

# DATA 621 Assignment 2

Darwhin Gomez

## Table of contents

1- Loading Data . . . . .	2
2- Table function . . . . .	3
3- Accuracy function . . . . .	3
4- Classification Error Rate . . . . .	4
5- Precision . . . . .	4
6- Sensitivity (Recall) . . . . .	5
7- Specificity . . . . .	5
8- F1 Score . . . . .	5
10- ROC Curve Function and plot . . . . .	6
11- All Metrics Comparison Caret . . . . .	7
12-13 pROC . . . . .	9

```
1 library(dplyr)
```

Attaching package: 'dplyr'

The following objects are masked from 'package:stats':

filter, lag

The following objects are masked from 'package:base':

intersect, setdiff, setequal, union

```
1 library(caret)
```

Loading required package: ggplot2

Loading required package: lattice

```
1 library(pROC)
```

Type 'citation("pROC")' for a citation.

Attaching package: 'pROC'

The following objects are masked from 'package:stats':

cov, smooth, var

```
1 library(ggplot2)
```

## 1- Loading Data

```
1 data <- read.csv("classification-output-data.csv")
```

```
1 head(data,5)
```

	pregnant	glucose	diastolic	skinfold	insulin	bmi	pedigree	age	class
1	7	124	70	33	215	25.5	0.161	37	0
2	2	122	76	27	200	35.9	0.483	26	0
3	3	107	62	13	48	22.9	0.678	23	1
4	1	91	64	24	0	29.2	0.192	21	0
5	4	83	86	19	0	29.3	0.317	34	0

	scored.class	scored.probability
1	0	0.32845226
2	0	0.27319044
3	0	0.10966039
4	0	0.05599835
5	0	0.10049072

The class column represents the actual class whilst the scored.class represents the predicted class.

## 2- Table function

```
1 conf_matrix<-table(data$class, data$scored.class)
2 conf_matrix
```

```
      0   1
0 119   5
1  30  27
```

Interpretation of the Values:

- True Positives (TP) = 119: The model *correctly* predicted Positive cases.
- True Negatives (TN) = 27: The model *correctly* predicted Negative cases.
- False Positives (FP) = 30: The model *incorrectly* predicted Positive (Type I Error or “False Alarm”). The actual case was Negative.
- False Negatives (FN) = 5: The model *incorrectly* predicted Negative (Type II Error or “Miss”). The actual case was Positive.

## 3- Accuracy function

```
1 acc_func <- function(df, actual, predicted){
2   cm <- table(df[[actual]], df[[predicted]])
3   tp <- cm[2,2]
4   tn <- cm[1,1]
5   fp <- cm[1,2]
6   fn <- cm[2,1]
7
8   (tp + tn) / (tp + tn + fp + fn)
9
10
11
12 }
13 computed_acc<- acc_func(data, "class", "scored.class")
14 print(paste("The calculated accuracy of the model is:",
15             round(acc_func(data, "class", "scored.class"),4)))
```

```
[1] "The calculated accuracy of the model is: 0.8066"
```

## 4- Classification Error Rate

```
1 error_func <- function(df, actual, predicted){
2   cm <- table(df[[actual]], df[[predicted]])
3   tp <- cm[2,2]
4   tn <- cm[1,1]
5   fp <- cm[1,2]
6   fn <- cm[2,1]
7
8   (fp + fn)/(tp + tn + fp +fn)
9
10
11
12 }
13 computed_err <- error_func(data, "class", "scored.class")
14 print(paste("The calculated classification error rate of the model is:",
15             round(error_func(data, "class", "scored.class"),4)))
```

```
[1] "The calculated classification error rate of the model is: 0.1934"
```

We can verify by adding the accuracy and the error rate which should sum to 1

```
1 print(computed_acc + computed_err)
```

```
[1] 1
```

## 5- Precision

```
1 precision_func <- function(df, actual, predicted) {
2   cm <- table(df[[actual]], df[[predicted]])
3   TP <- cm[2,2]; FP <- cm[1,2]
4   TP / (TP + FP)
5 }
6 calculated_prec<- precision_func(data, "class", "scored.class")
7 print(paste("The calculated precision of the model is:",
8             round(precision_func(data, "class", "scored.class"),4)))
```

```
[1] "The calculated precision of the model is: 0.8438"
```

## 6- Sensitivity (Recall)

```
1 sensitivity_func <- function(df, actual, predicted) {
2   cm <- table(df[[actual]], df[[predicted]])
3   TP <- cm[2,2]; FN <- cm[2,1]
4   TP / (TP + FN)
5 }
6 calculated_sensitivity<-sensitivity_func(data, "class", "scored.class")
7 print(paste("The calculated Sensitivity (recall rate) of the model is:",
8             round(sensitivity_func(data, "class", "scored.class"),4)))
```

```
[1] "The calculated Sensitivity (recall rate) of the model is: 0.4737"
```

