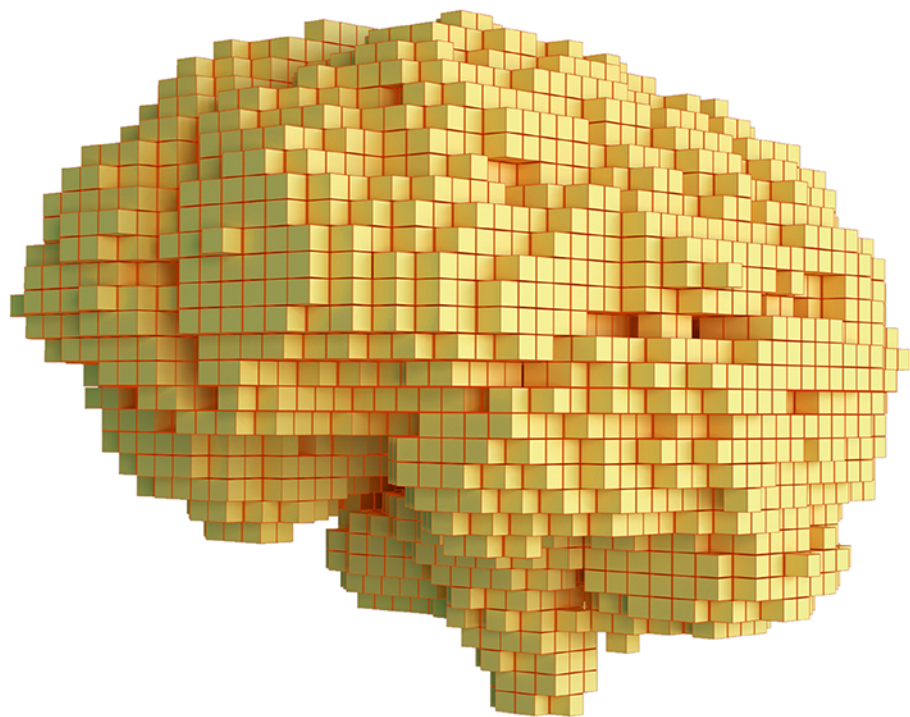


ALEX J. GUTMAN, JORDAN GOLDMEIER



# ANALITYK DANYCH

**PRZEWODNIK  
PO DATA SCIENCE, STATYSTYCE  
I UCZENIU MASZYNOWYM**

Helion 

Tytuł oryginału: Becoming a Data Head: How to Think, Speak and Understand  
Data Science, Statistics and Machine Learning

Tłumaczenie: Grzegorz Werner

ISBN: 978-83-289-0215-2

Copyright © 2021 by John Wiley & Sons, Inc., Indianapolis, Indiana  
All Rights Reserved. This translation published under license with the  
original publisher John Wiley & Sons, Inc.

Translation copyright © 2023 by Helion S.A.

No part of this publication may be reproduced, stored in a retrieval system,  
or transmitted in any form or by any means, electronic, mechanical,  
photocopying, recording, scanning, or otherwise, without either the prior  
written permission of the Publisher.

Wiley and the Wiley logo are trademarks or registered trademarks of John  
Wiley & Sons, Inc. and/or its affiliates, in the United States and other countries,  
and may not be used without written permission. All other trademarks are the  
property of their respective owners. John Wiley & Sons, Inc. is not associated  
with any product or vendor mentioned in this book.

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

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

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

Drogi Czytelniku!

Jeżeli chcesz ocenić tę książkę, zajrzyj pod adres

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

Możesz tam wpisać swoje uwagi, spostrzeżenia, recenzję.

Helion S.A.

ul. Kościuszki 1c, 44-100 Gliwice

tel. 32 230 98 63

e-mail: [helion@helion.pl](mailto:helion@helion.pl)

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

Printed in Poland.

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

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

# Spis treści

O redaktorach technicznych	ix
Podziękowania	xi
Przedmowa	xxi
Wprowadzenie	xxv

## CZĘŚĆ I

---

### Myśl jak spec od danych

#### ROZDZIAŁ 1

#### **Na czym polega problem? 3**

Pytania, które powinien zadawać spec od danych	4
<i>Dlaczego problem jest ważny?</i>	4
<i>Na kogo wpływa ten problem?</i>	6
<i>Co, jeśli nie mamy właściwych danych?</i>	7
<i>Kiedy projekt się zakończy?</i>	7
<i>Co, jeśli nie spodobają nam się rezultaty?</i>	7
Dlaczego projekty związane z danymi kończą się niepowodzeniem?	8
<i>Wrażenia klientów</i>	8
<i>Omówienie</i>	10
Praca nad problemami, które mają znaczenie	11
Podsumowanie rozdziału	11

#### ROZDZIAŁ 2

#### **Czym są dane? 13**

Dane a informacje	13
<i>Przykładowy zbiór danych</i>	14
Typy danych	15
Jak gromadzi się dane i jaką mają strukturę?	16
<i>Dane obserwacyjne i eksperymentalne</i>	16
<i>Dane ustrukturyzowane i nieustrukturyzowane</i>	17
Podstawowe statystyki zbiorcze	18
Podsumowanie rozdziału	19

**ROZDZIAŁ 3****Przygotowanie do myślenia statystycznego 21**

Zadawaj pytania	22
Wszystko jest zmienne	23
<i>Scenariusz: wrażenia klientów (kontynuacja)</i>	24
<i>Studium przypadku: zachorowalność na raka nerki</i>	26
Prawdopodobieństwo i statystyka	28
<i>Prawdopodobieństwo a intuicja</i>	29
<i>Odkrywanie informacji za pomocą statystyki</i>	31
Podsumowanie rozdziału	33

**CZĘŚĆ II****Mów jak spec od danych****ROZDZIAŁ 4****Polemizuj z danymi 37**

Co byś zrobił(a)?	38
<i>Katastrofa spowodowana brakiem danych</i>	39
Jaka jest historia pochodzenia danych?	43
<i>Kto zebrał dane?</i>	44
<i>Jak zebrano dane?</i>	44
Czy dane są reprezentatywne?	45
<i>Czy poprawnie dobrano próbę?</i>	46
<i>Co zrobiono z wartościami odstającymi?</i>	46
Jakich danych nie widzę?	47
<i>Jak rozwiązano problem brakujących wartości?</i>	47
<i>Czy dane mogą zmierzyć to, co ma być mierzone?</i>	48
Polemizuj z danymi każdej wielkości	48
Podsumowanie rozdziału	49

**ROZDZIAŁ 5****Eksploruj dane 51**

Eksploracyjna analiza danych i Ty	52
Przyjmij nastawienie eksploracyjne	52
<i>Pytania naprowadzające</i>	53
<i>Scenariusz</i>	53
Czy dane mogą odpowiedzieć na pytanie?	54
<i>Określ oczekiwania i użyj zdrowego rozsądku</i>	54
<i>Czy wartości mają intuicyjny sens?</i>	54
<i>Uważaj! Wartości odstające i brakujące</i>	58

