterek produkcyjnych czy też automatycznego usuwania elementów odstających w danych uczących przed ich przekazaniem algorytmowi uczenia maszynowego. Taki system jest trenowany za pomocą standardowych przykładów, a w przypadku zaprezentowania nowego przykładu jest w stanie stwierdzić, czy na przykład jest on standardowy, czy też stanowi anomalię (rysunek 1.10).

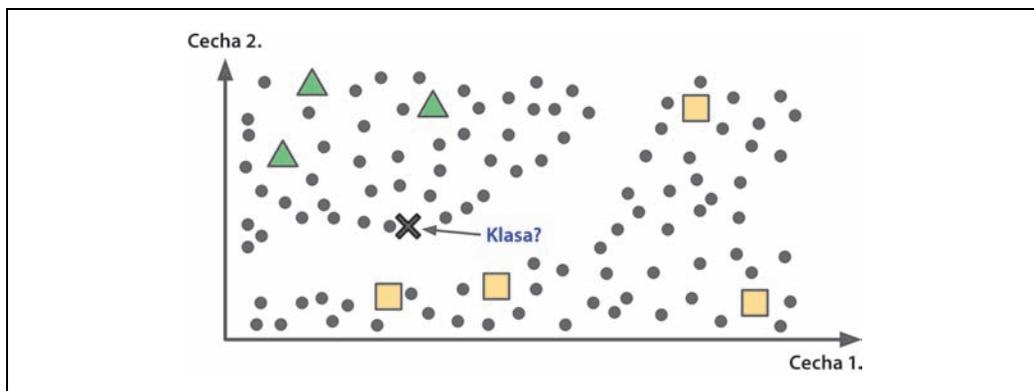


Rysunek 1.10. Wykrywanie anomalii

Jeszcze jednym powszechnie stosowanym zadaniem uczenia nienadzorowanego jest **uczenie przy użyciu reguł asocjacyjnych** (ang. *association rule learning*), którego celem jest analiza ogromnej ilości danych i wykrycie interesujących zależności pomiędzy atrybutami. Załóżmy, że jesteś właścicielką/właścicielem supermarketu. Przetworzenie danych dotyczących obrotów za pomocą algorytmu reguł asocjacyjnych pozwoliłoby nam odkryć, że osoby kupujące keczup i czipsy zazwyczaj zaopatrują się również w steki. Dzięki tej wiedzy możesz chcieć, na przykład, umieścić te produkty niedaleko siebie.

Uczenie półnadzorowane

Niektóre algorytmy są w stanie przetwarzać częściowo oznakowane dane uczące, najczęściej składające się z większości nieoznakowanych i tylko odrobiny oznakowanych przykładów. Taką sytuację nazywamy **uczeniem półnadzorowanym** (ang. *semisupervised learning*; rysunek 1.11).



Rysunek 1.11. Uczenie półnadzorowane

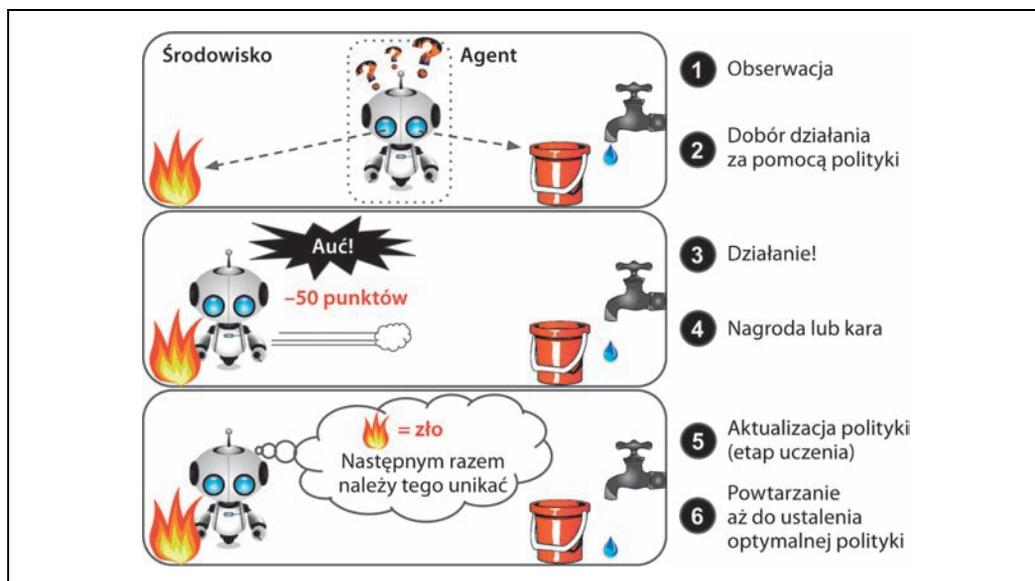
Dobrym przykładem zastosowania tego typu metod są niektóre usługi przechowywania obrazów, np. Zdjęcia Google. Po przesłaniu do takiej usługi wszystkich zdjęć rodzinnych rozpoznaje ona automatycznie, że na zdjęciach 1., 5. i 11. pojawia się osoba A, natomiast do fotografii 2., 5. i 7. pozuje osoba B. Na tym etapie mamy do czynienia z uczeniem nienadzorowanym (analizą skupień). Teraz wystarczy, że podasz imiona poszczególnych osób. Każda osoba potrzebuje tylko jednej etykiety⁶, co wystarczy do oznakowania wszystkich osób na każdym zdjęciu i znacznie ułatwia wyszukiwanie zdjęć.

Większość algorytmów półnadzorowanych stanowi kombinację algorytmów uczenia nadzorowanego i nienadzorowanego. Na przykład **głębokie sieci przekonań** (ang. *deep belief networks* — DBN) bazują na ułożonych warstwowo elementach nienadzorowanych, zwanych **ograniczonymi maszynami Boltzmann** (ang. *restricted Boltzmann machines* — RBM). Maszyny te są uczone sekwencyjnie w nienadzorowany sposób, a następnie cały system zostaje dostrojony przy użyciu technik uczenia nadzorowanego.

Uczenie przez wzmacnianie

Uczenie przez wzmacnianie (ang. *reinforcement learning*) to zupełnie inna bajka. System uczący, zwany w tym kontekście **agentem**, może obserwować środowisko, dobierać i wykonywać czynności, a także odbierać **nagrody** (lub **kary** będące ujemnymi nagrodami; patrz rysunek 1.12). Musi następnie samodzielnie nauczyć się najlepszej strategii (nazywanej **polityką**; ang. *policy*), aby uzyskiwać jak największą nagrodę. Polityka definiuje rodzaj działania, jakie agent powinien wybrać w danej sytuacji.

⁶ Tak wygląda wyidealizowana sytuacja. W praktyce usługa ta często tworzy kilka skupień dla jednej osoby, a czasami myli dwie podobne do siebie osoby, dlatego należy każdej osobie przydzielić kilka etykiet oraz własnoręcznie oczyścić niektóre skupienia.



Rysunek 1.12. Uczenie przez wzmacnianie

Na przykład metody uczenia przez wzmacnianie są używane w wielu modelach robotów do nauki chodzenia. Również program AlphaGo firmy DeepMind był trenowany przez wzmacnianie: zrobiło się o nim głośno w marcu 2016 roku, gdy pokonał dotychczasowego mistrza świata w grę *Go* — Lee Sedola. Aplikacja ta zdefiniowała politykę gwarantującą zwycięstwo, analizując miliony rozgrywek i rozgrywając mnóstwo partii przeciwko samej sobie. Warto zauważyć, że algorytmy uczące zostały wyłączone podczas partii z ludzkim przeciwnikiem; program AlphaGo korzystał jedynie z wyuczonoj polityki.

Uczenie wsadowe/przyrostowe

Innym kryterium stosowanym do klasyfikowania systemów uczenia maszynowego jest możliwość trenowania przyrostowego przy użyciu strumienia nadsyłanych danych.

Uczenie wsadowe

W mechanizmie **uczenia wsadowego** (ang. *batch learning*) system nie jest w stanie trenować przyrostowo: do jego nauki muszą wystarczyć wszystkie dostępne dane. Wymaga to zazwyczaj poświęcenia dużej ilości czasu i zasobów, dlatego najczęściej proces ten jest przeprowadzany w trybie offline. System najpierw jest uczony, a następnie zostaje wdrożony do cyklu produkcyjnego i już więcej nie jest trenowany; korzysta jedynie z dotychczas zdobytych informacji. Zjawisko to bywa nazywane **uczeniem offline** (ang. *offline learning*).

Jeśli chcesz, aby system uczenia wsadowego brał pod uwagę nowe dane (np. nowe rodzaje spamu), musisz od podstaw wytrenować nową wersję systemu przy użyciu wszystkich dostępnych danych (nowych i starych), następnie wyłączyć stary system i zastąpić go nowym.

Na szczęście cały proces trenowania, oceniania i uruchamiania systemu uczenia maszynowego można dość łatwo zautomatyzować (co widać na rysunku 1.3), dlatego nawet układ uczenia wsadowego jest w stanie dostosowywać się do zmian. Wystarczy zaktualizować dane i z odpowiednią częstotliwością trenować system od podstaw.

Rozwiązanie to jest proste i często skuteczne, ale trenowanie za pomocą pełnego zbioru danych uczących nieraz trwa wiele godzin, dlatego zazwyczaj trenujemy nowy system raz na dobę, a czasami nawet raz w tygodniu. Jeśli Twój system musi dostosowywać się do szybko zmieniających się danych (np. podczas prognozowania cen akcji giełdowych), będziesz potrzebować dynamiczniejszego mechanizmu.

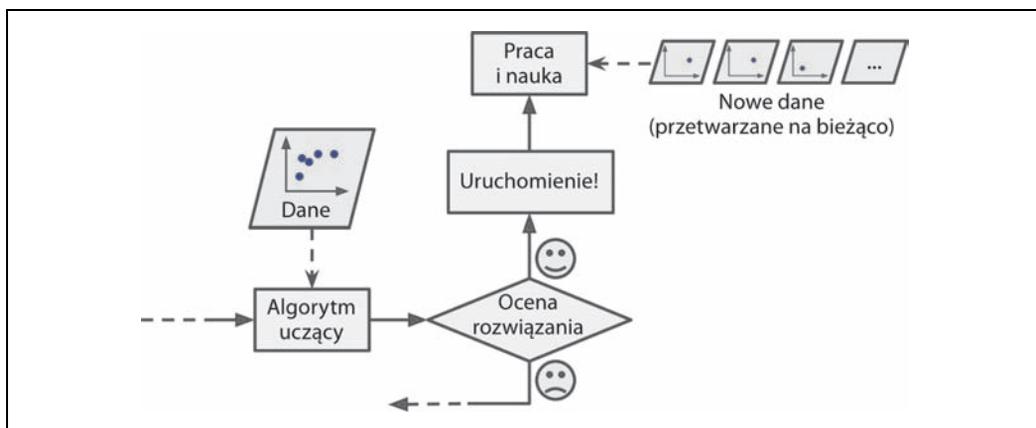
Ponadto uczenie się na pełnym zbiorze danych pochłania mnóstwo zasobów obliczeniowych (mocy procesora, miejsca i szybkości dysku, szybkości sieci itd.). Jeśli masz dużo danych i trenujesz system codziennie od podstaw, będzie Cię to kosztowało wiele pieniędzy. W przypadku danych o naprawdę pokaźnych rozmiarach korzystanie z algorytmu uczenia wsadowego może być wręcz niemożliwe.

Zwróćmy na koniec uwagę, że jeśli system ma być zdolny do autonomicznej nauki przy użyciu ograniczonych zasobów (np. aplikacja na smartfona albo marsjański łazik), to magazynowanie olbrzymich ilości danych i przeznaczanie zasobów na codzienne wielogodzinne trenowanie systemu mogą zwykle uniemożliwiać normalną pracę urządzenia.

Na szczęście istnieje lepsze rozwiązanie dla wymienionych powyżej przykładów: algorytmy zdolne do uczenia przyrostowego.

