

Przemysław Biecek

Analiza danych z programem R

Modele liniowe
z efektami stałymi,
losowymi i mieszanymi



W Y D A W N I C T W O N A U K O W E P W N

Analiza danych z programem R

Przemysław Biecek

Analiza danych z programem R

Modele liniowe
z efektami stałymi,
losowymi i mieszanymi



WYDAWNICTWO NAUKOWE PWN
WARSZAWA 2013

Projekt okładki i stron tytułowych **Agnieszka Łydźba**

Ilustracja na okładce **serknor/Shutterstock**

Redaktor inicjujący **Berenika Grajkowska**

Redaktor **Izabela Ewa Mika**

Koordynator produkcji **Mariola Grzywacka**

Łamanie **FixPoint, Warszawa**

Książka, którą nabyłeś, jest dziełem twórcy i wydawcy. Prosimy, abyś przestrzegał praw, jakie im przysługują. Jej zawartość możesz udostępnić nieodpłatnie osobom bliskim lub osobiście znanym. Ale nie publikuj jej w internecie. Jeśli cytujesz jej fragmenty, nie zmieniaj ich treści i koniecznie zaznacz, czyje to dzieło. A kopiując jej część, rób to jedynie na użytek osobisty.

Szanujmy cudzą własność i prawo

Więcej na www.legalnakultura.pl

Copyright © by Wydawnictwo Naukowe PWN SA
Warszawa 2011, 2013

ISBN 978-83-01-17453-8

Wydanie II

Wydawnictwo Naukowe PWN SA
tel. 22 69 54 321; faks 22 69 54 288
infolinia 801 33 33 88
e-mail: pwn@pwn.com.pl; www.pwn.pl
Druk i oprawa: OSDW Azymut Sp z o.o

Spis treści

Przedmowa	IX
1. Modele liniowe — wprowadzenie, podstawowe twierdzenia i wzory	1
1.1. Wprowadzenie	1
1.2. Model	3
1.3. Estymatory najmniejszych kwadratów i największej wiarygodności	7
1.3.1. Metoda najmniejszych kwadratów	7
1.3.2. Metoda największej wiarygodności	9
1.4. Rozkłady estymatorów	11
1.4.1. Asymptotyczny rozkład estymatorów największej wiarygodności	12
1.4.2. Rozkład estymatorów oparty na metodach permutacyjnych i metodzie bootstrap	12
1.5. Testy i przedziały ufności	14
1.5.1. Przedział ufności dla β_i	14
1.5.2. Test dla hipotezy brzegowej dotyczącej β_i	14
1.5.3. Przedziały ufności dla zbioru współczynników wektora β	15
1.5.4. Test dla hipotezy dotyczącej podzbioru współczynników β	16
1.5.5. Przedział ufności dla σ^2	18
1.5.6. Przedział ufności dla y_i	19
1.5.7. Ortogonalność macierzy modelu	19
1.5.8. Permutacyjne testy dla parametrów modelu β i σ^2	22
1.5.9. Bootstrapowe przedziały ufności dla parametrów modelu	23
1.6. Inne metody estymacji współczynników w modelu liniowym	25
2. Przykładowe modele liniowe i ich zastosowania	29
2.1. Regresja prosta	30
2.1.1. Wprowadzenie do regresji prostej	31
2.1.2. Przykład: zależność pomiędzy wzrostem żony a męża	33
2.1.3. Przykład: zależność pomiędzy współczynnikiem GC a wielkością genomu	38
2.1.4. Zagadnienie: diagnostyka modelu liniowego	43
2.1.5. Zagadnienie: transformacje zmiennej objaśnianej	58

