

SLLD - Module 1

Classification

L.Emer, F.Chiaromonte

Sant'Anna School of Advanced Study - Pisa

12/2/2026

Libraries

```
library(dplyr)
library(caret)
library(MASS)
library(klaR)
library(mvtnorm)
```

Data

We simulate data as follows

```
set.seed(123)
n <- 500
p <- 5
X <- matrix(rnorm(n*p), nrow = n, ncol = p)
X <- data.frame(X)
colnames(X) <- paste("X", 1:p, sep="")
eta <- -0.5 - 1*X[,1] - 2*X[,2] + X[,3] + 3*X[,4] + 0.5*X[,5]
piy <- (exp(eta)) / (1 + exp(eta))
mean(piy)
```

```
## [1] 0.5653752
```

```
Y<-rbinom(n, 1, piy)
summary(Y)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0.00    0.00    1.00    0.56    1.00    1.00
```

```
df<-cbind(Y,X)

set.seed(123)
training_samples <- df$Y %>%
  caret::createDataPartition(p = 0.8, list = FALSE)
# p indicates the percentage of training data
train <- df[training_samples, ]
test  <- df[-training_samples, ]
```

Logistic Regression

Simple Logistic Regression

The **glm()** function can be used to fit many types of generalized linear models, including logistic regression.

It is similar to **lm()**, except that we must linear model pass in the argument **family = binomial** in order to tell R to run a logistic regression rather than some other type of generalized linear model.

Using the training set, we build a simple logistic regression model using X_1 as the only explanatory variable of the response Y

```
simple_glm <- glm(Y ~ X1,  
                 data = train, family = 'binomial')  
summary(simple_glm)
```

```
##  
## Call:  
## glm(formula = Y ~ X1, family = "binomial", data = train)  
##  
## Coefficients:  
##             Estimate Std. Error z value Pr(>|z|)  
## (Intercept)  0.2409      0.1028   2.343  0.0191 *  
## X1          -0.4338      0.1109  -3.910 9.22e-05 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## (Dispersion parameter for binomial family taken to be 1)  
##  
##      Null deviance: 549.22  on 399  degrees of freedom  
## Residual deviance: 533.02  on 398  degrees of freedom  
## AIC: 537.02  
##  
## Number of Fisher Scoring iterations: 4
```

Let's see predictive power of our model in terms of **accuracy** -- i.e. the proportion of correct predictions, both true positives and true negatives, among the total number of cases examined

```
# Test for accuracy: predict test data
predict_1 <- predict(simple_glm, newdata = test,type = 'response')

# round up the predictions
predict_1 <- ifelse(predict_1>0.5, 1, 0)

# calculate accuracy
accuracySim <- mean(predict_1==test$Y)
accuracySim
```

```
## [1] 0.71
```

```
ConfMat_SR = confusionMatrix(as.factor(predict_1),
                             as.factor(test$Y))
ConfMat_SR$table
```

```
##           Reference
## Prediction  0  1
##           0 22  8
##           1 21 49
```

Note: the "event" is 0!

Multiple Logistic Regression

```
glm_complete <- glm(Y ~ ., data=train, family = 'binomial')
summary(glm_complete)
```

```
##
## Call:
## glm(formula = Y ~ ., family = "binomial", data = train)
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   0.5152     0.1759   2.929  0.00340 **
## X1            -1.2380     0.2220  -5.577 2.44e-08 ***
## X2            -2.4090     0.2807  -8.582  < 2e-16 ***
## X3             1.0528     0.1935   5.441 5.29e-08 ***
## X4             3.0775     0.3355   9.173  < 2e-16 ***
## X5             0.5673     0.1755   3.232  0.00123 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 549.22  on 399  degrees of freedom
## Residual deviance: 231.04  on 394  degrees of freedom
## AIC: 243.04
##
## Number of Fisher Scoring iterations: 6
```

Let's see its predictive accuracy

```
# Test for accuracy: predict test data
predict_1 <- predict(glm_complete, newdata =
                      test, type = 'response')

# round up the predictions
predict_1 <- ifelse(predict_1 > 0.5, 1, 0)

# calculate accuracy
accuracySat <- mean(predict_1 == test$Y)
accuracySat
```

```
## [1] 0.92
```

```
ConfMat_LR = confusionMatrix(as.factor(predict_1),
                              as.factor(test$Y))
ConfMat_LR$table
```

```
##           Reference
## Prediction  0  1
##           0 38  3
##           1  5 54
```