Uczenie przyrostowe

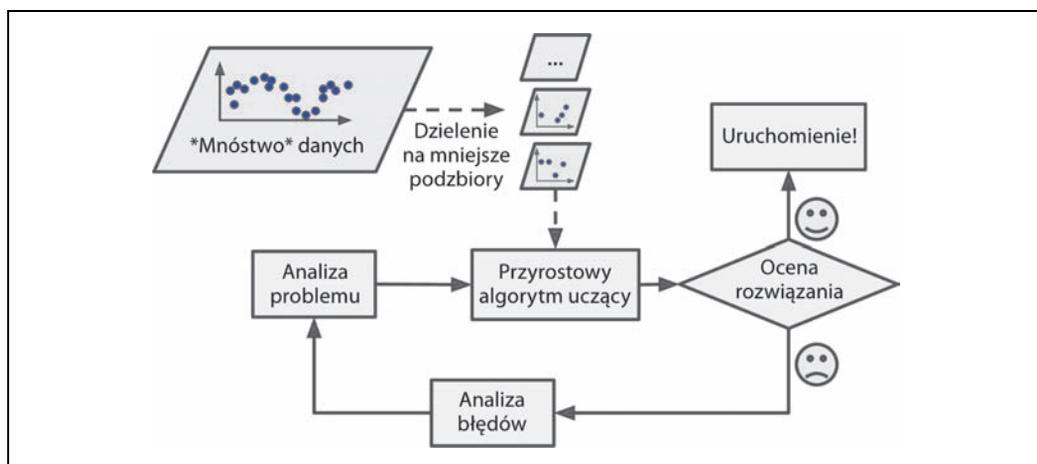
W procesie **uczenia przyrostowego** (ang. *online learning*) system jest trenowany na bieżąco poprzez sekwencyjne dostarczanie danych; mogą być one pojedyncze lub przyjmować postać tzw. **minipakietów/minigrup** (ang. *mini-batches*), czyli niewielkich zbiorów. Każdy krok uczący jest szybki i niezbyt kosztowny, dlatego system jest w stanie uczyć się na bieżąco przy użyciu nowych danych, gdy tylko się pojawiają (patrz rysunek 1.13).



Rysunek 1.13. Uczenie przyrostowe

Uczenie przyrostowe sprawdza się znakomicie w systemach odbierających ciągle strumień nowych danych (np. ceny akcji) oraz wymagających szybkiego lub autonomicznego dostosowywania się do nowych warunków. Metody te przydają się również przy ograniczonych zasobach obliczeniowych: dane użyte do nauki nie są już potrzebne i można się ich pozbyć (chyba że chcesz zachować możliwość cofnięcia się do wcześniejszego stanu i „odtworzyć ponownie” dane). W ten sposób możesz mocno zaoszczędzić na przestrzeni dyskowej.

Algorytmy uczenia przyrostowego są w stanie również trenować systemy za pomocą dużych zbiorów danych niemieszczących się w pamięci urządzenia (jest to tzw. **uczenie pozakorowe** — ang. *out-of-core learning*). Następuje wczytanie części danych i trenowanie systemu za ich pomocą, po czym cały proces zostanie powtórzony dla całego zbioru danych uczących (rysunek 1.14).



Rysunek 1.14. Przetwarzanie dużych ilości danych za pomocą algorytmu uczenia przyrostowego

Jednym z najważniejszych parametrów systemów uczenia przyrostowego jest szybkość, z jaką dostosowują się do zmieniających się danych: jest to tak zwany **współczynnik uczenia** (ang. *learning rate*). Wysoka wartość tego współczynnika oznacza, że system będzie szybko dostosowywał się do nowych danych, ale jednocześnie dość szybko zapominał o starych danych (raczej nie chcesz, aby filtr spamu oznaczał jedynie najnowsze rodzaje spamu). Z drugiej strony, system o niskim współczynniku uczenia będzie cechował się większą inercją; oznacza to, że będzie uczył się wolniej, ale jednocześnie będzie bardziej odporny na zaszumienie nowych danych oraz na sekwencje niereprezentatywnych punktów danych.

Z kolei dużym problemem uczenia przyrostowego jest stopniowy spadek wydajności systemu w przypadku dostarczenia mu nieprawidłowych danych. Klienci korzystający z danego systemu na pewno to zauważą. Przykładowo, nieprawidłowe dane mogą pochodzić z uszkodzonego czujnika w robocie lub od osoby zasypującej silnik przeglądarki danymi w celu podbicia swojego rankingu w wynikach wyszukiwania. Aby zmniejszyć to ryzyko, musisz uważnie obserwować swój system i świadomie wyłączyć proces uczenia (a także ewentualnie przywrócić stan do wcześniejszego stanu) po wykryciu spadku szybkości algorytmu. Możesz także śledzić przychodzące dane i reagować na wszelkie nieprawidłowości (np. za pomocą algorytmu wykrywania anomalii).

Uczenie z przykładów/z modelu

Możemy również podzielić systemy uczenia maszynowego pod względem **uogólniania**. Większość zadań uczenia maszynowego polega na sporządzaniu prognoz. Oznacza to, że na podstawie określonej liczby próbek uczących system musi być w stanie generalizować wyniki na niewidziane wcześniej przykłady. Uzyskanie dobrej szybkości algorytmu wobec danych uczących jest pożądane, ale niewystarczające; prawdziwy cel stanowi dobra wydajność wobec nowych przykładów.

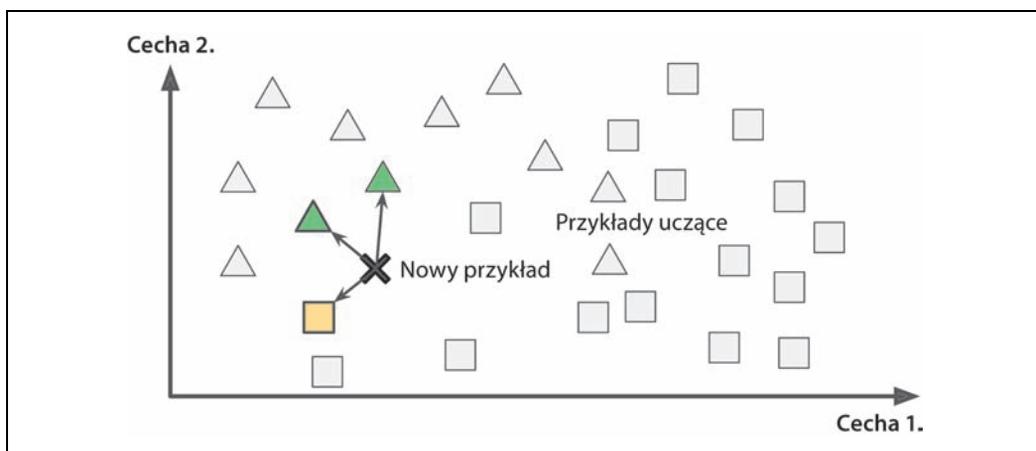
Znamy dwa główne mechanizmy uogólniania: uczenie z przykładów i uczenie z modelu.

Uczenie z przykładów

Prawdopodobnie najprostszą formą nauki jest „wykuwanie na blachę”. Gdybyśmy stworzyli w ten sposób filtr spamu, oznaczałby po prostu wszystkie wiadomości identyczne z oznakowanymi przez użytkowników — rozwiązanie nie najgorsze, ale zarazem nie najlepsze.

Zamiast oznaczać wiadomości identyczne z rozpoznanym wcześniej spamem, możemy zaprogramować filtr w taki sposób, aby znakował również wiadomości bardzo podobne do powszechnie rozpoznawanych wzorców. Potrzebna staje się jakaś **miara podobieństwa** dwóch wiadomości e-mail. Może ją stanowić (w wielkim uproszczeniu) porównanie liczby takich samych słów występujących w obydwu wiadomościach. Filtr mógłby oznaczać wiadomość jako spam, jeśli będzie miała wiele słów wspólnych ze znanym spamem.

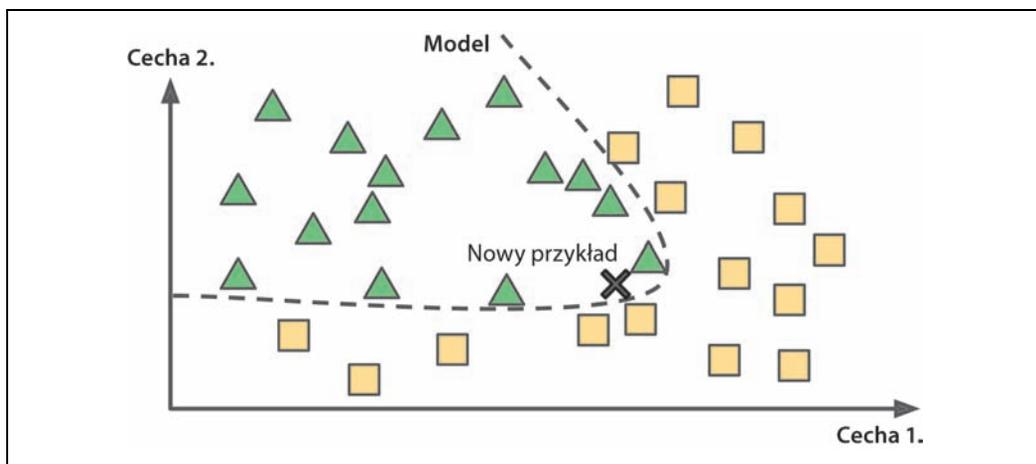
Mamy tu do czynienia z **uczeniem z przykładów** (ang. *instance-based learning*): system uczy się przykładów „na pamięć”, a następnie porównuje (generalizuje) je z nowymi przypadkami za pomocą miary podobieństwa (rysunek 1.15).



Rysunek 1.15. *Uczenie z przykładów*

Uczenie z modelu

Innym sposobem uogólniania wyników uzyskiwanych z danych uczących jest stworzenie modelu z tych przykładów i użycie go do **przewidywania** (**prognozowania**; ang. *prediction*). Jest to **uczenie z modelu** (ang. *model-based learning*; rysunek 1.16).



Rysunek 1.16. Uczenie z modelu

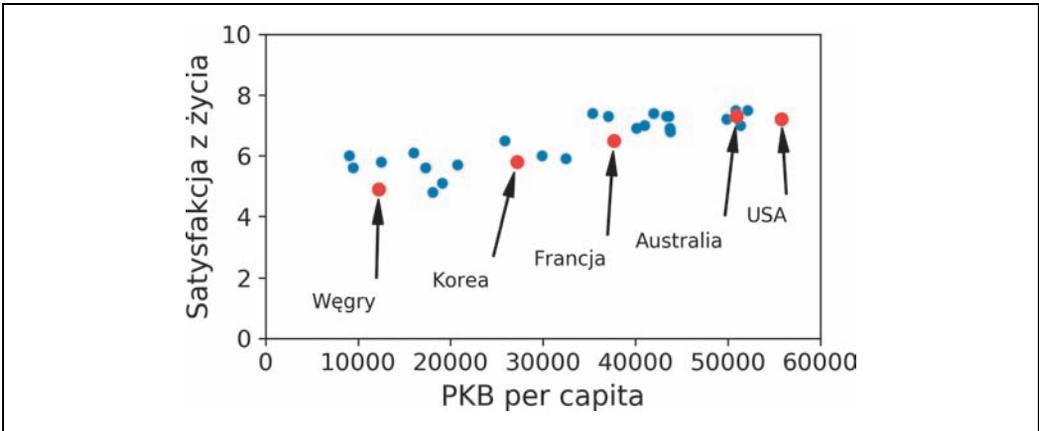
Załóżmy, na przykład, że chcemy sprawdzić, czy pieniądze dają szczęście; w tym celu moglibyśmy pobrać dane *Better Life Index* ze strony Organizacji Współpracy Gospodarczej i Rozwoju (OECD, <http://stats.oecd.org/index.aspx?DataSetCode=BLI>), a także statystyki PKB per capita z witryny Międzynarodowego Funduszu Walutowego (IMF, <http://goo.gl/j1MSKe>). Teraz wystarczyłoby połączyć obydwie tabele i posortować dane dotyczące PKB. Tabela 1.1 prezentuje fragment uzyskanych informacji.

Tabela 1.1. Czy pieniądze przynoszą szczęście?

Państwo	PKB per capita (w dolarach amerykańskich)	Satysfakcja z życia
Węgry	12 240	4,9
Korea	27 195	5,8
Francja	37 675	6,5
Australia	50 962	7,3
Stany Zjednoczone	55 805	7,2

Wygenerujmy teraz wykres danych dla kilku przykładowych krajów (rysunek 1.17).

Możemy dostrzec tu pewien trend! Pomimo tego, że dane są **zaszumione** (tzn. częściowo losowe), wygląda na to, że poziom satysfakcji z życia wzrasta w sposób mniej więcej liniowy wraz ze wzrostem produktu krajowego brutto na osobę. Postanawiasz zatem stworzyć model satysfakcji z życia jako funkcji liniowej wobec parametru PKB per capita. Etap ten nazywamy **doborem modelu**: wybraliśmy **model liniowy** satysfakcji z życia wykorzystujący tylko jeden atrybut — PKB per capita (równanie 1.1).