2.2.	Jednokierunkowa analiza wariancji	61
2.2.1.	Wprowadzenie do jednokierunkowej analizy wariancji	62
2.2.2.	Przykład: ostra białaczka szpikowa	64
2.2.3.	Przykład: najmniejsza efektywna dawka	71
2.2.4.	Zagadnienie: testy <i>post hoc</i>	74
2.2.5.	Zagadnienie: testowanie jednorodności wariancji w grupach	81
2.2.6.	Zagadnienie: analiza kontrastów	82
2.3.	Analiza wariancji dwu- i wielokierunkowa	88
2.3.1.	Wprowadzenie do dwukierunkowej analizy wariancji	88
2.3.2.	Przykład: genetyczne podłoże schizofrenii	92
2.3.3.	Zagadnienie: model addytywny a model z interakcją	92
2.4.	Hierarchiczna analiza wariancji	103
2.4.1.	Wprowadzenie do hierarchicznej analizy wariancji	103
2.4.2.	Przykład: badanie ECAP	106
2.5.	Analiza kowariancji	111
2.5.1.	Wprowadzenie do analizy kowariancji	111
2.5.2.	Przykład: badanie endometriozy	113
2.6.	Regresja liniowa z wieloma zmiennymi objaśniającymi	117
2.6.1.	Wprowadzenie do regresji liniowej z wieloma zmiennymi objaśniającymi	117
2.6.2.	Zagadnienie: kolejność testowania	118
2.6.3.	Zagadnienie: wybór zmiennych w modelu	121
2.6.4.	Zagadnienie: modele z p bliskim n	124
2.6.5.	Zagadnienie: współliniowość zmiennych objaśniających	126
2.6.6.	Przykład: zależność pomiędzy pracą nerki, poziomem elastazy a innymi zmiennymi zależnymi	127
2.6.7.	Przykład: zależność ceny metra kwadratowego mieszkania od parametrów tego mieszkania	130
2.6.8.	Przykład: zależność pomiędzy genotypem a kątem zwinięcia ssawki u muszek owocowych	135
2.6.9.	Zagadnienie: strategie przeszukiwania listy modeli w poszukiwaniu najlepszego	137
3.	Modele mieszane — wprowadzenie, podstawowe twierdzenia i wzory	144
3.1.	Wprowadzenie	144
3.2.	Model	145
3.3.	Metoda największej wiarygodności ML i metoda resztowej największej wiarygodności REML	147
3.4.	Estymatory największej wiarygodności i resztowej wiarygodności	148
3.4.1.	Metoda estymacji z użyciem algorytmu Newtona–Rapshona	149
3.4.2.	Metoda estymacji z wykorzystaniem operacji na macierzach rzadkich	152
3.4.3.	Szczególna postać macierzy V	155
3.5.	Równania Hendersona i rozkłady estymatorów	156
3.5.1.	Równania Hendersona	157
3.5.2.	Rozkłady ocen efektów β i \tilde{u}	158
3.5.3.	Rozkład estymatora parametru θ	159
3.6.	Testy dla efektów losowych i stałych	159
3.6.1.	Testy dla efektów stałych	159
3.6.2.	Testy dla komponentów wariancyjnych	161

