# ADS 503 Final Project

## April Chia

## 2025-06-05

## Packages used

```
library(caret)
library(tidyverse)
library(ggplot2)
library(DataExplorer)
library(reshape2)
library(corrplot)
library(Hmisc)
library(mlbench)
library(e1071)
library(randomForest)
library(gt)
library(pls)
library(elasticnet)
library(pROC)
```

## Import dataset

```
cancer_data <- read.csv("breast-cancer.csv")
head(cancer_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
```

## EDA

```
summary(cancer_data)
```

```
##        id               diagnosis          radius_mean       texture_mean
##  Min.   :     8670   Length:569         Min.   : 6.981   Min.   : 9.71
##  1st Qu.:   869218   Class :character   1st Qu.:11.700   1st Qu.:16.17
##  Median :   906024   Mode  :character   Median :13.370   Median :18.84
##  Mean   : 30371831                      Mean   :14.127   Mean   :19.29
##  3rd Qu.:  8813129                      3rd Qu.:15.780   3rd Qu.:21.80
##  Max.   :911320502                      Max.   :28.110   Max.   :39.28
##  perimeter_mean     area_mean      smoothness_mean   compactness_mean
##  Min.   : 43.79   Min.   : 143.5   Min.   :0.05263   Min.   :0.01938
##  1st Qu.: 75.17   1st Qu.: 420.3   1st Qu.:0.08637   1st Qu.:0.06492
##  Median : 86.24   Median : 551.1   Median :0.09587   Median :0.09263
##  Mean   : 91.97   Mean   : 654.9   Mean   :0.09636   Mean   :0.10434
##  3rd Qu.:104.10   3rd Qu.: 782.7   3rd Qu.:0.10530   3rd Qu.:0.13040
##  Max.   :188.50   Max.   :2501.0   Max.   :0.16340   Max.   :0.34540
##  concavity_mean     concave.points_mean symmetry_mean    fractal_dimension_mean
##  Min.   :0.00000   Min.   :0.00000     Min.   :0.1060   Min.   :0.04996
##  1st Qu.:0.02956   1st Qu.:0.02031     1st Qu.:0.1619   1st Qu.:0.05770
##  Median :0.06154   Median :0.03350     Median :0.1792   Median :0.06154
##  Mean   :0.08880   Mean   :0.04892     Mean   :0.1812   Mean   :0.06280
##  3rd Qu.:0.13070   3rd Qu.:0.07400     3rd Qu.:0.1957   3rd Qu.:0.06612
```

