

# Eksploracja danych

*true*

*2019-04-16*



# Spis treści

<b>Wstęp</b>	<b>5</b>
O książce . . . . .	5
Zakres przedmiotu . . . . .	5
Zakres technik stosowanych w data mining . . . . .	5
Etapy eksploracji danych . . . . .	7
<b>1 Import danych</b>	<b>9</b>
<b>2 Przygotowanie danych</b>	<b>17</b>
2.1 Identyfikacja braków danych . . . . .	18
2.2 Zastępowanie braków danych . . . . .	21
<b>3 Podział metod data mining</b>	<b>35</b>
3.1 Rodzaje wnioskowania . . . . .	35
3.2 Modele regresyjne . . . . .	36
3.3 Modele klasyfikacyjne . . . . .	37
3.4 Modele grupujące . . . . .	37
<b>4 Drzewa decyzyjne</b>	<b>39</b>
4.1 Węzły i gałęzie . . . . .	39
4.2 Rodzaje reguł podziału . . . . .	39
4.3 Algorytm budowy drzewa . . . . .	41
4.4 Kryteria zatrzymania . . . . .	42
4.5 Reguły podziału . . . . .	42
4.6 Przycinanie drzewa decyzyjnego . . . . .	43
4.7 Obsługa braków danych . . . . .	45
4.8 Zalety i wady . . . . .	45
4.9 Inne algorytmy budowy drzew decyzyjnych implementowane w <b>R</b> . . . . .	53
<b>5 Pochodne drzew decyzyjnych</b>	<b>61</b>
5.1 Bagging . . . . .	61
5.2 Lasy losowe . . . . .	71
5.3 Boosting . . . . .	72
<b>6 Klasyfikatory liniowe</b>	<b>75</b>
6.1 Reprezentacja progowa . . . . .	75
6.2 Reprezentacja logitowa . . . . .	76
6.3 Wady klasyfikatorów liniowych . . . . .	76
<b>7 Regresja logistyczna</b>	<b>77</b>
7.1 Model . . . . .	77
7.2 Estymacja parametrów modelu . . . . .	77
7.3 Interpretacja . . . . .	78

<b>8</b>	<b>Analiza dyskryminacyjna</b>	<b>81</b>
8.1	Liniowa analiza dyskryminacyjna Fisher’a . . . . .	81
8.2	Liniowa analiza dyskryminacyjna - podejście probabilistyczne . . . . .	87
8.3	Analiza dyskryminacyjna metodą częściowych najmniejszych kwadratów . . . . .	95
8.4	Regularyzowana analiza dyskryminacyjna . . . . .	97
8.5	Analiza dyskryminacyjna mieszana . . . . .	98
8.6	Elastyczna analiza dyskryminacyjna . . . . .	98
<b>9</b>	<b>Klasyfikatory bayesowskie</b>	<b>101</b>
9.1	Klasyfikator maximum a posteriori (MAP) . . . . .	101
9.2	Klasyfikator największej wiarygodności (ML) . . . . .	102
9.3	Naiwny klasyfikator Bayesa (NB) . . . . .	102
9.4	Zalety i wady . . . . .	106
<b>10</b>	<b>Metoda <math>k</math> najbliższych sąsiadów</b>	<b>107</b>
<b>11</b>	<b>Uogólnione modele addytywne</b>	<b>113</b>
11.1	Przypadek jednowymiarowy . . . . .	113
11.2	Przypadek wielowymiarowy . . . . .	116
11.3	Uogólnione modele addytywne . . . . .	117

# Wstęp

## O książce

Niniejsza książka powstała na bazie doświadczeń autora, a głównym jej celem jest przybliżenie czytelnikowi podstaw z dziedziny *Data mining* studentom kierunku *Matematyka* Politechniki Lubelskiej. Będzie łączyć w sobie zarówno treści teoretyczne związane z przedstawianymi etapami eksploracji danych i budową modeli, jak i praktyczne wskazówki dotyczące budowy modeli w środowisku **R** (R Core Team 2018). Podane zostaną również wskazówki, jak raportować wyniki analiz i jak dokonać właściwych ilustracji wyników. Bardzo użyteczny w napisaniu książki były pakiety programu R: **bookdown** (Xie 2018a), **knitr** (Xie 2018b) oraz pakiet **rmarkdown** (Allaire et al. 2018).

## Zakres przedmiotu

Przedmiot *Eksploracja danych* będzie obejmował swoim zakresem eksplorację i wizualizację danych oraz uczenie maszynowe. Eksploracja danych ma na celu pozyskiwanie i systematyzację wiedzy pochodzącej z danych. Odbywa się ona głównie przy użyciu technik statystycznych, rachunku prawdopodobieństwa i metod z zakresu baz danych. Natomiast uczenie maszynowe, to gałąź nauki (obejmuje nie tylko statystykę, choć to na niej się głównie opiera) dotyczącej budowy modeli zdolnych do rozpoznawania wzorców, przewidywania wartości i klasyfikacji obiektów. Data mining to szybko rosnąca grupa metod analizy danych rozwijana nie tylko przez statystyków ale również przez biologów, genetyków, cybernetyków, informatyków, ekonomistów, osoby pracujące nad rozpoznawaniem obrazów i wiele innych grup zawodowych. W dzisiejszych czasach trudno sobie wyobrazić życie bez sztucznej inteligencji. Towarzyszy ona nam w codziennym, życiu kiedy korzystamy z telefonów komórkowych, wyszukiwarek internetowych, robotów sprzątających, automatycznych samochodów, nawigacji czy gier komputerowych. Lista ta jest niepełna i stale się wydłuża.

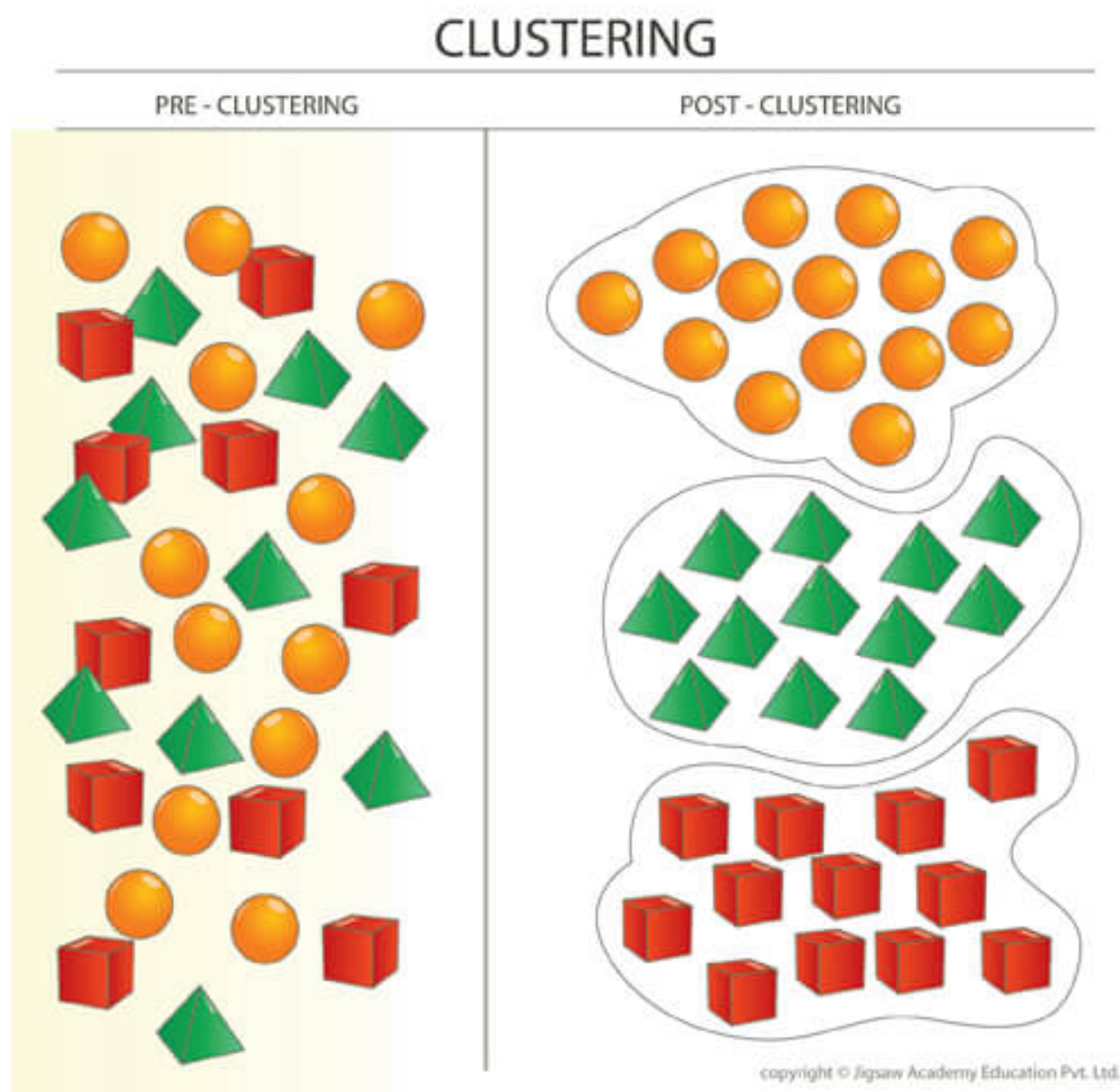
href="https://twitter.com/i/status/1091069356367200256">January 31, 2019

## Zakres technik stosowanych w data mining

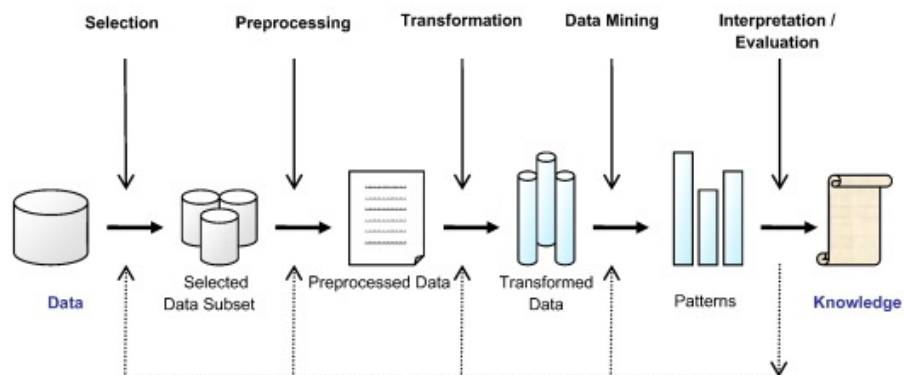
- statystyka opisowa
- wielowymiarowa analiza danych
- analiza szeregów czasowych
- analiza danych przestrzennych
- reguły asocjacji
- uczenie maszynowe<sup>1</sup>, w tym:
  - klasyfikacja
  - predykcja
  - analiza skupień
  - *text mining*
- i wiele innych

---

<sup>1</sup>ang. *machine learning*



Rysunek 1: Przykład nienadzorowanego uczenia maszynowego. Źródło: <https://analyticstraining.com/cluster-analysis-for-business/>



Rysunek 2: Etapy eksploracji danych (Kavakiotis et al. 2017)

href="https://twitter.com/i/status/1097199751072690176">Ferbruary 17, 2019

## Etapy eksploracji danych

1. Czyszczenie danych - polega na usuwaniu braków danych, usuwaniu stałych zmiennych, imputacji braków danych oraz przygotowaniu danych do dalszych analiz.
2. Integracja danych - łączenie danych pochodzących z różnych źródeł.
3. Selekcja danych - wybór z bazy tych danych, które są potrzebne do dalszych analiz.
4. Transformacja danych - przekształcenie i konsolidacja danych do postaci przydatnej do eksploracji.
5. Eksploracja danych - zastosowanie technik wymienionych wcześniej w celu odnalezienia wzorców<sup>2</sup> i zależności.
6. Ewaluacja modeli - ocena poprawności modeli oraz wzorców z nich uzyskanych.
7. Wizualizacja wyników - graficzne przedstawienie odkrytych wzorców.
8. Wdrażanie modeli - zastosowanie wyznaczonych wzorców.

---

<sup>2</sup>ang. *patterns*





# Rozdział 1

## Import danych

Środowisko **R** pozwala na import i export plików o różnych rozszerzeniach (`txt`, `csv`, `xls`, `xlsx`, `sav`, `xpt`, `dta`, itd.)<sup>1</sup>. W tym celu czasami trzeba zainstalować pakiety rozszerzające podstawowe możliwości R-a. Najnowsza<sup>2</sup> wersja programu RStudio (v. 1.1.463)<sup>3</sup> pozwala na wczytanie danych z popularnych źródeł za pomocą GUI.

Jeśli dane są zapisane w trybie tekstowym (np. `txt`, `csv`), to wczytujemy je w następujący sposób

```
dane1 <- read.table("data/dane1.txt", header = T)
head(dane1)
```

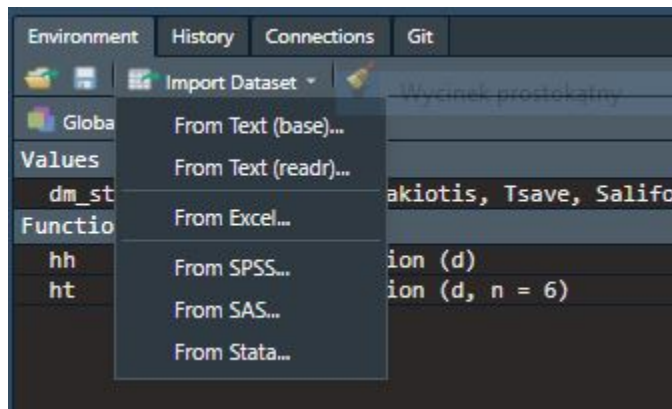
```
##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1         5.1         3.5         1.4         0.2   setosa
## 2         4.9         3.0         1.4         0.2   setosa
## 3         4.7         3.2         1.3         0.2   setosa
## 4         4.6         3.1         1.5         0.2   setosa
## 5         5.0         3.6         1.4         0.2   setosa
## 6         5.4         3.9         1.7         0.4   setosa
```

```
dane2 <- read.csv2("data/dane1.csv", header = T)
head(dane2)
```

<sup>1</sup>póki co nie jest mi znana funkcja pozwalająca na import plików programu Statistica

<sup>2</sup>na dzień 19.02.2019

<sup>3</sup>istnieją również nowsze wersje deweloperskie



Rysunek 1.1: Narzędzie do importu plików programu RStudio

	A	B	C	D	E	F
1	Długość kielic	Szerokość	Długość p	Szerokość	Gatunki	
2	5,1	3,5	1,4	0,2	setosa	
3	4,9	3	1,4	0,2	setosa	
4	4,7	3,2	1,3	0,2	setosa	
5	4,6	3,1	1,5	0,2	setosa	
6	5	3,6	1,4	0,2	setosa	
7	5,4	3,9	1,7	0,4	setosa	
8	-	-	1,4	0,3	setosa	
9	5	3,4	1,5	0,2	BD	
10	4,4	2,9	1,4	0,2	setosa	
11	4,9	3,1	1,5	0,1	setosa	
12	5,4	3,7	1,5	0,2	setosa	
13	4,8	3,4	1,6	0,2	setosa	

Rysunek 1.2: Fragment pliku Excel

```
## Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1 5.1 3.5 1.4 0.2 setosa
## 2 4.9 3.0 1.4 0.2 setosa
## 3 4.7 3.2 1.3 0.2 setosa
## 4 4.6 3.1 1.5 0.2 setosa
## 5 5.0 3.6 1.4 0.2 setosa
## 6 5.4 3.9 1.7 0.4 setosa
```

```
# funkcja pakietu readr wczytuje plik jako ramkę danych w formacie tibble
# pakiet readr jest częścią większego pakietu tidyverse,
# który został wczytany wcześniej
dane3 <- read_csv2("data/dane1.csv")
dane3
```

```
## # A tibble: 150 x 5
## Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## <dbl> <dbl> <dbl> <dbl> <chr>
## 1 5.1 3.5 1.4 0.2 setosa
## 2 4.9 3 1.4 0.2 setosa
## 3 4.7 3.2 1.3 0.2 setosa
## 4 4.6 3.1 1.5 0.2 setosa
## 5 5 3.6 1.4 0.2 setosa
## 6 5.4 3.9 1.7 0.4 setosa
## 7 4.6 3.4 1.4 0.3 setosa
## 8 5 3.4 1.5 0.2 setosa
## 9 4.4 2.9 1.4 0.2 setosa
## 10 4.9 3.1 1.5 0.1 setosa
## # ... with 140 more rows
```

Jeśli dane są przechowywane w pliku Excel (np. `xlsx`), to importujemy je za pomocą funkcji `read_excel` pakietu `readxl`. Domyślnie jest wczytywany arkusz pierwszy ale jeśli zachodzi taka potrzeba, to można ustalić, który arkusz pliku Excel ma być wczytany za pomocą paramteru `sheet`, np. `sheet=3`, co oznacza, że zostanie wczytany trzeci arkusz pliku.

Ponieważ w pliku `dane1.xlsx` braki danych zostały zakodowane znakami BD oraz -, to należy ten fakt przekazać funkcji, aby poprawnie wczytać braki danych. W przeciwnym przypadku zmienne zawierające braki tak kodowane, będą wczytane jako zmienne znakowe.

```
library(readxl)
dane4 <- read_excel("data/dane1.xlsx", na = c("BD", "-"))
dane4

## # A tibble: 150 x 5
##   `Długość kielic` `Szerokość kielic` `Długość płatka` `Szerokość płat~
##   <dbl>          <dbl>          <dbl>          <dbl>
## 1         5.1         3.5         1.4         0.2
## 2         4.9         3         1.4         0.2
## 3         4.7         3.2         1.3         0.2
## 4         4.6         3.1         1.5         0.2
## 5          5         3.6         1.4         0.2
## 6         5.4         3.9         1.7         0.4
## 7          NA         NA         1.4         0.3
## 8          5         3.4         1.5         0.2
## 9         4.4         2.9         1.4         0.2
## 10        4.9         3.1         1.5         0.1
## # ... with 140 more rows, and 1 more variable: Gatunki <chr>
```

Istnieją oczywiście jeszcze wiele innych formatów danych, charakterystycznych dla programów, w których są traktowane jako domyślne.<sup>4</sup> W szczególny sposób należy zwrócić uwagę na pliki o rozszerzeniu `RData` lub `rda`<sup>5</sup> oraz pliki `rds`. Pliki `rda` służą do przechowywania obiektów programu R. Mogą to być pliki danych ale również obiekty graficzne (typu wyniki funkcji `ggplot`), modele (np. wynik funkcji `lm()`), zdefiniowane funkcje i wszystkie inne obiekty, które da się zapisać w środowisku R. Ponadto pliki `rda` pozwalają na zapisanie wielu obiektów w jednym pliku. Pliki o rozszerzeniu `rds` mają podobną funkcję z tym, że pozwalają na przechowywanie tylko jednego obiektu.

```
# wszystkie wczytane wcześniej pliki zapisuje w jednym pliku
save(dane1, dane2, dane3, dane4, file = "data/dane.rda")
# plik rda został zapisany
list.files(path = "data/")

## [1] "algae.csv"      "Analysis.txt"  "dane.rda"      "dane1.csv"
## [5] "dane1.txt"      "dane1.xlsx"    "dane4.rds"     "dane4.sav"

# usuwam dane ze środowiska R
rm(dane1, dane2, dane3, dane4)
# sprawdzam co jest wczytane do R
ls()

## character(0)

# wczytuję plik rda
load("data/dane.rda")
# jeszcze raz sprawdzam co jest wczytane do R
ls()

## [1] "dane1" "dane2" "dane3" "dane4"
```

Zapisując obiekty jako oddzielne pliki, można przy wczytywaniu nadawać im nazwy.

```
rm(dane1, dane2, dane3)
ls()
```

<sup>4</sup>do ich wczytywania stosujemy funkcje pakietu `foreign`

<sup>5</sup>oznaczają to samo

```
## [1] "dane4"
```

```
saveRDS(dane4, file = "data/dane4.rds")
nowe_dane <- readRDS("data/dane4.rds")
nowe_dane
```

```
## # A tibble: 150 x 5
##   `Długość kielic~` `Szerokość kielic~` `Długość płatka~` `Szerokość płat~`
##   <dbl>           <dbl>           <dbl>           <dbl>
## 1         5.1         3.5         1.4         0.2
## 2         4.9         3         1.4         0.2
## 3         4.7         3.2         1.3         0.2
## 4         4.6         3.1         1.5         0.2
## 5         5         3.6         1.4         0.2
## 6         5.4         3.9         1.7         0.4
## 7         NA         NA         1.4         0.3
## 8         5         3.4         1.5         0.2
## 9         4.4         2.9         1.4         0.2
## 10        4.9         3.1         1.5         0.1
## # ... with 140 more rows, and 1 more variable: Gatunki <chr>
```

Oprócz wielu zalet takiego sposobu importu i eksportu danych jest jedna poważna wada, pliki te można odczytać jedynie za pomocą R. Osobiście polecam stosować do importu i eksportu danych plików w takich formatach, które mogą przeczytać wszyscy. Jak dotąd widać do importu różnych formatów danych potrzebujemy różnych funkcji, czasami nawet z różnych pakietów. Istnieje rozwiązanie tej niedogodności

```
library(rio)
dane1 <- import("data/dane1.txt")
head(dane1)
```

```
##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1         5.1         3.5         1.4         0.2 setosa
## 2         4.9         3.0         1.4         0.2 setosa
## 3         4.7         3.2         1.3         0.2 setosa
## 4         4.6         3.1         1.5         0.2 setosa
## 5         5.0         3.6         1.4         0.2 setosa
## 6         5.4         3.9         1.7         0.4 setosa
```

```
dane2 <- import("data/dane1.csv", dec = ",")
# dane1.csv miały , jako znak rozdzielający cechę i mantysę liczb
# dlatego włączamy parametr dec
head(dane2)
```

```
##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1         5.1         3.5         1.4         0.2 setosa
## 2         4.9         3.0         1.4         0.2 setosa
## 3         4.7         3.2         1.3         0.2 setosa
## 4         4.6         3.1         1.5         0.2 setosa
## 5         5.0         3.6         1.4         0.2 setosa
## 6         5.4         3.9         1.7         0.4 setosa
```

```
dane3 <- import("data/dane1.xlsx", na=c("BD", "-"))
head(dane3)
```

```
##   Długość kielicha Szerokość kielicha Długość płatka Szerokość płatka
## 1         5.1         3.5         1.4         0.2
## 2         4.9         3.0         1.4         0.2
## 3         4.7         3.2         1.3         0.2
```

```
## 4      4.6      3.1      1.5      0.2
## 5      5.0      3.6      1.4      0.2
## 6      5.4      3.9      1.7      0.4
##  Gatunki
## 1  setosa
## 2  setosa
## 3  setosa
## 4  setosa
## 5  setosa
## 6  setosa
```

```
dane4 <- import("data/dane4.rds")
dane4
```

```
## # A tibble: 150 x 5
##   `Długość kielic~` `Szerokość kiel~` `Długość płatka` `Szerokość płat~`
##           <dbl>           <dbl>           <dbl>           <dbl>
## 1           5.1           3.5           1.4           0.2
## 2           4.9           3           1.4           0.2
## 3           4.7           3.2           1.3           0.2
## 4           4.6           3.1           1.5           0.2
## 5           5           3.6           1.4           0.2
## 6           5.4           3.9           1.7           0.4
## 7           NA           NA           1.4           0.3
## 8           5           3.4           1.5           0.2
## 9           4.4           2.9           1.4           0.2
## 10          4.9           3.1           1.5           0.1
## # ... with 140 more rows, and 1 more variable: Gatunki <chr>
```

Lista możliwości jaką daje nam pakiet `rio` (Chan and Leeper 2018) jest niemal nieograniczona:<sup>6</sup>

- Comma-separated data (.csv), using `fread` or, if `fread = FALSE`, `read.table` with `row.names = FALSE` and `stringsAsFactors = FALSE`
- Pipe-separated data (.psv), using `fread` or, if `fread = FALSE`, `read.table` with `sep = '|'`, `row.names = FALSE` and `stringsAsFactors = FALSE`
- Tab-separated data (.tsv), using `fread` or, if `fread = FALSE`, `read.table` with `row.names = FALSE` and `stringsAsFactors = FALSE`
- SAS (.sas7bdat), using `read_sas`.
- SAS XPORT (.xpt), using `read_xpt` or, if `haven = FALSE`, `read_xport`.
- SPSS (.sav), using `read_sav`. If `haven = FALSE`, `read_spss` can be used.
- Stata (.dta), using `read_dta`. If `haven = FALSE`, `read_dta` can be used.
- SAS XPORT (.xpt), using `read_xport`.
- SPSS Portable Files (.por), using `read_por`.
- Excel (.xls and .xlsx), using `read_excel`. Use which to specify a sheet number. For .xlsx files, it is possible to set `readxl = FALSE`, so that `read.xlsx` can be used instead of `readxl` (the default).
- R syntax object (.R), using `dget`
- Saved R objects (.RData,.rda), using `load` for single-object .Rdata files. Use which to specify an object name for multi-object .Rdata files. This can be any R object (not just a data frame).
- Serialized R objects (.rds), using `readRDS`. This can be any R object (not just a data frame).
- Epiinfo (.rec), using `read.epiinfo`
- Minitab (.mtp), using `read.mtp`
- Systat (.syd), using `read.systat`
- "XBASE" database files (.dbf), using `read.dbf`
- Weka Attribute-Relation File Format (.arff), using `read.arff`

<sup>6</sup>fragment pliku `help` funkcji `import`

- Data Interchange Format (.dif), using read.DIF
- Fortran data (no recognized extension), using read.fortran
- Fixed-width format data (.fwf), using a faster version of read.fwf that requires a widths argument and by default in rio has stringsAsFactors = FALSE. If readr = TRUE, import will be performed using read\_fwf, where widths should be: NULL, a vector of column widths, or the output of fwf\_empty, fwf\_widths, or fwf\_positions.
- gzip comma-separated data (.csv.gz), using read.table with row.names = FALSE and stringsAsFactors = FALSE
- CSVY (CSV with a YAML metadata header) using read\_csvy.
- Feather R/Python interchange format (.feather), using read\_feather
- Fast storage (.fst), using read.fst
- JSON (.json), using fromJSON
- Matlab (.mat), using read.mat
- EViews (.wfl), using readEViews
- OpenDocument Spreadsheet (.ods), using read\_ods. Use which to specify a sheet number.
- Single-table HTML documents (.html), using read\_html. The data structure will only be read correctly if the HTML file can be converted to a list via as\_list.
- Shallow XML documents (.xml), using read\_xml. The data structure will only be read correctly if the XML file can be converted to a list via as\_list.
- YAML (.yaml), using yaml.load
- Clipboard import (on Windows and Mac OS), using read.table with row.names = FALSE
- Google Sheets, as Comma-separated data (.csv)

**Przykład 1.1.** Poniższa ilustracja przedstawia fragment pliku danych `Analysis.txt` zawierającego pewne błędy, które należy naprawić na etapie importu danych. Po pierwsze brakuje w nim nazw zmiennych (choć nie widać tego na rysunku). Poszczególne kolumny nazywają się następująco: `season`, `size`, `speed`, `mxPH`, `mn02`, `C1`, `N03`, `NH4`, `oP04`, `P04`, `Ch1a`, `a1`, `a2`, `a3`, `a4`, `a5`, `a6`, `a7`. Naszym zadaniem jest import tego pliku z jednoczesną obsługą braków (braki danych są zakodowane przez `XXXXXXX`) oraz nadaniem nagłówków kolumn. Plik `Analysis.txt` jest umieszczony w katalogu `data/`. Z racji, że plik dotyczy glonów, to dane zapiszemy pod nazwą `algae`.

```
algae <- import('data/Analysis.txt', header=F,
               dec='.',
               col.names=c('season','size','speed','mxPH','mn02','C1',
                           'N03','NH4','oP04','P04','Ch1a','a1','a2',
                           'a3','a4','a5','a6','a7'),
               na.strings=c('XXXXXXX'))
head(algae)
```

```
##  season size speed mxPH mn02    C1    N03    NH4    oP04    P04 Ch1a
## 1 winter small medium 8.00  9.8 60.800  6.238 578.000 105.000 170.000 50.0
## 2 spring small medium 8.35  8.0 57.750  1.288 370.000 428.750 558.750  1.3
## 3 autumn small medium 8.10 11.4 40.020  5.330 346.667 125.667 187.057 15.6
## 4 spring small medium 8.07  4.8 77.364  2.302  98.182  61.182 138.700  1.4
## 5 autumn small medium 8.06  9.0 55.350 10.416 233.700  58.222  97.580 10.5
## 6 winter small  high 8.25 13.1 65.750  9.248 430.000  18.250  56.667 28.4
##      a1  a2  a3  a4  a5  a6  a7
## 1  0.0  0.0  0.0  0.0 34.2  8.3  0.0
## 2  1.4  7.6  4.8  1.9  6.7  0.0  2.1
## 3  3.3 53.6  1.9  0.0  0.0  0.0  9.7
## 4  3.1 41.0 18.9  0.0  1.4  0.0  1.4
## 5  9.2  2.9  7.5  0.0  7.5  4.1  1.0
## 6 15.1 14.6  1.4  0.0 22.5 12.6  2.9
```

```
summary(algae)
```

Analysis.txt — Notatnik								
Plik	Edycja	Format	Widok	Pomoc				
winter	small	high	8.30000	7.70000	50.00000	8.54300	76.00000	264
spring	small	high	8.30000	8.80000	54.14300	7.83000	51.42900	270
winter	small	high	8.40000	13.40000	69.75000	4.55500	37.50000	10
spring	small	high	8.30000	12.50000	87.00000	4.87000	22.50000	27
autumn	small	high	8.00000	12.10000	66.30000	4.53500	39.00000	16
winter	small	low	XXXXXXX	12.60000	9.00000	0.23000	10.00000	5
spring	small	medium	7.60000	9.60000	15.00000	3.02000	40.00000	
autumn	small	medium	7.29000	11.21000	17.75000	3.07000	35.00000	
winter	small	medium	7.60000	10.20000	32.30000	4.50800	192.50000	
summer	small	medium	8.00000	7.90000	27.23300	1.65100	28.33300	
winter	small	high	7.90000	11.00000	6.16700	1.17200	18.33300	7
spring	small	high	7.90000	9.00000	5.27300	0.91000	33.63600	9
winter	small	high	6.60000	10.80000	XXXXXXX	3.24500	10.00000	1
spring	small	medium	5.60000	11.80000	XXXXXXX	2.22000	5.00000	
autumn	small	medium	5.70000	10.80000	XXXXXXX	2.55000	10.00000	
spring	small	high	6.60000	9.50000	XXXXXXX	1.32000	20.00000	1
summer	small	high	6.60000	10.80000	XXXXXXX	2.64000	10.00000	2
autumn	small	medium	6.60000	11.30000	XXXXXXX	4.17000	10.00000	

Rysunek 1.3: Fragment pliku danych Analysis.txt

```
##      season              size              speed              mxPH
## Length:200      Length:200      Length:200      Min.   :5.600
## Class :character  Class :character  Class :character  1st Qu.:7.700
## Mode  :character  Mode  :character  Mode  :character  Median :8.060
##                                     Mean  :8.012
##                                     3rd Qu.:8.400
##                                     Max.   :9.700
##                                     NA's   :1
##      mnO2              Cl              NO3              NH4
## Min.   : 1.500      Min.   : 0.222      Min.   : 0.050      Min.   : 5.00
## 1st Qu.: 7.725      1st Qu.: 10.981      1st Qu.: 1.296      1st Qu.: 38.33
## Median : 9.800      Median : 32.730      Median : 2.675      Median : 103.17
## Mean   : 9.118      Mean   : 43.636      Mean   : 3.282      Mean   : 501.30
## 3rd Qu.:10.800      3rd Qu.: 57.824      3rd Qu.: 4.446      3rd Qu.: 226.95
## Max.   :13.400      Max.   :391.500      Max.   :45.650      Max.   :24064.00
## NA's   :2          NA's   :10          NA's   :2          NA's   :2
##      oP04              P04              Chla              a1
## Min.   : 1.00      Min.   : 1.00      Min.   : 0.200      Min.   : 0.00
## 1st Qu.: 15.70      1st Qu.: 41.38      1st Qu.: 2.000      1st Qu.: 1.50
## Median : 40.15      Median :103.29      Median : 5.475      Median : 6.95
## Mean   : 73.59      Mean   :137.88      Mean   :13.971      Mean   :16.92
## 3rd Qu.: 99.33      3rd Qu.:213.75      3rd Qu.:18.308      3rd Qu.:24.80
## Max.   :564.60      Max.   :771.60      Max.   :110.456      Max.   :89.80
## NA's   :2          NA's   :2          NA's   :12
##      a2              a3              a4              a5
## Min.   : 0.000      Min.   : 0.000      Min.   : 0.000      Min.   : 0.000
## 1st Qu.: 0.000      1st Qu.: 0.000      1st Qu.: 0.000      1st Qu.: 0.000
```

```
## Median : 3.000   Median : 1.550   Median : 0.000   Median : 1.900
## Mean    : 7.458   Mean    : 4.309   Mean    : 1.992   Mean    : 5.064
## 3rd Qu.:11.375   3rd Qu.: 4.925   3rd Qu.: 2.400   3rd Qu.: 7.500
## Max.    :72.600   Max.    :42.800   Max.    :44.600   Max.    :44.400
##
##          a6          a7
## Min.    : 0.000   Min.    : 0.000
## 1st Qu.: 0.000   1st Qu.: 0.000
## Median : 0.000   Median : 1.000
## Mean    : 5.964   Mean    : 2.495
## 3rd Qu.: 6.925   3rd Qu.: 2.400
## Max.    :77.600   Max.    :31.600
##
```

```
export(algae, file = "data/algae.csv")
```



## Rozdział 2

# Przygotowanie danych

Dane, które importujemy z zewnętrznego źródła najczęściej nie spełniają formatów obowiązujących w **R**. Często zmienne zawierają niedopuszczalne znaki szczególne, odstępy w nazwach, powtórzone nazwy kolumn, nazwy zmiennych zaczynające się od liczby, czy puste wiersze lub kolumny. Przed przystąpieniem do analizy zbioru należy rozważyć ewentualne poprawki nazw zmiennych, czy usunięcie pustych kolumn i wierszy. Niektórych czynności można dokonać już na etapie importu danych, stosując pewne pakiety oraz nowe funkcjonalności środowiska **RStudio**. W większości przypadków uchroni nas to od żmudnego przekształcania typów zmiennych. Oczywiście wszystkie te czynności czyszczenia danych można również dokonać już po imporcie danych, za pomocą odpowiednich komend **R**.

```
## przykładowe niepożądane nazwy zmiennych
test_df <- as.data.frame(matrix(rnorm(18),ncol = 6))
names(test_df) <- c("hIgHlo", "REPEAT VALUE", "REPEAT VALUE",
                    "% successful (2009)", "abc@!*", "")
test_df

##           hIgHlo REPEAT VALUE REPEAT VALUE % successful (2009)      abc@!*
## 1  0.62816036   -0.5201609    1.1043470      0.1738823 -0.08412172
## 2 -0.28213680   -0.6262217   -0.1659695      1.2445501  1.15968770
## 3  0.09870219   -0.6941247   -0.4375282      0.1012135  1.02700292
##
## 1  2.3730675
## 2  0.4643571
## 3 -0.6335834
```

```
## do poprawy nazw zmiennych użyjemy funkcji make.names
names(test_df) <- make.names(names(test_df))
test_df

##           hIgHlo REPEAT.VALUE REPEAT.VALUE X..successful..2009.      abc...
## 1  0.62816036   -0.5201609    1.1043470      0.1738823 -0.08412172
## 2 -0.28213680   -0.6262217   -0.1659695      1.2445501  1.15968770
## 3  0.09870219   -0.6941247   -0.4375282      0.1012135  1.02700292
##           X
## 1  2.3730675
## 2  0.4643571
## 3 -0.6335834
```

Efekt końcowy choć skuteczny, to nie jest zadowalający. Czyszczenia nazw zmiennych można też dokonać stosując funkcję `clean_names` pakietu **janitor** (Firke 2018). Pozwala on również na usuwanie pustych wierszy i kolumn, znajdowanie zduplikowanych rekordów, itp.

```
library(janitor)
test_df %>% # aby na stałe zmienić nazwy zmiennych trzeba podstawienia
  clean_names()

##      h_ig_hlo repeat_value repeat_value_2 x_successful_2009      abc
## 1  0.62816036  -0.5201609      1.1043470      0.1738823 -0.08412172
## 2 -0.28213680  -0.6262217      -0.1659695      1.2445501  1.15968770
## 3  0.09870219  -0.6941247      -0.4375282      0.1012135  1.02700292
##
##      x
## 1  2.3730675
## 2  0.4643571
## 3 -0.6335834

# przykładowe dane
x <- data.frame(w1=c(1,4,2,NA),w2=c(NA,2,3,NA), w3=c(1,NA,1,NA))
x

##    w1 w2 w3
## 1   1 NA  1
## 2   4  2 NA
## 3   2  3  1
## 4  NA NA NA

x %>% remove_empty("rows")

##    w1 w2 w3
## 1   1 NA  1
## 2   4  2 NA
## 3   2  3  1
```

## 2.1 Identyfikacja braków danych

Zanim usuniemy jakiegokolwiek braki w zbiorze, powinniśmy je najpierw zidentyfikować, określić ich charakter, a dopiero potem ewentualnie podjąć decyzję o uzupełnianiu braków.

```
algae <- rio::import("data/algae.csv")

# najprościej jest wywołać summary
summary(algae)
```

```
##      season      size      speed      mxPH
## Length:200    Length:200    Length:200    Min.   :5.600
## Class :character Class :character Class :character 1st Qu.:7.700
## Mode  :character Mode  :character Mode  :character Median :8.060
##                                     Mean  :8.012
##                                     3rd Qu.:8.400
##                                     Max.   :9.700
##                                     NA's    :1
##      mnO2      Cl      NO3      NH4
## Min.   : 1.500 Min.   : 0.222 Min.   : 0.050 Min.   :  5.00
## 1st Qu.: 7.725 1st Qu.:10.981 1st Qu.: 1.296 1st Qu.: 38.33
## Median : 9.800 Median :32.730 Median : 2.675 Median :103.17
## Mean   : 9.118 Mean   :43.636 Mean   : 3.282 Mean   :501.30
## 3rd Qu.:10.800 3rd Qu.:57.824 3rd Qu.: 4.446 3rd Qu.:226.95
## Max.   :13.400 Max.   :391.500 Max.   :45.650 Max.   :24064.00
## NA's    :2      NA's   :10      NA's   :2      NA's   :2
```

```
##          oP04          P04          Chla          a1
## Min.    : 1.00    Min.    : 1.00    Min.    : 0.200    Min.    : 0.00
## 1st Qu.: 15.70    1st Qu.: 41.38    1st Qu.: 2.000    1st Qu.: 1.50
## Median : 40.15    Median :103.29    Median : 5.475    Median : 6.95
## Mean    : 73.59    Mean    :137.88    Mean    : 13.971    Mean    :16.92
## 3rd Qu.: 99.33    3rd Qu.:213.75    3rd Qu.: 18.308    3rd Qu.:24.80
## Max.    :564.60    Max.    :771.60    Max.    :110.456    Max.    :89.80
## NA's    :2        NA's    :2        NA's    :12
##          a2          a3          a4          a5
## Min.    : 0.000    Min.    : 0.000    Min.    : 0.000    Min.    : 0.000
## 1st Qu.: 0.000    1st Qu.: 0.000    1st Qu.: 0.000    1st Qu.: 0.000
## Median : 3.000    Median : 1.550    Median : 0.000    Median : 1.900
## Mean    : 7.458    Mean    : 4.309    Mean    : 1.992    Mean    : 5.064
## 3rd Qu.:11.375    3rd Qu.: 4.925    3rd Qu.: 2.400    3rd Qu.: 7.500
## Max.    :72.600    Max.    :42.800    Max.    :44.600    Max.    :44.400
##
##          a6          a7
## Min.    : 0.000    Min.    : 0.000
## 1st Qu.: 0.000    1st Qu.: 0.000
## Median : 0.000    Median : 1.000
## Mean    : 5.964    Mean    : 2.495
## 3rd Qu.: 6.925    3rd Qu.: 2.400
## Max.    :77.600    Max.    :31.600
##
## wyświetl niekompletne wiersze
algae[!complete.cases(algae),] %>% head()
```

```
##      season size speed mxPH mnO2   Cl   NO3 NH4 oP04 P04 Chla  a1  a2  a3
## 28 autumn small  high  6.8 11.1 9.00 0.630 20  4.0  NA  2.7 30.3 1.9 0.0
## 38 spring small  high  8.0  NA  1.45 0.810 10  2.5 3.0  0.3 75.8 0.0 0.0
## 48 winter small  low   NA 12.6 9.00 0.230 10  5.0 6.0  1.1 35.5 0.0 0.0
## 55 winter small  high  6.6 10.8  NA 3.245 10  1.0 6.5   NA 24.3 0.0 0.0
## 56 spring small medium 5.6 11.8  NA 2.220  5  1.0 1.0   NA 82.7 0.0 0.0
## 57 autumn small medium 5.7 10.8  NA 2.550 10  1.0 4.0   NA 16.8 4.6 3.9
##      a4 a5 a6 a7
## 28  0.0 2.1 1.4 2.1
## 38  0.0 0.0 0.0 0.0
## 48  0.0 0.0 0.0 0.0
## 55  0.0 0.0 0.0 0.0
## 56  0.0 0.0 0.0 0.0
## 57 11.5 0.0 0.0 0.0
##
## policz niekompletne wiersze
nrow(algae[!complete.cases(algae),])

## [1] 16
## sprawdzenie liczby braków w wierszach
apply(algae, 1, function(x) sum(is.na(x)))

##      1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18
##      0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
##     19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36
##      0  0  0  0  0  0  0  0  0  1  0  0  0  0  0  0  0  0
##     37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54
```

```
## 0 1 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0
## 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72
## 2 2 2 2 2 2 2 6 1 0 0 0 0 0 0 0 0 0
## 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90
## 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108
## 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125 126
## 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0
## 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144
## 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## 145 146 147 148 149 150 151 152 153 154 155 156 157 158 159 160 161 162
## 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0
## 163 164 165 166 167 168 169 170 171 172 173 174 175 176 177 178 179 180
## 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198
## 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0
## 199 200
## 6 0
```