4. Przykładowe modele mieszane i ich zastosowania	162
4.1. Model mieszany z jednym komponentem wariacyjnym	162
4.1.1. Wprowadzenie do modelu z jednym komponentem wariacyjnym, jedna zmienna grupująca	163
4.1.2. Przykład: mleczność krów	164
4.1.3. Przykład: efekt stały genu i jeden komponent wariacyjny	169
4.1.4. Przykład: interakcja efektów środowiskowego i genetycznego a badania mikromacierzowe	177
4.2. Model mieszany z dwoma komponentami wariacyjnymi, dwie zmienne grupujące	187
4.2.1. Wprowadzenie do modelu z dwoma komponentami wariacyjnymi, dwie zmienne grupujące	188
4.2.2. Przykład: EUNOMIA <i>study</i> — zależność pomiędzy liczbą hospitalizacji a stanem pacjenta	190
4.3. Model mieszany z dwoma komponentami wariacyjnymi, jedna zmienna grupująca	205
4.3.1. Wprowadzenie do modelu z dwoma komponentami wariacyjnymi, jedna zmienna grupująca	205
4.3.2. Przykład: metaanaliza danych dotyczących ośpienia	207
4.4. Hierarchiczny model mieszany	211
4.4.1. Wprowadzenie do modelu hierarchicznego z dwoma komponentami wariacyjnymi	211
4.4.2. Przykład: badanie EDEN — efekt trybu leczenia i efekt lekarza badającego	214
4.5. Model mieszany w analizie pomiarów powtarzanych w czasie (ang. <i>longitudinal data</i>)	222
4.5.1. Wprowadzenie do analizy danych z pomiarami powtarzonymi w czasie	223
4.5.2. Przykład: funkcjonowanie nerki po przeszczepie	225
4.6. Model mieszany i zadane struktury macierzy kowariancji	237
4.6.1. Wprowadzenie do modelu mieszanego z zadanymi strukturami kowariancji	237
4.6.2. Przykład: parametry biomechaniczne mięśni	239
4.6.3. Przykład: badanie cen mieszkań w powiązaniu z lokalizacją przestrzenną (efekt przestrzenny)	245
5. Lista funkcji programu R do analizy modeli liniowych	258
5.1. Formuły	258
5.2. Modele liniowe z efektami stałymi i losowymi	265
5.2.1. Czas działania funkcji do estymacji parametrów w modelu	272
5.2.2. Szczegółowy opis funkcji <code>lm()</code> , <code>glm()</code> i <code>aov()</code>	275
5.2.3. Szczegółowy opis funkcji <code>lme()</code> , <code>lmer()</code> , <code>lmeKin()</code>	276
6. Charakterystyki zbiorów danych użytych w tej książce	279
6.1. Badanie wzrostu w małżeństwie	279
6.2. Badanie zależności między procentową zawartością GC a wielkością genomu	279
6.3. Badanie wpływu analogów witaminy D ₃ na ostrą białaczkę szpikową	281
6.4. Badanie wpływu dawki leku na reakcję organizmu	282
6.5. Badanie genetycznego podłoża schizofrenii	284
6.6. Badanie Epidemiologii Chorób Alergicznych w Polsce (ECAP)	285
6.7. Badanie ekspresji receptorów α i β u pacjentek chorych na endometriozę	287
6.8. Badanie zależności funkcji nerki od poziomu elastazy	288
6.9. Badanie czynników wpływających na cenę metra kwadratowego mieszkania	289
6.10. Badanie mleczności krów	290

6.11. Badanie efektu chłodu i linii komórkowej w eksperymentach mikromacierzowych	292
6.12. Badanie EUNOMIA i poziom psychotyczności	296
6.13. Badanie wpływu wieku i płci na występowanie otępienia	301
6.14. Badanie EDEN i efektywność oddziałów dziennych	304
6.15. Badanie funkcji nerki po przeszczepie	306
6.16. Badanie parametrów biomechanicznych mięśnia udowego	310
Dodatek	312
D.1. Uogólniona odwrotność	312
D.2. Dekompozycja na wartości osobliwe (ang. <i>singular value decomposition</i>)	313
D.3. Dekompozycja LU (ang. <i>LU decomposition</i>)	313
D.4. Dekompozycja Choleskiego (ang. <i>Cholesky decomposition</i>)	314
D.5. Dekompozycja LDM (ang. <i>LDM decomposition</i>)	314
D.6. Dekompozycja LDL (ang. <i>LDL decomposition</i>)	314
D.7. Dekompozycja QR (ang. <i>QR decomposition</i>)	314
D.8. Dekompozycja spektralna (na wartości własne i wektory własne)	315
D.9. Iloczyn Kroneckera	315
Bibliografia	316
Skorowidz	319

Przedmowa

Zanim opowiem, o czym i dla kogo jest ta książka oraz czego można z niej się nauczyć, wyjaśnię, dlaczego w ogóle ją napisałem.