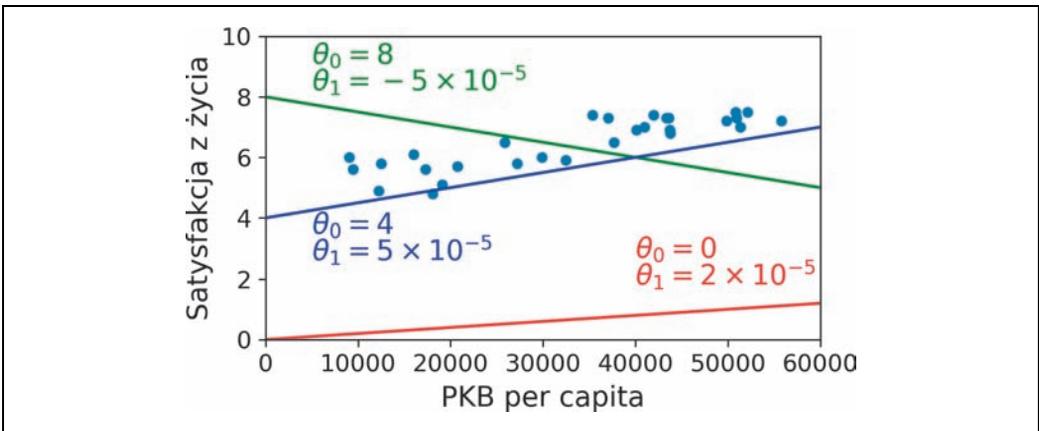


Rysunek 1.17. Czy widzisz tu trend?

Równanie 1.1. Prosty model liniowy

$$\text{satysfakcja_z_życia} = \theta_0 + \theta_1 \times \text{PKB_per_capita}$$

Model ten zawiera dwa **parametry**, θ_0 i θ_1 ⁷. Poprzez modyfikowanie tych parametrów możesz za ich pomocą uzyskać dowolną funkcję liniową, co zostało ukazane na rysunku 1.18.



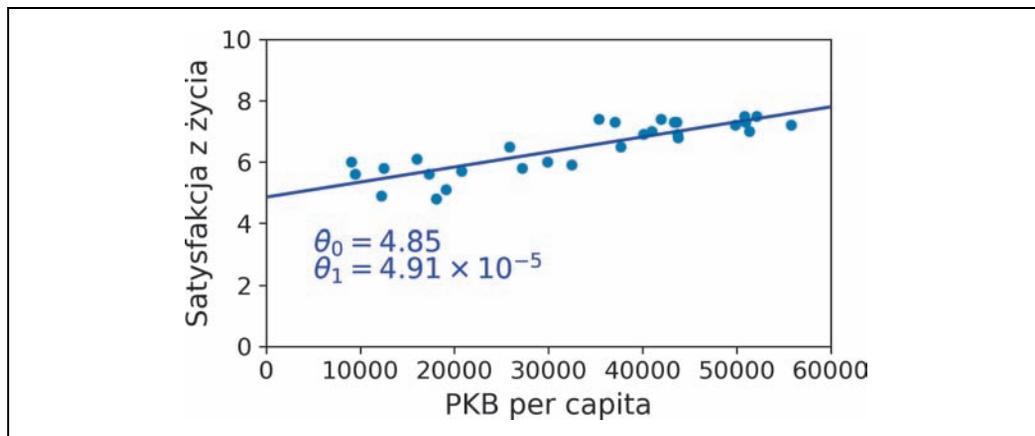
Rysunek 1.18. Kilka wybranych modeli liniowych

Zanim zaczniesz korzystać z modelu, musisz zdefiniować wartości jego parametrów θ_0 i θ_1 . Skąd mamy wiedzieć, które wartości nadają się najlepiej dla danego modelu? Aby móc odpowiedzieć na to pytanie, należy zdefiniować miarę szybkości. Możesz wyznaczyć **funkcję użyteczności** (zwaną także **funkcją dopasowania**), mówiącą nam, jak **dobry** jest dany model, lub **funkcję kosztu**, mającą przeciwne zastosowanie. W zagadnieniach wykorzystujących regresję liniową zazwyczaj jest stosowana funkcja kosztu mierząca odległości pomiędzy przewidywaniami modelu liniowego a przykładami uczącymi; naszym zadaniem jest zminimalizowanie tego dystansu.

⁷ Zgodnie z konwencją grecka litera θ (teta) często symbolizuje parametry modelu.

W tym właśnie miejscu przydaje się algorytm regresji liniowej: dostarczamy mu dane uczące, a on określa parametry najlepiej pasujące do danego modelu liniowego. Proces ten nosi nazwę uczenia (trenowania) modelu. W naszym przykładzie algorytm regresji liniowej wyznacza następujące optymalne wartości parametrów: $\theta_0 = 4,85$ i $\theta_1 = 4,91 \times 10^{-5}$.

Teraz nasz model jest maksymalnie dopasowany do danych uczących (jak na model liniowy), o czym możemy się przekonać na rysunku 1.19.



Rysunek 1.19. Model liniowy najlepiej dopasowany do danych uczących

Możemy w końcu uruchomić model, aby przeprowadził prognozy. Załóżmy, na przykład, że chcemy się dowiedzieć, jak szczęśliwi są Cypryjczycy; informacji tej nie znajdziemy w bazie danych OECD. Na szczęście, możemy posłużyć się modelem, aby uzyskać wiarygodne przewidywania: sprawdzamy wartość PKB per capita dla Cypru (22 587 dolarów), a po uruchomieniu modelu dowiadujemy się, że współczynnik satysfakcji z życia oscyluje wokół wartości $4,85 + 22\,587 \times 4,91 \times 10^{-5} = 5,96$.

Zaostrzymy Ci apetyt, pokazując w przykładzie 1.1 kod Python służący do wczytywania danych, ich przygotowania⁸, utworzenia wykresu, a także wyczerpania modelu liniowego i prognozowania wyników⁹.

Przykład 1.1. Uczenie i uruchamianie modelu liniowego za pomocą biblioteki Scikit-Learn

```
import matplotlib
import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import sklearn

# Wczytuje dane
oecd_bli = pd.read_csv("oecd_bli_2015.csv", thousands=',')
```

⁸ Zakładamy tu, że funkcja `prepare_country_stats()` została już zdefiniowana: łączy ona dane dotyczące PKB i satysfakcji z życia w jeden obiekt `DataFrame` biblioteki Pandas.

⁹ Nie szkodzi, jeżeli nie do końca rozumiesz ukazany kod; w następnych rozdziałach zajmiemy się omówieniem biblioteki Scikit-Learn.

```

gdp_per_capita = pd.read_csv("gdp_per_capita.csv",thousands=',',delimiter='\t',
                             encoding='latin1', na_values="bd")

# Przygotowuje dane
country_stats = prepare_country_stats(oecd_bli, gdp_per_capita)
X = np.c_[country_stats["PKB per capita"]]
y = np.c_[country_stats["Satysfakcja z życia"]]

# Wizualizuje dane
country_stats.plot(kind='scatter', x="PKB per capita", y='Satysfakcja z życia')
plt.show()

# Dobiera model liniowy
model = sklearn.linear_model.LinearRegression()

# Trenuje model
model.fit(X, y)

# Prognozuje wyniki dla Cypru
X_new = [[22587]] # PKB per capita dla Cypru
print(model.predict(X_new)) # generuje wynik [[ 5.96242338]]

```



Gdybyś użyła/użył algorytmu uczenia z przykładów, zauważyłabyś/zauważyłbyś, że wartość PKB najbardziej zbliżoną do Cypru ma Słowenia (20 732 dolary), a skoro dzięki danym OECD wiemy, że współczynnik satysfakcji z życia wynosi u Słowenców 5,7, moglibyśmy w drodze analogii stwierdzić, że jego wartość dla Cypryjczyków powinna być bardzo podobna. Po nieznacznym „oddaleniu skali” i przyjrzeniu się dwóm kolejnym państwom o zbliżonej wartości satysfakcji z życia okaże się, że są to Portugalia (5,1) i Hiszpania (6,5). Średnia tych trzech wartości wynosi 5,77, co stanowi wartość całkiem zbliżoną do wyliczonej przez nasz model. Wspomniany tu prosty algorytm nosi nazwę regresji **k-najbliższych sąsiadów** (ang. *k-nearest neighbors*); w tym przykładzie $k = 3$.

Aby zastąpić algorytm regresji liniowej algorytmem k-najbliższych sąsiadów, wystarczy podmienić wiersz:

```
model = sklearn.linear_model.LinearRegression()
```

na następujący:

```
model = sklearn.neighbors.KNeighborsRegressor(n_neighbors=3)
```

Jeśli wszystko zostało dobrze wykonane, Twój model powinien wyliczać dobre prognozy. W przeciwnym wypadku może być konieczne wykorzystanie większej liczby atrybutów (stopy zatrudnienia, zdrowia społeczeństwa, zanieczyszczenia powietrza itd.), zdobycie więcej danych uczących lub danych lepszej jakości, ewentualnie skorzystanie z wydajniejszego modelu (np. algorytmu regresji wielomianowej).

Podsumowując:

- Przeanalizowaliśmy dane.
- Wybraliśmy model.
- Wytrenowaliśmy go na danych uczących (np. algorytm uczący wyszukał wartości parametrów pozwalających na zminimalizowanie funkcji kosztu).

- Na końcu wykorzystaliśmy wytrenowany model do prognozowania wyników dla nowych przypadków (jest to tzw. **wnioskowanie**) z nadzieją, że będzie on skutecznie generalizował zdobytą wiedzę.

Tak właśnie wygląda typowy projekt uczenia maszynowego. W rozdziale 2. doświadczysz tego na własnej skórze podczas przygotowywania projektu od podstaw.

Poruszyliśmy do tej pory wiele zagadnień: wiesz już, czym jest uczenie maszynowe, dlaczego okazuje się przydatne i jakie są podstawowe kategorie tej dziedziny, a także znasz już podstawowy cykl pracy nad projektem uczenia maszynowego. Spójrzmy teraz, co może pójść nieprawidłowo w procesie uczenia i uniemożliwić uzyskiwanie dokładnych prognoz.

Główne problemy uczenia maszynowego

Skoro naszym głównym zadaniem jest dobór algorytmu uczącego i wyuczenie go za pomocą określonych danych, to dwoma największymi zagrożeniami są „zły algorytm” i „złe dane”. Zacznijmy od przyjrzenia się niewłaściwym danym.

Niedobór danych uczących

Aby nauczyć berbecia, czym jest jabłko, wystarczy wskazać mu ten owoc i powiedzieć „jabłko” (najprawdopodobniej musisz to jeszcze kilkakrotnie powtórzyć). Dziecko teraz będzie w stanie rozpoznawać jabłka w najróżniejszych kolorach i kształtach. Mały geniusz.

Niestety, techniki uczenia maszynowego znajdują się pod tym względem daleko w tyle; żeby większość algorytmów działała prawidłowo, należy im zapewnić mnóstwo danych. Nawet w przypadku bardzo prostych problemów należy zazwyczaj dostarczyć tysiące przykładów, natomiast bardziej zaawansowane zagadnienia, takie jak rozpoznawanie mowy, wymagają nawet milionów próbek (chyba że jesteś w stanie wykorzystać fragmenty już istniejącego modelu).

Niedorzeczna efektywność danych

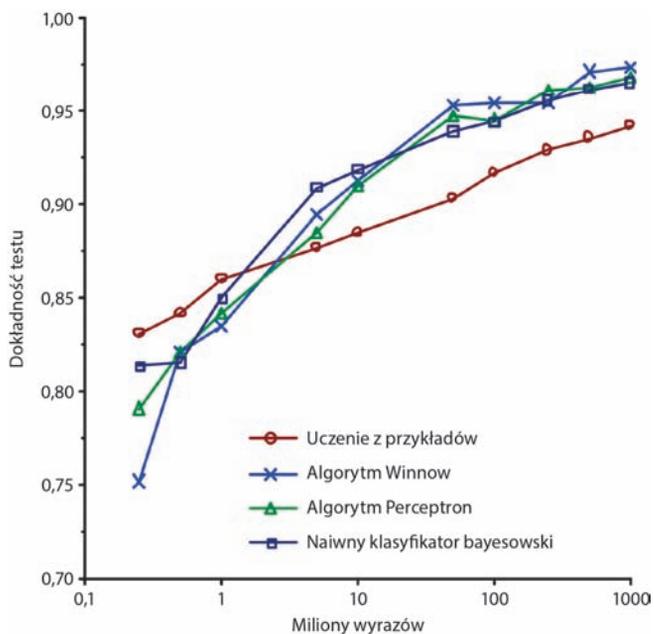
W słynnym dokumencie (<http://www.aclweb.org/anthology/P01-1005>) z 2001 roku badacze z firmy Microsoft, Michele Banko i Eric Brill, udowodnili, że różne algorytmy uczenia maszynowego, włącznie z najprostszymi, po przesłaniu im odpowiedniej ilości danych osiągają niemal identyczną wydajność wobec złożonych problemów ujednoznaczeń w języku naturalnym¹⁰ (co zaprezentowano na rysunku 1.20).