Wiele ciekawych funkcji do eksploracji danych znajduje się w pakiecie **DMwR** (Torgo 2013), który został przygotowany przy okazji publikacji książki *Data Mining with R*.

```
## poszukiwanie wierszy zawierających wiele braków
## w tym przypadku próg wyświetlania ustawiony jest na 0.2
## czyli 20% wszystkich kolumn
library(DMwR)
manyNAs(algae)
```

```
## 62 199
## 62 199
```

```
## tworzenie zbioru pozbawionego wierszy zawierających wiele braków
algae2 <- algae[-manyNAs(algae), ]

## sprawdzamy liczbę wybrakowanych wierszy które pozostały
nrow(algae2[!complete.cases(algae2),])
```

```
## [1] 14
```

```
## usuwamy wszystkie wiersze z brakami
algae3 <- na.omit(algae)

## wyświetl wiersze z brakami
algae3[!complete.cases(algae3),] %>% head()
```

```
## [1] season size speed mxPH mn02 Cl N03 NH4 oP04 P04
## [11] Chla a1 a2 a3 a4 a5 a6 a7
## <0 rows> (or 0-length row.names)
```

```
## liczba pozostałych wybrakowanych wierszy
nrow(algae3[!complete.cases(algae3),])
```

```
## [1] 0
```

```
## można oczywiście też ręcznie usuwać wiersze (nie polecam)
algae4 <- algae[-c(62,199),]
```

Można też zbudować funkcję, która będzie usuwała braki danych wg naszego upodobania.

```
## najpierw budujemy funkcję i ją kompilujemy aby R mógł ja stosować
## parametr prog ustala próg odcięcia wierszy
czysc.dane <- function(dt, prog = 0){
  licz.braki <- apply(dt, 1, function(x) sum(is.na(x)))
  czyste.dt <- dt[!(licz.braki/ncol(dt)>prog), ]
  return(czyste.dt)
}

## potem ją możemy stosować
algae4 <- czysc.dane(algae)
nrow(algae4[!complete.cases(algae4),])

## [1] 0

## czyścimy wiersze, których liczba braków przekracza 20% wszystkich kolumn
algae5 <- czysc.dane(algae, prog = 0.2)
nrow(algae5[!complete.cases(algae5),])

## [1] 14
```

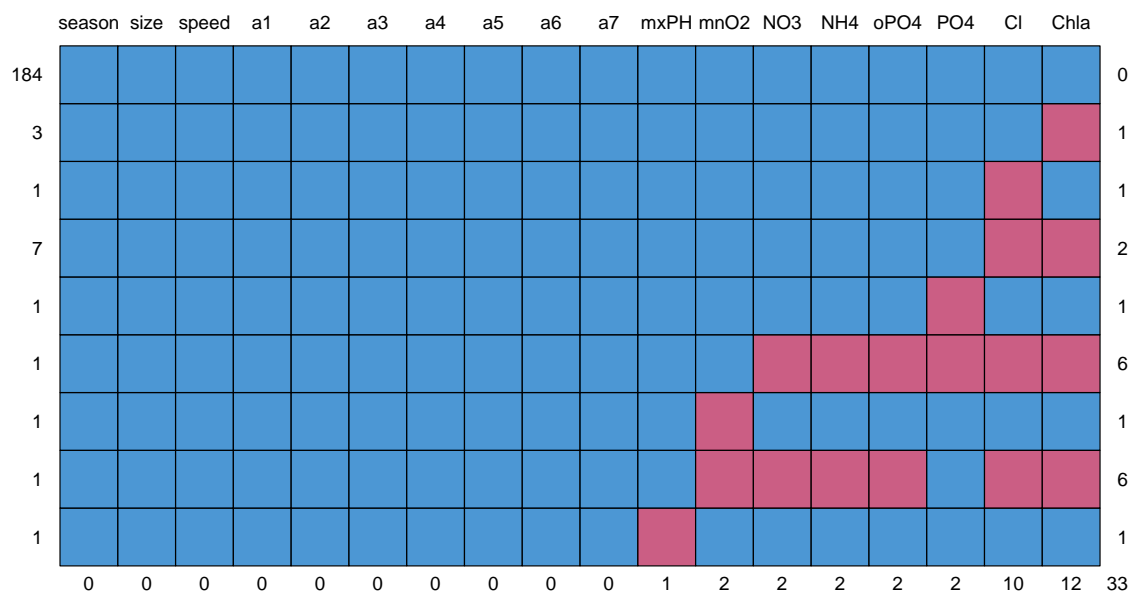
Bardzo ciekawym narzędziem do znajdowania braków danych jest funkcja `md.pattern` pakietu **mice** (van Buuren and Groothuis-Oudshoorn 2018). Wskazuje on ile braków występuje w ramach każdej zmiennej.

```
library(mice)
md.pattern(algae)
```

```
##      season size speed a1 a2 a3 a4 a5 a6 a7 mxPH mn02 N03 NH4 oP04 P04 C1
## 184      1    1    1  1  1  1  1  1  1  1    1    1    1    1    1    1
## 3        1    1    1  1  1  1  1  1  1  1    1    1    1    1    1    1
## 1        1    1    1  1  1  1  1  1  1  1    1    1    1    1    1    0
## 7        1    1    1  1  1  1  1  1  1  1    1    1    1    1    1    0
## 1        1    1    1  1  1  1  1  1  1  1    1    1    1    1    0    1
## 1        1    1    1  1  1  1  1  1  1  1    1    1    0    0    0    0
## 1        1    1    1  1  1  1  1  1  1  1    1    0    1    1    1    1
## 1        1    1    1  1  1  1  1  1  1  1    1    0    0    0    0    0
## 1        1    1    1  1  1  1  1  1  1  1    0    1    1    1    1    1
## 1        1    1    1  1  1  1  1  1  1  1    0    1    1    1    1    1
##      0    0    0  0  0  0  0  0  0  0    1    2    2    2    2    10
##      Chla
## 184      1    0
## 3        0    1
## 1        1    1
## 7        0    2
## 1        1    1
## 1        0    6
## 1        1    1
## 1        0    6
## 1        1    1
##      12 33
```

## 2.2 Zastępowanie braków danych

Zastępowanie braków danych (zwane także *imputacją danych*) jest kolejnym etapem procesu przygotowania danych do analiz. Nie można jednak wyróżnić uniwersalnego sposobu zastępowania braków dla wszystkich możliwych sytuacji. Wśród statystyków panuje przekonanie, że w przypadku wystąpienia braków danych



Rysunek 2.1: Na czerwono zaznaczone są zmienne, które zawierają braki danych. Liczba w wierszu po lewej stronie wykresu wskazuje ile wierszy w bazie ma daną charakterystykę, a liczba po prawej oznacza ile zmiennych było *wybrakowanych*

można zastosować trzy strategie:

- nie robić z brakami - co wydaje się niedorzeczne ale wcale takie nie jest, ponieważ istnieje wiele modeli statystycznych (np. drzewa decyzyjne), które świetnie radzą sobie w sytuacji braków danych. Niestety nie jest to sposób, który można stosować zawsze, ponieważ są również modele wymagające kompletności danych jak na przykład sieci neuronowe.
- usuwać braki wierszami<sup>1</sup> - to metoda, która jest stosowana domyślnie w przypadku kiedy twórca modelu nie zdecydował o innym sposobie obsługi luk. Metoda ta ma swoją niewątpliwą zaletę w postaci jasnej i prostej procedury, ale szczególnie w przypadku niewielkich zbiorów może skutkować obciążeniem estymatorów. Nie wiemy bowiem jaka wartość faktycznie jest przypisana danej cesze. Jeśli jest to wartość bliska np. średniej, to nie wpłynie znacząco na obciążenie estymatora wartości oczekiwanej. W przypadku, gdy różni się ona znacznie od średniej tej cechy, to estymator może już wykazywać obciążenie. Jego wielkość zależy również od liczby usuniętych elementów. Nie jest zalecane usuwanie wielu wierszy ze zbioru danych i na podstawie okrojonego zbioru wyciąganie wniosków o populacji, ponieważ próba jest wówczas znacząco inna niż populacja. Dodatkowo jeśli estymatory są wyznaczane na podstawie zbioru wyraźnie mniej liczego, to precyzja estymatorów wyrażona wariancją spada. Reasumując, jeśli liczba wierszy z brakującymi danymi jest niewielka w stosunku do całego zbioru, to usuwanie wierszy jest sensownym rozwiązaniem.
- uzupełnianie braków - to procedura polegająca na zastępowaniu braków różnymi technikami. Jej niewątpliwą zaletą jest fakt posiadania kompletnych danych bez konieczności usuwania wierszy. Niestety wiąże się to również z pewnymi wadami. Zbiór posiadający wiele braków uzupełnianych nawet bardzo wyrafinowanymi metodami może cechować się zaniżoną wariancją poszczególnych cech oraz tzw. przeuczeniem<sup>2</sup>.

1. Uzupełnianie średnią - braki w zakresie danej zmiennej uzupełniamy średnią tej zmiennej przypadków uzupełnionych.

```
algae[is.na(algae$mxPH), ]
```

```
##   season  size speed mxPH mn02 C1  N03 NH4 oP04 P04 Chla  a1 a2 a3 a4 a5
## 48 winter small  low  NA 12.6 9 0.23 10    5  6  1.1 35.5 0 0 0 0
##    a6 a7
## 48  0  0
```

```
m <- mean(algae$mxPH, na.rm = T)
algae[is.na(algae$mxPH), "mxPH"] <- m
algae[is.na(algae$mxPH), ]
```

```
## [1] season size  speed mxPH  mn02  C1    N03   NH4   oP04   P04
## [11] Chla   a1    a2    a3    a4    a5    a6    a7
## <0 rows> (or 0-length row.names)
```

2. Uzupełnianie medianą - braki w zakresie danej zmiennej uzupełniamy medianą tej zmiennej przypadków uzupełnionych.

```
algae %>% filter(is.na(Chla)) %>% head
```

```
##   season  size speed mxPH mn02 C1  N03 NH4 oP04 P04 Chla  a1 a2 a3
## 1 winter small  high  6.6 10.8 NA 3.245 10    1  6.5  NA 24.3 0.0 0.0
## 2 spring small medium 5.6 11.8 NA 2.220  5    1  1.0  NA 82.7 0.0 0.0
## 3 autumn small medium 5.7 10.8 NA 2.550 10    1  4.0  NA 16.8 4.6 3.9
## 4 spring small  high  6.6  9.5 NA 1.320 20    1  6.0  NA 46.8 0.0 0.0
## 5 summer small  high  6.6 10.8 NA 2.640 10    2 11.0  NA 46.9 0.0 0.0
## 6 autumn small medium 6.6 11.3 NA 4.170 10    1  6.0  NA 47.1 0.0 0.0
##    a4 a5 a6 a7
```

<sup>1</sup> polega na usuwaniu wierszy zawierających braki

<sup>2</sup> więcej o zjawisku przeuczenia w dalszej części książki

```
## 1  0.0  0 0.0  0
## 2  0.0  0 0.0  0
## 3 11.5  0 0.0  0
## 4 28.8  0 0.0  0
## 5 13.4  0 0.0  0
## 6  0.0  0 1.2  0
```

```
algae[is.na(algae$Chla), "Chla"] <- median(algae$Chla, na.rm = T)
```

3. Wypełnianie zmiennych typu wyliczeniowego, logicznego lub znakowego odbywa się najczęściej przez dobranie w miejsce brakującej wartości, elementu powtarzającego się najczęściej wśród obiektów obserwowanych. W pakiecie **DMwR** istnieje funkcja `centralImputation`, która wypełnia braki wartością centralną (w przypadku zmiennych typu liczbowego - medianą, a dla wartości logicznych, wyliczeniowych lub tekstowych - modą).

```
algae[48, "season"]
```

```
## [1] "winter"
```

```
algae[48, "season"] <- NA
algae.uzup <- centralImputation(algae)
algae.uzup[48,]
```

```
##   season size speed   mxPH mn02 Cl  N03 NH4 oP04 P04 Chla  a1 a2 a3
## 48 winter small   low 8.011734 12.6 9 0.23 10    5  6  1.1 35.5 0  0
##   a4 a5 a6 a7
## 48  0  0  0  0
```

4. Jeszcze innym sposobem imputacji danych są algorytmy oparte o metodę  $k$ -najbliższych sąsiadów. Algorytm opiera się na prostej zasadzie, uzupełniania brakujących wartości medianą (w przypadku zmiennych ilościowych) lub modą (w przypadku zmiennych jakościowych) elementów, które są  $k$ -tymi najbliższymi sąsiadami w metryce

$$d(x, y) = \sqrt{\sum_{i=1}^p \delta_i(x_i, y_i)}, \quad (2.1)$$

gdzie  $\delta_i$  jest odległością pomiędzy dwoma elementami ze względu na  $i$ -tą cechę, określoną następująco

$$\delta_i(v_1, v_2) = \begin{cases} 1, & \text{jeśli zmienna jest jakościowa i } v_1 \neq v_2 \\ 0, & \text{jeśli zmienna jest jakościowa i } v_1 = v_2 \\ (v_1 - v_2)^2, & \text{jeśli zmienna jest ilościowa.} \end{cases} \quad (2.2)$$

Odległości są mierzone dla zmiennych standaryzowanych. Istnieje też odmiana z wagami, które maleją wraz ze wzrostem odległości pomiędzy sąsiadem a uzupełnianym elementem (np.  $w(d) = \exp(-d)$ ).

```
algae[48, ]
```

```
##   season size speed   mxPH mn02 Cl  N03 NH4 oP04 P04 Chla  a1 a2 a3
## 48  <NA> small   low 8.011734 12.6 9 0.23 10    5  6  1.1 35.5 0  0
##   a4 a5 a6 a7
## 48  0  0  0  0
```

```
algae <- algae %>%
  mutate_if(is.character, as.factor)
algae.uzup <- knnImputation(algae, k = 5, scale = F, meth = "median")
algae.uzup[48,]
```



```
##      season  size speed      mxPH mnO2 Cl  N03 NH4 oP04 P04 Ch1a  a1 a2 a3
## 48 summer small   low 8.011734 12.6 9 0.23 10    5  6  1.1 35.5 0 0
##      a4 a5 a6 a7
## 48  0  0  0  0
```

Istnieją również dużo bardziej złożone algorytmy imputacji danych oparte na bardziej wyrafinowanych technikach, takich jak: predykcja modelami liniowymi, nieliniowymi, analiza dyskryminacyjna, drzewa klasyfikacyjne. Dwa najbardziej znane pakiety zawierające funkcje do imputacji w sposób złożony, to **Amelia** i **mice**.

Imputacja danych z zastosowaniem pakietu **mice** wymaga podjęcia kilku decyzji przed przystąpieniem do uzupełniania danych:

1. Czy dane są MAR (ang. *Missing At Random*) czy MNAR (ang. *Missing Not At Random*), co oznacza, że musimy się zastanowić jakie mogły być źródła braków danych, przypadkowe czy systematyczne?
2. Należy się zdecydować na formę imputacji, określając strukturę zależności pomiędzy cechami oraz rozkład błędu danej cechy?
3. Wybrać zbiór danych, który posłuży nam za predyktory w imputacji (nie mogą zawierać braków).
4. Określenie, które niepełne zmienne są funkcjami innych wybrakowanych zmiennych.
5. Określić w jakiej kolejności dane będą imputowane.
6. Określić parametry startowe imputacji (liczbę iteracji, warunek zbieżności).
7. Określić liczbę imputowanych zbiorów.

Ad 1. Wyróżniamy następujące rodzaje braków danych:

- MCAR (ang. *Missing Completely At Random*) - z definicji to braki, których pojawienie się jest kompletnie losowe. Przykładowo gdy osoba poproszona o wypełnienie wieku w ankiecie będzie rzucać monetą czy wypełnić tą zmienną.
- MAR - oznacza, że obserwowane wartości i wybrakowane mają inne rozkłady ale da się je oszacować na podstawie danych obserwowanych. Przykładowo ciśnienie tętnicze u osób, które nie wypełniły tej wartości jest wyższe niż u osób, które wpisały swoje ciśnienie. Okazuje się, że osoby starsze z nadciśnieniem nie wypełniały ankiety w tym punkcie.
- MNAR - jeśli nie jest spełniony warunek MCAR i MAR, wówczas brak ma charakter nielosowy. Przykładowo respondenci osiągający wyższe zarobki sukcesywnie nie wypełniają pola "zarobki" i dodatkowo nie ma w ankiecie zmiennych, które pozwoliłyby nam ustalić, jakie to osoby.

Ad 2. Decyzja o algorytmie imputacji wynika bezpośrednio ze skali w jakiej jest mierzona dana zmienna. Ze względu na rodzaj cechy używać będziemy następujących metod:

Każdy z czterech typów danych ma swój domyślny algorytm przeznaczony do imputacji:

- zmienna ilościowa - **pmm**
- zmienna dychotomiczna (stany 0 lub 1) - **logreg**
- zmienna typu wyliczeniowego (nieuporządkowana) - **polyreg**
- zmienna typu wyliczeniowego (uporządkowana) - **polr**

Niewątpliwą zaletą metody **pmm** jest to, że wartości imputowane są ograniczone jedynie do obserwowanych wartości. Metody **norm** i **norm.nob** uzupełniają brakujące wartości w oparciu o model liniowy. Są one szybkie i efektywne w przypadku gdy reszty modelu są zbliżone rozkładem do normalności. Druga z tych technik nie bierze pod uwagę niepewności związanej z modelem imputującym. Metoda **2L.norm** opiera się na dwupoziomym heterogenicznym modelu liniowym (skupienia są włączone jako efekt do modelu). Technika **polyreg** korzysta z funkcji **multinom** pakietu **nnet** tworzącej model wielomianowy. **polr** opiera się o proporcjonalny model logitowy z pakietu **MASS**. **lda** to model dyskryminacyjny klasyfikujący obiekty na podstawie prawdopodobieństw *a posteriori*. Metoda **sample** zastępuje braki losowa wybranymi wartościami spośród wartości obserwowanych.

Ad 3. Do ustalenia predyktorów w modelu **mice** służy funkcja **predictorMatrix**. Po pierwsze wyświetla ona domyślny układ predyktorów włączanych do modelu. Można go dowolnie zmienić i podstawić do modelu imputującego dane parametrem **predictorMatrix**. Zera występujące w kolejnych wierszach macierzy

predyktorów oznaczają pominięcie tej zmiennej przy imputacji innej zmiennej. Jeśli dodatkowo chcemy by jakaś zmienna nie była imputowana, to oprócz usunięcia jej z listy predyktorów, należy wymazać ją z listy metod predykcji (`method`).

Ogólne zalecenia co do tego jakie zmienne stosować jako predyktory jest takie, żeby brać ich jak najwięcej. Spowoduje to, że bardziej prawdopodobny staje się brak typu MAR a nie MNAR. Z drugiej jednak strony, nierzadko zbiory zawierają olbrzymią liczbę zmiennych i włączanie ich wszystkich do modelu imputującego nie będzie miało sensu.

Zalecenia doboru zmiennych są następujące:

- weź wszystkie te zmienne, które są włączane do modelu właściwego, czyli tego za pomocą którego chcesz poznać strukturę zależności;
- czasem do modelu imputującego należy też włączyć interakcje zmiennych z modelu właściwego;
- dodaj zmienne, które mogą mieć wpływ na wybrakowane cechy;
- włącz zmienne istotnie podnoszące poziom wyjaśnionej wariancji modelu;
- na koniec usuń te zmienne spośród predyktorów, które same zawierają zbyt wiele braków.

Ad 4-7. Decyzje podejmowane w tych punktach zależą istotnie od analizowanego zbioru i będą przedmiotem oddzielnych analiz w kontekście rozważanych zbiorów i zadań.

**Przykład 2.1.** Dokonamy imputacji zbioru `airquality` z wykorzystaniem pakietów `mice` i `VIM` (Templ et al. 2019)

```
data <- airquality
summary(data)
```

```
##      Ozone      Solar.R      Wind      Temp
## Min.   : 1.00   Min.   : 7.0   Min.   : 1.700   Min.   :56.00
## 1st Qu.: 18.00   1st Qu.:115.8   1st Qu.: 7.400   1st Qu.:72.00
## Median : 31.50   Median :205.0   Median : 9.700   Median :79.00
## Mean   : 42.13   Mean   :185.9   Mean   : 9.958   Mean   :77.88
## 3rd Qu.: 63.25   3rd Qu.:258.8   3rd Qu.:11.500   3rd Qu.:85.00
## Max.   :168.00   Max.   :334.0   Max.   :20.700   Max.   :97.00
## NA's   :37      NA's   :7
##      Month      Day
## Min.   :5.000   Min.   : 1.0
## 1st Qu.:6.000   1st Qu.: 8.0
## Median :7.000   Median :16.0
## Mean   :6.993   Mean   :15.8
## 3rd Qu.:8.000   3rd Qu.:23.0
## Max.   :9.000   Max.   :31.0
##
```

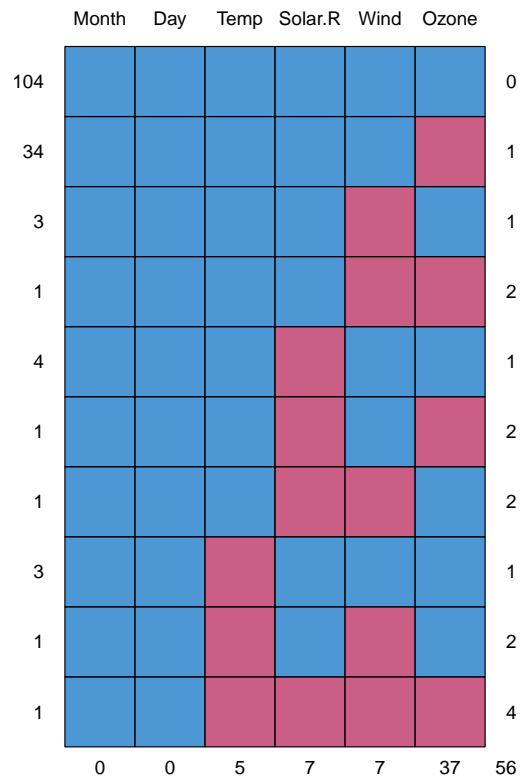
```
# tworzymy dodatkowe braki danych
```

```
data[4:10,3] <- rep(NA,7)
data[1:5,4] <- NA
summary(data)
```

```
##      Ozone      Solar.R      Wind      Temp
## Min.   : 1.00   Min.   : 7.0   Min.   : 1.700   Min.   :57.00
## 1st Qu.: 18.00   1st Qu.:115.8   1st Qu.: 7.400   1st Qu.:73.00
## Median : 31.50   Median :205.0   Median : 9.700   Median :79.00
## Mean   : 42.13   Mean   :185.9   Mean   : 9.806   Mean   :78.28
## 3rd Qu.: 63.25   3rd Qu.:258.8   3rd Qu.:11.500   3rd Qu.:85.00
## Max.   :168.00   Max.   :334.0   Max.   :20.700   Max.   :97.00
## NA's   :37      NA's   :7      NA's   :7      NA's   :5
##      Month      Day
```

```
## Min. :5.000 Min. : 1.0
## 1st Qu.:6.000 1st Qu.: 8.0
## Median :7.000 Median :16.0
## Mean :6.993 Mean :15.8
## 3rd Qu.:8.000 3rd Qu.:23.0
## Max. :9.000 Max. :31.0
##
```

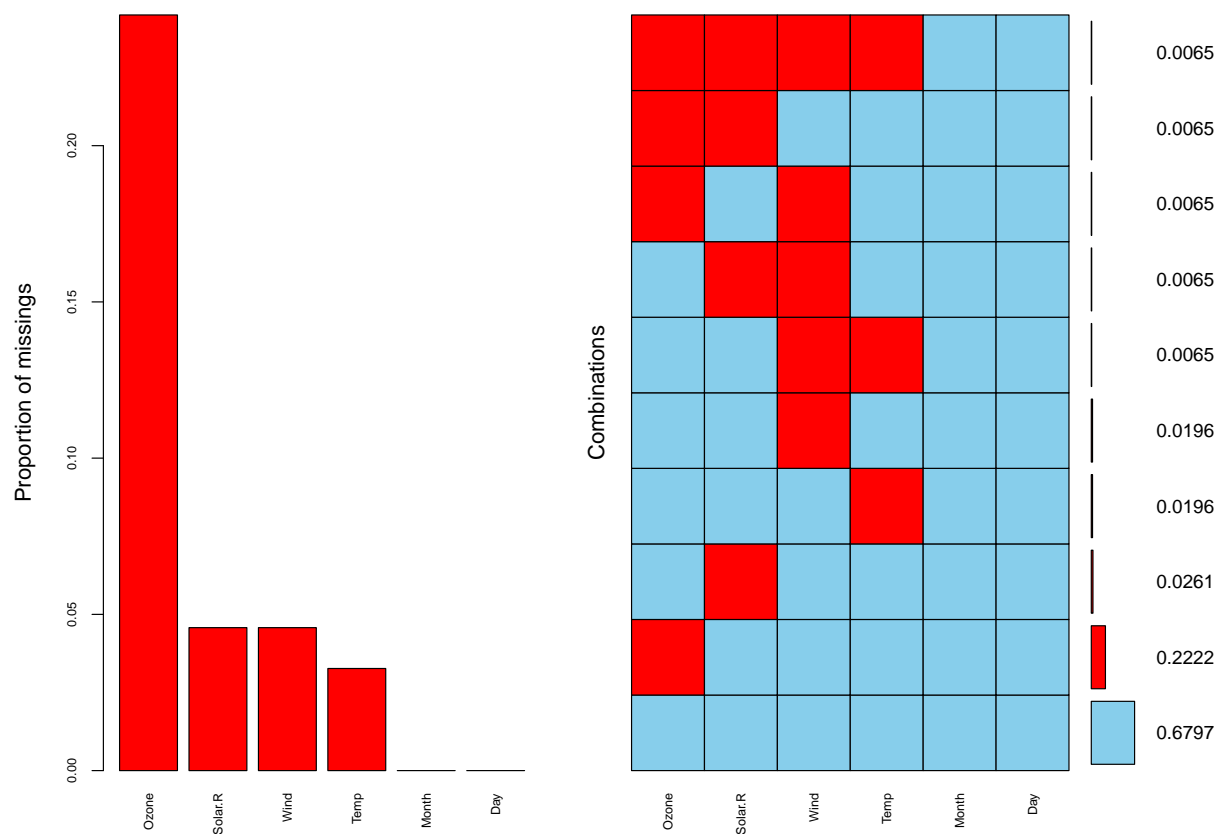
```
md.pattern(data)
```



```
##      Month Day Temp Solar.R Wind Ozone
## 104      1  1  1      1    1    1  0
## 34       1  1  1      1    1    0  1
## 3        1  1  1      1    0    1  1
## 1        1  1  1      1    0    0  2
## 4        1  1  1      0    1    1  1
## 1        1  1  1      0    1    0  2
## 1        1  1  1      0    0    1  2
## 3        1  1  0      1    1    1  1
## 1        1  1  0      1    0    1  2
## 1        1  1  0      0    0    0  4
##         0  0  5      7    7   37 56
```

Do ilustracji braków danych można zastosować funkcje pakietu **VIM**.

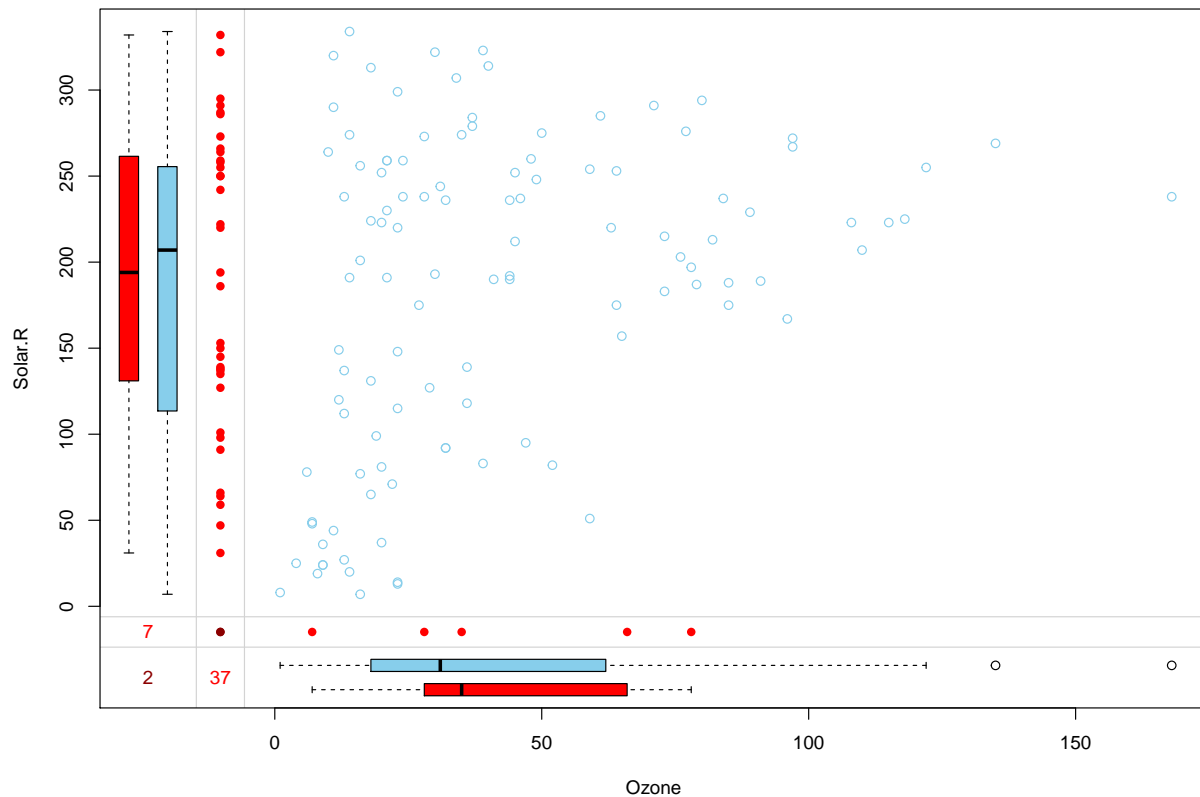
```
library(VIM)
aggr(data, numbers=TRUE,
      sortVars=TRUE,
      labels=names(data),
      cex.axis=.7)
```



```
##
## Variables sorted by number of missings:
## Variable      Count
## Ozone 0.24183007
## Solar.R 0.04575163
## Wind 0.04575163
## Temp 0.03267974
## Month 0.00000000
## Day 0.00000000
```

Tak przedstawia się wykres rozrzutu zmiennych `Ozone` i `Solar.R` z uwzględnieniem położenia braków danych.

```
marginplot(data[c(1,2)])
```



Dokonyamy imputacji metodą pmm.

```
tempData <- mice(data,
  maxit=50,
  meth='pmm',
  seed=44,
  printFlag = F)
summary(tempData)
```

```
## Class: mids
## Number of multiple imputations: 5
## Imputation methods:
##   Ozone Solar.R   Wind   Temp   Month   Day
##   "pmm"  "pmm"   "pmm"  "pmm"    ""    ""
## PredictorMatrix:
##       Ozone Solar.R Wind Temp Month Day
## Ozone      0      1   1   1   1   1
## Solar.R     1      0   1   1   1   1
## Wind        1      1   0   1   1   1
## Temp        1      1   1   0   1   1
## Month       1      1   1   1   0   1
## Day         1      1   1   1   1   0
```

Ponieważ, funkcja `mice` domyślnie dokonuje 5 kompletnych imputacji, możemy się przekonać jak bardzo różnią się poszczególne imputacje i zdecydować się na jedną z nich.

```
head(tempData$imp$Ozone)
```

```
##      1  2  3  4  5
## 5   21 20  7 36 13
## 10  21 16 44 22 21
## 25  14 14 14  6  8
## 26  23 18  8 19 14
## 27  37 23 21  7  9
## 32  63 23  7 52 39
```

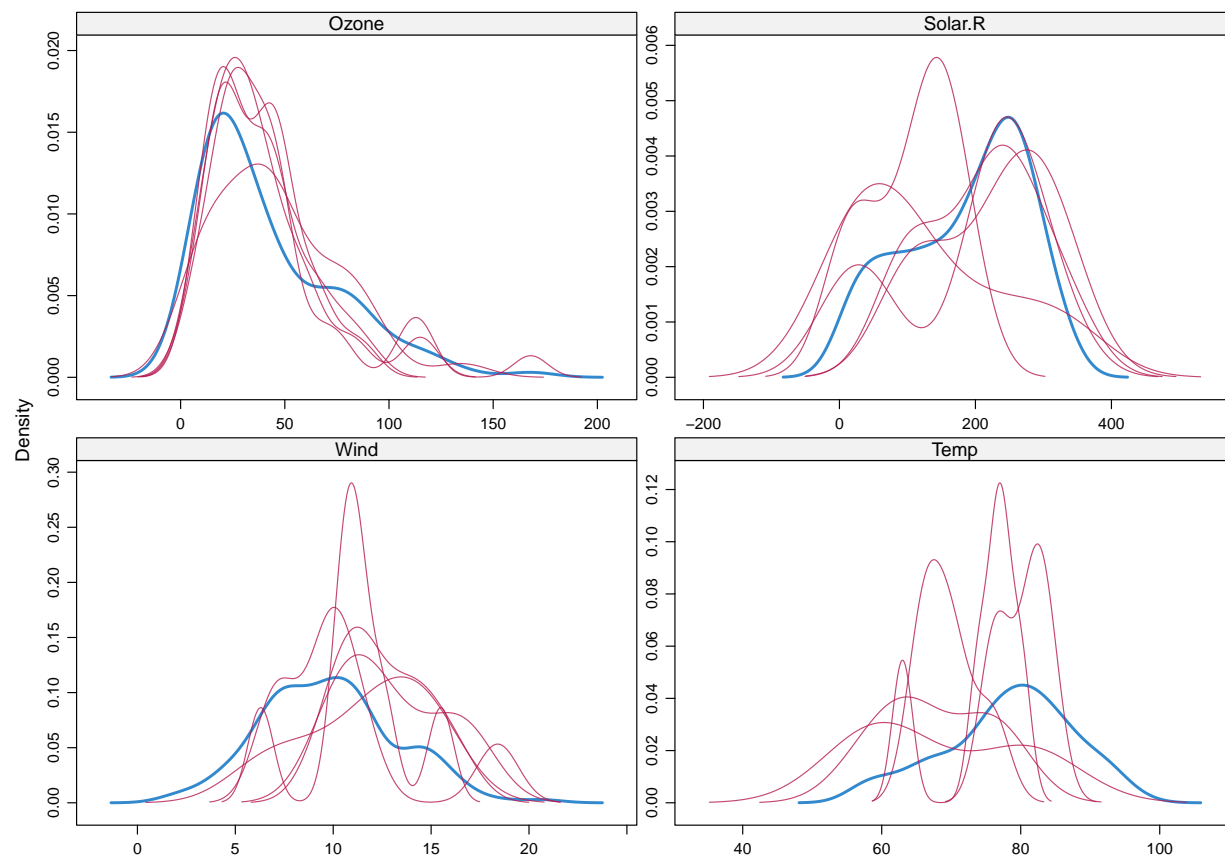
Ostatecznie imputacji dokonujemy wybierając jeden z zestawów danych uzupełniających (np. pierwszy).

```
completedData <- mice::complete(tempData, 1)
summary(completedData)
```

```
##      Ozone      Solar.R      Wind      Temp
## Min.   : 1.0   Min.   : 7.0   Min.   : 1.700   Min.   :57.00
## 1st Qu.: 20.0   1st Qu.:115.0   1st Qu.: 7.400   1st Qu.:73.00
## Median : 32.0   Median :212.0   Median : 9.700   Median :79.00
## Mean   : 42.5   Mean   :187.9   Mean   : 9.931   Mean   :78.14
## 3rd Qu.: 59.0   3rd Qu.:259.0   3rd Qu.:11.500   3rd Qu.:85.00
## Max.   :168.0   Max.   :334.0   Max.   :20.700   Max.   :97.00
##      Month      Day
## Min.   :5.000   Min.   : 1.0
## 1st Qu.:6.000   1st Qu.: 8.0
## Median :7.000   Median :16.0
## Mean   :6.993   Mean   :15.8
## 3rd Qu.:8.000   3rd Qu.:23.0
## Max.   :9.000   Max.   :31.0
```

Za pomocą funkcji pakietu `mice` możemy również przedstawić graficznie gdzie i jak zostały uzupełnione dane.

```
densityplot(tempData, ~Ozone+Solar.R+Wind+Temp)
```



```
stripplot(tempData, Ozone+Solar.R+Wind+Temp~.imp, pch = 20, cex = 1.2)
```

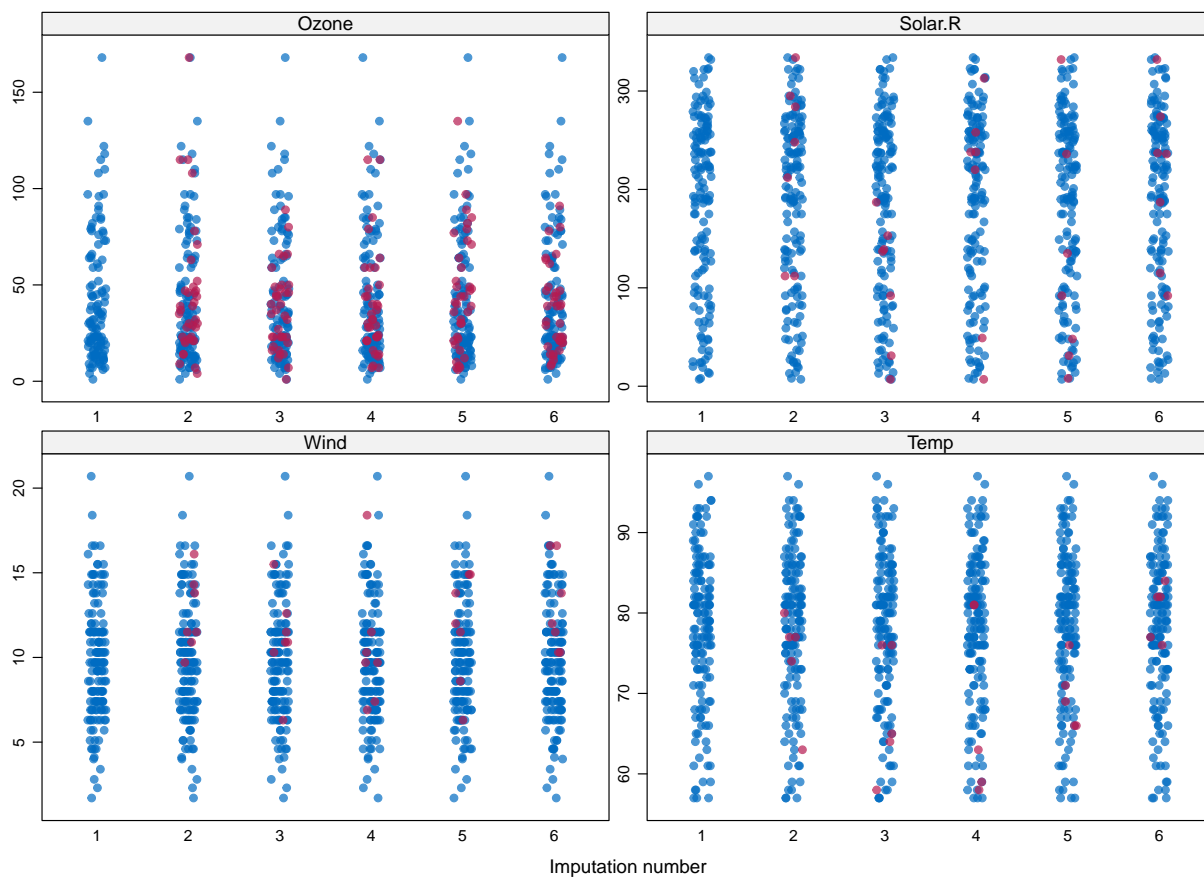




Tabela 2.1: Zestaw metod imputacji danych stosowanych w pakiecie **mice**

method	type	description
pmm	any	Predictive.mean.matching
midastouch	any	Weighted predictive mean matching
sample	any	Random sample from observed values
cart	any	Classification and regression trees
rf	any	Random forest imputations
mean	numeric	Unconditional mean imputation
norm	numeric	Bayesian linear regression
norm.nob	numeric	Linear regression ignoring model error
norm.boot	numeric	Linear regression using bootstrap
norm.predict	numeric	Linear regression, predicted values
quadratic	numeric	Imputation of quadratic terms
ri	numeric	Random indicator for nonignorable data
logreg	binary	Logistic regression
logreg.boot	binary	Logistic regression with bootstrap
polr	ordered	Proportional odds model
polyreg	unordered	Polytomous logistic regression
lda	unordered	Linear discriminant analysis
2l.norm	numeric	Level-1 normal heteroscedastic
2l.lmer	numeric	Level-1 normal homoscedastic, lmer
2l.pan	numeric	Level-1 normal homoscedastic, pan
2l.bin	binary	Level-1 logistic, glmer
2lonly.mean	numeric	Level-2 class mean
2lonly.norm	numeric	Level-2 class normal
2lonly.pmm	any	Level-2 class predictive mean matching



## Rozdział 3

# Podział metod data mining

### 3.1 Rodzaje wnioskowania

*Data mining* to zestaw metod pozyskiwania wiedzy na podstawie danych. Ową wiedzę zdobywamy w procesie wnioskowania na podstawie modeli. Wnioskowanie możemy podzielić na dedukcyjne i indukcyjne. I tak z wnioskowaniem dedukcyjnym mamy do czynienia wówczas, gdy na podstawie obecnego stanu wiedzy potrafimy odpowiedzieć na postawione pytanie dotyczące nowej wiedzy, stosując reguły wnioskowania. O wnioskowaniu indukcyjnym powiemy, że jest to metoda pozyskiwania wiedzy na podstawie informacji ze zbioru uczącego. Znajduje ono szerokie zastosowanie w data mining i charakteryzuje się omylnością, ponieważ nawet najlepiej nauczony model na zbiorze uczącym nie zapewnia nam prawdziwości odpowiedzi w przypadku nowych danych, a jedynie je uprawdopodobnia. Esencją wnioskowania indukcyjnego w zakresie data mining, jest poszukiwanie na podstawie danych uczących modelu charakteryzującego się najlepszymi właściwościami predykcyjnymi i dającego się zastosować do zupełnie nowego zbioru danych.

