

ISLR | Chapter 9 Exercises

Marshall McQuillen

11/23/2018

Conceptual

1

- A & B.

```
suppressPackageStartupMessages(library(gridExtra))
suppressPackageStartupMessages(library(ggplot2))

# data setup
x <- matrix(rnorm(1000*2), ncol = 2)

# define hyperplane 1 (part A)
hyperplane1 <- 1 + 3*x[, 1] - x[, 2]
y <- ifelse(hyperplane1 > 0, 1, -1)

# define hyperplane 2 (part b)
hyperplane2 <- -2 + x[, 1] + 2*x[, 2]
y2 <- ifelse(hyperplane2 > 0, 1, -1)

par(mfrow = c(1, 2))

# plot hyperplanes
plot1 <- ggplot(data.frame(x = x,
  y = factor(y, levels = c(-1, 1))), aes(x.1, x.2, color=y)) +
  geom_point(show.legend = FALSE) +
  geom_abline(intercept = 1,
    slope = 3,
    linetype = 'dashed') +
  ggtitle("Binary Class Data | Hyperplane A") +
  xlab("X1") +
  ylab("X2") +
  labs(color = "Class")

plot2 <- ggplot(data.frame(x = x,
  y = factor(y2, levels = c(-1, 1))), aes(x.1, x.2, color=y)) +
  geom_point() +
  geom_abline(intercept = 1,
    slope = -0.5,
    linetype = 'dashed') +
  ggtitle("Binary Class Data | Hyperplane B") +
  xlab("X1") +
  ylab("X2") +
  labs(color = "Class")
```

```
grid.arrange(plot1, plot2, ncol = 2)
```



2

- **A, B & C.** The hyperplane is the circle encompassing the pink/orange data points below, where those data points are the ones whose value, when plugged into the equation $f(X_1, X_2) = (1 + X_1)^2 + (2 - X_2)^2 - 4$, will be negative. The blue data points output value of the above equation would be positive. The 4 data points plotted in black are those requested in part C, and it is clear which class they would fall into.

```
# data setup
x_1 <- runif(2500, min = -4, max = 2)
x_2 <- runif(2500, min = -1, max = 5)

# define hyperplane
hyperplane3 <- (1 + x_1)^2 + (2 - x_2)^2 - 4
y3 <- factor(ifelse(hyperplane3 > 0, ">4", "<4"))

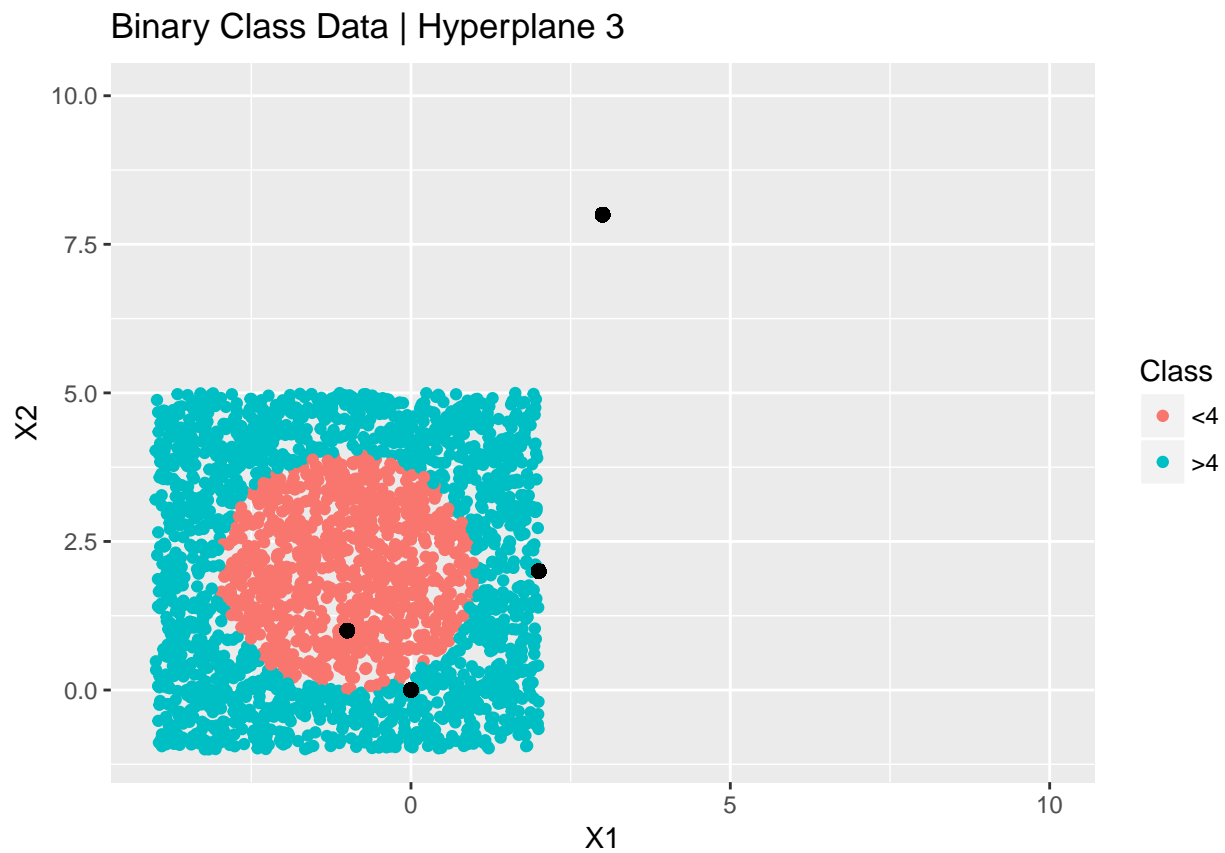
df <- data.frame(x1 = x_1,
                 x2 = x_2,
                 y = y3)

# plot hyperplane with additional points in black
ggplot(df, aes(x1, x2, color=y)) +
```

```

geom_point() +
geom_point(x = 0,
           y = 0,
           col = 'black',
           cex = 2) +
geom_point(x = -1,
           y = 1,
           col = 'black',
           cex = 2) +
geom_point(x = 2,
           y = 2,
           col = 'black',
           cex = 2) +
geom_point(x = 3,
           y = 8,
           col = 'black',
           cex = 2) +
ggtitle("Binary Class Data | Hyperplane 3") +
scale_y_continuous(limits = c(-1, 10)) +
scale_x_continuous(limits = c(-4, 10)) +
xlab("X1") +
ylab("X2") +
labs(color = "Class")

```



- **D.** One can see that the equation given in the text is non-linear with regard to X_1 and X_2 . However, when expanded and refactored, it is clear that the hyperplane is linear with regard to X_1, X_2, X_1^2 and X_2^2 .