Autorzy ujmują to następująco: „Wyniki te sugerują, że warto zastanowić się nad przeniesieniem środka ciężkości w kompromisie pomiędzy poświęceniem czasu i pieniędzy na tworzenie algorytmów a przeznaczaniem zasobów na rozwój korpusu językowego”.

Koncepcja przewagi znaczenia danych nad algorytmami w rozwiązywaniu złożonych problemów została spopularyzowana przez Petera Nordviga i in. w dokumencie *The Unreasonable Effectiveness of Data* (<https://research.google.com/pubs/archive/35179.pdf>), opublikowanym w 2009 roku¹¹. Należy jednak zauważyć, że ciągle są popularne niewielkie i średniej wielkości zbiory danych, a pozyskanie ich większej ilości często bywa bardzo trudne lub kosztowne, dlatego nie bagatelizujemy znaczenia algorytmów.

¹⁰ Na przykład rozróżnianie znaczeń wyrazów „buk”, „Bug” i „bóg” w zależności od kontekstu.

¹¹ *The Unreasonable Effectiveness of Data*, Peter Norvig i in. (2009).

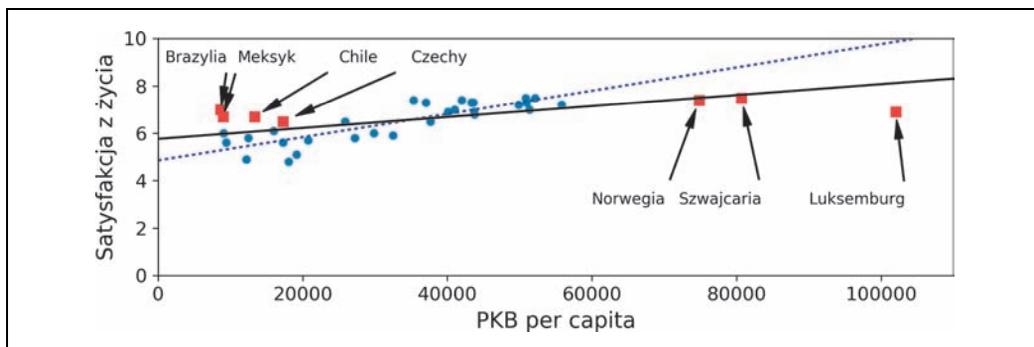


Rysunek 1.20. Porównanie znaczenia danych i algorytmów¹²

Niereprezentatywne dane uczące

Aby proces uogólniania przebiegał skutecznie, dane uczące muszą być reprezentatywne wobec nowych danych, wobec których nastąpi generalizacja. Jest to prawdziwe zarówno dla uczenia z przykładów, jak i dla uczenia z modelu.

Na przykład zbiór krajów użyty przez nas do trenowania modelu liniowego nie był doskonale reprezentatywny; brakuje w nim kilku państw. Na rysunku 1.21 widzimy, jak będą wyglądać dane po dodaniu tych kilku brakujących krajów.



Rysunek 1.21. Bardziej reprezentatywny zbiór przykładów uczących

¹² Rysunek umieszczony za pozwoleniem Banko i Brilla (2001), *Learning Curves for Confusion Set Disambiguation*.

Model wyuczony za pomocą tych danych został ukazany w postaci linii ciągłej, natomiast stary model jest reprezentowany przez linię przerywaną. Jak widać, dodanie kilku brakujących krajów nie tylko w znaczny sposób modyfikuje dany model, ale pozwala nam również uzmysłowić sobie, że prosty model liniowy raczej nigdy nie będzie bardzo dokładny. Wygląda na to, że mieszkańcy bardzo bogatych państw nie są szczęśliwsi od ich odpowiedników w średnio zamożnych krajach (wydaje się wręcz, że są bardziej nieszczęśliwi), a z kolei ludność wielu ubogich państw wydaje się znacznie szczęśliwsza od populacji bogatych krajów.

Za pomocą niereprezentatywnego zbioru danych wyuczyliliśmy model, który raczej nie będzie generował dokładnych prognoz, zwłaszcza w przypadku krajów ubogich i bardzo zamożnych.

Jest niezwykle istotne, aby zbiór danych uczących był reprezentatywny wobec generalizowanych przypadków. Często jest to trudniej osiągnąć, niż się wydaje: jeżeli przykład będzie zbyt mały, musisz liczyć się z **zaszumieniem próbkowania** (ang. *sampling noise*; niereprezentatywne dane pojawiają się tu w wyniku przypadku), ale nawet olbrzymie zbiory danych mogą nie być reprezentatywne, jeśli użyta zostanie niewłaściwa metoda próbkowania. W tym przypadku mamy do czynienia z **obciążeniem próbkowania** (ang. *sampling bias*).

Słynny przykład obciążenia próbkowania

Prawdopodobnie najsłynniejszy przykład wpływu obciążenia próbkowania można było zaobserwować podczas wyborów na prezydenta Stanów Zjednoczonych w 1936 roku, w których przeciwko Rooseveltowi stanął Landon: magazyn „Literary Digest” przeprowadził ankietę, wysyłając pocztą ulotki do około 10 milionów obywateli. Zostało odesłanych 2,4 miliona odpowiedzi, dzięki którym z dużą dozą pewności prognozowano, że Landon uzyska 57% głosów.

Ostatecznie okazało się, że to Roosevelt zebrał 62% głosów. Problem polegał na użytej metodzie próbkowania:

- Po pierwsze, magazyn „Literary Digest” zebrał adresy osób, którym została wysłana ankieta, z książek telefonicznych, list abonentów i członków klubów itd. We wszystkich tych źródłach dominowały majątne osoby, które chętniej głosowały na republikanów (czyli na Landona).
- Po drugie, odpowiedzi odesłało mniej niż 25% osób. Również ten czynnik odpowiada za obciążenie próbkowania, ponieważ nie są brane pod uwagę osoby apolityczne, przeciwnicy magazynu „Literary Digest” i inne kluczowe grupy. Jest to specjalny rodzaj obciążenia próbkowania, zwany **błędem braku odpowiedzi** (ang. *nonresponse bias*).

Weźmy pod uwagę inny przykład: powiedzmy, że chcesz stworzyć system rozpoznający teledyski z muzyką funkową. Jednym ze sposobów przygotowania zbioru danych uczących jest wyszukanie wyników „funk music” w serwisie YouTube i skorzystanie z nich. Zakładamy jednak w tym przypadku, że silnik wyszukiwania zwróci zestaw teledysków stanowiący przekrój wszystkich piosenek funkowych umieszczonych w serwisie. W rzeczywistości wyniki wyszukiwania są często tendencyjnie nastawione wobec najpopularniejszych artystów (natomiast jeśli mieszkasz w Brazylii, zostanie wyświetlonych wiele teledysków z nurtu „funk carioca”, który w ogóle nie brzmi jak James Brown). Z drugiej strony, jak inaczej można pozyskać duży zbiór danych uczących?

Dane kiepskiej jakości

Jest oczywiste, że jeśli dane uczące zawierają mnóstwo błędów, elementów odstających i szumu (np. z powodu niskiej jakości pomiarów), to systemowi będzie znacznie trudniej wykryć wzorce, przez co nie osiągnie optymalnej wydajności. Warto często poświęcić czas na oczyszczenie takich danych. Prawda jest taka, że większość analityków danych poświęca na tę czynność wiele czasu. Na przykład:

- jeżeli niektóre elementy wyraźnie odstają od reszty przykładów, wystarczy je po prostu odrzucić albo spróbować własnoręcznie poprawić błędy,
- jeśli niektórym przykładom brakuje kilku cech (np. 5% klientów nie podało wieku), musisz zdecydować, czy należy takie cechy całkowicie ignorować, ignorować te przykłady, wypełnić brakujące wartości (np. korzystając z mediany wieku), wytrenować jeden model przy użyciu tej cechy, a drugi bez niej itd.

Nieistotne cechy

Zgodnie ze słynnym powiedzeniem: z gipsu tortu nie ulepisz. Twój system będzie w stanie uczyć się jedynie za pomocą danych zawierających wystarczającą liczbę istotnych cech i niezaśmieconych nadmiarem cech nieistotnych. Elementem krytycznym dla sukcesu w projekcie uczenia maszynowego jest wybór dobrego zbioru cech uczących. Proces ten, zwany **inżynierią cech** (ang. *feature engineering*), składa się z następujących elementów:

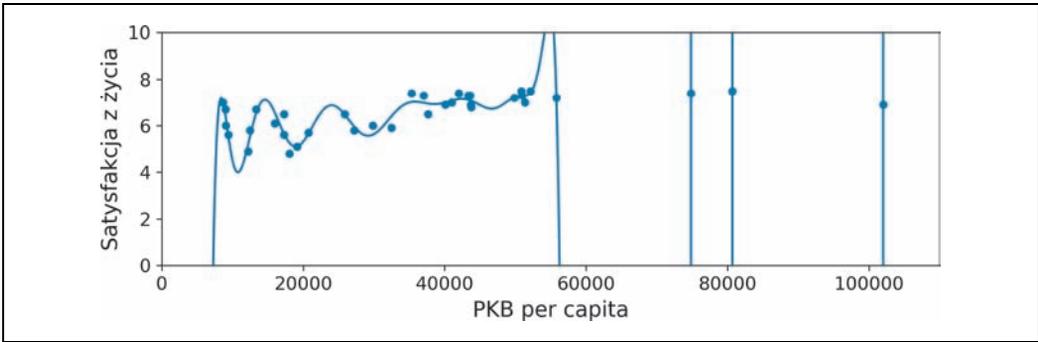
- **dobór cech** (ang. *feature selection*): dobór najprzydatniejszych spośród dostępnych cech,
- **odkrywanie cech** (ang. *feature extraction*): łączenie ze sobą istniejących cech w celu uzyskania przydatniejszej cechy (jak już wiemy, pomóc nam w tym może algorytm redukcji wymiarowości),
- uzyskiwanie nowych cech z nowych danych.

Skoro już wiemy, jak wyglądają złe dane, przeanalizujemy kwestię nieprawidłowych algorytmów.

Przetrenowanie danych uczących

Załóżmy, że na zagranicznej wycieczce zostałeś okradziony/zostałaś okradziona przez taksówkarza. Być może na Twoje usta ciśnie się stwierdzenie, że *wszyscy* taksówkarze w danym kraju są złodziejami. Ludzie często mają tendencję do nadmiernego generalizowania i, niestety, maszyny również łatwo wpadają w tę pułkę, jeśli nie zachowamy ostrożności. Zjawisko to w terminologii uczenia maszynowego nosi nazwę **przetrenowania** albo **nadmiernego dopasowania** (ang. *overfitting*). Termin ten oznacza, że model sprawdza się w przypadku danych uczących, ale sam proces uogólniania nie sprawuje się zbyt dobrze.

Na rysunku 1.22 widzimy przykład wielomianowego modelu satysfakcji z życia o olbrzymim stopniu przetrenowania za pomocą danych uczących. Nawet jeśli sprawuje się on znacznie lepiej wobec danych uczących niż zwykły model liniowy, to czy moglibyśmy zawierzyć jego prognozą?