Czy odkryliście jakieś związki?	58
<i>Korelacja</i>	59
<i>Uważaj! Błędne interpretowanie korelacji</i>	60
<i>Uważaj! Korelacja nie implikuje przyczynowości</i>	62
Czy znaleźliście w danych nowe możliwości?	63
Podsumowanie rozdziału	63

## ROZDZIAŁ 6

### **Badaj prawdopodobieństwa** **65**

Zgadnij odpowiedź	66
Reguły gry	67
<i>Notacja</i>	67
<i>Prawdopodobieństwo warunkowe i zdarzenia niezależne</i>	69
<i>Prawdopodobieństwo wielu zdarzeń</i>	69
Dwie rzeczy, które zdarzają się razem	69
Jedno albo drugie	71
Ćwiczenie myślowe z zakresu prawdopodobieństwa	72
<i>Następne kroki</i>	73
Uważaj z zakładaniem niezależności	74
<i>Nie popełniaj błędu hazardzisty</i>	75
Wszystkie prawdopodobieństwa są warunkowe	75
<i>Nie przestawiaj zależności</i>	76
<i>Twierdzenie Bayesa</i>	77
Upewnij się, że prawdopodobieństwa mają znaczenie	80
<i>Kalibracja</i>	80
<i>Rzadkie zdarzenia mogą się zdarzać i się zdarzają</i>	80
Podsumowanie rozdziału	81

## ROZDZIAŁ 7

### **Kwestionuj statystyki** **83**

Krótkie lekcje o wnioskowaniu	83
<i>Zostaw sobie trochę przestrzeni</i>	84
<i>Więcej danych, więcej dowodów</i>	84
<i>Kwestionuj status quo</i>	85
<i>Dowody na twierdzenie przeciwne</i>	86
<i>Równoważenie błędów decyzyjnych</i>	88
Proces wnioskowania statystycznego	89
Pytania, które pomogą Ci kwestionować statystyki	90
<i>Jaki jest kontekst tych statystyk?</i>	90
<i>Jaki jest rozmiar próby?</i>	91
<i>Co testujecie?</i>	91

<i>Jaka jest hipoteza zerowa?</i>	92
<i>Zakładanie równoważności</i>	93
<i>Jaki jest poziom istotności?</i>	93
<i>Ile przeprowadzacie testów?</i>	94
<i>Czy mogę zobaczyć przedziały ufności?</i>	94
<i>Czy jest to praktycznie istotne?</i>	96
<i>Czy zakładacie przyczynowość?</i>	96
Podsumowanie rozdziału	97

## CZĘŚĆ III

### Przybornik specjalisty data science

#### ROZDZIAŁ 8

#### **W poszukiwaniu ukrytych grup** **101**

Uczenie nienadzorowane	102
Redukcja wymiarowości	102
<i>Tworzenie cech złożonych</i>	103
Analiza składowych głównych	105
<i>Składowe główne zdolności sportowych</i>	105
<i>Podsumowanie PCA</i>	108
<i>Potencjalne pułapki</i>	109
Klasteryzacja	109
Klasteryzacja metodą k-średnich	111
<i>Klasteryzacja sklepów detalicznych</i>	111
<i>Potencjalne pułapki</i>	113
Podsumowanie rozdziału	114

#### ROZDZIAŁ 9

#### **Model regresji** **117**

Uczenie nadzorowane	117
Jak działa regresja liniowa?	119
<i>Regresja metodą najmniejszych kwadratów:</i>	
<i>nie tylko pomysłowa nazwa</i>	120
Regresja liniowa: co Ci daje?	123
<i>Rozszerzanie modelu na wiele cech</i>	124
Regresja liniowa: jakie powoduje nieporozumienia?	125
<i>Pominięte zmienne</i>	126
<i>Współliniowość</i>	126
<i>Przeciek danych</i>	127
<i>Błędy ekstrapolacji</i>	128
<i>Relacje nie zawsze są liniowe</i>	128

Wyjaśniasz czy przewidujesz?	128
Skuteczność regresji	130
Inne modele regresji	131
Podsumowanie rozdziału	131

## ROZDZIAŁ 10

### **Model klasyfikacji** **133**

Wprowadzenie do klasyfikacji	133
<i>Czego się nauczysz?</i>	134
<i>Przykładowy problem klasyfikacji</i>	135
Regresja logistyczna	135
<i>Regresja logistyczna — i co z tego?</i>	138
Drzewa decyzyjne	139
Metody zespołowe	142
<i>Lasy losowe</i>	143
<i>Drzewa wzmacniane gradientowo</i>	143
<i>Interpretowalność modeli zespołowych</i>	145
Strzeż się pułapek	145
<i>Złe podejście do problemu</i>	146
<i>Przeciek danych</i>	146
<i>Brak podziału danych</i>	146
<i>Wybór odpowiedniego progu decyzyjnego</i>	147
Błędne rozumienie dokładności	147
<i>Macierze błędów</i>	148
Podsumowanie rozdziału	150

## ROZDZIAŁ 11

### **Analiza tekstu** **151**

Oczekiwania wobec analizy tekstu	151
Jak tekst staje się liczbami	153
<i>Wielki worek słów</i>	153
<i>N-gramy</i>	157
<i>Osadzenia słów</i>	158
Modelowanie tematyczne	160
Klasyfikacja tekstu	162
<i>Naiwny klasyfikator bayesowski</i>	164
<i>Analiza odczuć</i>	166
Kwestie praktyczne podczas pracy z tekstem	167
<i>Giganci technologiczni mają przewagę</i>	168
Podsumowanie rozdziału	169

**ROZDZIAŁ 12****Uczenie głębokie 171**

Sieci neuronowe	172
<i>Pod jakimi względami sieci neuronowe przypominają ludzki mózg?</i>	172
<i>Prosta sieć neuronowa</i>	173
<i>Jak uczy się sieć neuronowa?</i>	174
<i>Nieco bardziej złożona sieć neuronowa</i>	175
Zastosowania uczenia głębokiego	178
<i>Korzyści z uczenia głębokiego</i>	179
<i>Jak komputery „widzą” obrazy?</i>	180
<i>Konwolucyjne sieci neuronowe</i>	181
<i>Uczenie głębokie w języku i sekwencjach</i>	183
Uczenie głębokie w praktyce	184
<i>Czy masz dane?</i>	184
<i>Czy Twoje dane są ustrukturyzowane?</i>	185
<i>Jak będzie wyglądać sieć?</i>	186
Sztuczna inteligencja i Ty	187
<i>Giganci technologiczni mają przewagę</i>	188
<i>Etyka w uczeniu głębokim</i>	188
Podsumowanie rozdziału	189