Mniej więcej 10–11 lat temu wybrałem się na seminarium „Statystyka w genetyce” prowadzone przez prof. Witolda Kloneckiego w Instytucie Matematyki Politechniki Wrocławskiej. Byłem już po kilku kursach statystyki, ale z jakiegoś powodu w żaden sposób mnie ta dziedzina nie pociągała. Wyglądała jak zbiór regulek, z których trzeba wybrać odpowiednią, żeby coś policzyć, i koniec. Seminarium prof. Kloneckiego wpłynęło na całkowitą zmianę tego poglądu, statystyka stała się narzędziem do podglądania świata. Analiza danych przypomina pracę detektywa, czasami nudną, gdy trzeba godzinami przygotowywać dane i poprawić literówki lub zamienić kropki na przecinki, ale w wielu momentach fascynującą, pełną zagadek, niespodzianek i trudnych pytań. Prawdziwa analiza danych, prowadzona wspólnie z ekspertem dziedzinowym, czy to lekarzem, biologiem, czy ekonomistą, może być źródłem świetnej zabawy.

W ciągu tych 10 lat zdarzyło mi się wielokrotnie pracować z tzw. „nie-matematykami”. Wielokrotnie stwierdzałem, że znacznie łatwiej byłoby mi z nimi rozmawiać, gdybym wyjaśniając, jak działa określone narzędzie statystyczne, miał pod ręką historyjkę przedstawiającą jego przykładowe wiarygodne zastosowanie. Stąd pomysł na napisanie zbioru opowiadań, które będą przedstawiały listę zagadnień. Każde opowiadanie naświetli inny problem, każde opowiadanie można przedstawić rzeczonemu „nie-matematykowi” i powiedzieć, na czym polega problem. Ta książka miała mieć podtytuł „Książka kucharska z elementami przygodowymi”, tyle że nie spodobał się on recenzentom. Tytuł byłby o tyle odpowiedni, że w zamierzeniu przedstawia ona zbiór zagadnień oraz informacje jakich funkcji z programu R należy użyć, by dane zagadnienie przeanalizować, a czasami też rozwiązać. Elementy przygodowe biorą się z przykładów rzeczywistych danych zbieranych przez polskie zespoły badawcze. Zamiast uczyć się, jak rozróżniać płatki irysa, będziemy rozwiązywać problemy, które kiedyś może pozwolą na wydłużenie czasu funkcjonowania

przeszczepu u pacjenta lub umożliwią skuteczniejszą walkę z białaczką, zwiększenie produkcji mleka lub obniżenie kosztów leczenia pacjentów przy zachowaniu skuteczności.

Dlaczego akurat modele liniowe z efektami stałymi i mieszanymi? Jest to wynikiem współpracy z firmą Netezza Polska. Jakiś czas temu w czasie wakacji pomagałem przy implementacji estymatora parametrów w modelu mieszanym w równoległym środowisku. Zaskoczyło mnie wtedy, jak niewiele jest w Polsce miejsc, gdzie modele mieszane z efektami losowymi są omawiane podczas kursów statystyki, i jak inaczej wygląda to na zachodzie, gdzie znajomość modeli mieszanych przez studenta statystyki jest czymś oczywistym. Zdziwił mnie również brak polskojęzycznej literatury poświęconej aplikacjom modeli mieszanych. Stwierdziłem więc, że napiszę kilkunastostronicowe wprowadzenie do modeli mieszanych. Z biegiem czasu i wskutek różnych rozmów to wprowadzenie przerodziło się w książkę.

Do kogo adresowana jest ta książka? Napisałem coś, co sam chciałem przeczytać kilka lat temu. Zakładałem, że czytelnik ma podstawową wiedzę o statystyce i programowaniu w programie R. Zakładałem też, że jest zainteresowany zastosowaniami i chce wiedzieć, kiedy jakiego narzędzia użyć, ale chce wiedzieć też, jak wykorzystywane narzędzia działają. Statystyka to nie czarna skrzynka, którą strach otworzyć, ale zbiór pomysłówych obserwacji pozwalających na kontrolowaną analizę danych.