Rysunek 1.22. Przetrenowanie danych uczących

Takie złożone modele, jak głębokie sieci neuronowe, są w stanie wykrywać subtelne wzorce w danych, ale jeżeli zbiór danych uczących jest zaszumiony lub zbyt mały (ryzyko zaszumienia próbkowania), to model ten prawdopodobnie będzie wykrywał wzorce nie w użytecznych danych, lecz w szumie. Oczywiście, szablonów tych nie da się generalizować na nowe próbki. Powiedzmy, że dostarczamy modelowi satysfakcji z życia wiele dodatkowych atrybutów, w tym takich nieprzydatnych, jak nazwy państw. W takiej sytuacji złożony model może wykrywać takie wzorce, jak np. wskaźnik satysfakcji z życia przekraczający wartość 7 w krajach mających w nazwie literę „w”: Nowa Zelandia (7,3), Norwegia (7,4), Szwecja (7,2) czy Szwajcaria (7,5). Czy masz pewność, że ta „reguła litery »w«” dotyczy również takich państw, jak Rwanda lub Zimbabwe? Oczywiście, wzorec ten wystąpił w danych uczących zupełnie przypadkowo, ale model nie jest w stanie stwierdzić, czy taki szablon jest rzeczywisty, czy stanowi wynik zaszumienia danych.



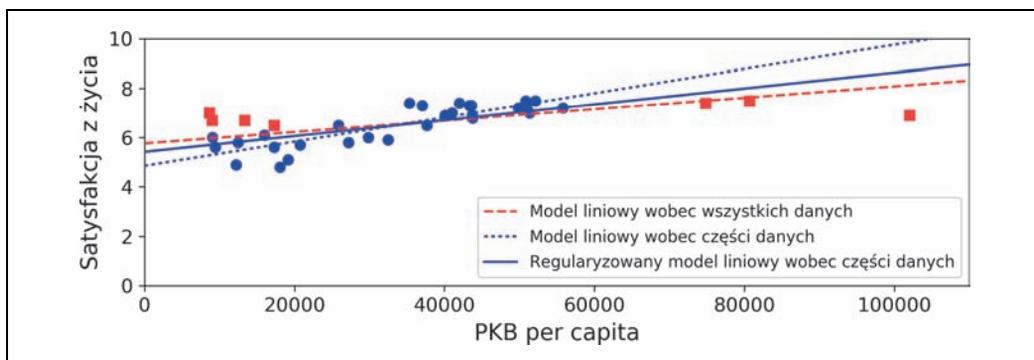
Zjawisko przetrenowania występuje, gdy model jest zbyt skomplikowany w porównaniu do ilości lub zaszumienia danych uczących. W takich przypadkach możliwe są następujące rozwiązania:

- uproszczenie modelu poprzez wybór zawierającego mniej parametrów (np. modelu liniowego w miejsce modelu wielomianowego), zmniejszenie liczby atrybutów w danych uczących lub ograniczenie modelu,
- uzyskanie większej ilości danych uczących,
- zmniejszenie zaszumienia danych uczących (np. poprzez usunięcie błędnych danych lub elementów odstających).

Ograniczenie modelu w celu jego uproszczenia i zmniejszenia ryzyka przetrenowania nosi nazwę **regularyzacji** (ang. *regularization*). Przykładowo, zdefiniowany przez nas wcześniej model liniowy zawiera dwa parametry: θ_0 i θ_1 . Algorytm uczący uzyskuje w ten sposób dwa **stopnie swobody**, za pomocą których możemy dostosowywać model do danych uczących: mamy wpływ zarówno na wysokość (θ_0), jak i nachylenie (θ_1) funkcji. Gdybyśmy wymusili wartość $\theta_1 = 0$, algorytm zostałby ograniczony do jednego stopnia swobody i jego prawidłowe wytrenowanie stałoby się znacznie trudniejsze: moglibyśmy jedynie zmieniać położenie funkcji wzdłuż osi y i starać się umieścić jej przebieg jak najbliżej punktów uczących, co w tym przypadku oznaczałoby określenie wartości średniej. Zaiste, bardzo prosty model! Jeżeli pozwolimy algorytmowi na nieznaczne zmiany wartości parametru θ_1 , to w istocie algorytm będzie znajdował się gdzieś pośrednio pomiędzy jednym a dwoma

stopniami swobody. Uzyskamy w ten sposób model prostszy od mającego dwa stopnie swobody, ale bardziej złożony od jednostopniowego. Musimy wyszukiwać odpowiednią równowagę pomiędzy perfekcyjnym dopasowywaniem danych a zachowaniem odpowiedniej prostoty modelu pozwalającej na generalizowanie wyników.

Rysunek 1.23 prezentuje trzy modele: linia kropkowana symbolizuje nasz oryginalny model wytrenowany pod nieobecność kilku krajów, linia kreskowana ukazuje nasz model wyuczony przy użyciu pełnego zakresu danych, natomiast linia ciągła reprezentuje model wytrenowany za pomocą takich samych danych, jak w pierwszym przypadku, ale przy wprowadzonej regularyzacji. Widzimy, że w wyniku regularyzacji funkcja ma trochę mniejsze nachylenie, co sprawia, że jest nieco mniej dopasowana do danych uczących, ale za to sprawuje się lepiej w procesie uogólniania nowych przykładów.



Rysunek 1.23. Regularyzacja zmniejsza ryzyko przetrenowania

Stopień regularyzacji przeprowadzanej na etapie nauki możemy kontrolować za pomocą **hiperparametrów** (ang. *hyperparameters*). Najprościej mówiąc, hiperparametrami nazywamy parametry algorytmu uczącego (nie całego modelu). Nie są one modyfikowane przez sam algorytm uczący; należy je wyznaczyć tuż przed rozpoczęciem procesu uczenia i w jego trakcie ich wartości pozostają niezmiennie. Jeśli wyznaczysz bardzo dużą wartość hiperparametru regularyzacji, uzyskana funkcja będzie niemal płaska (będzie miała niemal zerowe nachylenie); algorytm uczący prawie na pewno nie ulegnie przetrenowaniu, ale jednocześnie trudniej będzie mu znaleźć prawidłowe rozwiązanie. Strojenie hiperparametrów stanowi istotną część tworzenia systemu uczenia maszynowego (w następnym rozdziale przekonasz się o tym na własnej skórze).

Niedotrenowanie danych uczących

Jak łatwo się domyślić, **niedotrenowanie** (ang. *underfitting*) stanowi przeciwieństwo przetrenowania: występuje ono wtedy, gdy model jest zbyt prosty, aby wyuczyc się struktur danych uczących. Na przykład, liniowy model satysfakcji z życia jest podatny na niedotrenowanie; rzeczywistość jest po prostu bardziej skomplikowana od modelu, zatem jego prognozy nie będą dokładne, nawet w przypadku danych uczących.

Głównymi sposobami rozwiązania tego problemu są:

- wybór potężniejszego modelu wykorzystującego większą liczbę parametrów,
- dołączenie większej liczby cech do algorytmu uczącego (inżyniera cech),
- zmniejszenie ograniczeń modelu (np. zredukowanie hiperparametru regularyzacji).

Podsumowanie

Wiemy już coraz więcej na temat uczenia maszynowego. Wspomnieliśmy jednak o tak wielu pojęciach, że możesz poczuć zagubienie, dlatego zatrzymajmy się na chwilę i sprawdźmy, co już wiemy:

- Uczenie maszynowe polega na ciągłym usprawnianiu wykonywania zadań przez urządzenie poprzez jego uczenie się z danych, a nie jawne dopisywanie nowych reguł do kodu.
- Istnieje wiele różnych metod uczenia maszynowego: nadzorowane, nienadzorowane, wsadowe, przyrostowe, z przykładów, z modelu itd.
- W projekcie uczenia maszynowego gromadzimy dane w zbiorze uczącym, który następnie dostarczamy algorytmowi uczącemu. Jeśli wykorzystujemy technikę uczenia z modelu, pewne parametry zostają dostrójone do zbioru danych uczących (np. w celu uzyskiwania dokładnych prognoz wobec samych danych uczących), dzięki czemu model ten może również dobrze przewidywać wyniki wobec nowych przykładów. Z kolei w metodach uczenia z przykładów próbki zostają po prostu wyuczone „na pamięć”, a model wykorzystuje miarę podobieństwa do uogólniania wyników na nowe przykłady.
- System nie będzie dobrze spełniał swojego zadania, jeżeli zbiór danych testowych jest zbyt mały, same dane są niereprezentatywne, zaszumione lub zaśmiecone nieistotnymi cechami (z gipsu tortu nie ulepisz). Model nie może być również zbyt prosty (niedotrenowanie) ani zbyt skomplikowany (przetrenowanie).

Musimy omówić jeszcze jedno istotne zagadnienie: po wyuczeniu modelu nie wystarczy sama nadzieja na skuteczne uogólnianie nowych próbek. Chcesz ocenić wydajność modelu i w razie potrzeby dostrójć go do swoich potrzeb. Dowiesz się teraz, jak tego dokonać.

Testowanie i ocenianie

Jedynym sposobem sprawdzenia, jak dobrze model generalizuje wyniki, jest wypróbowanie go na nowych danych. Możemy, na przykład, zaimplementować go w środowisku produkcyjnym i monitorować jego wydajność. Jest to skuteczne rozwiązanie, ale jeśli model okaże się beznadziejny, użytkownicy będą narzekać — a tego lepiej unikać.

Lepszym pomysłem jest podział danych na dwa zestawy: **zbiór uczący** (ang. *training set*) i **zbiór testowy** (ang. *test set*). Zgodnie z tymi nazwami, model trenujemy za pomocą zbioru uczącego, a sprawdzamy przy użyciu zbioru testowego. Współczynnik błędów uzyskiwany dla nowych przykładów nosi nazwę **błędów uogólniania** (lub **błędów generalizacji**), a dzięki zbiorowi testowemu możemy oszacować jego wartość. Parametr ten mówi nam również, jak dobrze model będzie się spisywał wobec nieznanych danych.

Jeśli wartość błędu uczenia jest niewielka (tzn. model rzadko się myli wobec zbioru uczącego), ale błąd uogólniania jest duży, oznacza to, że model jest przetrenowany.



Zazwyczaj na zbiór uczący składa się 80% danych, a pozostałe 20% *przechowujemy* w celu testowania.

Zatem ocena modelu nie jest wcale taka trudna: wystarczy użyć zbioru testowego. Załóżmy teraz, że wahamy się przy wyborze jednego z dwóch modeli (np. pomiędzy wyborem modelu liniowego a wielomianowego): w jaki sposób mamy zdecydować? Jednym z rozwiązań jest wyuczenie obydwu modeli i sprawdzenie, jak sobie radzą z uogólnianiem wobec zbioru testowego.

Przyjmijmy teraz, że model liniowy radzi sobie lepiej z uogólnianiem, ale chcesz wprowadzić regularyzację w celu uniknięcia przetrenowania. Rodzi się w tym momencie pytanie: w jaki sposób dobrać wartość hiperparametru regularyzacji? Możesz na przykład wytrenować 100 różnych modeli, z których każdy będzie miał inną wartość tego hiperparametru. Załóżmy, że znajdziesz optymalną wartość tego hiperparametru, dającą model o najniższym błędzie generalizacji (np. rzędu 5%).

Zatem wprowadzasz ten model do środowiska produkcyjnego, ale niestety, sprawuje się on poniżej oczekiwań i generuje 15% błędów. Co się właściwie stało?

Problem polega na tym, że wielokrotnie mierzyłaś/mierzyłeś błąd uogólniania wobec zbioru testowego i dostosowałaś/dostosowałeś go *do tych danych*. Oznacza to, że teraz prawdopodobnie model ten nie będzie sobie radził zbyt dobrze z nowymi próbkami.

Powszechnie stosowanym rozwiązaniem tego problemu jest odłożenie jeszcze jednego zbioru danych, zwanego **zbiorem walidacyjnym** (ang. *validation set*). Trenujemy wiele modeli mających różne wartości hiperparametrów za pomocą zbioru uczącego, dobieramy model i hiperparametry najlepiej sprawujące się wobec zbioru walidacyjnego, a gdy uzyskiwane wyniki będą Cię zadowalać, wystarczy przeprowadzić ostatni sprawdzian wobec zbioru testowego, aby oszacować wartość błędu uogólniania.

Aby uniknąć „zużywania” zbyt dużej ilości danych uczących na zbiór walidacyjny, możesz skorzystać z techniki zwanej **sprawdzianem krzyżowym** lub **krosvalidacją** (ang. *cross-validation*): zbiór uczący zostaje rozdzielony na wzajemnie uzupełniające się podzbiory; każdy model jest uczony za pomocą różnych kombinacji tych podzbiorów i oceniany przy użyciu pozostałych, nieużytych podzestawów. Po dobraniu rodzaju i hiperparametrów modelu ostateczny model zostaje wytrenowany wobec pełnego zbioru uczącego, a błąd uogólniania określamy przy użyciu zbioru testowego.