Każdy proces uczenia z wykorzystaniem wnioskowania indukcyjnego składa się z następujących elementów.

#### 3.1.1 Dziedzina

*Dziedzina* to zbiór wszystkich obiektów pozostających w zainteresowaniu badacza, będących przedmiotem wnioskowania, oznaczana najczęściej przez  $X$ . Przykładowo mogą to być zbiory osób, transakcji, urzędów, instytucji, itp.

#### 3.1.2 Obserwacja

Każdy element dziedziny  $x \in X$  nazywamy obserwacją. Obserwacją nazywać będziemy zarówno rekordy danych ze zbioru uczącego, jak i ze zbioru testowego.

#### 3.1.3 Atrybuty obserwacji

Każdy obiekt z dziedziny  $x \in X$  można opisać zestawem cech (atrybutów), które w notacji matematycznej oznaczymy przez  $a : X \rightarrow A$ , gdzie  $A$  jest przestrzenią wartości atrybutów. Każda obserwacja  $x$  posiadająca  $k$  cech da się wyrazić wektorowo jako  $(a_1(x), a_2(x), \dots, a_k(x))$ . Dla większości algorytmów uczenia maszynowego wyróżnia się trzy typy atrybutów:

- *nominalne* - posiadające skończoną liczbę stanów, które posiadają porządku;
- *porządkowe* - posiadające skończoną liczbę stanów z zachowaniem porządku;
- *ciągłe* - przyjmujące wartości numeryczne.

Często jeden z atrybutów spełnia specjalną rolę, ponieważ stanowi realizację cechy, którą traktujemy jako wyjściową (ang. *target value attribute*). W tym przypadku powiemy o **nadzorowanym uczeniu maszy-**

**nowym.** Jeśli zmiennej wyjściowej nie ma dziedzinie, to mówimy o **nienadzorowanym uczeniu maszynowym**.

### 3.1.4 Zbiór uczący

Zbiorem uczącym  $T$  (ang. *training set*) nazywamy podzbiór  $D$  dziedziny  $X$  (czyli  $T \subseteq D \subseteq X$ ), gdzie zbiór  $D$  stanowi ogół dostępnych obserwacji z dziedziny  $X$ . Zbiór uczący zawiera informacje dotyczące badanego zjawiska, na podstawie których, dokonuje się doboru modelu, selekcji cech istotnych z punktu widzenia własności predykcyjnych lub jakości klasyfikacji, budowy modelu oraz optymalizacji jego parametrów. W przypadku uczenia z nauczycielem (nadzorowanego) zbiór  $T$  zawiera informacje o wartościach atrybutów zmiennej wynikowej.

### 3.1.5 Zbiór testowy

Zbiór testowy  $T'$  (ang. *test set*) będący dopełnieniem zbioru uczącego do zbioru  $D$ , czyli  $T' = D \setminus T$ , stanowi zestaw danych służący do oceny poprawności modelu nadzorowanego. W przypadku metod nienadzorowanych raczej nie stosuje się zbiorów testowych.

### 3.1.6 Model

Model to narzędzie pozyskiwania wiedzy na podstawie zbioru uczącego. Nauczony model jest zbiorem reguł  $f$ , którego zadaniem jest oszacowanie wielkości wartości wynikowej lub odpowiednia klasyfikacja obiektów. W zadaniu grupowania obiektów (ang. *clustering task*), celem modelu jest podanie grup możliwie najbardziej jednorodnych przy zadanym zestawie zmiennych oraz ustalonej liczbie skupień (czasami wyznaczenie liczby skupień jest również częścią zadania stawianego przed modelem).

### 3.1.7 Jakość dopasowania modelu

Do oceny jakości dopasowania modelu wykorzystuje się, w zależności od zadania, wiele współczynników (np. dla zadań regresyjnych są to błąd średnio-kwadratowy - ang. *Mean Square Error*, a dla zadań klasyfikacyjnych - trafność - ang. *Accuracy*). Możemy mówić dwóch rodzajach dopasowania modeli:

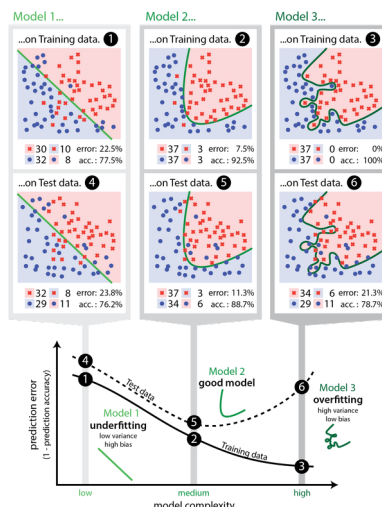
- poziom dopasowania na zbiorze uczącym
- poziom dopasowania na zbiorze testowym (oczywiście z punktu widzenia użyteczności modelu ten współczynnik jest ważniejszy).

W sytuacji, w której model wykazuje dobre charakterystyki jakości dopasowania na zbiorze uczącym ale słabe na testowym, mówimy o zjawisku przeuczenia modelu (ang. *overfitting*). Oznacza to, że model wskazuje predykcję poprawnie jedynie dla zbioru treningowego ale ma słabe własności generalizacyjne nowe przypadki danych. Takie model nie przedstawiają znaczącej wartości w odkrywaniu wiedzy w sposób indukcyjny.

Z drugiej strony parametry dopasowania modelu mogą pokazywać słabe dopasowanie, zarówno na zbiorze uczącym, jak i testowym. Wówczas również model nie jest użyteczny w pozyskiwaniu wiedzy na temat badanego zjawiska, a sytuację taką nazywamy niedouczeniem (ang. *underfitting*).

## 3.2 Modele regresyjne

Jednym z rodzajów zadań bazującym na wnioskowaniu indukcyjnym jest model regresyjny. Należy on do grupy metod nadzorowanych, których celem jest oszacowanie wartości cechy wyjściowej (która jest ilościowa) na podstawie zestawu predyktorów, które mogą być ilościowe i jakościowe. Uczenie takich modeli odbywa się poprzez optymalizację funkcji celu (np. *MSE*) na podstawie zbioru uczącego.



Rysunek 3.1: Przykłady niedouczenia (wykresy 1 i 4), poprawnego modelu (2 i 5) i przeuczenia (3 i 6). Pierwszy wiersz wykresów pokazuje klasyfikację na podstawie modelu na zbiorze uczącym, a drugi na zbiorze testowym. Wykres na dole pokazuje związek pomiędzy złożonością modelu a wielkością błędu predykcji. Źródło: <https://cambridgecoding.wordpress.com/2016/03/24/misleading-modelling-overfitting-cross-validation-and-the-bias-variance-trade-off/>

### 3.3 Modele klasyfikacyjne

Podobnie jak modele regresyjne, modele klasyfikacyjne należą do grupy metod nadzorowanego uczenia maszynowego. Ich zadaniem jest właściwa klasyfikacja obiektów na podstawie wielkości predyktorów. Odpowiedzią modelu jest zawsze cecha typu jakościowego, natomiast predyktory mogą mieć dowolny typ. Wyróżnia się klasyfikację dwu i wielostanową. Lista modeli realizujących klasyfikację binarną jest nieco dłuższa niż w przypadku modeli z wielostanową cechą wynikową. Proces uczenia modelu klasyfikacyjnego również opiera się na optymalizacji funkcji celu. Tym razem są to zupełnie inne miary jakości dopasowania (np. trafność, czyli odsetek poprawnych klasyfikacji).

### 3.4 Modele grupujące

Bardzo szeroką gamę modeli nienadzorowanych stanowią metody analizy skupień. Ich zadaniem jest grupowanie obiektów w możliwie najbardziej jednorodne grupy, na podstawie wartości atrybutów poddanych analizie. Ponieważ są to metody “bez nauczyciela”, to ocena ich przydatności ma nieco inny charakter i choć istnieją różne wskaźniki jakości grupowania, to trudno tu o obiektywne wskazanie najlepszego rozwiązania.



# Rozdział 4

## Drzewa decyzyjne

*Drzewo decyzyjne*<sup>1</sup> jest strukturą hierarchiczną przedstawiającą model klasyfikacyjny lub regresyjny. Stosowane są szczególnie często wówczas, gdy funkcyjna postać związku pomiędzy predyktorami a zmienną wynikową jest nieznana lub ciężka do ustalenia. Każde drzewo decyzyjne składa się z korzenia (ang. *root*), węzłów (ang. *nodes*) i liści (ang. *leaves*). Korzeniem nazywamy początkowy węzeł drzewa, z którego poprzez podziały (ang. *splits*) powstają kolejne węzły potomne. Końcowe węzły, które nie podlegają podziałom nazywamy liśćmi, a linie łączące węzły nazywamy gałęziami (ang. *branches*).

Jeśli drzewo służy do zadań klasyfikacyjnych, to liście zawierają informację o tym, która klasa w danym ciągu podziałów jest najbardziej prawdopodobna. Natomiast, jeśli drzewo jest regresyjne, to liście zawierają warunkowe miary tendencji centralnej (najczęściej średnią) wartości zmiennej wynikowej. Warunek stanowi szereg podziałów doprowadzający do danego węzła terminalnego (liścia). W obu przypadkach (klasyfikacji i regresji) drzewo “dąży” do takiego podziału by kolejne węzły, a co za tym idzie również liście, były ja najbardziej jednorodne ze względu na zmienną wynikową.

### 4.1 Węzły i gałęzie

Każdy podział rozdziela dziedzinę  $X$  na dwa lub więcej podobszarów dziedziny i wówczas każda obserwacja węzła nadrzędnego jest przyporządkowana węzłom potomnym. Każdy odchodzący węzeł potomny jest połączony gałęzią, która to wiąże się ściśle z możliwymi wynikami podziału. Każdy  $n$ -ty węzeł można opisać jako podzbiór dziedziny w następujący sposób

$$X_n = \{x \in X | t_1(x) = r_1, t_2(x) = r_2, \dots, t_k(x) = r_k\}, \quad (4.1)$$

gdzie  $t_1, t_2, \dots, t_k$  są podziałami, które przeprowadzają  $x$  w obszary  $r_1, r_2, \dots, r_k$ . Przez

$$S_{n,t=r} = \{x \in S | t(x) = r\} \quad (4.2)$$

rozumiemy, że dokonano takiego ciągu podziałów zbioru  $S$ , że jego wartości znalazły się w  $n$ -tym węźle.

### 4.2 Rodzaje reguł podziału

Najczęściej występujące reguły podziału w drzewach decyzyjnych są jednowymiarowe, czyli warunek podziału jest generowany na podstawie jednego atrybutu. Istnieją podziały wielowymiarowe ale ze względu na złożoność obliczeniową są rzadziej stosowane.

---

<sup>1</sup>wyglądem przypomina odwrócone drzewo, stąd nazwa





### 4.2.1 Podziały dla atrybutów ze skali nominalnej

Istnieją dwa typy reguł podziału dla skali nominalnej:

- oparte na wartości atrybutu (ang. *value based*) - wówczas funkcja testowa przyjmuje postać  $t(x) = a(x)$ , czyli podział generują wartości atrybutu;
- oparte na równości (ang. *equality based*) - gdzie funkcja testowa jest zdefiniowana jako

$$t(x) = \begin{cases} 1, & \text{gdzie } a(x) = \nu \\ 0, & \text{w przeciwnym przypadku,} \end{cases} \quad (4.3)$$

gdzie  $\nu \in A$  i  $A$  jest zbiorem możliwych wartości  $a$ . W tym przypadku podział jest dychotomiczny, albo obiekt ma wartość atrybutu równą  $\nu$ , albo go nie ma.

### 4.2.2 Podziały dla atrybutów ze skali ciągłej

Reguły podziału stosowane do skali ciągłej, to:

- oparta na nierównościach (ang. *inequality based*) - zdefiniowana jako

$$t(x) = \begin{cases} 1, & \text{gdzie } a(x) \leq \nu \\ 0, & \text{w przeciwnym przypadku,} \end{cases} \quad (4.4)$$

gdzie  $\nu \in A$ ;

- przedziałowa (ang. *interval based*) - zdefiniowana jako

$$t(x) = \begin{cases} 1, & \text{gdzie } a(x) \in I_1 \\ 2, & \text{gdzie } a(x) \in I_2 \\ \vdots & \\ k, & \text{gdzie } a(x) \in I_k \end{cases} \quad (4.5)$$

gdzie  $I_1, I_2, \dots, I_k \subset A$  stanowią rozłączny podział (przedziałami) przeciwdziedziny  $A$ .

### 4.2.3 Podziały dla atrybutów ze skali porządkowej

Podziały te mogą wykorzystywać oba wcześniej wspomniane typy, w zależności od potrzeb.

## 4.3 Algorytm budowy drzewa

1. stwórz początkowy węzeł (korzeń) i oznacz go jako *otwarty*;
2. przypisz wszystkie możliwe rekordy do węzła początkowego;
3. **dopóki** istnieją otwarte węzły **wykonuj**:
  - wybierz węzeł  $\mathbf{n}$ , wyznacz potrzebne statystyki opisowe zmiennej zależnej dla tego węzła i przypisz wartość docelową;
  - **jeśli** kryterium zatrzymania podziału jest spełnione dla węzła  $\mathbf{n}$ , **to** oznacz go za **zamknięty**;
  - **w przeciwnym przypadku** wybierz podział  $r$  elementów węzła  $\mathbf{n}$ , i dla każdego podzbioru podziału stwórz węzeł niższego rzędu (potomka)  $\mathbf{n}_r$  oraz oznacz go jako *otwarty*;
  - następnie przypisz wszystkie przypadki generowane podziałem  $r$  do odpowiednich węzłów potomków  $\mathbf{n}_r$ ;
  - oznacz węzeł  $\mathbf{n}$  jako *zamknięty*.

Sposób przypisywania wartości docelowej wiąże się ściśle z rodzajem drzewa. W drzewach regresyjnych chodzi o wyliczenie średniej lub mediany dla obserwacji ujętych w danym węźle. Natomiast w przypadku

drzewa klasyfikacyjnego, wyznacza się wartości prawdopodobieństw przynależności obserwacji znajdującej się w danym węźle do poszczególnych klas

$$\mathbb{P}(d|\mathbf{n}) = \mathbb{P}_{T_{\mathbf{n}}}(d) = \frac{|T_{\mathbf{n}}^d|}{|T_{\mathbf{n}}|}, \quad (4.6)$$

gdzie  $T_{\mathbf{n}}$  oznaczają obserwacje zbioru uczącego znajdujące się w węźle  $\mathbf{n}$ , a  $T_{\mathbf{n}}^d$  oznacza dodatkowo podzbiór zbioru uczącego w  $\mathbf{n}$  węźle, które należą do klasy  $d$ . Oczywiście klasyfikacja na podstawie otrzymanych prawdopodobieństw w danym węźle jest dokonana przez wybór klasy charakteryzującej się najwyższym prawdopodobieństwem.

## 4.4 Kryteria zatrzymania

Kryterium zatrzymania jest warunkiem, który decyduje o tym, że dany węzeł uznajemy za zamknięty i nie dokonujemy dalszego jego podziału. Wyróżniamy następujące kryteria zatrzymania:

- jednorodność węzła - w przypadku drzewa klasyfikacyjnego może zdarzyć się sytuacja, że wszystkie obserwacje węzła będą pochodziły z jednej klasy. Wówczas nie ma sensu dokonywać dalszego podziału węzła;
- węzeł jest pusty - zbiór przypisanych obserwacji zbioru uczącego do  $\mathbf{n}$ -tego węzła jest pusty;
- brak reguł podziału - wszystkie reguły podziału zostały wykorzystane, zatem nie da się stworzyć potomnych węzłów, które charakteryzowałyby się większą homogenicznością;

Warunki ujęte w pierwszych dwóch kryteriach mogą być nieco złagodzone, poprzez zatrzymanie podziałów wówczas, gdy prawdopodobieństwo przynależenia do pewnej klasy przekroczy ustalony próg lub gdy liczebność węzła spadnie poniżej ustalonej wartości.

W literaturze tematu istnieje jeszcze jedno często stosowane kryterium zatrzymania oparte na wielkości drzewa. Węzeł potomny ustala się jako zamknięty, gdy długość ścieżki dojścia do niego przekroczy ustaloną wartość.

## 4.5 Reguły podziału

Ważnym elementem algorytmu tworzenia drzewa regresyjnego jest *reguła podziału*. Dobierana jest w taki sposób aby zmaksymalizować zdolności generalizacyjne drzewa. Złożoność drzewa mierzona jest najczęściej przeciętną liczbą podziałów potrzebnych do dotarcia do liścia zaczynając od korzenia. Liście są najczęściej tworzone wówczas gdy dyspersja wartości wynikowej jest stosunkowo mała lub węzeł zawiera w miarę homogeniczne obserwacje ze względu na przynależność do klasy zmiennej wynikowej. W przypadku drzew regresyjnych zmienność na poziomie węzłów jest dobrą miarą służącą do definiowania podziału w węźle. I tak, jeśli pewien podział generuje nam stosunkowo małe dyspersje wartości docelowych w węzłach potomnych, to można ten podział uznać za właściwy. Jeśli  $T_n$  oznacza zbiór rekordów należących do węzła  $n$ , a  $T_{n,t=r}$  są podzbiorami generowanymi przez podział  $r$  w węzłach potomnych dla  $n$ , to dyspersję wartości docelowej  $f$  będziemy oznaczali następująco

$$\text{disp}_{T_{n,t=r}}(f). \quad (4.7)$$

Regułę podziału możemy określać poprzez minimalizację średniej ważonej dyspersji wartości docelowej następującej postaci

$$\text{disp}_n(f|t) = \sum_{r \in R_t} \frac{|T_{n,t=r}|}{|T_n|} \text{disp}_{T_{n,t=r}}(f), \quad (4.8)$$

gdzie  $||$  oznacza moc zbioru, a  $R_t$  zbiór wszystkich możliwych wartości reguły podziału. Czasami wygodniej będzie maksymalizować przyrost dyspersji (lub spadek)

$$\Delta \text{disp}_n(f|t) = \text{disp}_n(f) - \sum_{r \in R_t} \frac{|T_{n,t=r}|}{|T_n|} \text{disp}_{T_{n,t=r}}(f). \quad (4.9)$$

Miarą heterogeniczności węzłów ze względu na zmienną wynikową (ang. *impurity*) w drzewach klasyfikacyjnych, która pozwala na tworzenie kolejnych podziałów węzła, są najczęściej wskaźnik Gini'ego i entropia (Breiman 1998).

Entropią podzbioru uczącego w węźle  $\mathbf{n}$ , wyznaczamy wg wzoru

$$E_{T_{\mathbf{n}}}(c|t) = \sum_{x \in R_t} \frac{|T_{\mathbf{n},t=r}|}{|T_{\mathbf{n}}|} E_{T_{\mathbf{n},t=r}}(c), \quad (4.10)$$

gdzie  $t$  jest podziałem (kandydatem),  $r$  potencjalnym wynikiem podziału  $t$ ,  $c$  jest oznaczeniem klasy zmiennej wynikowej, a

$$E_{T_{\mathbf{n},t=r}}(c) = \sum_{d \in C} -\mathbb{P}_{T_{\mathbf{n},t=r}}(c = d) \log \mathbb{P}_{T_{\mathbf{n},t=r}}(c = d), \quad (4.11)$$

przy czym

$$\mathbb{P}_{T_{\mathbf{n},t=r}}(c = d) = \mathbb{P}_{T_{\mathbf{n}}}(c = d | t = r). \quad (4.12)$$

Podobnie definiuje się indeks Gini'ego

$$Gi_{T_{\mathbf{n}}}(c|t) = \sum_{x \in R_t} \frac{|T_{\mathbf{n},t=r}|}{|T_{\mathbf{n}}|} Gi_{T_{\mathbf{n},t=r}}(c), \quad (4.13)$$

gdzie

$$Gi_{T_{\mathbf{n},t=r}}(c) = \sum_{d \in C} \mathbb{P}_{T_{\mathbf{n},t=r}}(c = d) \cdot (1 - \mathbb{P}_{T_{\mathbf{n},t=r}}(c = d)) = 1 - \sum_{d \in C} \mathbb{P}_{T_{\mathbf{n},t=r}}^2(c = d). \quad (4.14)$$

Dla tak zdefiniowanych miar “nieczystości” węzłów, podziału dokonujemy w taki sposób, aby zminimalizować współczynnik Gini'ego lub entropię. Im niższe miary nieczystości, tym bardziej obserwacje znajdujące się w węźle są monokulturą<sup>2</sup>. Nierzadko korzysta się również z współczynnika przyrostu informacji (ang. *information gain*)

$$\Delta E_{T_{\mathbf{n}}}(c|t) = E_{T_{\mathbf{n}}}(c) - E_{T_{\mathbf{n}}}(c|t). \quad (4.15)$$

Istnieje również jego odpowiednik dla indeksu Gini'ego. W obu przypadkach optymalnego podziału szukamy poprzez maksymalizację przyrostu informacji.

## 4.6 Przycinanie drzewa decyzyjnego

Uczenie drzewa decyzyjnego wiąże się z ryzykiem przeuczenia modelu (podobnie jak to się ma w przypadku innych modeli predykcyjnych). Wcześniej przytoczone reguły zatrzymania (np. głębokość drzewa czy zatrzymanie przy osiągnięciu jednorodności na zadanym poziomie) pomagają kontrolować poziom generalizacji drzewa ale czasami będzie dodatkowo potrzebne przycięcie drzewa, czyli usunięcie pewnych podziałów, a co za tym idzie, również liści (węzłów).

### 4.6.1 Przycinanie redukujące błąd

Jedną ze strategii przycinania drzewa jest przycinanie redukujące błąd (ang. *reduced error pruning*). Polega ono na porównaniu błędów (najczęściej używana jest miara odsetka błędnych klasyfikacji lub MSE) liścia  $\mathbf{l}$  i węzła do którego drzewo przycinamy  $\mathbf{n}$  na całym nowym zbiorze uczącym  $R$ . Niech  $e_R(\mathbf{l})$  i  $e_R(\mathbf{n})$  oznaczają odpowiednio błędy na zbiorze  $R$  liścia i węzła. Przez błąd węzła rozumiemy błąd pod-drzewa o korzeniu  $\mathbf{n}$ . Wówczas jeśli zachodzi warunek

$$e_R(\mathbf{l}) \leq e_R(\mathbf{n}), \quad (4.16)$$

to zaleca się zastąpić węzeł  $\mathbf{n}$  liściem  $\mathbf{l}$ .

---

<sup>2</sup>prawie wszystkie są w jednej klasie

### 4.6.2 Przycinanie minimalizujące błąd

Przycinanie minimalizujące błąd opiera się na spostrzeżeniu, że błąd drzewa przyciętego charakteryzuje się zbyt pesymistyczną oceną i dlatego wymaga korekty. Węzeł drzewa klasyfikacyjnego  $\mathbf{n}$  zastępujemy liściem  $\mathbf{l}$ , jeśli

$$\hat{e}_T(\mathbf{l}) \leq \hat{e}_T(\mathbf{n}), \quad (4.17)$$

gdzie

$$\hat{e}_T(\mathbf{n}) = \sum_{\mathbf{n}' \in N(\mathbf{n})} \frac{|T_{\mathbf{n}'}|}{|T_{\mathbf{n}}|} \hat{e}_T(\mathbf{n}'), \quad (4.18)$$

a  $N(\mathbf{n})$  jest zbiorem wszystkich możliwych węzłów potomnych węzła  $\mathbf{n}$  i

$$\hat{e}_T(\mathbf{l}) = 1 - \frac{|\{x \in T_1 | c(x) = d_1\}| + mp}{|T_1| + m}, \quad (4.19)$$

gdzie  $p$  jest prawdopodobieństwem przynależności do klasy  $d_1$  ustalona na podstawie zewnętrznej wiedzy (gdy jej nie posiadamy przyjmujemy  $p = 1/|C|$ ).

W przypadku drzewa regresyjnego znajdujemy wiele analogii, ponieważ jeśli dla pewnego zbioru rekordów  $T$  spełniony jest warunek

$$\text{mse}_T(\mathbf{l}) \leq \text{mse}_T(\mathbf{n}), \quad (4.20)$$

gdzie  $\mathbf{l}$  i  $\mathbf{n}$  oznaczają odpowiednio liść i węzeł, to wówczas zastępujemy węzeł  $\mathbf{n}$  przez liść  $\mathbf{l}$ .

Estymatory wyznaczone na podstawie niewielkiej próby, mogą być obciążone znaczącym błędem. Wyliczanie błędu średnio-kwadratowego dla podzbioru nowych wartości może się charakteryzować takim obciążeniem. Dlatego stosuje się statystyki opisowe z poprawką, której pochodzenie może mieć trzy źródła: wiedza merytoryczna na temat szukanej wartości, założeń modelu lub na podstawie wyliczeń opartych o cały zbiór wartości.

Skorygowany estymator błędu średnio-kwadratowego ma następującą postać

$$\widehat{\text{mse}}_T(\mathbf{l}) = \frac{\sum_{x \in T} (f(x) - m_{\mathbf{l}, m, m_0}(f))^2 + m S_0^2}{|T_1| + m}, \quad (4.21)$$

gdzie

$$m_{\mathbf{l}, m, m_0}(f) = \frac{\sum_{x \in T_1} f(x) + m m_0}{|T_1| + m}, \quad (4.22)$$

a  $m_0$  i  $S_0^2$  są średnią i wariancją wyznaczonymi na całej próbie uczącej. Błąd średnio-kwadratowy węzła  $\mathbf{n}$  ma postać

$$\widehat{\text{mse}}_T(\mathbf{n}) = \sum_{\mathbf{n}' \in N(\mathbf{n})} \frac{|T_{\mathbf{n}'}|}{|T_{\mathbf{n}}|} \widehat{\text{mse}}_T(\mathbf{n}'). \quad (4.23)$$

Wówczas kryterium podcięcia można zapisać w następujący sposób

$$\widehat{\text{mse}}_T(\mathbf{l}) \leq \widehat{\text{mse}}_T(\mathbf{n}) \quad (4.24)$$

### 4.6.3 Przycinanie ze względu na współczynnik złożoności drzewa

Przycinanie ze względu na współczynnik złożoności drzewa (ang. *cost-complexity pruning*) polega na wprowadzeniu “kary” za zwiększoną złożoność drzewa. Drzewa klasyfikacyjne przycinamy gdy spełniony jest warunek

$$e_T(\mathbf{l}) \leq e_T(\mathbf{n}) + \alpha C(\mathbf{n}), \quad (4.25)$$

gdzie  $C(\mathbf{n})$  oznacza złożoność drzewa mierzona liczbą liści, a  $\alpha$  parametrem wagi kary za złożoność drzewa.

Wspomniane kryterium przycięcia dla drzew regresyjnych bazuje na względnym błędzie średnio-kwadratowym (ang. *relative square error*), czyli

$$\widehat{\text{rse}}_T(\mathbf{n}) = \frac{|T|\widehat{\text{mse}}_T(\mathbf{n})}{(|T| - 1)S_T^2(f)}, \quad (4.26)$$

gdzie  $T$  oznacza podzbiór  $X$ ,  $S_T^2$  wariancję na zbiorze  $T$ . Wówczas kryterium podcięcia wygląda następująco

$$\widehat{\text{rse}}_T(\mathbf{l}) \leq \widehat{\text{rse}}_T(\mathbf{n}) + \alpha C(\mathbf{n}). \quad (4.27)$$

## 4.7 Obsługa braków danych

Drzewa decyzyjne wyjątkowo dobrze radzą sobie z obsługą zbiorów z brakami. Stosowane są głównie dwie strategie:

- udziałów obserwacji (ang. *fractional instances*) - rozważane są wszystkie możliwe podziały dla brakującej obserwacji i przypisywana jest im odpowiednia waga lub prawdopodobieństwo, w oparciu o zaobserwowany rozkład znanych obserwacji. Te same wagi są stosowane do predykcji wartości na podstawie drzewa z brakami danych.
- podziałów zastępczych (ang. *surrogate splits*) - jeśli wynik podziału nie może być ustalony dla obserwacji z brakami, to używany jest podział zastępczy (pierwszy), jeśli i ten nie może zostać ustalony, to stosuje się kolejny. Kolejne podziały zastępcze są generowane tak, aby wynik podziału możliwie najbardziej przypominał podział właściwy.

## 4.8 Zalety i wady

### 4.8.1 Zalety

- łatwe w interpretacji;
- nie wymagają żmudnego przygotowania danych (brak standaryzacji, wprowadzania zmiennych binarnych, dopuszcza występowanie braków danych);
- działa na obu typach zmiennych - jakościowych i ilościowych;
- dopuszcza nieliniowość związku między zmienną wynikową a predyktorami;
- odporny na odstępstwa od założeń;
- pozwala na obsługę dużych zbiorów danych.

### 4.8.2 Wady

- brak jawnej postaci zależności;
- zależność struktury drzewa od użytego algorytmu;
- przegrywa jakością predykcji z innymi metodami nadzorowanego uczenia maszynowego.

**Przykład 4.1.** Przykładem zastosowania drzew decyzyjnych będzie klasyfikacja irysów na podstawie długości i szerokości kielicha i płatk.

Przykładem zastosowania drzew decyzyjnych będzie klasyfikacja irysów na podstawie długości i szerokości kielicha i płatk.

```
library(tidyverse)
library(rpart) # pakiet do tworzenia drzew typu CART
library(rpart.plot) # pakiet do rysowania drzew
```

Każde zadanie uczenia maszynowego zaczynamy od czyszczenia danych i odpowiedniego ich przygotowania ale w tym przypadku skupimy się jedynie na budowie, optymalizacji i ewaluacji modelu.

**Podział zbioru na próbę uczącą i testową**

```

set.seed(44)
dt.train <- iris %>%
  sample_frac(size = 0.7)
dt.test <- setdiff(iris, dt.train)
str(dt.train)

## 'data.frame': 105 obs. of 5 variables:
## $ Sepal.Length: num 6.4 4.4 6.6 5.4 5 5.4 5.6 4.4 5.4 6.1 ...
## $ Sepal.Width : num 2.7 3.2 3 3 3.6 3.4 2.9 2.9 3.9 2.9 ...
## $ Petal.Length: num 5.3 1.3 4.4 4.5 1.4 1.7 3.6 1.4 1.3 4.7 ...
## $ Petal.Width : num 1.9 0.2 1.4 1.5 0.2 0.2 1.3 0.2 0.4 1.4 ...
## $ Species : Factor w/ 3 levels "setosa","versicolor",...: 3 1 2 2 1 1 2 1 1 2 ...

str(dt.test)

## 'data.frame': 45 obs. of 5 variables:
## $ Sepal.Length: num 4.7 4.6 5.4 4.8 5.8 5.1 5.1 5.1 5 5.2 ...
## $ Sepal.Width : num 3.2 3.1 3.9 3.4 4 3.8 3.7 3.3 3 3.5 ...
## $ Petal.Length: num 1.3 1.5 1.7 1.6 1.2 1.5 1.5 1.7 1.6 1.5 ...
## $ Petal.Width : num 0.2 0.2 0.4 0.2 0.2 0.3 0.4 0.5 0.2 0.2 ...
## $ Species : Factor w/ 3 levels "setosa","versicolor",...: 1 1 1 1 1 1 1 1 1 1 ...

```

### Budowa drzewa

Budowy drzewa dokonujemy za pomocą funkcji **rpart** pakietu **rpart** (Therneau and Atkinson 2018) stosując zapis formuły zależności. Drzewo zostanie zbudowane z uwzględnieniem kilku kryteriów zatrzymania:

- minimalna liczebność węzła, który może zostać podzielony to 10 - ze względu na małą liczebność zbioru uczącego;
- minimalna liczebność liścia to 5 - aby nie dopuścić do przeuczenia modelu;
- maksymalna głębokość drzewa to 4 - aby nie dopuścić do przeuczenia modelu.

```

mod.rpart <- rpart(Species~., data = dt.train,
  control = rpart.control(minsplit = 10,
    minbucket = 5,
    maxdepth = 4))
summary(mod.rpart)

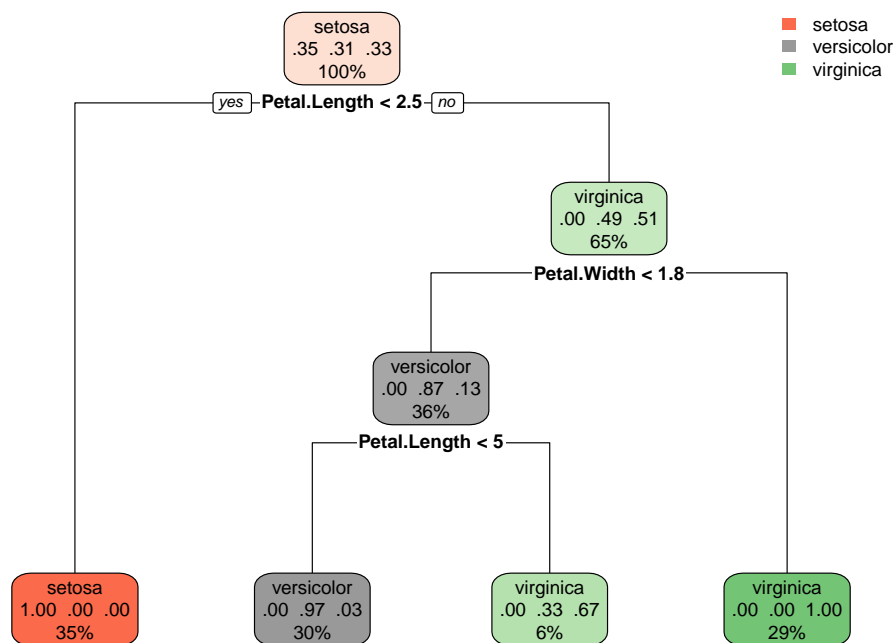
## Call:
## rpart(formula = Species ~ ., data = dt.train, control = rpart.control(minsplit = 10,
##   minbucket = 5, maxdepth = 4))
##   n= 105
##
##           CP nsplit  rel error    xerror    xstd
## 1 0.51470588     0 1.00000000 1.1764706 0.06418173
## 2 0.41176471     1 0.48529412 0.6617647 0.07457243
## 3 0.02941176     2 0.07352941 0.1029412 0.03758880
## 4 0.01000000     3 0.04411765 0.1029412 0.03758880
##
## Variable importance
##   Petal.Width Petal.Length Sepal.Length  Sepal.Width
##           35           33           21           12
##
## Node number 1: 105 observations,    complexity param=0.5147059
##   predicted class=setosa    expected loss=0.647619  P(node) =1
##   class counts:    37    33    35
##   probabilities: 0.352 0.314 0.333

```

```

## left son=2 (37 obs) right son=3 (68 obs)
## Primary splits:
##   Petal.Length < 2.45 to the left, improve=35.95322, (0 missing)
##   Petal.Width < 0.8 to the left, improve=35.95322, (0 missing)
##   Sepal.Length < 5.45 to the left, improve=25.39467, (0 missing)
##   Sepal.Width < 3.35 to the right, improve=12.69596, (0 missing)
## Surrogate splits:
##   Petal.Width < 0.8 to the left, agree=1.000, adj=1.000, (0 split)
##   Sepal.Length < 5.45 to the left, agree=0.924, adj=0.784, (0 split)
##   Sepal.Width < 3.35 to the right, agree=0.819, adj=0.486, (0 split)
##
## Node number 2: 37 observations
## predicted class=setosa expected loss=0 P(node) =0.352381
## class counts: 37 0 0
## probabilities: 1.000 0.000 0.000
##
## Node number 3: 68 observations, complexity param=0.4117647
## predicted class=virginica expected loss=0.4852941 P(node) =0.647619
## class counts: 0 33 35
## probabilities: 0.000 0.485 0.515
## left son=6 (38 obs) right son=7 (30 obs)
## Primary splits:
##   Petal.Width < 1.75 to the left, improve=25.286380, (0 missing)
##   Petal.Length < 4.75 to the left, improve=24.879360, (0 missing)
##   Sepal.Length < 5.75 to the left, improve= 6.713875, (0 missing)
##   Sepal.Width < 3.25 to the left, improve= 1.336180, (0 missing)
## Surrogate splits:
##   Petal.Length < 4.75 to the left, agree=0.882, adj=0.733, (0 split)
##   Sepal.Length < 6.15 to the left, agree=0.721, adj=0.367, (0 split)
##   Sepal.Width < 3.15 to the left, agree=0.618, adj=0.133, (0 split)
##
## Node number 6: 38 observations, complexity param=0.02941176
## predicted class=versicolor expected loss=0.1315789 P(node) =0.3619048
## class counts: 0 33 5
## probabilities: 0.000 0.868 0.132
## left son=12 (32 obs) right son=13 (6 obs)
## Primary splits:
##   Petal.Length < 4.95 to the left, improve=4.0800440, (0 missing)
##   Petal.Width < 1.45 to the left, improve=1.2257490, (0 missing)
##   Sepal.Width < 2.65 to the right, improve=0.6168705, (0 missing)
##   Sepal.Length < 5.95 to the left, improve=0.4736842, (0 missing)
## Surrogate splits:
##   Petal.Width < 1.55 to the left, agree=0.868, adj=0.167, (0 split)
##
## Node number 7: 30 observations
## predicted class=virginica expected loss=0 P(node) =0.2857143
## class counts: 0 0 30
## probabilities: 0.000 0.000 1.000
##
## Node number 12: 32 observations
## predicted class=versicolor expected loss=0.03125 P(node) =0.3047619
## class counts: 0 31 1
## probabilities: 0.000 0.969 0.031
##

```



Rysunek 4.2: Obraz drzewa klasyfikacyjnego.

```
## Node number 13: 6 observations
## predicted class=virginica expected loss=0.3333333 P(node) =0.05714286
## class counts:      0      2      4
## probabilities: 0.000 0.333 0.667
```

```
rpart.plot(mod.rpart)
```

Powyższy wykres przedstawia strukturę drzewa klasyfikacyjnego. Kolorami są oznaczone klasy, które w danym węźle dominują. Nasycenie barwy decyduje o sile tej dominacji. W każdym węźle podana jest klasa, do której najprawdopodobniej należą jego obserwacje. Ponadto podane są proporcje przynależności do klas zmiennej wynikowej oraz procent obserwacji zbioru uczącego należących do danego węzła. Pod każdym węzłem podana jest reguła podziału.

### Przycinanie drzewa

Zanim przystąpimy do przycinania drzewa należy sprawdzić, jakie są zdolności generalizacyjne modelu. Oceny tej dokonujemy najczęściej sprawdzając macierz klasyfikacji.

```
pred.prob <- predict(mod.rpart,
                     newdata = dt.test)
pred.prob[10:20,]
```

```
##      setosa versicolor virginica
## 10         1    0.00000    0.00000
```



```
## 11      1      0.00000  0.00000
## 12      1      0.00000  0.00000
## 13      1      0.00000  0.00000
## 14      0      0.96875  0.03125
## 15      0      0.96875  0.03125
## 16      0      0.96875  0.03125
## 17      0      0.96875  0.03125
## 18      0      0.96875  0.03125
## 19      0      0.96875  0.03125
## 20      0      0.00000  1.00000
```

```
pred.class <- predict(mod.rpart,
                      newdata = dt.test,
                      type = "class")
pred.class
```

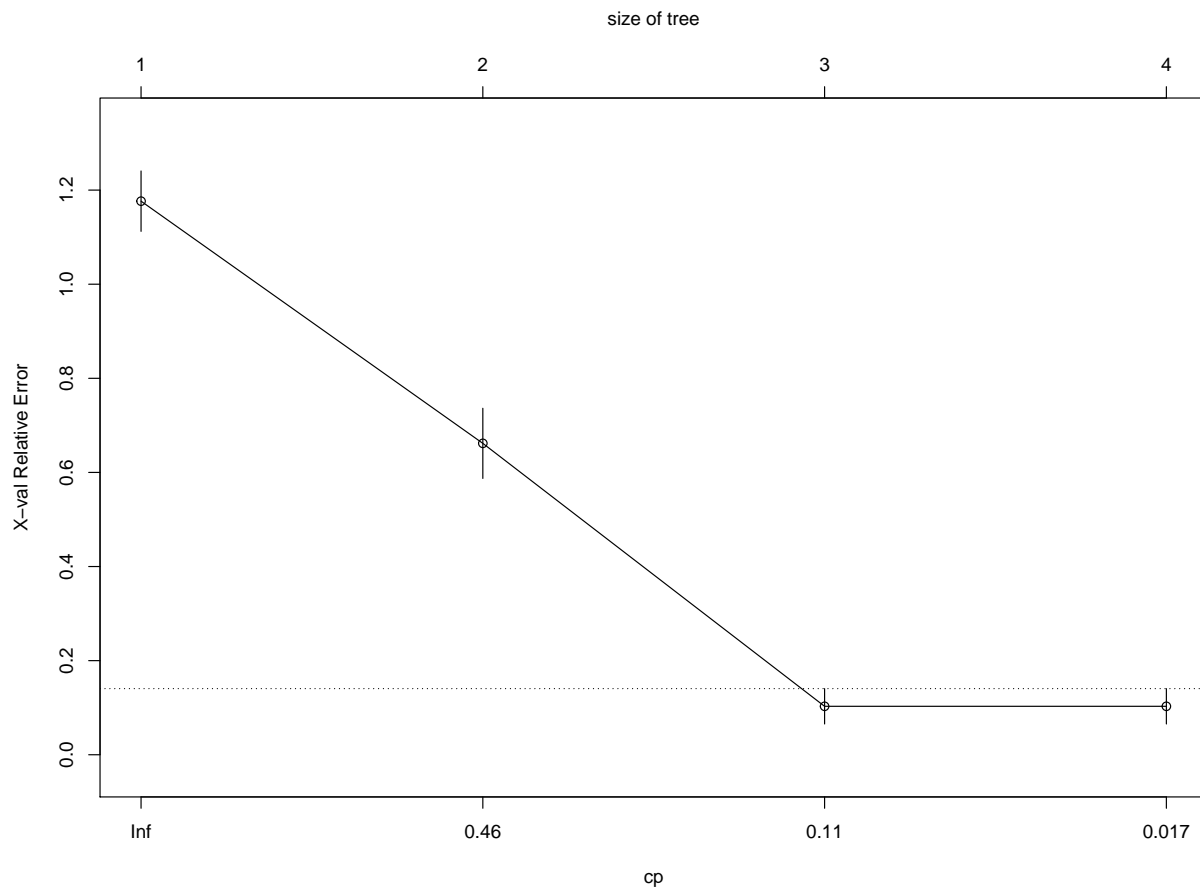
```
##          1          2          3          4          5          6
##      setosa      setosa      setosa      setosa      setosa      setosa
##          7          8          9         10         11         12
##      setosa      setosa      setosa      setosa      setosa      setosa
##         13         14         15         16         17         18
##      setosa versicolor versicolor versicolor versicolor versicolor
##         19         20         21         22         23         24
## versicolor virginica versicolor versicolor versicolor versicolor
##         25         26         27         28         29         30
## versicolor versicolor versicolor versicolor versicolor versicolor
##         31         32         33         34         35         36
## virginica virginica virginica virginica virginica virginica
##         37         38         39         40         41         42
## virginica virginica virginica virginica virginica virginica
##         43         44         45
## virginica virginica virginica
## Levels: setosa versicolor virginica
```

```
tab <- table(predykcja = pred.class, obserwacja = dt.test$Species)
tab
```

```
##          obserwacja
## predykcja  setosa versicolor virginica
##   setosa      13          0          0
## versicolor    0         16          0
## virginica     0          1         15
```

Jak widać z powyższej tabeli, model całkiem dobrze radzi sobie z poprawną klasyfikacją obserwacji do odpowiednich kategorii. Tylko jedna obserwacja została błędnie zaklasyfikowana.

