

Formatowanie spamu Data Mining

Michał Maj i Anna Mieszkalska
album 256556 i 255699

21 kwietnia 2023

Spis treści

1 Analiza opisowa i wizualizacja.

2

1 Analiza opisowa i wizualizacja.

W naszym projekcie będziemy analizować dane o nazwie *Spambase* z biblioteki *kernlab*. Zestaw danych *spambase* jest zbiorem wiadomości e-mail, które zostały przeanalizowane i sklasyfikowane jako spam lub non-spam. Celem tego zbioru danych jest dostarczenie użytecznych materiałów potrzebnych do analiz i eksploracji w tym zakresie. W tym projekcie użyjemy różnych metod i technik pozyskiwania wiedzy, aby przeanalizować dane *spambase* w celu opracowania modelu klasyfikującego wiadomości e-mail jako spam lub non-spam. Modele opracowane w tym projekcie mogą być przydatne w rzeczywistych serwisach poczt e-mailowych, gdzie problem dostarczania niechcianych wiadomości jest nam powszechnie znany.

```
#format

#Pierwsze 48 zmiennych zawiera częstotliwość występowania nazwy zmiennej (np. biznes)
#oznacza to częstotliwość odpowiadającej jej liczby (np. 650). Zmienne 49-54 wskazują
#Zmienne 55-57 zawierają średnią, najdłuższą i całkowitą długość wielkich liter Zmiennych

#dane
data(spam)

#pierwsze 10 wyników
head(spam)
```

##	make	address	all	num3d	our	over	remove	internet	order	mail	receive	will
## 1	0.00	0.64	0.64	0	0.32	0.00	0.00	0.00	0.00	0.00	0.00	0.64
## 2	0.21	0.28	0.50	0	0.14	0.28	0.21	0.07	0.00	0.94	0.21	0.79
## 3	0.06	0.00	0.71	0	1.23	0.19	0.19	0.12	0.64	0.25	0.38	0.45
## 4	0.00	0.00	0.00	0	0.63	0.00	0.31	0.63	0.31	0.63	0.31	0.31
## 5	0.00	0.00	0.00	0	0.63	0.00	0.31	0.63	0.31	0.63	0.31	0.31
## 6	0.00	0.00	0.00	0	1.85	0.00	0.00	1.85	0.00	0.00	0.00	0.00
##	people	report	addresses	free	business	email	you	credit	your	font	num000	
## 1	0.00	0.00	0.00	0.32	0.00	1.29	1.93	0.00	0.96	0	0.00	
## 2	0.65	0.21	0.14	0.14	0.07	0.28	3.47	0.00	1.59	0	0.43	
## 3	0.12	0.00	1.75	0.06	0.06	1.03	1.36	0.32	0.51	0	1.16	
## 4	0.31	0.00	0.00	0.31	0.00	0.00	3.18	0.00	0.31	0	0.00	
## 5	0.31	0.00	0.00	0.31	0.00	0.00	3.18	0.00	0.31	0	0.00	
## 6	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0	0.00	
##	money	hp	hpl	george	num650	lab	labs	telnet	num857	data	num415	num85
## 1	0.00	0	0	0	0	0	0	0	0	0	0	0
## 2	0.43	0	0	0	0	0	0	0	0	0	0	0
## 3	0.06	0	0	0	0	0	0	0	0	0	0	0
## 4	0.00	0	0	0	0	0	0	0	0	0	0	0
## 5	0.00	0	0	0	0	0	0	0	0	0	0	0
## 6	0.00	0	0	0	0	0	0	0	0	0	0	0
##	technology	num1999	parts	pm	direct	cs	meeting	original	project	re	edu	
## 1	0	0.00	0	0	0.00	0	0	0.00	0	0.00	0.00	
## 2	0	0.07	0	0	0.00	0	0	0.00	0	0.00	0.00	

```
## 3      0      0.00      0 0      0.06 0      0      0.12      0 0.06 0.06
## 4      0      0.00      0 0      0.00 0      0      0.00      0 0.00 0.00
## 5      0      0.00      0 0      0.00 0      0      0.00      0 0.00 0.00
## 6      0      0.00      0 0      0.00 0      0      0.00      0 0.00 0.00
##      table conference charSemicolon charRoundbracket charSquarebracket
## 1      0      0      0.00      0.000      0
## 2      0      0      0.00      0.132      0
## 3      0      0      0.01      0.143      0
## 4      0      0      0.00      0.137      0
## 5      0      0      0.00      0.135      0
## 6      0      0      0.00      0.223      0
##      charExclamation charDollar charHash capitalAve capitalLong capitalTotal type
## 1      0.778      0.000      0.000      3.756      61      278 spam
## 2      0.372      0.180      0.048      5.114      101      1028 spam
## 3      0.276      0.184      0.010      9.821      485      2259 spam
## 4      0.137      0.000      0.000      3.537      40      191 spam
## 5      0.135      0.000      0.000      3.537      40      191 spam
## 6      0.000      0.000      0.000      3.000      15      54 spam

#liczba maili
nrow(spam)

## [1] 4601

View(spam)
#ile spamów i non-spamów
table(spam$type)

##
## nonspam      spam
##      2788      1813
```