Twierdzenie o nieistnieniu darmowych obiadów

Model stanowi uproszczoną postać obserwacji. Uproszczenia wynikają z konieczności unikania niepotrzebnych szczegółów, które nie ułatwiają procesu uogólniania. Podczas dobierania i odrzucania danych musimy jednak przyjmować **założenia**. Na przykład w modelu liniowym zakładamy, że dane są, nomen omen, liniowe, a odległość pomiędzy przykładami i wykresem funkcji stanowi wyłącznie szum, który możemy bezpiecznie zignorować.

W słynnej publikacji z 1996 roku¹³ (<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.390.9412&rep=rep1&type=pdf>) David Wolpert udowodnił, że jeśli nie przyjmimy jakichkolwiek założeń dotyczących danych, to okaże się, że żaden model nie będzie lepszy od pozostałych. Jest to tak zwane **twierdzenie o nieistnieniu darmowych obiadów** (ang. *No Free Lunch Theorem* — NFL). Dla pewnych zbiorów danych najlepiej nadaje się model liniowy, natomiast dla innych — sieci neuronowe. Nie istnieje żaden model, który z założenia będzie działał lepiej (stąd nazwa twierdzenia). Jedynie poprzez ocenę działania każdego modelu możemy przekonać się, który z nich będzie sprawował się najlepiej. Jest to niewykonalne, dlatego w praktyce przyjmujemy rozsądne założenia dotyczące danych i oceniamy działanie tylko kilku rozsądnie dobranych modeli. Na przykład wobec prostych zadań możemy ocenić działanie modeli liniowych różniących się stopniem regularyzacji, a bardziej skomplikowane problemy możemy przetestować przy użyciu sieci neuronowych.

Ćwiczenia

W rozdziale tym poznaliśmy część najważniejszych koncepcji opisujących dziedzinę uczenia maszynowego. W kolejnych rozdziałach skoncentrujemy się na szczegółach i pisaniu właściwego kodu, zanim jednak przejdiesz dalej, upewnij się, że jesteś w stanie odpowiedzieć na poniższe pytania:

1. Podaj definicję uczenia maszynowego.
2. Wymień cztery rodzaje problemów, z których rozwiązaniem najlepiej sobie radzi proces uczenia maszynowego.
3. Czym jest oznakowany zbiór danych uczących?
4. Jakie są dwa najczęstsze zastosowania uczenia nadzorowanego?
5. Wymień cztery najpowszechniejsze zastosowania uczenia nienadzorowanego.
6. Jakiego rodzaju algorytmu uczenia maszynowego użyłabyś/użyłbyś w robocie przeznaczonym do poruszania się po nieznanym terenie?
7. Jakiego rodzaju algorytmu uczenia maszynowego użyłabyś/użyłbyś do rozdzielania klientów na kilka różnych grup?
8. Czy problem wykrywania spamu stanowi część mechanizmu uczenia nadzorowanego lub nienadzorowanego?
9. Czym jest system uczenia przyrostowego?
10. Czym jest uczenie pozakorowe?
11. W jakim algorytmie uczenia maszynowego jest wymagana miara podobieństwa do uzyskiwania prognoz?
12. Wyjaśnij różnicę pomiędzy parametrem modelu a hiperparametrem algorytmu uczącego.
13. Czego poszukują algorytmy uczenia z modelu? Z jakiej strategii najczęściej korzystają? W jaki sposób przeprowadzają prognozy?

¹³ *The Lack of A Priori Distinctions Between Learning Algorithms*, D. Wolpert (1996).

14. Wymień cztery główne problemy związane z uczeniem maszynowym.
15. Z czym mamy do czynienia, jeżeli model sprawuje się znakomicie wobec danych uczących, ale nie radzi sobie z uogólnianiem wobec nowych próbek? Wymień trzy możliwe rozwiązania.
16. Czym jest zbiór testowy i dlaczego powinniśmy z niego korzystać?
17. Do czego służy zbiór walidacyjny?
18. Co nam grozi w przypadku strojenia hiperparametru wobec zbioru testowego?
19. Czym jest sprawdzian krzyżowy oraz dlaczego warto go stosować wobec zbioru walidacyjnego?

Odpowiedzi znajdziesz w dodatku A.

A

- AdaBoost, 197
- adaptacyjny współczynnik uczenia, 297
- agent, 33, 432
- agregacja, 187, 191
 - przykładów wstępnych, 190
- aktualizacje
 - asynchroniczne, 346
 - synchroniczne, 345
- AlexNet, 364
- algorytmy
 - AdaBoost, 199
 - AdaGrad, 296
 - Adam, 298
 - dynamicznego rozmieszczania, 318
 - genetyczne, 434
 - Isomap, 226
 - iteracji wartości, 447
 - LLE, 224
 - lokalnie liniowego zanurzenia, 224
 - Nesterova, 295
 - normalizacji wsadowej, 282
 - pozapolityczne, 451
 - Q-uczenia, 449, 451
 - REINFORCE, 441
 - RMSProp, 298
 - stochastycznego spadku wzdłuż gradientu, 130
 - uczące
 - CART, 179
 - dostarczanie danych, 241
 - uczenia maszynowego, 77
 - zachłanne, 179
- analiza
 - błędów, 111
 - głównych składowych, PCA, 215
 - najlepszych modeli, 91
 - PCA, 215, 409
 - główne składowe, 215
 - jądrowa, 221
 - liczba wymiarów, 218
 - losowa, 220
 - przyrostowa, 220
 - rzutowanie, 217
 - współczynnik wariancji wyjaśnionej, 217
 - zachowanie wariancji, 215
 - zastosowania, 219
 - projektu, 55
 - sentymentów, 377
 - skupień, 30
 - hierarchiczna, 30
 - słowotwórcza, 117
- aplikacja Jupyter, 62
- architektura
 - CNN
 - AlexNet, 364
 - GoogLeNet, 366, 368
 - LeNet-5, 362
 - ResNet, 369
 - kory wzrokowej, 352
 - architektury sieci neuronowych, 505
 - asynchroniczna komunikacja, 328
 - atrybut, 29
 - atrybuty kategorialne, 80
 - autokodery, 407, 427
 - antagonistyczne, 427
 - generatywne, 423
 - kurczliwe, 427
 - liniowe, 409
 - odszumiające, 419
 - implementacja, 420
 - stosowe, 419

autokodery
 probabilistyczne, 423
 rzadkie
 implementacja, 422
 stosowe, 410
 implementacja, 411
 nienadzorowane uczenie wstępne, 417
 splotowe, 427
 uczenie pojedynczo, 413
 symetryczne, 412
 wariacyjne, 423
 WTA, 427
automatyczne różniczkowanie, 239

B

biblioteka
 Caffe, 233
 CUDA, 315
 cuDNN, 315
 Deeplearning4j, 233
 H2O, 233
 LIBLINEAR, 163
 LIBSVM, 163
 MXNet, 233
 Scikit-Learn, 40
 TensorFlow, 233
 Theano, 233, 288
 Torch, 233
biblioteki uczenia głębokiego, 233
błędy
 analiza, 111
 braku odpowiedzi, 44
 MSE, 87, 243, 390
 przeciwobrazu rekonstrukcji, 223
 rekonstrukcji, 219
 resztowe, 200
 RMSE, 85
 uogólniania, 48
bramka
 wejściowa, 398
 wyjściowa, 398
 zapominająca, 398
bramkowana jednostka rekurencyjna, 400
bufor protokołu, 325

C

cecha, 29
cechy podobieństwa, 160
cel uczenia, 166

charakterystyka robocza odbiornika, 106
CNN, convolutional neural networks, 351
cykl życia wartości w węźle, 236
człon regularyzacyjny, 139
czułość, 101
czytniki wieloklasowe, 336

D

dane
 kiepskiej jakości, 45
 zaszumione, 38
DBN, deep belief network, 509
dekoder, 408
długa pamięć krótkotrwała, 397
dobór
 cechy, 45
 liczby drzew, 203
 modelu, 38
dogenerowanie danych, 308
dominanta, 191
drzewa
 binarne, 177
 decyzyjne, 87, 175
 granice decyzyjne, 178
 niestabilność, 184
 regresyjne, 182
 szacowanie prawdopodobieństw, 178
 uczenie, 175
 wizualizowanie, 175
 wyliczanie prognoz, 176
 z agregacją, 192
 złożoność obliczeniowa, 180
 klasyfikacyjne, 179
 regresyjne, 179
 wzmacniane gradientowo, 200
dualność, 493
dynamiczne rozmieszczanie operacji, 318
dywergencja Kullbacka-Leiblera, 152

E

efektywność parametrów, 271
ELU, exponential linear unit, 280
entropia, 180
 krzyżowa, 151, 152
estymatory, 79

F

filtry, 354
FNN, feedforward neural network, 264
FTRL, follow the regularized leader, 300
funkcje
 aktywacji, 263, 273
 ELU, 280
 logistyczna, 278
 nasylenie, 277
 nienasycające, 278
 problem śmierci ReLU, 278
 przeciekająca ReLU, 279
 ReLU, 264, 278
 tangens hiperboliczny, 264, 278
 typu softmax, 264
błędu MSE, 390
decyzyjne, 102, 165
dopasowania, 39
dynamic_rnn(), 384
dynamicznego rozmieszczania, 320
get_variable(), 251
kosztu, 146
kosztu algorytmu CART, 179
Lagrange'a, 493
logistyczne, 146
make(), 436
pełności, 105
podobieństwa, 160
pomocnicze, 338
prognozy, 165
RBF, 161
relu(), 251
sigmoidalna, 146
skokowa Heaviside'a, 259
softmax, 151
static_rnn(), 382
straty, 147
straty zawiasowa, 173
tf.layers.conv1d(), 372
tf.layers.conv2d_transpose(), 373
tf.layers.conv3d(), 372
tf.layers.separable_conv2d(), 373
tf.nn.atrous_conv2d(), 372
tf.nn.depthwise_conv2d(), 373
typu softmax, 150
uczące, 146
użyteczności, 39

G

gaussowskie jądro RBF, 161
generatywna sieć stochastyczna, 427
generowanie cyfr, 426
globalne uśrednione łączenia, 368
głęboka sieć neuronowa, GSN, 262
głębokie
 Q-uczenie, 452
 sieci przekonań, DBN, 509
 sieci rekurencyjne, RSN, 393
głosowanie
 miękkie, 190
 większościowe, 186, 188
główne składowe, 215
GoogLeNet, 366
gradienty
 eksplodujące, 275
 obcinanie, 285
 obliczanie, 238
 polityk, 441
 polityk/reguła, 431
 proste, 124, 297
 proste wsadowe, 127
 przedawnione, 346
 wzmacnianie drzew, 200
 zanikające, 275
grafy, 234
 obliczeniowe, 231
 tworzenie, 234
 wizualizacja, 243
 zarządzanie, 236
granice decyzyjne, 148, 149, 214
 predyktorów, 198
grupy zadań, 323

H

harmonogramy uczenia, 131, 300
 potęgowe, 301
 wydajnościowe, 301
 wykładnicze, 301
hesjan, 299
hiperparametry
 regularyzacyjne, 181
 sieci neuronowej, 271
hiperpłaszczyzna, 215
hipoteza, 58
 zerowa, 181

histogram

- atrybutu numerycznego, 66
- kategorii dochodów, 70

I

implementacja

- autokodera stosowego, 411
- autokoderów odszumiających, 420
- metody gradientu prostego, 238
- normalizacji wsadowej, 283
- sprawdzianu krzyżowego, 99

inicjacja

- Glorota, 277
- losowa, 124
- wag, 276
- Xaviera, 277

inspekcja, 79

instalacja modułu TensorFlow, 234

interfejs API TF.Learn, 232

inżynieria cech, 45

istotność cech, 195

iteracja Q-wartości, 448

J

jakobian, 299

jądra

- drugiego stopnia, 165
- gaussowskie RBF, 171
- liniowe, 171
- łańcuchowe, 162
- łączące, 360
- podsekwencji łańcucha znaków, 162
- RBF, 161
- sigmoidalne, 171
- wielomianowe, 159, 160
 - drugiego stopnia, 171

jądrowa analiza PCA, 221

jeden

- przeciw jednemu, 109
- przeciw reszcie, 109
- przeciw wszystkim, 109

jednokierunkowa sieć neuronowa, FNN, 264

jednolitość, 79

jednostka

- GRU, 401
- liniowa, 247
- ReLU, 247, 252
- rezydualna, 370
- wykładniczo-liniowa, 280

