



Stuart
Russell

Peter
Norvig

Sztuczna inteligencja

Nowe spojrzenie

Wydanie IV. Tom 2



Tytuł oryginału: Artificial Intelligence: A Modern Approach, 4th Edition

Tłumaczenie: Andrzej Grażyński

ISBN: 978-83-283-7773-8

Authorized translation from the English language edition, entitled ARTIFICIAL INTELLIGENCE: A MODERN APPROACH, 4th Edition by STUART RUSSELL; PETER NORVIG, published by Pearson Education, Inc, publishing as Pearson, Copyright © 2021, 2010, 2003 by Pearson Education, Inc. or its affiliates, 221 River Street, Hoboken, NJ 07030.

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 Pearson Education, Inc.

Polish language edition published by Helion S.A., Copyright © 2022.

PEARSON, ALWAYS LEARNING is an exclusive trademark owned by Pearson Education, Inc. or its affiliates in the U.S. and/or other countries.

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

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

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

Helion S.A.

ul. Kościuszki 1c, 44-100 Gliwice

tel. 32 231 22 19, 32 230 98 63

e-mail: helion@helion.pl

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

Drogi Czytelniku!

Jeżeli chcesz ocenić tę książkę, zajrzyj pod adres

<https://helion.pl/user/opinie/szti42>

Możesz tam wpisać swoje uwagi, spostrzeżenia, recenzję.

Printed in Poland.

- Kup książkę
- Poleć książkę
- Oceń książkę

- Księgarnia internetowa
- Lubię to! » Nasza społeczność

SPIS TREŚCI

V **Uczenie maszynowe**

19	Uczenie maszynowe z przykładowych danych	7
19.1.	Formy uczenia maszynowego	8
19.2.	Uczenie nadzorowane	10
19.3.	Drzewa decyzyjne w uczeniu maszynowym	14
19.4.	Selekcja modelu i optymalizacja	23
19.5.	Teoria uczenia maszynowego	31
19.6.	Regresja liniowa i klasyfikacja	35
19.7.	Modele nieparametryczne	48
19.8.	Uczenie zespołowe	59
19.9.	Budowanie systemów uczenia maszynowego	68
	Podsumowanie	79
	Bibliografia i uwagi historyczne	80
20	Uczenie modeli probabilistycznych	86
20.1.	Uczenie statystyczne	86
20.2.	Uczenie z kompletnych danych	90
20.3.	Uczenie z ukrytymi zmiennymi: algorytm EM	104
	Podsumowanie	114
	Bibliografia i uwagi historyczne	114
21	Głębokie uczenie	117
21.1.	Proste sieci ze sprzężeniem w przód	118
21.2.	Grafy obliczeniowe dla głębokiego uczenia	124
21.3.	Konwolucyjne sieci neuronowe (CNN)	128
21.4.	Algorytmy głębokiego uczenia	134
21.5.	Generalizacja	137
21.6.	Rekurencyjne sieci neuronowe (RNN)	142
21.7.	Nienadzorowane uczenie transferowe	145
21.8.	Zastosowania	152
	Podsumowanie	154
	Bibliografia i uwagi historyczne	155

22	Uczenie ze wzmacnianiem	159
22.1.	Uczenie się dla nagród	159
22.2.	Pasywne uczenie ze wzmacnianiem	161
22.3.	Aktywne uczenie ze wzmacnianiem	167
22.4.	Generalizacja w uczeniu ze wzmacnianiem	174
22.5.	Wyszukiwanie polityki	181
22.6.	Uczenie praktykanckie i odwrotne uczenie ze wzmacnianiem	183
22.7.	Zastosowania uczenia ze wzmacnianiem	186
	Podsumowanie	189
	Bibliografia i uwagi historyczne	191
VI	Komunikacja, percepcja i działanie	
23	Przetwarzanie języka naturalnego	195
23.1.	Modele językowe	196
23.2.	Gramatyka	208
23.3.	Parsowanie	210
23.4.	Gramatyki augmentowane	216
23.5.	Komplikacje języków naturalnych	221
23.6.	Zadania NLP	225
	Podsumowanie	227
	Bibliografia i uwagi historyczne	228
24	Głębokie uczenie w przetwarzaniu języka naturalnego	233
24.1.	Embeddingi słów	233
24.2.	Rekurencyjne sieci neuronowe w NLP	238
24.3.	Modele „sekwencja na sekwencję”	242
24.4.	Architektura transformerów	246
24.5.	Trenowanie wstępne i uczenie transferowe	249
24.6.	Obecny stan sztuki	253
	Podsumowanie	257
	Bibliografia i uwagi historyczne	257
25	Widzenie komputerowe	260
25.1.	Wstęp	260
25.2.	Formowanie obrazów	261
25.3.	Podstawowe cechy obrazów	268
25.4.	Klasyfikowanie obrazów	275
25.5.	Wykrywanie obiektów	279

25.6.	Rzeczywistość 3D	282
25.7.	Widzenie komputerowe w praktyce	287
	Podsumowanie	302
	Bibliografia i uwagi historyczne	302
26	Robotyka	308
26.1.	Wstęp	308
26.2.	Sprzęt robotów	309
26.3.	Jakie rodzaje problemów rozwiązywać może robotyka?	314
26.4.	Percepcja robotów	315
26.5.	Planowanie i sterowanie	323
26.6.	Planowanie ruchu w warunkach niepewności	343
26.7.	Uczenie ze wzmacnianiem w robotyce	346
26.8.	Ludzie i roboty	348
26.9.	Alternatywne frameworki robotyczne	356
26.10.	Domeny zastosowań robotyki	358
	Podsumowanie	362
	Bibliografia i uwagi historyczne	364
VII	Konkluzje	
27	Bezpieczeństwo oraz etyczne i filozoficzne aspekty sztucznej inteligencji	369
27.1.	Granice sztucznej inteligencji	369
27.2.	Czy maszyny mogą naprawdę myśleć?	373
27.3.	Sztuczna inteligencja a etyka	375
	Podsumowanie	397
	Bibliografia i uwagi historyczne	397
28	Przyszłość sztucznej inteligencji	404
28.1.	Komponenty sztucznej inteligencji	404
28.2.	Architektury sztucznej inteligencji	411
	Bibliografia	417
	Skorowidz	459

GŁĘBOKIE UCZENIE W PRZETWARZANIU JĘZYKA NATURALNEGO

Głębokie sieci neuronowe, które zrewolucjonizowały uczenie maszynowe, nie mogły rzecz jasna ominąć kluczowej gałęzi tego uczenia, jaką jest przetwarzanie języka naturalnego. W tym rozdziale pokazujemy, jak przyczyniają się one do uchwycenia struktury języka naturalnego, zapewnienia jego płynności i wydajnego usprawnienia wielu innych zadań.

W rozdziale 23. omawialiśmy kluczowe elementy języka naturalnego, między innymi gramatykę i semantykę. Systemy oparte na parsingu i analizie semantycznej okazały się znakomite w realizacji wielu zadań, lecz ich wydajność jest ograniczona ze względu na stopień złożoności języków naturalnych i zachodzące w rzeczywistym świecie zjawiska i zdarzenia wpływające na ich ewolucję. Dysponując ogromną ilością dostępnego tekstu w formie czytelnej dla maszyn, moglibyśmy wysunąć hipotezę, że podejścia oparte na uczeniu sterowanym danymi mogą być bardziej efektywne; tę właśnie hipotezę przeanalizujemy dokładniej w tym rozdziale, posługując się narzędziami wykorzystywanymi przez systemy głębokiego uczenia, omawiane w rozdziale 21.

Rozpocniemy od pokazania (w podrozdziale 24.1), jak można usprawnić uczenie maszynowe, traktując słowa języka jako punkty w wysokowymiarowej przestrzeni, a nie jako niepodzielne (atomowe) wartości. Przedmiotem podrozdziału 24.2 jest zastosowanie rekurencyjnych sieci neuronowych do rozpoznawania znaczenia i długodystansowego kontekstu w miarę sekwencyjnego przetwarzania tekstu. Treść podrozdziału 24.3 koncentruje się na tłumaczeniach maszynowych, jako jednym z najbardziej pomyślnych obszarów zastosowań głębokiego uczenia w NLP. Podrozdziały 24.4 i 24.5 poświęcone są modelom, które można trenować za pomocą dużych porcji nieetykietowanego tekstu, a następnie wykorzystywać do realizacji specyficznych zadań, ze zdumiewającym nieraz powodzeniem. Na zakończenie, w podrozdziale 24.6, przedstawiamy obecny stan wiedzy w zakresie zastosowań głębokiego uczenia w NLP i przewidywane kierunki rozwoju w tej dziedzinie.

24.1. Embeddingi słów

Potrzebujemy takiej reprezentacji słów, która nie wymaga ręcznej inżynierii cech, ale umożliwia generalizację powiązań między słowami — powiązań syntaktycznych („bezbarwny” i „idealny” to przymiotniki), semantycznych („kot” i „kotek” to przedstawiciele rodziny kotowatych), tematycznych („słonecznie” i „gołoledź” to zjawiska pogodowe), wyrażających sentyment („fantastyczny” to sentyment przeciwny do „żenujący”) itp.

Jak więc mamy zakodować słowo w formie wektora \mathbf{x} jako wejściowego dla sieci neuronowej? Mógłby to być **wektor z gorącą jedyńką** (ang. *one-hot vector*), czyli taki, w którym i -ty element równy jest 1, a pozostałe są zerowe, gdzie i jest numerem pozycji słowa w słowniku (patrz sekcja 21.2.1). Taka jednak reprezentacja w żaden sposób nie odzwierciedla podobieństw między słowami.

Wzorując się na maksymie językoznawcy Johna R. Firtha (1957) „Poznasz słowo po jego towarzystwie”¹ możemy reprezentować dane słowo w formie wektora, którego k -ty element jest liczbą całkowitą odzwierciedlającą jego występowanie w n -gramach. Numerujemy wszystkie możliwe n -gramy (dla ustalonego n), a w wektorze reprezentującym dane słowo k -ty element równy jest liczbie wystąpień tego słowa w k -tym n -gramie. Operowanie „surowymi” licznikami byłoby jednak nieporęczne — dla 100 000 słów liczba możliwych 5-gramów wynosi $100\,000^5 = 10^{25}$, choć znakomita większość wektorów w tej 10^{25} -wymiarowej przestrzeni to wektory rzadkie, w których większość elementów jest zerowa. Bardziej praktyczną generalizację otrzymamy, posługując się wektorami, których liczba elementów jest znacznie mniejsza, zredukowana do (powiedzmy) kilkuset. Taki zredukowany wektor, reprezentujący wystąpienie danego słowa w określonym kontekście, nazywamy jest (z angielska) **embeddingiem**² tego słowa. Embeddingi poszczególnych słów budowane są automatycznie w procesie uczenia, na podstawie danych (w jaki sposób — pokażemy to w dalszej części rozdziału). Fizycznie są one *wektorami liczbowymi* o ustalonym rozmiarze, na przykład

„aardvark” = $[-0,7, +0,2, -3,2, \dots]$
 „abacus” = $[+0,5, +0,9, -1,3, \dots]$
 „zyzzyzus”³ = $[-0,1, +0,8, -0,4, \dots]$

Co do konkretnych wartości elementów tych wektorów, to nie mają one samoistnego znaczenia, ale wynikają z wymogu, by podobne słowa transformowały się na podobne wektory. Zasadę tę ilustruje rysunek 24.1, na którym widać, jak embeddingi słów pokrewnych grupują się w odrębne klaster, zawierające nazwy państw, relacji rodzinnych, środków transportu i żywności.

Jak się okazuje, z powodów, których do końca nie rozumiemy, embeddingi słów wykazują jeszcze inne właściwości oprócz bliskości do podobnych słów. Załóżmy na przykład, że poszukujemy embeddingu \mathbf{A} , odzwierciedlającego słowo „Ateny”, i embeddingu \mathbf{B} , odzwierciedlającego słowo „Grecja”. Intuicyjnie można podejrzewać, że wektor stanowiący różnicę embeddingów $\mathbf{B} - \mathbf{A}$ będzie *zakodowaną postacią związku* „państwo-stolica” między słowami. A więc innym parom słów pozostających w takim związku — „Francja-Paryż”, „Rosja-Moskwa”, „Zambia-Lusaka” — powinny odpowiadać (zasadniczo) takie same wektory różnicowe.

Możemy wykorzystywać tę własność do rozwiązywania problemów polegających na analogiach związków między słowami, na przykład „»Ateny« i »Grecja« pozostają ze sobą w takim samym związku, jak »Oslo« i [co]?”. Niech \mathbf{C} będzie embeddingiem słowa „Oslo”, a \mathbf{D} niech oznacza embedding szukanego słowa; jako że domniemyamy równość

$$\mathbf{B} - \mathbf{A} = \mathbf{D} - \mathbf{C}$$

więc embedding szukanego słowa określony będzie przez relację

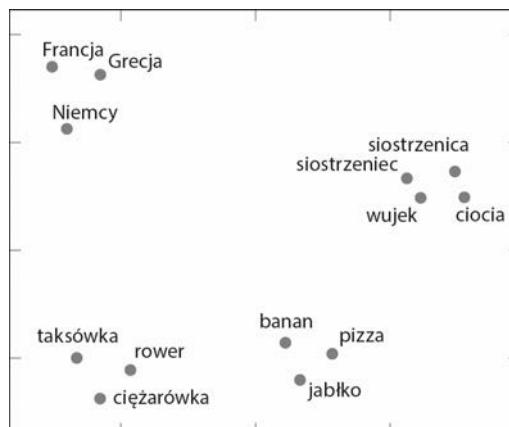
$$\mathbf{D} = \mathbf{C} + (\mathbf{B} - \mathbf{A})$$

I gdy już obliczymy embedding \mathbf{D} , zauważymy, że najbliższym z odpowiadających mu słów jest „Norwegia”. W tabeli 24.1 pokazujemy przykłady funkcjonowania tej arytmetyki wektorowej dla kilku innych relacji.

¹ „You shall know a word by the company it keeps” — *przyp. tłum.*

² W języku polskim termin ten jest często tłumaczony zarówno jako czynność „osadzania” słowa, przez analogię do „łączenia i osadzania obiektów” (ang. OLE — *Object Linking and Embedding*), jak i jako „osadzenie” będące rezultatem tej czynności. Tymczasem w znaczeniu OLE „osadzanie” i „osadzenie” odnosi się w istocie do *integrowania obiektów* (na przykład włączania obrazka do dokumentu tekstowego), natomiast *embedding* słowa jest wektorem stanowiącym wynik pewnej jego *transformacji* (z postaci znakowej na wektor liczbowy), niekojarzącej się z jakimś integrowaniem. Ponadto z perspektywy języka kontrowersyjne wydaje się użycie rzeczownika *odczasownikowego* na określenie obiektu, jakim jest wspomniany wektor. W związku z tym pozostawiam występujące w oryginale słowo „embedding” na oznaczenie owego wektora, zaś samą czynność jego tworzenia nazywał będę „embedowaniem” — *przyp. tłum.*

³ Patrz https://hrviki.net/Wikipedia_Zyzzyzus — *przyp. tłum.*



RYSUNEK 24.1. Embeddingi słów wyliczone przez algorytm GloVe wytrenowany na zbiorze liczącym 6 miliardów słów. Na rysunku widoczny jest efekt rzutowania 100-wymiarowych wektorów słów na dwuwymiarową płaszczyznę; wektory reprezentujące podobne słowa znajdują się blisko siebie

TABELA 24.1. W niektórych przypadkach model embeddingów może dostarczać odpowiedzi na pytania w rodzaju „A ma się tak do B, jak C do [czego]?” za pomocą prostej arytmetyki wektorowej: dysponując embeddingami słów A, B i C, możemy wyliczyć nieznaną embedding D z zależności $D = C + (B - A)$ i następnie poszukiwać słowa, którego embedding jest najbardziej zbliżony do D (odpowiedź widoczna w kolumnie D wyliczona została automatycznie przez model, opisy w kolumnie *Relacja* zostały dodane ręcznie). Zaadaptowano z: Mikolov i in., (2013, 2014)

A	B	C	$D = C + (B - A)$	Relacja
Ateny	Grecja	Oslo	Norwegia	<i>Stolica</i>
Astana	Kazachstan	Harare	Zimbabwe	<i>Stolica</i>
Angola	kwanza	Iran	rial	<i>Waluta</i>
Miedź	Cu	Złoto	Au	<i>Symbol pierwiastka</i>
Microsoft	Windows	Google	Android	<i>System operacyjny</i>
Nowy Jork	New York Times	Baltimore	Baltimore Sun	<i>Gazeta</i>
Berlusconi	Silvio	Obama	Barack	<i>Imię</i>
Szwajcaria	Szwajcar	Kambodża	Kambodżanin	<i>Narodowość</i>
Einstein	naukowiec	Picasso	malarz	<i>Profesja</i>
brat	siostra	wnuczek	wnuczka	<i>Pokrewieństwo</i>
Chicago	Illinois	Stockton	Kalifornia	<i>Stan</i>
możliwe	niemożliwe	etyczne	nieetyczne	<i>Negacja</i>
mysz	myszy	dolar	dolary	<i>Liczba mnoga</i>
łatwe	najłatwiejsze	szczęśliwy	najszcześniejszy	<i>Stopień najwyższy</i>
idę	szedłem	pływam	pływałem	<i>Czas przeszły</i>

Nie ma jednak gwarancji, że konkretny algorytm konstruowania embeddingów, uruchomiony dla konkretnego korpusu, prawidłowo uchwyci konkretne relacje semantyczne. Technika embeddingów stała się popularna bynajmniej nie dlatego, by skrywała w sobie jakąś magiczną zdolność bezbłędnego wychwytywania analogii relacyjnych między słowami korpusu, ale po prostu dlatego, że udowodniła swą wartość jako znakomita reprezentacja słów dla zadań NLP typu *downstream*, czyli operowania na gotowych tekstach: odpowiadania na pytania, tłumaczenia czy generowania streszczeń.

Wykorzystywanie embeddingów zamiast reprezentacji „z gorącą jedyneką” okazuje się wielce pomocne w odniesieniu do niemal wszelkich zastosowań głębokiego uczenia w zadaniach NLP. W przypadku bowiem wielu z tych zadań można posługiwać się generycznymi **wstępnie wytrenowanymi** zbiorami embeddingów, oferowanymi przez wielu dostawców na potrzeby konkretnych zastosowań i pozwalającymi zaoszczędzić sporo czasu. Gdy piszemy te słowa, do najbardziej popularnych słowników tego typu należą WORD2VEC, GloVe (od *Global Vectors*) i FASTTEXT, oferujące embeddingi słów w 157 językach. Powrócimy do tej kwestii w sekcji 24.5.1.