## 7- Specificity

```
1 specificity_func <- function(df, actual, predicted) {
2   cm <- table(df[[actual]], df[[predicted]])
3   TN <- cm[1,1]; FP <- cm[1,2]
4   TN / (TN + FP)
5 }
6 calculated_specif<-specificity_func(data, "class", "scored.class")
7 print(paste("The calculated Specificity of the model is:",
8             round(specificity_func(data, "class", "scored.class"),4)))
```

```
[1] "The calculated Specificity of the model is: 0.9597"
```

## 8- F1 Score

```
1 f1_score_func <- function(df, actual, predicted) {
2   p <- precision_func(df, actual, predicted)
3   r <- sensitivity_func(df, actual, predicted)
4   2 * (p * r) / (p + r)
5 }
6 calculated_f1<-f1_score_func(data, "class", "scored.class")
7 print(paste("The calculated F1 score of the model is:",
8             round(f1_score_func(data, "class", "scored.class"),4)))
```

```
[1] "The calculated F1 score of the model is: 0.6067"
```

The F1-score will always be between 0 and 1 because it is based on precision and recall, which are both proportions that range from 0 to 1. When either precision or recall is 0, the F1-score is also 0. When both are 1, the F1-score reaches its maximum value of 1. Therefore, since it is the harmonic mean of two values that cannot be less than 0 or greater than 1, the F1-score will always fall between 0 and 1.

## 10- ROC Curve Function and plot

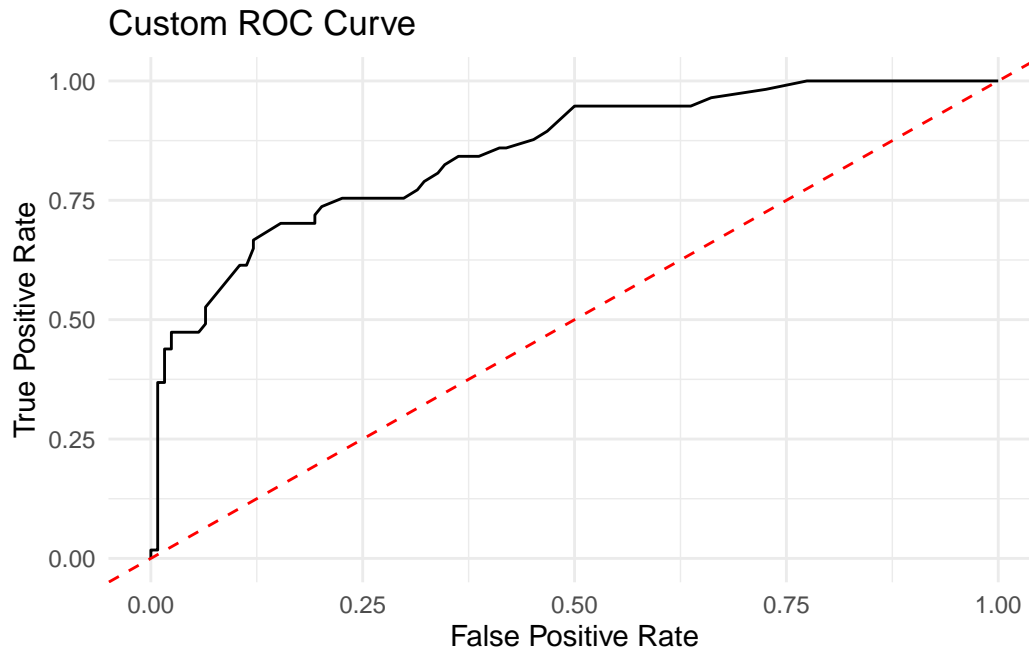
```
1 roc_curve_custom <- function(df, actual, prob_col) {
2   thresholds <- seq(1, 0, -0.01)
3
4   TPR <- FPR <- numeric(length(thresholds))
5
6   for (i in seq_along(thresholds)) {
7     t <- thresholds[i]
8     df$pred <- ifelse(df[[prob_col]] >= t, 1, 0)
9
10
11     cm <- table(factor(df[[actual]], levels = c(0,1)),
12                factor(df$pred, levels = c(0,1)))
13
14     TP <- cm[2,2]; TN <- cm[1,1]; FP <- cm[1,2]; FN <- cm[2,1]
15     TPR[i] <- TP / (TP + FN)
16     FPR[i] <- FP / (FP + TN)
17   }
18
19   auc <- sum(diff(FPR) * (head(TPR, -1) + tail(TPR, -1)) / 2)
20
21   p<-ggplot(data.frame(FPR, TPR), aes(x = FPR, y = TPR)) +
22     geom_line(color = "black") +
23     geom_abline(linetype = "dashed", color = "red") +
24     labs(title = "Custom ROC Curve", x = "False Positive Rate",
25          y = "True Positive Rate") +
26     theme_minimal()
27
28
29   return(list( Plot=p,AUC = auc))
30 }
```

```

31
32 roc_curve_custom(data, "class", "scored.probability")