W dalszej kolejności sprawdzimy, czy nie jest konieczne przycięcie drzewa. Jednym z kryteriów przycinania drzewa jest przycinanie ze względu na złożoność drzewa. W tym przypadku jest wyrażony parametrem `cp`. Istnieje powszechnie stosowana reguła jednego odchylenia standardowego, która mówi, że drzewo należy przyciąć wówczas, gdy błąd oszacowany na podstawie sprawdzianu krzyżowego (`xerror`), pierwszy raz zejdzie poniżej poziomu wyznaczonego przez najniższą wartość błędu powiększonego o odchylenie standardowe tego błędu (`xstd`). Na podstawie poniższej tabeli można ustalić, że poziomem odcięcia jest wartość  $0.10294 + 0.037589 = 0.140529$ . Pierwszy raz błąd przyjmuje wartość mniejszą od 0.140529 po drugim podziale (`nsplit=2`). Temu poziomowi odpowiada `cp` o wartości 0.029412 i to jest złożoność drzewa, którą powinniśmy przyjąć do przycięcia drzewa.



Rysunek 4.3: Na wykresie błędów punkt odcięcia zaznaczony jest linią przerywaną

```
printcp(mod.rpart)
```

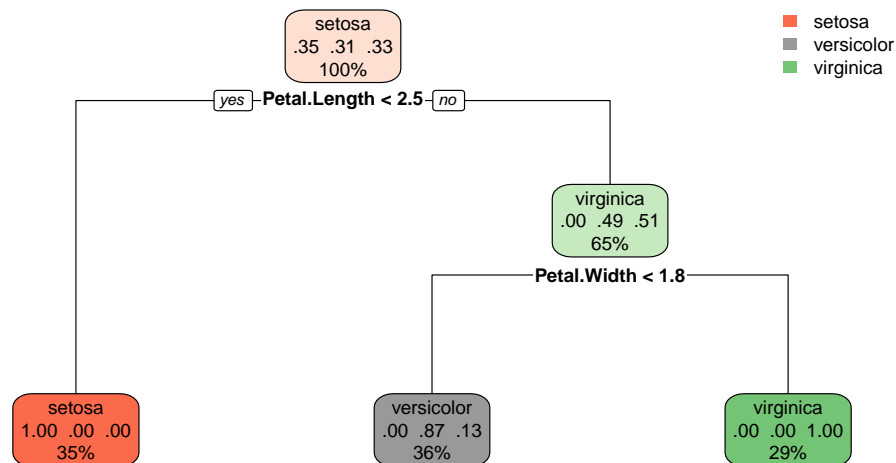
```
##
## Classification tree:
## rpart(formula = Species ~ ., data = dt.train, control = rpart.control(minsplit = 10,
##   minbucket = 5, maxdepth = 4))
##
## Variables actually used in tree construction:
## [1] Petal.Length Petal.Width
##
## Root node error: 68/105 = 0.64762
##
## n= 105
##
##      CP nsplit rel error  xerror    xstd
## 1 0.514706     0  1.000000 1.17647 0.064182
## 2 0.411765     1  0.485294 0.66176 0.074572
## 3 0.029412     2  0.073529 0.10294 0.037589
## 4 0.010000     3  0.044118 0.10294 0.037589
```

```
plotcp(mod.rpart)
```

Przycięte drzewo wygląda następująco:

```
mod.rpart2 <- prune(mod.rpart, cp = 0.029412)
summary(mod.rpart2)
```

```
## Call:
## rpart(formula = Species ~ ., data = dt.train, control = rpart.control(minsplit = 10,
##   minbucket = 5, maxdepth = 4))
##   n= 105
##
##           CP nsplit  rel error   xerror   xstd
## 1 0.5147059      0 1.00000000 1.1764706 0.06418173
## 2 0.4117647      1 0.48529412 0.6617647 0.07457243
## 3 0.0294120      2 0.07352941 0.1029412 0.03758880
##
## Variable importance
##   Petal.Width Petal.Length Sepal.Length  Sepal.Width
##           35           31           22           12
##
## Node number 1: 105 observations,   complexity param=0.5147059
##   predicted class=setosa   expected loss=0.647619  P(node) =1
##   class counts:    37    33    35
##   probabilities: 0.352 0.314 0.333
##   left son=2 (37 obs) right son=3 (68 obs)
##   Primary splits:
##     Petal.Length < 2.45 to the left,   improve=35.95322, (0 missing)
##     Petal.Width < 0.8 to the left,   improve=35.95322, (0 missing)
##     Sepal.Length < 5.45 to the left,   improve=25.39467, (0 missing)
##     Sepal.Width < 3.35 to the right,   improve=12.69596, (0 missing)
##   Surrogate splits:
##     Petal.Width < 0.8 to the left,   agree=1.000, adj=1.000, (0 split)
##     Sepal.Length < 5.45 to the left,   agree=0.924, adj=0.784, (0 split)
##     Sepal.Width < 3.35 to the right,   agree=0.819, adj=0.486, (0 split)
##
## Node number 2: 37 observations
##   predicted class=setosa   expected loss=0  P(node) =0.352381
##   class counts:    37      0      0
##   probabilities: 1.000 0.000 0.000
##
## Node number 3: 68 observations,   complexity param=0.4117647
##   predicted class=virginica expected loss=0.4852941  P(node) =0.647619
##   class counts:      0    33    35
##   probabilities: 0.000 0.485 0.515
##   left son=6 (38 obs) right son=7 (30 obs)
##   Primary splits:
##     Petal.Width < 1.75 to the left,   improve=25.286380, (0 missing)
##     Petal.Length < 4.75 to the left,   improve=24.879360, (0 missing)
##     Sepal.Length < 5.75 to the left,   improve= 6.713875, (0 missing)
##     Sepal.Width < 3.25 to the left,   improve= 1.336180, (0 missing)
##   Surrogate splits:
##     Petal.Length < 4.75 to the left,   agree=0.882, adj=0.733, (0 split)
##     Sepal.Length < 6.15 to the left,   agree=0.721, adj=0.367, (0 split)
```



Rysunek 4.4: Drzewo klasyfikacyjne po przycięciu

```
##      Sepal.Width < 3.15 to the left,  agree=0.618, adj=0.133, (0 split)
##
## Node number 6: 38 observations
##   predicted class=versicolor  expected loss=0.1315789  P(node) =0.3619048
##   class counts:      0      33      5
##   probabilities: 0.000 0.868 0.132
##
## Node number 7: 30 observations
##   predicted class=virginica   expected loss=0  P(node) =0.2857143
##   class counts:      0      0      30
##   probabilities: 0.000 0.000 1.000
rpart.plot(mod.rpart2)
```

### Ocena dopasowania modelu

Na koniec budowy modelu należy sprawdzić jego jakość na zbiorze testowym.

```
pred.class2 <- predict(mod.rpart2,
                      newdata = dt.test,
                      type = "class")
tab2 <- table(predykcja = pred.class2, obserwacja = dt.test$Species)
tab2
```

```
##          obserwacja
## predykcja  setosa versicolor virginica
##   setosa      13         0         0
##   versicolor  0         16         0
##   virginica   0          1        15
```

Mimo przycięcia drzewa, klasyfikacja pozostaje na niezmiennym poziomie. Odsetek poprawnych klasyfikacji możemy oszacować za pomocą

```
round(sum(diag(tab2))/sum(tab2)*100,1)
```

```
## [1] 97.8
```

## 4.9 Inne algorytmy budowy drzew decyzyjnych implementowane w R

Oprócz najbardziej znanego algorytmu CART implementowanego w postaci funkcji pakietu **rpart**, istnieją również inne algorytmy, które znalazły swoje implementacje w R. Są to:

- *CHAID*<sup>3</sup> - algorytm przeznaczony do budowy drzew klasyfikacyjnych, gdzie zarówno zmienna wynikowa, jak i zmienne niezależne muszą być ze skali jakościowej. Główną różnicą w stosunku do drzew typu CART jest sposób budowy podziałów, oparty na teście niezależności  $\chi^2$  Pearsona. Wyboru reguły podziału dokonuje się poprzez testowanie niezależności zmiennej niezależnej z predyktorami. Reguła o największej wartości statystyki  $\chi^2$  jest stosowana w pierwszej kolejności. Implementacja tego algorytmu znajduje się w pakiecie **CHAID**<sup>4</sup> (funkcja do tworzenia drzewa o tej samej nazwie **chaid**) (Team 2015).
- *Ctree*<sup>5</sup> - algorytm zbliżony zasadą działania do CHAID, ponieważ również wykorzystuje testowanie do wyboru reguły podziału. Różni się jednak tym, że może być stosowany do zmiennych dowolnego typu oraz tym, że może być zarówno drzewem klasyfikacyjnym jak i regresyjnym. Implementację R-ową można znaleźć w pakietach **party** (Hothorn, Hornik, and Zeileis 2006) lub **partykit** (Hothorn and Zeileis 2015) - funkcją do tworzenia modelu jest **ctree**.
- *C4.5* - algorytm stworzony przez Quinlan (1993) w oparciu, o również jego autorstwa, algorytm ID3. Służy jedynie do zadań klasyfikacyjnych. W dużym uproszczeniu, dobór reguł podziału odbywa się na podstawie przyrostu informacji (patrz Reguły podziału). W przeciwieństwie do pierwotnego algorytmu ID3, C4.5 nie raczej nie przeucza drzew. Implementacja R-owa znajduje się w pakiecie **RWeka** (Hornik, Buchta, and Zeileis 2009) - funkcją do budowy drzewa to **J48**.
- *C5.0* - kolejny algorytm autorstwa Kuhn and Quinlan (2018) jest usprawnieniem algorytmu C4.5, generującym mniejsze drzewa automatycznie przycinane na podstawie złożoności drzewa. Służy jedynie do zadań klasyfikacyjnych. Jest szybszy od poprzednika i pozwala na zastosowanie metody *boosting*<sup>6</sup>. Implementacja R-owa znajduje się w pakiecie *C50*, a funkcją do budowy drzewa to **C5.0**.

**Przykład 4.2.** W celu porównania wyników klasyfikacji na podstawie drzew decyzyjnych o różnych algorytmach, zostaną nauczone modele w oparciu o funkcje **ctree**, **J48** i **C5.0** dla tego samego zestawu danych co w przykładzie wcześniejszym 4.1.

### • Drzewo **ctree**

Na początek ustalamy parametry ograniczające rozrost drzewa podobne jak w poprzednim przykładzie.

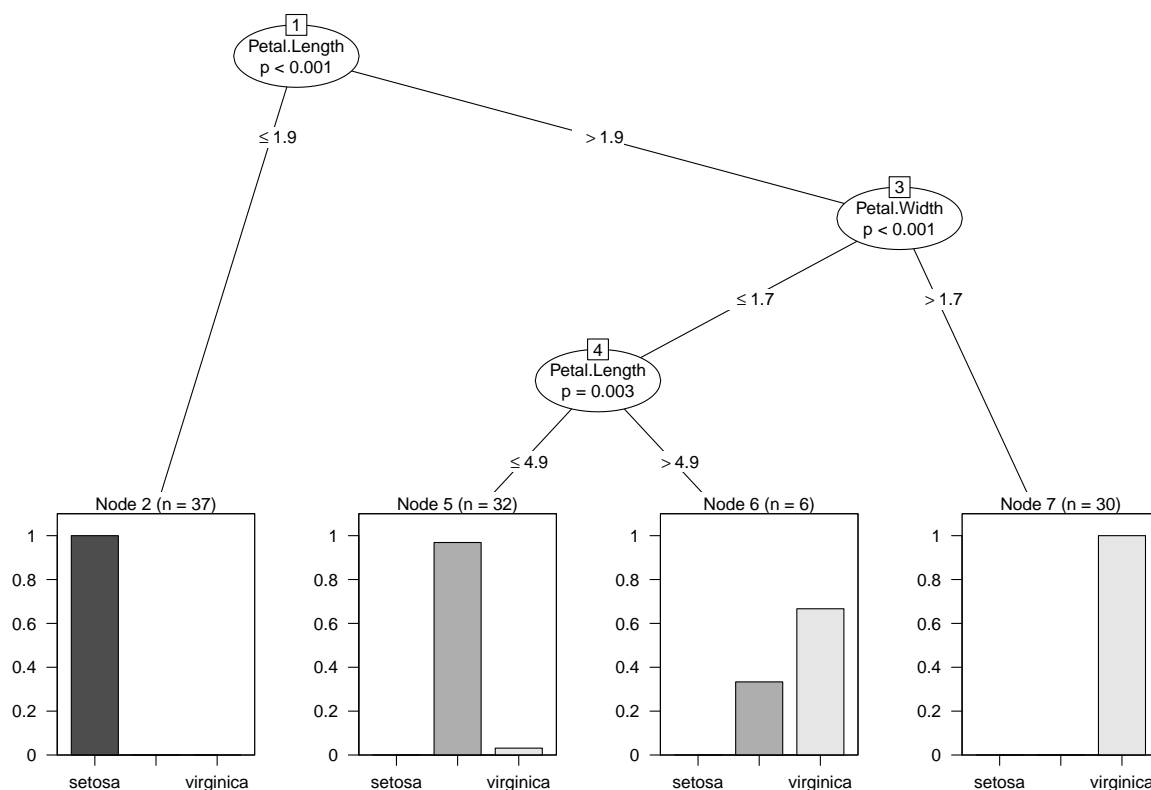
```
library(partykit)
tree2 <- ctree(Species~., data = dt.train,
               control = ctree_control(minsplit = 10,
```

<sup>3</sup>Chi-square automatic interaction detection

<sup>4</sup>brak w oficjalnej dystrybucji CRAN

<sup>5</sup>Conditional Inference Trees

<sup>6</sup>budowa klasyfikatora w oparciu o proces iteracyjny, w którym kolejne w kolejnych iteracjach budowane są proste drzewa i przypisywane są im wagi - im gorszy klasyfikator, tym większa waga - po to aby nauczyć drzewo klasyfikować "trudne" przypadki



Rysunek 4.5: Wykres drzewa decyzyjnego zbudowanego metodą ctree

```

minbucket = 5,
maxdepth = 4))

tree2

##
## Model formula:
## Species ~ Sepal.Length + Sepal.Width + Petal.Length + Petal.Width
##
## Fitted party:
## [1] root
## |   [2] Petal.Length <= 1.9: setosa (n = 37, err = 0.0%)
## |   [3] Petal.Length > 1.9
## |   |   [4] Petal.Width <= 1.7
## |   |   |   [5] Petal.Length <= 4.9: versicolor (n = 32, err = 3.1%)
## |   |   |   [6] Petal.Length > 4.9: virginica (n = 6, err = 33.3%)
## |   |   [7] Petal.Width > 1.7: virginica (n = 30, err = 0.0%)
##
## Number of inner nodes:    3
## Number of terminal nodes: 4

plot(tree2)

```

Wydaje się, że drzewo nie jest optymalne, ponieważ w węźle 6 obserwacje z grup *versicolor* i *virginica* są nieco pomieszane. Ostateczne oceny dokonujemy na podstawie próby testowej.

```
pred2 <- predict(tree2, newdata = dt.test)
tab <- table(predykcja = pred2, obserwacja = dt.test$Species)
tab
```

```
##              obserwacja
## predykcja   setosa versicolor virginica
## setosa      13         0          0
## versicolor   0        16          0
## virginica    0         1         15
```

Dopiero ocena jakości klasyfikacji na podstawie próby testowej pokazuje, że model zbudowany za pomocą ctree daje podobną precyzję jak rpart przycięty.

- Drzewo J48

W tym przypadku model sam poszukuje optymalnego rozwiązania przycinając się automatycznie.

```
library(RWeka)
tree3 <- J48(Species~., data = dt.train)
tree3

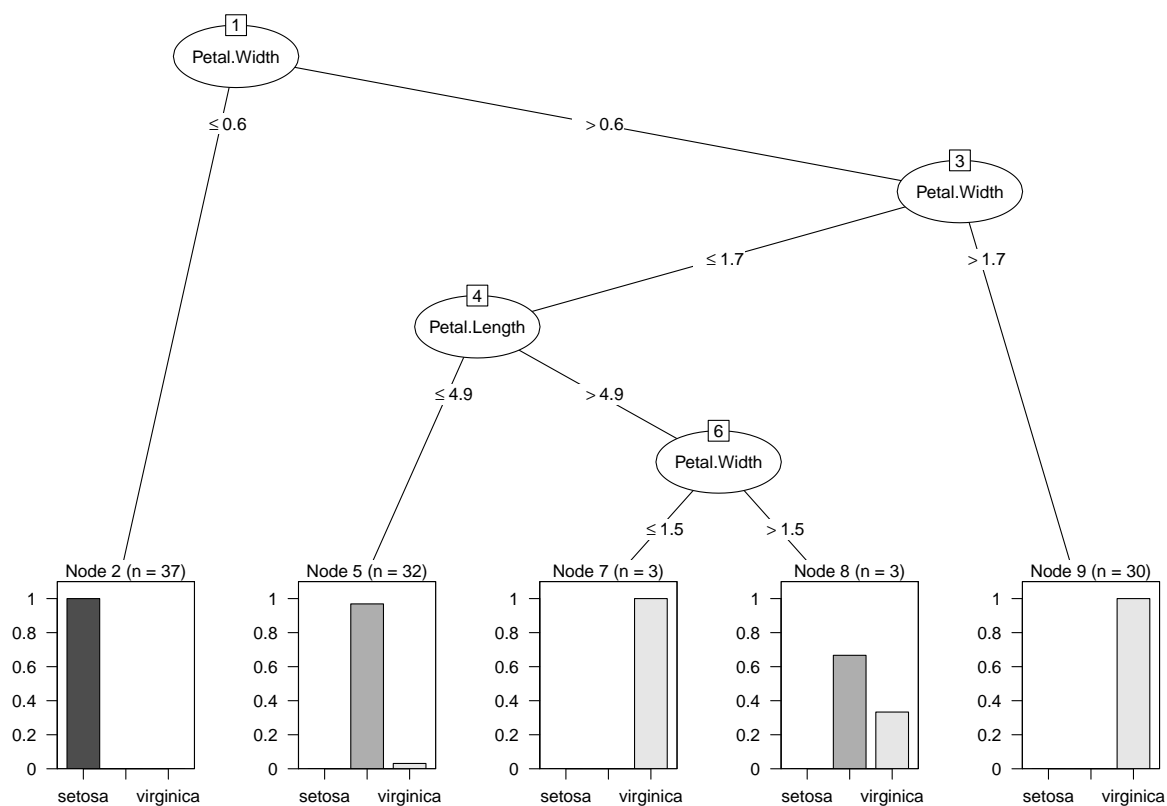
## J48 pruned tree
## -----
##
## Petal.Width <= 0.6: setosa (37.0)
## Petal.Width > 0.6
## |   Petal.Width <= 1.7
## |   |   Petal.Length <= 4.9: versicolor (32.0/1.0)
## |   |   Petal.Length > 4.9
## |   |   |   Petal.Width <= 1.5: virginica (3.0)
## |   |   |   Petal.Width > 1.5: versicolor (3.0/1.0)
## |   Petal.Width > 1.7: virginica (30.0)
##
## Number of Leaves   :   5
##
## Size of the tree   :   9

plot(tree3)
```

Drzewo jest nieco bardziej rozbudowane niż tree2 i mod.rpart2.

```
summary(tree3)

##
## === Summary ===
##
## Correctly Classified Instances      103          98.0952 %
## Incorrectly Classified Instances     2           1.9048 %
## Kappa statistic                     0.9714
## Mean absolute error                  0.0208
## Root mean squared error              0.1019
## Relative absolute error              4.6776 %
## Root relative squared error          21.628 %
## Total Number of Instances           105
##
## === Confusion Matrix ===
##
##  a  b  c  <-- classified as
```



Rysunek 4.6: Wykres drzewa decyzyjnego zbudowanego metodą J48



```
## 37 0 0 | a = setosa
## 0 33 0 | b = versicolor
## 0 2 33 | c = virginica
```

Podsumowanie dopasowania drzewa na próbie uczącej jest bardzo dobre, bo poprawnych klasyfikacji jest ponad 98%. Oceny dopasowania i tak dokonujemy na zbiorze testowym.

```
pred3 <- predict(tree3, newdata = dt.test)
tab <- table(predykacja = pred3, obserwacja = dt.test$Species)
tab
```

```
##           obserwacja
## predykacja  setosa versicolor virginica
## setosa      13         0         0
## versicolor  0         16        0
## virginica   0         1         15
```

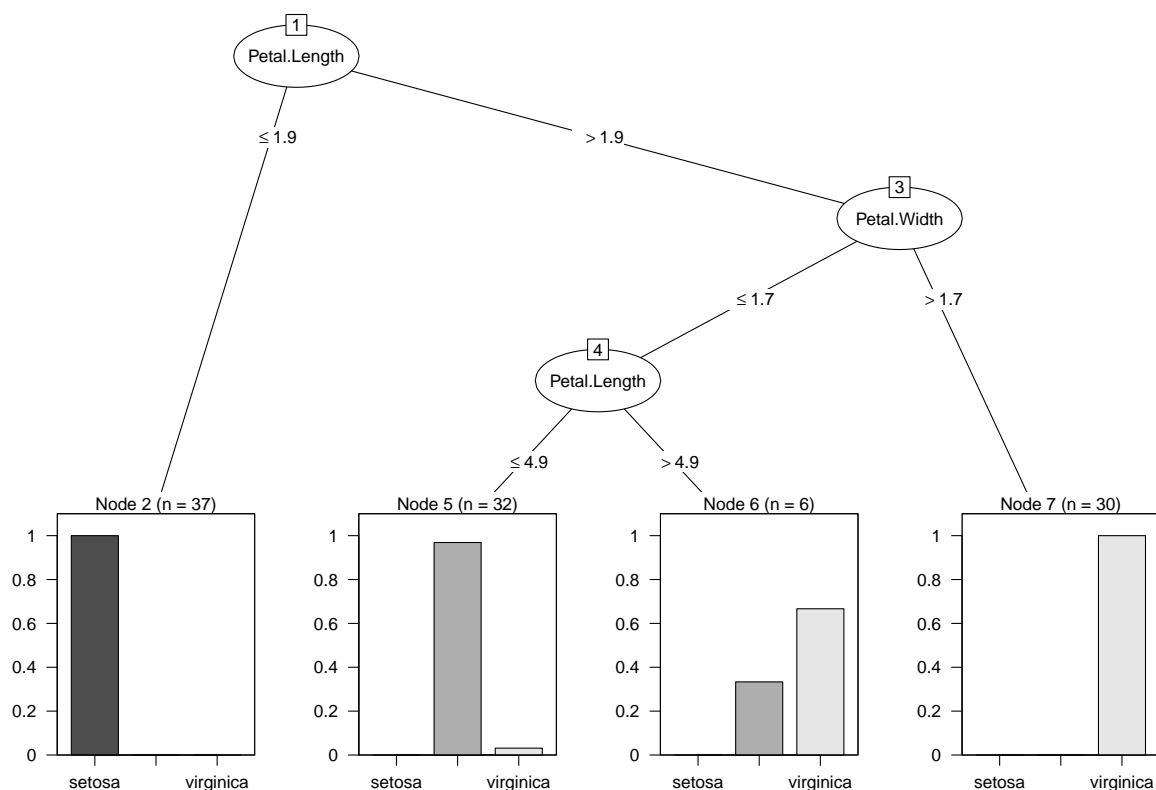
Otrzymujemy identyczną macierz klasyfikacji jak w poprzednich przypadkach.

- **Drzewo C50**

Tym razem również nie trzeba ustawiać parametrów drzewa, ponieważ algorytm działa tak aby zapobiec rozrostowi drzewa przy jednoczesnej wysokiej poprawności klasyfikacji.

```
library(C50)
tree4 <- C5.0(Species~., data = dt.train)
summary(tree4)
```

```
##
## Call:
## C5.0.formula(formula = Species ~ ., data = dt.train)
##
##
## C5.0 [Release 2.07 GPL Edition]      Tue Apr 16 12:14:35 2019
## -----
##
## Class specified by attribute `outcome'
##
## Read 105 cases (5 attributes) from undefined.data
##
## Decision tree:
##
## Petal.Length <= 1.9: setosa (37)
## Petal.Length > 1.9:
## :...Petal.Width > 1.7: virginica (30)
##   Petal.Width <= 1.7:
##     :...Petal.Length <= 4.9: versicolor (32/1)
##     Petal.Length > 4.9: virginica (6/2)
##
##
## Evaluation on training data (105 cases):
##
##      Decision Tree
##      -----
##      Size      Errors
##
##      4      3( 2.9%)  <<
```



Rysunek 4.7: Wykres drzewa decyzyjnego zbudowanego metodą C5.0

```
##
##      (a)  (b)  (c)    <-classified as
##      ----  ----  ----
##      37                    (a): class setosa
##                31    2    (b): class versicolor
##                1    34    (c): class virginica
##
##
## Attribute usage:
##
## 100.00% Petal.Length
##  64.76% Petal.Width
##
##
## Time: 0.0 secs
```

Otrzymujemy identyczne drzewo jak w przypadku zastosowania algorytmu `ctree`.

```
plot(tree4)
```

Dla pewności przeprowadzimy sprawdzenie na zbiorze testowym.

```
pred4 <- predict(tree4, newdata = dt.test)
tab <- table(predycja = pred4, obserwacja = dt.test$Species)
tab
```

```
##          obserwacja
## predykcja  setosa versicolor virginica
##   setosa      13         0         0
## versicolor   0         16         0
##  virginica    0          1        15
```



## Rozdział 5

# Pochodne drzew decyzyjnych

Przykład zastosowania drzew decyzyjnych na zbiorze `iris` w poprzednich przykładach może skłaniać do przypuszczenia, że drzewa decyzyjne zawsze dobrze radzą sobie z predykcją wartości wynikowej. Niestety w przykładach nieco bardziej skomplikowanych, gdzie chociażby klasy zmiennej wynikowej nie są tak wyraźnie separowalne, drzewa decyzyjne wypadają gorzej w porównaniu z innymi modelami nadzorowanego uczenia maszynowego.

I tak u podstaw metod bazujących na prostych drzewach decyzyjnych stał pomysł, że skoro jedno drzewo nie ma wystarczających własności predykcyjnych, to może zastosowanie wielu drzew połączonych w pewien sposób poprawi je. Tak powstały metody *bagging*, *random forest* i *boosting*<sup>1</sup>. Należy zaznaczyć, że metody znajdują swoje zastosowanie również w innych modelach nadzorowanego uczenia maszynowego.

### 5.1 Bagging

Technika ta została wprowadzona przez Breiman (1996) i ma na celu zmniejszenie wariancji modelu pojedynczego drzewa. Podobnie jak technika *bootstrap*, w której statystyki są wyliczane na wielu próbach pobranych z tego samego rozkładu (próby), w metodzie bagging losuje się wiele prób ze zbioru uczącego (najczęściej poprzez wielokrotne losowanie próby o rozmiarze zbioru uczącego ze zwracaniem), a następnie dla każdej próby bootstrapowej buduje się drzewo. W ten sposób otrzymujemy  $B$  drzew decyzyjnych  $\hat{f}^1(x), \hat{f}^2(x), \dots, \hat{f}^B(x)$ . Na koniec poprzez uśrednienie otrzymujemy model charakteryzujący się większą precyzją

$$\hat{f}_{bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^b(x). \quad (5.1)$$

Ponieważ podczas budowy drzew na podstawie prób bootstrapowych nie kontrolujemy złożoności, to w rezultacie każde z drzew może charakteryzować się dużą wariancją. Poprzez uśrednianie wyników pojedynczych drzew otrzymujemy mniejsze obciążenie ale również przy dostatecznie dużej liczbie prób ( $B$  często liczy się w setkach, czy tysiącach) zmniejszamy wariancję “średniej” predykcji z drzew. Oczywiście metodę tą trzeba dostosować do zadań klasyfikacyjnych, ponieważ nie istnieje średnia klasyfikacji z wielu drzew. W miejsce średniej stosuje się modę, czyli wartość dominującą.

Przyjrzyjmy się jak maszyna losuje obserwacje ze zwracaniem

```
n <- NULL
m <- NULL
for(i in 1:1000){
  x <- sample(1:500, size = 500, replace = T)
```

<sup>1</sup>chyba tylko dla drugiej metody istnieją dobre polskie tłumaczenie nazwy - las losowy

```

y <- setdiff(1:500, x)
z <- unique(x)
n[i] <- length(z)
m[i] <- length(y)
}
mean(n)/500*100

```

```
## [1] 63.2574
```

```
mean(m)/500*100
```

```
## [1] 36.7426
```

Faktycznie uczenie modelu metodą bagging odbywa się średnio na 2/3 obserwacji zbioru uczącego wylosowanych do prób bootstrapowych, a pozostała 1/3 (ang. *out-of-bag*) jest wykorzystana do oceny jakości predykcji.

Niewątpliwą zaletą drzew decyzyjnych była ich łatwa interpretacja. W przypadku metody bagging jest ona znacznie utrudniona, ponieważ jej wynik składa się z agregacji wielu drzew. Można natomiast ocenić ważność predyktorów (ang. *variable importance*). I tak, przez obserwację spadku *RSS* dla baggingu regresyjnego przy zastosowaniu danego predyktora w podziałach drzewa i uśrednieniu wyniku otrzymamy wskaźnik ważności predyktora dużo lepszy niż dla pojedynczego drzewa. W przypadku baggingu klasyfikacyjnego w miejsce *RSS* stosujemy indeks Gini’ego.

Implementacja R-owa metody bagging znajduje się w pakiecie **ipred**, a funkcja do budowy modelu nazywa się **bagging** (Peters and Hothorn 2018). Można również stosować funkcję **randomForest** pakietu **randomForest** (Liaw and Wiener 2002) - powody takiego działania wyjaśnia się w podrozdziale Lasy losowe.

**Przykład 5.1.** Tym razem cel zadania jest regresyjny i polega na ustaleniu miary tendencji centralnej ceny mieszkań w Bostonie na podstawie zmiennych umieszczonych w zbiorze **Boston** pakietu **MASS** (Venables and Ripley 2002). Zmienną zależną będzie mediana cen mieszkań na przedmieściach Bostonu (**medv**).

```

library(MASS)
head(Boston)

```

```

##      crim zn indus chas   nox    rm  age    dis rad tax ptratio  black
## 1 0.00632 18  2.31    0 0.538 6.575 65.2 4.0900   1  296    15.3 396.90
## 2 0.02731  0  7.07    0 0.469 6.421 78.9 4.9671   2  242    17.8 396.90
## 3 0.02729  0  7.07    0 0.469 7.185 61.1 4.9671   2  242    17.8 392.83
## 4 0.03237  0  2.18    0 0.458 6.998 45.8 6.0622   3  222    18.7 394.63
## 5 0.06905  0  2.18    0 0.458 7.147 54.2 6.0622   3  222    18.7 396.90
## 6 0.02985  0  2.18    0 0.458 6.430 58.7 6.0622   3  222    18.7 394.12
##    lstat medv
## 1   4.98 24.0
## 2   9.14 21.6
## 3   4.03 34.7
## 4   2.94 33.4
## 5   5.33 36.2
## 6   5.21 28.7

```

```

set.seed(2019)
boston.train <- Boston %>%
  sample_frac(size = 2/3)
boston.test <- setdiff(Boston, boston.train)

```

Aby móc porównać wyniki predykcji z metody bagging, najpierw zostanie zbudowane jedno drzewo decyzyjne w oparciu o algorytm CART.

```
library(rpart)
library(rpart.plot)
boston.rpart <- rpart(medv ~ ., data = boston.train)
x <- summary(boston.rpart)
```

```
## Call:
## rpart(formula = medv ~ ., data = boston.train)
##   n= 337
##
##              CP nsplit rel error   xerror   xstd
## 1  0.43506104      0 1.0000000 1.0037495 0.10496568
## 2  0.21114710      1 0.5649390 0.6856438 0.07732133
## 3  0.05641774      2 0.3537919 0.4393220 0.05974589
## 4  0.04154842      3 0.2973741 0.3726563 0.05716622
## 5  0.02707678      4 0.2558257 0.3520312 0.05569786
## 6  0.01489117      5 0.2287489 0.3238915 0.05681943
## 7  0.01202564      6 0.2138578 0.2922610 0.05311293
## 8  0.01057622      7 0.2018321 0.2889364 0.05318206
## 9  0.01031677      8 0.1912559 0.2838433 0.05152251
## 10 0.01006729      9 0.1809391 0.2838187 0.05152098
## 11 0.01000000     10 0.1708718 0.2815210 0.05152993
##
## Variable importance
##   lstat      nox      indus      crim      tax      rm      age      dis ptratio
##      24       13       13       13       11       10       10       2      2
##   rad   black
##      1       1
##
## Node number 1: 337 observations,      complexity param=0.435061
##   mean=22.61157, MSE=79.33004
##   left son=2 (186 obs) right son=3 (151 obs)
##   Primary splits:
##     lstat < 10.02      to the right, improve=0.4350610, (0 missing)
##     rm    < 6.8375     to the left,  improve=0.4305766, (0 missing)
##     indus < 6.66      to the right, improve=0.2914821, (0 missing)
##     ptratio < 19.15   to the right, improve=0.2608119, (0 missing)
##     nox    < 0.5125   to the right, improve=0.2169607, (0 missing)
##   Surrogate splits:
##     indus < 7.625     to the right, agree=0.846, adj=0.656, (0 split)
##     nox    < 0.519     to the right, agree=0.828, adj=0.616, (0 split)
##     crim   < 0.12995   to the right, agree=0.786, adj=0.523, (0 split)
##     age    < 63.9      to the right, agree=0.777, adj=0.503, (0 split)
##     tax    < 377       to the right, agree=0.769, adj=0.483, (0 split)
##
## Node number 2: 186 observations,      complexity param=0.05641774
##   mean=17.31828, MSE=19.86042
##   left son=4 (58 obs) right son=5 (128 obs)
##   Primary splits:
##     crim < 5.84803     to the right, improve=0.4083024, (0 missing)
##     dis  < 2.0754      to the left,  improve=0.3684093, (0 missing)
##     lstat < 14.405     to the right, improve=0.3516672, (0 missing)
##     nox  < 0.657       to the right, improve=0.3255969, (0 missing)
##     age  < 84.9        to the right, improve=0.2247741, (0 missing)
##   Surrogate splits:
```

```

##      rad   < 16      to the right, agree=0.855, adj=0.534, (0 split)
##      tax   < 551.5   to the right, agree=0.839, adj=0.483, (0 split)
##      nox   < 0.657   to the right, agree=0.828, adj=0.448, (0 split)
##      dis   < 2.0754  to the left,  agree=0.801, adj=0.362, (0 split)
##      lstat < 19.055   to the right, agree=0.796, adj=0.345, (0 split)
##
## Node number 3: 151 observations,      complexity param=0.2111471
##   mean=29.13179, MSE=75.5574
##   left son=6 (120 obs) right son=7 (31 obs)
##   Primary splits:
##     rm      < 7.127   to the left,  improve=0.4947648, (0 missing)
##     lstat   < 4.495   to the right, improve=0.4054324, (0 missing)
##     nox     < 0.574   to the left,  improve=0.1389706, (0 missing)
##     ptratio < 14.75   to the right, improve=0.1349232, (0 missing)
##     age     < 89.45   to the left,  improve=0.1133301, (0 missing)
##   Surrogate splits:
##     lstat   < 3.21    to the right, agree=0.841, adj=0.226, (0 split)
##     ptratio < 14.15   to the right, agree=0.828, adj=0.161, (0 split)
##     tax     < 207     to the right, agree=0.808, adj=0.065, (0 split)
##     nox     < 0.639   to the left,  agree=0.801, adj=0.032, (0 split)
##
## Node number 4: 58 observations
##   mean=13.08793, MSE=14.14485
##
## Node number 5: 128 observations,      complexity param=0.01489117
##   mean=19.23516, MSE=10.66681
##   left son=10 (61 obs) right son=11 (67 obs)
##   Primary splits:
##     lstat   < 14.405   to the right, improve=0.2915760, (0 missing)
##     dis     < 1.99235  to the left,  improve=0.2280873, (0 missing)
##     age     < 84.15    to the right, improve=0.1950219, (0 missing)
##     ptratio < 20.95    to the right, improve=0.1349341, (0 missing)
##     rm      < 5.706    to the left,  improve=0.1194638, (0 missing)
##   Surrogate splits:
##     age     < 91.15    to the right, agree=0.758, adj=0.492, (0 split)
##     dis     < 2.0418   to the left,  agree=0.664, adj=0.295, (0 split)
##     nox     < 0.607    to the right, agree=0.633, adj=0.230, (0 split)
##     indus   < 18.84    to the right, agree=0.625, adj=0.213, (0 split)
##     rm      < 5.703    to the left,  agree=0.617, adj=0.197, (0 split)
##
## Node number 6: 120 observations,      complexity param=0.04154842
##   mean=26.02417, MSE=34.39883
##   left son=12 (98 obs) right son=13 (22 obs)
##   Primary splits:
##     lstat   < 5.145    to the right, improve=0.2690898, (0 missing)
##     dis     < 2.0891   to the right, improve=0.2163813, (0 missing)
##     rm      < 6.543    to the left,  improve=0.2036454, (0 missing)
##     age     < 89.45    to the left,  improve=0.1796977, (0 missing)
##     tax     < 548      to the left,  improve=0.1751322, (0 missing)
##   Surrogate splits:
##     zn      < 92.5     to the left,  agree=0.833, adj=0.091, (0 split)
##     nox     < 0.4035   to the right, agree=0.833, adj=0.091, (0 split)
##     indus   < 1.495    to the right, agree=0.825, adj=0.045, (0 split)
##     dis     < 1.48495  to the right, agree=0.825, adj=0.045, (0 split)

```



```

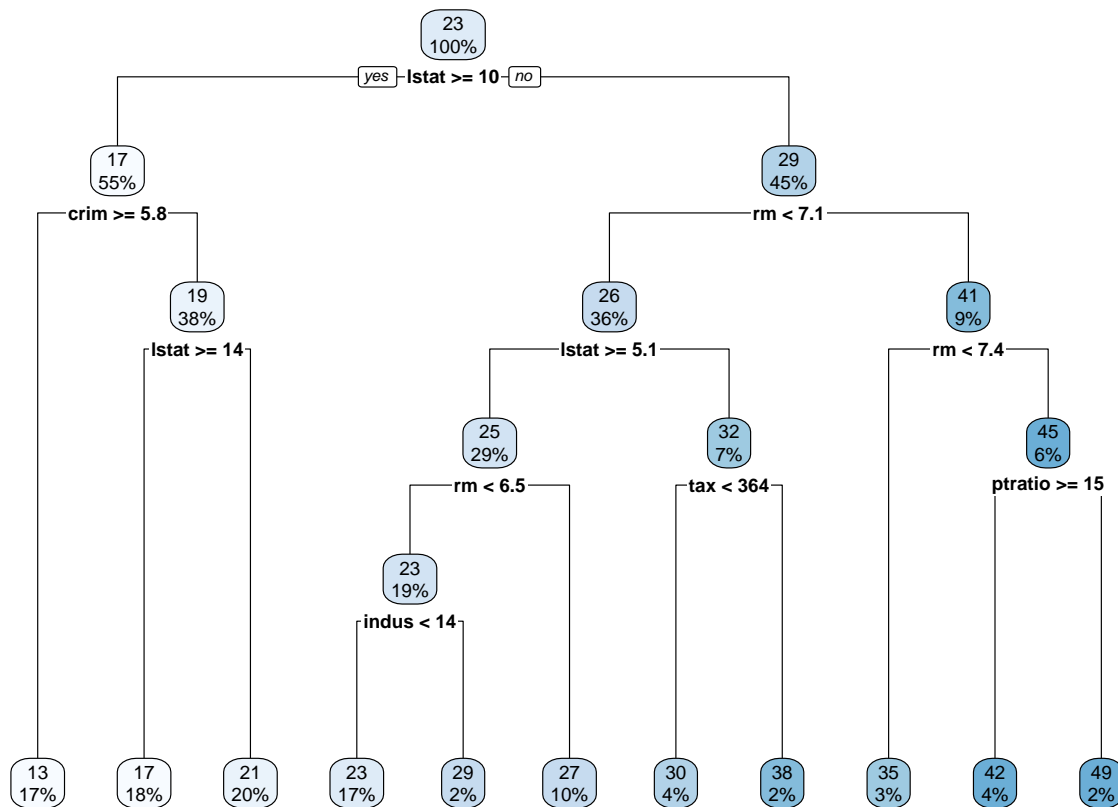
##
## Node number 7: 31 observations,      complexity param=0.02707678
##   mean=41.16129, MSE=52.78882
##   left son=14 (11 obs) right son=15 (20 obs)
##   Primary splits:
##       rm      < 7.437      to the left,  improve=0.4423448, (0 missing)
##       lstat    < 5.185      to the right, improve=0.3125696, (0 missing)
##       ptratio  < 15.05     to the right, improve=0.1896089, (0 missing)
##       black    < 392.715   to the right, improve=0.1133472, (0 missing)
##       age      < 37.6      to the right, improve=0.0737298, (0 missing)
##   Surrogate splits:
##       lstat < 4.635      to the right, agree=0.774, adj=0.364, (0 split)
##       indus < 2.32      to the left,  agree=0.742, adj=0.273, (0 split)
##       dis    < 5.9736    to the right, agree=0.710, adj=0.182, (0 split)
##       black < 390.095    to the right, agree=0.710, adj=0.182, (0 split)
##       crim   < 0.10593   to the left,  agree=0.677, adj=0.091, (0 split)
##
## Node number 10: 61 observations
##   mean=17.38689, MSE=8.122779
##
## Node number 11: 67 observations
##   mean=20.91791, MSE=7.041172
##
## Node number 12: 98 observations,      complexity param=0.01202564
##   mean=24.58265, MSE=20.9745
##   left son=24 (64 obs) right son=25 (34 obs)
##   Primary splits:
##       rm      < 6.543      to the left,  improve=0.1564077, (0 missing)
##       black    < 364.385   to the right, improve=0.1331323, (0 missing)
##       age      < 89.45     to the left,  improve=0.1241124, (0 missing)
##       tax      < 223.5     to the right, improve=0.1204819, (0 missing)
##       dis      < 4.46815   to the right, improve=0.1048755, (0 missing)
##   Surrogate splits:
##       dis      < 3.6589    to the right, agree=0.704, adj=0.147, (0 split)
##       rad      < 6.5       to the left,  agree=0.704, adj=0.147, (0 split)
##       age      < 68.9     to the left,  agree=0.694, adj=0.118, (0 split)
##       indus    < 1.605     to the right, agree=0.673, adj=0.059, (0 split)
##       nox      < 0.4045    to the right, agree=0.673, adj=0.059, (0 split)
##
## Node number 13: 22 observations,      complexity param=0.01031677
##   mean=32.44545, MSE=43.70884
##   left son=26 (15 obs) right son=27 (7 obs)
##   Primary splits:
##       tax      < 364       to the left,  improve=0.2868266, (0 missing)
##       lstat    < 3.855     to the right, improve=0.2413545, (0 missing)
##       age      < 31.85     to the left,  improve=0.1598075, (0 missing)
##       dis      < 5.4085    to the right, improve=0.1258591, (0 missing)
##       black    < 381.59    to the right, improve=0.1052855, (0 missing)
##   Surrogate splits:
##       crim     < 2.6956    to the left,  agree=0.773, adj=0.286, (0 split)
##       indus     < 14       to the left,  agree=0.773, adj=0.286, (0 split)
##       nox       < 0.5875    to the left,  agree=0.773, adj=0.286, (0 split)
##       age       < 89.65     to the left,  agree=0.773, adj=0.286, (0 split)
##       dis       < 2.3371    to the right, agree=0.773, adj=0.286, (0 split)