Oczywiście możliwe jest samodzielne trenowanie embeddingów — może się ono odbywać niejako przy okazji trenowania sieci neuronowej pod kątem określonych zadań, dla których przydatność wspomnianych generycznych embeddingów jest ograniczona czy wręcz żadna. Weźmy na przykład tagowanie części mowy, opisywane w sekcji 23.1.6, polegające na przewidywaniu przynależności poszczególnych słów zdania do określonych kategorii gramatycznych. Mimo iż jest to zadanie stosunkowo proste, to bynajmniej nie można uważać go za banalne, ponieważ dla wielu słów nie istnieje jedna „prawidłowa” kategoria — popularne słowo *cut* może być czasownikiem w czasie teraźniejszym — przechodnim (*I cut my hair* — ostrzygłem się) lub nieprzechodnym (*He cut a lecture* — opuścił wykład), czasownikiem w czasie przeszłym prostym (*he cut off smoking* — on rzucił palenie), bezokolicznikiem (ciąć, ścinać, odcinać, itp.), imiesłowem czasu przeszłego (*I was cut off while talking* — przerwano mi połączenie), przymiotnikiem (ścięty, rozcięty) lub rzeczownikiem (*price cut* — obniżka ceny). Jeśli pobliski przysówek czasu odnosi się do przeszłości, sugeruje to, że to konkretne wystąpienie *cut* jest czasownikiem w czasie przeszłym; a zatem możemy mieć nadzieję, że embedding uchwyci aspekt przysłówkowy, odnoszący się do przeszłości.

Tagowanie części mowy to bardzo wdzięczny temat wprowadzający w zastosowania głębokiego uczenia w NLP, jest bowiem wolny od komplikacji charakterystycznej dla problemów bardziej zaawansowanych, na przykład odpowiadania na pytania, którym to problemem zajmujemy się w sekcji 24.5.3. Dla danego korpusu zdań etykietowanych tagami części mowy, możemy realizować równoczesne uczenie dwóch rzeczy: parametrów dla embeddingów słów i rozpoznawania części mowy. Proces ten przebiega w następujących krokach:

1. Wybieramy nieparzystą liczbę w jako szerokość tzw. okna przewidywania, służącego do tagowania słów; zwykle przyjmuje się $w = 5$, co oznacza, że tag danego słowa przewidywany jest na podstawie tegoż słowa oraz dwóch słów sąsiadujących z nim z lewej i dwóch z prawej strony. Dzielimy każde zdanie korpusu na nakładające się okna o szerokości w ; każde z tych okien traktowane jest jako przykład treningowy w postaci ciągu w słów, etykietowany przewidywaną kategorią POS środkowego słowa.
2. Tworzymy słownik wszystkich unikalnych słów-tokenów występujących w danych treningowych więcej niż (powiedzmy) 5 razy; oznaczmy przez v liczbę wszystkich słów w tak utworzonym słowniku.
3. Sortujemy zawartość tego słownika w dowolnej kolejności (na przykład alfabetycznej).
4. Wybieramy d jako rozmiar embeddingu dla każdego słowa.
5. Tworzymy macierz wagową \mathbf{E} o rozmiarze v wierszy $\times d$ kolumn; i -ty wiersz tej macierzy stanie się (po zakończeniu uczenia) embeddingiem i -tego słowa w słowniku. Inicjujemy tę macierz losowymi wartościami (albo, jeśli dysponujemy wstępnie wytrenowanymi embeddingami, zapisujemy te embeddingi jako wiersze macierzy).

6. Budujemy sieć neuronową, której wyjściem będzie przewidywana etykieta części mowy, jak pokazaliśmy na rysunku 24.2. Pierwsza warstwa tej sieci składa się z w kopii wspomnianej macierzy kontekstów. W sieci tej występują ponadto dwie warstwy ukryte \mathbf{z}_1 i \mathbf{z}_2 , z macierzami wagowymi (odpowiednio) \mathbf{W}_1 i \mathbf{W}_2 , po których następuje warstwa softmax zwracająca rozkład prawdopodobieństwa $\hat{\mathbf{y}}$ możliwych kategorii wyjściowych:

$$\begin{aligned} \mathbf{z}_1 &= \sigma(\mathbf{W}_1 \mathbf{x}) \\ \mathbf{z}_2 &= \sigma(\mathbf{W}_2 \mathbf{z}_1) \\ \hat{\mathbf{y}} &= \text{softmax}(\mathbf{W}_{out} \mathbf{z}_2) \end{aligned}$$



RYSUNEK 24.2. Model ze sprzężeniem w przód dla tagowania części mowy. Na podstawie składającego się z 5 słów okna jako wejścia model ten przewiduje kategorię gramatyczną (część mowy) dla środkowego słowa — w tym przypadku *cut*. Model zdolny jest do uwzględniania pozycji słowa w zdaniu, ponieważ każdy z 5 embeddingów mnożony jest przez różne części pierwszej ukrytej warstwy. Wartości parametrów dla embeddingów i dla trzech warstw są wynikiem zrównoległego treningu (*Yesterday they cut the rope* — „wczoraj [oni] przecięli linę”)

7. Aby zakodować sekwencję w słów w postaci wektora wejściowego, po prostu wyszukujemy embeddingi dla tych słów i konkatenujemy je, tworząc wektor \mathbf{x} złożony z $w \cdot d$ liczb rzeczywistych. Mimo iż embedding danego słowa nie zależy od pozycji, jaką słowo to zajmuje w sekwencji w słów, to od pozycji tej zależy mnożnik w warstwie \mathbf{z}_1 , przez który mnożony jest wspomniany embedding; w ten oto sposób w wyniku kodowania słów uwzględniana jest ich relatywna pozycja w sekwencji.
8. Trenujemy macierz \mathbf{E} oraz macierze wagowe \mathbf{W}_1 , \mathbf{W}_2 i \mathbf{W}_{out} , wykorzystując metodę spadku gradientowego. Gdy wszystko przebiegnie pomyślnie, środkowe słowo *cut* zaetykietowane zostanie jako czasownik w czasie przeszłym (*past tense verb*) — za sprawą innych słów w oknie przewidywania, między innymi określenia temporalnego *yesterday* („wczoraj”) oraz zaimka *they* („oni”, ew. „one”) w 3. osobie liczby mnogiej, występującego bezpośrednio przed *cut*.

Alternatywą dla modelu embeddingów może być **model znakowy**, w ramach którego informacją wejściową jest sekwencja znaków, z których każdy kodowany jest jako wektor „z gorącą jedynką”. Model taki umożliwia uczenie się sposobów formowania znaków w słowa. Zdecydowanie jednak większość badań z zakresu NLP realizowana jest raczej na poziomie słów niż na poziomie oddzielnych znaków.

24.2. Rekurencyjne sieci neuronowe w NLP

Dysponujemy już dobrymi reprezentacjami pojedynczych, izolowanych słów, ale przecież język to uporządkowane sekwencje słów, dla których istotny jest **kontekst** wynikający ze słów w bliższym lub dalszym otoczeniu. Dla prostych zadań NLP, takich jak tagowanie części mowy, wystarczający okazuje się zwykle kontekst wyznaczany przez okna przewidywań o ustalonym, niewielkim rozmiarze (jak pięcioelementowe okna z poprzedniej sekcji).

Zadania bardziej ambitne, jak odpowiadanie na pytania czy rozwiązywanie odwołań (ang. *reference resolution*), mogą jednak wymagać uwzględnienia znacznie szerszego kontekstu. Na przykład w zdaniu *Eduardo told me that Miguel was very sick so I took him to the hospital* („Eduardo powiedział mi, że Miguel jest bardzo chory, więc zabrałam go do szpitala”) rozpoznanie, że **him** odnosi się do *Miguel*, a nie do *Eduardo*, wymaga kontekstu rozciągającego się na wszystkie słowa 15-wyrazowego zdania.

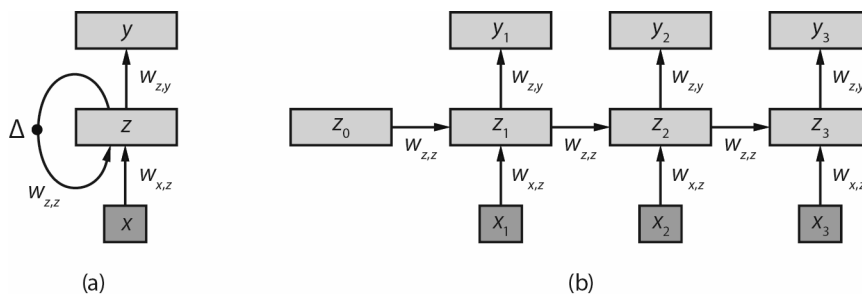
24.2.1. Modele językowe wykorzystujące sieci RNN

Rozpocznijmy od problemu stworzenia **modelu językowego** zapewniającego wystarczający kontekst. Jak pamiętamy, model językowy to w istocie rozkład prawdopodobieństwa względem sekwencji słów, pozwalający na przewidywanie kolejnego słowa w tekście na podstawie wszystkich poprzednich słów. Modele takie stanowią integralne elementy rozwiązań wielu złożonych problemów NLP.

Budowanie modelu językowego w oparciu zarówno o n -gramy (patrz podrozdział 23.1), jak i o sieci ze sprzężeniem w przód (patrz podrozdział 21.1) z oknem przewidywania o ustalonej szerokości wiąże się z tą trudnością, że albo szerokość ta jest zbyt mała dla wymaganego kontekstu, albo liczba wymaganych parametrów modelu jest nadmiernie duża (jedno zresztą nie wyklucza drugiego).

Ponadto sieci ze sprzężeniem w przód cechują się swoistą **asymetrią**: czegokolwiek nauczy się sieć odnośnie (powiedzmy) słowa *him* występującego na 12. pozycji w zdaniu, okaże się to raczej nieprzydatne w odniesieniu do takiego samego słowa występującego w innym miejscu tegoż zdania; powodem tego jest oczywiście różnica wag na różnych pozycjach.

W podrozdziale 21.6 omawialiśmy podstawy **rekurencyjnych sieci neuronowych** (RNN), zaprojektowanych z myślą o przetwarzaniu szeregów czasowych, po jednym elemencie szeregu w danej chwili. Natychmiast nasuwa się analogia z przetwarzaniem tekstu na poziomie kolejnych słów — analogia sugerująca przydatność sieci RNN do tego celu. Rysunek 24.3, przedstawiający schemat podstawowej sieci RNN, jest repliką rysunku 21.8.



RYСУNEK 24.3. (a) Schematyczny diagram podstawowej sieci RNN, w której ukryta warstwa z zawiera rekurencyjne połączenia; symbol Δ reprezentuje opóźnienie. Każdy wektor wejściowy x jest embeddingiem następnego słowa w zdaniu. (b) Ta sama sieć rozwinięta na trzy kroki czasowe, do postaci sieci ze sprzężeniem w przód. Zauważmy, że wagi współdzielone są przez wszystkie kroki czasowe

W modelu językowym opartym na sieci RNN każde słowo wejściowe zakodowane zostaje w postaci embeddingu \mathbf{x}_t . Ukryta warstwa \mathbf{z}_t tej sieci przekazuje wejście z jednego kroku czasowego do następnego. Jesteśmy zainteresowani klasyfikacją wieloklasową, rolę klas pełnią słowa zawarte w słowniku. Wyjście \mathbf{y}_t jest więc rozkładem prawdopodobieństwa softmax względem możliwych postaci następnego słowa w zdaniu.

Architektura RNN zapewnia rozwiązanie problemu zbyt dużej liczby parametrów: ich liczba w macierzach wagowych $w_{*,z,z}$, $w_{*,x,z}$ i $w_{*,z,y}$ pozostaje stała, niezależnie od liczby słów — czyli jest rzędu $O(1)$. Kontrastuje to z sieciami ze sprzężeniem w przód, wymagającymi $O(n)$ parametrów, oraz modelami n -gramowymi, w których liczba wymaganych parametrów jest rzędu $O(v^n)$, gdzie v oznacza rozmiar słownika.

Architektura RNN rozwiązuje także problem asymetrii, ponieważ wagi są takie same dla wszystkich pozycji słów w zdaniu.

W niektórych przypadkach architektura RNN może także rozwiązywać problem ograniczenia kontekstu. Teoretycznie nie istnieje ograniczenie na długość dystansu, na jaki model może spoglądać wstecz, na wcześniejsze słowa wejściowe. Każda aktualizacja warstwy wejściowej \mathbf{z}_t ma dostęp zarówno do bieżącego słowa wejściowego \mathbf{x}_t , jak i poprzedniej warstwy ukrytej \mathbf{z}_{t-1} , co oznacza, że informacja na temat dowolnego słowa wejściowego może być przechowywana w warstwie ukrytej w nieskończoność, oraz kopiowana (być może wraz z odpowiednią modyfikacją) między danym krokiem czasowym i krokiem następnym. W praktyce oczywiście owa „nieskończoność” ograniczona jest rozmiarem dostępnej pamięci dla warstwy \mathbf{z} , nie jest więc możliwe przechowywanie wszystkiego na temat wszystkich słów.

Modele RNN sprawdzają się bardzo dobrze w rozwiązywaniu wielu problemów, ale nie wszystkich — i trudne może być przewidywanie, czy rozwiązywanie konkretnego problemu zakończy się powodzeniem. Jednym z czynników zwiększających szansę takiego powodzenia jest fakt, że proces treningowy zachęca sieć do alokowania pamięci w warstwie \mathbf{z} przede wszystkim na potrzeby tych aspektów wejściowych, których faktyczna użyteczność została już wykazana.

Proces trenowania modelu RNN opisywaliśmy w sekcji 21.6.1. Wektory wejściowe \mathbf{x}_t reprezentują słowa odczytywane z korpusu treningowego, a obserwowanymi wyjściami są te same słowa przesunięte o 1. Na przykład dla tekstu treningowego „witaj świecie” pierwszym wektorem wejściowym \mathbf{x}_1 będzie embedding słowa „witaj”, a pierwszym wyjściem \mathbf{y}_1 będzie embedding słowa „świecie”. Trenujemy nasz model w celu przewidywania następnego słowa i spodziewamy się, że będzie on ten cel realizował wykorzystując ukrytą warstwę do reprezentowania użytecznych informacji. Jak wyjaśniliśmy w sekcji 21.6.1, obliczamy różnicę między wyjściem obserwowanym a wyjściem prognozowanym przez sieć i propagujemy tę różnicę wstecz, utrzymując te same wagi dla wszystkich kroków czasowych.

Gdy proces treningowy zostanie zakończony, możemy generować losowe teksty za pomocą wytrenowanego modelu. Przekazujemy temu modelowi początkowy wektor wejściowy \mathbf{x}_1 i obserwujemy wyprodukowane wyjście \mathbf{y}_1 , które jest rozkładem softmax względem możliwych słów. Próbkujemy pojedyncze słowo z tego rozkładu, rejestrujemy to słowo jako wyjście dla kroku czasowego t i dostarczamy je na wejście sieci jako słowo wejściowe \mathbf{x}_2 . Powtarzamy ten proces dowolną liczbę razy. Próbkując z rozkładu \mathbf{y}_1 możemy dokonać wyboru strategii próbkowania: możemy mianowicie konsekwentnie wybierać słowo najbardziej prawdopodobne, możemy faworyzować poszczególne słowa proporcjonalnie do ich prawdopodobieństwa, albo też „przepróbkowywać” słowa mniej prawdopodobne w celu większego urozmaicenia generowanego wyjścia. Wybrana strategia jest hiperparametrem modelu i należy ją konsekwentnie stosować we wszystkich krokach.

Oto próbka losowego tekstu wygenerowanego w powyższy sposób przez model wytrenowany na bazie dzieł Szekspira (Karpathy, 2015):

*Marry, and will, my lord, to weep in such a one were prettiest;
Yet now I was adopted heir
Of the world's lamentable day,
To watch the next way with his father with his face?*

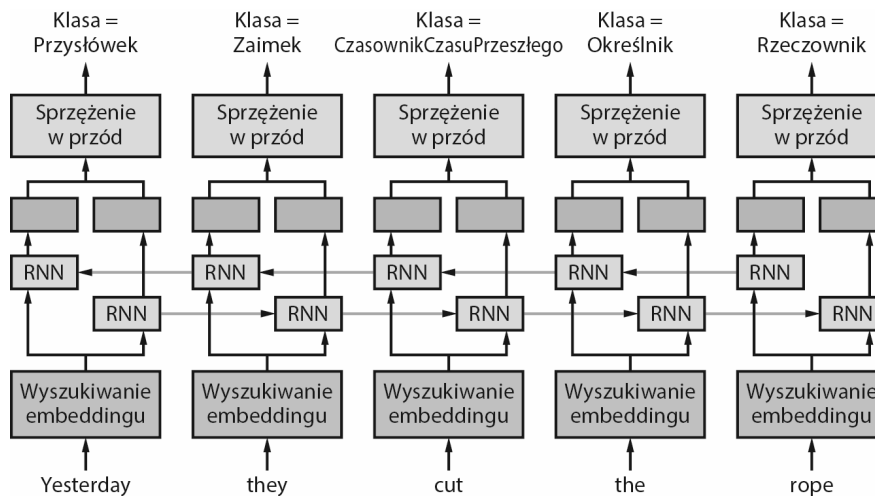
24.2.2. Klasyfikacja z użyciem sieci RNN

Rekurencyjne sieci neuronowe przydają się także w innych zadaniach NLP, między innymi tagowaniu części mowy i rozwiązywaniu koreferencji. W obu tych przypadkach warstwa wejściowa i warstwy ukryte pozostają niezmienione, lecz w przypadku tagera części mowy wyjściowy rozkład softmax dotyczy tagów POS, a w przypadku rozwiązywania koreferencji — możliwych poprzedników. Gdy na przykład sieć otrzyma na wejściu zdanie *Eduardo told me that Miguel was very sick so I took **him** to the hospital*, powinna przypisać wysokie prawdopodobieństwo słowu *Miguel*.

Trenowanie sieci RNN pod kątem takich zadań przebiega niemal tak samo, jak w przypadku modelu językowego RNN, z tą jednak różnicą, że dane treningowe *wymagają etykiet* — tagów części mowy lub wskaźników odwołań. Kolekcjonowanie takich danych jest znacznie trudniejsze w porównaniu z danymi nieetykietowanymi, które dla trenowania modelu językowego są wszystkim, co potrzebne.

Zadaniem modelu językowego jest przewidywanie n -tego słowa na podstawie poprzednich słów, natomiast w zadaniach klasyfikacyjnych nie ma żadnego powodu, byśmy ograniczali się wyłącznie do słów poprzedzających, bo równie owocne może być przeglądanie tekstu w przód. W naszym przykładzie obiekt koreferencji, do którego odnosi się zaimek *him*, byłby inny, gdyby zakończenie zdania brzmiało *to see Miguel*, a nie *to the hospital*, co doskonale ilustruje znaczenie przeglądania w przód. Liczne badania dowiodły zresztą, że człowiek analizujący wzrokowo tekst nie czyni tego *strictly* od lewej do prawej.

By uchwycić prawostronny kontekst słowa, możemy posłużyć się **dwukierunkową siecią RNN** (ang. *bidirectional RNN*) konkatenującą w sobie dwa oddzielne modele — „od lewej do prawej” i „od prawej do lewej”. Efekt zastosowania takiej sieci do tagowania części mowy widoczny jest na rysunku 24.4.



RYSUNEK 24.4. Dwukierunkowa sieć RNN dla tagowania części mowy

W wielowarstwowej sieci RNN z_t będzie ukrytym wektorem ostatniej warstwy. W dwukierunkowej sieci RNN wektor z_t wybierany jest zazwyczaj jako kombinacja wektorów obu modeli: „od lewej do prawej” i „od prawej do lewej”.