$$(1 + X_1)^2 + (2 - X_2)^2 = 4 \quad (1)$$

$$(1 + X_1)^2 + (2 - X_2)^2 - 4 = 0 \quad (2)$$

$$(1 + 2X_1 + X_1^2) + (4 - 4X_2 + X_2^2) - 4 = 0 \quad (3)$$

$$1 + 2X_1 + X_1^2 - 4X_2 + X_2^2 = 0 \quad (4)$$

$$(5)$$

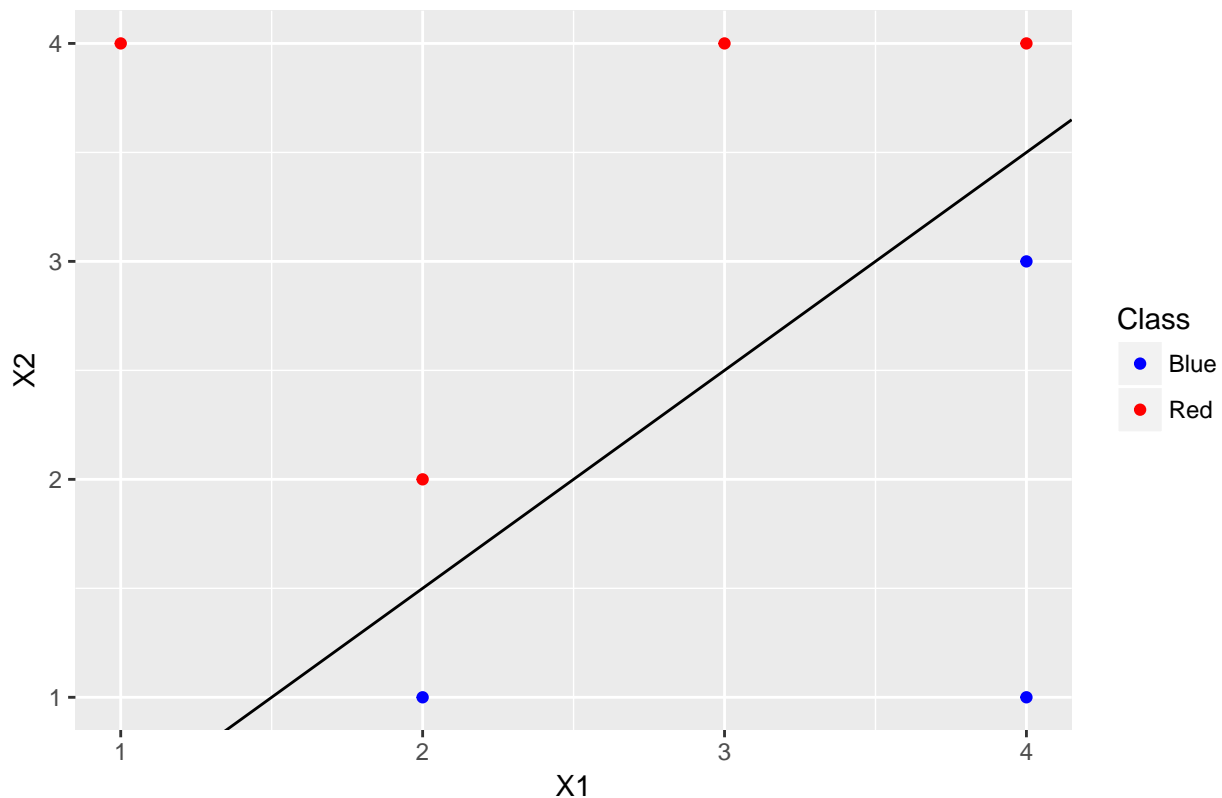
3

- **A & B.** The separating hyperplane has the equation $X_1 - X_2 - 0.5 = 0$.

```
# data setup
df <- data.frame(x_1 = c(3,2,4,1,2,4,4),
                 x_2 = c(4,2,4,4,1,3,1),
                 y = c(rep("Red", 4), rep("Blue", 3)))

# plot MMC
ggplot(df, aes(x_1, x_2, color = y)) +
  geom_point() +
  geom_abline(intercept = -0.5,
              slope = 1,
              linetype = 'solid') +
  scale_color_manual(values = c("Red" = 'red', "Blue" = 'blue')) +
  ggtitle("Maximal Marginal Classifier") +
  xlab("X1") +
  ylab("X2") +
  labs(color = 'Class')
```

Maximal Marginal Classifier



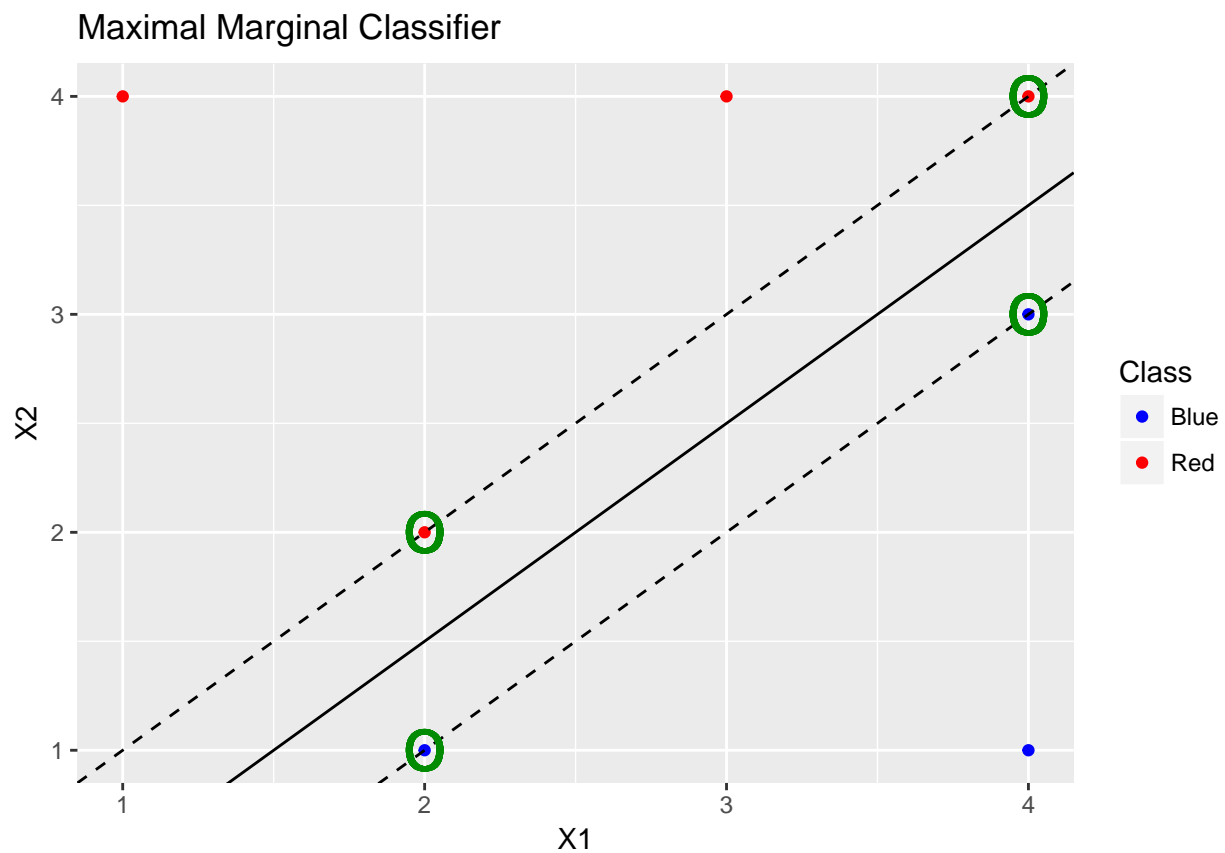
- C. Classify a test observation x_t to Red if $(1)X_1 + (-1)X_2 - 0.5 < 0$, otherwise classify to Blue.
- D. The support vectors are those data points circled in green, the edge of the margins are the dashed lines.

```
# illustrate support vectors and margins
ggplot(df, aes(x_1, x_2, color = y)) +
  geom_point() +
  geom_abline(intercept = -0.5,
              slope = 1,
              linetype = 'solid') +
  geom_abline(intercept = -0,
              slope = 1,
              linetype = 'dashed') +
  geom_abline(intercept = -1,
              slope = 1,
              linetype = 'dashed') +
  geom_point(x = 2,
            y = 2,
            color = 'green4',
            size = 10,
            shape = "o") +
  geom_point(x = 2,
            y = 1,
            color = 'green4',
```

```

      size = 10,
      shape = "o") +
geom_point(x = 4,
  y = 3,
  color = 'green4',
  size = 10,
  shape = "o") +
geom_point(x = 4,
  y = 4,
  color = 'green4',
  size = 10,
  shape = "o") +
scale_color_manual(values = c("Red" = 'red', "Blue" = 'blue')) +
ggtitle("Maximal Marginal Classifier") +
xlab("X1") +
ylab("X2") +
labs(color = 'Class')

