

Projet 1/5 Traitement-Des-Donnees

Nicolas SALVAN - Alexandre CORRIOU

2024-05-17

Lecture des données

Notre jeu de données contient des informations sur des patientes atteintes d'un cancer du sein. Nous allons commencer par lire les données et les afficher pour mieux les comprendre. Ce fichier contient le code pour *pré-traiter les données*.

Importation du dataset

```
data <- read.csv("data/breast-cancer.csv", header = TRUE, sep = ",")
```

Affichage des premières lignes

```
head(data)
```

```
##           id diagnosis radius_mean texture_mean perimeter_mean area_mean
## 1    842302          M      17.99       10.38         122.80      1001.0
## 2    842517          M      20.57       17.77         132.90      1326.0
## 3  84300903          M      19.69       21.25         130.00      1203.0
## 4  84348301          M      11.42       20.38          77.58       386.1
## 5  84358402          M      20.29       14.34         135.10      1297.0
## 6    843786          M      12.45       15.70          82.57       477.1
## smoothness_mean compactness_mean concavity_mean concave.points_mean
## 1         0.11840         0.27760         0.3001         0.14710
## 2         0.08474         0.07864         0.0869         0.07017
## 3         0.10960         0.15990         0.1974         0.12790
## 4         0.14250         0.28390         0.2414         0.10520
## 5         0.10030         0.13280         0.1980         0.10430
## 6         0.12780         0.17000         0.1578         0.08089
## symmetry_mean fractal_dimension_mean radius_se texture_se perimeter_se
## 1         0.2419         0.07871      1.0950      0.9053         8.589
## 2         0.1812         0.05667      0.5435      0.7339         3.398
## 3         0.2069         0.05999      0.7456      0.7869         4.585
## 4         0.2597         0.09744      0.4956      1.1560         3.445
## 5         0.1809         0.05883      0.7572      0.7813         5.438
## 6         0.2087         0.07613      0.3345      0.8902         2.217
## area_se smoothness_se compactness_se concavity_se concave.points_se
```

```
## 1 153.40      0.006399      0.04904      0.05373      0.01587
## 2  74.08      0.005225      0.01308      0.01860      0.01340
## 3  94.03      0.006150      0.04006      0.03832      0.02058
## 4  27.23      0.009110      0.07458      0.05661      0.01867
## 5  94.44      0.011490      0.02461      0.05688      0.01885
## 6  27.19      0.007510      0.03345      0.03672      0.01137
##      symmetry_se fractal_dimension_se radius_worst texture_worst perimeter_worst
## 1      0.03003              0.006193          25.38          17.33          184.60
## 2      0.01389              0.003532          24.99          23.41          158.80
## 3      0.02250              0.004571          23.57          25.53          152.50
## 4      0.05963              0.009208          14.91          26.50           98.87
## 5      0.01756              0.005115          22.54          16.67          152.20
## 6      0.02165              0.005082          15.47          23.75          103.40
##      area_worst smoothness_worst compactness_worst concavity_worst
## 1      2019.0              0.1622              0.6656              0.7119
## 2      1956.0              0.1238              0.1866              0.2416
## 3      1709.0              0.1444              0.4245              0.4504
## 4       567.7              0.2098              0.8663              0.6869
## 5      1575.0              0.1374              0.2050              0.4000
## 6       741.6              0.1791              0.5249              0.5355
##      concave.points_worst symmetry_worst fractal_dimension_worst
## 1              0.2654              0.4601              0.11890
## 2              0.1860              0.2750              0.08902
## 3              0.2430              0.3613              0.08758
## 4              0.2575              0.6638              0.17300
## 5              0.1625              0.2364              0.07678
## 6              0.1741              0.3985              0.12440
```

On observe qu'il y a une variable qualitative "diagnosis" qui correspond au diagnostic de la patiente. Toutes les autres variables sont quantitatives, et décrivent les caractéristiques du cancer.