Staralem się, by ta książka była użyteczna zarówno dla bardziej zmatematyzowanych lekarzy czy biologów, jak i dla statystyków, matematyków czy informatyków, którzy lubią analizę danych. Realizacja tego planu wymagała różnych kompromisów. Nie jest to klasyczny podręcznik statystyki matematycznej. Nie wyprowadzam tutaj dowodów twierdzeń, nie przytaczam nawet twierdzeń (za wyjątkiem kilku najważniejszych). Teoria matematyczna związana z modelami liniowymi zarówno o efektach stałych, jak stałych i mieszanych jest przedstawiona w dwóch rozdziałach 1 i 3. Prezentując teorię, koncentruję się raczej na aspekcie obliczeniowym umożliwiającym stosowanie algorytmów dla dużych zbiorów danych. Osoby zainteresowane przykładami, lecz niespecjalnie zainteresowane teorią, mogą te rozdziały pominąć przy pierwszym czytaniu. Zakładałem, że Czytelnik ma podstawową wiedzę o statystyce i dlatego nie tłumaczę takich pojęć jak estymator, test czy poziom istotności. Nawet jeżeli nie jest biegły zaznajomiony z tymi pojęciami, powinien nabyć intuicji z kontekstu opisywanych wyników.

One difficulty in constructing lectures or a book at this level is that most genuine data offer too many statistical problems simultaneously and often also involve biology that is unfamiliar to a first-year student. Finney, podestane przez T.C.

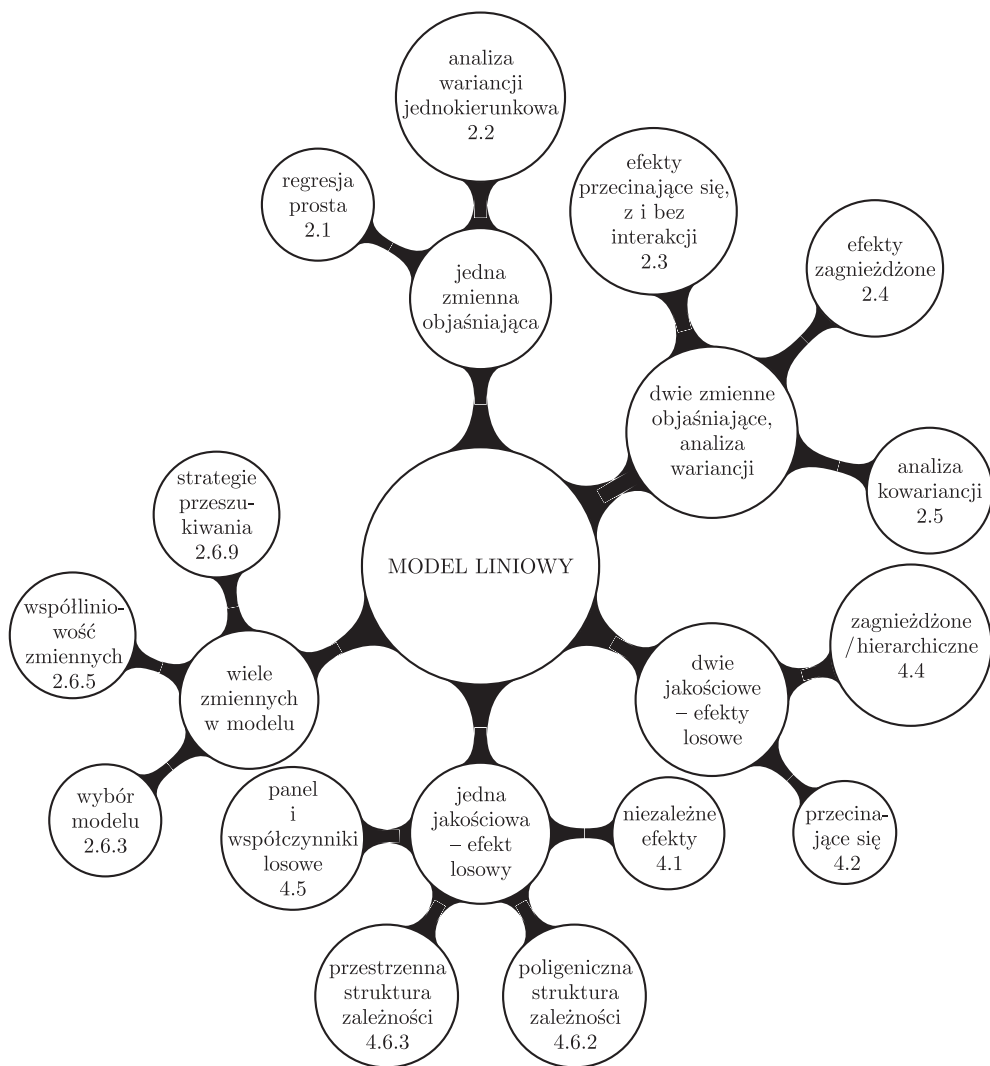
intuicji z kontekstu opisywanych wyników.

