

Exercise 10.2

Andrew Jordan

2022-05-21

Question 1

```
library(foreign)

surgery_df <- read.arff("data/ThoracicSurgery.arff")

surgery_glm <- glm(Risk1Yr ~ DGN + PRE4 + PRE5 + PRE6 + PRE7 + PRE8 + PRE9 + PRE10 + PRE11 + PRE14 + PRE17 + PRE19 + PRE25 + PRE30 + PRE32 + AGE, family = binomial, data = surgery_df)
summary(surgery_glm)
```



```
##
## Call:
## glm(formula = Risk1Yr ~ DGN + PRE4 + PRE5 + PRE6 + PRE7 + PRE8 +
##      PRE9 + PRE10 + PRE11 + PRE14 + PRE17 + PRE19 + PRE25 + PRE30 +
##      PRE32 + AGE, family = binomial, data = surgery_df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.6084  -0.5439  -0.4199  -0.2762   2.4929
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -1.655e+01  2.400e+03  -0.007  0.99450
## DGNDGN2      1.474e+01  2.400e+03   0.006  0.99510
## DGNDGN3      1.418e+01  2.400e+03   0.006  0.99528
## DGNDGN4      1.461e+01  2.400e+03   0.006  0.99514
## DGNDGN5      1.638e+01  2.400e+03   0.007  0.99455
## DGNDGN6      4.089e-01  2.673e+03   0.000  0.99988
## DGNDGN8      1.803e+01  2.400e+03   0.008  0.99400
## PRE4         -2.272e-01  1.849e-01  -1.229  0.21909
## PRE5         -3.030e-02  1.786e-02  -1.697  0.08971 .
## PRE6PRZ1     -4.427e-01  5.199e-01  -0.852  0.39448
## PRE6PRZ2     -2.937e-01  7.907e-01  -0.371  0.71030
## PRE7T        7.153e-01  5.556e-01   1.288  0.19788
## PRE8T        1.743e-01  3.892e-01   0.448  0.65419
## PRE9T        1.368e+00  4.868e-01   2.811  0.00494 **
## PRE10T       5.770e-01  4.826e-01   1.196  0.23185
## PRE11T       5.162e-01  3.965e-01   1.302  0.19295
## PRE14OC12    4.394e-01  3.301e-01   1.331  0.18318
## PRE14OC13    1.179e+00  6.165e-01   1.913  0.05580 .
```

```
## PRE14OC14      1.653e+00  6.094e-01   2.713  0.00668 **
## PRE17T         9.266e-01  4.445e-01   2.085  0.03709 *
## PRE19T        -1.466e+01  1.654e+03  -0.009  0.99293
## PRE25T        -9.789e-02  1.003e+00  -0.098  0.92227
## PRE30T         1.084e+00  4.990e-01   2.172  0.02984 *
## PRE32T        -1.398e+01  1.645e+03  -0.008  0.99322
## AGE           -9.506e-03  1.810e-02  -0.525  0.59944
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 395.61  on 469  degrees of freedom
## Residual deviance: 341.19  on 445  degrees of freedom
## AIC: 391.19
##
## Number of Fisher Scoring iterations: 15
```

According to the above summary, it appears that the variables “Pre9”, those with Dyspnoea before the surgery, and the “OC14” variant of “Pre14” variable that list the size of the original tumor, have the greatest effect on survival rate. “Pre30” and “Pre17” also having a large impact as well, showing that these variables respectively tracking smoking and type 2 diabetes also have a large impact. The p-values of “Pre9” at .005 and “Pre14OC14” at .007 indicate a high statistical significance. The p-values of “Pre30” at .03 and “Pre17” at .04 also indicate statistical significance, but not nearly as much as “Pre9” and “Pre14OC14”.

```
surgery_res <- predict(surgery_glm, type = "response")
surgeryPrediction <- table(Actual = surgery_df$Risk1Yr, Predicted = surgery_res > .5)

(surgeryPrediction[[1,1]] + surgeryPrediction[[2,2]]) / sum(surgeryPrediction)
```

```
## [1] 0.8361702
```

The accuracy of the model is approximately 84%.

Question 2

```
binary_df <- read.csv("data/binary-classifier-data.csv")
head(binary_df)
```

```
##   label      x      y
## 1     0 70.88469 83.17702
## 2     0 74.97176 87.92922
## 3     0 73.78333 92.20325
## 4     0 66.40747 81.10617
## 5     0 69.07399 84.53739
## 6     0 72.23616 86.38403
```

```
binary_glm <- glm(label ~ x+y, data = binary_df, family = binomial)
summary(binary_glm)
```

```
##
## Call:
## glm(formula = label ~ x + y, family = binomial, data = binary_df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.3728  -1.1697  -0.9575   1.1646   1.3989
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  0.424809   0.117224   3.624  0.00029 ***
## x           -0.002571   0.001823  -1.411  0.15836
## y           -0.007956   0.001869  -4.257  2.07e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 2075.8  on 1497  degrees of freedom
## Residual deviance: 2052.1  on 1495  degrees of freedom
## AIC: 2058.1
##
## Number of Fisher Scoring iterations: 4

binary_res <- predict(binary_glm, type = "response")
binaryPrediction <- table(Actual = binary_df$label, Predicted = binary_res > .5)

(binaryPrediction[[1,1]] + binaryPrediction[[2,2]]) / sum(binaryPrediction)

## [1] 0.5834446
```

The accuracy of the logistic regression classifier is approximately 58%.