```



- **F.** Not only can the seventh observation be removed from the data set without affecting the hyperplane, **any observation that is not a support vector can be removed without distorting the hyperplane**, as shown in the plot below.

```

# remove non support vectors
sub_df <- df[-c(1, 4, 7), ]

```

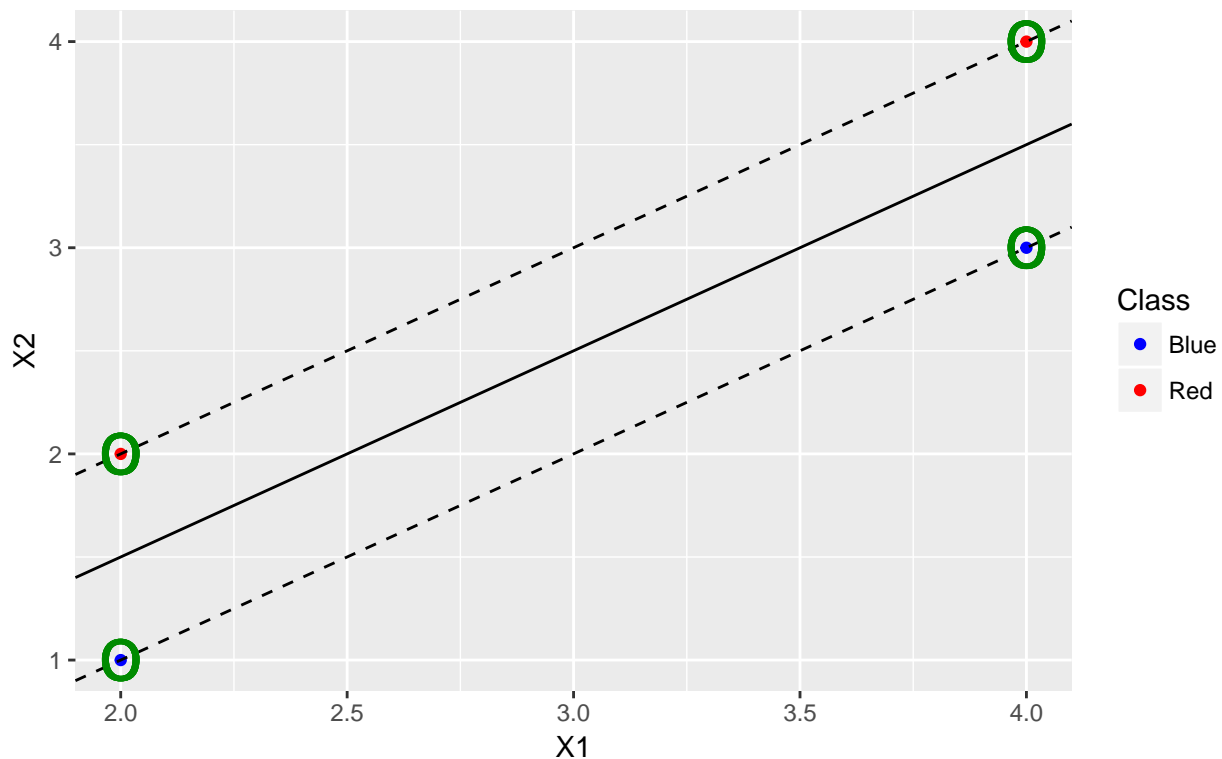
```

# remove non support vectors
ggplot(sub_df, aes(x_1, x_2, color = y)) +
  geom_point() +
  geom_abline(intercept = -0.5,
              slope = 1,
              linetype = 'solid') +
  geom_abline(intercept = -0,
              slope = 1,
              linetype = 'dashed') +
  geom_abline(intercept = -1,
              slope = 1,
              linetype = 'dashed') +
  geom_point(x = 2,
            y = 2,
            color = 'green4',
            size = 10,
            shape = "o") +
  geom_point(x = 2,
            y = 1,
            color = 'green4',
            size = 10,
            shape = "o") +
  geom_point(x = 4,
            y = 3,
            color = 'green4',
            size = 10,
            shape = "o") +
  geom_point(x = 4,
            y = 4,
            color = 'green4',
            size = 10,
            shape = "o") +
  scale_color_manual(values = c("Red" = 'red', "Blue" = 'blue')) +
  ggtitle("Maximal Marginal Classifier",
          "Non Support Vectors Removed") +
  xlab("X1") +
  ylab("X2") +
  labs(color = 'Class')

```

Maximal Marginal Classifier

Non Support Vectors Removed



- **G.** The below hyperplane, with the equation $(1)X_1 + (-1)X_2 - 0.25$ would not be the **maximal** separating hyperplane because the margin between the red support vectors and the hyperplane is smaller than the margin between the blue support vectors and the machine.