**CZĘŚĆ IV****Droga do sukcesu****ROZDZIAŁ 13****Strzeż się pułapek 193**

Tendencyjność i dziwne zjawiska w danych	194
<i>Błąd przeżywalności</i>	194
<i>Regresja do średniej</i>	195
<i>Paradoks Simpsona</i>	195
<i>Błąd potwierdzenia</i>	197
<i>Błąd utopionych kosztów</i>	197
<i>Dyskryminacja algorytmiczna</i>	197
<i>Nieskatyzowane przejawy tendencji</i>	198
Wielka lista pułapek	199
<i>Pułapki związane ze statystyką i uczeniem maszynowym</i>	199
<i>Pułapki związane z projektem</i>	200
Podsumowanie rozdziału	202



**ROZDZIAŁ 14****Znaj ludzi i osobowości 203**

Siedem scenariuszy fiaska komunikacyjnego	204
<i>Post mortem</i>	204
<i>Wieczorynka</i>	205
<i>Głuchy telefon</i>	206
<i>W gąszczu szczegółów</i>	206
<i>Konfrontacja z rzeczywistością</i>	207
<i>Wrogie przejęcie</i>	207
<i>Egocentryk</i>	207
Osobowości w świecie danych	208
<i>Entuzjasta</i>	208
<i>Cynik</i>	209
<i>Spec od danych</i>	209
Podsumowanie rozdziału	209

**ROZDZIAŁ 15****Co dalej? 211**



# Eksploruj dane

„Jeśli powiesz specjalście data science, żeby poszedł łowić ryby...  
to zasługujesz na to, co otrzymasz, czyli kiepską analizę”<sup>1</sup>.

– Thomas C. Redman, „Data Doc”,  
współpracownik „Harvard Business Review”

**P**rojekty związane z danymi nigdy nie są takie proste, jakie wydają się podczas prezentacji w sali posiedzeń zarządu. Interesariusze zwykle widzą dopieszczoną prezentację PowerPointa, która przebiega według sztywno ustalonego skryptu, od pytania, poprzez dane, do odpowiedzi. Tym, czego nie widać w tej historii, są jednak wszystkie idee, które nie przeszły przez sito — ważne decyzje i założenia, które po drodze podjął zespół ds. danych, aby uzyskać odpowiedź. Dobry zespół ds. danych podąża nie prostą, ale krętą ścieżką, dostosowując się do dokonywanych odkryć. W miarę podróży wraca do wcześniejszych pomysłów i zauważa, że w rezultacie otworzyło się wiele nowych dróg.

Ten proces iteracji, odkryć i przyglądania się danym nosi nazwę **eksploracyjnej analizy danych** (ang. *exploratory data analysis*, EDA). Został sformułowany przez statystyka Johna Tukeya w latach 70. jako sposób na wstępne zrozumienie danych poprzez zbiorcze statystyki i wizualizacje przed zastosowaniem bardziej skomplikowanych metod<sup>2</sup>. Tukey postrzegał EDA jako pracę detektywistyczną. W danych ukryte są wskazówki, a właściwa eksploracja może zasugerować

<sup>1</sup> Cytat autorstwa Amy Gallo zaczerpnięty z *Understand Regression Analysis*, rozdział 10. przewodnika *HBR Guide to Data Analytics Basics for Managers*, HBR Guide Series.

<sup>2</sup> J.W. Tukey, *Exploratory data analysis*, 1977, vol. 2, s.131 – 160.

następne kroki. W rzeczywistości EDA jest kolejnym sposobem „polemizowania” z danymi. Jest to fundamentalna część pracy z danymi, która wyznacza i zmienia kierunek projektu w zależności od dokonywanych odkryć.

## **EKSPLORACYJNA ANALIZA DANYCH I TY**

---

Eksploracyjna analiza danych może być dla niektórych niekomfortowa — ujawnia ona subiektywną naturę (sztukę?) pracy z danymi. Dwa zespoły, otrzymawszy ten sam problem i dane, mogą wybrać dwie różne ścieżki analizy, czasem dochodząc do tych samych wniosków. A czasem nie. Po prostu po drodze jest do podjęcia zbyt wiele decyzji, żeby dwa zespoły (lub dwie osoby) zrobiły wszystko tak samo. Każda osoba będzie miała inne kompetencje, pomysły i narzędzia potrzebne do rozwiązania problemu.

Dlatego w tym rozdziale prezentujemy EDA jako ciągły proces, który należy do obowiązków każdego specja od danych, bez względu na to, czy bezpośrednio pracuje z danymi, czy jest członkiem zarządu firmy. Nauczysz się, jakie zadawać pytania i na co zwracać uwagę podczas eksplorowania danych.

### **Czy jesteś menedżerem, czy liderem?**

Jeśli jesteś interesariuszem, menedżerem lub ekspertem merytorycznym, staraj się być dostępny dla swojego zespołu ds. danych. Prowadź otwarty dialog i oczekuj iteracji. Współpracuj z zespołem, aby ustalić poprawne założenia. Nie pozwól, aby zespół udał się na ryby bez odpowiedniego kontekstu biznesowego. W przeciwnym razie może podążać ścieżkami, które mają większy sens statystyczny niż praktyczny. Jedno błędne założenie może zagrozić wszystkiemu, co nastąpi potem.

Rozumiemy, że menedżerowie nie mogą angażować się w szczegóły projektu w takim samym stopniu jak osoby bezpośrednio pracujące z danymi. Jest tu jednak miejsce na poprawę. Nie musisz stosować mikrozarządzania. Po prostu nie możesz ignorować<sup>3</sup>.

## **PRZYMIJ NASTAWIENIE EKSPLORACYJNE**

---

Dziesiątki dostępnych narzędzi i języków programowania mogą pomóc zespołom ds. danych szybko i niedrogo eksplorować dane poprzez zbiorcze statystyki i wizualizacje. Nie należy jednak myśleć o EDA jak o przyborniku z narzędziami

---

<sup>3</sup> Żeby wszystko było jasne, interesariusze powinni powstrzymać się od mikrozarządzania. Konieczny jest pewien stopień zaufania między pracownikami biznesowymi a zespołami ds. danych.

albo liście kontrolnej. Jest to raczej nastawienie mentalne wplecione w każdą fazę pracy z danymi, które możesz przyjąć nawet bez zaplecza analitycznego.

## Pytania naprowadzające

Aby pomóc Ci w przyjęciu nastawienia eksploracyjnego, przeprowadzimy Cię przez krótki scenariusz oparty na popularnym zbiorze danych utworzonym do celów edukacyjnych: *Ames Housing Data*<sup>4</sup>. Zapewni Ci to wgląd w proces EDA.

Choć nie ma jednej właściwej ścieżki, którą należy podążać, jest kilka pytań, które możesz zadać, aby pomóc zespołowi w dojściu do użytecznych wniosków:

- Czy dane mogą odpowiedzieć na pytanie?
- Czy odkryliśmy jakieś związki?
- Czy znaleźliśmy w danych nowe możliwości?

Teraz nakreślimy ramy naszego scenariusza i omówimy każde z tych trzech pytań, aby wyjaśnić, dlaczego warto je zadać, i wskazać potencjalne problemy, które możesz napotkać.