Sieci RNN mogą być także używane do klasyfikacji na poziomie całych zdań (a nawet na poziomie dokumentów), kiedy to wyjściem jest pojedyncza wartość, zamiast strumienia wartości wyjściowych, po jednej dla każdego kroku czasowego. Przykładem klasyfikacji tego rodzaju jest **analiza sentymentu**, polegająca na ocenie tekstu jako *Pozytywnego*

albo *Negatywnego* — na przykład recenzja *This movie was poorly written and poorly acted* („Ten film to kiepski scenariusz i słaba gra aktorów”) ma wydzźwięk zdecydowanie *Negatywny*. (Niektóre systemy analizy sentymentu wykorzystują więcej niż dwie kategorie albo dokonują oceny w postaci liczbowej).

Wykorzystywanie modelu RNN *na poziomie zdań* jest nieco bardziej skomplikowane, musimy bowiem uzyskać reprezentację \mathbf{y} całego zdania jako *agregację* reprezentacji wyjściowych \mathbf{y}_t poszczególnych słów tego zdania. Najprostszym sposobem uzyskania takiej agregacji jest wykorzystanie ukrytego stanu sieci odpowiadającego *ostatniemu słowu* zdania wejściowego — bieżący krok czasowy zawiera wówczas informację będącą funkcją całego zdania. Tak otrzymana agregacja jest jednak wyraźnie obciążona w tym sensie, że w większym stopniu odzwierciedla końcówkę zdania niż jego początek. Bardziej sprawiedliwą alternatywą jest pooling wszystkich ukrytych wektorów; najprostszą jego odmianą jest **pooling uśredniający** (ang. *averaging pooling*), którego wynikowa agregacja jest średnią wektorów składowych:

$$\tilde{\mathbf{z}} = \frac{1}{s} \sum_{t=1}^s \mathbf{z}_t$$

Wynikowy, d -wymiarowy wektor $\tilde{\mathbf{z}}$ może zostać następnie podany na jedną lub więcej warstw sprzężenia w przód, zanim pojawi się w warstwie wyjściowej.

24.2.3. Pamięci LSTM w zadaniach NLP

Jak wcześniej stwierdziliśmy, sieci RNN zdolne są w wielu przypadkach do niwelowania problemu niedostatecznego kontekstu. Teoretycznie można w nich przekazywać dowolną informację od jednej ukrytej warstwy do następnej, dla dowolnej liczby kroków czasowych. W praktyce jednak przekazywana informacja może zostać utracona lub zniekształcona, mniej więcej tak, jak w grze w „głuchy telefon”: gracze ustawiają się w kolejkę, pierwszy szepcze drugiemu do ucha jakiś komunikat, drugi gracz szeptem przekazuje ów komunikat trzeciemu, itd.; zazwyczaj komunikat otrzymany przez ostatniego z graczy wyraźnie różni się od oryginału wymyślonego przez pierwszego gracza. W sieciach RNN problem ten jest analogią zjawiska **zanikającego gradientu**, opisywanego w sekcji 21.1.2 — z tą różnicą, że poszczególne warstwy odpowiadają teraz kolejnym krokom czasowym, a nie kolejnym stopniom zagłębienia.

W sekcji 21.6.2 przedstawiliśmy podstawy modelu **długotrwałej pamięci krótkoterminowej (LSTM)**. Jest to odmiana sieci RNN, której bramki wolne są od problemu zniekształcania komunikatu przy kopiowaniu go między krokami czasowymi. Sieć RNN może wybierać do *zapamiętywania* niektóre części wejścia, kopiować je do następnego kroku czasowego i zapominać o innych częściach. Rozważmy przykładowe zdanie:

The athletes, who all won their local qualifiers and advanced to the finals in Tokyo, now ...

(„Lekkoatleci, którzy wygrali lokalne kwalifikacje i awansowali do finałów w Tokio, teraz ...”). Jeśli spytamy model, które ze słów⁴: *competes* czy *compete* powinno wystąpić w miejscu wielokropka, spodziewamy się otrzymać odpowiedź *compete* jako zgodną z podmiotem *The athletes*. LSTM może nauczyć się tworzyć ukrytą cechę dla określonej osoby i liczby, i kopiować w przód tę cechę bez zmian aż do momentu, gdy okaże się konieczne dokonanie opisanego wyboru. Warto zaznaczyć, że „regularne” sieci RNN (a także modele n -gramowe) często mylą się w przypadku długich zdań, w których podmiot oddzielony jest od orzeczenia długim ciągiem słów.

⁴ *Competes* („rywalizuje”) to forma specyficzna dla 3. osoby liczby pojedynczej, w zdaniu mamy do czynienia z 3. osobą liczby mnogiej (*they*), więc musi być *compete* — *przyp. tłum.*

24.3. Modele „sekwencja na sekwencję”

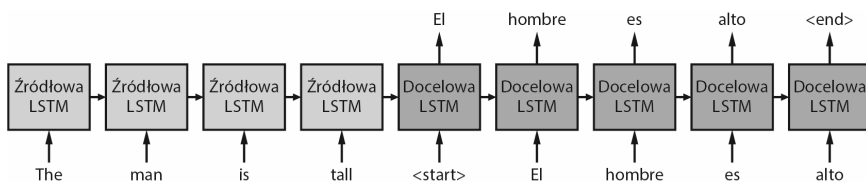
Jednym z najintensywniej badanych obszarów NLP jest **tłumaczenie maszynowe** (ang. MT — *Machine Translation*), którego celem jest tłumaczenie zdania zapisanego w języku źródłowym (na przykład hiszpańskim) na równoważne zdanie wyrażone w **języku docelowym** (na przykład angielskim). Trenowanie modeli MT odbywa się zwykle w oparciu o duże korpusy zawierające pary równoważnych zdań wyrażonych w obu językach, a jego celem jest nauczenie modelu dokładnego tłumaczenia nowych zdań, niewystępujących w zbiorach treningowych.

Czy możliwe jest budowanie systemów MT w oparciu o sieci RNN? Oczywiście możemy zakodować zdanie źródłowe w formie sieci RNN i gdyby tłumaczenie zdań polegało wyłącznie na odwzorowywaniu „jeden-do-jednego” słów języka źródłowego na słowa języka docelowego, zadanie MT sprowadzałoby się do prostego tagowania słów — hiszpańskie słowo *perro* („pies”) zostałoby otagowane przez angielski odpowiednik *dog*. W rzeczywistości tłumaczenie zdań (i w ogóle sekwencji słów) nie daje się zdekomponować na tłumaczenie oddzielnych słów; przykładowo hiszpańskiemu określeniu *caballo de mar* („konik morski”) odpowiada pojedyncze słowo angielskie *seahorse*, a w tłumaczeniu *perro grande* → *big dog* („duży pies”) odpowiadające sobie słowa występują w odwrotnej kolejności. Zmiana kolejności słów w wyniku tłumaczenia może być bardziej spektakularna; w języku angielskim podmiot zwykle rozpoczyna zdanie, ale na przykład w języku, którym posługują się mieszkańcy Republiki Fidżi, występuje on przeważnie na końcu zdania. Skoro więc nie odwzorowanie na poziomie słów, to w jaki sposób mamy generować zdania w języku docelowym?

Wygląda na to, że powinniśmy generować jedno słowo na raz, jednocześnie śledząc kontekst, zapamiętując te części źródła, które nie zostały jeszcze przetłumaczone i śledząc te, które już przetłumaczone zostały (byśmy nie musieli się powtarzać). W przypadku niektórych zdań może być i tak, że będziemy musieli przeczytać całe zdanie źródłowe przed wygenerowaniem pierwszego słowa zdania docelowego.

Uwidacznia to związek generowania tekstu w ramach MT ze standardowym modelem językowym RNN, opisywanym w podrozdziale 24.2. Istotnie, model RNN wytrenowany pod kątem języka angielskiego, wygeneruje raczej sekwencję *big dog* niż *dog big*. Nie zapominajmy jednak, że mamy do czynienia z *tłumaczeniem* — nie chodzi nam więc o generowanie losowych zdań w języku docelowym, lecz zdań stanowiących odpowiedniki zdań tekstu źródłowego. Najprostszym sposobem realizacji tego zadania jest użycie dwóch sieci RNN, z których jedna (nazwijmy ją „źródłową”) wytrenowana jest w języku źródłowym, a druga (docelowa) w języku docelowym. Podajemy następnie zdanie źródłowe na wejście źródłowej RNN, po czym wykorzystujemy jej finalny ukryty stan jako początkowy ukryty stan docelowej RNN. Tak oto każde słowo docelowe zostaje samoczynnie uzależnione zarówno od całego zdania źródłowego, jak i od wygenerowanych poprzednio słów docelowych.

Taką architekturę sieci neuronowej nazywamy podstawowym modelem „sekwencja na sekwencję” (ang. *sequence-to-sequence model*) — przykład takowego przedstawiamy na rysunku 24.5. Modele takie wykorzystywane są najczęściej w tłumaczeniu maszynowym, ale również w kilku innych zadaniach, takich jak automatyczne generowanie podpisów (tytułów) obrazów, streszczenia oraz skracanie tekstu, bez zmiany jego znaczenia.



RYСУNEK 24.5. Podstawowy model „sekwencja na sekwencję”. Każdy blok reprezentuje jeden krok czasowy LSTM (dla prostoty pominięliśmy embeddingi i warstwy wyjściowe). W kolejnych krokach przekazujemy do sieci kolejne słowa zdania źródłowego *The man is tall* („Ten mężczyzna jest wysoki”), zakończonego ogranicznikiem *<start>*, oznaczającym, że sieć może rozpocząć generowanie zdania docelowego. Końcowy ukryty stan na końcu zdania źródłowego wykorzystywany jest jako ukryty stan na starcie zdania docelowego. Następnie każde słowo zdania docelowego w kroku t wykorzystywane jest jako wejście w kroku $t + 1$ aż do momentu, gdy sieć wygeneruje znacznik *<end>* oznaczający, że generowanie zdania docelowego zakończyło się

Podstawowe modele „sekwencja na sekwencję” okazały się przełomem w przetwarzaniu języka naturalnego, a w szczególności w tłumaczeniu maszynowym. Według artykułu Wu i in. (2016b) przyczyniły się one do 60% redukcji błędów w stosunku do poprzednich metod MT. Nie są one jednak wolne od mankamentów, z których trzy najważniejsze to:

- **Polaryzacja w kierunku najbliższego kontekstu.** Wszystko, co sieć RNN zdecyduje się zapamiętywać odnośnie przeszłości, musi być zgodne z jej ukrytym stanem. Jeśli na przykład sieć przetwarza aktualnie 57. słowo (czyli krok czasowy $t = 57$) sekwencji liczącej 70 słów, to ukryty stan sieci prawdopodobnie zawierał więcej informacji na temat kroku $t = 56$ niż na temat kroku $t = 5$. Jest tak dlatego, że każdorazowo, gdy ukryty wektor jest aktualizowany, pewna część istniejącej informacji musi zostać zastąpiona przez nową informację. Takie zachowanie jest zamierzoną konsekwencją projektu modelu i jest szczególnie użyteczne w zadaniach NLP, gdzie najbliższy kontekst jest przeważnie ważniejszy od dalszego. W sytuacji jednak, gdy dalszy kontekst jest równie ważny, czy nawet kluczowy dla zadania, może zostać zmarginalizowany lub wręcz zagubiony przez model RNN i nawet LSTM może napotkać trudności nie do przewyżczenia.
- **Ustalony limit rozmiaru kontekstu.** W modelu tłumaczeniowym RNN całe zdanie źródłowe skompresowane zostaje do postaci pojedynczego ukrytego wektora stanu o ustalonej wymiarowości. W pamięciach LSTM wykorzystywanych w systemach NLP z górnej półki wymiarowość ta wynosi typowo 1024, jeśli więc chcemy reprezentować w jej ramach zdanie zawierające (powiedzmy) 64 słowa, to wypada nam $1024/64 = 16$ wymiarów na słowo — stanowczo zbyt mało dla skomplikowanych zdań. Powiększanie rozmiaru wspomnianego wektora skutkuje spowalnianiem treningu i stwarza ryzyko przeuczenia modelu.
- **Sekwencyjność spowalniająca przetwarzanie.** Jak wyjaśnialiśmy w podrozdziale 21.3, sieci neuronowe przyczyniają się do znacznego wzrostu wydajności, ponieważ przetwarzają dane treningowe partiami, by efektywnie wykorzystywać sprzętowe możliwości zrównoleglenia obliczeń oferowane przez procesory macierzowe. Tak się jednak składa, że sieci RNN ograniczone są raczej do sekwencyjnego przetwarzania słowa po słowie i możliwości zrównoleglenia są w ich przypadku co najwyżej szczątkowe.

24.3.1. Mechanizmy uwagi w modelach NLP

Docelowa RNN może być uzależniona od *wszystkich* ukrytych wektorów źródłowej RNN, nie tylko od ostatniego. Nie występują wtedy problemy polaryzacji kontekstu ani ograniczonego rozmiaru kontekstu, model ma wówczas równoprawny dostęp do każdego poprzedniego słowa. Jednym ze sposobów zapewnienia takiego dostępu jest konkatencja wszystkich ukrytych wektorów źródłowej RNN, jednak takie rozwiązanie skutkuje gwałtownym wzrostem liczby wag, z konsekwencjami w postaci wydłużenia czasu obliczeń i zaistnienia ryzyka przeuczenia modelu. Rozwiązanie bardziej efektywne zasadza się na spostrzeżeniu, że docelowa RNN generuje po jednym słowie na raz, a więc w danej chwili (kroku czasowym) prawdopodobnie tylko niewielka część wejścia będzie rzeczywiście istotna (relewantna) dla danego słowa docelowego.

Co najważniejsze, owe „relewantne części” informacji źródłowej mogą być różne dla poszczególnych słów docelowych. Załóżmy, że wytrenowaliśmy sieć pod kątem tłumaczenia z języka angielskiego na hiszpański, a na jej wejściu pojawiła się sekwencja *The front door is red* („Drzwi frontowe są czerwone”), po której następuje znacznik końca zdania, uprawniający do rozpoczęcia generowania hiszpańskich słów. Sieć powinna więc zwrócić uwagę na słowo *The* i wygenerować słowo *La*, następnie zwrócić uwagę na słowo *door* i wygenerować *puerta*, i tak dalej.

Owa koncepcja „zwracania uwagi” odzwierciedlona została w projekcie sieci w postaci jej komponentu zwanego (po prostu) **uwagą** (ang. *attention*), za pomocą którego można kreować „kontekstowe podsumowanie” zdania źródłowego w postaci wektora o ustalonej wymiarowości. Ów kontekstowy wektor c_i używany jest następnie jako dodatkowe wejście do docelowej RNN. Model „sekwencja na sekwencję” wykorzystujący taki komponent nazywany jest **uwagowym modelem „sekwencja na sekwencję”** (ang. *attentional sequence-to-sequence model*). Jeśli standardową docelową sieć RNN zapiszemy w postaci

$$\mathbf{h}_t = RNN(\mathbf{h}_{t-1}, \mathbf{x}_t)$$

to docelową sieć RNN dla wspomnianego modelu zapisać można w tej konwencji jako

$$\mathbf{h}_i = RNN(\mathbf{h}_{i-1}, [\mathbf{x}_i; \mathbf{c}_i])$$

gdzie $[\mathbf{x}_i; \mathbf{c}_i]$ oznacza konkatencję wektorów wejściowego (\mathbf{x}_i) i kontekstowego (\mathbf{c}_i), przy czym ten ostatni definiowany jest następująco:

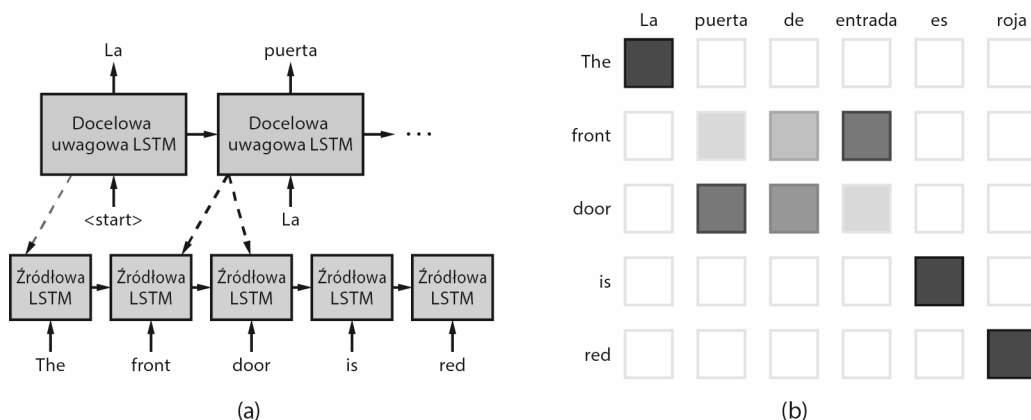
$$r_{ij} = \mathbf{h}_{i-1} \cdot \mathbf{s}_j$$

$$a_{ij} = \frac{e^{r_{ij}}}{\sum_k e^{r_{ik}}}$$

$$\mathbf{c}_i = \sum_j a_{ij} \cdot \mathbf{s}_j$$

gdzie \mathbf{h}_{i-1} jest wektorem docelowej RNN, który będzie użyty do przewidywania słowa docelowego w kroku czasowym i , a \mathbf{s}_j jest wektorem wyjściowym źródłowej RNN dla słowa źródłowego j (czyli kroku czasowego j). Zarówno \mathbf{h}_{i-1} jak i \mathbf{s}_j są wektorami d -wymiarowymi, gdzie d jest ukrytym rozmiarem. Wartość r_{ij} jest więc surową „punktacją uwagową” (ang. *attention score*) między bieżącym stanem docelowym a j -tym słowem źródłowym. Punktacje te są następnie normalizowane do prawdopodobieństwa a_{ij} według rozkładu softmax dla wszystkich słów źródłowych. Ostatecznie prawdopodobieństwa te wykorzystywane są do generowania średniej ważonej wektorów źródłowej sieci RNN, czyli \mathbf{c}_i (która to średnia również jest wektorem d -wymiarowym).

Przykład zastosowania uwagowego modelu „sekwencja na sekwencję” widoczny jest w części (a) rysunku 24.6, zawierającego kilka istotnych szczegółów. Po pierwsze, komponent uwagi nie posiada sam z siebie wyuczonych wag i obsługuje sekwencje o zmiennej długości, zarówno po stronie źródłowej, jak i docelowej. Po drugie, analogicznie do wielu innych modeli opartych na sieciach neuronowych, komponent ten jest całkowicie ukryty. Programista nie może dyktować modelowi, jakich ten ma używać informacji — model uczy się tego samodzielnie. Komponent uwagi może być łączony z wielowarstwowymi sieciami RNN, zwykle komponenty takie znajdują się wówczas w każdej warstwie.



RYСУNEK 24.6. (a) Uwagowy model „sekwencja na sekwencję” dla tłumaczenia z języka angielskiego na hiszpański (*The front door is red* — „Drzwi frontowe są czerwone”). (b) Przykład uwagowej macierzy prawdopodobieństw dla dwujęzycznej pary zdań, ciemniejsze kwadraty reprezentują większe wartości a_{ij} . Wartości prawdopodobieństw uwagowych sumują się do 1 w każdej kolumnie

Sformułowanie funkcjonalności komponentu uwagi w kategoriach rozkładu softmax służy trojakiemu celowi. Po pierwsze, komponent ten staje się różniczkowalny, co jest konieczne do użycia go w propagacji wstecz; propagacja ta wciąż ma miejsce w źródłowej i docelowej RNN, mimo iż jej komponenty uwagi nie posiadają własnych wyuczonych wag. Po drugie, formuła probabilistyczna umożliwia modelowi uchwycenie pewnych typów długodystansowych kontekstualizacji, które inaczej pozostawałyby niezauważalne dla źródłowej RNN, ponieważ komponent uwagi może rozważać całą sekwencję źródłową naraz i nauczyć się wybierać z niej to, co istotne. Po trzecie wreszcie, probabilistyczny komponent uwagi pozwala sieci na reprezentowanie niepewności — jeśli sieć nie wie dokładnie, jakie słowo źródłowe ma przetłumaczyć w następnej kolejności, może uwzględnić kilka opcji, opatrując je prawdopodobieństwami, i następnie wybrać właściwe słowo, wykorzystując docelową RNN.