LDA and QDA

```
set.seed(123)
library(mvtnorm)
Nj      <- c(15, 25, 20)
Sigma   <- matrix(c(18,-2, -2,13), byrow=TRUE, ncol=2)
mu1     <- c(-4, 4)
mu2     <- c( 3, 3)
mu3     <- c( 1, -1)
df1     <- rmvnorm(Nj[1], mean=mu1, sigma=Sigma)
df2     <- rmvnorm(Nj[2], mean=mu2, sigma=Sigma)
df3     <- rmvnorm(Nj[3], mean=mu3, sigma=Sigma)
df      <- rbind(df1, df2, df3)
Y       <- factor(rep(1:length(Nj), Nj))
df      <- data.frame(Y, X1=df[, 1], X2=df[, 2])

set.seed(123)
training_samples <- df$Y %>%
  caret::createDataPartition(p = 0.8, list = FALSE)
train <- df[training_samples, ]
test  <- df[-training_samples, ]
```

We use the function **lda** to perform linear discriminant analysis

```
lda.Y <- lda(Y~ ., data=train)
lda.Y
```

```
## Call:
## lda(Y ~ ., data = train)
##
## Prior probabilities of groups:
##          1          2          3
## 0.2500000 0.4166667 0.3333333
##
## Group means:
##          X1          X2
## 1 -3.7122490  4.072981
## 2  3.3231770  3.871667
## 3  0.7695065 -0.823288
##
## Coefficients of linear discriminants:
##          LD1          LD2
## X1 -0.2348374 -0.1279290
## X2 -0.1849840  0.2605693
##
## Proportion of trace:
##    LD1    LD2
## 0.5384 0.4616
```

Here is the model accuracy on training data.

```
predmodel.train.lda = predict(lda.Y, data=train_transformed)
ConfMat_LDAttrain = confusionMatrix(as.factor(
  predmodel.train.lda$class),
                                   as.factor(train$Y))
ConfMat_LDAttrain$table
```

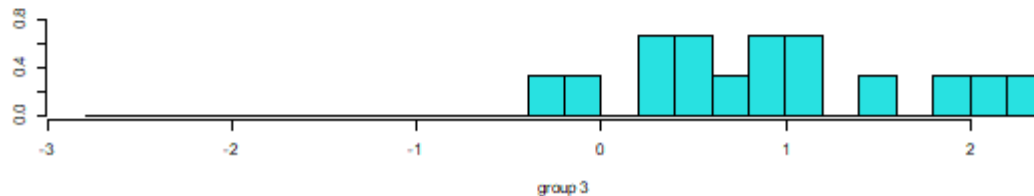
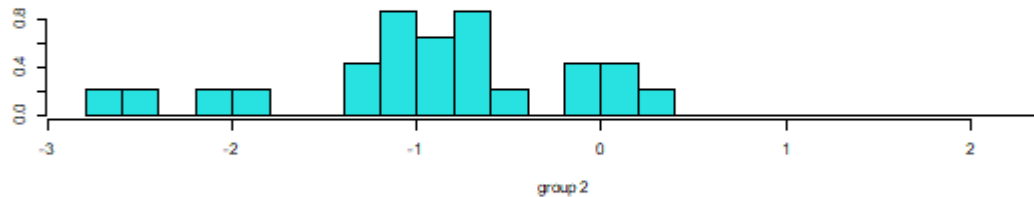
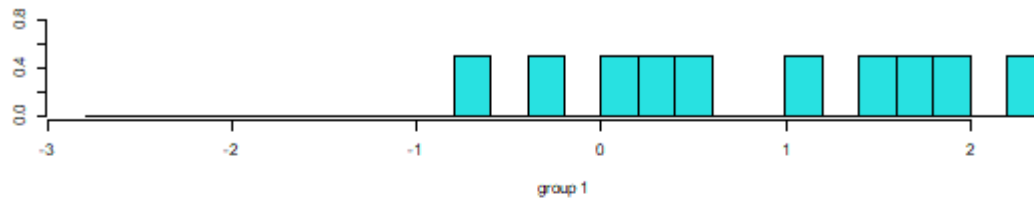
```
##           Reference
## Prediction  1  2  3
##           1  7  3  0
##           2  3 16  4
##           3  2  1 12
```

```
ConfMat_LDAttrain$byClass[1:3,1:2]
```

```
##           Sensitivity Specificity
## Class: 1    0.5833333    0.9166667
## Class: 2    0.8000000    0.7500000
## Class: 3    0.7500000    0.9062500
```