## Scenariusz

Pracujesz w start-upie działającym na rynku nieruchomości i musisz zwiększyć ruch w swojej witrynie. Trudno jednak odciągnąć klientów od technologicznych gigantów, takich jak amerykański serwis Zillow.com. Jego słynne narzędzie do wyceny nieruchomości, Zestimate®, przyciąga ludzi (i zyski) do marki Zillow<sup>5</sup>. Aby skutecznie konkurować z potentatami, Twoja firma potrzebuje własnego narzędzia predykcyjnego. Otrzymałeś(-łaś) więc zadanie zbudowania *modelu*, który przyjmuje informacje o domu jako *dane wejściowe* i generuje szacowaną cenę sprzedaży jako *dane wyjściowe*.

Na dobry początek szef wysłał Ci zbiór danych. Ma on 80 kolumn opisujących różne aspekty setek domów mieszkalnych sprzedanych w Ames w stanie Iowa w latach 2006 – 2011.

Tak duża ilość danych może wydawać się przytłaczająca. Jednak wymienione poprzednio pytania pomogą Ci przygotować wstępny plan pracy z danymi.

Omówmy je po kolei.

---

<sup>4</sup> D. De Cock, *Ames, Iowa: Alternative to the Boston housing data as an end of semester regression project*, „Journal of Statistics Education” 2013, vol. 19(3). Dane możesz pobrać pod adresem [www.kaggle.com/c/house-prices-advanced-regression-techniques](http://www.kaggle.com/c/house-prices-advanced-regression-techniques).

<sup>5</sup> Firma Zillow traktuje swoje narzędzie Zestimate® bardzo poważnie. W 2019 roku przyznała milion dolarów nagrody zespołowi specjalistów data science za zwiększenie dokładności szacunków Zestimate®: [venturebeat.com/2019/01/30/zillow-awards-1-million-to-team-that-reduced-home-valuation-algorithm-error-to-below-4](http://venturebeat.com/2019/01/30/zillow-awards-1-million-to-team-that-reduced-home-valuation-algorithm-error-to-below-4).

## CZY DANE MOGĄ ODPOWIEDZIEĆ NA PYTANIE?

---

Choć możesz odczuwać pokusę, aby natychmiast przepuścić dane przez obecnie modny algorytm (na przykład uczenie głębokie opisane w rozdziale 12.), najpierw musisz zapytać: „Czy dane mogą odpowiedzieć na pytanie?”. A żeby znaleźć odpowiedź, często wystarczy po prostu przyrzeć się danym.

### Określ oczekiwania i użyj zdrowego rozsądku

Powinieneś (powinnaś) dość dobrze wiedzieć, jakie informacje są potrzebne do oszacowania ceny sprzedaży domu: rozmiar, liczba sypialni, liczba łazienek, rok budowy itd. Są to popularne cechy, które potencjalni kupujący wyszukują w Twojej witrynie. Przewidywanie ceny sprzedaży bez tych informacji nie wydawałoby się rozsądne.

Nazwy kolumn i typy danych są widoczne od razu po otwarciu pliku. Obecne są tam zdroworozsądkowe informacje, których oczekujesz, a także pomocne dane porządkowe (ogólna jakość domu w skali 1 – 10, przy czym 10 to jakość „doskonała”), dane nominalne (dzielnica) i wiele innych cech. Dane przeszły zatem wstępny test.

Następnie powinieneś (powinnaś) zbadać wartości zmiennych. Czy obejmują one scenariusze, które chcesz analizować? Jeśli na przykład odkryjesz, że zmienna „Typ nieruchomości” uwzględnia tylko domy jednorodzinne, ale nie apartamenty, bliźniaki albo mieszkania, Twój model będzie miał ograniczony zakres w porównaniu z tym, którego używa Zillow. Zestimate® może przewidzieć cenę sprzedaży mieszkania, ale jeśli nie masz historycznych danych dotyczących mieszkań, Twoja firma nie będzie mogła precyzyjnie szacować ich cen.

Lekcja: unikaj wypraw na ryby, przed którymi ostrzegał Cię cytat na początku tego rozdziału. Upewnij się, że dane mają sens z perspektywy tego, jak będą wykorzystywane.

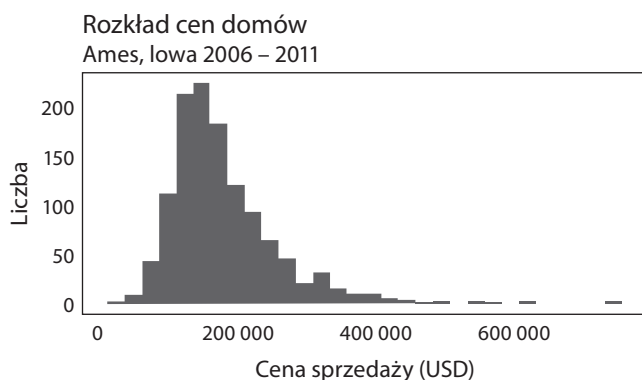
### Czy wartości mają intuicyjny sens?

Możesz użyć oprogramowania, aby wygenerować mnóstwo zbiorczych statystyk. Twoim zadaniem jest umieścić dane w kontekście. Sprawdź, czy zbiorcze statystyki odpowiadają Twojemu intuicyjnemu rozumieniu problemu. Równie ważnym elementem EDA są wizualizacje — wykorzystaj je do odkrywania anomalii i innych dziwactw w danych.

## Powtórka z wizualizacji danych

Omówmy kilka przykładów EDA z histogramami oraz wykresami skrzynkowymi, słupkowymi i punktowymi. Możesz pominąć tę ramkę, jeśli rozumiesz te wykresy i wiesz, co mogą Ci odpowiedzieć.

Możesz poznać kształt albo rozkład ciągłych danych liczbowych poprzez przyjrzenie się histogramowi. Rozważ histogram cen sprzedaży pokazany na rysunku 5.1. Widać na nim przykład około 125 domów w przedziale 200 000 dol. i długi „ogon” po prawej stronie, który reprezentuje najdroższe domy. Ogon ten przeciąga średnią cenę sprzedaży (181 000 dol.) poza medianę (163 000 dol.). Kilka drogich domów sprawia, że średnia jest większa niż mediana.