W przeciwieństwie do większości komponentów sieci neuronowych, prawdopodobieństwa komponentu uwagi są intuicyjnie zrozumiałe dla ludzi. Przykładowo, w przypadku tłumaczenia maszynowego prawdopodobieństwa te mogą być odnoszone do kojarzenia na poziomie pojedynczych słów, stosowanego często przy ręcznym tłumaczeniu; przykład takiego kojarzenia przedstawiamy w części (b) rysunku 24.6.

Modele „sekwencja na sekwencję” w naturalny sposób przystosowane są do tłumaczenia maszynowego, ale niemal każde zadanie NLP można zakodować w postaci problemu rozwiązywalnego przez taki model. Przykładowo, system odpowiadania na pytania można trenować na zbiorze par, których pierwszym elementem jest pytanie, a drugim prawidłowa odpowiedź ograniczona specjalnym znacznikiem końca.

24.3.2. Dekodowanie

W czasie treningu model „sekwencja na sekwencję” dąży do maksymalizacji prawdopodobieństwa wystąpienia każdego słowa w docelowej sekwencji treningowej, w zależności od źródła i wszystkich poprzednich słów docelowych. Po zakończeniu treningu dostarczamy modelowi zdanie źródłowe i oczekujemy wygenerowania odpowiadającego mu zdania docelowego. Jak wynika to z rysunku 24.6, możemy generować po jednym słowie docelowym na raz, po czym propagujemy wstecz słowo, które wygenerowaliśmy w następnym kroku czasowym. Procedura ta nosi nazwę **dekodowania** (ang. *decoding*).

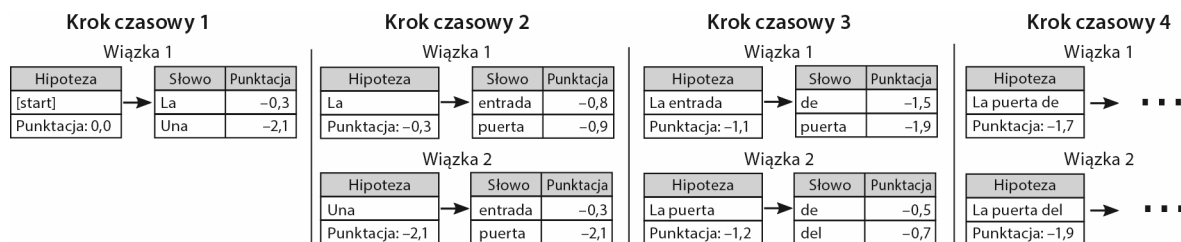
Najprostszą formą dekodowania jest wybieranie, w każdym kroku czasowym, słowa o największym prawdopodobieństwie i przekazywanie tego słowa jako wejścia do następnego kroku czasowego. Taką strategię nazywamy **zachłannym dekodowaniem** (ang. *greedy decoding*), ponieważ po wygenerowaniu każdego ze słów docelowych system bezgranicznie polega na hipotezie, którą sformułował do tej pory. Problem jednak w tym, że celem dekodowania jest maksymalizacja prawdopodobieństwa *całej sekwencji docelowej*, a tego zachłanna strategia nie jest w stanie zagwarantować. Rozpatrzmy na przykład użycie zachłannego dekodera do przetłumaczenia cytowanego wcześniej zdania *The front door is red*.

Poprawnym tłumaczeniem jest *La puerta de entrada es roja* („Drzwi wejściowe są czerwone”). Załóżmy, że docelowa RNN poprawnie wygeneruje pierwsze słowo *La* jako odpowiednik *The*, po czym zachłanny dekodery w charakterze drugiego słowa zaproponuje *entrada* jako odpowiednik *front*. I tu pojawi się błąd — zgodnie z hiszpańskim szykiem wyrazów, podmiot *puerta* powinien wystąpić przed przydawką *entrada*. Zachłanne dekodowanie jest co prawda szybkie — w każdym kroku czasowym dokonywany jest tylko jeden wybór — ale niestety, model nie dysponuje mechanizmem umożliwiającym korygowanie popełnionych błędów.

Moglibyśmy spróbować ulepszyć mechanizm uwagi tak, by zawsze zwracał uwagę na właściwe słowo i za każdym razem poprawnie je odgadywał. To jest jednak niemożliwe w przypadku niektórych zdań, bo nie sposób prawidłowo odgadnąć początkowe wyrazy zdania bez wiedzy o tym, co znajduje się na jego końcu.

Lepszym podejściem jest poszukiwanie optymalnego (lub przynajmniej — dobrego) dekodowania przy użyciu jednego z algorytmów wyszukiwania opisywanych w rozdziałach 3. i 4. Często wykorzystuje się w tej roli **lokalne wyszukiwanie skupione** (patrz sekcja 4.1.3). W kontekście dekodowania MT, przeszukiwanie to utrzymuje zazwyczaj k najlepszych hipotez na każdym etapie, rozszerzając każdą o jedno słowo, na k sposobów przy użyciu k najlepszych wyborów słów, a następnie wybiera najlepsze k z wynikowych k^2 nowych hipotez. Gdy każda z hipotez w wiązce wygeneruje specjalny token <end>, algorytm wyprowadza hipotezę o najwyższej punktacji.

Rysunek 24.7 przedstawia wizualizację przykładowego wyszukiwania skupionego. W miarę jak model głębokiego uczenia staje się bardziej dokładny, możemy przeważnie ograniczyć rozmiar wiązki; w nowoczesnych modelach neuronowych MT rozmiar ten waha się między 4 a 8, podczas gdy starsze modele statystycznego MT posługiwały się wiązkami o rozmiarach 100 i więcej.



RYСУNEK 24.7. Wyszukiwanie skupione z rozmiarem wiązki $b = 2$. Punktacja każdego słowa równa jest logarytmicznemu prawdopodobieństwu wynikającemu z rozkładu softmax generowanego przez sieć RNN, a punktacja każdej hipotezy jest sumą punktacji poszczególnych słów. W kroku czasowym $t = 3$ hipoteza o najwyższej punktacji *La entrada* może generować jedynie kontynuacje o niskim prawdopodobieństwie, więc plasuje się poza wiązką.

24.4. Architektura transformerów

Vaswani i in. w swym wpływowym artykule z 2018 roku wprowadzają architekturę **transformera** (ang. *transformer*), wykorzystującą mechanizm **samouwagi** (ang. *self-attention*) zdolny do modelowania długodystansowego kontekstu, bez zależności sekwencyjnych.

24.4.1. Samouwaga

W modelach „sekwencja na sekwencję” uwaga zwrócona była ze strony docelowej RNN na źródłową RNN. Mechanizm **samouwagi** rozszerza tę funkcjonalność w ten sposób, że każda sekwencja ukrytych stanów (źródłowa i docelowa) zwraca uwagę również na samą siebie. Pozwala to modelowi dodatkowo uchwycić kontekst zarówno długodystansowy, jak i pobliski w każdej sekwencji.

Najprostszym wariantem zastosowania samouwagi jest sytuacja, w której macierz uwagi tworzona jest bezpośrednio przez iloczyn skalarny wektorów wejściowych. Przypadek taki jest jednak o tyle problematyczny, że iloczyn skalarny wektora przez samego siebie zawsze jest duży, więc każdy ukryty stan będzie preferował zwracanie uwagi na samego siebie. Transformer rozwiązuje ten problem, rozpoczynając od rzutowania danych wejściowych na trzy różne reprezentacje, przy użyciu trzech różnych macierzy wagowych⁵:

- **wektor zapytania** $\mathbf{q}_i = \mathbf{W}_q \mathbf{x}_i$ jest tym *obserwującym*, podobnie jak strona docelowa w standardowym mechanizmie uwagi;
- **wektor klucza** $\mathbf{k}_i = \mathbf{W}_k \mathbf{x}_i$ jest tym *obserwowanym*, podobnie jak strona źródłowa w standardowym mechanizmie uwagi;
- **wektor wartości** $\mathbf{v}_i = \mathbf{W}_v \mathbf{x}_i$ jest kontekstem generowania.

⁵ Patrzyć także <https://ichi.pro/pl/google-deepmind-s-rfa-approximating-softmax-attention-mechanism-in-transformers-235867597421592> — przyp. tłum.

W standardowym mechanizmie uwagi sieci odpowiadające kluczowi i wartości są identyczne, ale intuicyjnie ma sens ich oddzielnie reprezentowanie. Rezultat kodowania i -tego słowa — \mathbf{c}_i — można obliczyć, stosując mechanizm uwagi do rzutowanych wektorów.

$$r_{ij} = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d}}$$

$$a_{ij} = \frac{e^{r_{ij}}}{\sum_m e^{r_{im}}}$$

$$\mathbf{c}_j = \sum_i a_{ij} \cdot \mathbf{v}_i$$

gdzie d jest wymiarowością wektorów \mathbf{k} i \mathbf{q} . Zauważmy, że i oraz j są indeksami słów w tym samym zdaniu, ponieważ kodujemy kontekst w konwencji samouwagi. W każdej warstwie transformera samouwaga wykorzystuje ukryte wektory poprzedniej warstwy — którą początkowo jest warstwa kontekstowa.

Kilka szczegółów w powyższych formułach wartych jest szczególnej uwagi. Przede wszystkim mechanizm samouwagi jest *asymetryczny*, ponieważ $r_{ij} \neq r_{ji}$. Po drugie, rolą czynnika skalującego \sqrt{d} jest poprawienie stabilności numerycznej algorytmu. Po trzecie, kodowanie poszczególnych słów w zdaniu może być prowadzone *równoległe* — fundamentem algorytmu są operacje macierzowe, dla których współczesne procesory oferują wydajne implementacje sprzętowe.

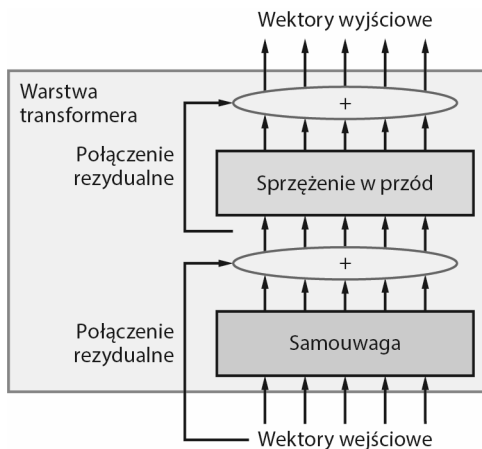
Wybór kontekstu do wykorzystania nie jest z góry określony, lecz całkowicie wyuczony na przykładach treningowych. Podsumowanie kontekstowe — \mathbf{c}_i — to suma wszystkich poprzednich pozycji w zdaniu. Teoretycznie wszelkie informacje zawarte w zdaniu mogą pojawić się w \mathbf{c}_i . Jednym ze sposobów zapobiegania temu zjawisku jest **wielouwaga** (ang. *multiheaded attention*): dzielimy zdanie na m równych części i do każdej z nich stosujemy model uwagi; każda część ma swój własny zbiór wag. Konkatenujemy następnie wyniki, otrzymując \mathbf{c}_i ; konkatenacja, w przeciwieństwie do sumowania, nie stwarza ryzyka utraty informacji cząstkowych.

24.4.2. Od samouwagi do transformera

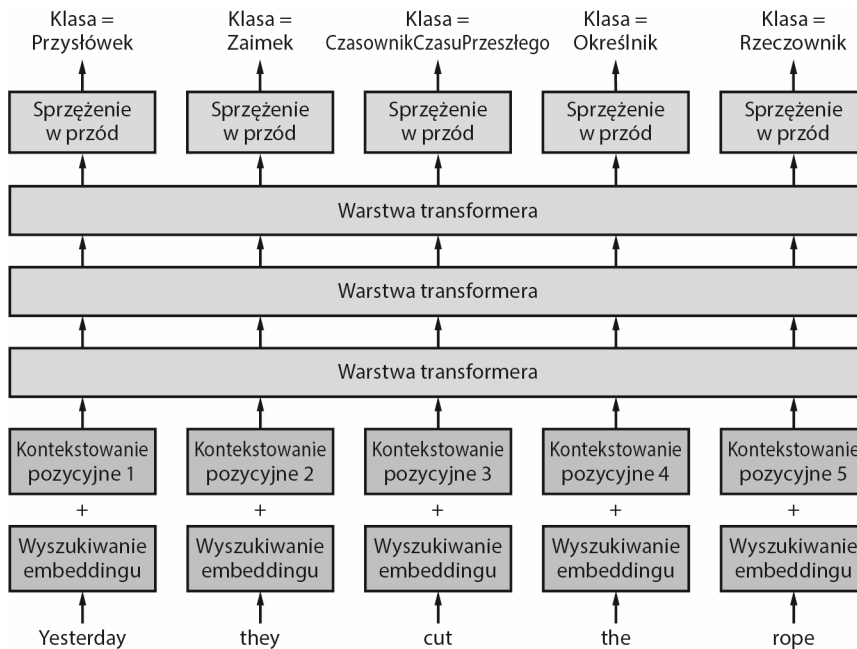
Samouwaga to tylko jeden z elementów modelu transformera. Każda warstwa transformera składa się z kilku podwarstw. Do każdej warstwy transformera najpierw stosowana jest samouwaga, następnie wyjście modułu uwagi kierowane jest na warstwy sprzężenia w przód, gdzie te same macierze wagowe sprzężenia stosowane są niezależnie do każdej pozycji. Nieliniowa funkcja aktywacyjna — zazwyczaj ReLU — stosowana jest po pierwszej warstwie sprzężenia. Aby rozwiązać potencjalny problem zanikającego gradientu, do warstwy transformera dodawane są dwa połączenia rezydualne. Jednowarstwowy transformer widoczny jest na rysunku 24.8. W praktyce modele transformatorów zawierają zwykle sześć lub więcej warstw. Podobnie jak w przypadku innych wcześniej opisywanych modeli, wyjście warstwy i jest podawane na wejście warstwy $i+1$.

Architektura transformera nie utrzymuje jawnej informacji o kolejności słów w sekwencji, ponieważ kontekst modelowany jest wyłącznie poprzez samouwagę, która kolejności tej nie zna. By jednak informację o tej kolejności utrzymywać, transformer wykorzystuje technikę **embedowania pozycyjnego** (ang. *positional embedding*) Jeśli wejściowa sekwencja ma maksymalną długość n , model uczy się n nowych embeddingów, po jednym dla każdej pozycji słowa. Wejściem do pierwszej warstwy transformera jest suma embeddingów słowa na pozycji t oraz embeddingu pozycyjnego odpowiadającego pozycji t .

Rysunek 24.9 przedstawia architekturę transformera dla tagowania części mowy (POS), w zastosowaniu do tego samego zdania, które widoczne jest u dołu rysunku 24.2. W dolnej części rysunku 24.9 embeddingi słów sumowane są z embeddingami pozycyjnymi, co daje w rezultacie wejście do trójwarstwowego transformera. Transformer produkuje jeden wektor dla każdego słowa, tak jak w tagowaniu POS opartym na RNN. Każdy taki wektor podawany jest na końcową warstw wyjściową i na warstwę softmax w celu obliczenia rozkładu prawdopodobieństwa poszczególnych tagów.



RYSUNEK 24.8. Jednowarstwowy transformer składający się z samouwagi, sieci ze sprzężeniem w przód i połączeń rezydualnych



RYSUNEK 24.9. Wykorzystanie architektury transformera do tagowania części mowy.

To, co dotychczas napisaliśmy o transformerach, to tylko część ich historii — ta sekcja ogranicza się do modelu zwanego **koderem transformera** (ang. *transformer encoder*), użytecznego w zadaniach klasyfikacji tekstu. Pełna architektura transformera została oryginalnie zaprojektowana jako model „sekwencja na sekwencję” na potrzeby tłumaczenia maszynowego. Oprócz kodera transformera istnieje więc także **dekoder transformera** (ang. *transformer decoder*). Koder i dekodek są niemal identyczne, z tą różnicą, że dekodek wykorzystuje wersję samouwagi, w której każde słowo może zwracać uwagę wyłącznie na słowa poprzedzające, ponieważ tekst jest generowany od strony lewej do prawej. Dekoder posiada też drugi moduł uwagi w każdej warstwie transformera, który to moduł zwraca uwagę na koder transformera.

24.5. Trenowanie wstępne i uczenie transferowe

Pozyskanie danych w ilości wystarczającej do zbudowania solidnego modelu może być prawdziwym wyzwaniem. W przypadku widzenia komputerowego (patrz rozdział 25.) wyzwanie to rozwiązano poprzez zgromadzenie dużych kolekcji obrazów (takich jak ImageNet) i ręczne ich zaetykietowanie.

Przetwarzanie języka naturalnego to przeważnie praca z nieetykietowanym tekstem — i znacznie trudniejsze jego etykietowanie: niewykształcony obserwator z łatwością może stwierdzić, że obraz przedstawia „kota” albo „zachód słońca”, ale adnotowanie tekstu za pomocą tagów części mowy lub drzew rozbioru syntaktycznego wymaga solidnego treningu. Specyfika przetwarzania tekstu wynika także z jego obfitości: każdego dnia w internecie pojawia się ponad 100 miliardów słów tekstu, w tym zdigitalizowane książki, wyselekcjonowane zasoby, takie jak Wikipedia oraz niere-dagowane posty w mediach społecznościowych.

Łatwy dostęp do takich danych zapewniają projekty takie jak Common Crawl. Do budowania modeli n -gramowych lub embeddingowych można użyć dowolnych strumieni tekstu, a struktura niektórych może być pomocna w rozmaitych zadaniach — istnieje na przykład wiele witryn z często zadawanymi pytaniami (ang. FAQ — *Frequently Asked Questions*), oferujących gotowe pary „pytanie-odpowiedź”; można je wykorzystywać do treningu systemów odpowiadających na pytania. Podobnie na wielu stronach sieci Web znaleźć można publikowane obok siebie warianty tekstu w różnych wersjach językowych, co oczywiście może stanowić cenny materiał treningowy dla systemów tłumaczenia maszynowego. Można też spotkać wiele tekstów etykietowanych w specyficzny sposób: na przykład recenzje, które autorzy opatrują ratingiem w skali (powiedzmy) od 0 do 5.

Wolelibyśmy oszczędzić sobie trudu tworzenia nowego zbioru danych dla każdego nowego modelu NLP. W tej sekcji przedstawiamy ideę **treningu wstępnego** (ang. *pretraining*) — formę **uczenia transferowego** (patrz sekcja 21.7.2), w której używamy dużej ilości współdzielonych danych ogólnej domeny językowej do trenowania początkowej wersji modelu NLP. Po takim wstępnym treningu możemy użyć mniejszej ilości danych specyficznych dla konkretnej domeny (być może danych częściowo etykietowanych) dla doskonalenia modelu — który następnie może uczyć się słownictwa, idiomów, struktur składniowych i innych osobliwości lingwistycznych specyficznych dla nowej domeny.