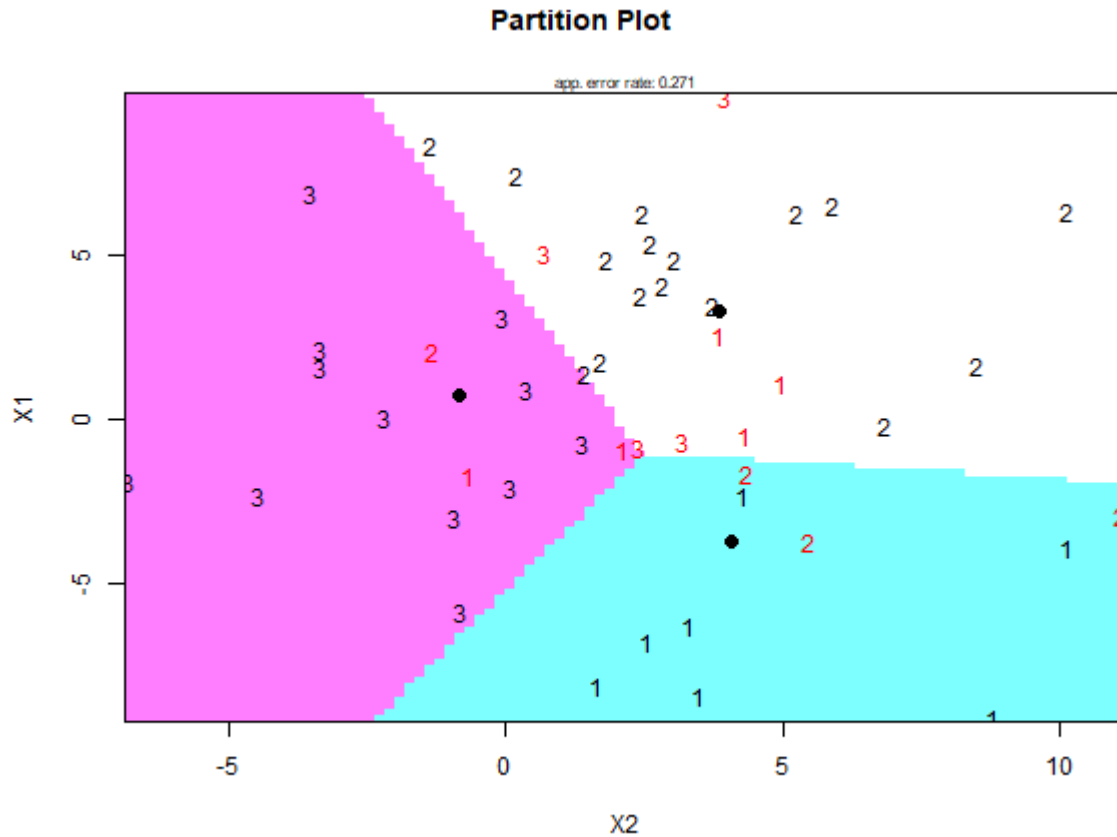
The plot below shows how the response class has been classified by the LDA classifier. The x -axis shows the value of the line defined by the coefficient of linear discriminant for LDA. Groups are the ones in the response classes.

```
ldahist(predmodel.train.lda$x[,1], g= predmodel.train.lda$class)
```



We use the function **partimat** to see geometric division

```
partimat(factor(Y)~ X1+X2,  
         data=train, method = "lda")
```



Now we check the model accuracy on test data.

```
predmodel.test.lda = predict(lda.Y, newdata=test)
ConfMat_LDAtest = confusionMatrix(as.factor(
  predmodel.test.lda$class),
  as.factor(test$Y))
ConfMat_LDAtest$table
```

```
##           Reference
## Prediction 1 2 3
##           1 1 0 0
##           2 0 3 1
##           3 2 2 3
```

```
mean(predmodel.test.lda$class==test$Y)
```

```
## [1] 0.5833333
```

```
ConfMat_LDAtest$byClass[1:2]
```

```
## [1] 0.3333333 0.6000000
```

Next we will fit the model through QDA. The command is similar to LDA and it outputs the prior probabilities and Group means. Note that “Prior Probabilities” and “Group Means” values are same as of LDA.

```
qda.Y <- qda(Y~ ., data=train)
qda.Y
```

```
## Call:
## qda(Y ~ ., data = train)
##
## Prior probabilities of groups:
##      1      2      3
## 0.2500000 0.4166667 0.3333333
##
## Group means:
##      X1      X2
## 1 -3.7122490 4.072981
## 2  3.3231770 3.871667
## 3  0.7695065 -0.823288
```

Now we check the model accuracy on test data.

```
predmodel.test.qda = predict(qda.Y, newdata=test)
ConfMat_QDAtest = confusionMatrix(as.factor(
  predmodel.test.qda$class),
  as.factor(test$Y))
ConfMat_QDAtest$table
```

```
##           Reference
## Prediction 1 2 3
##           1 1 0 0
##           2 0 3 1
##           3 2 2 3
```

```
mean(predmodel.test.qda$class==test$Y)
```

```
## [1] 0.5833333
```

```
ConfMat_QDAtest$byClass[1:2]
```