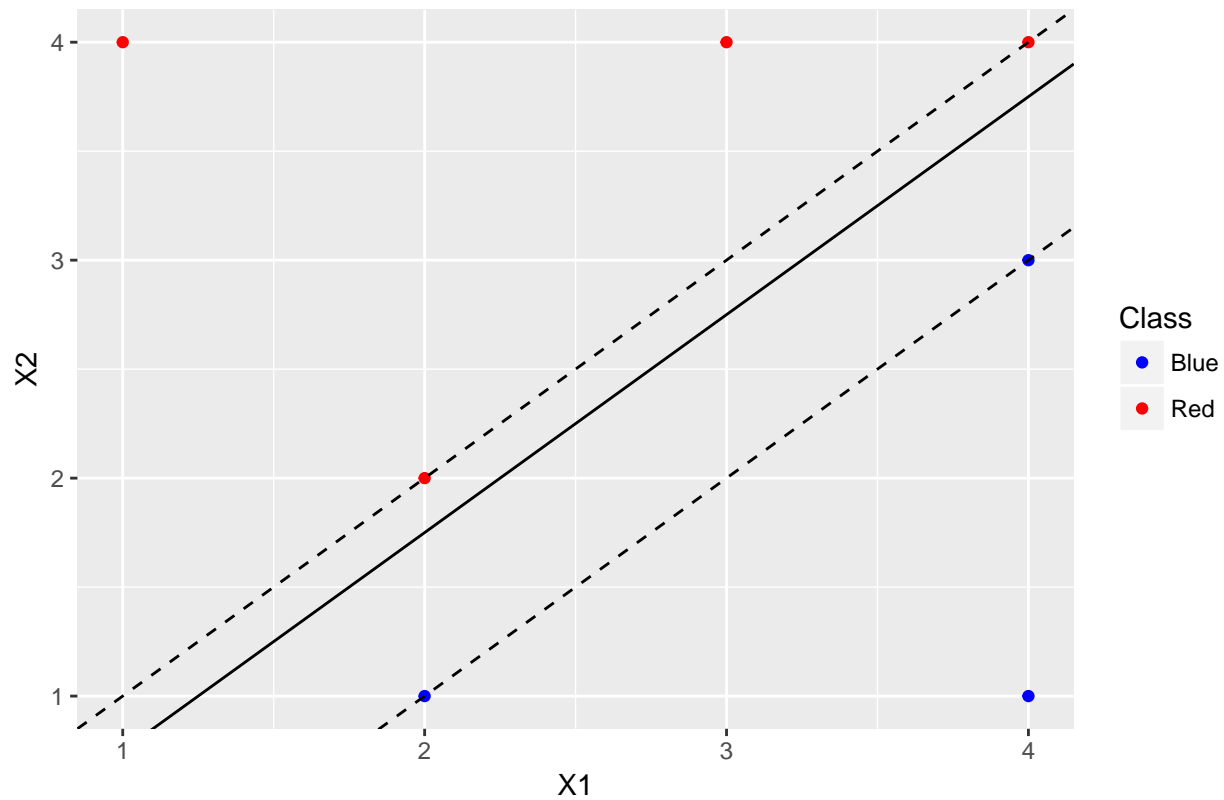
The reason that this is less desirable than the Maximal Marginal Classifier is because the hyperplane, *in this scenario*, is favoring the blue data points, **with no “supporting” evidence (pun intended)**. It is giving them (the blue support vectors) a wider berth than is necessary, at the cost of the berth to the red support vectors.

```
# illustrate non-maximal margin classifier
ggplot(df, aes(x_1, x_2, color = y)) +
  geom_point() +
  geom_abline(intercept = -0.25,
              slope = 1,
              linetype = 'solid') +
  geom_abline(intercept = -0,
              slope = 1,
              linetype = 'dashed') +
  geom_abline(intercept = -1,
              slope = 1,
              linetype = 'dashed') +
  scale_color_manual(values = c("Red" = 'red', "Blue" = 'blue')) +
  ggtitle("Non-Maximal Marginal Classifier") +
  xlab("X1") +
  ylab("X2") +
```



```
labs(color = 'Class')
```

Non-Maximal Marginal Classifier



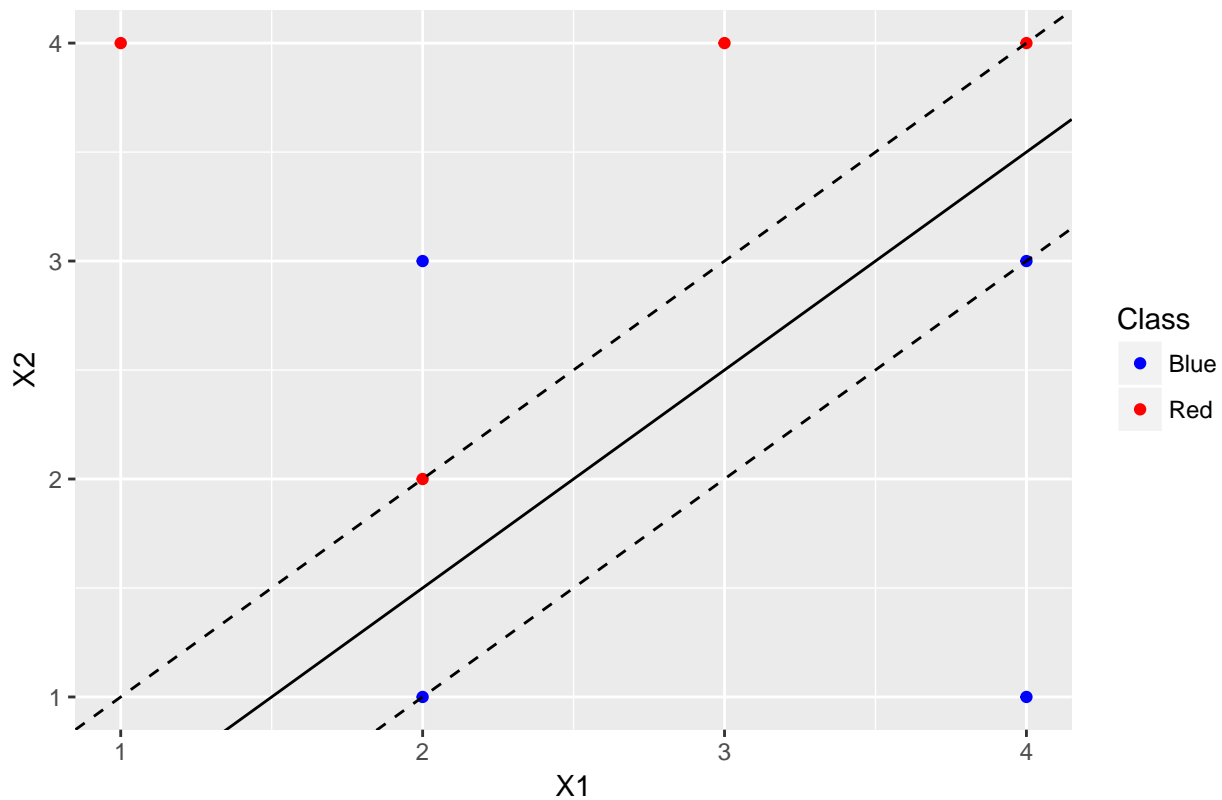
- **H.** A data set where the Maximal Marginal Classifier does not exist.

```
new.point <- data.frame(x_1 = 2,
                        x_2 = 3,
                        y = "Blue")
df <- rbind(df, new.point)

# illustrate the data set where a non-maximal is not possible
ggplot(df, aes(x_1, x_2, color = y)) +
  geom_point() +
  geom_abline(intercept = -0.5,
              slope = 1,
              linetype = 'solid') +
  geom_abline(intercept = -0,
              slope = 1,
              linetype = 'dashed') +
  geom_abline(intercept = -1,
              slope = 1,
              linetype = 'dashed') +
  scale_color_manual(values = c("Red" = 'red', "Blue" = 'blue')) +
  ggtitle("Maximal Marginal Classifier") +
  xlab("X1") +
```

```
ylab("X2") +
labs(color = 'Class')
```