24.5.1. Wstępny trening embeddingów

W podrozdziale 24.1 przedstawiliśmy pokrótce koncepcję kontekstowania słów, w szczególności pokazując, jak podobieństwo słów — na przykład „jabłko” i „banan” — przekłada się na podobieństwo ich embeddingów i jak dzięki temu można rozwiązywać problemy z analogią słów, posługując się jedynie dodawaniem i odejmowaniem wektorów. Oznacza to, że embeddingi słów wychwytyują istotne informacje na ich temat.

W tej sekcji zagłębimy się w szczegóły tworzenia embeddingów słów przy użyciu całkowicie nienadzorowanego procesu zastosowanego do dużego korpusu tekstu. Różni się to od kontekstowania opisywanego w podrozdziale 24.1, stanowiącego nadzorowany element tagowania części mowy i tym samym wymagającego tagów POS, które są wynikiem pracochłonnego ręcznego adnotowania.

Skoncentrujemy się na jednym konkretnym modelu embeddingowym — GloVe (ang. *Global Vectors*). Model rozpoczyna swe działanie od zliczenia, ile razy każde ze słów pojawia się w oknie innego słowa, podobnie jak w modelach skipgramowych (patrz sekcja 23.1.3). Najpierw wybieramy rozmiar okna (powiedzmy 5 słów). Niech X_{ij} oznacza liczbę wystąpień słów i oraz j w tym samym oknie, a X_i niech oznacza liczbę wystąpień słowa i we wspólnym oknie z jakimkolwiek innym słowem; wówczas $P_{ij} = X_{ij}/X_i$ jest prawdopodobieństwem wystąpienia słowa j w kontekście słowa i . Tak jak poprzednio, oznaczmy przez \mathbf{E}_i embedding słowa i .

Intuicyjny charakter modelu GloVe bierze się w dużej mierze stąd, że związek między dwoma słowami można najlepiej uchwycić przez porównywanie ich z innymi słowami. Rozważmy na przykład słowa *ice* („lód”) i *steam* („para”) i stosunek prawdopodobieństw ich współwystępowania z pewnym słowem w , czyli

$$P_{w, ice}/P_{w, steam}$$

Ponieważ lód jest ciałem stałym (ang. *solid*), a para wodna jest gazem (ang. *gas*), więc $P_{solid, ice}$ będzie wyraźnie większe od $P_{solid, steam}$, a $P_{gas, ice}$ będzie znacznie mniejsze od $P_{gas, steam}$ — w obu przypadkach stosunek wymienionych prawdopodobieństw będzie wyraźnie odbiegał od 1. Gdy natomiast w będzie słowem pozbawionym treści, na przykład *the*, stosunek $P_{w, ice}/P_{w, steam}$ będzie bliski 1; tak samo będzie w przypadku, gdy w będzie słowem jednakowo relevantnym dla *ice* i *steam* (na przykład *water* — „woda”) albo jednakowo dla nich nierelwantnym (na przykład *fashion* — „moda”).

Model GloVe, opierając się na tym intuicyjnym spostrzeżeniu, wykonuje pewne operacje matematyczne (Pennington i in., 2014), które przekształcają stosunki prawdopodobieństw na różnice wektorowe i iloczyny skalarne, co ostatecznie prowadzi do zależności

$$\mathbf{E}_i \cdot \mathbf{E}'_j = \log(P_{ij})$$

Innymi słowy, iloczyn skalarny dwóch embeddingów jest równy logarytmicznemu prawdopodobieństwu ich współwystępowania. Ma to intuicyjny sens: dwa prawie prostopadłe wektory mają iloczyn skalarny bliski zeru, a dwa prawie identyczne znormalizowane⁶ wektory mają iloczyn skalarny bliski 1. Istnieje pewna techniczna komplikacja polegająca na tym, że model GloVe tworzy dwa embeddingi \mathbf{E}_i i \mathbf{E}'_i dla każdego słowa i na końcu dodaje je do siebie; intencją tego zabiegu jest zmniejszenie ryzyka przeuczenia modelu.

Uczenie modelu takiego jak GloVe jest zazwyczaj znacznie tańsze niż uczenie standardowej sieci neuronowej: nowy model można wytrenować na korpusie zawierającym miliardy słów w ciągu kilku godzin przy użyciu standardowego procesora domowego komputera.

Zdarza się, że trenowanie modelu w kontekstowaniu słów z określonej domeny doprowadza do odkrycia wiedzy głęboko tkwiącej w tej domenie, ale jeszcze powszechnie nieuświadomionej. Tshitoyan i in. (2019) wykorzystali 3,3 miliona abstraktów publikacji naukowych z dziedziny materiałoznawstwa jako korpus treningowy. Widzieliśmy wcześniej, jak model embeddingów potrafi przewidywać właściwe słowa na zasadzie analogii („Ateny mają się go Grecji tak, jak Oslo do...”); na tej samej zasadzie wspomniany zespół postanowił uzyskać odpowiedź na pytanie „NiFe tak się ma do ferromagnetyków, jak IrMn do...” — i GloVe udzielił odpowiedzi „antyferromagnetyków”.

Wydaje się, że ich model nie polegał wyłącznie na statystyce współwystępowania słów, lecz dysponował dodatkowo pewną wiedzą domenową. Zapytany o to, jakie związki chemiczne są „termoelektrykami” i „izolatorami topologicznymi”, był w stanie poprawnie odpowiedzieć, mimo iż wzór chemiczny CsAgGa₂Se₄ nigdzie w korpusie nie pojawia się w pobliżu słowa „termoelektryk”, ale występuje wraz ze słowami „chalkogenki”, „przerwa energetyczna” i „optoelektryki”, które są wskazówkami pozwalającymi na zaklasyfikowanie go jako podobnego do „termoelektryka”.

Co więcej, po wytrenowaniu wyłącznie na bazie abstraktów nieprzekraczających 2008 roku i zażądaniu wybrania związków, które są „termoelektryczne”, wśród pięciu czołowych zwróconych przez model znalazły się trzy, które za termoelektryki uznane zostały dopiero w literaturze opublikowanej w latach 2009 – 2019. Tak oto model lingwistyczny antycypował o wiele lat oficjalne odkrycie na gruncie doświadczalnym.

⁶ „Wektor znormalizowany” to wektor o długości 1. Iloczyn skalarny dwóch wektorów równy jest iloczynowi ich długości oraz kosinusa kąta między nimi. Dla wektorów prawie równoległych kąt ten jest bliski 0, a jego kosinus jest bliski 1 — *przyj. tłum.*

24.5.2. Wstępny trening reprezentacji kontekstu zdaniowego

Embeddingi w opisywanej dotychczas formie są *embeddingami absolutnymi* w tym sensie, że nie odzwierciedlają *znaczenia* słów wynikającego z *kontekstu zdań*, w których zostały użyte, co staje się problemem w przypadku słów wieloznacznych (polisemicznych). Przykładowo, słowo *rose* może być rzeczownikiem „róża”, ale może też być formą czasu przeszłego czasownika *rise* („pojawić się”, „wzrastać”, „wschodzić”). Spodziewamy się zatem odnaleźć co najmniej dwa całkowicie różne klasterki kontekstów zdaniowych dla *rose*: jeden skupiający wyrazy „kwiatopodobne” — *dahlia* („georginia”), *freesia* („frezja”), *geranium* („pelargonia”) — i drugi skupiający wyrazy bliskoznaczne do *rise* (na przykład *upsurge* — „ożywienie, rozkwit”); nie jest możliwe odzwierciedlenie obu tych znaczeń w formie pojedynczego embeddingu. *Rose* to tylko jeden przykład słowa o (co najmniej) dwóch całkowicie różnych znaczeniach, lecz język angielski zawiera wiele bardziej subtelnych wieloznaczności, silnie zależnych od kontekstu znaczeniowego: słowo *need*, kojarzące się z „potrzebą”, ma nieco odmienne znaczenie w zdaniach *you need to see this movie* („powinieneś zobaczyć ten film”) i *humans need oxygen to survive* („ludzie potrzebują tlenu do przeżycia”). Jeszcze ciekawszym przykładem wieloznaczności są konstrukcje idiomatyczne — wyrażenie *break the bank* („rozbić bank”) ma znaczenie zupełnie inne niż sugerowane przez każde z jego słów składowych.

Zamiast więc trenować model w oparciu o tablicę odwzorowań „słowo → embedding”, powinniśmy prowadzić trening oparty na reprezentacji odzwierciedlającej **kontekst znaczeniowy** poszczególnych słów tworzących zdanie. Taka reprezentacja danego słowa, zwana jego **embeddingiem znaczeniowym**, to embedding stanowiący wynik mapowania słowa *wraz z otaczającym je kontekstem znaczeniowym*; innymi słowy, jeśli przekażemy modelowi słowo *rose* i kontekst *the gardener planted a rose bush* („ogrodnik zasadził krzew róży”), model ten powinien wyprodukować embedding znaczeniowy podobny do tego odpowiadającego kontekstowi *the cabbage rose had an unusual fragrance* („róża stulistna pachniała niezwykle”), ale zupełnie inny niż dla kontekstu zdania *the river rose five feet* („poziom wody w rzece podniósł się o pięć stóp”).

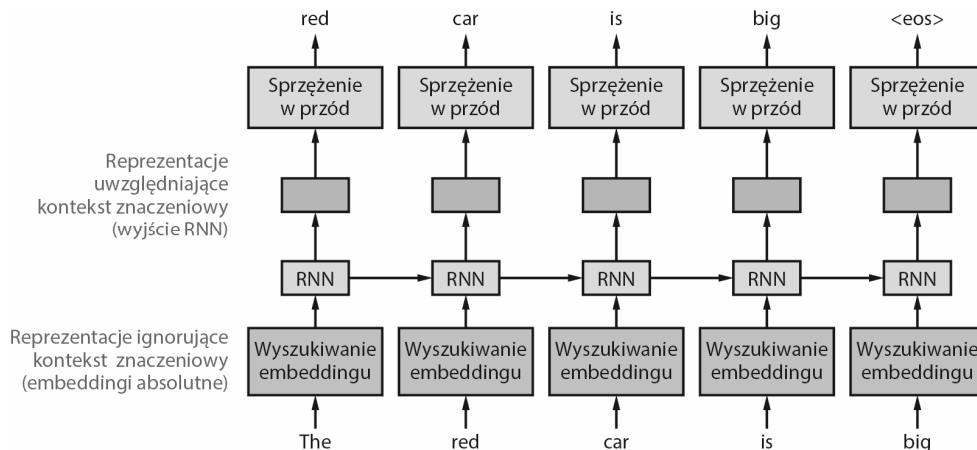
Rysunek 24.10 przedstawia sieć RNN, która tworzy embeddingi znaczeniowe słów zdania *The red car is big* („czerwony samochód jest duży”) — embeddingi te przedstawione są w postaci nieetykietowanych prostokątów. Zakładamy, że zbudowaliśmy już kolekcję embeddingów absolutnych. Podajemy jedno słowo na raz i żądamy od modelu, aby przewidział następane słowo. W przykładzie na rysunku w punkcie, w którym dotarliśmy do słowa *car*, węzeł RNN w tym kroku czasowym otrzyma dwie dane wejściowe: embedding absolutny słowa *car* oraz kontekst, który koduje informacje z poprzednich słów: *The red*. Węzeł RNN wygeneruje następną kontekstową reprezentację dla *car*, a sieć jako całość wyprowadza przewidywanie dla następnego słowa: *is*. Potem aktualizujemy wagi sieci, tak by zminimalizować błąd między przewidywaniem a faktycznie występującym następnym słowem.

Ten model jest podobny do modelu tagowania części mowy (POS) z rysunku 24.4, z dwiema istotnymi różnicami. Po pierwsze, model ten jest jednokierunkowy (od lewej do prawej), natomiast model POS jest dwukierunkowy. Po drugie, zamiast przewidywać tagi POS dla bieżącego słowa, model ten przewiduje następane słowo przy użyciu poprzedniego kontekstu. Po zbudowaniu modelu możemy go używać do wyszukiwania reprezentacji słów i przekazywania ich do innych zadań; nie musimy już przewidywać następnego słowa. Zauważmy, że obliczanie reprezentacji kontekstu zdaniowego zawsze wymaga dwóch danych wejściowych — bieżącego słowa i kontekstu jego występowania.

24.5.3. Maskowane modele językowe

Słabością standardowych modeli językowych, takich jak modele *n*-gramowe, jest to, że kontekstualizacja każdego słowa opiera się tylko na słowach poprzedzających — przewidywania tworzone są od lewej do prawej. Niekiedy jednak kontekst danego słowa staje się zrozumiały dopiero po wyjaśnieniach dostarczanych przez późniejsze zdania, na przykład w zdaniu „Należy pozamykać okna na pulpicie” wyjaśnienia takiego dostarcza słowo „pulpit”.

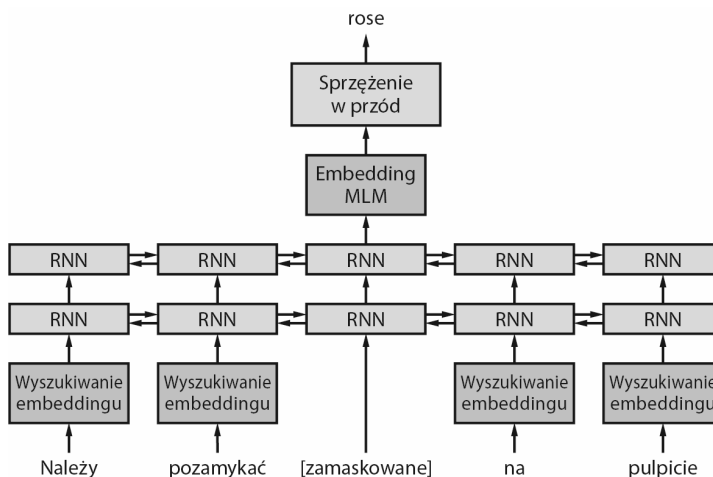
Jednym z prostych sposobów obejścia tego problemu jest wytrenowanie oddzielnego modelu języka, zapisywanego od strony prawej do lewej, który to model kontekstualizuje każde słowo na podstawie następných słów w zdaniu; potem można połączyć oba modele w celu otrzymania połączonej reprezentacji dla obu kierunków. Nie da się jednak w ten sposób zintegrować materiałów dowodowych wnoszonych przez każdy z kierunków.



RYSUNEK 24.10. Trenowanie reprezentacji kontekstu zdaniowego w kierunku od lewej do prawej (*The red car is big* — „czerwony samochód jest duży”).

Zamiast tego możemy użyć **modelu języka maskowanego** (ang. MLM — *Masked Language Model*). Modele MLM trenowane są przez maskowanie (ukrywanie) wybranych słów w danych wejściowych, a następnie żądanie od modelu, by owe zamaskowane słowa odgadywał. Do realizacji tego zadania można użyć głębokiej dwukierunkowej sieci RNN lub transformera w odniesieniu do zdania wejściowego. Na przykład w zdaniu „Należy pozamykać okna na pulpicie” można zamaskować „okna”, uzyskując „należy pozamykać ___ na pulpicie”, i zażądać od modelu wypełnienia powstałej luki.

Przykład zastosowania modelu języka maskowanego przedstawia rysunek 24.11.



RYSUNEK 24.11. Modelowanie języka maskowanego: wstępny trening modelu dwukierunkowego — na przykład wielowarstwowej sieci RNN — przez maskowanie niektórych słów i następnie ich odgadywanie

Finalne ukryte wektory odpowiadające zamaskowanym tokenom są następnie używane do przewidywania zamaskowanych słów — w tym przykładzie „okna”. W procesie trenowania jedno zdanie może być użyte wiele razy, z różnymi pozycjami maskowania. Piękno tego podejścia polega na tym, że nie wymaga ono uprzedniego etykietowania danych — samo zdanie stanowi etykietę dla zamaskowanego słowa. Jeśli przeprowadza się taki trening na dużym korpusie tekstowym, można otrzymać zbiór wstępnie wytrenowanych reprezentacji, użytecznych w wielu zadaniach NLP — tłumaczeniu maszynowym, odpowiadaniu na pytania, podsumowaniach treści, ocenie poprawności gramatycznej i wielu innych.

24.6. Obecny stan sztuki

Głębokie uczenie i uczenie transferowe wydatnie przyczyniły się do podniesienia wiedzy w zakresie NLP na wyższy poziom — do tego stopnia, że jeden z komentatorów w 2018 roku oświadczył: „Nadszedł moment ImageNet NLP” (Ruder, 2018). Stwierdzenie to stanowi zapewne wspomnienie punktu zwrotnego w historii sztucznej inteligencji, gdy w 2012 roku zastosowanie głębokiego uczenia w systemach widzenia komputerowego przyniosło tym systemom zaskakująco dobre osiągnięcia w konkursie ImageNet. Podobny przełom dało się zauważyć właśnie w 2018 roku, gdy odkryto zadziwiająco użyteczność uczenia transferowego w zadaniach NLP. Dostępny jest bowiem gotowy ogólny model języka, który można pobrać i dostosować do specyfiki konkretnego zadania.

Zaczął się od prostych systemów embeddingowych — WORD2VEC w 2013 roku i GloVe w 2014 roku. Badacze mogą pobierać gotowe modele tego typu bądź szybko trenować własne, i to bez konieczności wykorzystywania superkomputerów. Z drugiej strony, wstępnie wytrenowane reprezentacje kontekstowe są o rząd wielkości bardziej pracochłonne w treningu.

Modele te stały się realizowalne dopiero po szerokim upowszechnieniu nowych rozwiązań sprzętowych, głównie procesorów GPU i TPU, dzięki którym badacze mogli pobierać gotowe modele zamiast zużywać zasoby na trenowanie własnych. Model transformera pozwolił na efektywne trenowanie znacznie większych i znacznie głębszych niż wcześniejszej sieci neuronowych — tym razem dzięki rozwojowi raczej oprogramowania niż sprzętu. Od 2018 roku nowe projekty NLP zwykle rozpoczynają się od wstępnie wytrenowanego modelu transformera.

Chociaż te modele transformera zostały wytrenowane pod kątem przewidywania następnego słowa w tekście, sprawują się zaskakująco dobrze w innych zadaniach językowych. Model ROBERTA, po pewnym dostrojeniu, osiąga znakomite rezultaty w testach odpowiadania na pytania i czytania ze zrozumieniem (Liu i in., 2019b). Model językowy GPT-2, podobny do transformera, posiadający 1,5 miliarda parametrów, wytrenowany z użyciem tekstu internetowego o objętości ok. 40 gigabajtów, osiąga dobre wyniki w tak różnorodnych zadaniach, jak tłumaczenie maszynowe z francuskiego na angielski i odwrotnie, wyszukiwanie odwołań w zależnościach długodystansowych i odpowiadanie na pytania z zakresu wiedzy ogólnej, wszystko bez dostrajania do specyfiki konkretnego zadania. Jak ilustruje tabela 24.2, GPT-2 może generować dość przekonujący tekst na podstawie zaledwie kilku słów promptu.