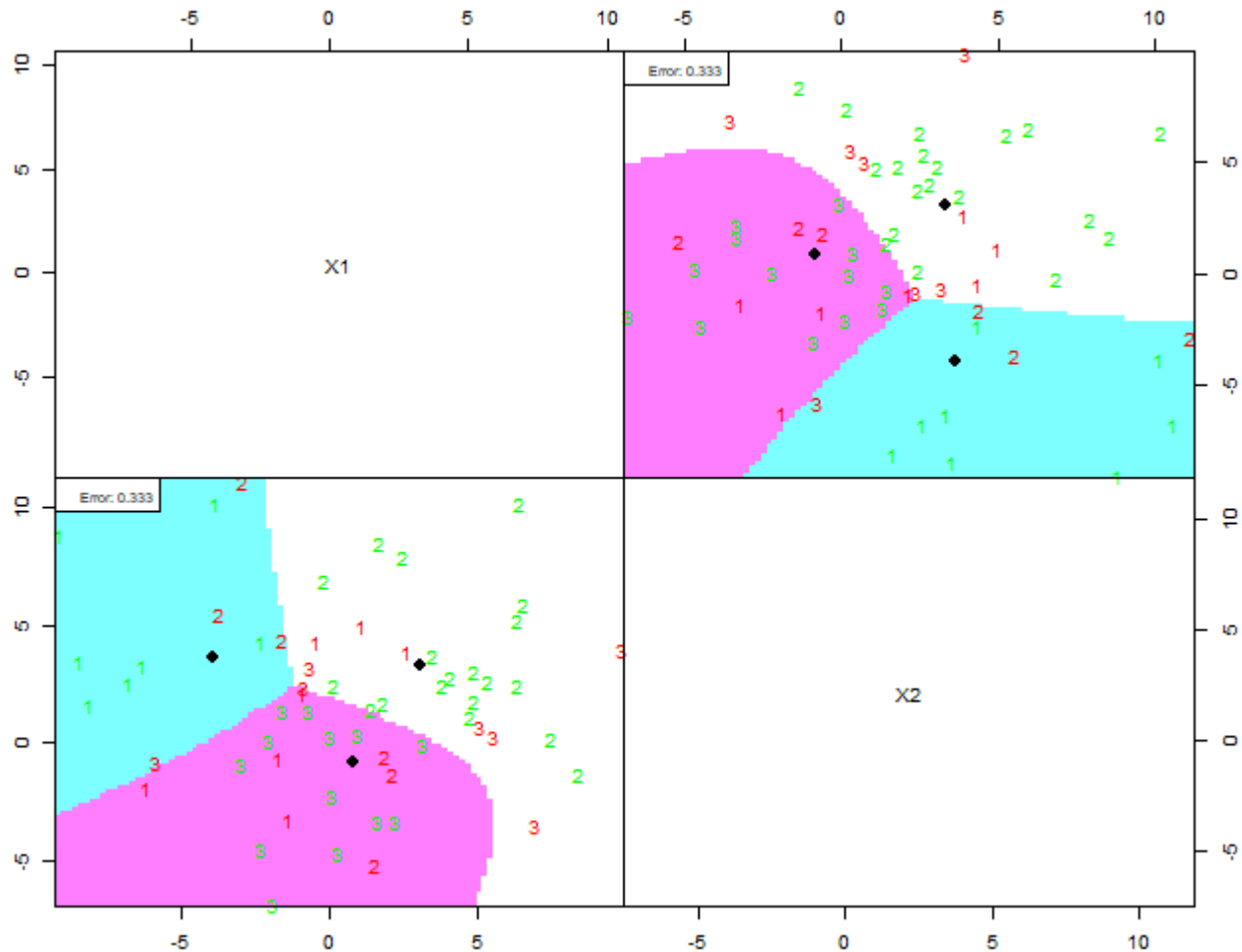
```
## [1] 0.3333333 0.6000000
```



```

partimat(factor(Y) ~ .,
data=df, method = "qda", plot.matrix=TRUE,
col.correct='green', col.wrong='red')

```



kNN

We are going to use the **knn3** function within caret package. Let's train the knn with $k = 1$, 10

```
knn_Y1 <- knn3(factor(Y) ~ ., data=train, k = 1)
knn_Y10 <- knn3(factor(Y) ~ ., data=train, k = 10)
predict_knn1 <- predict(knn_Y1, test, type='class')
ConfMat_knn1 = confusionMatrix(predict_knn1, as.factor(test$Y))
ConfMat_knn1$overall[1] # the accuracy
```

```
## Accuracy
## 0.5833333
```

```
predict_knn10 <- predict(knn_Y10, test, type='class')
ConfMat_knn10 = confusionMatrix(predict_knn10, as.factor(test$Y))
ConfMat_knn10$overall[1] # the accuracy
```

```
## Accuracy
## 0.5
```

Now it's your turn!!!