Nacisk położyłem na zainteresowanie tematem i zrozumienie. Chciałbym przekonać Czytelnika, że rozróżnienie między efektem stałym a losowym ma znaczenie, że diagnostyka to rzecz potrzebna, że transformacje zmiennych to chleb powszedni, że modele konstruuje

się w sposób iteracyjny. Rzadko się zdarza, że analiza danych to wykonanie testu t -Studenta w dwóch grupach albo zastosowanie modelu mieszanego do gotowego zbioru danych. Detektyw spędzi wiele godzin, czuwając w zaroślach, przemierza-

jąc kanały pod miastem, zanim pojawi się odpowiedni moment do przeprowadzenia akcji. Znaczna część praktyków uważa, że 80% czasu poświęconego na analizę danych zabiera przygotowanie danych, samo stosowanie metod statystycznych pochłania jedynie 20% czasu. Pozostała część praktyków uważa, że na przygotowanie danych trzeba poświęcić przynajmniej 90% czasu.

Na poniższym diagramie są przedstawione zagadnienia opisane w tej książce. Numery obok zagadnień to numery rozdziałów, w których są one omówione.



Wszystkie przykłady prezentowane w tej książce są wykonane z użyciem programu R. Jest to świetny, darmowy program do analiz statystycznych, zdobywający wciąż coraz większą popularność. Zakładam, że Czytelnik zna podstawy R.

W Internecie można znaleźć wiele stron z materiałami do nauki programu R, linki do wybranych materiałów umieszczone są na stronie <http://biecek.pl/R/>. W języku polskim dostępnych jest kilka książek poświęconych samemu językowi R lub też analizie danych z użyciem programu R, na przykład do poznania programu R od podstaw do zaawansowanych zastosowań może służyć książka [Biecek 2011]. Wszystkie dane, z których korzystałem w tej książce, znajdują się w opracowanym przeze mnie pakiecie PBI`misc`. Pakiet ten jest dostępny na CRAN-ie i można go zainstalować z konsoli programu R poleceniem

```
install.packages("PBImisc")
```