Maximal Marginal Classifier



Applied

4

As shown below, a linear decision boundary is not possible with the data provided. Knowing this, it is unlikely that a linear SVM would outperform a more complex model. This intuition is confirmed in the following three plots, showing that the radial SVM (with $\gamma = 2$) is the model that fits the training data the best.

Not surprisingly, the radial SVM outperforms the linear and polynomial SVM's on the test set as well. Interestingly, the radial SVM also has fewer support vectors, therefore making the model more computationally efficient to store.

```
suppressPackageStartupMessages(library(e1071))
suppressPackageStartupMessages(library(caret))
set.seed(2)

# data setup
x_1 <- rnorm(100)
x_2 <- 5*(x_1^2) + rnorm(100)
idx <- sample(100, 50)
```

```

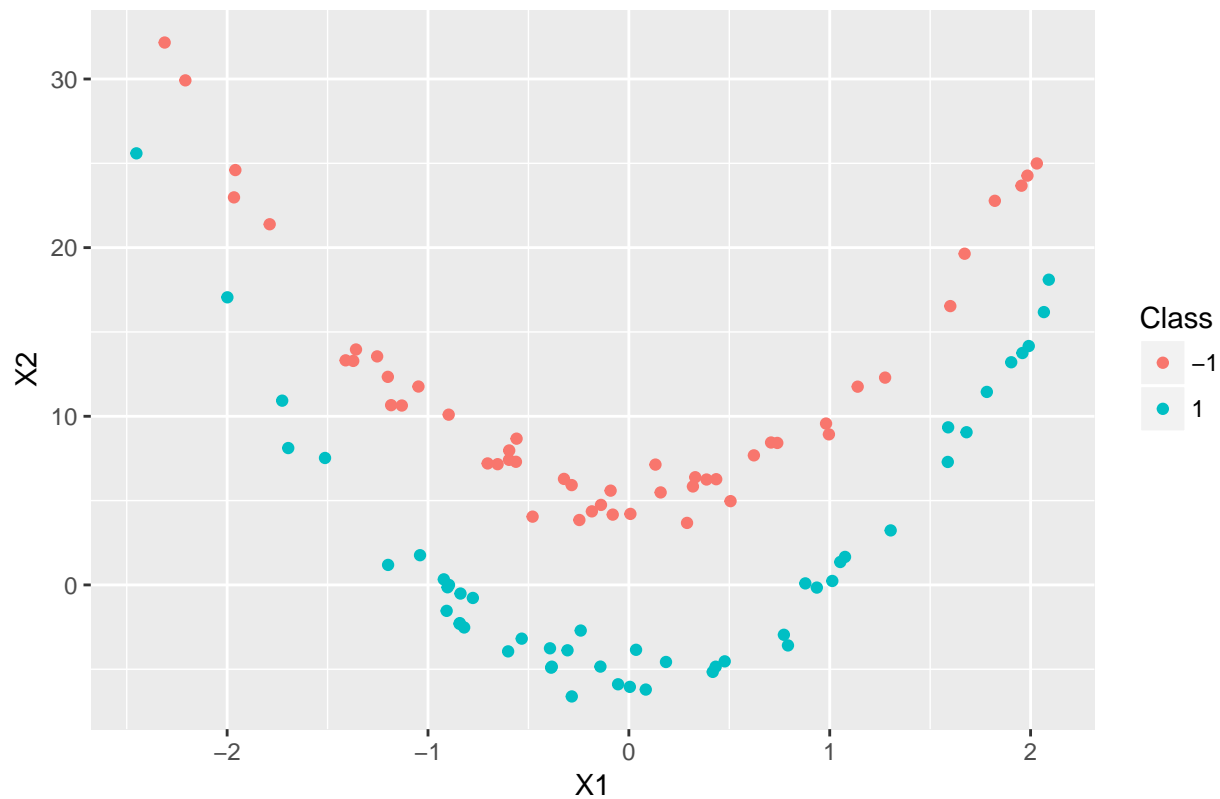
x_2[idx] <- x_2[idx] + 5
x_2[-idx] <- x_2[-idx] - 5
y <- factor(ifelse(x_2 > 5*(x_1^2) + rnorm(100), -1, 1), levels = c(-1, 1))

df <- data.frame(x.1 = x_1,
                 x.2 = x_2,
                 y = y)

# plot data
ggplot(df, aes(x.1, x.2, color = y)) +
  geom_point() +
  ggtitle("Maximal Marginal Classifier NOT Possible") +
  xlab("X1") +
  ylab("X2") +
  labs(color = 'Class')