2

```
## Max.    :0.42680  Max.    :0.20120   Max.    :0.3040   Max.    :0.09744
##    radius_se        texture_se       perimeter_se        area_se
## Min.   :0.1115   Min.   :0.3602   Min.   : 0.757   Min.   :  6.802
## 1st Qu.:0.2324   1st Qu.:0.8339   1st Qu.: 1.606   1st Qu.: 17.850
## Median :0.3242   Median :1.1080   Median : 2.287   Median : 24.530
## Mean   :0.4052   Mean   :1.2169   Mean   : 2.866   Mean   : 40.337
## 3rd Qu.:0.4789   3rd Qu.:1.4740   3rd Qu.: 3.357   3rd Qu.: 45.190
## Max.   :2.8730   Max.   :4.8850   Max.   :21.980   Max.   :542.200
## smoothness_se      compactness_se      concavity_se     concave.points_se
## Min.   :0.001713   Min.   :0.002252   Min.   :0.00000   Min.   :0.000000
## 1st Qu.:0.005169   1st Qu.:0.013080   1st Qu.:0.01509   1st Qu.:0.007638
## Median :0.006380   Median :0.020450   Median :0.02589   Median :0.010930
## Mean   :0.007041   Mean   :0.025478   Mean   :0.03189   Mean   :0.011796
## 3rd Qu.:0.008146   3rd Qu.:0.032450   3rd Qu.:0.04205   3rd Qu.:0.014710
## Max.   :0.031130   Max.   :0.135400   Max.   :0.39600   Max.   :0.052790
##   symmetry_se       fractal_dimension_se  radius_worst     texture_worst
## Min.   :0.007882   Min.   :0.0008948   Min.   : 7.93   Min.   :12.02
## 1st Qu.:0.015160   1st Qu.:0.0022480   1st Qu.:13.01   1st Qu.:21.08
## Median :0.018730   Median :0.0031870   Median :14.97   Median :25.41
## Mean   :0.020542   Mean   :0.0037949   Mean   :16.27   Mean   :25.68
## 3rd Qu.:0.023480   3rd Qu.:0.0045580   3rd Qu.:18.79   3rd Qu.:29.72
## Max.   :0.078950   Max.   :0.0298400   Max.   :36.04   Max.   :49.54
## perimeter_worst     area_worst       smoothness_worst  compactness_worst
## Min.   : 50.41   Min.   : 185.2   Min.   :0.07117   Min.   :0.02729
## 1st Qu.: 84.11   1st Qu.: 515.3   1st Qu.:0.11660   1st Qu.:0.14720
## Median : 97.66   Median : 686.5   Median :0.13130   Median :0.21190
## Mean   :107.26   Mean   : 880.6   Mean   :0.13237   Mean   :0.25427
## 3rd Qu.:125.40   3rd Qu.:1084.0   3rd Qu.:0.14600   3rd Qu.:0.33910
## Max.   :251.20   Max.   :4254.0   Max.   :0.22260   Max.   :1.05800
## concavity_worst  concave.points_worst symmetry_worst    fractal_dimension_worst
## Min.   :0.0000   Min.   :0.00000   Min.   :0.1565   Min.   :0.05504
## 1st Qu.:0.1145   1st Qu.:0.06493   1st Qu.:0.2504   1st Qu.:0.07146
## Median :0.2267   Median :0.09993   Median :0.2822   Median :0.08004
## Mean   :0.2722   Mean   :0.11461   Mean   :0.2901   Mean   :0.08395
## 3rd Qu.:0.3829   3rd Qu.:0.16140   3rd Qu.:0.3179   3rd Qu.:0.09208
## Max.   :1.2520   Max.   :0.29100   Max.   :0.6638   Max.   :0.20750
```

```r
# Data types
str(cancer_data)
```

```
## 'data.frame':    569 obs. of  32 variables:
##  $ id                : int  842302 842517 84300903 84348301 84358402 843786 844359 84458202 8449
##  $ 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 ...
```

```r
# Missing values
sum(is.na(cancer_data))
```
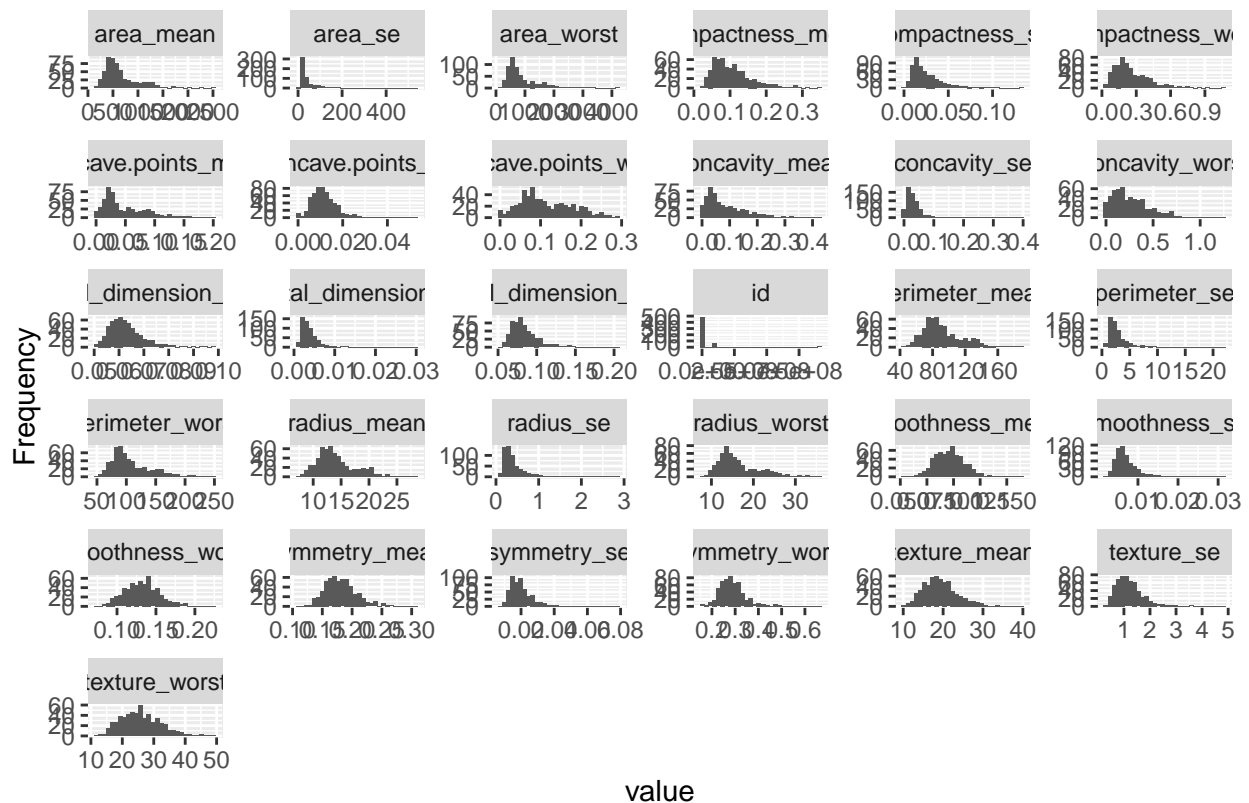
```
## [1] 0
```

```r
# Duplicates
sum(duplicated(cancer_data))
```