Affichage des dimensions

```
dim(data)
```

```
## [1] 569 32
```

Notre jeu de données contient 569 observations et 32 variables.

Affichage des types des variables

```
str(data)
```

```
## 'data.frame': 569 obs. of 32 variables:
## $ id : int 842302 842517 84300903 84348301 84358402 843786 844359 84458202 844...
## $ diagnosis : chr "M" "M" "M" "M" ...
## $ radius_mean : num 18 20.6 19.7 11.4 20.3 ...
## $ texture_mean : num 10.4 17.8 21.2 20.4 14.3 ...
## $ perimeter_mean : num 122.8 132.9 130 77.6 135.1 ...
```

```
## $ area_mean : num 1001 1326 1203 386 1297 ...
## $ smoothness_mean : num 0.1184 0.0847 0.1096 0.1425 0.1003 ...
## $ compactness_mean : num 0.2776 0.0786 0.1599 0.2839 0.1328 ...
## $ concavity_mean : num 0.3001 0.0869 0.1974 0.2414 0.198 ...
## $ concave.points_mean : num 0.1471 0.0702 0.1279 0.1052 0.1043 ...
## $ symmetry_mean : num 0.242 0.181 0.207 0.26 0.181 ...
## $ fractal_dimension_mean : num 0.0787 0.0567 0.06 0.0974 0.0588 ...
## $ radius_se : num 1.095 0.543 0.746 0.496 0.757 ...
## $ texture_se : num 0.905 0.734 0.787 1.156 0.781 ...
## $ perimeter_se : num 8.59 3.4 4.58 3.44 5.44 ...
## $ area_se : num 153.4 74.1 94 27.2 94.4 ...
## $ smoothness_se : num 0.0064 0.00522 0.00615 0.00911 0.01149 ...
## $ compactness_se : num 0.049 0.0131 0.0401 0.0746 0.0246 ...
## $ concavity_se : num 0.0537 0.0186 0.0383 0.0566 0.0569 ...
## $ concave.points_se : num 0.0159 0.0134 0.0206 0.0187 0.0188 ...
## $ symmetry_se : num 0.03 0.0139 0.0225 0.0596 0.0176 ...
## $ fractal_dimension_se : num 0.00619 0.00353 0.00457 0.00921 0.00511 ...
## $ radius_worst : num 25.4 25 23.6 14.9 22.5 ...
## $ texture_worst : num 17.3 23.4 25.5 26.5 16.7 ...
## $ perimeter_worst : num 184.6 158.8 152.5 98.9 152.2 ...
## $ area_worst : num 2019 1956 1709 568 1575 ...
## $ smoothness_worst : num 0.162 0.124 0.144 0.21 0.137 ...
## $ compactness_worst : num 0.666 0.187 0.424 0.866 0.205 ...
## $ concavity_worst : num 0.712 0.242 0.45 0.687 0.4 ...
## $ concave.points_worst : num 0.265 0.186 0.243 0.258 0.163 ...
## $ symmetry_worst : num 0.46 0.275 0.361 0.664 0.236 ...
## $ fractal_dimension_worst : num 0.1189 0.089 0.0876 0.173 0.0768 ...
```