```

Maximal Marginal Classifier NOT Possible



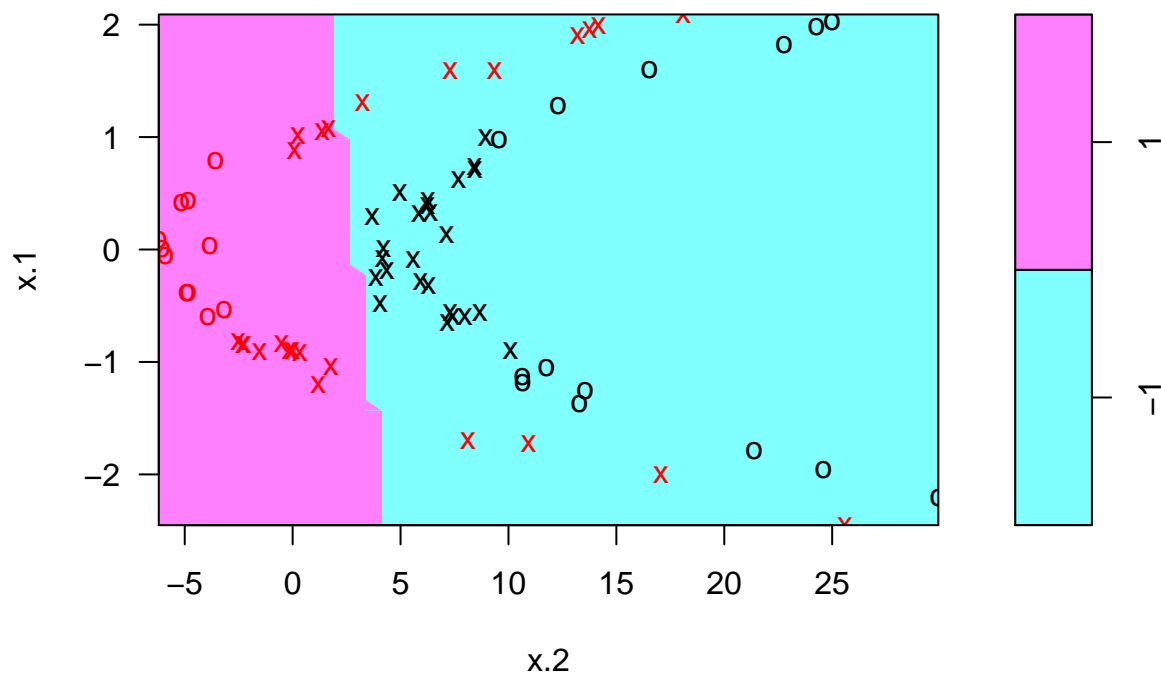
```

set.seed(5)
train <- sample(100, 75)

# linear support vector machine
linear.svm <- svm(y ~ .,
                  data = df[train,],
                  kernel = 'linear')
plot(linear.svm, df[train,])

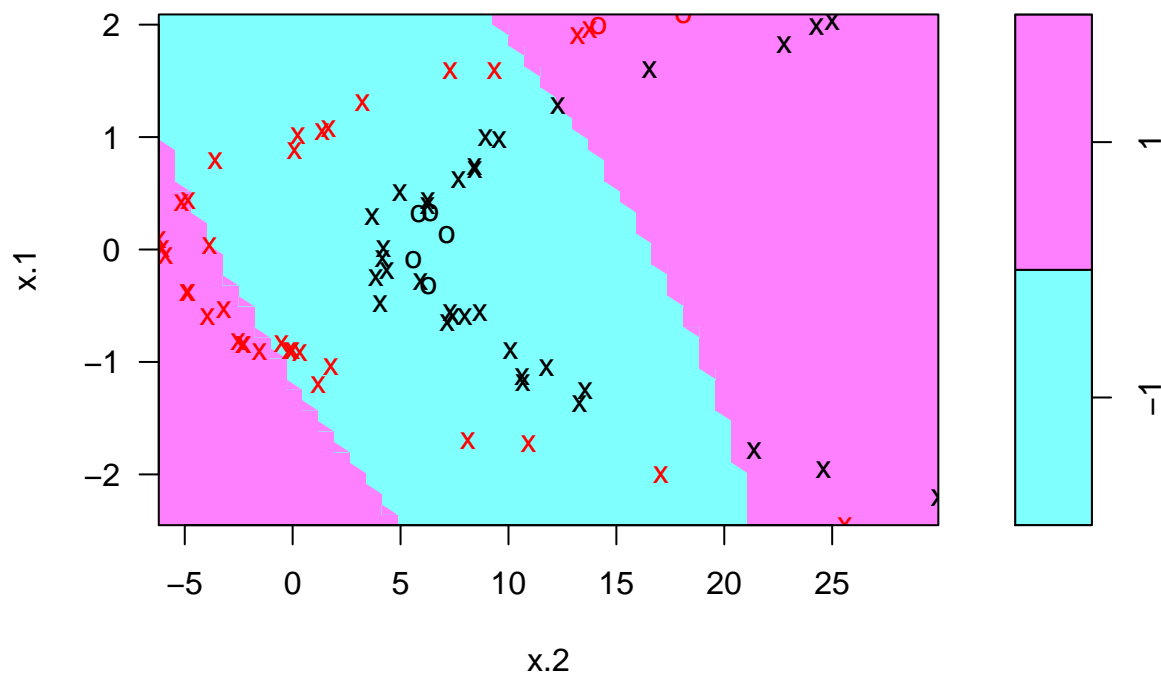
```

SVM classification plot



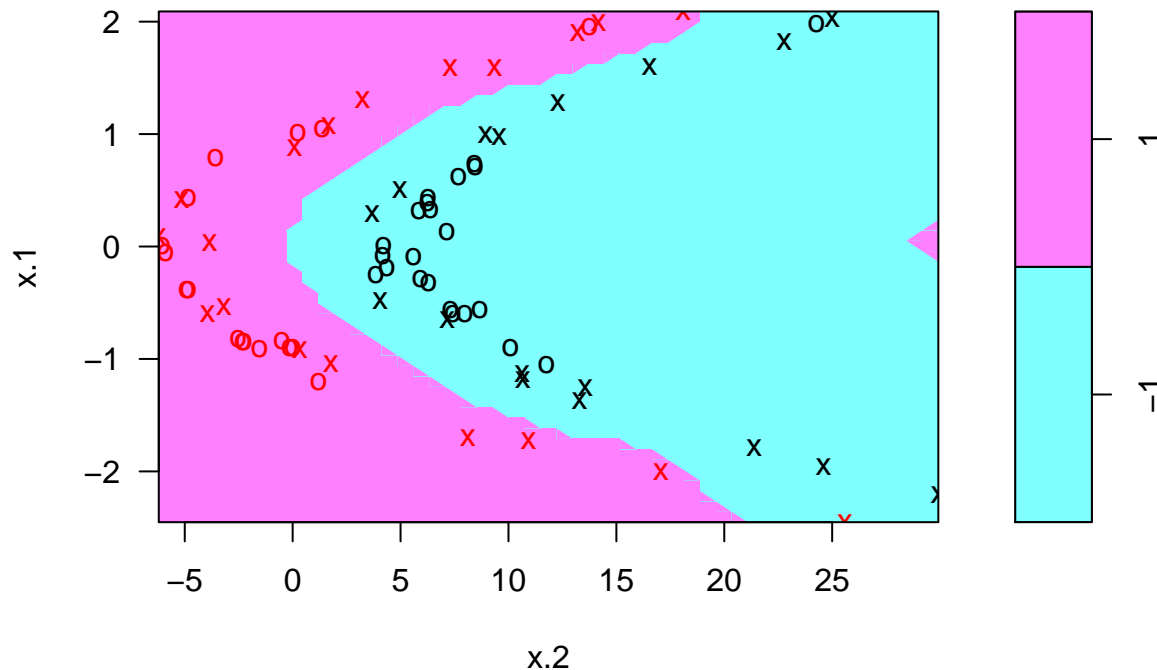
```
# polynomial support vector machine
poly.svm <- svm(y ~ .,
  data = df[train,],
  kernel = 'polynomial',
  degree = 2)
plot(poly.svm, df[train,])
```