```

```

##
## Node number 14: 11 observations
##   mean=34.64545, MSE=3.304298
##
## Node number 15: 20 observations,   complexity param=0.01057622
##   mean=44.745, MSE=43.81147
##   left son=30 (12 obs) right son=31 (8 obs)
##   Primary splits:
##       ptratio < 15.4      to the right, improve=0.3226860, (0 missing)
##       rad      < 6        to the right, improve=0.2170243, (0 missing)
##       tax      < 270      to the right, improve=0.1545997, (0 missing)
##       age      < 71.85    to the right, improve=0.1331209, (0 missing)
##       zn       < 10       to the left,  improve=0.1328727, (0 missing)
##   Surrogate splits:
##       zn < 10             to the left,  agree=0.80, adj=0.500, (0 split)
##       nox < 0.541         to the left,  agree=0.80, adj=0.500, (0 split)
##       age < 86.7          to the left,  agree=0.80, adj=0.500, (0 split)
##       dis < 2.5813        to the right, agree=0.80, adj=0.500, (0 split)
##       crim < 0.45114      to the left,  agree=0.75, adj=0.375, (0 split)
##
## Node number 24: 64 observations,   complexity param=0.01006729
##   mean=23.2625, MSE=21.96891
##   left son=48 (57 obs) right son=49 (7 obs)
##   Primary splits:
##       indus < 14.48       to the left,  improve=0.19142190, (0 missing)
##       crim < 0.841845     to the left,  improve=0.17407590, (0 missing)
##       black < 374.635     to the right, improve=0.14590640, (0 missing)
##       dis < 2.6499        to the right, improve=0.13374910, (0 missing)
##       age < 79.85        to the left,  improve=0.08856433, (0 missing)
##   Surrogate splits:
##       crim < 1.163695     to the left,  agree=0.984, adj=0.857, (0 split)
##       nox < 0.589         to the left,  agree=0.984, adj=0.857, (0 split)
##       age < 84.35         to the left,  agree=0.984, adj=0.857, (0 split)
##       dis < 2.28545       to the right, agree=0.969, adj=0.714, (0 split)
##       black < 361.635     to the right, agree=0.969, adj=0.714, (0 split)
##
## Node number 25: 34 observations
##   mean=27.06765, MSE=9.646894
##
## Node number 26: 15 observations
##   mean=30.02667, MSE=14.56062
##
## Node number 27: 7 observations
##   mean=37.62857, MSE=66.76776
##
## Node number 30: 12 observations
##   mean=41.675, MSE=48.28521
##
## Node number 31: 8 observations
##   mean=49.35, MSE=1.7575
##
## Node number 48: 57 observations
##   mean=22.54386, MSE=10.87053
##

```



Rysunek 5.1: Drzewo regresyjne pełne

```
## Node number 49: 7 observations
## mean=29.11429, MSE=73.89265
```

```
rpart.plot(boston.rpart)
```

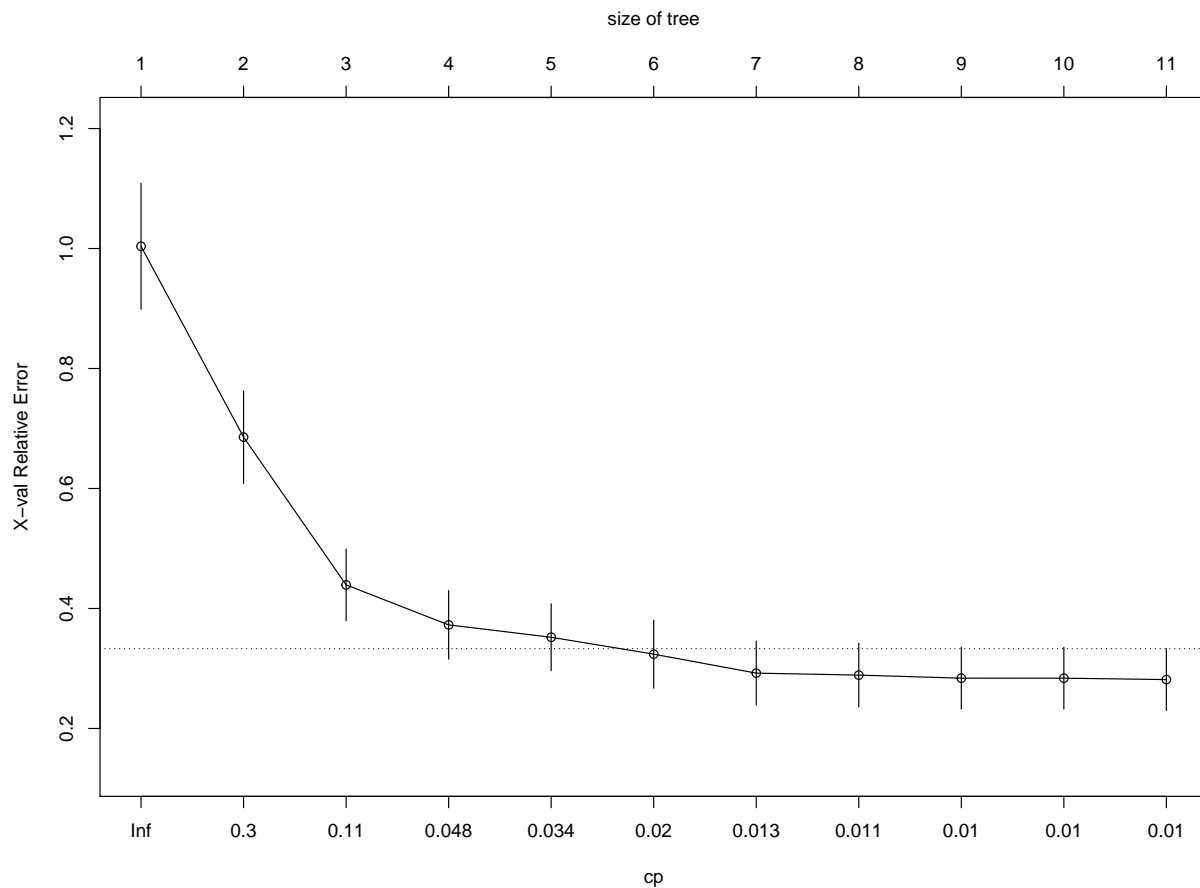
Przycinamy drzewo...

```
printcp(boston.rpart)
```

```
##
## Regression tree:
## rpart(formula = medv ~ ., data = boston.train)
##
## Variables actually used in tree construction:
## [1] crim   indus  lstat  ptratio rm    tax
##
## Root node error: 26734/337 = 79.33
##
## n= 337
##
##      CP nsplit rel error  xerror   xstd
## 1  0.435061      0  1.00000 1.00375 0.104966
## 2  0.211147      1  0.56494 0.68564 0.077321
```

```
## 3 0.056418      2 0.35379 0.43932 0.059746
## 4 0.041548      3 0.29737 0.37266 0.057166
## 5 0.027077      4 0.25583 0.35203 0.055698
## 6 0.014891      5 0.22875 0.32389 0.056819
## 7 0.012026      6 0.21386 0.29226 0.053113
## 8 0.010576      7 0.20183 0.28894 0.053182
## 9 0.010317      8 0.19126 0.28384 0.051523
## 10 0.010067     9 0.18094 0.28382 0.051521
## 11 0.010000    10 0.17087 0.28152 0.051530
```

```
plotcp(boston.rpart)
```



```
boston.rpart2 <- prune(boston.rpart, cp = 0.012026)
```

```
rpart.plot(boston.rpart2)
```

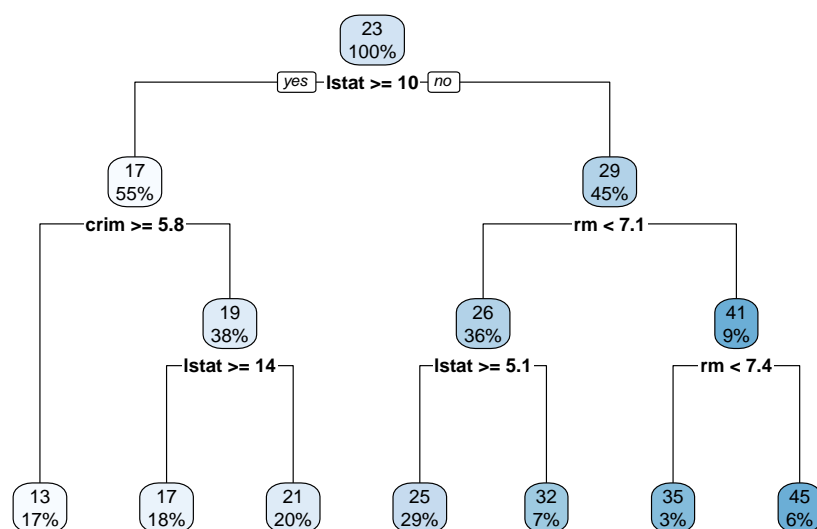
Predykcja na podstawie drzewa na zbiorze testowym.

```
boston.pred <- predict(boston.rpart2, newdata = boston.test)
rmse <- function(pred, obs) sqrt(1/length(pred)*sum((pred-obs)^2))
rmse(boston.pred, boston.test$medv)
```

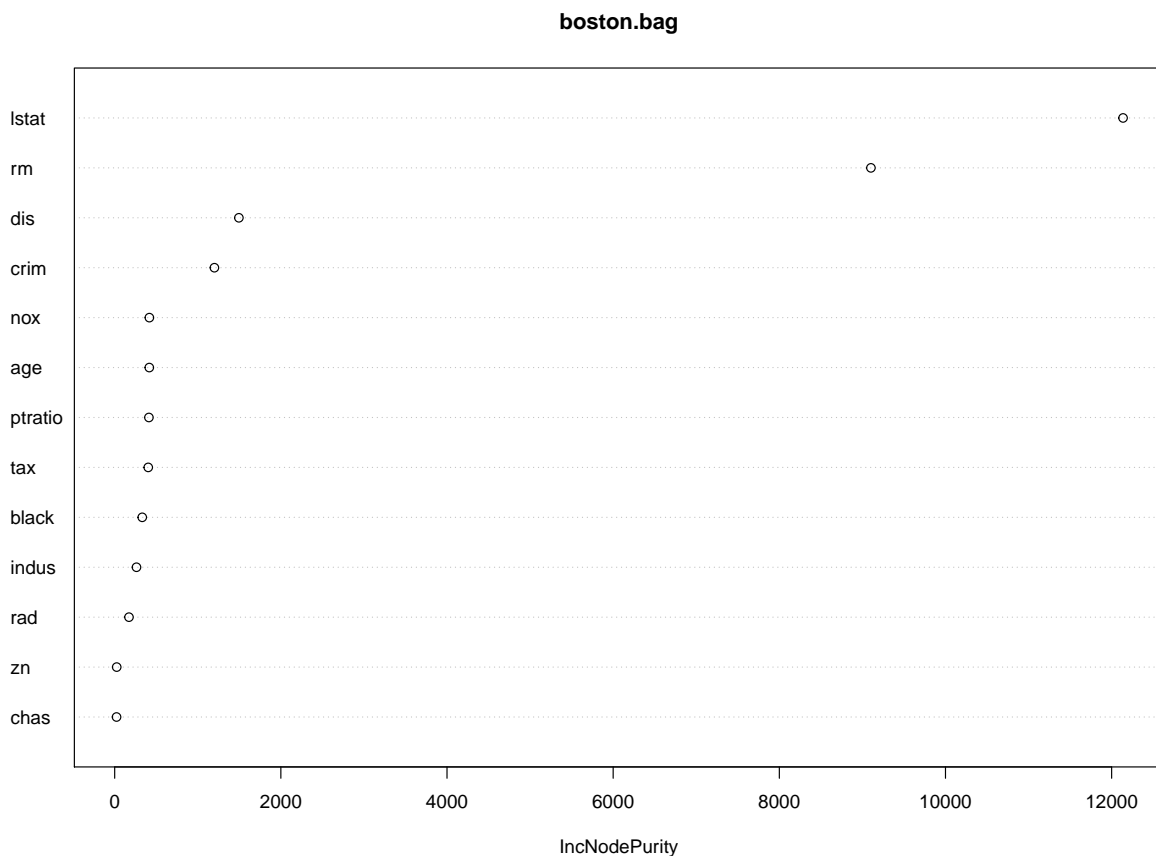
```
## [1] 4.825862
```

Teraz zbudujemy model metodą bagging.

```
library(randomForest)
boston.bag <- randomForest(medv~., data = boston.train,
```



Rysunek 5.2: Drzewo regresyjne przycięte



Rysunek 5.3: Wykres ważności predyktorów

```

                                mtry = ncol(boston.train)-1)
boston.bag

##
## Call:
## randomForest(formula = medv ~ ., data = boston.train, mtry = ncol(boston.train) - 1)
##               Type of random forest: regression
##               Number of trees: 500
## No. of variables tried at each split: 13
##
##               Mean of squared residuals: 13.06701
##               % Var explained: 83.53

Predykcja na podstawie modelu
boston.pred2 <- predict(boston.bag, newdata = boston.test)
rmse(boston.pred2, boston.test$medv)

## [1] 3.039308

```

Zatem predykcja na podstawie modelu bagging jest nieco lepsza niż z pojedynczego drzewa. Dodatkowo możemy ocenić ważność zmiennych użytych w budowie drzew.

```
varImpPlot(boston.bag)
```

```
importance(boston.bag)
```

```
##          IncNodePurity
## crim          1200.11828
## zn             24.17836
## indus          262.33396
## chas           22.27133
## nox            417.32236
## rm            9102.58339
## age           416.48170
## dis           1494.79734
## rad            171.92103
## tax            403.66309
## ptratio        411.88528
## black          331.58495
## lstat         12137.38999
```

```
x$variable.importance
```

```
##      lstat      nox      indus      crim      tax      rm
## 15197.8587 8683.8225 8325.2431 8074.7200 6991.0756 6768.5423
##      age      dis      ptratio      rad      black      zn
## 6538.5039 1305.3786 1193.2073 853.4309 323.8576 242.3521
```

W porównaniu do ważności zmiennych dla pojedynczego drzewa widać pewne różnice.

## 5.2 Lasy losowe

Lasy losowe są uogólnieniem metody bagging, polegającą na losowaniu dla każdego drzewa wchodzącego w skład lasu  $m$  predyktorów spośród  $p$  dostępnych, a następnie budowaniu drzew z wykorzystaniem tylko tych predyktorów (Ho 1995). Dzięki temu za każdy raz drzewo jest budowane w oparciu o nowy zestaw cech (najczęściej przyjmujemy  $m = \sqrt{p}$ ). W przypadku modeli bagging za każdym razem najsilniejszy predyktor wchodził w skład zbioru uczącego, a co za tym idzie również uczestniczył w tworzeniu reguł podziału. Wówczas wiele drzew zawierało reguły stosujące dany atrybut, a wtedy predykcje otrzymywane za pomocą drzew były skorelowane. Dlatego nawet duża liczba prób bootstrapowych nie zapewniała poprawy precyzji. Implementacja tej metody znajduje się w pakiecie **randomForest**.

**Przykład 5.2.** Kontynuując poprzedni przykład 5.1 możemy zbudować las losowy aby przekonać się czy nastąpi poprawa predykcji zmiennej wynikowej.

```
boston.rf <- randomForest(medv ~ ., data = boston.train)
boston.rf
```

```
##
## Call:
## randomForest(formula = medv ~ ., data = boston.train)
##              Type of random forest: regression
##              Number of trees: 500
## No. of variables tried at each split: 4
##
##              Mean of squared residuals: 13.09902
##              % Var explained: 83.49
```

Porównanie MSE na próbach uczących pomiędzy lasem losowym i modelem bagging wypada nieco na korzyść

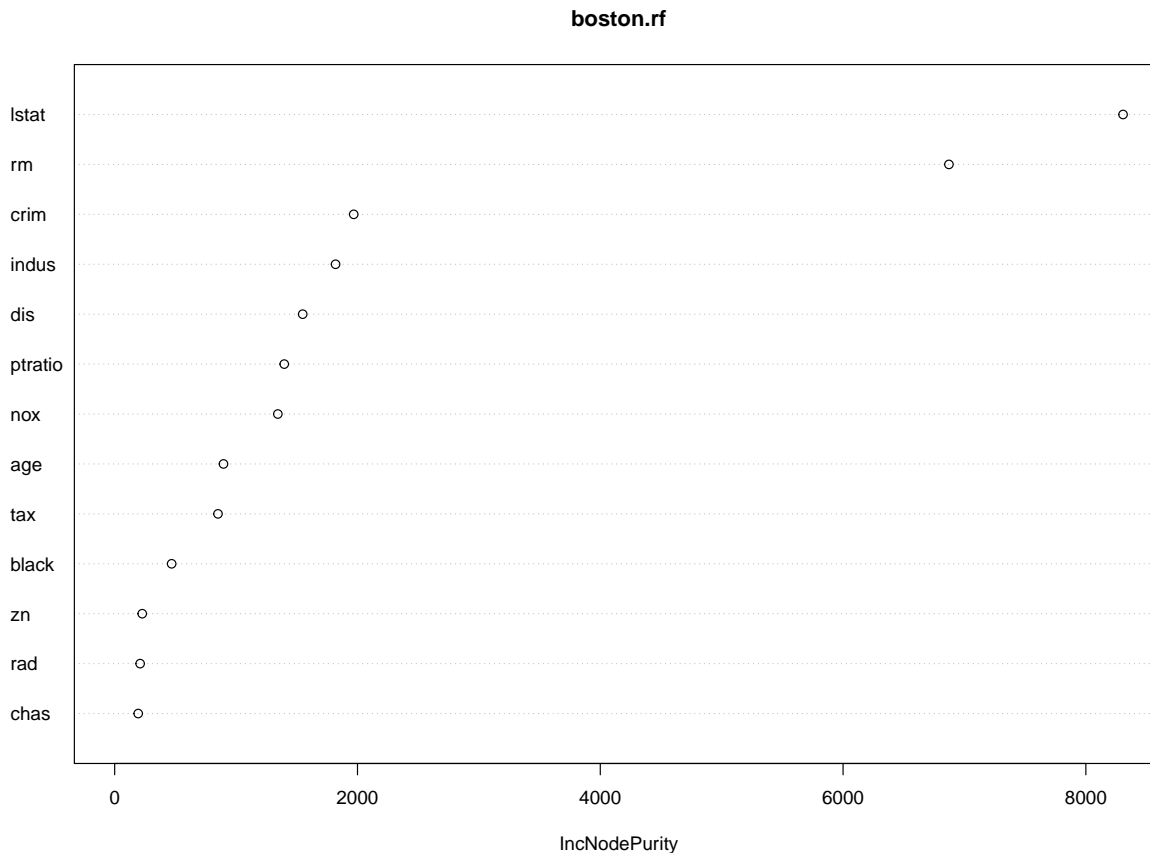
bagging.

```
boston.pred3 <- predict(boston.rf, newdata = boston.test)
rmse(boston.pred3, boston.test$medv)
```

```
## [1] 3.418302
```

Ważność zmiennych również się nieco różni.

```
varImpPlot(boston.rf)
```



### 5.3 Boosting

Rozważania na temat metody *boosting* zaczęły się od pytań postawionych w publikacji Kearns and Valiant (1989), czy da się na podstawie na podstawie zbioru słabych modeli stworzyć jeden dobry? Odpowiedzi pozytywnej na nie udzielili, najpierw Schapire (1990), a potem Breiman (1998). W metodzie boosting nie stosuje się prób bootstrapowych ale odpowiednio modyfikuje się drzewo wyjściowe w kolejnych krokach na tym samym zbiorze uczącym. Algorytm dla drzewa regresyjnego jest następujący:

1. Ustal  $\hat{f}(x) = 0$  i  $r_i = y_i$  dla każdego  $i$  w zbiorze uczącym.
2. Dla  $b = 1, 2, \dots, B$  powtarzaj:
  - a) naucz drzewo  $\hat{f}^b$  o  $d$  regułach podziału (czyli  $d + 1$  liściach) na zbiorze  $(X_i, r_i)$ ,
  - b) zaktualizuj drzewo do nowej “skurczonej” wersji

$$\hat{f}(x) \leftarrow \hat{f}(x) + \lambda \hat{f}^b(x), \quad (5.2)$$



c) zaktualizuj reszty

$$r_i \leftarrow r_i - \lambda \hat{f}^b(x_i). \quad (5.3)$$

3. Wyznacz boosted model

$$\hat{f}(x) = \sum_{b=1}^B \lambda \hat{f}^b(x) \quad (5.4)$$

Uczenie drzew klasyfikacyjnego metoda boosting przebiega w podobny sposób. Wynik uczenia drzew metodą boosting zależy od trzech parametrów:

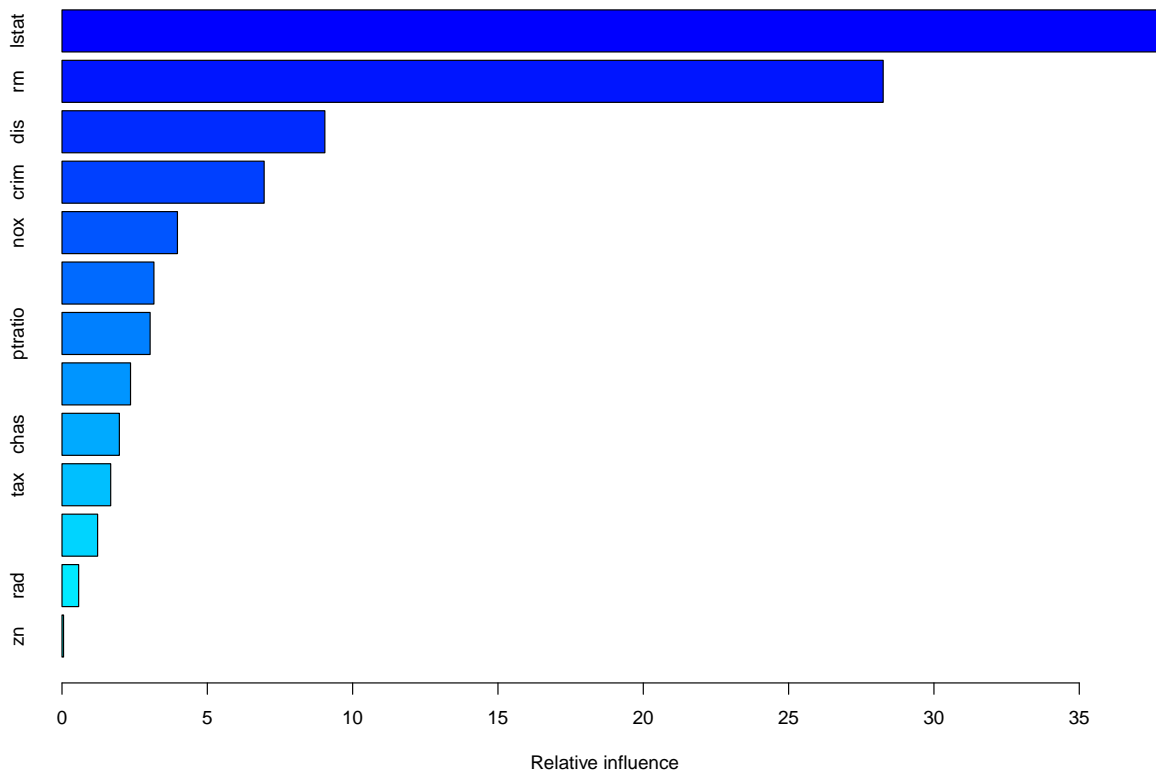
1. Liczby drzew  $B$ . W przeciwieństwie do metody bagging i lasów losowych, zbyt duże  $B$  może doprowadzić do przeuczenia modelu.  $B$  ustala się najczęściej na podstawie walidacji krzyżowej.
2. Parametru “kurczenia” (ang. *shrinkage*)  $\lambda$ . Kontroluje on szybkość uczenia się kolejnych drzew. Typowe wartości  $\lambda$  to 0.01 lub 0.001. Bardzo małe  $\lambda$  może wymagać dobrania większego  $B$ , aby zapewnić dobrą jakość predykcyjną modelu.
3. Liczby podziałów w drzewach  $d$ , która decyduje o złożoności drzewa. Bywa, że nawet  $d = 1$  daje dobre rezultaty, ponieważ model wówczas uczy się powoli.

Implementację metody boosting można znaleźć w pakiecie **gbm** (Greenwell et al. 2019)

**Przykład 5.3.** Metodę boosting zastosujemy do zadania predykcji ceny mieszkań na przedmieściach Bostonu. Dobór parametrów modelu będzie arbitralny, więc niekoniecznie model będzie najlepiej dopasowany.

```
library(gbm)
boston.boost <- gbm(medv ~ ., data = boston.train,
                    distribution = "gaussian",
                    n.trees = 5000,
                    interaction.depth = 2,
                    shrinkage = 0.01)
boston.boost
```

```
## gbm(formula = medv ~ ., distribution = "gaussian", data = boston.train,
##      n.trees = 5000, interaction.depth = 2, shrinkage = 0.01)
## A gradient boosted model with gaussian loss function.
## 5000 iterations were performed.
## There were 13 predictors of which 13 had non-zero influence.
summary(boston.boost)
```



```
##          var      rel.inf
## lstat    lstat 37.72235740
## rm       rm   28.25340805
## dis      dis   9.04378958
## crim     crim  6.95484787
## nox      nox   3.97210067
## black    black 3.16250916
## ptratio  ptratio 3.03202030
## age      age   2.35790500
## chas     chas  1.97366108
## tax      tax   1.67544858
## indus    indus 1.22537648
## rad      rad   0.57329299
## zn       zn    0.05328283
```

Predykcja na podstawie metody boosting

```
boston.pred4 <- predict(boston.boost, newdata = boston.test, n.trees = 5000)
rmse(boston.pred4, boston.test$medv)
```

```
## [1] 3.06509
```

*RMSE* jest w tym przypadku mniejsze niż w lasach losowych ale nieco większe niż w metodzie bagging. Wszystkie metody wzmacnianych drzew dają wyniki lepsze niż pojedyncze drzewa.

# Rozdział 6

## Klasyfikatory liniowe

Obszerną rodzinę klasyfikatorów stanowią modele liniowe (ang. *linear classification models*). Klasyfikacji w tej rodzinie technik dokonuje się na podstawie modeli funkcji kombinacji liniowej predyktorów. Jest to ujęcie parametryczne, w którym klasyfikacji nowej wartości dokonujemy na podstawie atrybutów obserwacji i wektora parametrów. Uczenie na podstawie zestawu treningowego polega na oszacowaniu parametrów modelu. W odróżnieniu od metod nieparametrycznych postać modelu tym razem jest znana. Każdy klasyfikator liniowy składa się z funkcji wewnętrznej (ang. *inner representation function*) i funkcji zewnętrznej (ang. *outer representation function*). Pierwsza jest funkcją rzeczywistą parametrów modelu i wartości atrybutów obserwacji

$$g(x) = F(\mathbf{a}(x), \mathbf{w}) = \sum_{i=0}^p w_i a_i(x) = \mathbf{w} \circ \mathbf{a}(x), \quad (6.1)$$

przyjmując, że  $a_0(x) = 1$ .

Funkcja zewnętrzna przyporządkowuje binarnie klasy na podstawie wartości funkcji wewnętrznej. Istnieją dwa główne typy tych klasyfikacji:

- brzegowa - przyjmujemy, że funkcje wewnętrzne tworzą granice zbiorów obserwacji różnych klas,
- probabilistyczna - bazująca na tym, że funkcje wewnętrzne mogą pośrednio wykazywać prawdopodobieństwo przynależności do danej klasy.

Pierwsza dzieli przestrzeń obserwacji za pomocą hiperpłaszczyzn na obszary jednorodne pod względem przynależności do klas. Druga jest próbą parametrycznej reprezentacji prawdopodobieństw przynależności do klas. Klasyfikacji na podstawie prawdopodobieństw można dokonać na różne sposoby, stosując:

- największe prawdopodobieństwo,
- funkcję najmniejszego kosztu błędnej klasyfikacji,
- krzywych ROC (ang. *Receiver Operating Characteristic* - o tym później).

Podejście brzegowe lub probabilistyczne prowadzi najczęściej do dwóch typów reprezentacji funkcji zewnętrznej:

- reprezentacji progowej (ang. *threshold representation*) - najczęściej przy podejściu brzegowym,
- reprezentacji logistycznej (ang. *logit representation*) - przy podejściu probabilistycznym.

### 6.1 Reprezentacja progowa

W przypadku klasyfikacji dwustanowej, dziedzina jest dzielona na dwa regiony (pozytywny i negatywny) poprzez porównanie funkcji zewnętrznej z wartością progową. Bez straty ogólności można sprawić, że będzie

to wartość 0

$$h(x) = H(g(x)) = \begin{cases} 1, & \text{jeśli } g(x) \geq 0 \\ 0, & \text{w przeciwnym przypadku.} \end{cases} \quad (6.2)$$

Czasami używa się parametryzacji  $\{-1, 1\}$ . Przez porównanie  $g(x)$  z 0 definiuje się hiperpłaszczyznę w  $p$  wymiarowej przestrzeni, która rozdziela dziedzinę na regiony pozytywne i negatywne. W tym ujęciu mówimy o liniowej separowalności obserwacji różnych klas, jeśli istnieje hiperpłaszczyzna je rozdziela.

## 6.2 Reprezentacja logitowa

Najbardziej popularną reprezentacją parametryczną stosowaną w klasyfikacji jest reprezentacja logitowa

$$\mathbb{P}(y = 1|x) = \frac{e^{g(x)}}{e^{g(x)} + 1}. \quad (6.3)$$

Wówczas  $g(x)$  nie reprezentuje bezpośrednio  $\mathbb{P}(y = 1|x)$  ale jego logit

$$g(x) = \text{logit}(\mathbb{P}(y = 1|x)), \quad (6.4)$$

gdzie  $\text{logit}(p) = \ln \frac{p}{1-p}$ . Dlatego właściwa postać reprezentacji jest następująca

$$\mathbb{P}(y = 1|x) = \text{logit}^{-1}(g(x)). \quad (6.5)$$

W ten sposób reprezentacja logitowa jest równoważna reprezentacji progowej, ponieważ

$$g(x) = \ln \frac{\mathbb{P}(y = 1|x)}{1 - \mathbb{P}(y = 1|x)} = \ln \frac{\mathbb{P}(y = 1|x)}{\mathbb{P}(y = 0|x)} > 0. \quad (6.6)$$

Jednak zaletą reprezentacji logitowej, w porównaniu do progowej, jest to, że można wyznaczyć prawdopodobieństwa przynależności do obu klas. W przypadku klasyfikacji wielostanowej uczymy tyle funkcji  $h$  ile jest klas.

## 6.3 Wady klasyfikatorów liniowych

- tylko w przypadku prostych funkcji wewnętrznych jesteśmy w stanie ocenić wpływ poszczególnych predyktorów na klasyfikację,
- jakość predykcji zależy od doboru funkcji wewnętrznej (liniowa w ścisłym sensie jest najczęściej niewystarczająca),
- nie jest w stanie klasyfikować poprawnie stanów (nie jest liniowo separowalna) w zagadnieniach typu XOR.

## Rozdział 7

# Regresja logistyczna

### 7.1 Model

Regresja logistyczna (ang. *logistic regression*) jest techniką z rodziny klasyfikatorów liniowych z reprezentacją logistyczną, a formalnie należy do rodziny uogólnionych modeli liniowych (GLM). Stosowana jest wówczas, gdy zmienna wynikowa posiada dwa stany (sukces i porażka), kodowane najczęściej za pomocą 1 i 0. W tej metodzie modelowane jest warunkowe prawdopodobieństwo sukcesu za pomocą kombinacji liniowej predyktorów  $X$ .

Ogólna postać modelu

$$Y \sim B(1, p) \quad (7.1)$$

$$p(X) = E(Y|X) = \frac{\exp(\beta X)}{1 + \exp(\beta X)}, \quad (7.2)$$

gdzie  $B(1, p)$  jest rozkładem dwumianowym o prawdopodobieństwie sukcesu  $p$ , a  $\beta X$  oznacza kombinację liniową parametrów modelu i wartości zmiennych niezależnych, przyjmując, że  $x_0 = 1$ . Jako funkcji łączącej (czyli opisującej związek między kombinacją liniową predyktorów i prawdopodobieństwem sukcesu) użyto *logitu*. Pozwala on na wygodną interpretację wyników w terminach szans.

Szansą (ang. *odds*) nazywamy stosunek prawdopodobieństwa sukcesu do prawdopodobieństwa porażki

$$o = \frac{p}{1 - p}. \quad (7.3)$$

Ponieważ będziemy przyjmowali, że  $p \in (0, 1)$ , to  $o \in (0, \infty)$ , a jej logarytm należy do przedziału  $(-\infty, \infty)$ .

Zatem logarytm szansy jest kombinacją liniową predyktorów

$$\log \left[ \frac{p(X)}{1 - p(X)} \right] = \beta_0 + \beta_1 x_1 + \dots + \beta_d x_d. \quad (7.4)$$

### 7.2 Estymacja parametrów modelu

Estymacji parametrów modelu logistycznego dokonujemy za pomocą metody największej wiarygodności. Funkcja wiarygodności w tym przypadku przyjmuje postać

$$L(X_1, \dots, X_n, \beta) = \prod_{i=1}^n p(X_i)_i^{Y_i} [1 - p(X_i)]^{1-Y_i}, \quad (7.5)$$

gdzie wektor  $\beta$  jest uwikłany w funkcji  $p(X_i)$ . Maksymalizacji dokonujemy raczej po nałożeniu na funkcję wiarygodności logarytmu, bo to ułatwia szukanie ekstremum.

$$\log L(X_1, \dots, X_n, \beta) = \sum_{i=1}^n (Y_i \log p(X_i) + (1 - Y_i) \log(1 - p(X_i))). \quad (7.6)$$

### 7.3 Interpretacja

Interpretacja (lat. *ceteris paribus* - "inne takie samo") poszczególnych parametrów modelu jest następująca:

- jeśli  $b_i > 0$  - to zmienna  $x_i$  ma wpływ stymulujący pojawienie się sukcesu,
- jeśli  $b_i < 0$  - to zmienna  $x_i$  ma wpływ ograniczający pojawienie się sukcesu,
- jeśli  $b_i = 0$  - to zmienna  $x_i$  nie ma wpływu na pojawienie się sukcesu.

Iloraz szans (ang. *odds ratio*) stosuje się w przypadku porównywania dwóch klas obserwacji. Jest on jak sama nazwa wskazuje ilorazem szans zajścia sukcesu w obu klasach

$$OR = \frac{p_1}{1 - p_1} \frac{1 - p_2}{p_2}, \quad (7.7)$$

gdzie  $p_i$  oznacza zajście sukcesu w  $i$ -tej klasie.

Interpretujemy go następująco:

- jeśli  $OR > 1$  - to w pierwszej grupie zajście sukcesu jest bardziej prawdopodobne,
- jeśli  $OR < 1$  - to w drugiej grupie zajście sukcesu jest bardziej prawdopodobne,
- jeśli  $OR = 1$  - to w obu grupach zajście sukcesu jest jednakowo prawdopodobne.

**Przykład 7.1.** Jako ilustrację działania regresji logistycznej użyjemy modelu dla danych ze zbioru `Default` pakietu `ISLR`.

```
library(ISLR)
head(Default)
```

```
##      default student  balance  income
## 1         No       No  729.5265 44361.625
## 2         No      Yes  817.1804 12106.135
## 3         No       No 1073.5492 31767.139
## 4         No       No  529.2506 35704.494
## 5         No       No  785.6559 38463.496
## 6         No      Yes  919.5885  7491.559
```

Zmienną zależną jest `default`, a pozostałe są predyktorami. najpierw dokonamy podziału próby na uczącą i testową, a następnie zbudujemy model.

```
set.seed(2019)
ind <- sample(1:nrow(Default), size = 2/3*nrow(Default))
dt.ucz <- Default[ind,]
dt.test <- Default[-ind,]
mod.logit <- glm(default~., dt.ucz, family = binomial("logit"))
summary(mod.logit)
```

```
##
## Call:
## glm(formula = default ~ ., family = binomial("logit"), data = dt.ucz)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.4481  -0.1470  -0.0597  -0.0226   3.6966
```

```
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.085e+01  5.896e-01 -18.409  <2e-16 ***
## studentYes  -4.970e-01  2.851e-01  -1.744  0.0812 .
## balance      5.604e-03  2.809e-04  19.949  <2e-16 ***
## income       7.933e-06  9.652e-06   0.822  0.4112
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1906.5  on 6665  degrees of freedom
## Residual deviance: 1059.8  on 6662  degrees of freedom
## AIC: 1067.8
##
## Number of Fisher Scoring iterations: 8
```

Tylko `income` nie ma żadnego wpływu na prawdopodobieństwo stanu `Yes` zmiennej `default`. Zmienna `balance` wpływa stymulująco na prawdopodobieństwo pojawienia się sukcesu. Natomiast jeśli badana osoba jest studentem (`studentYes`), to ma wpływ ograniczający na pojawienie się sukcesu. Chcąc porównać dwie grupy obserwacji, przykładowo studentów z nie-studentami, możemy wykorzystać iloraz szans.

```
exp(cbind(OR = coef(mod.logit), confint(mod.logit))) %>%
  kable(digits = 4)
```

	OR	2.5 %	97.5 %
(Intercept)	0.0000	0.0000	0.0001
studentYes	0.6083	0.3485	1.0668
balance	1.0056	1.0051	1.0062
income	1.0000	1.0000	1.0000

Z powyższej tabeli wynika, że bycie studentem zmniejsza szanse na `Yes` w zmiennej `default` o około 40% (w stosunku do nie-studentów). Natomiast wzrost zmiennej `balance` przy zachowaniu pozostałych zmiennych na tym samym poziomie skutkuje wzrostem szans na `Yes` o około 0.6%.

Chcąc przeprowadzić predykcję na podstawie modelu dla ustalonych wartości cech (np. `student = Yes`, `balance = $1000` i `income = $40000`) postępujemy następująco

```
dt.new <- data.frame(student = "Yes", balance = 1000, income = 40000)
predict(mod.logit, newdata = dt.new, type = "response")
```

```
##              1
## 0.004367692
```

Otrzymany wynik jest oszacowanym prawdopodobieństwem warunkowym wystąpienia sukcesu (`default = Yes`). Widać zatem, że poziomy badanych cech sprzyja raczej porażce.

Jeśli chcemy sprawdzić jakość klasyfikacji na zbiorze testowym, to musimy ustalić na jakim poziomie prawdopodobieństwa będziemy uznawać obserwację za sukces. W zależności od tego, na predykcji jakiego stanu zależy nam bardziej, możemy różnie dobierać ten próg (bez żadnych dodatkowych przesłanek najczęściej jest to 0.5).

```
pred <- predict(mod.logit, newdata = dt.test, type = "response")
pred.class <- ifelse(pred > 0.5, "Yes", "No")
(tab <- table(pred.class, dt.test$default))
```

```
##  
## pred.class   No   Yes  
##           No 3204   76  
##           Yes  13   41  
  
(acc <- sum(diag(prop.table(tab))))
```

```
## [1] 0.9733053
```

Klasyfikacja na poziomie 97% wskazuje na dobre dopasowanie modelu.