Nowoczesny system NLP — ARISTO (Clark i in., 2019) — uzyskał wynik 91,6% w teście na poziomie ośmioklasisty (patrz tabela 24.3). System ten składa się z zespołu solverów, z których część wyszukuje informację (podobnie do wyszukiwania w sieci Web), część realizuje wynikanie tekstowe i wnioskowanie jakościowe, niektóre natomiast używają dużych modeli transformerowych. Okazało się, że sam model ROBERTA osiągnął w tym teście wynik 88,2%. W treści bardziej zaawansowanej, na poziomie maturzysty, ARISTO osiąga 83%; wynik 65% oznacza „spełnia normy” (*meeting the standards*), a wynik 85% „spełnia normy z wyróżnieniem” (*meeting the standards with distinction*).

TABELA 24.3. Przykład testu szkolnego na poziomie ośmioklasisty, który system ARISTO potrafi poprawnie rozwiązać wykorzystując zespół metod, z modelem ROBERTA na czele. Poprawne rozwiązanie takiego testu wymaga znajomości języka naturalnego, posiadania tzw. wiedzy ogólnej i elementarnej wiedzy naukowej oraz zasad konstruowania testów tego rodzaju

1. Za pomocą czego najlepiej rozdzielić mieszaninę opiłków żelaza i czarnego pieprzu?

(a) magnesu (b) papierowego filtru (c) wagi laboratoryjnej⁷ (d) woltomierza

2. Jaką formę energii wytwarza drgająca gumka-recepturka?

(a) chemiczną (b) świetlną (c) elektryczną (d) akustyczną

3. Miedź jest metalem, zatem

(a) jest cieczą w temperaturze pokojowej (b) nie reaguje z innymi substancjami
(c) jest kiepskim przewodnikiem prądu (d) jest dobrym przewodnikiem ciepła

4. Który proces zachodzący w jabłoni jest najważniejszym skutkiem podziału komórkowego?

(a) wzrost (b) fotosynteza (c) wymiana gazowa (d) wydalanie produktów przemiany materii

ARISTO ma jednak ograniczenia: radzi sobie tylko z testami wyboru spośród gotowych odpowiedzi, nie z zadaniami „otwartymi” (opisowymi), nie potrafi też czytać ani generować diagramów⁸.

System T5 (*Text-to-Text Transfer Transformer*) przeznaczony jest do generowania odpowiedzi tekstowych na wejściowe dane tekstowe różnego rodzaju. Zawiera standardowy model kodera-dekodera transformerowego, wstępnie wytrenowany na 35 miliardach słów pochodzących z 750-gigabajtowego korpusu C4 (*Colossal Clean Crawled Corpus*). Ten nieetykietowany trening zaprojektowany został w celu wyposażenia modelu w wiedzę językową, którą można generalizować i która będzie przydatna do wielu specyficznych zadań. T5 jest następnie trenowany dla każdego zadania z użyciem danych tekstowych, których każda pozycja ma postać

nazwa zadania : treść

na przykład polecenie

translate English to German : That is good

daje w wyniku *Das ist gut*. W przypadku niektórych zadań dane wejściowe są znakowane, na przykład w tekście *Winograd Schema Challenge* zadanie polegało na określeniu podmiotu, do którego odnosi się zaimek *they* w zdaniu

*referent : The city councilmen refused the demonstrators a permit because **they** feared violence*

(„radni miejscy odmówili demonstrantom zezwolenia, ponieważ obawiali się przemocy”). Odpowiedzią jest *councilmen* („radni”), ale już w bardzo podobnym zdaniu

*referent : The city councilmen refused the demonstrators a permit because **they** advocated violence*

(„radni miejscy odmówili demonstrantom zezwolenia, ponieważ ci nawoływali do przemocy”) *they* odnosi się do „demonstrantów”. Różnica ta nie wynika w żaden sposób z gramatyki języka, lecz z ogólnej wiedzy o otaczającym świecie

⁷ W oryginale *triple beam balance* — „waga przesuwnikowa jednoszalkowa trójramienna” — *przyp. tłum.*

⁸ Zauważono, że w przypadku niektórych testów wielokrotnego wyboru możliwe jest uzyskanie dobrego wyniku nawet *bez związku z pytaniami*, bo w błędnych wariantach wyboru mogą skrywać się elementy demaskujące ich zwodniczy charakter (Gururangan i in., 2018). Wydaje się, że dotyczy to również odpowiedzi na pytania wizualne (Chao i in., 2018).

— to nie radni chcieli manifestować poparcie dla przemocy, skoro zgody na manifestację nie wydali, zresztą, jak pokazuje doświadczenie, jeśli przemoc ma miejsce, to rodzi się przeważnie ze spontanicznych manifestacji, a nie odgórnych zarządzeń władz municypalnych, a w ogóle w zdaniu jest mowa o demonstrantach. Komputerom jednak z natury obca jest wiedza tego rodzaju, dlatego wzmiankowany test postrzegany jest jako najważniejsza alternatywa dla oryginalnego testu Turinga⁹.

Na polu ulepszania systemów NLP wciąż wiele pozostaje do zrobienia. Jednym z podstawowych problemów większości modeli transformerów jest ograniczoność kontekstu — opierają się one na wąskim kontekście, nieprzekraczającym kilkuset słów. Próba złagodzenia tego limitu jest przedmiotem kilku podejść eksperymentalnych, między innymi system Reformer (Kitaev i in., 2020) zdolny jest do obsługi kontekstu o długości do miliona słów.

Ostatnie wyniki eksperymentów pokazały, że użycie większej ilości danych treningowych owocuje lepszymi modelami — na przykład model ROBERTA osiągnął imponujące wyniki po wytrenowaniu na zbiorze liczącym 2,2 biliona słów. A skoro tak zachęcające wyniki uzyskuje się dzięki danym wyłącznie tekstowym, to co można by osiągnąć, sięgając po dane innego rodzaju — ustrukturyzowane bazy danych, dane liczbowe, obrazy i wideo? Niewątpliwie do efektywnego spożytkowania tego panoptikum informacji konieczny byłby kolejny przełom w możliwościach sprzętu obliczeniowego, głównie w szybkości procesorów i pojemności pamięci. Stosowny przełom musiałby się także dokonać w metodologiach sztucznej inteligencji.

Wnikliwi Czytelnicy zapewne zapytają w tym momencie, z nutką zdziwienia, dlaczego w poprzednim rozdziale tak szczegółowo omawialiśmy gramatyki, analizę syntaktyczną i interpretację semantyczną, skoro w świetle treści niniejszego rozdziału deprecjonujemy te koncepcje na rzecz modeli opartych wyłącznie na danych? Szybka odpowiedź powinna eksponować bezsprzeczne zalety systemów sterowanych danymi: łatwość budowania i rozwijania, wspaniałe osiągi w standardowych benchmarkach — pod tym względem ustępują im wyraźnie systemy budowane ręcznie, przy użyciu rozsądnego wysiłku ludzkiego, oparte na podejściach opisywanych w rozdziale 23. Może być tak dlatego, że modele transformerów i ich pochodne uczą się ukrytych reprezentacji wychwytyjących te same podstawowe idee co gramatyki i informacje semantyczne. Nie można też wykluczyć, że w tych ogromnych modelach dzieje się coś zupełnie innego, o czym po prostu (jeszcze) nie wiemy. Wiemy natomiast, że system wytrenowany przy użyciu danych tekstowych jest łatwiejszy w utrzymaniu i przystosowywaniu do nowych domen i nowych języków naturalnych niż system, który opiera się na ręcznie definiowanych cechach.

Niewykluczone, że przyszłe przełomy w jawnym modelowaniu gramatycznym i semantycznym spowodują odwrócenie tej tendencji, bardziej jednak prawdopodobne jest ukształtowanie się podejść hybrydowych, integrujących najlepsze koncepcje z obu wymienionych światów. Kitaev i Klein (2018) wykorzystali mechanizm uwagi do ulepszenia tradycyjnego parsera składnikowego, osiągając najlepszy wynik, jaki zarejestrowano kiedykolwiek przy użyciu zestawu testowego Penn Treebank. Podobnie Ringgaard i in. (2017) pokazują, w jaki sposób można ulepszyć parser zależności za pomocą embeddingów i rekurencyjnej sieci neuronowej. Ich system — SLING — przeprowadza parsowanie bezpośrednio do reprezentacji w postaci ramek semantycznych, łagodząc problem kumulowania się błędów w tradycyjnym systemie potokowym.

Systemy NLP, mimo niekwestionowanych osiągnięć, wciąż ustępują człowiekowi pod względem wydajności w wielu zadaniach, chociaż każdy z nich trenowany jest przy użyciu takiej objętości tekstu, jakiej człowiek nie byłby w stanie przeczytać w ciągu całego swojego życia. Stanowi to inspirujące pole do nowych pomysłów, odkryć i wynalazków dla językoznawców, psychologów i oczywiście badaczy NLP.

⁹ Patrz <https://pol.phonocardsoftware.com/s2468-winograd-schema-challenge-a-turing-test-alternative-puts-ai-programs-to-shame> — mimo nieco dziwnego (automatycznego) tłumaczenia tekst jest całkiem zrozumiały — *przyp. tłum.*

Podsumowanie

W tym rozdziale omawialiśmy następujące kluczowe aspekty związków głębokiego uczenia z przetwarzaniem języka naturalnego:

- Reprezentacje ciągłe słów, z embeddingami na czele, są bardziej użyteczne od dyskretnych reprezentacji atomowych, można je bowiem trenować przy użyciu nieetykietowanych danych tekstowych.
- Lokalny i długodystansowy kontekst można efektywnie modelować za pomocą rekurencyjnych sieci neuronowej, przechowując w jej ukrytych wektorach relewantne informacje.
- Modele „sekwencja na sekwencję” okazują się efektywne w rozwiązywaniu problemów tłumaczenia maszynowego i generowania tekstu.
- Modele transformerów wykorzystujące samouagę mogą być wykorzystywane do modelowania kontekstu (lokalnego i długodystansowego) z dużą efektywnością dzięki możliwości wykorzystywania operacji macierzowych implementowanych w nowoczesnym sprzęcie obliczeniowym.
- Uczenie transferowe wykorzystujące wytrenowane wstępnie embeddingi słów umożliwia budowanie modeli w oparciu o ogromne nieetykietowane korpusy i wykorzystywanie tych modeli do rozwiązywania różnorodnych zadań. Modele wytrenowane wstępnie do przewidywania brakujących słów okazują się użyteczne także w realizacji innych zadań, takich jak odpowiadanie na pytania i wynikanie tekstowe, po dostrojeniu do specyfiki docelowych domen.

Bibliografia i uwagi historyczne

Rozkład słów i fraz w języku naturalnym jest zgodny z **prawem Zipfa** (Zipf, 1935, 1949), zgodnie którym słowo zajmujące n -tą pozycję w rankingu występowania pojawia się z częstotliwością w przybliżeniu odwrotnie proporcjonalną do n . Oznacza to, że mamy problem ze słowami rzadko występującymi: analizując nawet ogromne, liczone w miliardach słów dane treningowe wciąż można napotykać nowe słowa i frazy nienapotkane wcześniej w tych danych.

Generalizację na nowe słowa i frazy ułatwiają reprezentacje hołdujące zasadzie, że słowa o podobnym znaczeniu pojawiają się najczęściej w podobnych kontekstach. Deerwester i in. (1990) dokonali rzutowania słów na wektory niskowymiarowe, dekomponując macierz współwystępowania utworzoną na podstawie słów oraz dokumentów, w których słowa te występowały. Inną możliwością jest potraktowanie otaczających słów, w odległości ograniczonej oknem o szerokości (powiedzmy) pięciu słów, jako kontekstu tychże słów.

Brown i in. (1992) pogrupowali słowa w hierarchiczne klastery zgodnie z kontekstem bigramowym tych słów, co okazało się skuteczne w przypadku zadań takich jak rozpoznawanie nazwanych encji (Turian i in., 2010). System WORD2VEC (Mikolov i in., 2013) był pierwszą znaczącą demonstracją zalet embeddingów budowanych w procesie trenowania sieci neuronowych. Embeddingi systemu GloVe (Pennington i in., 2014) uzyskano, operując bezpośrednio na macierzy współwystępowania słów zbudowanej na podstawie miliardów słów tekstu. Levy i Goldberg (2014) wyjaśniają, skąd bierze się zdolność embeddingów do odzwierciedlania regularności językowych.

Bengio i in. (2003) to autorzy pionierskiego pomysłu związanego z zastosowaniem sieci neuronowych do modeli językowych, polegającego na (cyt.) „połączeniu (1) rozproszonej reprezentacji każdego słowa oraz (2) funkcji prawdopodobieństwa dla sekwencji słów, wyrażonej w kategoriach tych reprezentacji”. Mikolov i in. (2010) zademonstrowali zastosowanie sieci RNN do modelowania lokalnego kontekstu w modelach językowych. Jozefowicz i in. (2016) pokazali, w jaki sposób sieć RNN wytrenowana na miliardzie słów może przewyższyć wydajnością starannie zaprojektowane „ręczne” modele n -gramowe. Kontekstowe reprezentacje słów zostały szczegółowo opisane przez Petersa i in. (2018), którzy opatrzyli je akronimem ELMO (*Embeddings from Language Models*).

Należy w tym miejscu zaznaczyć, że niektórzy autorzy porównują modele językowe pod względem ich **powikłania** (ang. *perplexity*). Powikłanie rozkładu prawdopodobieństwa wynosi 2^H , gdzie H jest entropią rozkładu (patrz sekcja 19.3.3). Spośród dwóch modeli lepszy jest ten o mniejszym powikłaniu, pod warunkiem identyczności pozostałych parametrów. W praktyce jednak bardzo rzadko daje się zapewnić ową „identyczność”, dlatego bardziej miarodajnym kryterium porównawczym jest wydajność modeli w zastosowaniu do rzeczywistych zadań.

Howard i Ruder (2018) opisują framework ULMFiT (ang. *Universal Language Model Fine-tuning*), ułatwiający dostrajanie wstępnie wytrenowanego modelu językowego, bez konieczności posługiwania się dużym zbiorem dokumentów z domeny docelowej. Artykuł Rudera i in. z 2019 roku to tutorial na temat zastosowań uczenia transferowego w zadaniach NLP.

Mikolov i in. (2010) przedstawili ideę wykorzystania sieci RNN w zadaniach NLP, a Sutskever i in. (2015) wprowadzili ideę uczenia „sekwencja na sekwencję” z wykorzystaniem głębokich sieci neuronowych. Zhu i in. (2017) oraz Liu i in. (2018b) wykazali, że podejście nienadzorowane faktycznie działa i znacznie ułatwia kolekcjonowanie danych. Bardzo szybko okazało się, że tego rodzaju modele potrafią zaskakująco dobrze radzić sobie z różnymi zadaniami, na przykład tworzeniem podpisów pod obrazami (Karpathy i Fei-Fei, 2015; Vinyals i in., 2017b).

Devlin i in. (2018) wykazali, że modele transformerów wstępnie wytrenowane pod kątem języka maskowanego mogą być bezpośrednio wykorzystywane do wielu zadań. Opisujący model opatrzyli oni akronimem BERT (ang. *Bidirectional Encoder Representations from Transformers*); wstępnie wytrenowane modele tego typu można dostosowywać do określonych domen i określonych zadań, w tym do odpowiadania na pytania, rozpoznawania nazwanych encji, klasyfikowania tekstu, analizy sentymentu i wnioskowania w języku naturalnym.

System XLNET (Yang i in., 2019) to ulepszenie modelu BERT, eliminujące rozbieżność między treningiem wstępnym a dostrajaniem. Framework ERNIE 2.0 (Sun i in., 2019) wydobywa więcej informacji z danych treningowych, dzięki uwzględnianiu kolejności zdań i obecności nazwanych encji, a nie tylko współwystępowania słów; wykazano, że przewyższa on zarówno BERT-a, jak i XLNET. Odpowiedzią twórców BERT-a było jego ponowne przeanalizowanie i ulepszenie — nowy system ROBERTA (Liu i in., 2019b) wykorzystywał więcej danych, różne hiperparametry oraz rozmaite procedury treningowe, dzięki czemu zdołał dorównać wydajnością systemowi XLNET. System Reformer (Kitaev i in., 2020) rozszerza zakres uwzględnianego kontekstu aż do miliona słów. Tymczasem ALBERT (ang. *A Lite BERT*) to krok w innym kierunku — zmniejszenia liczby parametrów ze 108 mln do 12 mln, by system zmieścił się na urządzeniu mobilnym, bez uszczerbku dla (wysokiej) dokładności.

System XLM (Lample i Conneau, 2019) to model transformera z danymi treningowymi w wielu językach. Wielojęzyczność ta jest oczywiście użyteczna w zadaniach tłumaczenia maszynowego, ale zapewnia również bardziej niezawodne reprezentacje w zadaniach jednojęzycznych. W niniejszym rozdziale opisałiśmy dwa inne ważne systemy, GPT-2 (Radford i in., 2019) oraz T5 (Raffel i in., 2019); w drugim z wymienionych artykułów opisano również korpus C4 (*Colossal Clean Crawled Corpus*) z 35 miliardami słów.

Zaproponowano różne obiecujące ulepszenia algorytmów wstępnego treningu (Yang i in., 2019; Liu i in., 2019b). Wstępnie wytrenowane modele kontekstowe opisywane są w artykułach Petersa i in. (2018) oraz Dai i Le (2016).

Benchmark GLUE (ang. *General Language Understanding Evaluation*) to zbiór wzorcowych zadań i narzędzi do ewaluacji systemów NLP; został skonstruowany przez Wang i in. (2018a). Wspomniane zadania obejmują odpowiadanie na pytania, analizę sentymentu, wynikania tekstowego, tłumaczenie maszynowe i parsowanie. Modele transformerów zdominowały ranking systemów NLP, przeciętne ludzkie możliwości porównywalne są z dziewiątą pozycją tego rankingu. GLUE zyskał godnego następcę o nazwie SUPERGLUE (Wang i in., 2019), którego wzorcowe zadania zaprojektowane zostały jako trudniejsze dla komputerów, lecz łatwe jak dotąd dla ludzi.

U schyłku 2019 roku T5 był ogólnym liderem wspomnianego rankingu, uzyskując wynik 89,3 — zaledwie 0,5 poniżej linii bazowej 89,8. W przypadku trzech z dziesięciu zadań T5 okazał się lepszy od człowieka. Pierwsze z tych zadań polegało na udzieleniu odpowiedzi typu „tak/nie” („Czy Francja należy do tej samej strefy czasowej co Wielka Brytania?”), dwa pozostałe sprowadzały się do czytania ze zrozumieniem, czyli przeczytania (odpowiednio) akapitu publikacji i artykułu prasowego, a następnie odpowiedzi na pytania związane z ich treścią.

Tłumaczenie maszynowe to główne zastosowanie modeli językowych. Już w 1933 roku Petr Troyanskii otrzymał patent na „maszynę tłumaczącą”, ale była to jedynie idea, której nie sposób było urzeczywistnić, z prostego powodu — nie było jeszcze komputerów. W 1947 Warren Weaver, opierając się na pracach z kryptografii i teorii informacji, napisał do Norberta Wienera:

Gdy patrzę na artykuł napisany w języku rosyjskim, myślę sobie: „Ten artykuł napisano w języku angielskim, ale został on zakodowany w postaci dziwacznych symboli, które muszą zdekodować”.