SVM classification plot



```
# radial support vector machine
radial.svm <- svm(y ~ .,
  data = df[train,],
  kernel = 'radial',
  gamma = 2)
plot(radial.svm, df[train,])
```

SVM classification plot



```
# prediction
linear.svm.yhat <- predict(linear.svm, newdata = df[-train, ])
poly.svm.yhat <- predict(poly.svm, newdata = df[-train, ])
radial.svm.yhat <- predict(radial.svm, newdata = df[-train, ])

linear.svm.acc <- confusionMatrix(df[-train, 'y'], linear.svm.yhat)$overall[1]
poly.svm.acc <- confusionMatrix(df[-train, 'y'], poly.svm.yhat)$overall[1]
radial.svm.acc <- confusionMatrix(df[-train, 'y'], radial.svm.yhat)$overall[1]

print(paste("Linear SVM test accuracy =",
            linear.svm.acc,
            "with",
            linear.svm$tot.nSV,
            "support vectors."))
```

```
## [1] "Linear SVM test accuracy = 0.84 with 50 support vectors."
```

```
print(paste("Polynomial SVM test accuracy =",
            poly.svm.acc,
            "with",
            poly.svm$tot.nSV,
            "support vectors."))
```

```
## [1] "Polynomial SVM test accuracy = 0.64 with 68 support vectors."
```

```
print(paste("Radial SVM test accuracy =",
            radial.svm.acc,
            "with",
            radial.svm$tot.nSV,
            "support vectors."))
```

```
## [1] "Radial SVM test accuracy = 1 with 36 support vectors."
```

5

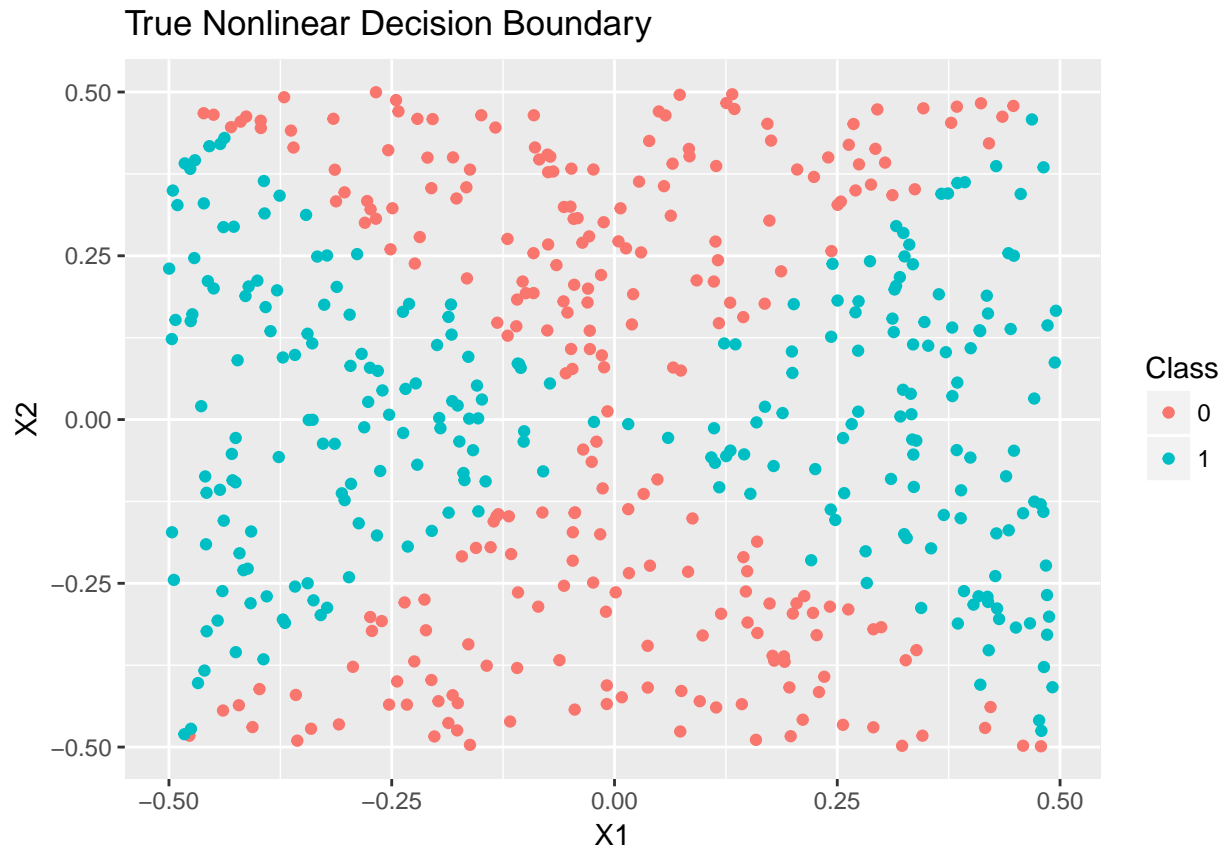
- A.

```
x1 <- runif(500) - 0.5
x2 <- runif(500) - 0.5
y <- 1 * (x1^2 - x2^2 > 0)
```

- B.

```
df <- data.frame(x1 = x1,
                 x2 = x2,
                 y = factor(y, levels = c(0, 1)))

ggplot(df, aes(x1, x2, color = y)) +
  geom_point() +
  ggtitle("True Nonlinear Decision Boundary") +
  xlab("X1") +
  ylab("X2") +
  labs(color = "Class")
```



- C.