Na koniec chciałbym podziękować osobom, bez których ta książka by nie powstała. Przede wszystkim mojej żonie Karolinie, która po raz kolejny przymknęła oko na dziwne hobby męża, wierząc jedynie, że kiedyś się to skończy. Jestem wdzięczny Witoldowi Kloneckiemu, który jako pierwszy próbował nauczyć mnie, czym są modele mieszane. Dziękuję wielu osobom, z którymi pracowałem i które miały znaczny wpływ na moje wyobrażenie statystyki: profesorowi Stanisławowi Cebrowi, Pawłowi Mackiewiczowi (z którymi pracowałem w zakładzie Genomiki Uniwersytetu Wrocławskiego), profesor Teresie Ledwinie (pod której kierunkiem miałem przyjemność pracować w Instytucie Matematycznym Polskiej Akademii Nauk), profesorowi Jerzemu Tiurynowi (który zwerbował mnie do pracy w grupie Biologii Obliczeniowej na Uniwersytecie Warszawskim), profesorowi Janowi Mielniczukowi, profesorowi Tadeuszowi Calińskiemu, profesorowi Jackowi Koronackiemu, profesorowi Joannie Szydzie, profesorowi Tomaszowi Burzykowskiego, doktorom Andrzejowi Dąbrowskiemu i Andrzejowi Michalskiemu (którzy zgodzili się przeczytać wstępną bardzo surową wersję pierwszych rozdziałów). Jednocześnie muszę zaznaczyć, że nie ponoszą oni żadnej odpowiedzialności za moje ewentualne wpadki w tej książce, ponieważ wielu ich dobrych rad nie posłuchałem, mając w głowie wizję książki kucharskiej z elementami przygodowymi. Osobne serdeczne podziękowania składam Adamowi Zagdańskiemu i Arturowi Suchwałce, z którymi współpracowałem we Wrocławiu i którzy mieli znaczny wpływ na spalenie mojego postrzegania statystyki. Chciałbym też podziękować wszystkim osobom, z którymi pracowałem nad analizą danych, w szczególności danych przedstawionych w tej książce — dużo się nauczyłem i też uświadomiłem sobie, ile jeszcze trzeba się nauczyć. Wiele też nauczyłem się, prowadząc kurs „Modele liniowe i mieszane na przykładach z biologii i medycyny” na wydziale MIMUW. Dziękuję studentom z tego kursu za wyrozumiałość i wiele uwag, którymi się podzielili, a w szczególności Oldze Kowalczuk, Aleksandrze Maj, Michałowi Lisowi i Janowi Matuszewskiemu, których krytyczne uwagi i komentarze były bardzo przydatne. Chciałbym również gorąco podziękować redaktor Izabeli Mice za włożony wysiłek w pracę nad książką.

Modele liniowe — wprowadzenie, podstawowe twierdzenia i wzory

1.1. Wprowadzenie

Modele liniowe to jedna z najstarszych i najpopularniejszych metod modelowania zależności między zbiorem zmiennych *objaśniających* a zmienną ilościową nazywaną zmienną *objaśnianą*. Zależność tę modeluje się zwykle po to, by móc szacować punktowo lub przedziałowo zmienną objaśnianą na podstawie zmiennych objaśniających, lub też po to, by lepiej zrozumieć zależności między obserwowanymi zmiennymi. Model liniowy można wykorzystać również do oceny, które zmienne objaśniające i w jaki sposób są zależne od zmiennej objaśnianej.

Za początek modeli liniowych uznać można prace z przełomu XVIII i XIX wieku dotyczące metody najmniejszych kwadratów, wykorzystywanej wówczas w nawigacji oraz astronomii. Nie ma jednoznacznej odpowiedzi na pytanie, kto pierwszy wymyślił tę metodę — do dziś problem pierwszeństwa rozpala dyskusję wśród historyków matematyki. Pierwsza opublikowana praca poświęcona tej metodzie pochodzi z roku 1805 i jest autorstwa francuskiego matematyka Adrien-Marie Legendre’a. Jednak Johann Carl Friedrich Gauss twierdził, że używał tej metody od ok. 1794 lub 1795 roku (Gauss miał wtedy 18 lat!) do obliczeń na potrzeby swoich prac dotyczących astronomii. W 1809 roku Gauss opisał metodę najmniejszych kwadratów w swojej książce, podając liczne przykłady zastosowań oraz udowadniając właściwości tej metody w sytuacji, gdy zakłócenie losowe ma rozkład normalny. Legendre wcale nie przyjął do wiadomości oświadczenia o pierwszeństwie Gaussa, a spory o to, kto wymyślił metodę najmniejszych kwadratów, trwają do dziś. O historii tego sporu można przeczytać np. w artykułach „C. F. Gauss and the Theory of Errors” [Sneytin 1979] lub „Gauss and the Invention of Least Squares” [Stigler 1981].