## Rozdział 8

# Analiza dyskryminacyjna

Analiza dyskryminacyjna (ang. *discriminant analysis*) jest grupą technik dyskryminacji obserwacji względem przynależności do klas. Część z nich należy do klasyfikatorów liniowych (choć nie zawsze w ścisłym sensie). Za autorów tej metody uważa się Fisher'a (1936) i Welch'a (1939). Każdy z nich prezentował nieco inne podejście do tematu klasyfikacji. Welch poszukiwał klasyfikacji minimalizującej prawdopodobieństwo błędnej klasyfikacji, znane jako klasyfikatory bayesowskie. Podejście Fisher'a skupiało się raczej na porównaniu zmienności międzygrupowej do zmienności wewnątrzgrupowej. Wychodząc z założenia, że iloraz tych wariancji powinien być stosunkowo duży przy różnych klasach, jeśli do ich opisu użyjemy odpowiednich zmiennych niezależnych. W istocie chodzi o znalezienie takiego wektora, w kierunku którego wspomniany iloraz wariancji jest największy.

### 8.1 Liniowa analiza dyskryminacyjna Fisher'a

#### 8.1.1 Dwie kategorie zmiennej grupującej

Niech  $\mathbf{D}$  będzie zbiorem zawierającym  $n$  punktów  $\{\mathbf{x}_i, y_i\}$ , gdzie  $\mathbf{x}_i \in \mathbb{R}^d$ , a  $y_i \in \{c_1, \dots, c_k\}$ . Niech  $\mathbf{D}_i$  oznacza podzbiór punktów zbioru  $\mathbf{D}$ , które należą do klasy  $c_i$ , czyli  $\mathbf{D}_i = \{\mathbf{x}_i | y_i = c_i\}$  i niech  $|\mathbf{D}_i| = n_i$ . Na początek założmy, że  $\mathbf{D}$  składa się tylko z  $\mathbf{D}_1$  i  $\mathbf{D}_2$ .

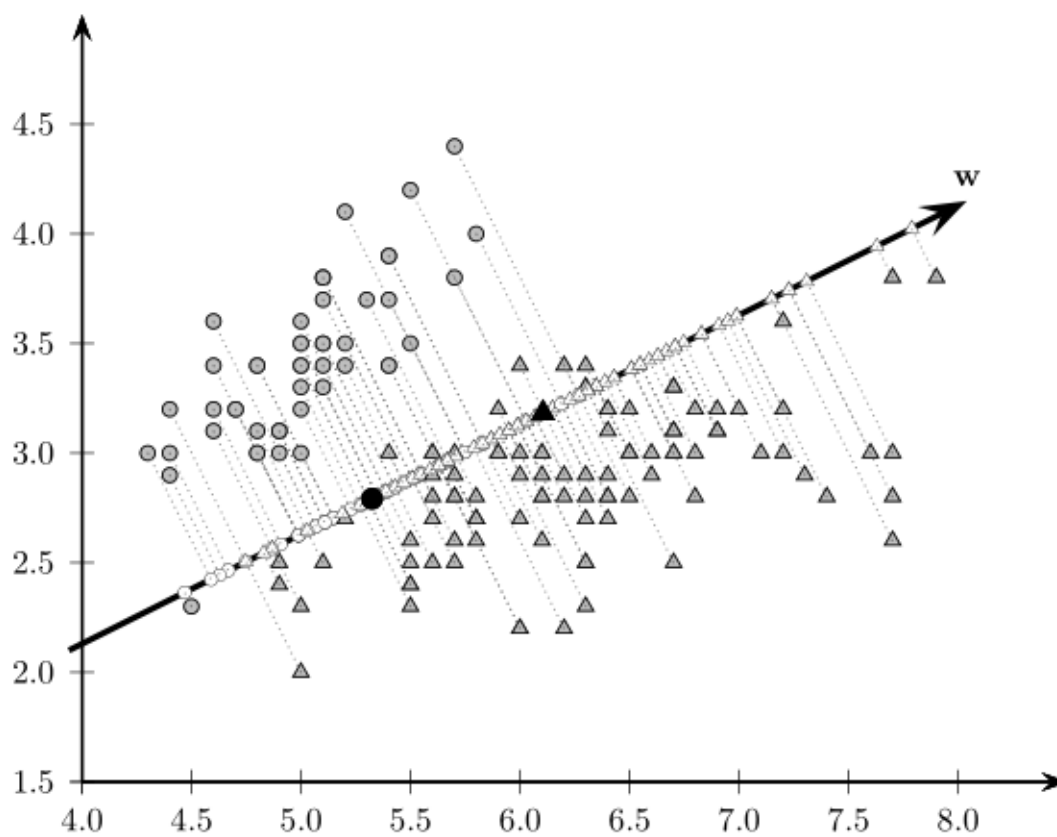
Niech  $\mathbf{w}$  będzie wektorem jednostkowym ( $\mathbf{w}'\mathbf{w} = 1$ ), wówczas rzut ortogonalny punktu  $\mathbf{x}_i$  na wektor  $\mathbf{w}$  można zapisać następująco

$$\tilde{\mathbf{x}}_i = \left( \frac{\mathbf{w}'\mathbf{x}_i}{\mathbf{w}'\mathbf{w}} \right) \mathbf{w} = (\mathbf{w}'\mathbf{x}_i) \mathbf{w} = a_i \mathbf{w}, \quad (8.1)$$

gdzie  $a_i$  jest współrzędną punktu  $\tilde{\mathbf{x}}_i$  w kierunku wektora  $\mathbf{w}$ , czyli

$$a_i = \mathbf{w}'\mathbf{x}_i. \quad (8.2)$$

Zatem  $(a_1, \dots, a_n)$  reprezentują odwzorowanie  $\mathbb{R}^d$  w  $\mathbb{R}$ , czyli z  $d$ -wymiarowej przestrzeni w przestrzeń generowaną przez  $\mathbf{w}$ .

Rysunek 8.1: Rzut ortogonalny punktów w kierunku wektora  $w$

Każdy punkt należy do pewnej klasy, dlatego możemy wyliczyć

$$m_1 = \frac{1}{n_1} \sum_{\mathbf{x}_i \in \mathcal{D}_1} a_i = \quad (8.3)$$

$$= \frac{1}{n_1} \sum_{\mathbf{x}_i \in \mathcal{D}_1} \mathbf{w}' \mathbf{x}_i = \quad (8.4)$$

$$= \mathbf{w}' \left( \frac{1}{n_1} \sum_{\mathbf{x}_i \in \mathcal{D}_1} \mathbf{x}_i \right) = \quad (8.5)$$

$$= \mathbf{w}' \boldsymbol{\mu}_1, \quad (8.6)$$

gdzie  $\boldsymbol{\mu}_1$  jest wektorem średnich punktów z  $\mathcal{D}_1$ . W podobny sposób można policzyć  $m_2 = \mathbf{w}' \boldsymbol{\mu}_2$ . Oznacza to, że średnia projekcji jest projekcją średnich.

Rozsądnym wydaje się teraz poszukać takiego wektora, aby  $|m_1 - m_2|$  była maksymalnie duża przy zachowaniu niezbyt dużej zmienności wewnątrz grup. Dlatego kryterium Fisher'a przyjmuje postać

$$\max_{\mathbf{w}} J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{ss_1^2 + ss_2^2}, \quad (8.7)$$

gdzie  $ss_j^2 = \sum_{\mathbf{x}_i \in \mathcal{D}_j} (a_i - m_j)^2 = n_j \sigma_j^2$ .

Zauważmy, że licznik w (8.7) da się zapisać jako

$$(m_1 - m_2)^2 = (\mathbf{w}'(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2))^2 = \quad (8.8)$$

$$= \mathbf{w}'((\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)')\mathbf{w} = \quad (8.9)$$

$$= \mathbf{w}' \mathbf{B} \mathbf{w} \quad (8.10)$$

gdzie  $\mathbf{B} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)'$  jest macierzą  $d \times d$ .

Ponadto

$$ss_j^2 = \sum_{\mathbf{x}_i \in \mathcal{D}_j} (a_i - m_j)^2 = \quad (8.11)$$

$$= \sum_{\mathbf{x}_i \in \mathcal{D}_j} (\mathbf{w}' \mathbf{x}_i - \mathbf{w}' \boldsymbol{\mu}_j)^2 = \quad (8.12)$$

$$= \sum_{\mathbf{x}_i \in \mathcal{D}_j} (\mathbf{w}'(\mathbf{x}_i - \boldsymbol{\mu}_j))^2 = \quad (8.13)$$

$$= \mathbf{w}' \left( \sum_{\mathbf{x}_i \in \mathcal{D}_j} (\mathbf{x}_i - \boldsymbol{\mu}_j)(\mathbf{x}_i - \boldsymbol{\mu}_j)' \right) \mathbf{w} = \quad (8.14)$$

$$= \mathbf{w}' \mathbf{S}_j \mathbf{w}, \quad (8.15)$$

gdzie  $\mathbf{S}_j = n_j \boldsymbol{\Sigma}_j$ . Zatem mianownik (8.7) możemy zapisać jako

$$ss_1^2 + ss_2^2 = \mathbf{w}'(\mathbf{S}_1 + \mathbf{S}_2)\mathbf{w} = \mathbf{w}' \mathbf{S} \mathbf{w}, \quad (8.16)$$

gdzie  $\mathbf{S} = \mathbf{S}_1 + \mathbf{S}_2$ . Ostatecznie warunek Fisher'a przyjmuje postać

$$\max_{\mathbf{w}} J(\mathbf{w}) = \frac{\mathbf{w}' \mathbf{B} \mathbf{w}}{\mathbf{w}' \mathbf{S} \mathbf{w}}. \quad (8.17)$$

Różniczkując (8.17) po  $\mathbf{w}$  otrzymamy warunek

$$\mathbf{B} \mathbf{w} = \lambda \mathbf{S} \mathbf{w}, \quad (8.18)$$

gdzie  $\lambda = J(\mathbf{w})$ . Maksimum (8.18) jest osiągnięte dla wektora  $\mathbf{w}$  równego wektorowi własnemu odpowiadającemu największej wartości własnej równania charakterystycznego  $|\mathbf{B} - \lambda\mathbf{S}| = 0$ . Jeśli  $\mathbf{S}$  nie jest osobiłwa, to rozwiązanie (8.18) otrzymujemy przez znalezienie największej wartości własnej macierzy  $\mathbf{B}\mathbf{S}^{-1}$  lub bez wykorzystania wartości i wektorów własnych.

Ponieważ  $\mathbf{B} = ((\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)') \mathbf{w}$  jest macierzą  $d \times d$  rzędu 1, to  $\mathbf{B}\mathbf{w}$  jest punktem na kierunku wyznaczonym przez wektor  $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$ , bo

$$\mathbf{B}\mathbf{w} = ((\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)') \mathbf{w} = \quad (8.19)$$

$$= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) ((\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)') \mathbf{w} = \quad (8.20)$$

$$= b(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2), \quad (8.21)$$

gdzie  $b = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \mathbf{w}$  jest skłarem.

Wówczas (8.18) zapiszemy jako

$$b(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) = \lambda \mathbf{S} \mathbf{w} \quad (8.22)$$

$$\mathbf{w} = \frac{b}{\lambda} \mathbf{S}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \quad (8.23)$$

A ponieważ  $b/\lambda$  jest liczbą, to kierunek najlepszej dyskryminacji grup wyznacza wektor

$$\mathbf{w} = \mathbf{S}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2). \quad (8.24)$$

### 8.1.2 $k$ -kategorii zmiennej grupującej

Uogólnieniem tej teorii na przypadek  $k$  klas otrzymujemy przez uwzględnienie  $k - 1$  funkcji dyskryminacyjnych. Zmienność wewnątrzgrupowa przyjmuje wówczas postać

$$\mathbf{S}_W = \sum_{i=1}^k \mathbf{S}_i, \quad (8.25)$$

gdzie  $\mathbf{S}_i$  jest zdefiniowane jak w (8.15). Niech średnia i rozrzut globalny będą dane wzorami

$$\mathbf{m} = \frac{1}{n} \sum_{i=1}^k n_i \mathbf{m}_i, \quad (8.26)$$

$$\mathbf{S}_T = \sum_{j=1}^k \sum_{\mathbf{x} \in D_j} (\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})' \quad (8.27)$$

gdzie  $\mathbf{m}_i$  jest określone jak w (8.6). Wtedy zmienność międzygrupową możemy wyrazić jako

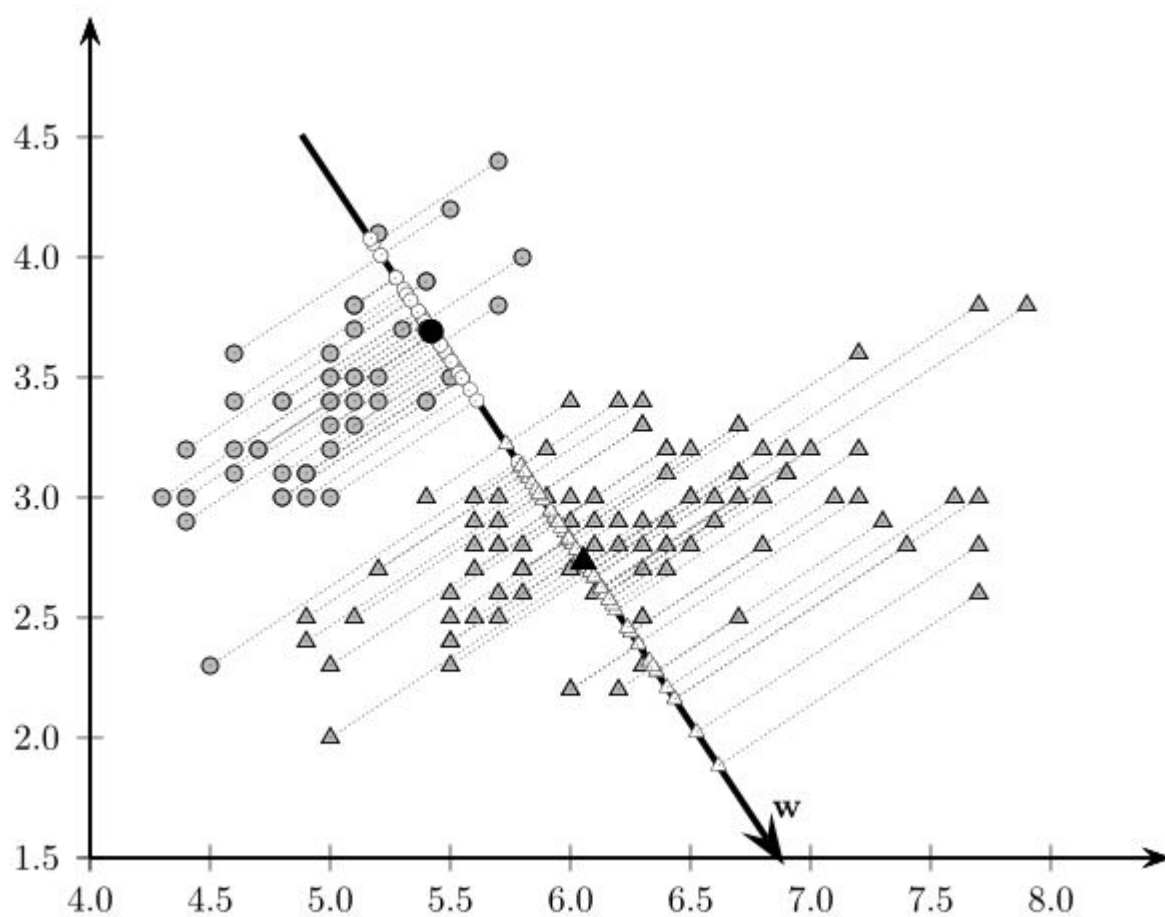
$$\mathbf{S}_B = \sum_{i=1}^k n_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})', \quad (8.28)$$

bo  $\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B$ . Określamy projekcję  $d$ -wymiarowej przestrzeni na  $k - 1$ -wymiarową przestrzeń za pomocą  $k - 1$  funkcji dyskryminacyjnych postaci

$$\mathbf{a}_j = \mathbf{w}_j' \mathbf{x}, \quad j = 1, \dots, k - 1. \quad (8.29)$$

Połączone wszystkie  $k - 1$  rzutów możemy zapisać jako

$$\mathbf{a} = \mathbf{W}' \mathbf{x}. \quad (8.30)$$



Rysunek 8.2: Rzut ortogonalny w kierunku wektora  $w$ , będącego najlepiej dyskryminującym obie grupy obserwacji

W nowej przestrzeni  $k - 1$ -wymiarowej możemy zdefiniować

$$\tilde{\mathbf{m}} = \frac{1}{n} \sum_{i=1}^k n_i \tilde{\mathbf{m}}_i, \quad (8.31)$$

gdzie  $\tilde{\mathbf{m}}_i = \frac{1}{n_i} \sum_{\mathbf{a} \in A_i} \mathbf{a}$ , a  $A_i$  jest projekcją obiektów z  $i$ -tej klasy w kierunku wektora  $\mathbf{W}$ . Dalej możemy zdefiniować zmienności między- i wewnątrzgrupowe dla obiektów przekształconych przez  $\mathbf{W}$

$$\tilde{S}_W = \sum_{i=1}^k \sum_{\mathbf{a} \in A_i} (\mathbf{a} - \tilde{\mathbf{m}})(\mathbf{a} - \tilde{\mathbf{m}})' \quad (8.32)$$

$$\tilde{S}_B = \sum_{i=1}^k n_i (\tilde{\mathbf{m}}_i - \tilde{\mathbf{m}})(\tilde{\mathbf{m}}_i - \tilde{\mathbf{m}})'. \quad (8.33)$$

Łatwo można zatem pokazać, że

$$\tilde{S}_W = \mathbf{W}' \mathbf{S}_W \mathbf{W} \quad (8.34)$$

$$\tilde{S}_B = \mathbf{W}' \mathbf{S}_B \mathbf{W}. \quad (8.35)$$

Ostatecznie warunek (8.7) w  $k$ -wymiarowym ujęciu można przedstawić jako

$$\max_{\mathbf{W}} J(\mathbf{W}) = \frac{\tilde{S}_W}{\tilde{S}_B} = \frac{\mathbf{W}' \mathbf{S}_W \mathbf{W}}{\mathbf{W}' \mathbf{S}_B \mathbf{W}}. \quad (8.36)$$

Maksimum można znaleźć poprzez rozwiązanie równania charakterystycznego

$$|\mathbf{S}_B - \lambda_i \mathbf{S}_W| = 0 \quad (8.37)$$

dla każdego  $i$ .

**Przykład 8.1.** Dla danych ze zbioru `iris` przeprowadzimy analizę dyskryminacji. Implementację metody LDA znajdziemy w pakiecie `MASS` w postaci funkcji `lda`.

Zaczynamy od standaryzacji zmiennych i podziału próby na uczącą i testową.

```
library(MASS)
library(tidyverse)
iris.std <- iris %>%
  mutate_if(is.numeric, scale)
set.seed(2019)
ind <- sample(nrow(iris.std), size = 100)
dt.ucz <- iris.std[ind,]
dt.test <- iris.std[-ind,]
```

Budowa modelu

```
mod.lda <- lda(Species~., data = dt.ucz)
mod.lda$prior
```

```
##      setosa versicolor virginica
##      0.36      0.31      0.33
```

Prawdopodobieństwa *a priori* przynależności do klas przyjęto na podstawie próby uczącej.

```
mod.lda$means

##      Sepal.Length Sepal.Width Petal.Length Petal.Width
## setosa      -1.0251463  0.8690229   -1.294839  -1.2527443
## versicolor   0.1619267 -0.5015842    0.316167   0.1786195
## virginica    0.9174351 -0.2636338    1.046883   1.0504160
```

W części `means` wyświetlone są średnie poszczególnych zmiennych niezależnych w podziale na grupy. Dzięki temu można określić położenia środków ciężkości poszczególnych klas w oryginalnej przestrzeni.

```
mod.lda$scaling
```

```
##                LD1        LD2
## Sepal.Length  1.0073378  0.211252
## Sepal.Width   0.4701094 -1.053135
## Petal.Length -4.0746585  1.488372
## Petal.Width  -2.5146178 -2.312201
```

Powyższa tabela zawiera współrzędne wektorów wyznaczających funkcje dyskryminacyjne. Na ich podstawie możemy określić, która z nich wpływa najmocniej na tworzenie się nowej przestrzeni.

Obiekt `svd` przechowuje pierwiastki z  $\lambda_i$ , dlatego podnosząc je do kwadratu i dzieląc przez ich sumę otrzymamy udział poszczególnych zmiennych w dyskryminacji przypadków. Jak widać pierwsza funkcja dyskryminacyjna w zupełności by wystarczyła.

```
mod.lda$svd^2/sum(mod.lda$svd^2)
```

```
## [1] 0.994875091 0.005124909
```

Klasyfikacja na podstawie modelu

```
pred.lda <- predict(mod.lda, dt.test)
```

Wynik predykcji przechowuje trzy rodzaje obiektów:

- klasy, które przypisał obiektom model (`class`);
- prawdopodobieństwa *a posteriori* przynależności do klas na podstawie modelu (`posterior`);
- współrzędne w nowej przestrzeni LD1, LD2 (`x`).

Sprawdzenie jakości klasyfikacji

```
tab <- table(pred = pred.lda$class, obs = dt.test$Species)
tab
```

```
##          obs
## pred      setosa versicolor virginica
## setosa      14         0         0
## versicolor   0        18         0
## virginica    0         1        17
```

```
sum(diag(prop.table(tab)))
```

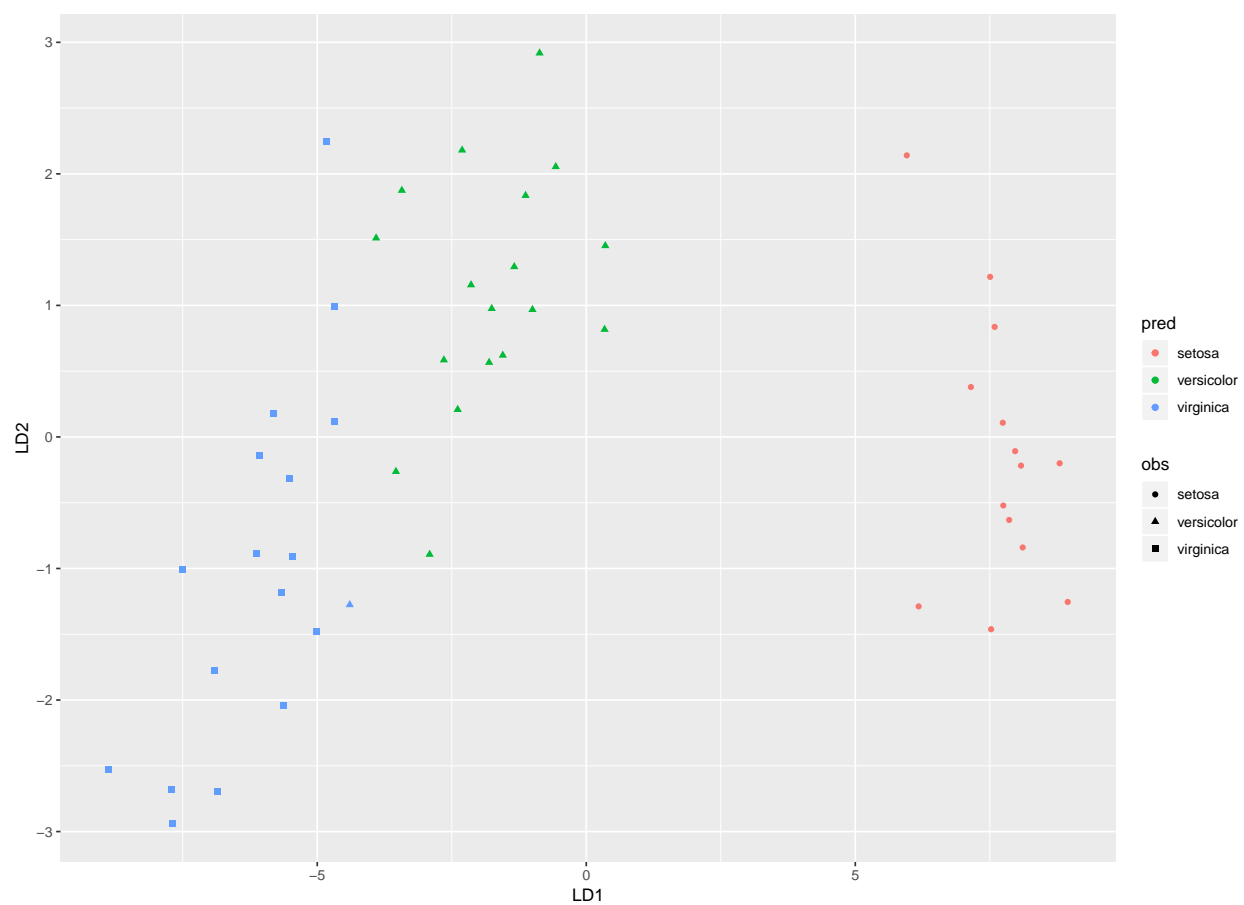
```
## [1] 0.98
```

Jak widać z powyższej tabeli model dobrze sobie radzi z klasyfikacją obiektów. Odsetek poprawnych klasyfikacji wynosi 98%.

```
cbind.data.frame(obs = dt.test$Species,
                  pred.lda$x,
                  pred = pred.lda$class) %>%
  ggplot(aes(x = LD1, y = LD2)) +
  geom_point(aes(color = pred, shape = obs))
```

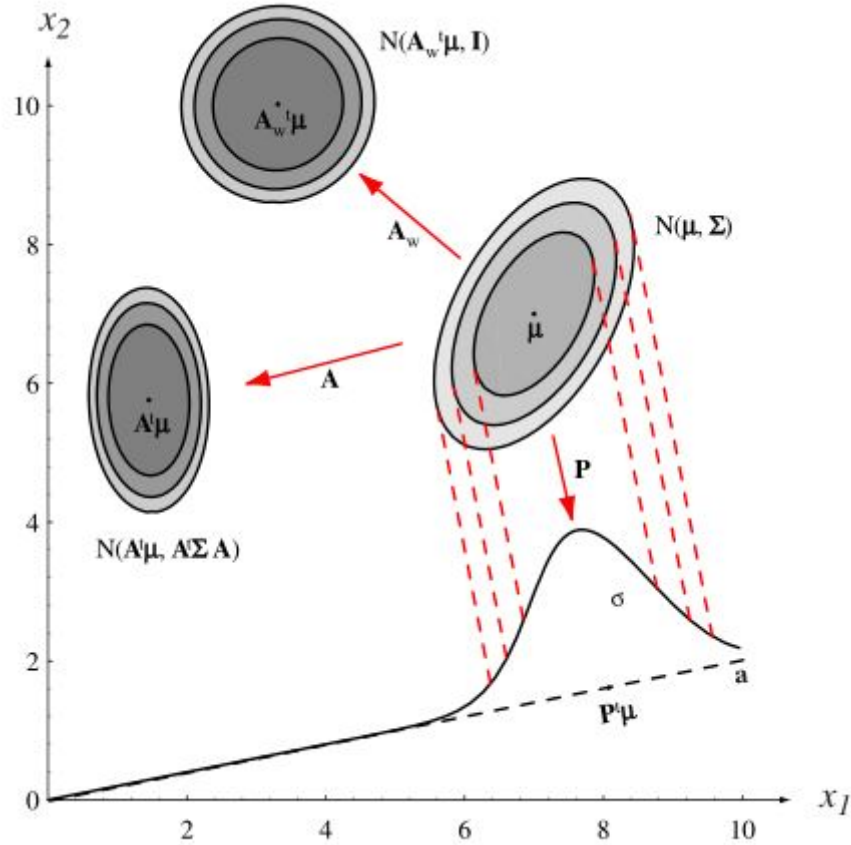
## 8.2 Liniowa analiza dyskryminacyjna - podejście probabilistyczne

Jak wspomniano na wstępie (patrz rozdział 8), podejście prezentowane przez Welch'a polegało na minimalizacji prawdopodobieństwa popełnienia błędu przy klasyfikacji. Cała rodzina klasyfikatorów Bayesa (patrz



Rysunek 8.3: Klasyfikacja w przestrzeni LD1, LD2 na podstawie modelu `mod.lda`





Rysunek 8.4: Transformacje rozkładu normalnego łącznego. Źródło: @duda2001

rozdział 9.1) polega na wyznaczeniu prawdopodobieństw *a posteriori*, na podstawie których dokonuje się decyzji o klasyfikacji obiektów. Tym razem dodajemy również założenie, że zmienne niezależne  $\mathbf{x} = (x_1, \dots, x_d)$  charakteryzują się wielowymiarowym rozkładem normalnym

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})' \Sigma (\mathbf{x} - \boldsymbol{\mu}) \right], \quad (8.38)$$

gdzie  $\boldsymbol{\mu}$  jest wektorem średnich  $\mathbf{x}$ , a  $\Sigma$  jest macierzą kowariancji  $\mathbf{x}$ .

*Uwaga.* Liniowa kombinacja zmiennych losowych o normalnym rozkładzie łącznym ma również rozkład łączny normalny. W szczególności, jeśli  $A$  jest macierzą wymiaru  $d \times k$  i  $\mathbf{y} = A'\mathbf{x}$ , to  $f(\mathbf{y}) \sim N(A'\boldsymbol{\mu}, A'\Sigma A)$ . Odpowiednia forma macierzy przekształcenia  $A_w$ , sprawia, że zmienne po transformacji charakteryzują się rozkładem normalnym łącznym o wariancji określonej przez  $I$ . Jeśli  $\Phi$  jest macierzą, której kolumny są ortogonalnymi wektorami własnymi macierzy  $\Sigma$ , a  $\Lambda$  macierzą diagonalną wartości własnych, to transformacja  $A_w = \Phi \Lambda^{-1}$  przekształca  $\mathbf{x}$  w  $\mathbf{y} \sim N(A'_w \boldsymbol{\mu}, I)$ .

**Definicja 8.1.** Niech  $g_i(\mathbf{x})$ ,  $i = 1, \dots, k$  będzie pewną funkcją dyskryminacyjną, wówczas obiekt  $\mathbf{x}$  należy zaklasyfikować do grupy  $c_i$  jeśli spełniony jest warunek

$$g_i(\mathbf{x}) > g_j(\mathbf{x}), \quad j \neq i. \quad (8.39)$$

W podejściu polegającym na minimalizacji prawdopodobieństwa błędnej klasyfikacji, przyjmuje się najczęściej, że

$$g_i(\mathbf{x}) = \mathbb{P}(c_i | \mathbf{x}), \quad (8.40)$$

czyli jako prawdopodobieństwo a posteriori. Wszystkie trzy poniższe postaci funkcji dyskryminacyjnych są dopuszczalne i równoważne ze względu na rezultat grupowania

$$g_i(\mathbf{x}) = \mathbb{P}(c_i|\mathbf{x}) = \frac{\mathbb{P}(\mathbf{x}|c_i)\mathbb{P}(c_i)}{\sum_{i=1}^k \mathbb{P}(\mathbf{x}|c_i)\mathbb{P}(c_i)}, \quad (8.41)$$

$$g_i(\mathbf{x}) = \mathbb{P}(\mathbf{x}|c_i)\mathbb{P}(c_i), \quad (8.42)$$

$$g_i(\mathbf{x}) = \log \mathbb{P}(\mathbf{x}|c_i) + \log \mathbb{P}(c_i) \quad (8.43)$$

W przypadku gdy  $\mathbf{x}|c_i \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ , to na podstawie (8.38)  $g_i$  danej jako (8.43) przyjmuje postać

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)' \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) - \frac{d}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| + \log \mathbb{P}(c_i). \quad (8.44)$$

W kolejnych podrozdziałach przeanalizujemy trzy możliwe formy macierzy kowariancji.

### 8.2.1 Przypadek gdy $\boldsymbol{\Sigma}_i = \sigma^2 I$

To najprostszy przypadek, zakładający niezależność zmiennych wchodzących w skład  $\mathbf{x}$ , których wariancje są stałe  $\sigma^2$ . Wówczas  $g_i$  przyjmuje postać

$$g_i(\mathbf{x}) = -\frac{\|\mathbf{x} - \boldsymbol{\mu}_i\|^2}{2\sigma^2} + \log \mathbb{P}(c_i), \quad (8.45)$$

gdzie  $\|\cdot\|$  jest normą euklidesową.

Rozpisując licznik równania (8.45) mamy

$$\|\mathbf{x} - \boldsymbol{\mu}_i\|^2 = (\mathbf{x} - \boldsymbol{\mu}_i)'(\mathbf{x} - \boldsymbol{\mu}_i). \quad (8.46)$$

Zatem

$$g_i(\mathbf{x}) = -\frac{1}{2\sigma^2}[\mathbf{x}'\mathbf{x} - 2\boldsymbol{\mu}_i'\mathbf{x} + \boldsymbol{\mu}_i'\boldsymbol{\mu}_i] + \log \mathbb{P}(c_i). \quad (8.47)$$

A ponieważ  $\mathbf{x}'\mathbf{x}$  nie zależy do  $i$ , to funkcję dyskryminacyjną możemy zapisać jako

$$g_i(\mathbf{x}) = \mathbf{w}_i'\mathbf{x} + w_{i0}, \quad (8.48)$$

gdzie  $\mathbf{w}_i = \frac{1}{\sigma^2}\boldsymbol{\mu}_i$ , a  $w_{i0} = \frac{-1}{2\sigma^2}\boldsymbol{\mu}_i'\boldsymbol{\mu}_i + \log \mathbb{P}(c_i)$ .

Na podstawie funkcji dyskryminacyjnych wyznaczamy hiperpłaszczyzny decyzyjne jako zbiory punktów dla których  $g_i(\mathbf{x}) = g_j(\mathbf{x})$ , gdzie  $g_i, g_j$  są największe. Możemy to zapisać w następujący sposób

$$\mathbf{w}'(\mathbf{x} - \mathbf{x}_0) = 0, \quad (8.49)$$

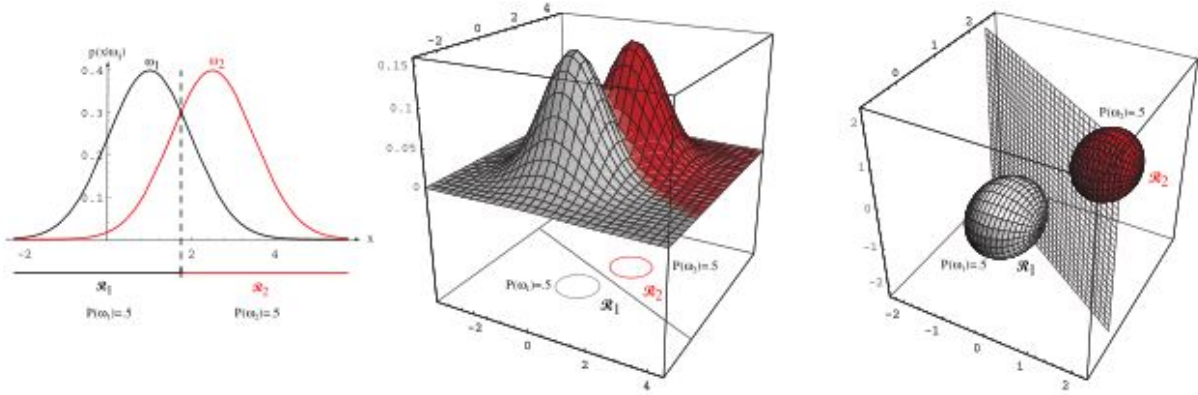
gdzie

$$\mathbf{w} = \boldsymbol{\mu}_i - \boldsymbol{\mu}_j, \quad (8.50)$$

oraz

$$\mathbf{x}_0 = \frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \frac{\sigma^2}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2} \log \frac{\mathbb{P}(c_i)}{\mathbb{P}(c_j)} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j). \quad (8.51)$$

Równanie (8.49) określa hiperpłaszczyznę przechodzącą przez  $\mathbf{x}_0$  i prostopadłą do  $\mathbf{w}$ .



Rysunek 8.5: Dyskrymiancja hiperpłaszczyznami w sygnacji dwóch klas. Wykres po lewej, to ujęcie jedno-wymiarowe, wykresy po środku - ujęcie 2-wymiarowe i wykresy po prawej, to ujęcie 3-wymiarowe. Źródło: @duda2001

### 8.2.2 Przypadek gdy $\Sigma_i = \Sigma$

Przypadek ten opisuje sytuację, gdy rozkłady  $\mathbf{x}$  są identyczne we wszystkich grupach ale zmienne w ich skład wchodzące nie są niezależne. W tym przypadku funkcje dyskryminacyjne przyjmują postać

$$g_i(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \log \mathbb{P}(c_i). \quad (8.52)$$

Jeśli  $\mathbb{P}(c_i)$  są identyczne dla wszystkich klas, to można je pominąć we wzorze (8.52). Metryka euklidesowa ze wzoru (8.45) została zastąpiona przez odległość Mahalanobis'a. Podobnie jak w przypadku gdy  $\boldsymbol{\Sigma}_i = \sigma^2 I$ , tak i teraz można uprościć (8.52) przez rozpisanie formy kwadratowej, aby otrzymać, że

$$g_i(\mathbf{x}) = \mathbf{w}'_i \mathbf{x} + w_{i0}, \quad (8.53)$$

gdzie  $\mathbf{w}_i = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i$ , a  $w_{i0} = -\frac{1}{2} \boldsymbol{\mu}'_i \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \log \mathbb{P}(c_i)$ .

Ponieważ funkcje dyskryminacyjne są liniowe, to hiperpłaszczyzny są ograniczeniami obszarów obserwacji każdej z klas

$$\mathbf{w}'(\mathbf{x} - \mathbf{x}_0) = 0, \quad (8.54)$$

gdzie

$$\mathbf{w} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j), \quad (8.55)$$

oraz

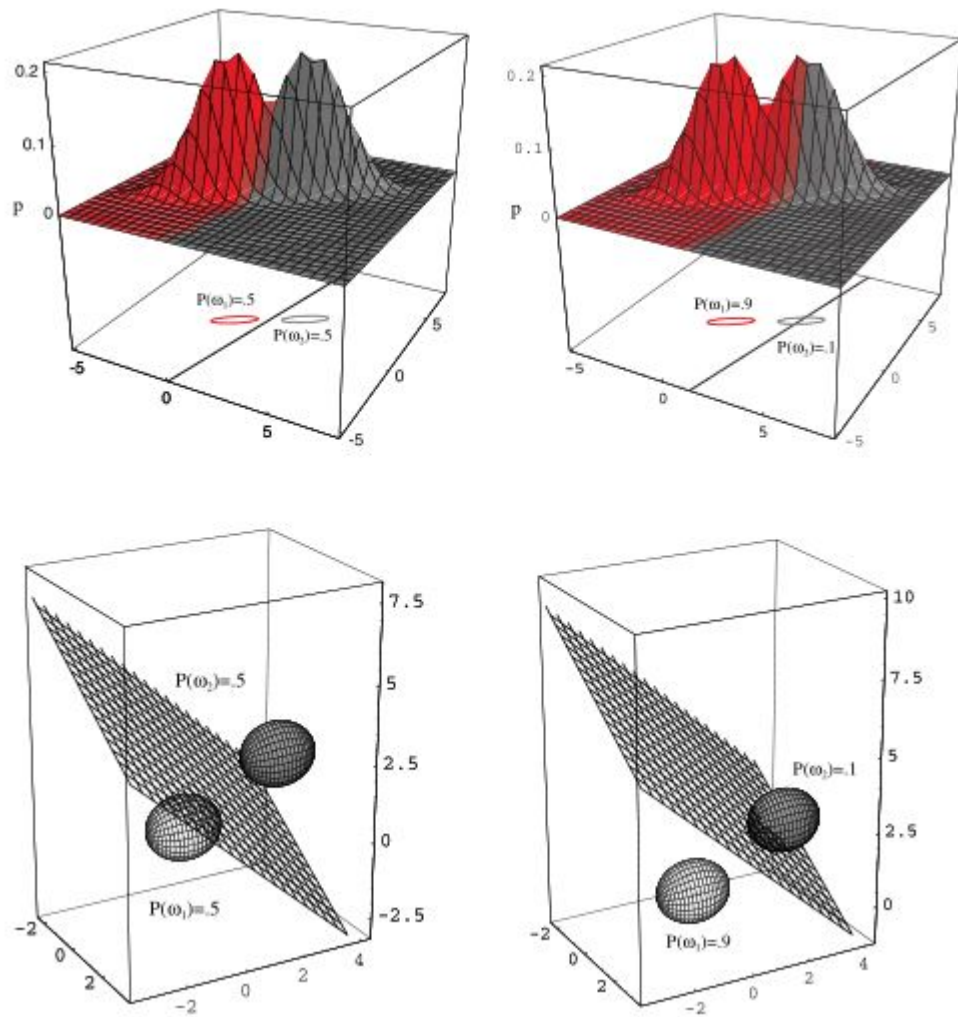
$$\mathbf{x}_0 = \frac{1}{2}(\boldsymbol{\mu}_i + \boldsymbol{\mu}_j) - \frac{\log[\mathbb{P}(c_i)/\mathbb{P}(c_j)]}{(\mathbf{x} - \boldsymbol{\mu}_i)' \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_i)}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j). \quad (8.56)$$

Tym razem  $\mathbf{w} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)$  nie jest wektorem w kierunku  $\boldsymbol{\mu}_i - \boldsymbol{\mu}_j$ , więc hiperpłaszczyzna rozdzielająca obiekty różnych klas nie jest prostopadła do wektora  $\boldsymbol{\mu}_i - \boldsymbol{\mu}_j$  ale przecina go w połowie (w punkcie  $\mathbf{x}_0$ ).

### 8.2.3 Przypadek gdy $\Sigma_i$ jest dowolnej postaci

Jest to najbardziej ogólny przypadek, kiedy nie nakłada się żadnych ograniczeń na macierze kowariancji grupowych. Postać funkcji dyskryminacyjnych jest następująca

$$g_i(\mathbf{x}) = \mathbf{x}' \mathbf{W}_i \mathbf{x} + \mathbf{w}'_i \mathbf{x} + w_{i0} \quad (8.57)$$



Rysunek 8.6: Hiperpłaszczyzna rozdzielająca obszary innych klas może być przesunięta w kierunku bardziej prawdopodobnej klasy, jeśli prawdopodobieństwa a priori są różne. Źródło: @duda2001

gdzie

$$\mathbf{W}_i = -\frac{1}{2}\boldsymbol{\Sigma}_i^{-1}, \quad (8.58)$$

$$\mathbf{w}_i = \boldsymbol{\Sigma}_i^{-1}\boldsymbol{\mu}_i, \quad (8.59)$$

$$w_{i0} = -\frac{1}{2}\boldsymbol{\mu}_i'\boldsymbol{\Sigma}_i^{-1}\boldsymbol{\mu}_i - \frac{1}{2}\log|\boldsymbol{\Sigma}_i| + \log\mathbb{P}(c_i). \quad (8.60)$$

Ograniczenia w ten sposób budowane są hiperpowierzchniami (nie koniecznie hiperpłaszczyznami). W literaturze ta metoda znana jest pod nazwą kwadratowej analizy dyskryminacyjnej (ang. *Quadratic Discriminant Analysis*).