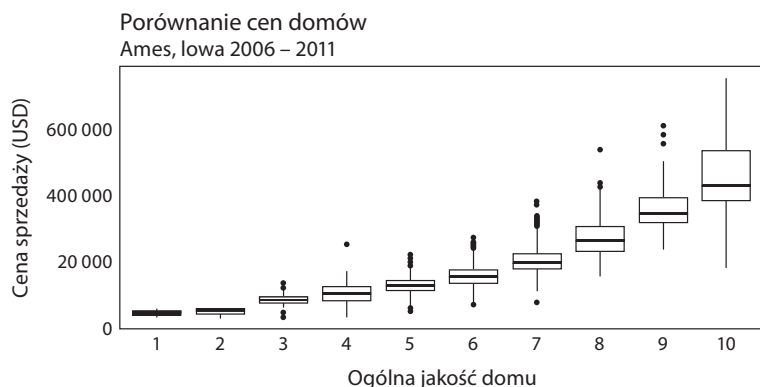


**RYSUNEK 5.1.** Histogram pokazujący kształt cen sprzedaży

Histogramy pomagają dostrzegać anomalie. Gdybyś zobaczył(a) wartości ujemne (ktoś otrzymuje zapłatę za kupienie domu?) albo przedziały z bardzo dużymi liczbami na skrajach rysunku 5.1, co zdarza się często, kiedy dane są „obcięte” (na przykład każda wartość powyżej 500 000 dol. jest wprowadzana jako 500 000 dol.), powinieneś (powinnaś) zadać kilka pytań.

Wykresy skrzynkowe<sup>6</sup> można wykorzystać do porównywania danych w kilku grupach. Na rysunku 5.2 pokazano wykres skrzynkowy dla ocen jakości domu, gdzie 1 to jakość bardzo niska, a 10 to doskonała.

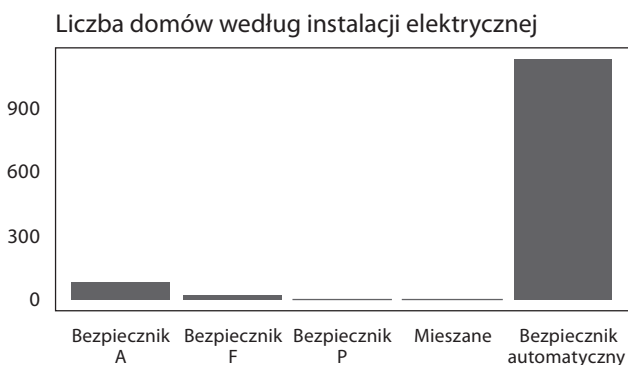
<sup>6</sup> Wykresy skrzynkowe nazywa się również wykresami typu „skrzynka i wąsy”. „Skrzynka” zawiera środkową połowę danych (wartości między 25. a 75. centylem), linia w środku to mediana, a „wąsy” pokazują zakres pozostałych punktów danych. Kropki pod wąsami to potencjalne wartości odstające.



**RYСУNEK 5.2.** Używanie wykresów skrzynkowych do porównywania cen przy różnych ocenach jakości

Związek między ogólną jakością a cenami domów wydaje się zgodny z intuicją. Domy wyższej jakości zwykle sprzedają się za wyższe kwoty. Dostrzegamy dom za 200 000 dol. z jakością ocenioną na 10 (dolny koniec linii), ale rozsądne wydaje się założenie, że sprzedano go za mniej niż inne idealne „dziesiątki” ze względu na inne czynniki. Jest to dokładnie ten rodzaj informacji, które powinien sprawdzać każdy, kto pracuje z danymi.

Wykresy słupkowe, takie jak pokazany na rysunku 5.3, zliczają dane kategoryczne.



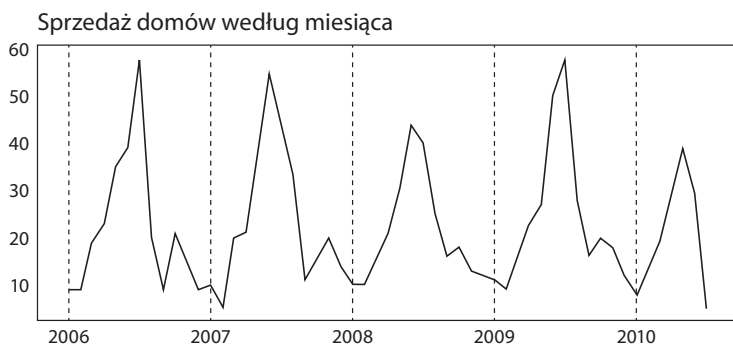
**RYСУNEK 5.3.** Wykres słupkowy pokazujący liczby domów według instalacji elektrycznej

Nie wszystkie obrazy są interesujące na pierwszy rzut oka. Warto jednak obejrzeć różne wizualizacje choćby po to, aby potwierdzić (a może zakwestionować) odpowiedź na poprzednie pytanie — czy dane mają intuicyjny sens? Na przykład rysunek 5.3 pokazuje, że niemal wszystkie domy mają



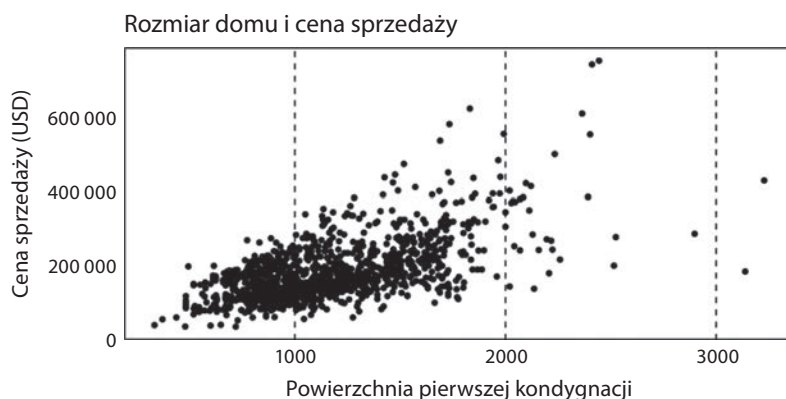
jednakową wartość tej cechy. Dla Twoich celów informacja ta jest jednak przydatna. Ponieważ większość domów ma jednakową wartość tej zmiennej, prawdopodobnie nie wpłynie ona znacząco na ceny sprzedaży domów.

Na rysunku 5.4 pokazano wykres liniowy z liczbą domów sprzedanych w poszczególnych miesiącach. Prezentuje on wzrosty sprzedaży latem i spadki zimą — przykład *sezonowości*. Wykresy liniowe pomagają dostrzegać takie trendy.



**RYСУNEK 5.4.** Wykres liniowy pokazujący liczbę domów sprzedanych w poszczególnych miesiącach

Następnie zbadajmy wykres punktowy pokazujący domy według ich rozmiaru (powierzchni pierwszej kondygnacji) i ceny sprzedaży (rysunek 5.5).



**RYСУNEK 5.5.** Wykres punktowy pokazujący powierzchnię i cenę sprzedaży

Z rysunku 5.5 wyłania się interesujący wzorzec. Większe domy zwykle sprzedają się za większe kwoty. Oczywiście, ta reguła nie zawsze się sprawdza. Czasem mniejsze domy kosztują więcej niż większe. Zawsze jest jakaś zmienność, ale ogólny trend jest oczywisty. Ponieważ zaś próbujemy przewidywać ceny sprzedaży, powierzchnia domu wydaje się bardzo cenną informacją.