Tabela 1: Database - pierwsze 14 rekordów.

make	address	all	num3d	our	over	remove	internet	order	mail	charRoundbracket	type
0.00	0.64	0.64	0	0.32	0.00	0.00	0.00	0.00	0.00	0.000	spam
0.21	0.28	0.50	0	0.14	0.28	0.21	0.07	0.00	0.94	0.132	spam
0.06	0.00	0.71	0	1.23	0.19	0.19	0.12	0.64	0.25	0.143	spam
0.00	0.00	0.00	0	0.63	0.00	0.31	0.63	0.31	0.63	0.137	spam
0.00	0.00	0.00	0	0.63	0.00	0.31	0.63	0.31	0.63	0.135	spam
0.00	0.00	0.00	0	1.85	0.00	0.00	1.85	0.00	0.00	0.223	spam

Tabela 1

```
summary(spam)

##      make      address      all      num3d
## Min.   :0.0000   Min.   : 0.000   Min.   :0.0000   Min.   : 0.00000
## 1st Qu.:0.0000   1st Qu.: 0.000   1st Qu.:0.0000   1st Qu.: 0.00000
```

##	Median :0.0000	Median : 0.000	Median :0.0000	Median : 0.00000
##	Mean :0.1046	Mean : 0.213	Mean :0.2807	Mean : 0.06542
##	3rd Qu.:0.0000	3rd Qu.: 0.000	3rd Qu.:0.4200	3rd Qu.: 0.00000
##	Max. :4.5400	Max. :14.280	Max. :5.1000	Max. :42.81000
##	our	over	remove	internet
##	Min. : 0.0000	Min. :0.0000	Min. :0.0000	Min. : 0.0000
##	1st Qu.: 0.0000	1st Qu.:0.0000	1st Qu.:0.0000	1st Qu.: 0.0000
##	Median : 0.0000	Median :0.0000	Median :0.0000	Median : 0.0000
##	Mean : 0.3122	Mean :0.0959	Mean :0.1142	Mean : 0.1053
##	3rd Qu.: 0.3800	3rd Qu.:0.0000	3rd Qu.:0.0000	3rd Qu.: 0.0000
##	Max. :10.0000	Max. :5.8800	Max. :7.2700	Max. :11.1100
##	order	mail	receive	will
##	Min. :0.00000	Min. : 0.0000	Min. :0.00000	Min. :0.0000
##	1st Qu.:0.00000	1st Qu.: 0.0000	1st Qu.:0.00000	1st Qu.:0.0000
##	Median :0.00000	Median : 0.0000	Median :0.00000	Median :0.1000
##	Mean :0.09007	Mean : 0.2394	Mean :0.05982	Mean :0.5417
##	3rd Qu.:0.00000	3rd Qu.: 0.1600	3rd Qu.:0.00000	3rd Qu.:0.8000
##	Max. :5.26000	Max. :18.1800	