Termin *regresja* został użyty po raz pierwszy przez Sir Francis Galtona jeszcze w XIX wieku. Francis Galton był człowiekiem renesansu, zajmującym się zarówno

antropologią, statystyką, jak i wieloma innymi naukami. W 1886 roku opublikował pracę „Regression towards mediocrity in hereditary stature” w *The Journal of the Anthropological Institute of Great Britain and Ireland* [Galton 1886], w której przedstawił wyniki badań nad dziedziczeniem wzrostu. Galton zauważył, że synowie bardzo wysokich ojców są średnio wyżsi niż synowie niższych ojców, ale ich średni wzrost jest niższy niż średni wzrost ojców. Podobnie wzrost synów niskich ojców jest bliższy średniej w populacji. To zjawisko nazwał „dążeniem do przeciętności” lub „tendencją do przeciętności” (w języku angielskim ten współczynnik nazywa się ang. *regression toward the mean*). Galton zaproponował równanie opisujące zależność między wzrostem synów i ojców lub równoważnie opisujące regresję wzrostu z pokolenia na pokolenie w kierunku wartości przeciętnej (dziś nazwalibyśmy to równanie równaniem regresji). Używając tego równania, Galton

Za słownikiem języka polskiego PWN: *regres* [fr. *regres*, łac. *regressus*, odejście], *cofanie się w rozwoju*; *przeciwnieństwo postępu*.

wyzaczył współczynnik dziedziczenia takiej cechy jak wzrost. Od tej pracy przyjęła się nazwa regresja na równania opisujące zależności między zmiennymi i nazwa ta upowszechniła się tak bardzo, że dziś jest stosowana nawet w sytuacjach, które nie mają żadnego związku z oryginalnym znaczeniem słowa *regres*.

Jako ciekawostkę można dodać, że zbiór danych zebrany przez Galtona jest dostępny w programie R pod nazwą `galton` w pakiecie `UsingR`. W tym zbiorze danych w jednej kolumnie jest średni wzrost rodziców liczony jako średnia wzrostu ojca i 1,08 wzrostu matki (ang. *midparent*), a w drugiej kolumnie wzrost dorosłego syna. Wartości z tego zbioru są przedstawione na rys. 1.1. Przerywana linia odpowiada prostej o równaniu $y = x$. Jak można zauważyć, w każdej grupie wzrostowej rodziców średni wzrost synów jest bliższy średniej populacyjnej.

Metody analizy modeli regresyjnych, głównie modeli o strukturze liniowej, zostały znacznie rozwinięte na początku XX wieku przez Karla Pearsona i jemu współczesnych. Dziś, wiek później, regresja jest jednym z najbardziej popularnych narzędzi modelowania statystycznego. Regresja jest popularna, ponieważ pozwala na opisanie związku między zmiennymi objaśniającymi a zmienną objaśnianą, oszacowanie średniej wartości zmiennej objaśnianej w zależności od zmiennych objaśniających, a także wybranie zmiennych istotnie zależnych od zmiennej objaśnianej. Postulowany w modelu związek między zmiennymi może mieć różnoraki charakter. Model regresji liniowej można relatywnie łatwo rozszerzyć na uogólniony model liniowy (w tej rodzinie jest bardzo popularna regresja logistyczna) lub na model nieliniowy (np. model GAM, modele krzywych sklepanych). W tej książce poświęcimy dużo uwagi modelom liniowym. Jeżeli zostaną one odpowiednio głęboko zrozumiane, to znacznie łatwiej będzie poznać bardziej elastyczne modele.

W podrozdziale 1.2 przedstawimy ogólną postać modelu liniowego, w podrozdziale 1.3 — najpopularniejsze metody estymacji współczynników w modelu liniowym, w podrozdziale 1.5 — wybrane metody testowania hipotez dotyczących współczynników w modelu liniowym oraz konstrukcji przedziałów ufności dla tych współczynników.