Był to tylko przedsmak informacji i spostrzeżeń, które możesz szybko poczynić poprzez naniesienie danych na wykres. Gdybyś chciał(a) dowiedzieć się więcej o wizualizacji w eksploracji danych, polecamy następujące tytuły:

- Stephen Few, *Now You See it: Simple Visualization Techniques for Quantitative Analysis*, Analytics Press, 2009.
- Edward Tufte, *The Visual Display of Quantitative Information*, Edward Tufte Graphics Press, 2011.

## Uważaj! Wartości odstające i brakujące

W każdym zbiorze danych znajdują się anomalie, wartości odstające i wartości brakujące. Sposób ich traktowania ma znaczenie.

Na przykład wykres skrzynkowy z rysunku 5.2 wykorzystuje praktyczne reguły, aby oznaczyć kilka punktów danych jako potencjalne wartości odstające. Jednak samo to, że jakaś grafika klasyfikuje pewne punkty jako „wartości odstające”, nie oznacza jeszcze, że możesz wyłączyć krytyczne myślenie i automatycznie usunąć te punkty, zakładając, że nie będą użyteczne. Zillow na pewno nie usuwa przydatnych informacji ze swoich zbiorów danych tylko dlatego, że wizualizacja opisuje je jako odstające. Zwróć uwagę na kontekst danych — domy, które kosztują znacznie więcej niż większość innych, to znana i częsta cecha danych dotyczących nieruchomości. Przypomnij sobie lekcje z poprzedniego rozdziału. Powinieneś (powinnaś) mieć dobre uzasadnienie biznesowe, aby usunąć wartości odstające. Czy masz je w tym przypadku?

A co z brakującymi wartościami? Czy brakująca wartość zmiennej „Rozmiar piwnicy” oznacza, że dom ma piwnicę, ale jej powierzchnia jest nieznana? Czy może nie ma piwnicy i wartość powinna wynosić 0?

Jeśli wydaje się, że mnożymy problemy, to właśnie taki był nasz zamiar. Osoby pracujące z danymi podejmują setki takich drobnych decyzji podczas każdego projektu. Ich łączny efekt może być znaczący. Pozostawieni samym sobie — i pozbawieni dostępu do wiedzy dziedzinowej — pracownicy zajmujący się danymi mogą po trochu je oskubywać, usuwając trudne i zniuansowane przypadki, aż dane staną się zbyt oddalone od rzeczywistości, żeby były użyteczne. Właśnie dlatego wszyscy, łącznie z menedżerami, powinni naprawdę rozumieć, co robią zespoły ds. danych.

## CZY ODKRYLIŚCIE JAKIEŚ ZWIĄZKI?

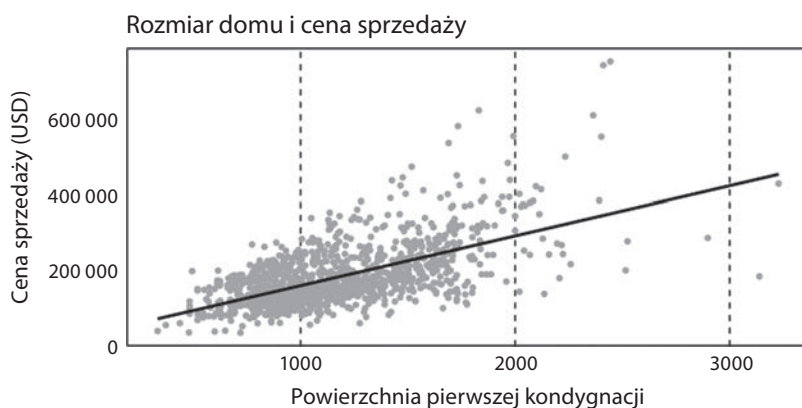
Wstępna analiza danych mieszkaniowych z wykorzystaniem zbiorczych statystyk i wizualizacji wydaje się zachęcająca i uważasz, że danych rzeczywiście będzie można użyć do zbudowania modelu przewidującego ceny sprzedaży, więc przechodzisz do następnego pytania: „Czy odkryliście jakieś związki?”

Wizualizacja danych dała Ci pierwsze wskazówki: wyższa ogólna jakość i większy metraż są, jak można było się spodziewać, związane z wyższymi cenami sprzedaży. Właśnie takich informacji zwrotnych szukasz w danych. Związki wydają się sensowne, a zmienne, które naniósł(-łaś) na wykres, pomogą Ci zbudować model do przewidywania ceny sprzedaży. Jakie inne zmienne mają związek z ceną sprzedaży?

W tym momencie na interesujące wzorce i relacje w danych mogą nakierować Cię statystyki zbiorcze, ponieważ generowanie każdego możliwego wykresu punktowego może nie być praktyczne. Relację znaną na wykresach punktowych można sprowadzić do statystyki zbiorczej zwanej **korelacją**, która wskazuje na związek między dwiema zmiennymi liczbowymi (choć go nie dowodzi).

## Korelacja

Korelacja jest miarą tego, jak powiązane są dwie zmienne. Najczęstszym typem korelacji używanym w biznesie jest **współczynnik korelacji Pearsona**, statystyka o wartości od  $-1$  do  $1$ , która mierzy liniowy związek (wyrażony zwykłymi liniami prostymi) między parami liczb na wykresie punktowym. Korelacja może być dodatnia, co oznacza, że zwiększenie jednej zmiennej powoduje wzrost drugiej: większe domy sprzedają się za wyższe kwoty. Może też być ujemna: cięższe samochody przejeżdżają mniej kilometrów na tej samej ilości benzyny. Korelacja między rozmiarem domu a ceną sprzedaży, pokazana na rysunku 5.6, wynosi  $0,62$ . Im „ciaśniej” skupione są punkty wokół linii trendu, tym wyższa korelacja<sup>7</sup>.



**RYSUNEK 5.6.** Metraż i cena sprzedaży mają korelację równą  $0,62$ , która mierzy, jak ciasno skupiają się punkty wokół linii trendu

<sup>7</sup> Korelacja nie oznacza „stromizny”. Wykresy dwóch ściśle skorelowanych zmiennych mogą wydawać się niemal płaskie (choć niezupełnie poziome).

Korelacja może Ci pomóc na dwa sposoby. Po pierwsze, znalezienie zmiennych skorelowanych z ceną sprzedaży pomogłoby ją przewidywać. Po drugie, korelacja może ograniczyć redundancję w danych, ponieważ dwie ściśle skorelowane zmienne zawierają mniej więcej te same informacje. Wyobraź sobie, że masz w danych dwie kolumny: powierzchnię domu w stopach kwadratowych i w metrach kwadratowych. Te dwie zmienne są idealnie skorelowane; do analizy potrzebna jest tylko jedna.