Max. :2.61000	Max. :9.6700
##	people	report	addresses	free
##	Min. :0.00000	Min. : 0.00000	Min. :0.0000	Min. : 0.0000
##	1st Qu.:0.00000	1st Qu.: 0.00000	1st Qu.:0.0000	1st Qu.: 0.0000
##	Median :0.00000	Median : 0.00000	Median :0.0000	Median : 0.0000
##	Mean :0.09393	Mean : 0.05863	Mean :0.0492	Mean : 0.2488
##	3rd Qu.:0.00000	3rd Qu.: 0.00000	3rd Qu.:0.0000	3rd Qu.: 0.1000
##	Max. :5.55000	Max. :10.00000	Max. :4.4100	Max. :20.0000
##	business	email	you	credit
##	Min. :0.0000	Min. :0.0000	Min. : 0.000	Min. : 0.00000
##	1st Qu.:0.0000	1st Qu.:0.0000	1st Qu.: 0.000	1st Qu.: 0.00000
##	Median :0.0000	Median :0.0000	Median : 1.310	Median : 0.00000
##	Mean :0.1426	Mean :0.1847	Mean : 1.662	Mean : 0.08558
##	3rd Qu.:0.0000	3rd Qu.:0.0000	3rd Qu.: 2.640	3rd Qu.: 0.00000
##	Max. :7.1400	Max. :9.0900	Max. :18.750	Max. :18.18000
##	your	font	num000	money
##	Min. : 0.0000	Min. : 0.0000	Min. :0.0000	Min. : 0.00000
##	1st Qu.: 0.0000	1st Qu.: 0.0000	1st Qu.:0.0000	1st Qu.: 0.00000
##	Median : 0.2200	Median : 0.0000	Median :0.0000	Median : 0.00000
##	Mean : 0.8098	Mean : 0.1212	Mean :0.1016	Mean : 0.09427
##	3rd Qu.: 1.2700	3rd Qu.: 0.0000	3rd Qu.:0.0000	3rd Qu.: 0.00000
##	Max. :11.1100	Max. :17.1000	Max. :5.4500	Max. :12.50000
##	hp	hpl	george	num650
##	Min. : 0.0000	Min. : 0.0000	Min. : 0.0000	Min. :0.0000
##	1st Qu.: 0.0000	1st Qu.: 0.0000	1st Qu.: 0.0000	1st Qu.:0.0000
##	Median : 0.0000	Median : 0.0000	Median : 0.0000	Median :0.0000
##	Mean : 0.5495	Mean : 0.2654	Mean : 0.7673	Mean :0.1248
##	3rd Qu.: 0.0000	3rd Qu.: 0.0000	3rd Qu.: 0.0000	3rd Qu.:0.0000
##	Max. :20.8300	Max. :16.6600	Max. :33.3300	Max. :9.0900
##	lab	labs	telnet	num857
##	Min. : 0.00000	Min. :0.0000	Min. : 0.00000	Min. :0.00000