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

```r
# Distribution of predictors
plot_histogram(cancer_data, nrow = 6, ncol = 6)
```

```r
# Distribution of diagnosis classes
table(cancer_data$diagnosis)
```

```
##
##   B   M
## 357 212
```

```r
prop.table(table(cancer_data$diagnosis))
```

```
##
##         B         M
## 0.6274165 0.3725835
```

```r
# Relationship between predictors and response
predictor_data <- cancer_data[, names(cancer_data) != "diagnosis"]

# Convert to long format
df_long <- data.frame(
  diagnosis = rep(cancer_data$diagnosis, times = ncol(predictor_data)),
  feature = rep(names(predictor_data), each = nrow(cancer_data)),
  value = as.vector(as.matrix(predictor_data))
)

ggplot(df_long, aes(x = value, fill = diagnosis)) +
  geom_density(alpha = 0.5) +
  facet_wrap(~ feature, scales = "free") +
  theme_minimal()
```

```r
# Predictors w/ near zero variance
degenerate <- nearZeroVar(predictor_data)
print(degenerate)
```

```
## integer(0)
```

```r
# Correlation between predictors
cor_matrix <- cor(predictor_data)
cor_long <- melt(cor_matrix)

ggplot(cor_long, aes(Var1, Var2, fill = value)) +
  geom_tile(color = "white") +
  scale_fill_gradient2(low = "blue", high = "red", mid = "white",
                       midpoint = 0, limit = c(-1, 1), space = "Lab",
                       name = "Correlation") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, vjust = 1, hjust = 1)) +
  coord_fixed() +
  labs(title = "Predictor Correlation Heatmap", x = "", y = "")
```

## Predictor Correlation Heatmap



```
# Skewness
apply(cancer_data[, -2], 2, skewness)
```

```
##                       id              radius_mean              texture_mean
##                6.4396595                0.9374168                 0.6470241
##            perimeter_mean                area_mean            smoothness_mean
##                0.9854334                1.6370654                 0.4539207
##          compactness_mean           concavity_mean        concave.points_mean
##                1.1838556                1.3938008                 1.1650124
##            symmetry_mean     fractal_dimension_mean                 radius_se
##                0.7217877                1.2976191                 3.0723468
##                texture_se               perimeter_se                   area_se
##                1.6377733                3.4254803                 5.4185001
##             smoothness_se             compactness_se               concavity_se
##                2.3022616                1.8922032                 5.0835502
##          concave.points_se               symmetry_se       fractal_dimension_se
##                1.4370701                2.1835728                 3.9033041
##              radius_worst              texture_worst             perimeter_worst
##                1.0973059                0.4956970                 1.1222227
##                area_worst           smoothness_worst           compactness_worst
##                1.8495814                0.4132383                 1.4657948
##           concavity_worst        concave.points_worst              symmetry_worst
##                1.1441794                0.4900213                 1.4263764
## fractal_dimension_worst
##                1.6538237
```

## Pre-processing

```r
# Remove uneccessary columns
df <- cancer_data[, -which(names(cancer_data) == "id")]
head(df)
```

```
##   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
```

```r
# Convert diagnosis to factor
df$diagnosis <- factor(df$diagnosis, levels = c("B", "M"))

# BoxCox Transformation
non_bct_cols <- c("smoothness_mean", "texture_worst", "smoothness_worst", "concave.points_worst")
bct_cols <- setdiff(names(df), non_bct_cols)

params <- preProcess(df[, bct_cols], method = "BoxCox")
df_transformed <- predict(params, df[, bct_cols])
df[, bct_cols] <- df_transformed
```

```r
# Confirm transformation
apply(df_transformed[, -1], 2, skewness)
```

```
##             radius_mean              texture_mean            perimeter_mean
##            -0.018084005              -0.013801528              -0.018259725
##               area_mean          compactness_mean            concavity_mean
##             0.283456808              -0.033906489               1.393800804
##      concave.points_mean            symmetry_mean     fractal_dimension_mean
##             1.165012377               0.001737667               0.150646585
##               radius_se                texture_se               perimeter_se
##             0.027176088               0.029036809               0.069227942
##                 area_se              smoothness_se             compactness_se
##             0.115303422              -0.024011982              -0.004019758
##             concavity_se          concave.points_se                symmetry_se
##             5.083550174               1.437070137               0.054910585
##      fractal_dimension_se              radius_worst            perimeter_worst
##             0.012191507               0.026399596               0.061225231
##              area_worst          compactness_worst            concavity_worst
##             0.067682043              -0.220675829               1.144179410
##           symmetry_worst fractal_dimension_worst
##            -0.056548989               0.047053460
```