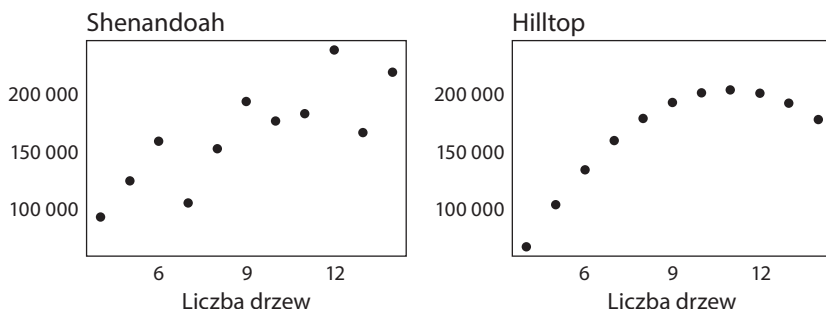
Choć większość z nas wie co nieco o korelacji i często posługuje się tą miarą, bywa ona zwodnicza. Sprawdźmy dlaczego.

### Uważaj! Błędne interpretowanie korelacji

Często zapomina się, że korelacja jest miarą trendu liniowego, a nie wszystkie trendy są liniowe.

Przypuśćmy na przykład, że analizujesz dwie dzielnice w zbiorze danych mieszkaniowych, każdą z 11 domami. Szybkie obliczenia ujawniają, że liczba drzew na działce jest ściśle skorelowana z ceną sprzedaży w tych dzielnicach. Korelacja wynosi 0,8: działki z większą liczbą drzew sprzedają się drożej.

Jednak wizualna kontrola danych ujawnia coś nieoczekiwanego. Na rysunku 5.7 dane po lewej stronie ukazują obraz, którego spodziewamy się po wysokiej korelacji: liniowy trend z nieco rozrzuconymi punktami danych. Ale wykres po prawej stronie pokazuje, że liczba drzew jest związana ze wzrostem ceny *tylko do pewnego punktu* (11 drzew). Następnie linia opada. Może w dzielnicy Hilltop na niektórych działkach rośnie tyle drzew, że przeszkadzają one w pielęgnowaniu trawnika.



**RYSUNEK 5.7.** Dwa zbiory danych z korelacją równą 0,8

Pora uczciwie przyznać: dane pokazane na rysunku 5.7 nie pochodzą z badanego przez nas zbioru danych Ames, ale z popularnego zestawu Anscombe's Quartet<sup>8</sup> składającego się z czterech zbiorów danych z iden-

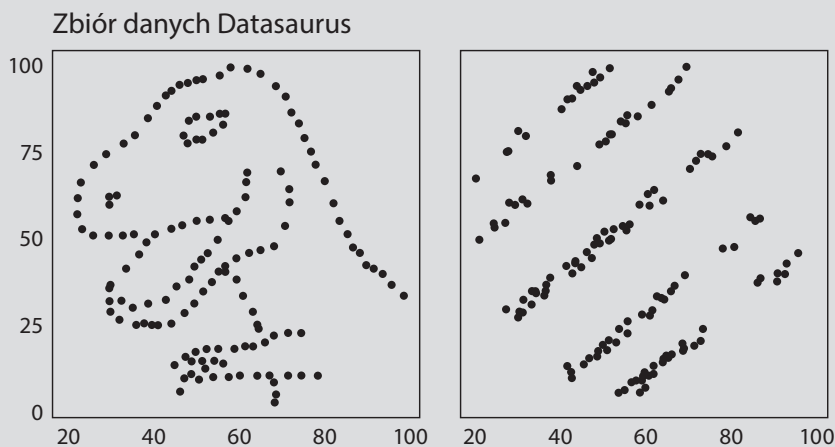
<sup>8</sup> F.J. Anscombe, *Graphs in statistical analysis*, „The American Statistician” 1973, vol. 27(1), s. 17 – 21. Pomnożyliśmy zmienną zależną przez 22 000, aby utworzyć realistyczne przykłady cen sprzedaży.

tycznymi statystykami zbiorczymi, ale wyraźnie różnymi wizualizacjami. (Pokazaliśmy tylko dwie i dostosowaliśmy dane, aby pozostać przy temacie sprzedaży nieruchomości).

Lekcja: używaj wizualizacji, aby zweryfikować godne uwagi korelacje w danych, ponieważ liniowy trend zidentyfikowany przez korelację może nie opowiadać pełnej historii.

## Nieskorelowane, ale wciąż interesujące

Na rysunku 5.8 pokazano dwa wykresy z niemal identycznymi, bliskimi zera współczynnikami korelacji. Nie oznacza to jeszcze, że nie ma w nich nic interesującego. I choć nie napotkasz wielu „danozaurów”, jak ten pokazany na lewym wykresie, możesz trafić na scenariusz z prawego wykresu: pięć grup skorelowanych liniowo danych, które jako całość w ogóle nie wykazują korelacji. Jest to tzw. paradoks Simpsona, o którym dowiesz się więcej w rozdziale 13.



**RYSUNEK 5.8.** Danozaur: możesz sam pobrać i zbadać dane<sup>9</sup>. Podobnie jak w przypadku Anscombe’s Quartet, oba pokazane tu zbiory danych mają identyczne statystyki zbiorcze

<sup>9</sup> Danozaur został stworzony przez Alberto Caira, a dane są dostępne na GitHubie: [github.com/lockedata/datasauRus](https://github.com/lockedata/datasauRus).

## Uważaj! Korelacja nie implikuje przyczynowości

Prawdopodobnie słyszałeś(-łaś) już wcześniej powiedzenie: „Korelacja nie implikuje przyczynowości”<sup>10</sup>. Warto je jednak tu powtórzyć, ponieważ jest często ignorowane, a nawet źle rozumiane.

Kiedy dwie zmienne są skorelowane, nawet bardzo ściśle, nie oznacza to, że jedna powoduje drugą. Ludzie jednak często wpadają w tę pułapkę, próbując zbudować jakąś narrację za każdym razem, kiedy dwie zmienne poruszają się razem. Istnieją typowe niemądre przykłady, których używają statystycy, aby pokazać, że korelacja nie oznacza przyczynowości. Sprzedaż lodów jest skorelowana z atakami rekinów (jedno i drugie rośnie w miesiącach letnich). Rozmiar buta jest skorelowany z umiejętnością czytania (jedno i drugie rośnie z biegiem czasu). Jednak sugerowanie, że zmniejszenie sprzedaży lodów ograniczyłoby liczbę ataków rekinów albo że kupno większych butów pomogłoby w czytaniu to oczywiście żart. W grę wchodzi inne czynniki — temperatura zewnętrzna w przykładzie z lodami, wiek w przykładzie z rozmiarem buta — które odgrywają oczywistą rolę w złudnym związku.

Kiedy jednak korelacje nie pojawiają się w kontekście żartów, a prawdziwe przyczyny nie są znane, mantra „korelacja nie implikuje przyczynowości” często idzie w zapomnienie.