K

kanał barw, 356

karty graficzne, 316

kernelizowane maszyny SVM, 170

k-krotny sprawdzian krzyżowy, 86

klasa

- AdaBoostClassifier, 200
- BaggingClassifier, 192, 194
- Coordinator, 336
- GradientBoostingRegressor, 201
- LabelBinazer, 81
- LinearSVC, 163
- OneVsOneClassifier, 110
- OneVsRestClassifier, 110
- PCA, 217
- PolynomialFeatures, 135
- QueueRunner, 336
- RandomForestClassifier, 107, 194
- ReLU, 250
- RMSPropOptimizer, 298
- StratifiedKfold, 99
- StratifiedShuffleSplit, 70
- SyncReplicasOptimizer, 349
- tensorflow, 266

klaster, 323

- TensorFlow, 323

klasyfikacja, 28, 95

- miękkiego marginesu, 156
- wieloetykietowa, 114
- wieloklasowa, 108
- wielowyjściowa, 115

klasyfikatory

- AdaBoost, 197
- binarne, 97
- głosujące, 187
- głosujące większościowo, 187
- RandomForestClassifier, 107, 110
- sekwencji, 387
- SGDClassifier, 103
- silne, 188
- słabe, 188
- SVM, 158

klasyfikowanie maksymalnego marginesu, 155

klątwa wymiarowości, 209, 210

k-najbliższych sąsiadów, 41

koder, 408

kodowanie gorącojedynkowe, 80

kolejca tokenów, 349

- kolejka, 328
 - PaddingFIFOQueue, 332
 - RandomShuffleQueue, 332
- kolejki
 - krotek, 330
 - umieszczanie danych, 329
 - usuwanie danych, 330
 - zamykanie, 331
- kombinacje atrybutów, 76
- komórka
 - GRU, 400
 - LSTM, 397
- komórki pamięci, 380
- kompozycja, 80
- kompresja, 219
- komunikacja asynchroniczna, 328
- koneksjonizm, 262
- konfiguracja głębokiej sieci neuronów, 309
- kontaminacja, 187, 204
- kontenery zasobów, 327
- korelacja, 73
- korzeń, 176
- krosvalidacja, 49, 86
- krzywe uczenia, 135, 137, 243
 - modelu wielomianowego, 138
- kurczliwość, 201
- kwantyzacja sieci neuronowej, 348
- kwartył, 65

L

- lagranżjan, 493
 - uogólniony, 494
- las losowy, 187, 194
- LASSO, 141
- LDA, linear discriminant analysis, 226
- LeNet-5, 362
- liczba
 - neuronów, 272
 - ukrytych warstw, 271
- liniowa
 - analiza dyskryminacyjna, LDA, 226
 - klasyfikacja SVM, 155
- liniowy zbiór danych, 122
- liść, 176
- LLE, locally linear embedding, 224
- lokalne pola recepcyjne, 352
- lokalnie liniowe zanurzenie, LLE, 224
- losowa analiza PCA, 220

- losowanie warstwowe, 69
- losowy las, 187, 194
- luz dopełniający, 494

Ł

- łańcuchy Markowa, 446
- łączenie prognoz, 205

M

- macierz
 - jednostkowa, 140
 - korelacji, 77
 - parametrów Θ , 150
 - pomyłek, 100
 - rzadka, 81
- MAE, 58
- mapa cech ϕ , 222
- mapy samoorganizujące, 511
- maszyny
 - Boltzmann, 506
 - ograniczone, 508
 - wektorów nośnych, SVM, 155, 493
- mądrość tłumu, 187
- MDS, multidimensional scaling, 226
- mechanizm
 - AdaBoost, 197
 - uwagi, 405
- metody
 - automatycznego wyliczania gradientów, 240
 - describe(), 65
 - elastycznej siatki, 143
 - export_graphviz(), 175
 - fit(), 83
 - gradientu prostego, 124, 238
 - hist(), 66
 - LASSO, 141, 143
 - mode(), 186
 - podprzestrzeni losowych, 194
 - porzucania, 395
 - predict_proba(), 192
 - rejonów losowych, 194
 - różnic czasowych, 449
 - toarray(), 81
 - uczenia zespołowego, 91, 187
 - wczesnego zatrzymywania, 203
- metryka wydajności, 57

miara
 podobieństwa, 37
 wydajności, 98
 zanieczyszczenia, 177
mikser, 204
minigrupy, 132
MLP, multi-layer perceptron, 262
mnożniki
 Karusha-Kuhna-Tuckera, 494
 Lagrange'a, 493
modele
 białej skrzynki, 178
 czarnej skrzynki, 178
 DecisionTreeRegressor, 86
 drzewa decyzyjnego, 87
 liniowe, 38, 40
 nieparametryczne, 181
 parametryczne, 181
 regresji wielomianowej, 135
 rzadkie, 142, 300
 tłumaczenia maszynowego, 403
 z innych środowisk, 288
 zrównolegane, 343
moduł
 Scikit-Learn, 69, 79, 101, 110
 TensorFlow, 231
modułowość, 247
moment, 294

N

nadmierne dopasowanie, 45
nagroda, 33
naruszanie marginesu, 156
narzędzie
 OpenAI gym, 435
 pip, 234
 TensorBoard, 233
nasylenie przepustowości, 347
nauka gry, 452
neurony
 biologiczne, 256
 obciążeniowe, 260
 sztuczne, 257
 rekurencyjne, 378
 wejściowe, 260
niedobór danych uczących, 42
niedorzeczna efektywność danych, 42
niedotrenowanie danych uczących, 47

nieistotne cechy, 45
nieliniowa klasyfikacja SVM, 158
nienadzorowane uczenie wstępne, 291
niereprezentatywne dane uczące, 43
nierozprzestrzenianie klas, 79
niestandardowe transformatory, 81
norma, 58
 euklidesowa, 58
 ℓ_k , 59
 taksówkowa, 59
normalizacja wsadowa, 281
 implementacja, 283

O

obciążalność
 wejściowa, 277
 wyjściowa, 277
obciążenie, 138
 próbekowania, 44
 związane z podglądaniem danych, 67
obcinanie gradientu, 285
obliczanie gradientów, 238
obsługa
 sekwencji wejściowych, 385
 sekwencji wyjściowych, 386
 tekstu, 80
ocena modelu, 85
ocenywanie, 48
oczyszczanie danych, 78
odchylenie standardowe σ , 65, 277, 423
odkrywanie cechy, 45
odległość
 geodezyjna, 226
 Levenstheina, 162
odsetek
 fałszywie pozytywnych, 106
 prawdziwie negatywnych, 106
 prawdziwie pozytywnych, 106
odwrotne różniczkowanie automatyczne, 240, 501
ograniczona propagacja wsteczna w czasie, 396
ograniczone maszyny Boltzmanna, 292, 508
określanie zakresu problemu, 55
OOB, out-of-bag instances, 193
operacje
 czytnikowe, 334
 na urządzeniach, 318
 TensorFlow, 313
 w wielu zadaniach, 325

- optymalizacja
 - Adam, 298
 - momentum, 294, 296
 - ograniczona, 167
- optymalizatory, 240, 293
- optymalna wartość stanu, 447
- otwieranie sesji, 324

P

- pamięć operacyjna karty graficznej, 316
- parametry polityki, 434
- PCA, Principal Component Analysis, 215
- pełność, 101
- percentyl, 65
- perceptron, 259
 - wielowarstwowy, MLP, 262
- pierwiastek błędu średniokwadratowego, 57
- plik zdarzeń, 244
- pobieranie danych, 62
- pochodne cząstkowe, 127
 - drugiego rzędu, 299
 - pierwszego rzędu, 299
- podążanie za regularyzowanym liderem, FTRL, 300
- podglądanie danych, 67
- podpróbkowanie, 360
- podprzestrzenie losowe, 194
- podprzestrzeń, 211
- podzbiór, 86
- polityka, 33, 433
- polityka poszukiwania, 449, 451
- połączenia
 - pomijające, 371
 - przezierne, 400
 - skrótowe, 369
- pomiar dokładności, 98
- porzucanie, 304, 395
- poszukiwanie, 451
- potoki, 56
 - transformujące, 83
- prawdopodobieństwo, 146
 - przynależności do klasy, 178
 - utrzymywania, 305
- prawo wielkich liczb, 189
- precyzja, 101
- predyktory, 79
 - miksujące, 205

- problem
 - dualny, 169
 - NP-zupełny, 179
 - przypisania zasługi, 440
 - śmierci ReLU, 278
- procesory tensorowe, 314
- procesy decyzyjne Markowa, 446
- prognozy
 - klasyfikatora regresji Softmax, 151
 - modelu regresji wielomianowej, 135
 - szeregu czasowego, 391
- program TensorBoard, 245
- programowanie
 - dynamiczne, 448
 - kwadratowe, 168
- projekt, 53
- propagacja wsteczna, 262
- prostowane jednostki liniowe, ReLU, 247
- protokół gRPC, 325
- próbkowana funkcja softmax, 405
- próg decyzyjny, 102
- przeciekająca funkcja ReLU, 279
- przeciwobraz, 222
- przestrzeń
 - cech, 221
 - przeszukiwania, 90
 - ukryta, 424
- przeszukiwanie
 - losowe, 90
 - siatki, 88
- przetrenowanie, 45, 302
- przetwarzanie
 - języka naturalnego, 351, 401
 - równoległe, 321, 339
- przewidywanie, 38
 - szeregów czasowych, 389
- przybliżający algorytm Q-uczenia, 451
- przykłady pozatreningowe, OOB, 193
- przypinanie operacji, 325
- przyrostowe
 - algorytmy PCA, 220
 - maszyny SVM, 172
- punkt
 - charakterystyczny, 160
 - obciążenia, 120
 - przebieg, 120

Q

Q-sieć, 454
Q-uczenie głębokie, 452
Q-wartość, 448

R

rachunek zdań, 256
radialna funkcja bazowa, 161
ramka, 378
redukowanie wymiarowości, 32, 209

- algorytm Isomap, 226
- algorytm LLE, 224
- analiza PCA, 215
- jądrowa analiza PCA, 221
- klątwa wymiarowości, 210
- liniowa analiza dyskryminacyjna, 226
- skalowanie wielowymiarowe, 226
- strategie, 211
- t-SNE, 226

regresja, 182

- grzbietowa, 139, 140
- k-najbliższych sąsiadów, 41
- liniowa, 120, 123
 - użycie modułu TensorFlow, 237
- logistyczna, 29, 145
 - wieloraka, 150
- metodą LASSO, 141
- softmax, 150
- SVM, 163
- wielomianowa, 133
 - wysokiego stopnia, 136
- z jedną zmienną, 57
- z wieloma zmiennymi, 57

regresyjne drzewo decyzyjne, 182
regularyzacja, 46, 303

- grzbietowa, 143
- metodą wczesnego zatrzymywania, 145
- przez porzucanie, 304
- Tichonowa, 139
- typu max-norm, 306

regularyzowane modele liniowe, 139
reguła Hebba, 260
rejony losowe, 194
rekonstrukcja, 409, 415
rekurencyjne sieci neuronowe, RSN, 377
ReLU, rectified linear unit, 247

replikacja

- międzygrafowa, 341
- wewnątrzgrafowa, 341

repliki zapasowe, 349
repozytoria modeli, 291
reprezentacje

- danych, 407
- wektorowe słów, 402

ResNet, Residual Network, 369
RMSE, 57
rozdzielanie zmiennych, 326
rozgłaszanie, 268
rozkład

- długoogonowy, 67
- Gausa, 65
- według wartości osobliwych, 216
- wykładniczy, 284

rozmaitości, 213
rozpoznawanie mowy, 351
rozproszenie, 75
rozsądne wartości domyślne, 80
równanie

- kwadratowe, 134
- normalne, 121
- optymalności Bellmana, 447

równoległe obliczenia, 232
różniczkowanie

- automatyczne, 233, 240, 497, 500
 - odwrotne, 240, 501
- numeryczne, 240, 499
- ręczne, 497
- symboliczne, 239, 240, 498

RNN, recurrent neural networks, 377
rzadkość, 420
rzutowanie, 211

- na d wymiarów, 217
- wyników, 390
- zbioru danych, 215

S

SAMME, 199
schodzenie po gradiencie z minigrupami, 132
sekwencja gradowa, 408
serwer parametrów, 323
sesja