Owa koncepcja „tłumaczenia przez dekodowanie”, sama z siebie interesująca, okazała się jednak nierealizowalna w praktyce, ze względu na brak wystarczających danych i możliwości ówczesnej infrastruktury obliczeniowej.

W latach siedemdziesiątych ubiegłego wieku ten stan rzeczy zaczął się to zmieniać, a system SYSTRAN (Toma, 1977) był pierwszym udanym komercyjnym systemem tłumaczenia maszynowego. Opierał się na regułach leksykalnych i gramatycznych opracowanych ręcznie przez lingwistów, a także na danych treningowych. W latach 80. społeczność NLP przyjęła modele czysto statystyczne, oparte na częstotliwości występowania słów i fraz (Brown i in., 1988; Koehn, 2009). Gdy zbiory treningowe osiągnęły miliardy lub biliony tokenów (Brants i in., 2007), pojawiły się systemy, które dawały wyniki zrozumiałe, ale jednak dalekie od płynności (Och i Ney, 2004; Zollmann i in., 2008). Och i Ney (2002) pokazują, jak trenowanie dyskryminatywne doprowadziło do postępów w tłumaczeniu maszynowym na początku XXI wieku.

Sutskever i in. (2015) po raz pierwszy wykazali, że jest możliwe uczenie kompleksowego modelu neuronowego „sekwencja na sekwencję” pod kątem tłumaczenia maszynowego. Bahdanau i in. (2015) wykazali zalety modelu, który uczy się dopasowywać zdania w języku źródłowym i docelowym oraz dokonywać tłumaczeń między tymi językami. Vaswani i in. (2018) wykazali, że neuronowe systemy tłumaczenia maszynowego można dalej ulepszać poprzez zastąpienie pamięci LSTM architekturą transformerów, wykorzystującą mechanizm uwagi do wychwytywania kontekstu. Te neuronowe systemy tłumaczenia szybko zdystansowały statystyczne metody oparte na frazach, a architektura transformerów wkrótce rozprzestrzeniła się na inne zadania NLP.

Wielce pomocnym w badaniach nad problemami **odpowiadania na pytania** okazało się stworzenie repozytorium SQuAD (ang. *Stanford Question Answering Dataset*)¹⁰ — pierwszego wielkoskalowego zestawu danych do trenowania i testowania takich systemów (Rajpurkar i in., 2016). Od tego czasu opracowano szereg modeli głębokiego uczenia dedykowanych temu zadaniu (Seo i in., 2017; Keskar i in., 2019). System ARISTO (Clark i in., 2019) wykorzystuje głębokie uczenie w połączeniu z zestawem innych taktyk. Począwszy od 2018 roku większość modeli odpowiadania na pytania korzysta z wytrenowanych wstępnie reprezentacji językowych, co prowadzi do zauważalnej poprawy w stosunku do wcześniejszych systemów tej kategorii.

Wnioskowanie na podstawie języka naturalnego polega na ocenie, czy dana hipoteza (*dogs need to eat* — „psy potrzebują pożywienia”) wynika z określonej przesłanki (*all animals need to eat* — „wszystkie zwierzęta potrzebują pożywienia”). Zadanie to stało się szczególnie popularne dzięki konkursom PASCAL Challenge (Dagan i in., 2005) i obecnie dostępne są dedykowane mu zbiory danych na dużą skalę (Bowman i in., 2015; Williams i in., 2018). Szczególnie dobrą wydajność zapewniają systemy oparte na wstępnie wytrenowanych modelach, takich jak ELMo i BERT.

Konferencja CoNLL (ang. *Conference on Computational Natural Language Learning*) koncentruje się — zgodnie z nazwą — na problemach uczenia maszynowego związanych z ogólnie pojętym przetwarzaniem języka naturalnego. Także wszelkie czasopisma i konferencje wspomniane w rozdziale 23. obejmują tematykę głębokiego uczenia, które obecnie jest dominującym kierunkiem badań na polu NLP.

¹⁰ Patrz <https://rajpurkar.github.io/SQuAD-explorer/> — przyp. tłum.

SKOROWIDZ

A

abstrakcja stanów, 193
adaptacyjne programowanie
 dynamiczne, 163, 189
addytywna dekompozycja, 181
agent
 uczenie, 7
 aktywne, 167
 aktywne różnic czasowych, 172
 na bazie modelu, 189
 niezależnie od modelu, 189
 pasywne, 161
 zachłanny, 168
aktuatory, 313, 362, 404
 hydrauliczne, 313
aktywne dalmierze optyczne, 311
algorytm
 „wewnątrz-zewnątrz”, 215
 ADABOOST, 64
 asymptotycznie optymalny, 334
 CYK, 211, 212, 227
 dowolnej chwili, 411
 dwukierunkowy, 112
 EM, 104, 112
 filtrowania cząstek, 319
 głębokiego uczenia, 134
 gradientowego TD, 192
 ignorowania nimaksymalności, 282
 klasy PAC, 31
 konstruowania modelu, 25
 pasywnego agenta, 164, 166
 probabilistyczny filtrowania, 363
 RRT, 333
 typu „poza polityką”, 174
 typu „według polityki”, 174

uczenia, 34
uczenia oparty na drzewie
 decyzyjnym, 17
 Viterbiego, 204
analiza
 błędów, 392
 drzewa błędów, 392
 sentymetu, 240
apertura kamery, 263
aproksymowanie
 funkcji, 174
 uczenia, 176
architektura
 odruchowa, 412
 oprogramowania, 368
 subsumpcyjna, 357
 sztucznej inteligencji, 411
atak kontrydaktoryjny, 139
atrybuty testowe, 19
augmentacja, 216, 227
 zbioru danych, 279
autokoder, 147
 wariacyjny, 148
autolokalizacja z jednoczesnym
 mapowaniem, 321
automatyczne różniczkowanie, 124
automatyzacja procesów
 biznesowych, 390
autonomiczne odkurzacze, 309

B

badania eksploracyjne danych, 72
bagging, 60, 80, 83
barwy podstawowe, 267
bayesowska regresja liniowa, 98

bayesowskie uczenie
 parametrów, 96
 ze wzmacnianiem, 171
bezpieczeństwo
 AI, 369, 401
 eksploracji, 170, 347
bezpieczna agregacja, 382
bezpośrednie estymowanie
 użyteczności, 189
bezrobocie technologiczne, 389
bierna percepcja, 260
big data, 408
bikony, 312
błąd pominięcia przykładu, 61
boosting, 63, 80
 gradientowy, 66
bramka
 wejściowa, 144
 wyjściowa, 144
 zapominania, 144
broń autonomiczna, 377, 400

C

całka po ścieżce, 335
CART, Classification And Regression
 Trees, 23
cechy, 71, 260
 obrazów, 268
certyfikacja, 387
certyfikowane usuwanie danych, 400
chwytyki trójpalczaste, 313
cienie, 265, 267
CNN, konwolucyjna sieć
 neuronowa, 128, 302,
 Patrz także sieci neuronowe
 czasy gramatyczne, 222
 częściowa obserwowalność, 362

czujniki, 262, 362, 404
 aktywne, 311
 bezwładnościowe, 313
 dotykowe, 312
 lokalizacji, 312
 pasywne, 311
 proprioceptywne, 313

D

dalmierze, 311
 dane
 kompletne, 90
 liniowo separowalne, 79
 data science, 408
 decyzja
 negatywna, 14
 pozytywna, 14
 deepfake, 297
 deidentyfikacja, 380
 deklaracja zgodności, 401
 dekodery obrotu, 313
 dekodowanie, 245
 zachłanne, 245
 dekompozycja komórkowa, 363
 demonstracja poklatkowa, 355
 deskryptory punktów, 305
 detektor
 kolizji, 331
 obiektów, 280
 diagramy Woronoja, 329, 366
 długotrwała pamięć krótkoterminowa,
 LSTM, 144, 241
 dokonywanie przewidywań, 350
 dolna granica dowodowa, 148
 domeny zastosowań robotyki, 358
 dopasowywanie cech, 185
 dostrajanie hiperparametrów, 29
 dotykowe sprzężenie zwrotne, 367
 downsampling, 131
 dropout, 141
 drzewo
 ekstremalnie randomizowane, 61
 decyzyjne, 14, 15, 79, 80
 przycinanie, 20
 zastosowania, 22
 k-d, 50

losowe
 błyskawicznie eksplorujące, 332
 regresyjne, 23
 rozbioru, 219
 syntaktycznego, 210
 dynamika odwrotna, 338
 dyskryminator, 149
 dysproporcja liczebności próbek, 385
 dywergencja Kullbacka-Leiblera, 126

E

efektory, 308
 efektywna liczba punktów, 106
 eksploracja, 168, 170
 eksploracyjne badanie, 72
 ekstradrzewa, 61
 ekstrakcja informacji, 231
 elementy inżynierii oprogramowania, 75
 embedding
 pozycyjny, 247
 słów, 153, 201, 233
 znaczeniowy, 251
 entropia, 19
 krzyżowa, 126
 wzajemna, 126
 epoka, 39
 estymowanie
 bezpośrednio użyteczności, 163
 gęstości, 90
 okien Parzena, 116
 prawdopodobieństwa, 102
 swego stanu, 315
 użyteczności, 175
 etykieta, 9

F

falszywe pominięcie, 74
 fałszywy alarm, 74
 Faster RCNN, 280, 281
 fiksacja, 284
 filtr
 gaussowski, 270
 Kalmana, 321
 rozszerzony Kalmana, 321

filtrowanie cząstek
 z Rao-Blackwellizacją, 366
 formowanie obrazów, 261
 frameworki robotyczne, 356
 funkcja
 aktywacyjna, 119
 eksploracji, 169
 ewaluacyjna, 174
 jakości, 161
 jądra, 53, 58, 103
 liniowa, 35
 logistyczna, 82
 najlepiej dopasowana, 10
 nagród, 406
 progowa, 44
 Q, 161
 regularyzacji, 29
 ReLU, 119
 sigmoidalna, 119
 softplus, 120
 straty, 27, 79, 125
 twardego progu, 46

G

GAN, generatywna sieć
 kontradiktoryjna, 149
 GBM, Gradient Boosting Machine, 83
 generalizacja, 20, 137, 174, 354, 381
 stosowa, 62
 generator, 149
 generowanie i preparowanie obrazów,
 295, 298
 geometria 3D, 282, 293, 302
 głęboka, 285
 ostrości, 265
 głęboki
 model autoregresywny, 149
 system uczenia ze wzmacnianiem,
 177
 GPS różnicowy, 312
 GPS, Global Positioning System, 312
 gra o niepełnej informacji, 349
 gradient, 66
 eksplodujący, 144
 empiryczny, 182
 polityki, 182

- graf
 obliczeniowy, 121, 124, 135
 przepływu danych, 121
 w pełni połączony, 122
 widoczności, 328, 366
 Woronoja, 329
- gramatyka, 196, 208
 augmentowana, 216
 bezkontekstowa, 227
 leksykalno-funkcjonalna, 229
 nadgeneratywna, 209
 semantyczna, 220
 zależnościowa, 213, 229
- granica decyzyjna, 43
 gry asystenckie, 394
- ## H
- haszowanie wrażliwe na lokalizację, 51
 heurystyka zysku informacyjnego, 79
 hierarchiczne uczenie ze
 wzmocnieniem, 178, 190, 406
 hiperparametry, 24, 29
 hipoteza, 10
 generalizująca, 10
 niedouczona, 11
 o maksymalnej wiarygodności, 89
 PAC, 31
 pasująca do danych, 11
 spójna, 10
 w przybliżeniu poprawna, 32
 zerowa, 21
- histogram pola wektorowego, 367
 hodometr, 313
- ## I
- identyfikacja w granicy, 81
 ignorowanie niemaksymalności, 282
 iloczyn skalarny, 204
 ILQR, Iterative Linear Quadratic
 Regulator, 343, 363
 indukcja, 9
 gramatyki, 230
 interakcje, 362
 interpolacja, 26
 interpretacja semantyczna, 218
- interpretowalność, 75
 inżynieria
 bezpieczeństwa, 392
 cech, 71
 oprogramowania, 75, 392
 sztucznej inteligencji, 414
- ## J
- jądro wielomianowe, 58
 jednostka, 119
 język naturalny, 195
 gramatyka, 208
 gramatyka augmentowana, 216
 komplikacje, 221
 parsowanie, 210
- ## K
- kamera, 285
 monitorująca, 380
 otworkowa, 262
 ToF, 311
- kamienie milowe, 332
 k-anonimowość, 381
 katastrofalne zapominanie, 176
 kategorie syntaktyczne, 208
 kinematyka odwrotna, 326
 klasa modeli, 10
 klasteryzacja nienadzorowana, 105
 klasyfikacja, 9, 79
 boolowska, 14
 liniowa
 z regresją logistyczną, 45
 z twardym progim, 43
 z użyciem sieci RNN, 240
 za pomocą maszyny SVM, 54
 klasyfikatory obrazów, 275, 302
 klonowanie zachowań, 192, 354
 kodeksy etyczne, 400
 kodowanie
 danych wejściowych, 124
 z gorącą jedynką, 71
 kojarzenie danych, 321
 kolekcja
 danych, 69
 zdań, 202
- kolokacja bezpośrednia, 342
 kolory, 267
 komórka, 330
 pamięci, 144
 komunikacja, 195
 kontaminacja danych, 62
 kontekst znaczeniowy, 251
 kontroler
 proporcjonalno-całkowo-
 różniczkowy, 340
 proporcjonalno-różniczkowy, 340
 proporcjonalny, 339
 reaktywny, 356
 ściśle stabilny, 340
- konwolucja, 128
 konwolucyjne sieci neuronowe, CNN,
 128, 302
 downsampling, 131
 klasyfikowanie obrazów, 277
 kumulacja, 131
 operacje tensorowe, 132
 warstwy, 129
- krawędzie, 268
 krzywa
 oceny jakości klasyfikatora, 74
 ROC, 74
 uczenia, 18, 79
- kształtowanie nagród, 178, 190
 kumulacja, 131
 maksymalizująca, 131
 uśredniająca, 131
 kwantyfikacja, 221
- ## L
- las losowy, 60, 83
 leksykon języka, 209
 lidary skanujące, 311
 linearyzacja, 318
 linia bazowa, 284
 liniowa logika temporalna, 407
 liniowy model Gaussa, 94
 listy decyzyjne, 33, 34
 lokalizacja, 316
 Markowa, 366
 Monte Carlo, 318, 319

lokalna regresja ważona, 79
 lokalne wyszukiwanie skupione, 245
 LOOCV, leave-one-out cross-validation,
 24

M

macierz
 bezwładności, 341
 pomyłek, 74
 wagowa, 121
 maksimum
 a posteriori, 114
 wiarygodności, 91, 114, 125
 maksymalizacja
 oczekiwań, 114
 wiarygodności logarytmicznej, 90
 manipulatory, 309, 362
 MAP, Maximum A Posteriori, 88
 mapa głębi, 293
 mapowanie, 316
 margines, 55
 maskowane modele językowe, 251
 maszyna
 boostingu gradientowego, GBM, 83
 wektorów nośnych, SVM, 54, 79, 82
 materiał dowodowy, 86
 metauczenie, 84, 367
 metawnioskowanie, 412
 metoda najbliższego sąsiedztwa, 79
 metody
 jądrowe, 79
 różnic czasowych, 189
 samonadzorowane, 323
 metonimia, 223
 metryka, 49
 mieszanka gaussowska, 106
 minimalna długość opisu, 81, 89
 miniwsad, 39
 model
 „sekwencja na sekwencję”, 242–245
 „worek słów”, 197, 228
 3D, 282
 akustyczny, 225
 autoregresywny, 149
 bayesowski naiwny, 93

bazowy, 59
 ciągły, 94
 danych, 10
 dynamiczny, 338
 dyskretny, 90
 dyskryminatywny, 94, 206
 Gaussa liniowy, 94
 generatywny, 94, 146
 językowy, 196, 225, 238
 maskowany, 252
 porównanie, 206
 k najbliższych sąsiadów, 48, 49, 82,
 103
 lasu losowego, 60
 mentalny, 224
 n-gramowy, 199
 wygładzanie, 200
 niejawni, 149
 nieparametryczny, 48, 79, 114
 obiektowy, 261
 regresyjny
 nieparametryczny, 52
 renderingowy, 261
 rzadki, 41
 świata, 224
 transformatora, 207
 ukryty Markowa, 111
 zespołowy, 59
 znakowy, 237

N

nadpróbkowanie, 71
 nagrody, 178
 naiwny model bayesowski, 93
 naśladownictwo, 354, 394
 nauka, 195
 nawigacja przybrzeżna, 345
 neokognitron, 156
 neuronauka obliczeniowa, 158
 n-gramy słów, 198
 niedopróbkowanie, 71
 niejednoznaczność
 leksykalna, 223
 syntaktyczna, 223
 znaczeniowa, 223

nienadzorowane
 tłumaczenia, 150
 uczenie transferowe, 145
 niepewność, 343
 niezmienniczość
 czasowa, 128
 przestrzenna, 128
 ze względu na przesunięcie, 305
 NLP, Natural Language Processing, 195
 norma, 49
 normalizacja, 49
 wsadowa, 136, 305

O

obciążenie, 11
 obiekty
 wykrywanie, 279
 obliczanie gradientów, 135
 obliczeniowa teoria uczenia, 31, 79
 obrazy
 dwuwymiarowe, 262
 generowanie, 295, 298
 klasyfikowanie, 275
 klasyfikowanie przy użyciu sieci
 CNN, 277
 krawędzie, 268
 preparowanie, 295
 przepływ optyczny, 273
 regiony, 274
 segmentowanie, 274
 tekstury, 272
 transformacja, 295
 oceny ryzyka, 383
 oczekiwanie cechy, 186
 odbicie
 lustrzane, 265
 rozproszone, 265
 wielokrotne, 267
 odbijalność, 267
 odbłaski, 266
 odległość
 Hamminga, 49
 Mahalanobisa, 49
 Minkowskiego, 49
 ze znakiem, 335, 363

odtworzenie doświadczeń, 177
 odwrotne uczenie ze wzmacnianiem, 407
 ognisko ekspansji, 285
 ogniskowa, 263
 ograniczenie cyklu, 296
 ograniczona optymalność, 413
 operacja
 konwolucji, 129
 tensorowa, 132
 optymalizacja, 24
 bayesowska, 30
 kwadratowa, 56
 polityki regionu zaufania, 193
 trajektorii, 334
 optymalność ograniczona, 413
 orientacja krawędzi, 271
 osadzanie w niższym wymiarze, 322