```

\$Plot



\$AUC

```
[1] 0.8488964
```

## 11- All Metrics Comparison Caret

```

1 data.frame(
2   Accuracy = computed_acc,
3   Error_Rate = computed_err,
4   Precision = calculated_prec,
5   Sensitivity = calculated_sensitivity,
6   Specificity = calculated_specif,
7   F1_Score = calculated_f1
8 )

```

	Accuracy	Error_Rate	Precision	Sensitivity	Specificity	F1_Score
1	0.8066298	0.1933702	0.84375	0.4736842	0.9596774	0.6067416

```

1 caret_conf <- confusionMatrix(
2   factor(data$scored.class),
3   factor(data$class),
4   positive = "1"
5 )
6 caret_conf

```

#### Confusion Matrix and Statistics

	Reference	
Prediction	0	1
0	119	30
1	5	27

Accuracy : 0.8066  
 95% CI : (0.7415, 0.8615)  
 No Information Rate : 0.6851  
 P-Value [Acc > NIR] : 0.0001712

Kappa : 0.4916

McNemar's Test P-Value : 4.976e-05

Sensitivity : 0.4737  
 Specificity : 0.9597  
 Pos Pred Value : 0.8438  
 Neg Pred Value : 0.7987  
 Prevalence : 0.3149  
 Detection Rate : 0.1492  
 Detection Prevalence : 0.1768  
 Balanced Accuracy : 0.7167

'Positive' Class : 1

The functions match output from the caret package indicating caret can be used to calculate these metrics reliably and easily.



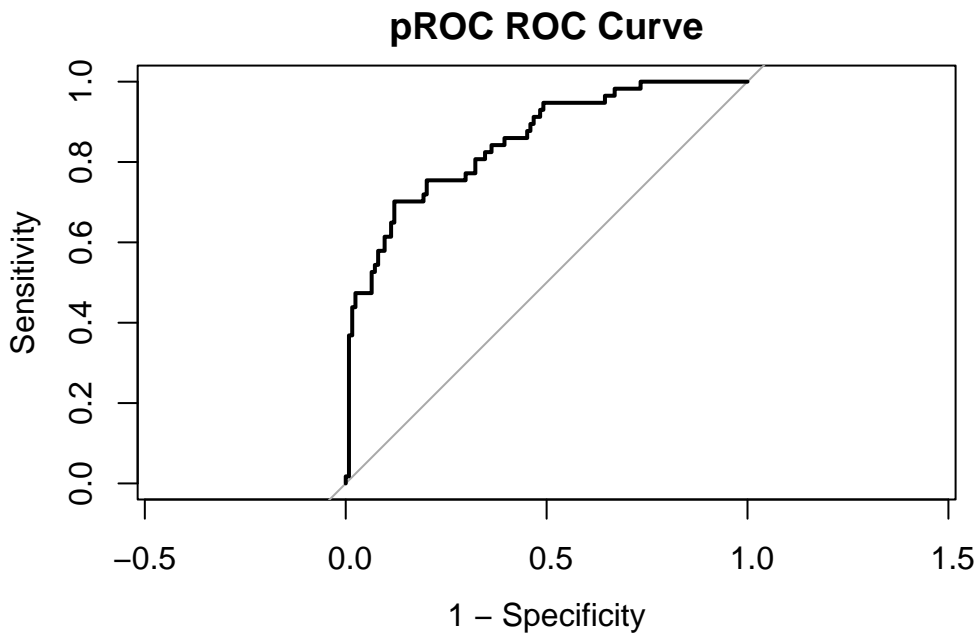
## 12-13 pROC

```
1 roc_obj <- roc(data$class, data$scored.probability)
```

Setting levels: control = 0, case = 1

Setting direction: controls < cases

```
1 plot(roc_obj, col="black", main="pROC ROC Curve", legacy.axes = TRUE)
```



```
1 auc(roc_obj)
```

Area under the curve: 0.8503

The computed ROC plot and the pROC plot display very similar curves, with the custom AUC at 0.8488 and the pROC AUC at 0.8503, indicating highly consistent results in the calculation of the area under the curve. Overall, both the ROC plots and AUC values demonstrate that the model is performing well, particularly in accurately identifying positive cases.

The packages work well and should be used to produce classification metrics efficiently and reliably.