##	1st Qu.: 0.00000	1st Qu.:0.0000	1st Qu.: 0.00000	1st Qu.:0.00000
##	Median : 0.00000	Median :0.0000	Median : 0.00000	Median :0.00000
##	Mean : 0.09892	Mean :0.1029	Mean : 0.06475	Mean :0.04705
##	3rd Qu.: 0.00000	3rd Qu.:0.0000	3rd Qu.: 0.00000	3rd Qu.:0.00000
##	Max. :14.28000	Max. :5.8800	Max. :12.50000	Max. :4.76000
##	data	num415	num85	technology
##	Min. : 0.00000	Min. :0.00000	Min. : 0.0000	Min. :0.00000
##	1st Qu.: 0.00000	1st Qu.:0.00000	1st Qu.: 0.0000	1st Qu.:0.00000
##	Median : 0.00000	Median :0.00000	Median : 0.0000	Median :0.00000
##	Mean : 0.09723	Mean :0.04784	Mean : 0.1054	Mean :0.09748
##	3rd Qu.: 0.00000	3rd Qu.:0.00000	3rd Qu.: 0.0000	3rd Qu.:0.00000
##	Max. :18.18000	Max. :4.76000	Max. :20.0000	Max. :7.69000
##	num1999	parts	pm	direct
##	Min. :0.000	Min. :0.0000	Min. : 0.00000	Min. :0.00000
##	1st Qu.:0.000	1st Qu.:0.0000	1st Qu.: 0.00000	1st Qu.:0.00000
##	Median :0.000	Median :0.0000	Median : 0.00000	Median :0.00000
##	Mean :0.137	Mean :0.0132	Mean : 0.07863	Mean :0.06483
##	3rd Qu.:0.000	3rd Qu.:0.0000	3rd Qu.: 0.00000	3rd Qu.:0.00000
##	Max. :6.890	Max. :8.3300	Max. :11.11000	Max. :4.76000
##	cs	meeting	original	project
##	Min. :0.00000	Min. : 0.0000	Min. :0.0000	Min. : 0.0000
##	1st Qu.:0.00000	1st Qu.: 0.0000	1st Qu.:0.0000	1st Qu.: 0.0000
##	Median :0.00000	Median : 0.0000	Median :0.0000	Median : 0.0000
##	Mean :0.04367	Mean : 0.1323	Mean :0.0461	Mean : 0.0792
##	3rd Qu.:0.00000	3rd Qu.: 0.0000	3rd Qu.:0.0000	3rd Qu.: 0.0000
##	Max. :7.14000	Max. :14.2800	Max. :3.5700	Max. :20.0000
##	re	edu	table	conference
##	Min. : 0.0000	Min. : 0.0000	Min. :0.000000	Min. : 0.00000
##	1st Qu.: 0.0000	1st Qu.: 0.0000	1st Qu.:0.000000	1st Qu.: 0.00000
##	Median : 0.0000	Median : 0.0000	Median :0.000000	Median : 0.00000
##	Mean : 0.3012	Mean : 0.1798	Mean :0.005444	Mean : 0.03187
##	3rd Qu.: 0.1100	3rd Qu.: 0.0000	3rd Qu.:0.000000	3rd Qu.: 0.00000
##	Max. :21.4200	Max. :22.0500	Max. :2.170000	Max. :10.00000
##	charSemicolon	charRoundbracket	charSquarebracket	charExclamation
##	Min. :0.00000	Min. :0.000	Min. :0.00000	Min. : 0.0000
##	1st Qu.:0.00000	1st Qu.:0.000	1st Qu.:0.00000	1st Qu.: 0.0000
##	Median :0.00000	Median :0.065	Median :0.00000	Median : 0.0000
##	Mean :0.03857	Mean :0.139	Mean :0.01698	Mean : 0.2691
##	3rd Qu.:0.00000	3rd Qu.:0.188	3rd Qu.:0.00000	3rd Qu.: 0.3150
##	Max. :4.38500	Max. :9.752	Max. :4.08100	Max. :32.4780
##	charDollar	charHash	capitalAve	capitalLong
##	Min. :0.00000	Min. : 0.00000	Min. : 1.000	Min. : 1.00
##	1st Qu.:0.00000	1st Qu.: 0.00000	1st Qu.: 1.588	1st Qu.: 6.00
##	Median :0.00000	Median : 0.00000	Median : 2.276	Median : 15.00
##	Mean :0.07581	Mean : 0.04424	Mean : 5.191	Mean : 52.17
##	3rd Qu.:0.05200	3rd Qu.: 0.00000	3rd Qu.: 3.706	3rd Qu.: 43.00
##	Max. :6.00300	Max. :19.82900	Max. :1102.500	Max. :9989.00
##	capitalTotal	type		

```
## Min.    :    1.0   nonspam:2788
## 1st Qu.:   35.0   spam    :1813
## Median :   95.0
## Mean    :  283.3
## 3rd Qu.:  266.0
## Max.    :15841.0
```