- lokalna, 327
- rozproszona, 327

SGD, Stochastic Gradient Descent, 98, 130

sieci neuronowe, 229, 255

- Hopfielda, 505
- funkcje aktywacji, 273
- głębokie, GSN, 262, 265, 275
 - funkcje aktywacji, 278
 - harmonogramowanie współczynnika uczenia, 300
 - konfiguracja, 309
 - normalizacja, 281
 - optymalizatory, 293
 - regularyzacja, 302
- jako polityki, 438
- jednokierunkowe, FNN, 264
- kwantyzacja, 348
- liczba neuronów, 272
- liczba ukrytych warstw, 271
- przekonań głębokie, DBN, 509
- rekurencyjne, RNN, 377
 - dynamiczne rozwijanie, 384
 - głęboka, 393
 - komórka GRU, 400
 - komórka LSTM, 397
 - neurony rekurencyjne, 378
 - przetwarzanie języka naturalnego, 401
 - przewidywanie szeregów czasowych, 389
 - sekwencje wejściowe o zmiennej długości, 385
 - sekwencje wejść i wyjść, 380
 - sekwencje wyjściowe o zmiennej długości, 386
 - statyczne rozwijanie, 382
 - twórcza, 392
 - uczenie, 386
- rezydualne, ResNet, 369
- rozdzielanie, 343, 344
- rozpoznawania, 408
- splotowe, CNN, 351
 - architektury, 362
 - warstwa łącząca, 360
 - warstwa splotowa, 353
- strojenie hiperparametrów, 221, 271
- sztuczne, SSN, 255, 262
- typu koder-dekoder, 403
- wielotaktowe, 396
- wielowarstwowe, MLP, 262
 - uczenie, 265
- skalowanie
 - cech, 82
 - wielowymiarowe, MDS, 226
- spadek przyśpieszony wzdłuż gradientu, 295
- specyficzność, 106
- specyfikacja klastra, 323
- splotowe sieci neuronowe, CNN, 351
- sprawdzian krzyżowy, 49, 86, 98
- SSN, artificial neural networks, 255
- statystyki uczenia, 245
- stochastyczne
 - wzmacnianie gradientowe, 204
 - zanurzanie sąsiadów przy użyciu rozkładu t , 226
- stochastyczny spadek wzdłuż gradientu, 98, 130
- stopa dyskontowa, 440
- stopnie swobody, 46
- stosy map cech, 356
- strata
 - rekonstrukcji, 409
 - rzadkości, 421
- strojenie hiperparametrów, 221, 271
- struktura danych, 63
- subrózniczka, 173
- SVM, support vector machine, 155
 - liniowa klasyfikacja, 155
 - maszyny kernelizowane, 170
 - maszyny przyrostowe, 172
 - mechanizm działania, 165
 - nieliniowa klasyfikacja, 158
 - cechy podobieństwa, 160
 - jądro RBF, 161
 - jądro wielomianowe, 159
 - złożoność obliczeniowa, 163
 - regresja, 163
- sygnał, 55
- symulowane wyżarzanie, 130
- system
 - monitorowanie, 92
 - uruchamianie, 92
 - utrzymywanie, 92
- systemy uczenia maszynowego, 27
- szacowanie prawdopodobieństwa, 146
- szeregi czasowe, 389
- sztuczka z funkcją jądra, 160
- sztuczne sieci neuronowe, SSN, 255
- szybkość uzyskania zbieżności, 129

Ś

- ścieżki algorytmów gradientu prostego, 133
- średni absolutny błąd, 58
- średnia harmoniczna, 102
- środowisko
 - CartPole, 436
 - izolowane, 60
 - symulowane, 435

T

- TensorFlow, 231
 - algorytm dynamicznego rozmieszczania, 318
 - biblioteki CUDA i cuDNN, 315
 - definiowanie sieci neuronowej, 439
 - implementacja autokoderów
 - odszumiających, 420
 - rzadkich, 422
 - stosowych, 411
 - instalacja modułu, 234
 - klaster, 323
 - metoda gradientu prostego, 238
 - modułowość, 247
 - normalizacja wsadowa, 283
 - obsługa kart graficznych, 316
 - przetwarzanie równoległe, 339
 - regresja liniowa, 237
 - rozdzielanie operacji, 313
 - sieci RSN, 381
 - standardowy interfejs, 266
 - tworzenie grafu, 234
 - uczenie głębokiej sieci neuronowej, 266
 - uczenie sieci MLP, 265
 - udostępnianie zmiennych, 250
 - warstwy splotowe, 357, 372
 - wielokrotne stosowanie modelu, 287
 - wizualizowanie
 - grafu, 243, 246
 - krzywych uczenia, 243
 - statystyk uczenia, 245
 - zakresy nazw, 246
 - zarządzanie grafami, 236
 - zrównoleglanie danych, 348
- teoria informacji Shannona, 180
- testowanie, 48
- tłumaczenie maszynowe, 403
- token EOS, 386
- transformator, 79, 81
 - DataFrameSelector, 84

- t-SNE, t-Distributed stochastic neighbor embedding, 226
- twierdzenie Mercera, 171
 - o nieistnieniu darmowych obiadów, 49
 - o zbieżności perceptronu, 261
- tworzenie grafu, 234

U

- uczenie
 - autokoderów stosowych, 413
 - dostarczanie danych algorytmowi, 241
 - drzewa decyzyjnego, 175
 - głębokiej sieci neuronowej, 13, 229, 233, 266, 275
 - faza konstrukcyjna, 266
 - faza wykonawcza, 270
 - harmonogramy, 300
 - klasyfikatora, 188
 - binarnego, 97
 - sekwencji, 387
 - maszynowe, 13
 - lista kontrolna projektu, 487
 - metodą różnic czasowych, 449
 - miksera, 206
 - modeli, 85, 119
 - modeli rzadkich, 300
 - nadzorowane, 28
 - nienadzorowane, 30
 - nienadzorowane wstępne, 417
 - offline, 34
 - pierwszej warstwy, 205
 - półnadzorowane, 33
 - przez wzmacnianie, 33, 431
 - przy użyciu reguł asocjacyjnych, 32
 - przyrostowe, 35, 98
 - rekurencyjnych sieci neuronowych, 386
 - resztowe, 369
 - rezydualne, 369
 - rozmaitościowe, 213
 - sieci MLP, 265
 - się optymalizowania nagród, 432
 - transferowe, 286
 - wsadowe, 34
 - wstępne
 - nienadzorowane, 291
 - za pomocą dodatkowego zadania, 293
 - z modelu, 38

- z przykładów, 37
- za pomocą agregacji/wklejania, 191
- zespołowe, 87, 187
- udostępnianie
 - stanu rozproszonych sesji, 327
 - zmiennych, 250
- urządzenia
 - na jednym komputerze, 314
 - na wielu serwerach, 323
- usługi
 - nadrzędne, 325
 - robocze, 325
- uśrednione kodowanie μ , 423
- użycie wielokrotne
 - modelu TensorFlow, 287
 - warstw, 286

W

- wariancja, 65, 138, 215
- warstwa, 69
 - górne
 - modyfikowanie, 291
 - usuwanie, 291
 - zastępowanie, 291
 - łączące, 360
 - w głębi, 366
 - neuronów rekurencyjnych, 378
 - niższe
 - zamrażanie, 289
 - ograniczające, 366
 - spłotowe, 353, 372
 - filtry, 354
 - stosy map cech, 356
 - zużycie pamięci operacyjnej, 359
 - ukryte, 262
 - zapamiętywanie, 290
 - wejściowe, 262
 - wielokrotne użycie, 286
 - wyjściowe, 262
- wartość stanu-czynności, 448
- wczesne zatrzymywanie, 144, 302
- wczytywanie
 - danych, 333
 - bezpośrednio z grafu, 334
 - do zmiennej, 333
 - modeli, 242
- wektor
 - momentu, 294
 - nośny, 156, 493

- wiązanie wag, 412
- wielowarstwowy zespół kontaminacji, 207
- wizualizowanie
 - cech, 416
 - danych, 71
 - danych geograficznych, 71
 - drzewa decyzyjnego, 175
 - grafu, 243, 246
 - rekonstrukcji, 415
 - statystyk uczenia, 245
 - t-SNE, 31
- wklejanie, 190, 191
- wnioskowanie, 42
- wsadowy gradient prosty, 127
- wskaźnik Giniego, 177, 180
- współczynnik
 - korelacji liniowej, 73
 - porzucenia, 304
 - uczenia, 36, 124, 221, 300
 - wariancji wyjaśnionej, 217
- wybór
 - jądra, 221
 - liczby wymiarów, 218
 - modelu, 85
 - podprzestrzeni rzutowania, 215
- wydobywanie
 - cech, 32
 - danych, 27
- wykres
 - funkcji logistycznej, 146
 - krzywej ROC, 106
 - mediany, 76
 - precyzji i pełności, 104
 - rozproszenia, 75
- wykrywanie anomalii, 32
- wyliczanie prognoz, 176
- wyszukiwanie polityki, 433, 434
- wzmocnianie, 187, 197
 - adaptacyjne, 197
 - gradientowe, 197, 200
 - hipotezy, 197

Z

- zachowanie wariancji, 215
- zadania, 323
- zakresy nazw, 246
- zależności sterujące, 322
- założenia, 49, 59
- założenie różnorodności, 213

zamrażanie niższych warstw, 289
zapamiętywanie warstw ukrytych, 290
zapisywanie modeli, 242
zarządzanie grafami, 236
zasada DRY, 248
zaszumienie próbkowania, 44
zawiasowa funkcja straty, 157, 173
zbieżność algorytmu, 124
zbiór
 danych
 liniowy, 122
 MNIST, 95
 Swiss roll, 212, 224
 wypaczony, 99
 testowy, 48, 67, 92
 uczący, 48
 walidacyjny, 49
zespół, 187
 Extra-Trees, 195
 GBRT, 201
złożoność obliczeniowa, 123, 163, 180
zmienna swobodna, 167
znormalizowana funkcja wykładnicza, 150
zrównoleglenie
 danych, 345
 modelu, 343

PROGRAM PARTNERSKI

GRUPY WYDAWNICZEJ HELION

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

Zmień swoją stronę WWW
w działający bankomat!

Dowiedz się więcej i dołącz już dzisiaj!

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

GRUPA WYDAWNICZA

 **Helion SA**

Już dziś zacznij tworzyć inteligentne systemy!

W ciągu ostatnich lat uczenie maszynowe stało się sercem wielu nowoczesnych produktów, takich jak zaawansowane techniki wyszukiwania w przeglądarkach, rozpoznawanie mowy w smartfonach czy proponowanie treści w zależności od indywidualnych preferencji użytkownika. Być może niedługo taki inteligentny system zastąpi nas za kierownicą samochodu. Uczenie głębokie wprowadziło nową jakość do uczenia maszynowego. Daje niesamowite możliwości, jednak wymaga olbrzymiej mocy obliczeniowej i potężnych ilości danych. Programiści implementujący takie rozwiązania są poszukiwanymi specjalistami i mogą liczyć na ekscytujące oferty!

Ta książka jest praktycznym podręcznikiem tworzenia systemów inteligentnych. Przedstawiono tu najważniejsze zagadnienia teoretyczne dotyczące uczenia maszynowego i sieci neuronowych. W zrozumiały sposób zaprezentowano koncepcje i narzędzia służące do tworzenia systemów inteligentnych. Opisano Scikit-Learn i TensorFlow – środowiska produkcyjne języka Python – i pokazano krok po kroku, w jaki sposób wykorzystuje się je do implementacji sieci neuronowych. Liczne praktyczne przykłady i ćwiczenia pozwolą na pogłębienie i utrwalenie zdobytej wiedzy. Jeśli tylko potrafisz posługiwać się Pythonem, dzięki tej przystępnie napisanej książce szybko zaczniesz implementować systemy inteligentne.

Aurélien Géron – specjalizuje się w dziedzinie uczenia maszynowego. Posiada również doświadczenie w branży finansowej, obronnej i medycznej. Był wykładowcą akademickim, pracował w Google, jest założycielem firmy Wifirst (wiodącego dostawcy bezprzewodowych usług internetowych we Francji). Napisał kilka książek o języku C++, sieciach Wi-Fi oraz architekturze sieci.

W tej książce między innymi:

- podstawowe koncepcje uczenia maszynowego, uczenia głębokiego i sieci neuronowych
- przygotowywanie zbiorów danych i zarządzanie nimi
- algorytmy uczenia maszynowego
- rodzaje architektury sieci neuronowych
- uczenie głębokich sieci neuronowych

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

Sprawdź nasze szkolenia!



AKADEMIA IT & BUSINESS
WWW.SZKOLENIA.HELION.PL

KOD KORZYŚCI
Śięgnij po więcej!



ISBN 978-83-283-4373-3



9 788328 343733