**Przykład 8.2.** Przeprowadzimy klasyfikację na podstawie zbioru `Smarket` pakietu `ILSR`. Dane zawierają kursy indeksu giełdowego S&P500 w latach 2001-2005. Na podstawie wartości waloru z poprzednich 2 dni będziemy chcieli przewidzieć czy ruch w kolejnym okresie czasu będzie w górę czy w dół.

```
library(ILSR)
head(Smarket)
```

```
##   Year  Lag1  Lag2  Lag3  Lag4  Lag5 Volume  Today Direction
## 1 2001  0.381 -0.192 -2.624 -1.055  5.010 1.1913  0.959      Up
## 2 2001  0.959  0.381 -0.192 -2.624 -1.055 1.2965  1.032      Up
## 3 2001  1.032  0.959  0.381 -0.192 -2.624 1.4112 -0.623     Down
## 4 2001 -0.623  1.032  0.959  0.381 -0.192 1.2760  0.614      Up
## 5 2001  0.614 -0.623  1.032  0.959  0.381 1.2057  0.213      Up
## 6 2001  0.213  0.614 -0.623  1.032  0.959 1.3491  1.392      Up
```

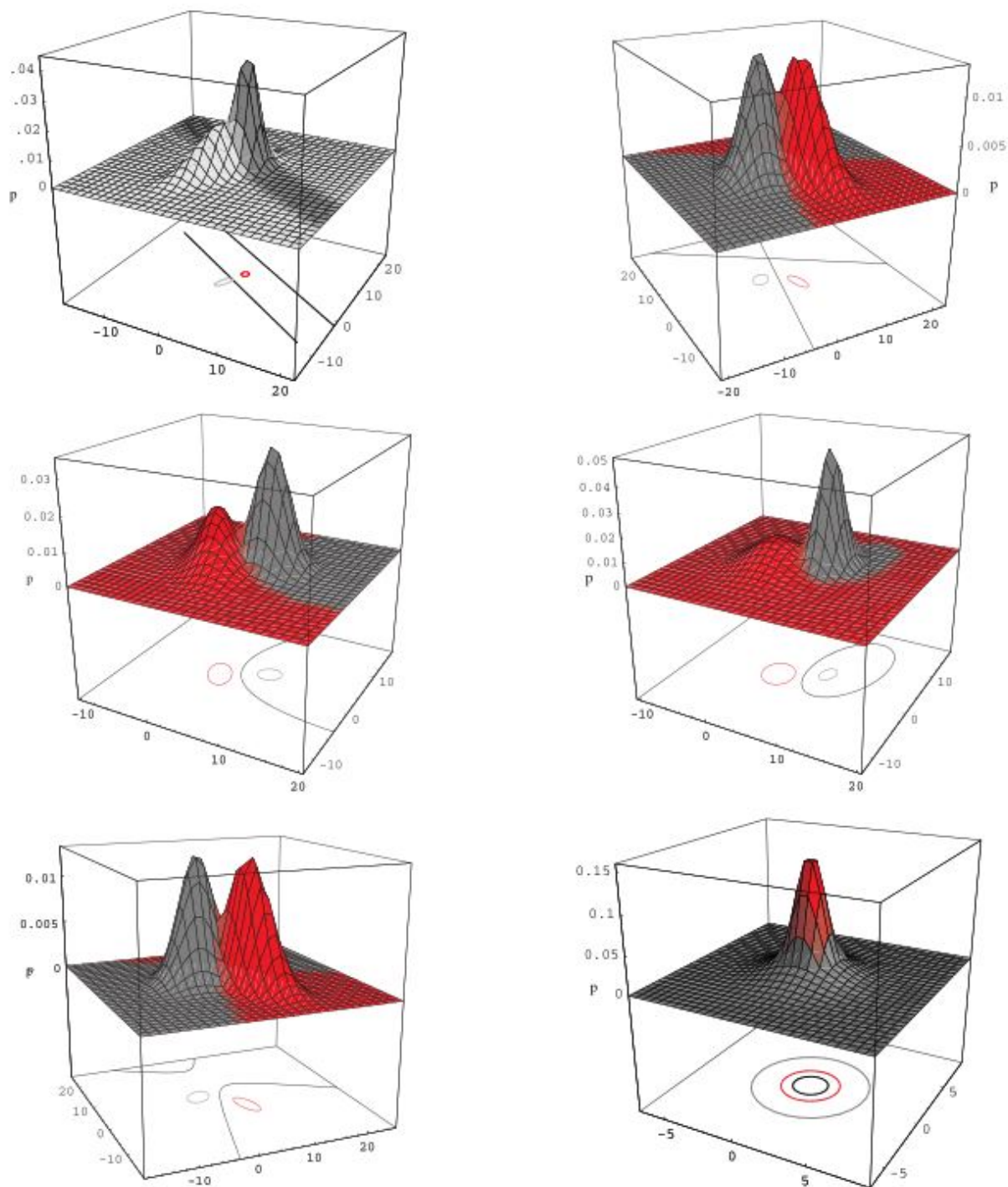
```
set.seed(2019)
dt.ucz <- Smarket %>%
  mutate_if(is.numeric, scale) %>%
  sample_frac(size = 2/3)
dt.test <- Smarket[~as.numeric(rownames(dt.ucz)),]
mod.qda <- qda(Direction~Lag1+Lag2, data = dt.ucz)
mod.qda
```

```
## Call:
## qda(Direction ~ Lag1 + Lag2, data = dt.ucz)
##
## Prior probabilities of groups:
##      Down      Up
## 0.4789916 0.5210084
##
## Group means:
##           Lag1           Lag2
## Down  0.01064899 -0.006153307
## Up    -0.08417078  0.030090639
```

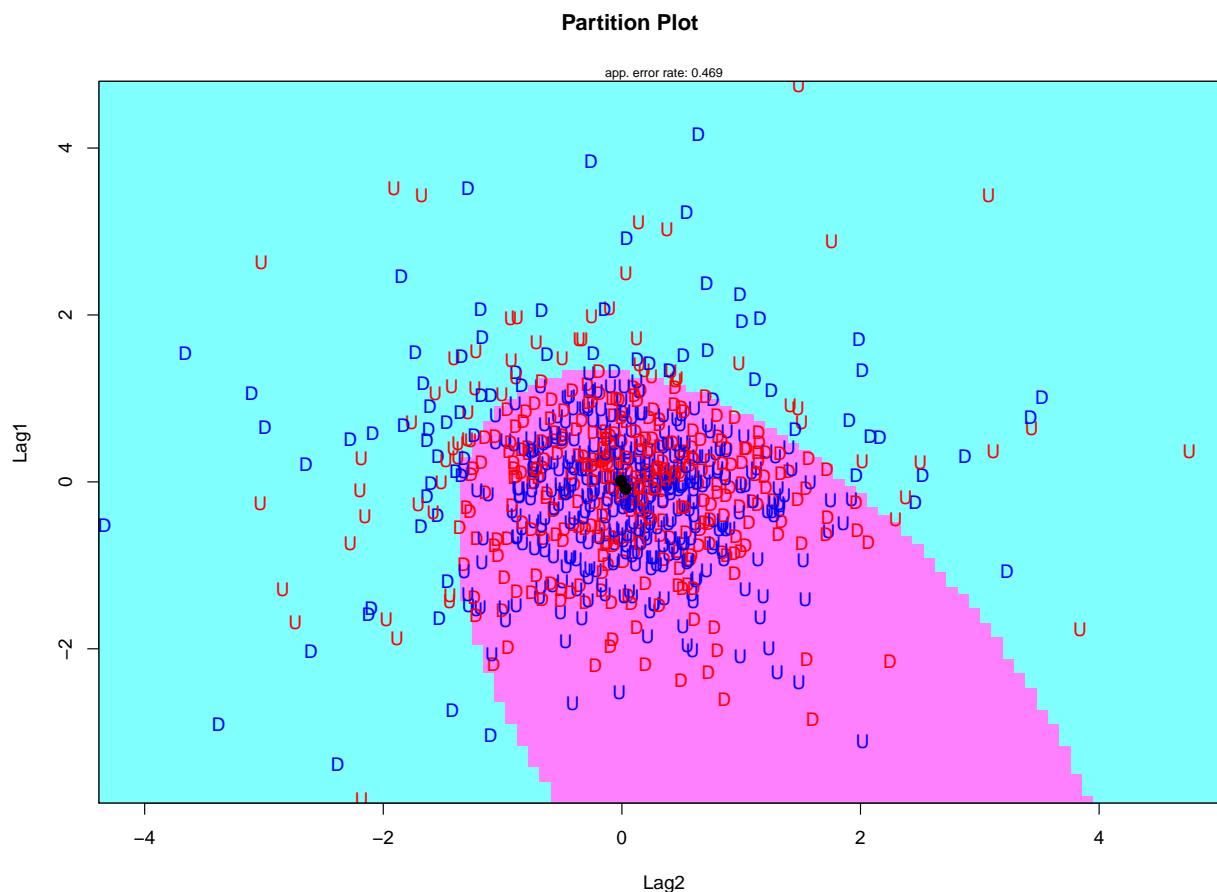
Ponieważ funkcje dyskryminacyjne mogą być nieliniowe, to podsumowanie modelu nie zawiera współczynników funkcji. Podsumowanie zawiera tylko prawdopodobieństwa a priori i średnie poszczególnych zmiennych niezależnych w klasach.

```
pred.qda <- predict(mod.qda, dt.test)
tab <- table(pred = pred.qda$class, dt.test$Direction)
tab

##
## pred  Down  Up
## Down   10  17
```



Rysunek 8.7: Przykład zastosowania kwadratowej analizy dyskryminacyjnej. Pokazane są dopuszczalne postaci zbiorów ograniczających. Źródło: @duda2001



Rysunek 8.8: Wykres klasyfikacji na podstawie QDA. Obserwacje zaznaczone kolorem niebieskim są prawidłowo zaklasyfikowane, a czerwonym źle

```
## Up 171 219
sum(diag(prop.table(tab)))

## [1] 0.5491607

library(klaR)
partimat(Direction ~ Lag1+Lag2,
  data = dt.ucz,
  method = "qda",
  col.correct='blue',
  col.wrong='red')
```

### 8.3 Analiza dyskryminacyjna metodą częściowych najmniejszych kwadratów

Analiza dyskryminacyjna metodą częściowych najmniejszych kwadratów (ang. *Partial Least Squares Discriminant Analysis*) jest wykorzystywana szczególnie w sytuacjach gdy zestaw predyktorów zawiera zmienne silnie ze sobą skorelowane. Jak wiadomo z wcześniejszych rozważań, metody dyskryminacji obserwacji są mało odporne na nadmiarowość zmiennych niezależnych. Stąd powstał pomysł zastosowania połączenia LDA z PLS (Partial Least Squares), której celem jest redukcja wymiaru przestrzeni jednocześnie maksymalizując



korelację zmiennych niezależnych ze zmienną wynikową.

Parametrem, który jest kontrolowany podczas budowy modelu jest liczba ukrytych zmiennych. Metoda PLS-DA ma kilka implementacji w R, ale najbardziej znana jest funkcja `plsda` z pakietu `caret` (Jed Wing et al. 2018).

**Przykład 8.3.** Kontynuując poprzedni przykład przeprowadzimy klasyfikację ruchu waloru korzystając z metody PLS-DA. W przeciwieństwie do poprzednich funkcji `plsda` potrzebuje przekazania zbioru predyktorów i wektora zmiennej wynikowej oddzielnie, a nie za pomocą formuły. Doboru liczby zmiennych latentnych dokonamy arbitralnie.

```
library(caret)
mod.plsda <- plsda(dt.ucz[, -c(1,7:9)],
                  as.factor(dt.ucz$Direction),
                  ncomp = 2)
mod.plsda$loadings
```

```
##
## Loadings:
##      Comp 1 Comp 2
## Lag1 -0.712  0.450
## Lag2  0.234  0.237
## Lag3  0.647
## Lag4  0.158  0.519
## Lag5      0.681
##
##              Comp 1 Comp 2
## SS loadings      1.008  1.001
## Proportion Var   0.202  0.200
## Cumulative Var   0.202  0.402
```

Dwie ukryte zmienne użyte do redukcji wymiaru przestrzeni wyjaśniają około 40% zmienności pierwotnych zmiennych. Ładunki (Loadings) pokazują kontrybucje poszczególnych zmiennych w tworzenie się zmiennych ukrytych.

```
pred.plsda <- predict(mod.plsda, dt.test[, -c(1,7:9)])
tab <- table(pred.plsda, dt.test$Direction)
tab
```

```
##
## pred.plsda Down  Up
##      Down   37  50
##      Up    144 186
```

```
sum(diag(prop.table(tab)))
```

```
## [1] 0.5347722
```

Ponieważ korelacje pomiędzy predyktorami w naszym przypadku nie były duże, to zastosowanie PLS-DA nie poprawiło znacząco klasyfikacji w stosunku do metody QDA.

```
cor(dt.ucz[, 2:6])
```

```
##      Lag1      Lag2      Lag3      Lag4      Lag5
## Lag1  1.000000000 -0.001713222  0.003820374  0.01583203  0.02504524
## Lag2 -0.001713222  1.000000000 -0.046611448 -0.02069792 -0.04105822
## Lag3  0.003820374 -0.046611448  1.000000000 -0.06142632 -0.03424691
## Lag4  0.015832026 -0.020697920 -0.061426325  1.00000000 -0.07102928
## Lag5  0.025045238 -0.041058218 -0.034246907 -0.07102928  1.00000000
```



## 8.4 Regularyzowana analiza dyskryminacyjna

Regularyzowana analiza dyskryminacyjna (ang. *Regularized Discriminant Analysis*) powstała jako technika równoważąca zalety i wady LDA i QDA. Ze względu na zdolności generalizacyjne model LDA jest lepszy od QDA (mniejsza wariancja modelu), ale jednocześnie QDA ma bardziej elastyczną postać hiperpowierzchni brzegowych rozdzielających obiekty różnych klas. Dlatego Friedman (1989) wprowadził technikę będącą kompromisem pomiędzy LDA i QDA poprzez odpowiednie określenie macierzy kowariancji

$$\tilde{\Sigma}_i(\lambda) = \lambda \Sigma_i + (1 - \lambda) \Sigma, \quad (8.61)$$

gdzie  $\Sigma_i$  jest macierzą kowariancji dla  $i$ -tej klasy, a  $\Sigma$  jest uśrednioną macierzą kowariancji wszystkich klas. Zatem odpowiedni dobór parametru  $\lambda$  decyduje czy poszukujemy modelu prostszego ( $\lambda = 0$  odpowiada LDA), czy bardziej elastycznego ( $\lambda = 1$  oznacza QDA).

Dodatkowo metoda RDA pozwala na elastyczny wybór pomiędzy postaciami macierzy kowariancji wspólnej dla wszystkich klas  $\Sigma$ . Może ona być macierzą jednostkową, jak w przypadku 8.2.1, co oznacza niezależność predyktorów modelu, może też być jak w przypadku 8.2.2, gdzie dopuszcza się korelacje między predyktorami. Dokonuje się tego przez odpowiedni dobór parametru  $\gamma$

$$\Sigma(\gamma) = \gamma \Sigma + (1 - \gamma) \sigma^2 I. \quad (8.62)$$

**Przykład 8.4.** Funkcja `rda` pakietu `klaR` jest implementacją powyższej metody. Ilustrują jej działania będzie klasyfikacja stanów z poprzedniego przykładu.

```
library(klaR)
mod.rda <- rda(Direction~Lag1+Lag2+Lag3+Lag4+Lag5, dt.ucz)
mod.rda

## Call:
## rda(formula = Direction ~ Lag1 + Lag2 + Lag3 + Lag4 + Lag5, data = dt.ucz)
##
## Regularization parameters:
##      gamma      lambda
## 0.33416870 0.03931045
##
## Prior probabilities of groups:
##      Down      Up
## 0.4789916 0.5210084
##
## Misclassification rate:
##      apparent: 43.337 %
## cross-validated: 45.524 %
```

Model został oszacowany z parametrami wyznaczonymi na podstawie sprawdzianu krzyżowego zastosowanego w funkcji `rda`.

```
pred.rda <- predict(mod.rda, dt.test)
(tab <- table(pred = pred.rda$class, dt.test$Direction))

##
## pred   Down  Up
##   Down   19  30
##    Up   162 206
##
## sum(diag(prop.table(tab)))

## [1] 0.5395683
```

Jakość klasyfikacji jest na zbliżonym poziomie jak przy poprzednich metodach.

## 8.5 Analiza dyskryminacyjna mieszana

Liniowa analiza dyskryminacyjna zakładała, że średnie (centroidy) w klasach są różne ale macierz kowariancji wszystkich klas jest jednakowa. Analiza dyskryminacyjna mieszana (ang. *Mixture Discriminant Analysis*) prezentuje jeszcze inne podejście ponieważ zakłada, że każda klasa może być charakteryzowana przez wiele wielowymiarowych rozkładów normalnych, których centroidy mogą się różnić, ale macierze kowariancji nie.

Wówczas rozkład dla danej klasy jest mieszaniną rozkładów składowych, a funkcja dyskryminacyjna dla  $i$ -tej klasy przyjmuje postać

$$g_i(\mathbf{x}) \propto \sum_{k=1}^{L_i} \phi_{ik} g_{ik}(\mathbf{x}), \quad (8.63)$$

gdzie  $L_i$  jest liczbą rozkładów składających się na  $i$ -tą klasę, a  $\phi_{ik}$  jest współczynnikiem proporcji estymowanych w czasie uczenia modelu.

**Przykład 8.5.** Funkcja `mda` pakietu `mda` (Trevor Hastie et al. 2017) jest implementacją tej techniki w R. Jej zastosowanie pokażemy na przykładzie danych giełdowych z poprzedniego przykładu. Użyjemy domyślnych ustawień funkcji (trzy rozkłady dla każdej klasy).

```
library(mda)
mod.mda <- mda(Direction~Lag1+Lag2+Lag3+Lag4+Lag5, dt.ucz)
mod.mda

## Call:
## mda(formula = Direction ~ Lag1 + Lag2 + Lag3 + Lag4 + Lag5, data = dt.ucz)
##
## Dimension: 5
##
## Percent Between-Group Variance Explained:
##      v1      v2      v3      v4      v5
## 48.45  88.33  94.80  99.68 100.00
##
## Degrees of Freedom (per dimension): 6
##
## Training Misclassification Error: 0.42737 ( N = 833 )
##
## Deviance: 1134.453

pred.mda <- predict(mod.mda, dt.test)
(tab <- table(pred = pred.mda, dt.test$Direction))

##
## pred   Down  Up
##   Down    23  38
##    Up   158 198

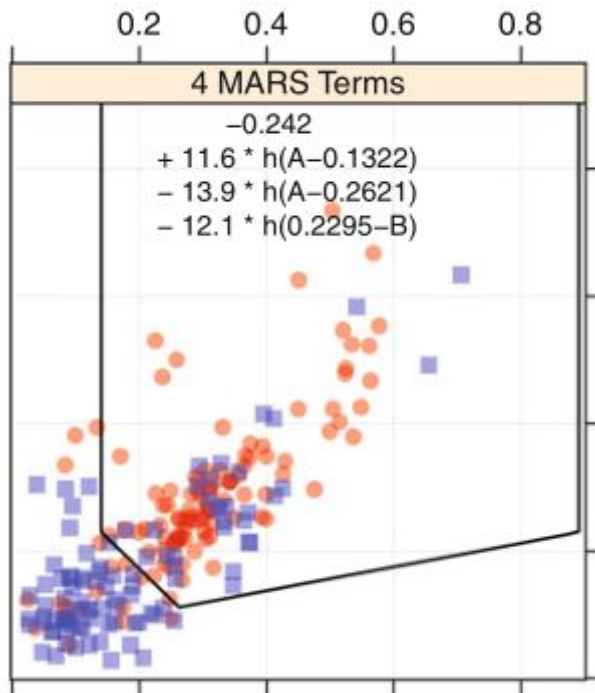
sum(diag(prop.table(tab)))

## [1] 0.529976
```

Kolejny raz model dyskryminacyjny charakteryzuje się podobną jakością klasyfikacji.

## 8.6 Elastyczna analiza dyskryminacyjna

zupełnie inne podejście w stosunku do wcześniejszych rozwiązań, prezentuje elastyczna analiza dyskryminacyjna (ang. *Flexible Discriminant Analysis*). Kodując klasy wynikowe jako zmienne dichotomiczne (dla każdej klasy jest odrębna zmienna wynikowa) dla każdej z nich budowanych jest  $k$  modeli regresji. Mogą to



Rysunek 8.9: Przykład klasyfikacji dwustanowej za pomocą metody FDA

być modele regresji penalizowanej, jak regresja grzbietowa lub LASSO, modele regresji wielomianowej albo modele regresji sklepanej (MARS), o których będzie mowa w dalszej części tego opracowania.

Przykładowo, jeśli modelem bazowym jest MARS, to funkcja dyskryminacyjna  $i$ -tej klasy może być postaci

$$g_i(\mathbf{x}) = \beta_0 + \beta_1 h(1 - x_1) + \beta_2 h(x_2 - 1) + \beta_3 h(1 - x_3) + \beta_4 h(x_1 - 1), \quad (8.64)$$

gdzie  $h$  są tzw. funkcjami bazowymi postaci

$$h(x) = \begin{cases} x, & x > 0 \\ 0, & x \leq 0. \end{cases} \quad (8.65)$$

Klasyfikacji dokonujemy sprawdzając znak funkcji dyskryminacyjnej  $g_i$ , jeśli jest dodatni, to funkcja przypisuje obiekt do klasy  $i$ -tej. W przeciwnym przypadku nie należy do tej klasy.

**Przykład 8.6.** Funkcja `fda` pakietu `mda` jest implementacją techniki FDA w R. Na podstawie danych z poprzedniego przykładu zostanie przedstawiona zasada działania. Przyjmiemy domyślne ustawienia funkcji, z wyjątkiem metody estymacji modelu, jako którą przyjmujemy MARS.

```
mod.fda <- fda(Direction~Lag1+Lag2, dt.ucz, method = mars)
mod.fda
```

```
## Call:
## fda(formula = Direction ~ Lag1 + Lag2, data = dt.ucz, method = mars)
##
## Dimension: 1
##
## Percent Between-Group Variance Explained:
## v1
## 100
##
## Training Misclassification Error: 0.43938 ( N = 833 )
```

Ponieważ, zmienna wynikowa jest dwustanowa, to powstała tylko jedna funkcja dyskryminacyjna. Parametry modelu są następujące

```
mod.fda$fit$coefficients
```

```
##           [,1]  
## [1,]  0.1129623  
## [2,] -0.5202437  
## [3,]  0.5462219
```

```
pred.fda <- predict(mod.fda, dt.test)  
(tab <- table(pred = pred.fda, dt.test$Direction))
```

```
##  
## pred   Down  Up  
##   Down  108 118  
##    Up    73 118
```

```
sum(diag(prop.table(tab)))
```

```
## [1] 0.5419664
```

Jakość klasyfikacji jest tylko nieco lepsza niż w przypadku poprzednich metod.

## Rozdział 9

# Klasyfikatory bayesowskie

Całą gamę klasyfikatorów opartych na twierdzeniu Bayesa nazywać będziemy bayesowskimi.

$$P(A|B) = \frac{P(A)P(B|A)}{P(B)}, \quad (9.1)$$

gdzie  $P(B) > 0$ .

Bayesowskie reguły podejmowania decyzji dały podstawy takich metod jak:

- liniowa analiza dyskryminacyjna;
- kwadratowa analiza dyskryminacyjna;

W ustaleniu klasyfikatora bayesowskiego będzie nam przyświecała cały czas ta sama reguła: *jeśli znam wartości cech charakteryzujących badane obiekty oraz klasy do których należą (w próbie uczącej), to na ich podstawie mogę wyznaczyć miary prawdopodobieństw a posteriori, które pomogą mi w ustaleniu klasy do której należy nowy testowy element.*

W dalszej części będziemy przyjmowali następujące oznaczenia:

- $T$  - zbiór danych uczących (treningowych),
- $T^j$  - zbiór danych uczących dla których przyjęliśmy decyzję o przynależności do  $j$ -tej klasy,
- $T_{a_i=v}^j$  - zbiór danych uczących o wartości atrybutu  $a_i$  równej  $v$  i klasy  $j$ -tej,
- $\mathbb{H}$  - przestrzeń hipotez,
- $P(h|a_1 = v_1, a_2 = v_2, \dots, a_p = v_p)$  - prawdopodobieństwo a posteriori, że prawdziwa jest hipoteza  $h \in \mathbb{H}$ , jeśli znamy atrybuty obiektu,
- $P(h)$  - prawdopodobieństwo a priori zajścia hipotezy  $h \in \mathbb{H}$ ,
- $c$  - prawdziwy stan obiektu.

### 9.1 Klasyfikator maximum a posteriori (MAP)

Na podstawie wiedzy o atrybutach obiektu  $x$  podejmujemy decyzję o klasyfikacji tego obiektu zgodnie z hipotezą  $h_{MAP} \in \mathbb{H}$ , która przyjmuje postać

$$h_{MAP} = \arg \max_{h \in \mathbb{H}} P(h|a_1 = v_1, a_2 = v_2, \dots, a_p = v_p) \quad (9.2)$$

$$= \arg \max_{h \in \mathbb{H}} P(a_1 = v_1, a_2 = v_2, \dots, a_p = v_p|h) \cdot P(h), \quad (9.3)$$

gdzie ostatnia równość wynika z twierdzenia Bayesa oraz faktu, że dla konkretnego obiektu  $x$  wielkości atrybutów nie zależą od postawionej hipotezy.

## 9.2 Klasyfikator największej wiarogodności (ML)

Na podstawie wiedzy o atrybutach obiektu  $x$  podejmujemy decyzję o klasyfikacji tego obiektu zgodnie z hipotezą  $h_{ML} \in \mathbb{H}$ , która przyjmuje postać

$$h_{ML} = \arg \max_{h \in \mathbb{H}} P(a_1 = v_1, a_2 = v_2, \dots, a_p = v_p | h). \quad (9.4)$$

*Uwaga.* Obie wspomniane metody wymagają znajomości prawdopodobieństwa  $P(a_1 = v_1, a_2 = v_2, \dots, a_p = v_p | h)$ , ale różnią się podejściem do wiedzy o prawdopodobieństwach a priori. W metodzie MAP brana pod uwagę jest wiedza o prawdopodobieństwie przynależności do poszczególnych klas, a w ML nie. Dla klasyfikacji, w których prawdopodobieństwa przynależności do klas są takie same, klasyfikatory MAP i ML są równoważne.

## 9.3 Naiwny klasyfikator Bayesa (NB)<sup>1</sup>

Największy problem w wyznaczeniu klasyfikatorów MAP i ML stanowi wyznaczenie rozkładu łącznego  $P(a_1 = v_1, a_2 = v_2, \dots, a_p = v_p | h)$ . W naiwnym klasyfikatorze Bayesa zakłada się niezależność warunkową poszczególnych atrybutów względem klasy do której ma należeć wg hipotezy obiekt. Założenie to często nie jest spełnione i stąd nazwa przymiotnik “naiwny”.

Definicja naiwnego klasyfikatora bayesowskiego różni się od klasyfikatora MAP tylko podejściem do prawdopodobieństwa a posteriori.

$$h_{NB} = \arg \max_{h_j \in \mathbb{H}} P(h_j) \prod_{i=1}^p P(a_i = v_i | h_j), \quad (9.5)$$

gdzie  $h_j$  oznacza hipotezę (decyzję), że badany obiekt należy do  $j$ -tej klasy.

Oczywiście zarówno prawdopodobieństwo a priori jak i a posteriori są wyznaczone na podstawie próby, i tak prawdopodobieństwo a priori wynosi

$$P(h_j) = P_T(h_j) = \frac{|T^j|}{|T|}, \quad (9.6)$$

gdzie  $|A|$  oznacza moc zbioru  $A$ .

Natomiast prawdopodobieństwo a posteriori dla  $i$ -tego atrybutu wynosi

$$P(a_i = v_i | h_j) = P_{T^j}(a_i = v_i) = \frac{|T_{a_i=v_i}^j|}{|T^j|}. \quad (9.7)$$

Na mocy powyższego możemy zauważyć, że jeżeli założenie o warunkowej niezależności jest spełnione, to klasyfikatory NB i MAP są równoważne.

Chcąc przypisać klasę nowemu obiektowi powstaje problem praktyczny, polegający na tym, że dla pewnych konfiguracji atrybutów nie ma odpowiedników w nauczonym modelu. Powodem takiego stanu rzeczy jest fakt, że takie kombinacje nie wystąpiły w próbie uczącej.

Istnieją dwa sposoby predykcji w takiej sytuacji:

1.

$$P(a_i = v_i | h_j) = \begin{cases} \frac{|T_{a_i=v_i}^j|}{|T^j|}, & T_{a_i=v_i}^j \neq \emptyset \\ \epsilon, & \text{w przeciwnym przypadku.} \end{cases} \quad (9.8)$$

W tym przypadku przyjmuje się, że  $\epsilon \ll 1/|T_j|$ .

---

<sup>1</sup>ang. *Naive Bayes Classifier*

2. Drugi sposób wykorzystuje estymację z poprawką

$$P(a_i = v_i | h_j) = \frac{|T_{a_i=v_i}^j| + mp}{|T^j| + mp}, \quad (9.9)$$

gdzie  $p$  oznacza prawdopodobieństwo a priori przyjęcia przez atrybut  $a$  wartości  $v$  (najczęściej  $p = 1/|A|$ ,  $A$  - zbiór wszystkich możliwych wartości atrybutu  $a$ ),  $m$  - waga (najczęściej  $m = |A|$ ).

W przypadku gdy atrybuty są mierzone na skali ciągłej najczęściej stosuje się dyskretyzację ich do zmiennych ze skali przedziałowej. Inna metoda stosowana w przypadku ciągłych atrybutów, to użycie gęstości  $g_i^j$  o rozkładzie normalnym w miejsce  $P(a_i = v_i | h_j)$ . Przy czym do obliczenia parametrów rozkładu stosujemy wzory

$$m_i^j = \frac{1}{|T^j|} \sum_{x \in T^j} a_i(x), \quad (9.10)$$

oraz

$$(s_i^j)^2 = \frac{1}{|T^j| - 1} \sum_{x \in T^j} (a_i(x) - m_i^j)^2. \quad (9.11)$$

Obsługa braków danych przez naiwny klasyfikator Bayesa jest dość prosta i opiera się na liczeniu prawdopodobieństw a posteriori wyłącznie dla obiektów, których wartości atrybutów są znane. Dlatego prawdopodobieństwa warunkowe liczy się wg wzoru

$$P(a_i = v_i | h_j) = \frac{|T_{a_i=v_i}^j|}{|T^j| - |T_{a_i=NA}^j|}. \quad (9.12)$$

Jeśli brakujące dane nie niosą w sobie istotnych informacji dotyczących klasyfikacji obiektów, to naiwny klasyfikator Bayesa będzie działał poprawnie.

Naiwny klasyfikator Bayesa jest implementowany w pakietach **e1071** (Meyer et al. 2019) i **klaR** (Weihs et al. 2005).

**Przykład 9.1.** Przeprowadzimy klasyfikację dla zbioru **Titanic**. W przypadku funkcji z pakietu **e1071** nie potrzeba zamieniać tabeli na przypadki. W pakiecie **klaR** istnieje inna funkcja budująca klasyfikator Bayesa **NaiveBayes**, ale w tym przypadku jeśli zbiór jest w formie tabeli, to należy go zamienić na ramkę danych z oddzielnymi przypadkami.

```
library(e1071)
```

```
Titanic
```

```
## , , Age = Child, Survived = No
```

```
##
```

```
##      Sex
```

```
## Class Male Female
```

```
## 1st      0      0
```

```
## 2nd      0      0
```

```
## 3rd     35     17
```

```
## Crew     0      0
```

```
##
```

```
## , , Age = Adult, Survived = No
```

```
##
```

```
##      Sex
```

```
## Class Male Female
```

```
## 1st    118      4
```

```
## 2nd    154     13
```

```
## 3rd    387     89
```

```
## Crew   670      3
```

```
##
```

```
## , , Age = Child, Survived = Yes
##
##      Sex
## Class  Male Female
##  1st      5      1
##  2nd     11     13
##  3rd     13     14
##  Crew      0      0
##
## , , Age = Adult, Survived = Yes
##
##      Sex
## Class  Male Female
##  1st     57     140
##  2nd     14      80
##  3rd     75      76
##  Crew    192      20
```

```
nb <- naiveBayes(Survived ~ ., data = Titanic)
nb$apriori
```

```
## Survived
##   No  Yes
## 1490  711
```

Poniższe tabele zawierają warunkowe prawdopodobieństwa przynależności do poszczególnych klas.

```
nb$tables
```

```
## $Class
##      Class
## Survived  1st      2nd      3rd      Crew
##      No  0.08187919 0.11208054 0.35436242 0.45167785
##      Yes 0.28551336 0.16596343 0.25035162 0.29817159
##
## $Sex
##      Sex
## Survived  Male      Female
##      No  0.91543624 0.08456376
##      Yes 0.51617440 0.48382560
##
## $Age
##      Age
## Survived  Child      Adult
##      No  0.03489933 0.96510067
##      Yes 0.08016878 0.91983122
```

```
dane <- as.data.frame(Titanic)
pred <- predict(nb, dane)
pred
```

```
## [1] Yes No  No  No  Yes Yes Yes Yes No  No  No  No  Yes Yes Yes Yes Yes
## [18] No  No  No  Yes Yes Yes Yes No  No  No  No  Yes Yes Yes Yes
## Levels: No Yes
```

```
tab <- table(pred, dane$Survived)
tab
```



```
##
## pred  No  Yes
##    No   7   7
##    Yes  9   9
```

```
sum(diag(prop.table(tab)))
```

```
## [1] 0.5
```

Naiwny klasyfikator spisał się bardzo słabo, ponieważ klasyfikacja na poziomie 0.5 jest taka jak przy rzucie monetą.

**Przykład 9.2.** Przeprowadzimy klasyfikację gatunków irysów na podstawie szerokości i długości kielicha i płatka.

```
library(klaR)
set.seed(2019)
uczaca <- sample(1:nrow(iris), 2*nrow(iris)/3)
pr.ucz <- iris[uczaca,]
pr.test <- iris[-uczaca,]
nb2 <- NaiveBayes(Species~., data = pr.ucz)
nb2$apriori
```

```
## grouping
##      setosa versicolor  virginica
##      0.36      0.31      0.33
```

Prawdopodobieństwa a priori zostały oszacowane na podstawie próby uczącej. Poniższe tabele zawierają średnie i odchylenia standardowe zmiennych w poszczególnych klasach.

```
nb2$tables
```

```
## $Sepal.Length
##           [,1]      [,2]
## setosa      4.994444 0.3438807
## versicolor  5.977419 0.5613731
## virginica   6.603030 0.7359029
##
## $Sepal.Width
##           [,1]      [,2]
## setosa      3.436111 0.3586903
## versicolor  2.838710 0.2996414
## virginica   2.942424 0.3211603
##
## $Petal.Length
##           [,1]      [,2]
## setosa      1.472222 0.1782632
## versicolor  4.316129 0.4576001
## virginica   5.606061 0.6269666
##
## $Petal.Width
##           [,1]      [,2]
## setosa      0.2444444 0.09394358
## versicolor  1.3354839 0.18537959
## virginica   2.0000000 0.25860201
```

```
pred <- predict(nb2, newdata = pr.test)
tab <- table(pred$class, pr.test$Species)
```

```
tab
```

```
##
##           setosa versicolor virginica
##  setosa         14          0          0
##  versicolor      0          18         1
##  virginica       0           1        16
```

```
sum(diag(prop.table(tab)))
```

```
## [1] 0.96
```

Klasyfikacja na podstawie modelu jest bardzo dobra (96%).

## 9.4 Zalety i wady

- Zalety:
  - prostota konstrukcji i prosty algorytm;
  - jeśli jest spełnione założenie warunkowej niezależności, to ten klasyfikator działa szybciej i czasem lepiej niż inne metody klasyfikacji;
  - nie potrzebuje dużych zbiorów danych do estymacji parametrów;
- Wady:
  - często nie spełnione założenie o warunkowej niezależności powoduje obciążenie wyników;
  - brak możliwości wprowadzania interakcji efektów kilku zmiennych;
  - potrzebuje założenia normalności warunkowych gęstości w przypadku ciągłych atrybutów;
  - często istnieją lepsze klasyfikatory.

## Rozdział 10

# Metoda $k$ najbliższych sąsiadów

Technika  $k$  najbliższych sąsiadów (ang. *k-Nearest Neighbors*) przewiduje wartość zmiennej wynikowej na podstawie  $k$  najbliższych obserwacji zbioru uczącego. W przeciwieństwie do wspomnianych wcześniej modeli liniowych, nie posiada ona jawnej formy i należy do klasy technik nazywanych czarnymi skrzynkami (ang. *black box*). Może być wykorzystywana, zarówno do zadań klasyfikacyjnych, jak i regresyjnych. W obu przypadkach predykcja dla nowych wartości predyktorów przebiega podobnie.

Niech  $\mathbf{x}_0$  będzie obserwacją, dla której poszukujemy wartości zmiennej wynikowej  $y_0$ . Na podstawie zbioru obserwacji  $\mathbf{x} \in T$  zbioru uczącego wyznacza się  $k$  najbliższych sąsiadów<sup>1</sup>, gdzie  $k$  jest z góry ustaloną wartością. Następnie, jeśli zadanie ma charakter klasyfikacyjny, to  $y_0$  przypisuje się modę zmiennej wynikowej obserwacji będących  $k$  najbliższymi sąsiadami. W przypadku zadań regresyjnych  $y_0$  przypisuje się średnią lub medianę.

Olbrzymie znaczenie dla wyników predykcji na podstawie metody  $kNN$  ma dobór metryki. Nie istnieje obiektywna technika wyboru najlepszej metryki, dlatego jej doboru dokonujemy metodą prób i błędów. Należy dodatkowo pamiętać, że wielkości mierzone  $\mathbf{x}$  mogą się różnić zakresami zmienności, a co za tym idzie, mogą znacząco wpłynąć na mierzone odległości pomiędzy punktami. Dlatego zaleca się standaryzację zmiennych przed zastosowaniem metody  $kNN$ .

Kolejnym parametrem, który ma znaczący wpływ na predykcję, jest liczba sąsiadów  $k$ . Wybór zbyt małej liczby  $k$  może doprowadzić do przeuczenia modelu jak to jest pokazane na rysunku 10.1

Z kolei zbyt duża liczba sąsiadów powoduje obciążenie wyników (patrz rysunek 10.2)

Dopiero dobór odpowiedniego  $k$  daje model o stosunkowo niskiej wariancji i obciążeniu. Najczęściej liczby  $k$  poszukujemy za pomocą próbkowania.

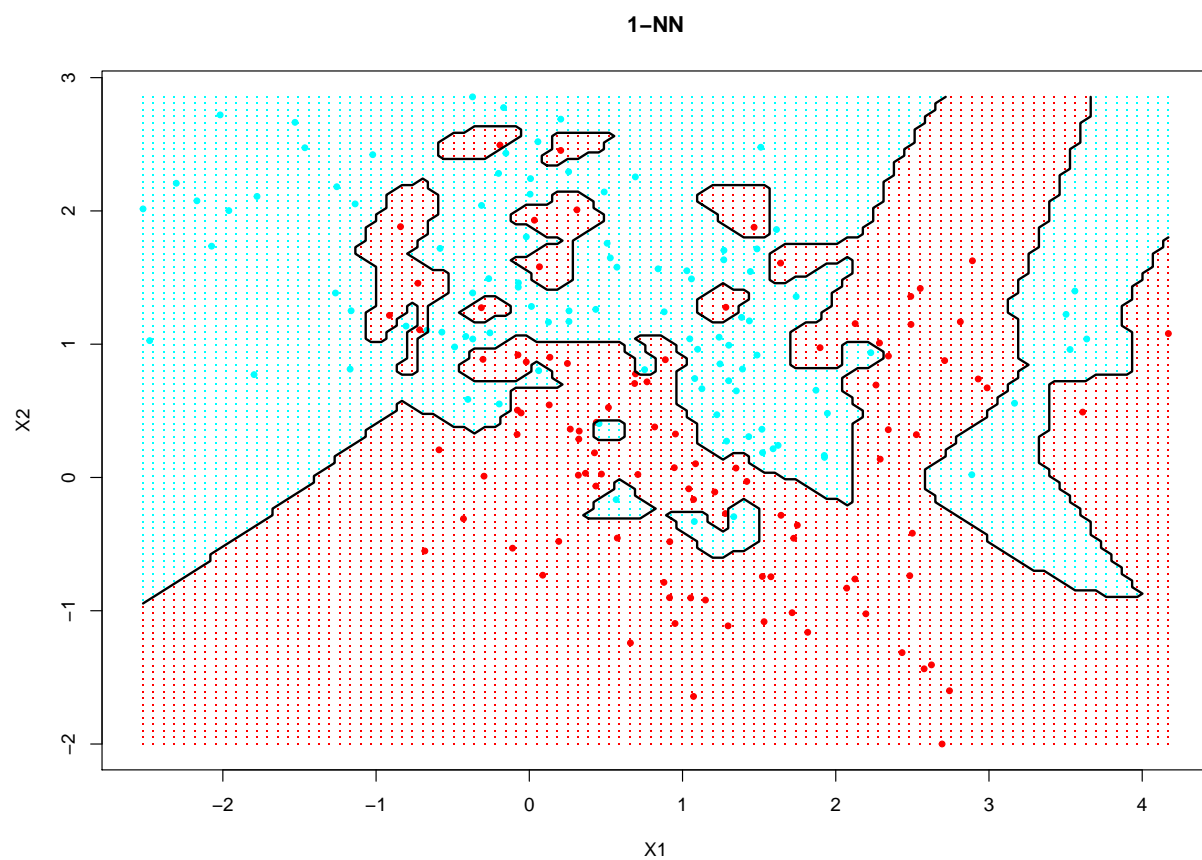
**Przykład 10.1.** Klasyfikację z wykorzystaniem metody  $kNN$  przeprowadzimy na przykładzie danych zbioru `spam` pakietu `ElemStatLearn`. Metoda  $kNN$  ma wiele implementacji R-owych ale na potrzeby przykładu wykorzystamy funkcję `knn3` pakietu `caret`.