Conversion des données qualitatives en factor

On observe que la variable “diagnosis” est de type “chr”. Nous allons la convertir en facteur pour faciliter l’analyse.

```
data$diagnosis <- as.factor(data$diagnosis)
str(data)
```

```
## 'data.frame': 569 obs. of 32 variables:
## $ id : int 842302 842517 84300903 84348301 84358402 843786 844359 84458202 844...
## $ diagnosis : Factor w/ 2 levels "B","M": 2 2 2 2 2 2 2 2 2 ...
## $ radius_mean : num 18 20.6 19.7 11.4 20.3 ...
## $ texture_mean : num 10.4 17.8 21.2 20.4 14.3 ...
## $ perimeter_mean : num 122.8 132.9 130 77.6 135.1 ...
## $ area_mean : num 1001 1326 1203 386 1297 ...
## $ smoothness_mean : num 0.1184 0.0847 0.1096 0.1425 0.1003 ...
## $ compactness_mean : num 0.2776 0.0786 0.1599 0.2839 0.1328 ...
## $ concavity_mean : num 0.3001 0.0869 0.1974 0.2414 0.198 ...
## $ concave.points_mean : num 0.1471 0.0702 0.1279 0.1052 0.1043 ...
## $ symmetry_mean : num 0.242 0.181 0.207 0.26 0.181 ...
## $ fractal_dimension_mean : num 0.0787 0.0567 0.06 0.0974 0.0588 ...
## $ radius_se : num 1.095 0.543 0.746 0.496 0.757 ...
## $ texture_se : num 0.905 0.734 0.787 1.156 0.781 ...
## $ perimeter_se : num 8.59 3.4 4.58 3.44 5.44 ...
## $ area_se : num 153.4 74.1 94 27.2 94.4 ...
```

```
## $ smoothness_se      : num  0.0064 0.00522 0.00615 0.00911 0.01149 ...
## $ compactness_se     : num  0.049 0.0131 0.0401 0.0746 0.0246 ...
## $ concavity_se       : num  0.0537 0.0186 0.0383 0.0566 0.0569 ...
## $ concave.points_se  : num  0.0159 0.0134 0.0206 0.0187 0.0188 ...
## $ symmetry_se        : num  0.03 0.0139 0.0225 0.0596 0.0176 ...
## $ fractal_dimension_se : num  0.00619 0.00353 0.00457 0.00921 0.00511 ...
## $ radius_worst       : num  25.4 25 23.6 14.9 22.5 ...
## $ texture_worst      : num  17.3 23.4 25.5 26.5 16.7 ...
## $ perimeter_worst    : num  184.6 158.8 152.5 98.9 152.2 ...
## $ area_worst         : num  2019 1956 1709 568 1575 ...
## $ smoothness_worst   : num  0.162 0.124 0.144 0.21 0.137 ...
## $ compactness_worst  : num  0.666 0.187 0.424 0.866 0.205 ...
## $ concavity_worst    : num  0.712 0.242 0.45 0.687 0.4 ...
## $ concave.points_worst : num  0.265 0.186 0.243 0.258 0.163 ...
## $ symmetry_worst     : num  0.46 0.275 0.361 0.664 0.236 ...
## $ fractal_dimension_worst: num  0.1189 0.089 0.0876 0.173 0.0768 ...
```

Tous les types de variables semblent corrects.