Na przykład w danych dotyczących nieruchomości widać, że miary edukacyjnych wyników szkół są skorelowane z wartością domów. Czy oznacza to, że lepsze szkoły powodują wzrost wartości domu? Można domniemywać, że dobre szkoły podnoszą atrakcyjność dzielnicy. A może przyczynowość działa w odwrotnym kierunku: wyższe ceny domów powodują poprawę edukacyjnych wyników szkół? Może większe przychody z podatków zapewniają szkołom więcej zasobów? A może przyczynowość działa w obu kierunkach, tworząc sprzężenie zwrotne? W większości przypadków po prostu tego nie wiemy. W grę wchodzi mnóstwo innych czynników i rzadko zdarza się, żebyś miał(a) wszystkie odpowiedzi w swoim zbiorze danych.

Bezpieczniej jest założyć, że „nie ma przyczynowości” między dwiema skorelowanymi zmiennymi, dopóki ktoś nie przeprowadzi eksperymentu, który jej dowiedzie. Nie należy jednak popadać w skrajność. Obaj autorzy widzieli, jak w środowiskach biznesowych, uniwersyteckich i medialnych zakładano przyczynowość tam, gdzie było to nieuzasadnione. Bywają też jednak przypadki, w których ważny związek jest natychmiast odrzucany jako rzekomo błędne założenie przyczynowości. (Przykład odrzucenia przyczynowości w sytuacji, w której nie należało tego robić, podajemy w ramce).

---

<sup>10</sup> Autorzy zastanawiali się nawet, czy można napisać książkę o danych i nie wspomnieć o tym, że „korelacja nie implikuje przyczynowości”. Jaki był wynik tej dyskusji, widzisz sam.

## Palenie a rak płuc

Ronald A. Fisher, jeden z czołowych statystyków XX wieku, który zresztą był współtwórcą wielu technik opisywanych w tej książce, podchodził dość sceptycznie do badań wiążących palenie tytoniu z rakiem.

Fisher obawiał się o wpływ zmiennych zakłócających. Co by było, gdyby na przykład niektórzy ludzie byli genetycznie predysponowani do zachorowania na raka płuc i palili, aby złagodzić dolegliwości? Według Fishera we wczesnych badaniach nad zagrożeniami związanymi z używaniem tytoniu popełniono „błąd [...] starego rodzaju, przechodząc od korelacji do związku przyczynowego”<sup>11</sup>.

A jednak obecnie wiemy, że ten związek jest niepodważalny. Choć powinniśmy wystrzegać się dostrzegania przyczynowości tam, gdzie jej nie ma, musimy również uważać, aby nie odrzucać jeszcze niedowiedzionych związków przyczynowych.

## CZY ZNALEZLIŚCIE W DANYCH NOWE MOŻLIWOŚCI?

EDA nie tylko pozwala lepiej zrozumieć dane i wyznaczyć drogę do rozwiązania problemu. Daje też szansę znalezienia w danych dodatkowych możliwości; nowych problemów, które mogą być cenne dla Twojej organizacji. Specjalista data science może dostrzec w zbiorze danych coś interesującego lub nietypowego, a następnie sformułować problem.

Nie będziesz jednak wiedział(a), czy ktokolwiek potrzebuje rozwiązania, które znalazłeś(-łaś), dopóki nie wykonasz czynności opisanych w rozdziale 1. „Na czym polega problem?”.

## PODSUMOWANIE ROZDZIAŁU

Aby zostać specem od danych, musisz polubić proces eksploracyjnej analizy danych. Pozwoli Ci to:

- znaleźć klarowniejszą ścieżkę do rozwiązania problemu;
- doprecyzować pierwotny problem biznesowy zgodnie z ograniczeniami odnalezionymi w danych;
- zidentyfikować nowe problemy, które można rozwiązać z wykorzystaniem dostępnych danych;
- anulować projekt. Kiedy wyniki są niezadowalające, analiza EDA zapobiega traceniu czasu i pieniędzy na projekty, które zabrnęły w ślepią uliczkę.

<sup>11</sup> R.A. Fisher, *Cancer and smoking*, „Nature” 1958, vol. 182 (4635), s. 596.

Zademonstrowaliśmy proces EDA na przykładzie zbioru danych dotyczących nieruchomości (do którego wrócimy w rozdziale 9., aby wreszcie zbudować ten model, o którym mówiliśmy) i opisaliśmy przeszkody, które możesz napotkać po drodze.

W rozdziale tym założyliśmy, że możesz być częścią procesu EDA od początku do końca. Czasem jest to niemożliwe, zwłaszcza w przypadku dyrektorów wyższego szczebla, którzy nadzorują wiele projektów. Ale nieobecność we wczesnych fazach projektu nie zwalnia speców od danych z obowiązku przyjęcia eksploracyjnego nastawienia. Jeśli dołączasz do projektu pod koniec, zapytaj, dlaczego zespół ds. danych wybrał konkretny sposób analizy i jakie napotkał trudności. Możesz odkryć założenia, których nie poczynił(a)byś sam(a).



# 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**

## Data science? Odsiejesz piasek od złota!

Musisz spojrzeć prawdzie w oczy: epoka danych to nie tylko imponujące możliwości, ale również obietnice bez pokrycia. Firmy wdrażają rozwiązania, które mają je wyręczać w podejmowaniu decyzji. Menedżerowie zatrudniają analityków, którzy nimi nie są. Specjaliści w dziedzinie data science są zatrudniani w organizacjach, które nie są na nich gotowe. Dyrektorzy wysłuchują technicznego żargonu i udają, że go rozumieją. Efekt? Pieniądze idą w błoto.

Oto praktyczny przewodnik po nauce o danych w miejscu pracy. Dowiesz się stąd wszystkiego, co ważne na początku Twojej drogi jako danologa: od osobowości, z którymi przyjdzie Ci pracować, przez detale analizy danych, po matematykę stojącą za algorytmami i uczeniem maszynowym. Nauczysz się myśleć krytycznie o danych i otrzymanych wynikach, będziesz też inteligentnie o tym mówić. Jednym zdaniem: zrozumiesz dane i związane z nimi wyzwania na głębszym, profesjonalnym poziomie.



**To książka dla każdego, kto chce przestawić firmę na tory data science.**



**Eric Weber**

*kierownik ds. eksperymentów  
i badań metrycznych, Yelp*

Naucz się:

- myśleć statystycznie i rozumieć rolę zmienności w podejmowaniu decyzji
- zadawać właściwe pytania na temat statystyk i wyników analiz
- sensownie korzystać z rozwiązań uczenia maszynowego i sztucznej inteligencji
- unikać typowych błędów podczas pracy z danymi i ich interpretowania

**Dr Alex J. Gutman** jest adiunktem w Instytucie Technicznym Wojsk Lotniczych, specjalistą z zakresu data science i instruktorem biznesowym.

**Jordan Goldmeier** jest światowej klasy ekspertem w dziedzinie analityki i wizualizacji danych, a także autorem książek. Od sześciu lat otrzymuje nagrodę Excel MVP. Jest również wolontariuszem ratownictwa medycznego.

**Helion**

[helion.pl](http://helion.pl)

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

**KOD KORZYŚCI**  
Sięgnij po więcej! ▶



ISBN 978-83-289-0215-2



9 788328 902152

Cena: 69,00 zł

**WILEY**