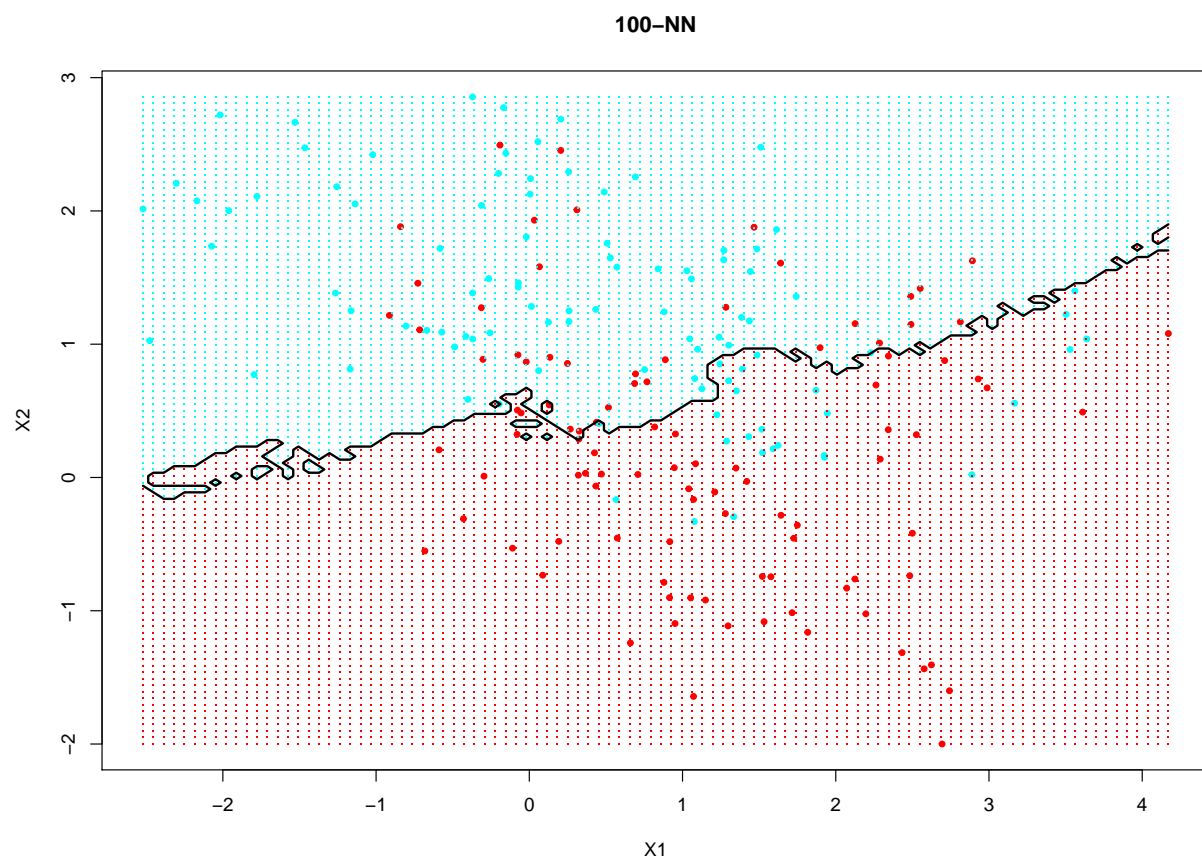
Najpierw dokonamy oszacowania optymalnego  $k$

```
library(ElemStatLearn)
library(tidyverse)

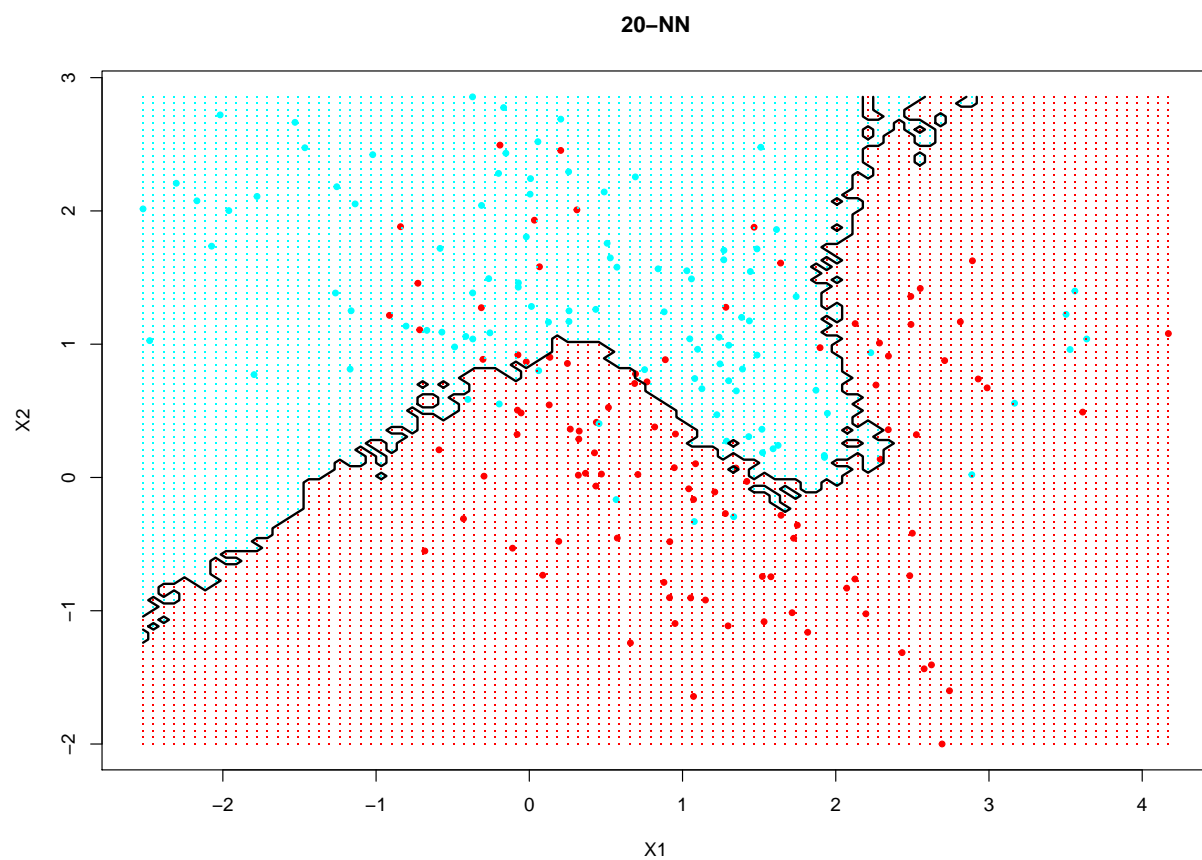
spam.std <- spam %>%
  mutate_if(is.numeric, scale)
set.seed(123)
ind <- sample(nrow(spam), size = nrow(spam)*2/3)
dt.ucz <- spam.std[ind,]
```

<sup>1</sup>metrykę można wybierać dowolnie, choć najczęściej jest to metryka euklidesowa

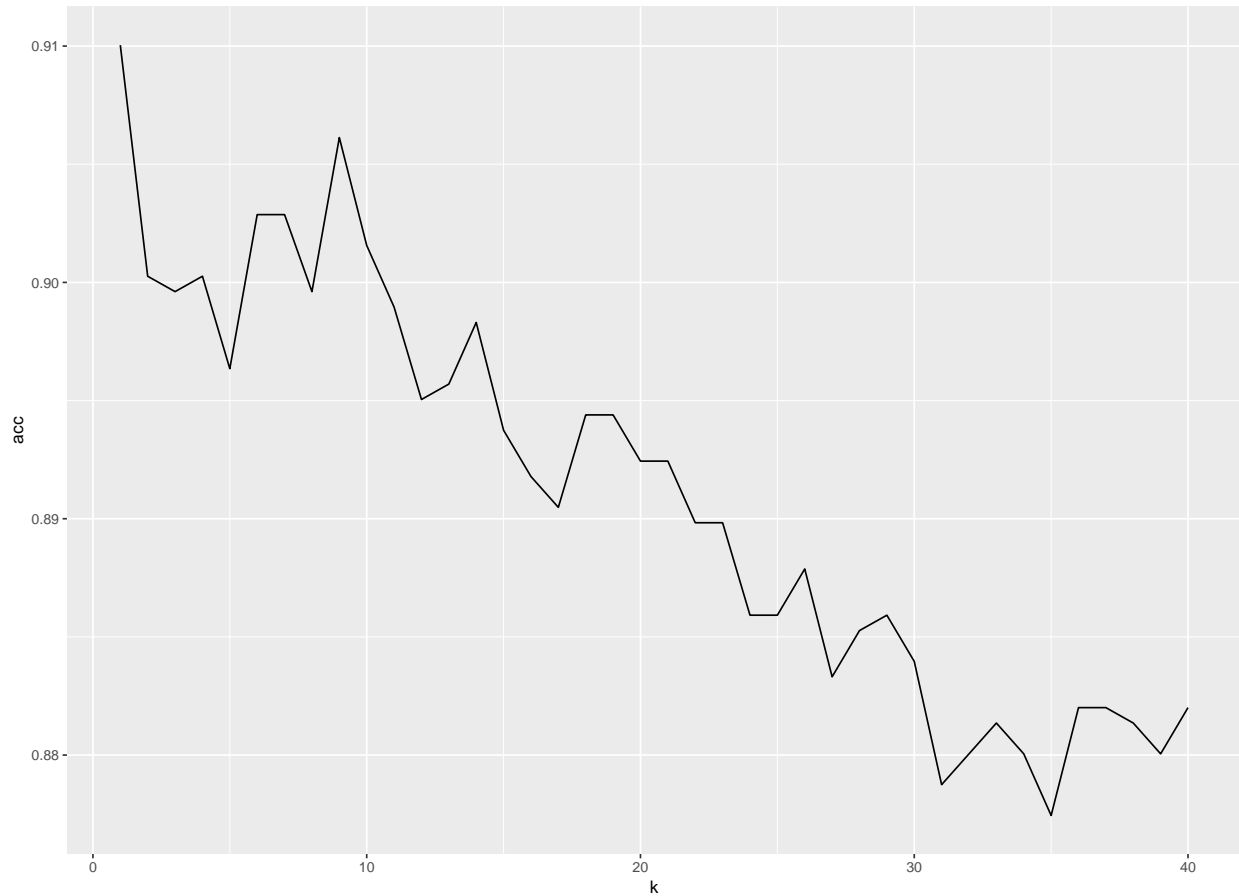
Rysunek 10.1: Przykład klasyfikacji dla  $k = 1$



Rysunek 10.2: Przykład zastosowania 100 sąsiadów



Rysunek 10.3: Model z optymalną liczbą sąsiadów



Rysunek 10.4: Ocena jakości dopasowania modelu dla różnej liczby sąsiadów

```
dt.test <- spam.std[-ind,]

acc <- function(pred, obs){
  tab <- table(pred,obs)
  acc <- sum(diag(prop.table(tab)))
  acc
}

1:40 %>%
  map(~knn3(spam~., data = dt.ucz, k = .x)) %>%
  map(~predict(.x, newdata = dt.test, type = "class")) %>%
  map_dbl(~acc(pred = .x, obs = dt.test$spam)) %>%
  tibble(k = 1:length(.), acc=.) %>%
  ggplot(aes(k, acc))+
  geom_line()
```

Biorąc pod uwagę wykres 10.4 można rozważać 1 lub 9 sąsiadów jako optymalne rozwiązanie, ponieważ wówczas poprawność klasyfikacji jest najwyższa. Proponuje unikać rozwiązania z 1 najbliższym sąsiadem ponieważ, będzie się ono charakteryzowało dużą zmiennością. Wybór  $k = 9$  wydaje się być optymalny.

```
mod.knn <- knn3(spam~., data = dt.ucz,
  k = 9)
```

```
mod.knn
```

```
## 9-nearest neighbor model
## Training set outcome distribution:
##
## email spam
## 1872 1195
```

Predykcyi dokonujemy w ten sam sposób co w innych modelach klasyfikacyjnych

```
pred.knn.class <- predict(mod.knn, newdata = dt.test, type = "class")
head(pred.knn.class)
```

```
## [1] spam email email spam spam email
## Levels: email spam
```

```
pred.knn <- predict(mod.knn, newdata = dt.test)
head(pred.knn)
```

```
##          email      spam
## [1,] 0.0000000 1.0000000
## [2,] 0.6666667 0.3333333
## [3,] 0.5555556 0.4444444
## [4,] 0.2222222 0.7777778
## [5,] 0.0000000 1.0000000
## [6,] 0.8888889 0.1111111
```

```
(tab <- table(pred.knn.class, dt.test$spam))
```

```
##
## pred.knn.class email spam
##          email    863    88
##          spam     53   530
```

```
sum(diag(prop.table(tab)))
```

```
## [1] 0.9080834
```



# Rozdział 11

## Uogólnione modele addytywne

Modele liniowe, jako techniki klasyfikacji i regresji, mają niewątpliwą zaletę - jawna postać zależności pomiędzy predyktorami i zmienną wynikową. Często w rzeczywistości tak uproszczony model nie potrafi oddać złożoności natury badanego zjawiska. Dlatego powstał pomysł aby w miejsce kombinacji liniowej predyktorów wstawić kombinację liniową ich funkcji, czyli

$$E(Y|X) = f(X) = \sum_{i=1}^M \beta_m h_m(X), \quad (11.1)$$

gdzie  $h_m : \mathbb{R}^d \rightarrow \mathbb{R}$  nazywana często funkcją bazową (ang. *linear basis expansion*). Wówczas w zależności od postaci funkcji bazowej otrzymujemy modele z różnymi poziomami elastyczności:

- gdy  $h_m(X) = X_m$ ,  $m = 1, \dots, M$ , to otrzymujemy model liniowy;
- gdy  $h_m(X) = X_j^2$  lub  $h_m(X) = X_j X_k$ , to otrzymujemy struktury wielomianowe, charakteryzujące się większą elastycznością modelu;
- gdy  $h_m(X) = \log X_j$  lub  $h_m(X) = \sqrt{X_j}$ , to uzyskujemy nieliniowość czynników wchodzących w skład kombinacji liniowej (11.1);
- dopuszczalne są również kawałkami liniowe funkcje postaci  $h_m(X) = I(l_m \leq X_k < u_m)$ , gdzie  $I$  jest funkcją charakterystyczną (ang. *indicator*) przedziału  $[l_m, u_m)$ .

Zbiory wszystkich funkcji bazowych definiowanych w ten sposób tworzy słownik funkcji bazowych  $\mathcal{D}$ . Aby kontrolować złożoność modeli, mając do dyspozycji tak zasobny słownik, wprowadza się następujące podejścia:

- ogranicza się klasę dostępnych funkcji bazowych

$$f(X) = \sum_{j=1}^d f_j(X_j) = \sum_{j=1}^d \sum_{m=1}^{M_j} \beta_{jm} h_{jm}(X_j), \quad (11.2)$$

- włącza się do modelu jedynie te funkcje ze słownika  $\mathcal{D}$ , które istotnie poprawiają dopasowanie modelu,
- używa się metod penalizowanych, czyli dopuszcza się stosowanie wszystkich funkcji bazowych ze słownika  $\mathcal{D}$ , ale współczynniki przy nich stojące są ograniczane.

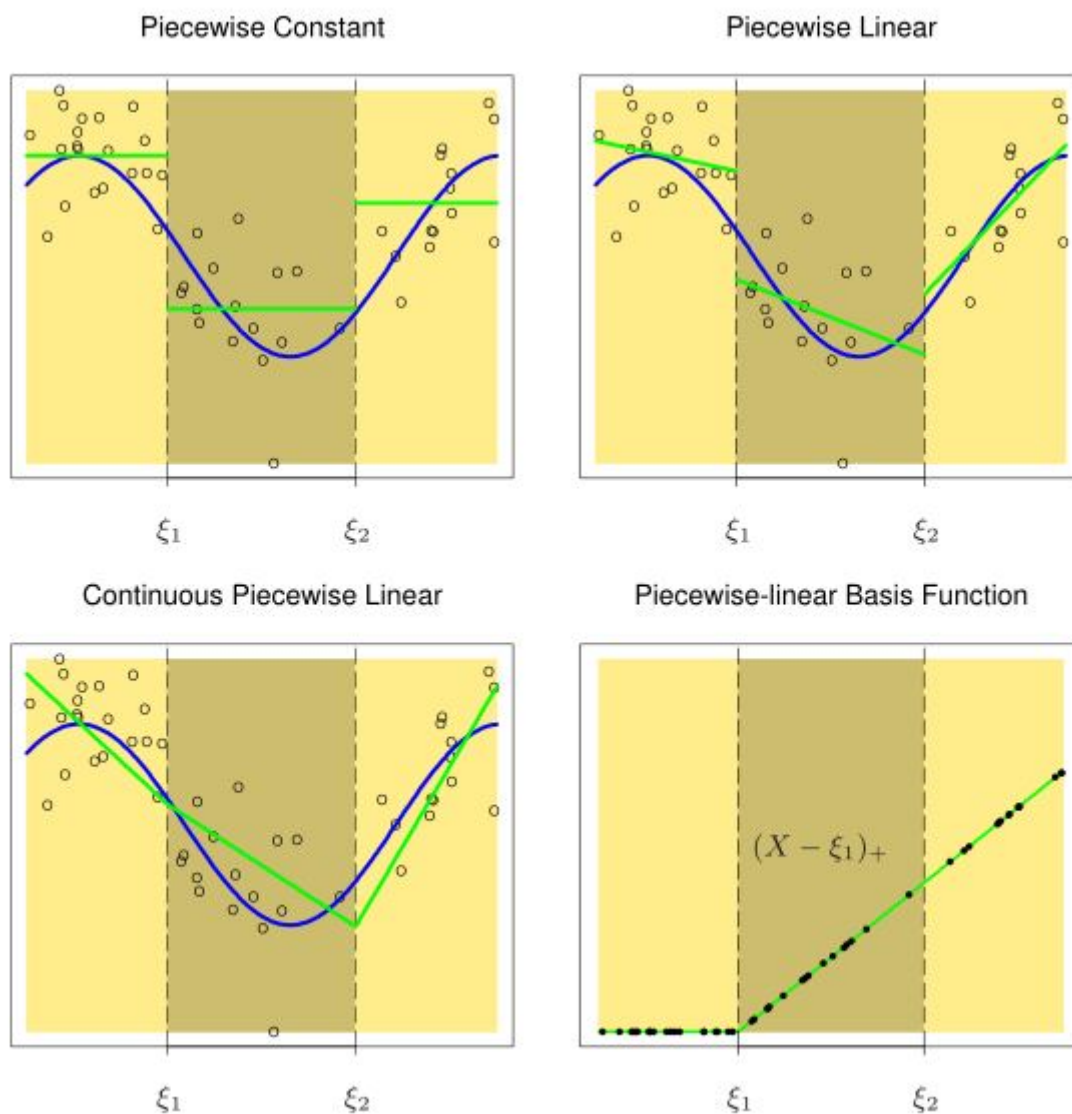
### 11.1 Przypadek jednowymiarowy

Dla uproszczenia rozważań przyjmiemy, że  $X$  jest jednowymiarowe.

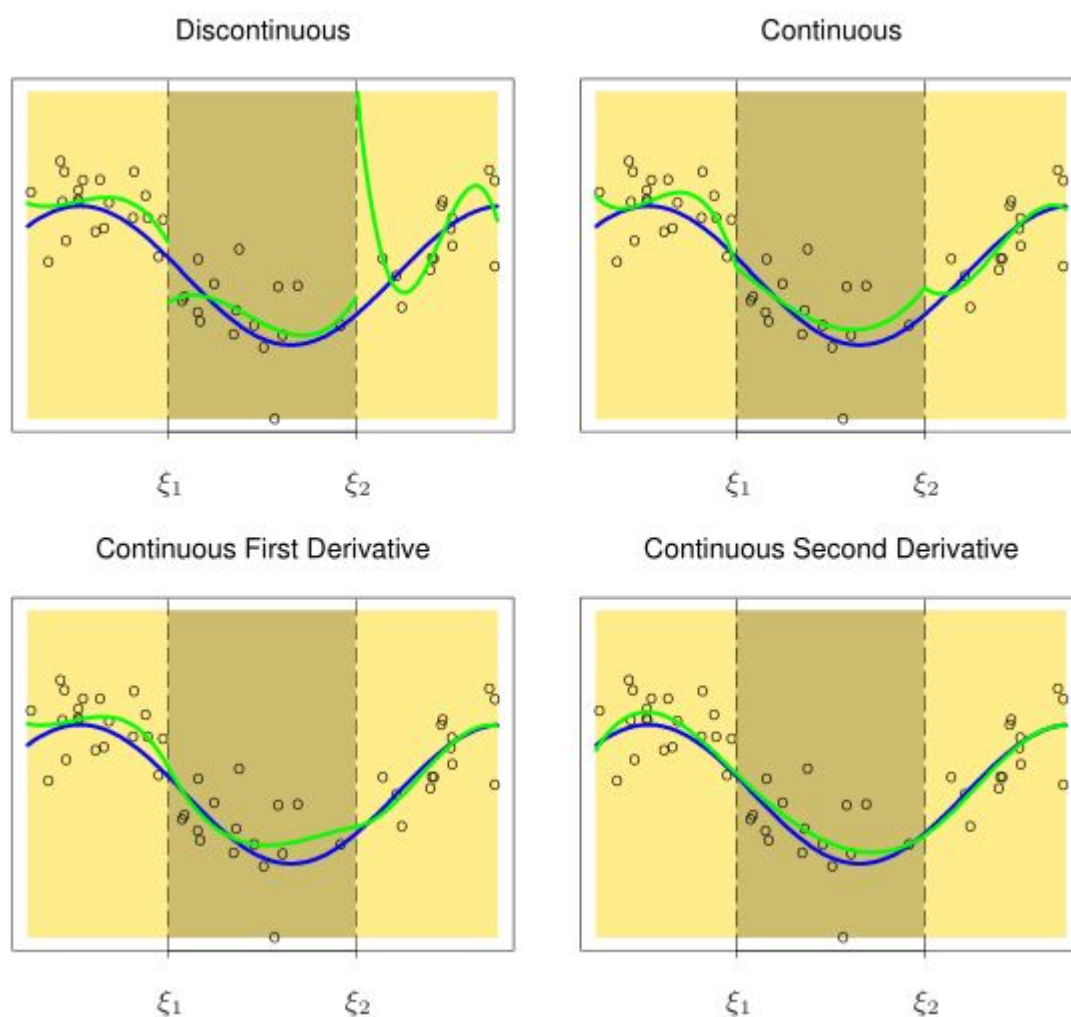
Przykładowo sześcienny splajn<sup>1</sup> dla dwóch punktów węzłowych składa się z następujących funkcji bazowych

---

<sup>1</sup>czyli funkcja łączona



Rysunek 11.1: Przykładowe zastosowanie kilku rodzajów funkcji bazowych. Wykres w lewym górnym rogu powstał ze stałych na przedziałach, wykres w górnym prawym rogu powstał z liniowych funkcji bazowych na przedziałach, w lewym dolnym rogu model powstał również z liniowych funkcji bazowych na przedziałach ale z założeniem ciągłości, a prawym dolnym rogu powstał z zastosowania funkcji bazowej  $\max(X - \xi_1, 0)$



Rysunek 11.2: Kolejne wykresy przedstawiają coraz bardziej gładkie modele będące efektem dodawania wielomianów trzeciego stopnia na przedziałach. Na każdym kolejnym modelu wymuszone zostały silniejsze założenia dotyczące gładkości

$$h_1(X) = 1, \quad h_3(X) = X^2, \quad h_5(X) = (X - \xi_1)_+^3 \quad (11.3)$$

$$h_2(X) = X, \quad h_4(X) = X^3, \quad h_6(X) = (X - \xi_2)_+^3. \quad (11.4)$$

Zachowanie wielomianów poza punktami węzłowymi jest czasami bardzo dziwne. Zdarza się, że charakteryzują się tam dużą zmiennością. Dlatego wprowadza się takie splajny aby w obszarach brzegowych zachowywały się przewidywalnie. Naturalny splajn sześcienny zakłada liniowość modelu poza węzłami brzegowymi. Dla  $K$  węzłów naturalny splajn sześcienny składa się z następujących funkcji bazowych

$$N_1(X) = 1, \quad N_2(X) = X, \quad N_{k+2}(X) = d_k(X) - d_{K-1}(X), \quad (11.5)$$

$$\text{gdzie } d_k(X) = \frac{(X - \xi_k)_+^3 - (X - \xi_K)_+^3}{\xi_K - \xi_k}.$$

Estymacji parametrów modelu dokonujemy metodą najmniejszych kwadratów, minimalizując

$$RSS(f, \lambda) = \sum_{i=1}^N (y_i - f(x_i))^2 + \lambda \int (f''(t))^2 dt, \quad (11.6)$$

gdzie  $\lambda$  jest parametrem wygładzania. Pierwsze wyrażenie po prawej stronie to ocena dopasowania, a drugie to kara za krzywoliniowość. Dla naturalnego splajna

$$f(x) = \sum_{j=1}^N N_j(x) \beta_j \quad (11.7)$$

minimalizujemy

$$RSS(\beta, \lambda) = (\mathbf{y} - \mathbf{N}\beta)'(\mathbf{y} - \mathbf{N}\beta) + \lambda \beta' \mathbf{\Omega} \beta, \quad (11.8)$$

gdzie  $\{\mathbf{N}\}_{ij} = N_j(x_i)$  i  $\{\mathbf{\Omega}\}_{jk} = \int N_j''(t) N_k''(t) dt$ . Rozwiązaniem zaganienia minimalizacji  $RSS(\beta, \lambda)$  jest

$$\hat{\beta} = ((\mathbf{N}'\mathbf{N}) + \lambda \mathbf{\Omega})^{-1} \mathbf{N}'\mathbf{y}. \quad (11.9)$$

## 11.2 Przypadek wielowymiarowy

W przypadku gdy  $X \in \mathbb{R}^d$  poszukujemy takiej  $d$ -wymiarowej regresji  $f(x)$ , która będzie minimalizowała wyrażenie

$$\min_f \sum_{i=1}^N (y_i - f(x_i))^2 + \lambda J(f), \quad (11.10)$$

gdzie  $J$  jest odpowiednią funkcją wyrażającą krzywoliniowość modelu. Dla  $X \in \mathbb{R}^2$  przyjmuje postać

$$J(f) = \iint_{\mathbb{R}^2} \left[ \left( \frac{\partial^2 f(x)}{\partial^2 x_1} \right)^2 + 2 \left( \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \right)^2 + \left( \frac{\partial^2 f(x)}{\partial^2 x_2} \right)^2 \right] dx_1 dx_2. \quad (11.11)$$

Rozwiązanie przyjmuje postać

$$f(x) = \beta_0 + \beta'x + \sum_{i=1}^N \alpha_i h_i(x), \quad (11.12)$$

gdzie  $h_i(x) = \|x - x_j\|^2 \log \|x - x_j\|$ .

## 11.3 Uogólnione modele addytywne

Przez uogólnione modele addytywne (ang. *Generalized Additive Models*) rozumiemy klasę modeli, które poprzez funkcję łączącą, opisują warunkową wartość zmiennej wynikowej w następujący sposób

$$g(E(Y|X)) = g(\mu(X)) = \alpha + f_1(X_1) + \dots + f_d(X_d), \quad (11.13)$$

gdzie  $g$  jest funkcją łączącą. Najczęściej stosowanymi funkcjami łączącymi są:

- $g(\mu) = \mu$  - stosowana w modelach, gdy zmienna wynikowa ma rozkład normalny;
- $g(\mu) = \text{logit } \mu$  - stosowana, gdy zmienna wynikowa ma rozkład dwumianowy;
- $g(\mu) = \text{probit } \mu$  - stosowana również w przypadku gdy zmienna ma rozkład dwumianowy, a  $\Phi^{-1}$  oznacza odwrotność dystrybuanty standaryzowanego rozkładu normalnego;
- $g(\mu) = \log \mu$  - stosowana, gdy zmienna wynikowa jest zmienną typu zliczeniowego (rozkład Poissona).

### 11.3.1 Algorytm uczenia modelu GAM

Algorytm uczenia wstecznego (ang. *backfitting*) przebiega wg następujących kroków:

1. Ustalamy wstępne oszacowania na  $\alpha = \bar{y}$  i  $\hat{f}_j = 0$ .
2. Dla  $j = 1, \dots, d, 1, \dots, d, 1, \dots$  powtarzamy szacowanie

$$\hat{f}_j \leftarrow \mathcal{S}_j \left[ (y_i - \hat{\alpha} - \sum_{k \neq j} \hat{f}_k(x_{ik}))_1^N \right], \quad (11.14)$$

$$\hat{f}_j \leftarrow \hat{f}_j - \frac{1}{N} \sum_{i=1}^N \hat{f}_j(x_{ij}) \quad (11.15)$$

dopóki  $\hat{f}_j$  osiągnie zbieżność. Funkcja  $\mathcal{S}_j \left[ (y_i - \hat{\alpha} - \sum_{k \neq j} \hat{f}_k(x_{ik}))_1^N \right]$  jest jednowymiarowym sześciennym splajnem o  $N$  węzłach. W jej miejsce można przyjąć również inne funkcje, takie jak: jednowymiarowe lokalne regresje wielomianowe (ang. *LOESS - locally estimated scatterplot smoothing*), regresje liniowe, wielomianowe.

**Przykład 11.1.** Dla zilustrowania zasady działania uogólnionych modeli addytywnych przeprowadzimy analizę stężenia ozonu ( $O_3$ ) w zależności od wybranych parametrów meteorologicznych. Do zbudowania modelu GAM wykorzystamy funkcję `gam` pakietu `mgcv` (Wood 2003).

```
library(faraway)
head(ozone)
```

```
##      03      vh wind humidity temp  ibh dpb ibt vis doy
## 1    3 5710     4      28   40 2693 -25  87 250  33
## 2    5 5700     3      37   45  590 -24 128 100  34
## 3    5 5760     3      51   54 1450  25 139  60  35
## 4    6 5720     4      69   35 1568  15 121  60  36
## 5    4 5790     6      19   45 2631 -33 123 100  37
## 6    4 5790     3      25   55  554 -28 182 250  38
```

```
library(mgcv)
mod.gam <- gam(03~s(temp, bs = "cr", m = 2)+s(ibh)+s(ibt), data = ozone)
summary(mod.gam)
```

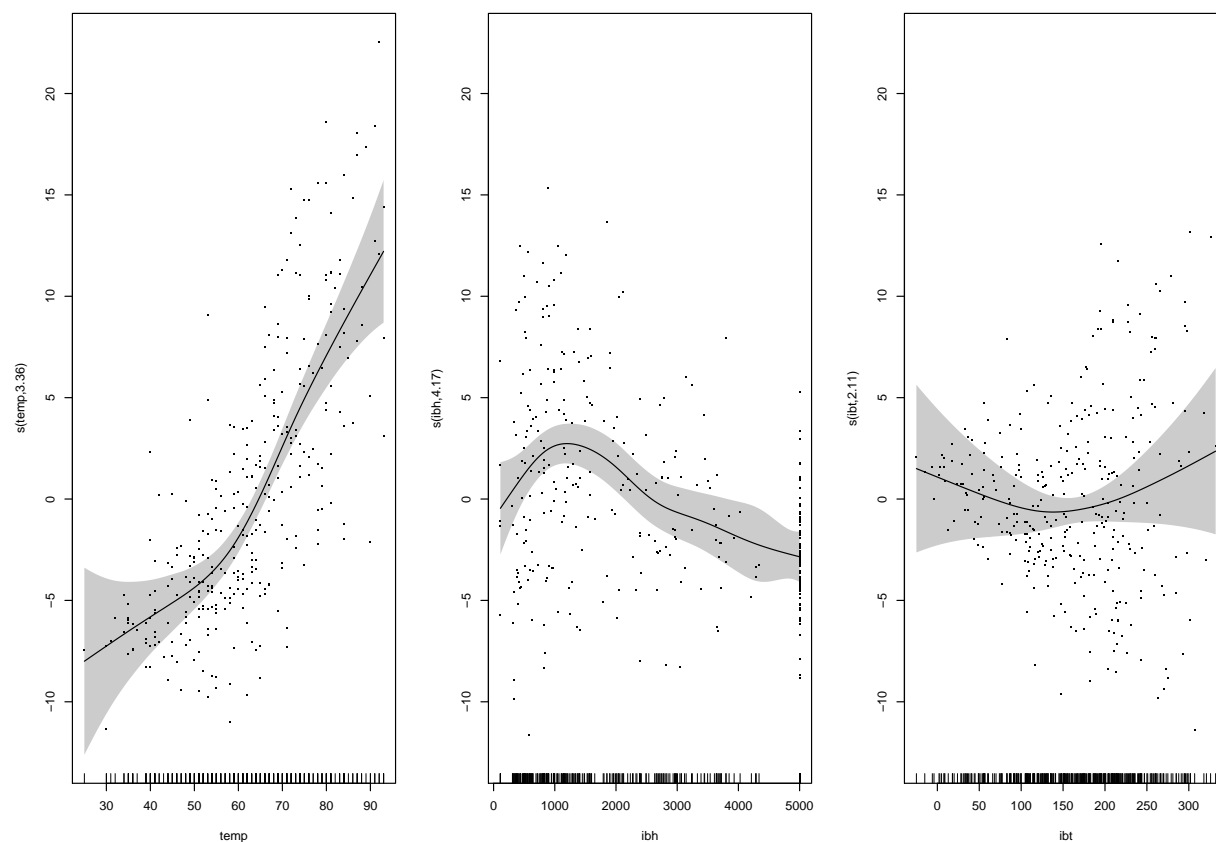
```
##
## Family: gaussian
## Link function: identity
##
## Formula:
```

```
## 03 ~ s(temp, bs = "cr", m = 2) + s(ibh) + s(ibt)
##
## Parametric coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  11.7758     0.2382  49.44  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Approximate significance of smooth terms:
##             edf Ref.df      F  p-value
## s(temp)  3.357  4.216 20.758 5.99e-16 ***
## s(ibh)   4.171  5.072  7.344 1.37e-06 ***
## s(ibt)   2.111  2.729  1.403  0.213
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## R-sq.(adj) =  0.708   Deviance explained = 71.7%
## GCV = 19.348   Scale est. = 18.724     n = 330
```

W powyższym modelu użyto splajnów jako funkcji  $f_i$ . W przypadku zmiennej `temp` był to sześcienny splajn z regularyzacją w postaci ciągłości drugiej pochodnej, a pozostałe są prostymi splajnami. Dopasowanie modelu sięga 71.7% a wartość uogólnionego sprawdzianu krzyżowego 19.35.

Poniższy wykres pokazuje rodzaje transformacji użyte przy dopasowaniu modelu.

```
par(mfrow = c(1,3))
plot(mod.gam, shade = T, residuals = T)
```



```
par(mfrow=c(1,1))
```

Allaire, JJ, Yihui Xie, Jonathan McPherson, Javier Luraschi, Kevin Ushey, Aron Atkins, Hadley Wickham, Joe Cheng, Winston Chang, and Richard Iannone. 2018. *Rmarkdown: Dynamic Documents for R*. <https://CRAN.R-project.org/package=rmarkdown>.

Breiman, Leo. 1996. “Bagging Predictors.” *Machine Learning* 24 (2): 123–40. <https://doi.org/10.1007/BF00058655>.

———. 1998. “Arcing Classifier (with Discussion and a Rejoinder by the Author).” *Ann. Statist.* 26 (3). The Institute of Mathematical Statistics: 801–49. <https://doi.org/10.1214/aos/1024691079>.

Chan, Chung-hong, and Thomas J. Leeper. 2018. *Rio: A Swiss-Army Knife for Data I/O*. <https://CRAN.R-project.org/package=rio>.

Firke, Sam. 2018. *Janitor: Simple Tools for Examining and Cleaning Dirty Data*. <https://CRAN.R-project.org/package=janitor>.

Fisher, R. A. 1936. “The Use of Multiple Measurements in Taxonomic Problems.” *Annals of Eugenics* 7 (2): 179–88. <https://doi.org/10.1111/j.1469-1809.1936.tb02137.x>.

Friedman, Jerome H. 1989. “Regularized Discriminant Analysis.” *Journal of the American Statistical Association* 84 (405): 165–75. <https://doi.org/10.2307/2289860>.

Greenwell, Brandon, Bradley Boehmke, Jay Cunningham, and GBM Developers. 2019. *Gbm: Generalized Boosted Regression Models*. <https://CRAN.R-project.org/package=gbm>.

Ho, Tin Kam. 1995. “Random Decision Forests.” In *Proceedings of 3rd International Conference on Document Analysis and Recognition*, 1:278–82. IEEE.

- Hornik, Kurt, Christian Buchta, and Achim Zeileis. 2009. "Open-Source Machine Learning: R Meets Weka." *Computational Statistics* 24 (2): 225–32. <https://doi.org/10.1007/s00180-008-0119-7>.
- Hothorn, Torsten, Kurt Hornik, and Achim Zeileis. 2006. "Unbiased Recursive Partitioning: A Conditional Inference Framework." *Journal of Computational and Graphical Statistics* 15 (3). Taylor & Francis: 651–74. <https://doi.org/10.1198/106186006X133933>.
- Hothorn, Torsten, and Achim Zeileis. 2015. "Partykit: A Modular Toolkit for Recursive Partytioning in R." *Journal of Machine Learning Research* 16: 3905–9. <http://jmlr.org/papers/v16/hothorn15a.html>.
- Jed Wing, Max Kuhn. Contributions from, Steve Weston, Andre Williams, Chris Keefer, Allan Engelhardt, Tony Cooper, Zachary Mayer, et al. 2018. *Caret: Classification and Regression Training*. <https://CRAN.R-project.org/package=caret>.
- Kavakiotis, Ioannis, Olga Tsave, Athanasios Salifoglou, Nicos Maglaveras, Ioannis Vlahavas, and Ioanna Chouvarda. 2017. "Machine Learning and Data Mining Methods in Diabetes Research." *Computational and Structural Biotechnology Journal* 15: 104–16.
- Kearns, M., and L. G. Valiant. 1989. "Cryptographic Limitations on Learning Boolean Formulae and Finite Automata." *Annual ACM Symposium on Theory of Computing*, 433. <http://search.ebscohost.com/login.aspx?direct=true&db=edb&AN=73725380&lang=pl&site=eds-live&scope=site>.
- Kuhn, Max, and Ross Quinlan. 2018. *C50: C5.0 Decision Trees and Rule-Based Models*. <https://CRAN.R-project.org/package=C50>.
- Liaw, Andy, and Matthew Wiener. 2002. "Classification and Regression by randomForest." *R News* 2 (3): 18–22. <https://CRAN.R-project.org/doc/Rnews/>.
- Meyer, David, Evgenia Dimitriadou, Kurt Hornik, Andreas Weingessel, and Friedrich Leisch. 2019. *E1071: Misc Functions of the Department of Statistics, Probability Theory Group (Formerly: E1071), Tu Wien*. <https://CRAN.R-project.org/package=e1071>.
- Peters, Andrea, and Torsten Hothorn. 2018. *Ipred: Improved Predictors*. <https://CRAN.R-project.org/package=ipred>.
- Quinlan, J Ross. 1993. *C4. 5: Programs for Machine Learning*. Morgan Kaufmann.
- R Core Team. 2018. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Schapire, Robert E. 1990. "The Strength of Weak Learnability." *Machine Learning* 5 (2): 197–227. <https://doi.org/10.1007/BF00116037>.
- Team, The FoRt Student Project. 2015. *CHAID: CHi-Squared Automated Interaction Detection*.
- Templ, Matthias, Alexander Kowarik, Andreas Alfons, and Bernd Prantner. 2019. *VIM: Visualization and Imputation of Missing Values*. <https://CRAN.R-project.org/package=VIM>.
- Therneau, Terry, and Beth Atkinson. 2018. *Rpart: Recursive Partitioning and Regression Trees*. <https://CRAN.R-project.org/package=rpart>.
- Torgo, Luis. 2013. *DMwR: Functions and Data for "Data Mining with R"*. <https://CRAN.R-project.org/package=DMwR>.
- Trevor Hastie, S original by, Robert Tibshirani. Original R port by Friedrich Leisch, Kurt Hornik, and Brian D. Ripley. 2017. *Mda: Mixture and Flexible Discriminant Analysis*. <https://CRAN.R-project.org/package=mda>.
- van Buuren, Stef, and Karin Groothuis-Oudshoorn. 2018. *Mice: Multivariate Imputation by Chained Equations*. <https://CRAN.R-project.org/package=mice>.
- Venables, W. N., and B. D. Ripley. 2002. *Modern Applied Statistics with S*. Fourth. New York: Springer. <http://www.stats.ox.ac.uk/pub/MASS4>.



Weihs, Claus, Uwe Ligges, Karsten Luebke, and Nils Raabe. 2005. “KlaR Analyzing German Business Cycles.” In *Data Analysis and Decision Support*, edited by D. Baier, R. Decker, and L. Schmidt-Thieme, 335–43. Berlin: Springer-Verlag.

Welch, B. L. 1939. “Note on Discriminant Functions.” *Biometrika* 31 (1/2): 218–20. <https://doi.org/10.2307/2334985>.

Wood, S. N. 2003. “Thin-Plate Regression Splines.” *Journal of the Royal Statistical Society (B)* 65 (1): 95–114.

Xie, Yihui. 2018a. *Bookdown: Authoring Books and Technical Documents with R Markdown*. <https://CRAN.R-project.org/package=bookdown>.

———. 2018b. *Knitr: A General-Purpose Package for Dynamic Report Generation in R*. <https://CRAN.R-project.org/package=knitr>.