P

PAC, probably approximately correct, 31
 pamięć LSTM, 241
 parser, 214
 tablicowy, 210, 227
 parsowanie, 210
 A*, 230
 częściowo nadzorowane, 215
 nienadzorowane, 215
 ukierunkowane na dane, 215
 zależnościowe, 213
 parytet demograficzny, 383
 PD-kontroler, 340
 pentaptery, 309
 percepcja, 322
 czynna, 260
 dotykowa, 302
 robotów, 315, 316
 węchowa, 302
 perceptron, 44, 79
 Gamby, 155
 PID-kontroler, 340, 363
 P-kontroler, 339
 planista, 332
 planowanie
 hierarchiczne, 178
 kinematyczne, 334
 ruchu, 314, 323, 327, 363, 366
 randomizowane, 332, 363
 w warunkach niepewności, 343
 z wieloma zapytaniami, 332
 zadań, 314
 płaszczyzna ogniskowa, 265
 podziały zbioru przykładów, 16
 pojazdy kosmiczne, 309
 pola uogólniające, 381
 polaryzacja, 11
 pole
 pod krzywą ROC, 74
 recepcyjne, 131
 polityka
 wyszukiwanie, 181
 ponowne planowanie online, 363
 pooling
 ROI, 281
 uśredniający, 241
 postać normalna Chomsky'ego, 211, 227
 poszukiwanie architektury neuronowej, NAS, 139
 prawa dla robotów, 391, 402
 prawda podstawowa, 10
 prawdopodobieństwo
 a posteriori, 88
 a priori hipotez, 87
 nieinformatywne a priori, 99
 sprzężone a priori, 97
 prawo
 kosinusowe Lamberta, 267
 Zipfa, 257
 predykat, 219
 probabilistyczna
 analiza głównych składowych, 146
 gramatyka bezkontekstowa, 208, 229
 mapa drogowa, 332
 probabilistyczne algorytmy filtrowania, 363
 problem
 aproksymowanego najbliższego sąsiedztwa, 51
 bandyty, 169
 eksploracyjny, 171
 króla Midasa, 394
 odpowiedniości, 354
 przypisania kredytu, 178
 tragarza fortepianu, 327, 366
 wyrównywania wartości, 394
 wyszukiwania, 327
 proces
 Dirichleta, 116
 gaussowski, 30, 116
 jednorodny czasowo, 142
 regulacji, 339
 programowanie
 dynamiczne adaptacyjne, 163
 dynamiczne różniczkowe, 367
 różniczkowalne, 409
 wizualne, 355
 prompt, 207
 propagacja
 wsteczna, 123, 154
 wsteczna w czasie, 144
 prosty planista, 332
 próbkowanie
 skorelowane, 183
 w dół, 131
 prymitywy ruchu, 348
 prywatność, 400
 różnicowa, 381
 przeczesywanie priorytetowe, 167
 przegląd k najbliższych sąsiadów, 48
 przejrzystość, 387
 przekleństwo wymiarowości, 50, 82
 przeniesienie stylu, 297
 przenośnia, 224
 przepływ optyczny, 273
 przestrzeń
 cech, 57
 hipotez, 10
 konfiguracyjna, 324, 328, 363
 połączonych stanów, 179
 robocza, 324, 328
 ścieżek, 328
 wagowa, 37
 przeszkoda konfiguracyjna, 325
 przeszukiwanie losowe, 30
 przetwarzanie języka naturalnego, NLP, 195, 233
 mechanizmy uwagi, 243
 pamięć LSTM, 241
 sieci RNN, 238
 przeuczenie, 12, 20, 89

przewidywanie, 364
 czynnika ludzkiego, 315
 dotyczące robotów, 351
 ludzkich działań, 349, 368
 przycinanie, 61
 drzewa decyzyjnego, 20
 przyszłość pracy, 402
 pseudodoświadczenie, 167
 pseudonagrody, 178
 pseudoodwrotność, 41
 punkty
 fizyczne, 363
 odstające, 71
 orientacyjne, 317
 zafiksowania, 284

Q

quadcoptery, 309
 Q-uczenie, 172, 190

R

rachunek wariacyjny, 335
 racjonalność Boltzmannna, 185
 radar, 312
 ramki zakotwiczeń, 280
 randomizacja domeny, 346
 randomizowana większość ważona, 67
 region, 274
 zainteresowania, 281
 regresja, 79
 k najbliższych sąsiadów, 53
 liniowa, 35, 79, 82, 95
 bayesowska, 98
 jednowymiarowa, 35
 wielozmienna, 40
 logistyczna, 45, 79, 82, 204
 lokalna ważona, 53
 nieparametryczna, 52
 ramki ograniczającej, 282
 regularyzacja, 29, 41, 42, 66
 regulator liniowo-kwadratowy, 342, 363
 reguła
 łańcuchowa, 122
 sterowania, 338
 Widrowa-Hoffa, 175

rekonstrukcja, 261, 291
 ścieżki, 293
 rekurencyjne sieci neuronowe, RNN, 142
 diagram, 238
 dwukierunkowe, 240
 klasyfikacja, 240
 modele językowe, 238
 przetwarzanie języka naturalnego,
 238
 trenowanie, 142
 relaksacja grafów, 322
 repozytorium
 ImageNet, 69
 SQuAD, 259
 reprezentacja
 czynnikowa, 9
 dualna, 55
 reprezentowanie
 słów, 201
 stanu świata, 405
 RGB, Red, Green, Blue, 268
 RNN, rekurencyjna sieć neuronowa,
 142, *Patrz także* sieci neuronowe
 roboty, 308
 kroczące, 309
 mobilne, 309, 362
 przemysłowe, 310
 uczące się ludzkich oczekiwań, 352
 ROC, Receiver Operating
 Characteristic, 74
 rodzina Dirichleta, 97
 rozbieżność
 dwuoczną, 283
 kątową, 285
 rozkład
 gęstości prawdopodobieństwa, 96
 normalny Wisharta, 97
 t-Studenta, 72
 rozkład-mieszanka, 105
 rozmycie przez ruch, 263
 rozpoznawanie, 261
 mowy, 225
 równania
 Yule'a-Walkera, 149
 Eulera-Lagrange'a, 336
 normalne, 41, 149
 Riccatiego, 343

równe szanse, 383
 równy wpływ, 383
 różnice czasowe, 176
 różniczkowanie w trybie odwrotnym,
 124
 RRT*, 334
 RRT, Rapidly-exploring Random Trees,
 332
 ruch
 planowanie, 327, 343
 randomizowane planowanie, 363
 sterowanie, 337
 rytm konwolucji, 129
 rzutowanie
 ortograficzne, 265
 perspektywiczne, 263

S

samochód autonomiczny, 351, 360
 samouwaga, 246
 SARSA, 173
 segmentacja, 274
 selekcja modeli, 24, 73, 79, 114
 semantyczne interpretowanie tekstu,
 227
 semantyka kompozycyjna, 218
 separator
 liniowy, 43
 o największym marginesie, 54, 55
 siatka zajętości, 365
 sieci neuronowe, 117, 154
 konwolucyjne, CNN, 128, 302
 downsampling, 131
 klasyfikowanie obrazów, 277
 kumulacja, 131
 operacje tensorowe, 132
 warstwy, 129
 rekurencyjne, RNN, 142
 diagram, 238
 dwukierunkowe, 240
 klasyfikacja, 240
 modele językowe, 238
 przetwarzanie języka
 naturalnego, 238
 trenowanie, 142
 z dwoma wejściami, 121

- sieć
 bayesowska
 uczenie się struktur, 100
 generatywna kontraduktoryjna,
 GAN, 149
 Hopfielda, 157
 jako złożona funkcja, 119
 przyczynowa, 116
 rezydualna, 132
 propozycji regionów, 280
 ze sprzężeniem w przód, 118
 siłownik, 313
 skalowane rzutowanie ortograficzne,
 265
 skrót perspektywiczny, 261
 słownik, 201
 słowo, 211
 soczewka, 264
 softmax, 126, 181
 spadek gradientowy, 37, 122
 deterministyczny, 39
 stochastyczny, 39
 wsadowy, 39
 spłot, 128, 270
 sprzężenie
 w przód, 237, 341, 363
 zwrotne, 9, 341
 stacjonarność przykładów, 23
 stacking, 62
 stan
 dynamiczny, 338
 przekonań, 315
 wyboru, 180
 stereopsja dwuoczną, 283
 sterowanie
 adaptacyjne, 366
 obliczanego momentu obrotowego,
 341
 odporne, 366
 optymalne, 341, 363
 ruchem, 299, 337
 śledzeniem trajektorii, 363, 366
 stochastyczne maszyny Boltzmana, 157
 stochastyczność, 362
 stopa błędu, 23, 27
 stopnie swobody, 325

 strata
 empiryczna, 28
 uogólniona, 28
 strategia „dziel i zwyciężaj”, 16
 strojenie ręczne, 30
 strukturalne estymowanie MDP, 192
 strukturalny algorytm EM, 113
 struktury sieci bayesowskich, 100
 suma kwadratów różnic, 273
 superpiksele, 275
 SVM, Support Vector Machines, 54
 syntetyzowanie mowy, 225
 systemy
 dialogu wizualnego, 291
 podpisywania, 290
 soczewek, 264
 tagowania, 290
 uczenia maszynowego, 68, 75
 interpretowalność, 74
 kolekcje danych, 69
 selekcja modeli, 73
 trenowanie, 73
 utrzymanie i konserwacja, 77
 wdrożenie, 77
 wyjaśnialność, 74
 zarządzanie, 69
 zaufanie, 74
 sztuczna inteligencja, 369
 architektury, 411
 argument
 „chińskiego pokoju”, 374
 nieformalności, 370
 ułomności, 371
 bezpieczeństwo, 379, 392, 401
 bezrobocie technologiczne, 389
 bezstronność, 382
 broń autonomiczna, 377
 czasu rzeczywistego, 411
 etyka, 375
 inżynieria, 414
 komponenty, 404
 nadzór, 379
 obiekcje natury matematycznej, 371
 ogólna, 413
 perspektywy, 414
 pomiar, 372
 prywatność, 379, 400
 przejrzystość, 387
 qualia, 374
 silna, 369, 397
 słaba, 369, 397
 systemy, 397
 świadomość, 374, 398
 ucieleśniona, 194
 uczciwość, 382
 wyjaśnialna, 387
 zaufanie, 387, 401
 szybkość uczenia, 38

Ś
 śmiertcionośna broń autonomiczna, 377
 świadomość, 398
 świat
 3D, 282
 W, 328
 światło, 265
 otoczenia, 267
 strukturalne, 311

T
 tabela przeglądowa, 48
 tablica, 210
 czujników, 318
 tagowanie części mowy, 202, 204, 248
 tangens hiperboliczny, 120
 technologiczna osobliwość, 395
 teksel, 272
 tekstury, 272
 tensor, 132
 teoria
 jednorodnej zbieżności, 82
 odpornego sterowania, 172
 optymalności, 228
 sterowania, 337
 trójchromatyczna, 267
 uczenia maszynowego, 31
 zintegrowanej informacji, 398
 testy
 cech i danych, 78
 infrastruktury uczenia
 maszynowego, 78
 istotności, 21

testy
 monitorujące dla uczenia
 maszynowego, 78
 punktu dzielącego, 22
 Turinga, 372, 398
 tworzenia modelu, 78
 tłumaczenie maszynowe, 225, 242, 259
 trajektoria, 323
 optymalizacja, 334
 transfer z symulacji do rzeczywistości,
 367
 transformacja obrazów, 295
 transformer, 246
 jednowarstwowy, 248
 transhumanizm, 396
 treebank, 227
 trening, 24
 kontradyktoryjny, 354
 modelu, 73
 wstępny, 249
 kontekstu zdaniowego, 251
 trik jądrowy, 54, 58
 t-SNE, 72
 twardy próg, 43
 twierdzenie
 Mercera, 58
 o uniwersalnej aproksymacji, 119

U

uczciwość
 grupowa, 383
 indywidualna, 383
 poprzez niewiedzę, 383
 uczenie
 agenta, 407
 bayesowskie, 87, 114
 bez żalu, 67
 częściowo nadzorowane, 68
 funkcji akcja-użyteczność, 161
 głębokie, 117
 algorytmy, 134
 grafy obliczeniowe, 124
 przetwarzanie języka
 naturalnego, 152, 233
 widzenie komputerowe, 152
 interfejsów, 354

kinestetyczne, 355
 kompleksowe, 124, 153, 347
 maszynowe, 7
 budowanie systemów, 68
 interpretowalne, 84
 zautomatyzowane, 84
 na podstawie różnic czasowych, 165
 nadzorowane, 9, 79
 naśladowcze, 184
 nienadzorowane, 9, 145
 online, 66, 80, 83
 PAC, 31, 33, 82
 parametrów, 90
 parametrów metodą maksimum
 wiarygodności, 90, 94
 parsera, 214
 populacyjne, 30
 praktykanckie, 183, 190, 192, 408
 predykcyjne, 409
 przez naśladownictwo, 354
 robota
 nadzorowane, 322
 nienadzorowane, 322
 sfederowane, 381
 się
 dla nagród, 159
 funkcji kosztu, 353
 gramatyk semantycznych, 220
 polityk, 354
 preferencji, 315
 przez naśladownictwo, 394
 przez perceptron, 44
 struktur sieci bayesowskich,
 100, 108, 113
 słabe, 64
 słabo nadzorowane, 409
 statystyczne, 86
 stopniowe, 215
 transferowe, 150, 249, 408
 transferowe nienadzorowane, 145
 wielkoskalowe, 29
 wielozadaniowe, 150, 151
 z kompletnych danych, 90
 z ukrytymi zmiennymi, 104
 za pomocą
 mieszanek gaussowskich, 105
 ukrytych modeli Markowa, 111

ze wzmacnianiem, 9, 153, 159
 aktywne, 167
 bayesowskie, 171
 generalizacja, 174
 głębokie, 177
 hierarchiczne, 178, 406
 na bazie modelu, 160, 346
 niezależnie od modelu, 161
 odwrotne, 183, 394, 407
 pasywne, 161
 relacyjne, 192
 w robotyce, 188, 346
 zastosowania do gier, 187
 zespołowe, 59
 ukryte
 indeksowanie semantyczne, 228
 modele Markowa, 111
 utajona alokacja Dirichleta, 228
 uwaga, 243

W

walidacja krzyżowa, 24, 79, 81
 wariacyjna granica dolna, 148
 wariacyjny rozkład a posteriori, 148
 wariancja, 28
 hipotezy, 12
 warstwa, 117
 gęstości mieszanej, 127
 kumulacyjna, 131
 wyjściowa, 122
 zbiorcza, 131
 warstwy
 sieci CNN, 129
 ukryte, 122, 127
 wyjściowe, 125
 wartość polityki, 182
 wczesne zatrzymywanie, 21
 wektor
 nośny, 56
 wag, 204
 z gorącą jędynką, 234
 weryfikacja i walidacja, 387
 wiarygodność, 87
 logarytmiczna danych, 90, 107

- widzenie
komputerowe, 287
stereoskopowe, 283, 311
- wiedza wstępna, 9
- wielkość kroku, 38
- wielostrzałowość, 342
- wirtualne liczniki, 97
- wizualizacja danych, 72
- wnioskowanie na podstawie języka naturalnego, 259
- wskazówki 3D, 285
- wspinaczka, 38
- współczynnik
wzmocnienia, 339
zysku informacyjnego, 22
- współtrening, 232
- wstępny trening embeddingów, 249
- wybór
akcji, 406
architektury sieci, 137
atrybutów testowych, 19
- wydajność
eksploracyjnego agenta, 170
zachłannego agenta, 168
- wygaszanie wag, 140
- wygładzanie, 112, 270
1-addytywne, 228
modeli n-gramowych, 200
przez interpolację liniową, 201
- wyjaśnialność, 76
- wykrywanie obiektów, 279
- wymiar VC, 82
- wyodrębnianie cech, 261
- wyszukiwanie
gridowe, 30
informacji, 226, 231
polityki, 161, 181, 190
skupione, 206
zachłanne, 205
- ## Z
- zachłanne
dekodowanie, 245
wyszukiwanie, 205
- zachłanność, 169
- zachłanny agent, 168
- zachowanie robota, 315
- założenie o stacjonarności przykładów, 23
- zanegowany logarytm wiarygodności, 125
- zanikający gradient, 123
- zasada minimalnej długości opisu, 29
- zasady etyki robotów, 398
- zasoby, 409
- zastosowania robotyki
eksploracja, 361
ochrona zdrowia, 359
opieka domowa, 358
przemysł, 362
- rozrywka, 361
- samochody autonomiczne, 360
- usługi, 360
- zbiory wstępnie wytrenowane, 236
- zbiór
MNIST, 73
testowy, 10, 24
treningowy, 10, 14, 18, 24, 79
dwuwymiarowy, 57
ważony, 63
walidacyjny, 24
- zespół, 59
- zleksykalizowane PCFG, 229
- złożoność
próby, 32
algorytmiczna, 81
Kołmogorowa, 81
obliczeniowa, 28
- zmiennie
ukryte, 108, 113
utajone, 104
wskaźnikowe, 107
- znaczniki części mowy, 203
- znajdowanie
ciągłej ścieżki, 327
najbliższych sąsiadów, 50
- znikający punkt, 264
- znormalizowany przekrój, 275
- zupełność probabilistyczna, 332
- zysk informacyjny, 345
- ## ε
- ε-kula, 32

PROGRAM PARTNERSKI

— GRUPY HELION —

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

Zmień swoją stronę WWW w działający bankomat!

Dowiedz się więcej i dołącz już dzisiaj!

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

GRUPA
Helion 

Sztuczna inteligencja: dokąd zmierzasz, technologio?

Praktycznie codziennie korzystamy z osiągnięć sztucznej inteligencji (ang. Artificial Intelligence – AI). Mimo to jej potencjał wciąż jest zagadką: nie wiemy, gdzie leżą granice jej rozwoju i jakie jeszcze technologie przyniesie nam ta relatywnie młoda dziedzina nauki. Równocześnie niektóre zastosowania sztucznej inteligencji budzą niepokój i zmuszają do zadawania trudnych pytań. Jakakolwiek próba odpowiedzi wymaga jednak wiedzy o tym, czym w istocie jest AI i jakie są jej ograniczenia.

To drugi tom klasycznego podręcznika wiedzy o sztucznej inteligencji. Podobnie jak w wypadku pierwszej części, lektura tej książki nie wymaga wybitnej znajomości tematu. Dzięki przejrzystości tekstu i umiejętnemu unikaniu nadmiernego formalizmu można w dość łatwy sposób zrozumieć kluczowe idee i koncepcje nauki o AI. Najnowsze technologiczne osiągnięcia zostały pokazane na tle rozwijającej się wiedzy, również z innych dziedzin inżynierii. Sporo miejsca poświęcono zagadnieniom, które budzą wątpliwości. Mowa tu o wyrafinowanych technikach uczenia maszynowego, modelach językowych czy widzeniu komputerowym, a także o sprawach, które już dziś wymagają najwyższej troski: o etycznych aspektach sztucznej inteligencji, bezpieczeństwie związanych z nią technologii i jej perspektywach.

W drugim tomie:

- różne modele i koncepcje uczenia maszynowego
- przetwarzanie języka naturalnego i modele językowe
- widzenie komputerowe, w tym generowanie obrazów
- roboty: percepcja, działanie, uczenie
- perspektywy sztucznej inteligencji

Stuart Russell jest znany ze swojego wkładu w naukę o sztucznej inteligencji. Jest profesorem na Uniwersytecie Kalifornijskim w Berkeley i założycielem Center for Human-Compatible AI. Laureat wielu prestiżowych nagród i autor kilkuset publikacji.

Peter Norvig jest wybitnym akademikiem, pracuje w Stanford Institute for Human-Centered AI. Był szefem Wydziału Nauk Obliczeniowych NASA Ames Research Center i dyrektorem do spraw badań w Google.

Helion 	KOD KORZYŚCI Sięgnij po więcej! ▶ 
 helion.pl	ISBN 978-83-283-7773-8 
 HELION SA ul. Kościuszki 1c 44-100 Gliwice tel.: 32 230 98 63 helion@helion.pl	9 788328 377738 Cena: 129,00 zł