Nettoyage

Certaines variables sont inutilisables, comme l'identifiant de la patiente. Nous allons les supprimer. Il nous faut également supprimer les NaNs pour éviter les erreurs dans les analyses.

```
# suppression des NaNs
data <- na.omit(data)

# suppression des colonnes inutiles : identifiant de la patiente
data <- data[,-c(1)]
head(data)
```

```
##      diagnosis radius_mean texture_mean perimeter_mean area_mean smoothness_mean
## 1           M      17.99      10.38      122.80      1001.0      0.11840
## 2           M      20.57      17.77      132.90      1326.0      0.08474
## 3           M      19.69      21.25      130.00      1203.0      0.10960
## 4           M      11.42      20.38       77.58       386.1      0.14250
## 5           M      20.29      14.34      135.10      1297.0      0.10030
## 6           M      12.45      15.70       82.57       477.1      0.12780
## compactness_mean concavity_mean concave.points_mean symmetry_mean
## 1      0.27760      0.3001      0.14710      0.2419
## 2      0.07864      0.0869      0.07017      0.1812
## 3      0.15990      0.1974      0.12790      0.2069
## 4      0.28390      0.2414      0.10520      0.2597
## 5      0.13280      0.1980      0.10430      0.1809
## 6      0.17000      0.1578      0.08089      0.2087
## fractal_dimension_mean radius_se texture_se perimeter_se area_se
## 1      0.07871      1.0950      0.9053      8.589 153.40
## 2      0.05667      0.5435      0.7339      3.398 74.08
## 3      0.05999      0.7456      0.7869      4.585 94.03
## 4      0.09744      0.4956      1.1560      3.445 27.23
## 5      0.05883      0.7572      0.7813      5.438 94.44
## 6      0.07613      0.3345      0.8902      2.217 27.19
## smoothness_se compactness_se concavity_se concave.points_se symmetry_se
```

```
## 1      0.006399      0.04904      0.05373      0.01587      0.03003
## 2      0.005225      0.01308      0.01860      0.01340      0.01389
## 3      0.006150      0.04006      0.03832      0.02058      0.02250
## 4      0.009110      0.07458      0.05661      0.01867      0.05963
## 5      0.011490      0.02461      0.05688      0.01885      0.01756
## 6      0.007510      0.03345      0.03672      0.01137      0.02165
## fractal_dimension_se radius_worst texture_worst perimeter_worst area_worst
## 1      0.006193      25.38      17.33      184.60      2019.0
## 2      0.003532      24.99      23.41      158.80      1956.0
## 3      0.004571      23.57      25.53      152.50      1709.0
## 4      0.009208      14.91      26.50      98.87      567.7
## 5      0.005115      22.54      16.67      152.20      1575.0
## 6      0.005082      15.47      23.75      103.40      741.6
## smoothness_worst compactness_worst concavity_worst concave.points_worst
## 1      0.1622      0.6656      0.7119      0.2654
## 2      0.1238      0.1866      0.2416      0.1860
## 3      0.1444      0.4245      0.4504      0.2430
## 4      0.2098      0.8663      0.6869      0.2575
## 5      0.1374      0.2050      0.4000      0.1625
## 6      0.1791      0.5249      0.5355      0.1741
## symmetry_worst fractal_dimension_worst
## 1      0.4601      0.11890
## 2      0.2750      0.08902
## 3      0.3613      0.08758
## 4      0.6638      0.17300
## 5      0.2364      0.07678
## 6      0.3985      0.12440
```

Exportation des données

Nous allons exporter les données nettoyées pour les utiliser dans les analyses suivantes.

```
write.csv(data, "data/data_cleaned.csv", row.names = FALSE)
```

Séparation des données

Nous allons séparer les données en deux parties : une partie pour l'apprentissage et une partie pour le test.

```
split_data <- function (data, train_ratio) {
  set.seed(123)
  n <- nrow(data)
  p <- ncol(data)-1
  test_ratio <- 1 - train_ratio
  n.test <- round(n*test_ratio)
  train_index <- sample(1:nrow(data), n.test)
  train_data <- data[-train_index,]
  test_data <- data[train_index,]
  return(list(train_data = train_data, test_data = test_data))
}
```

```
data_split <- split_data(data, 0.8) # 1/5 des données pour le test
train_data <- data_split$train_data
test_data <- data_split$test_data
```

Exportation des données d'apprentissage et de test

```
write.csv(train_data, "data/train_data.csv", row.names = FALSE)
write.csv(test_data, "data/test_data.csv", row.names = FALSE)
```

Il faut noter qu'il faudra convertir la colonne "diagnosis" en facteur dans les données d'apprentissage et de test, mais aussi dans les données cleaned.

Données d'entraînement équilibrées

```
train_data_balanced <- rbind(train_data[train_data$diagnosis == "M",], train_data[train_data$diagnosis == "B",])
table(train_data_balanced$diagnosis)
```

```
##
##      B      M
## 179 179
```

Nous avons maintenant des données d'entraînement équilibrées.

Exportation des données d'apprentissage équilibrées

```
write.csv(train_data_balanced, "data/train_data_balanced.csv", row.names = FALSE)
```