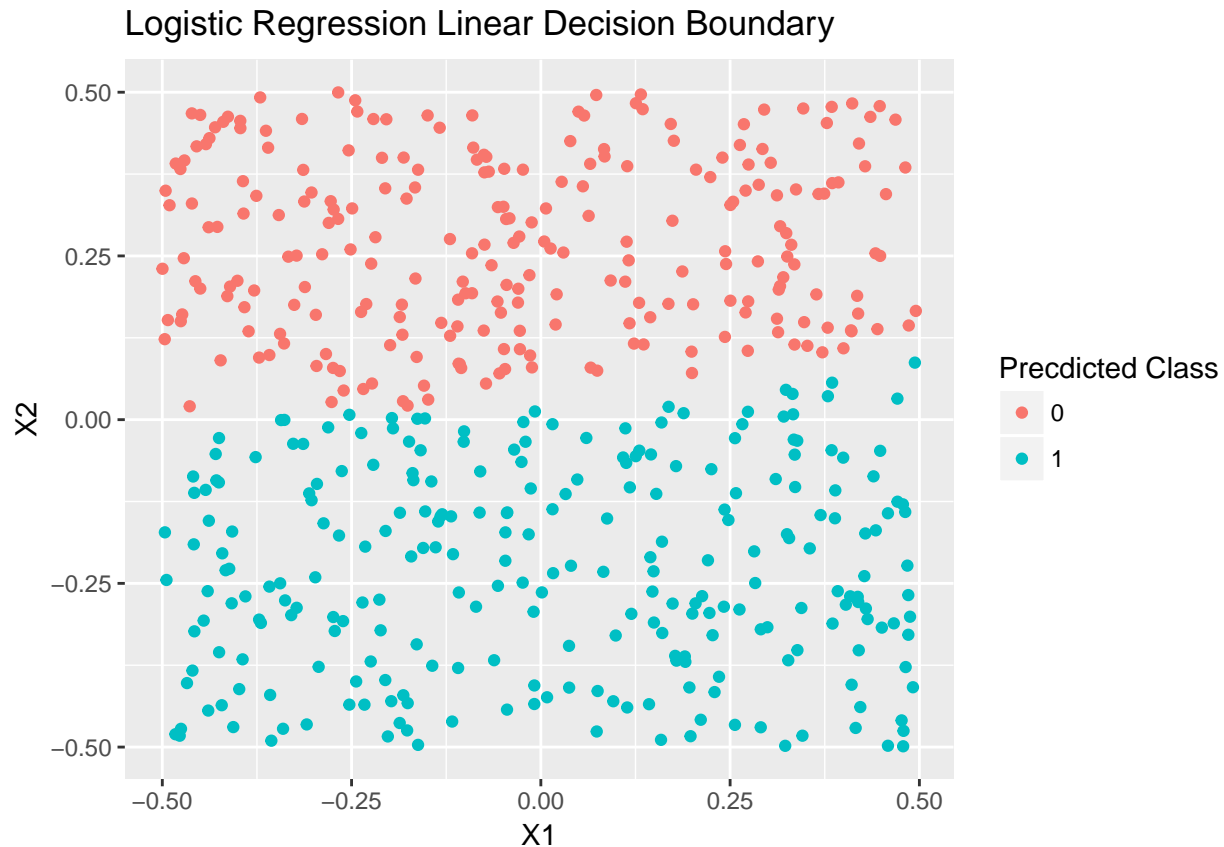
```
log.mod <- glm(y ~ .,
               data = df,
               family = 'binomial')
```

- D.

```
log.mod.y.prob <- predict(log.mod, newdata = df, type = "response")
threshold <- 0.5
log.mod.y.hat <- ifelse(log.mod.y.prob >= threshold, 1, 0)

df$y_hat <- factor(log.mod.y.hat, levels = c(0, 1))

ggplot(df, aes(x1, x2, color = y_hat)) +
  geom_point() +
  ggtitle("Logistic Regression Linear Decision Boundary") +
  xlab("X1") +
  ylab("X2") +
  labs(color = "Predicted Class")
```

• E.

```
non.linear.log.mod <- glm(y ~ poly(x1, 2) + poly(x2, 2) + x1*x2,
  data = df,
  family = 'binomial')

## Warning: glm.fit: algorithm did not converge
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

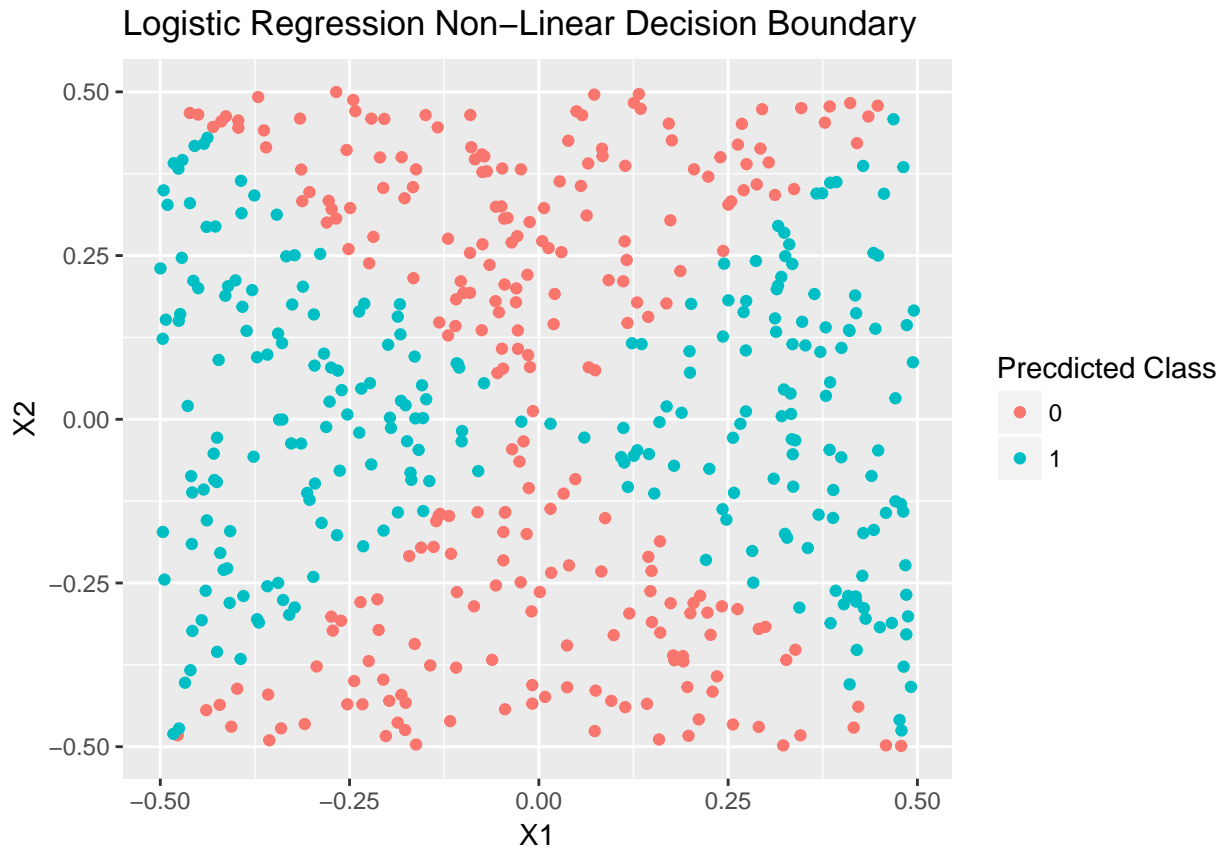
• F.

```
non.linear.log.mod.y.prob <- predict(non.linear.log.mod,
  newdata = df,
  type = 'response')

## Warning in predict.lm(object, newdata, se.fit, scale = 1, type =
## ifelse(type == : prediction from a rank-deficient fit may be misleading
threshold <- 0.5
non.linear.log.mod.y.hat <- ifelse(non.linear.log.mod.y.prob >= threshold,
  1, 0)

df$non_lin_y_hat <- factor(non.linear.log.mod.y.hat, levels = c(0, 1))
```

```
ggplot(df, aes(x1, x2, color = non_lin_y_hat)) +
  geom_point() +
  ggtitle("Logistic Regression Non-Linear Decision Boundary") +
  xlab("X1") +
  ylab("X2") +
  labs(color = "Predicted Class")
```

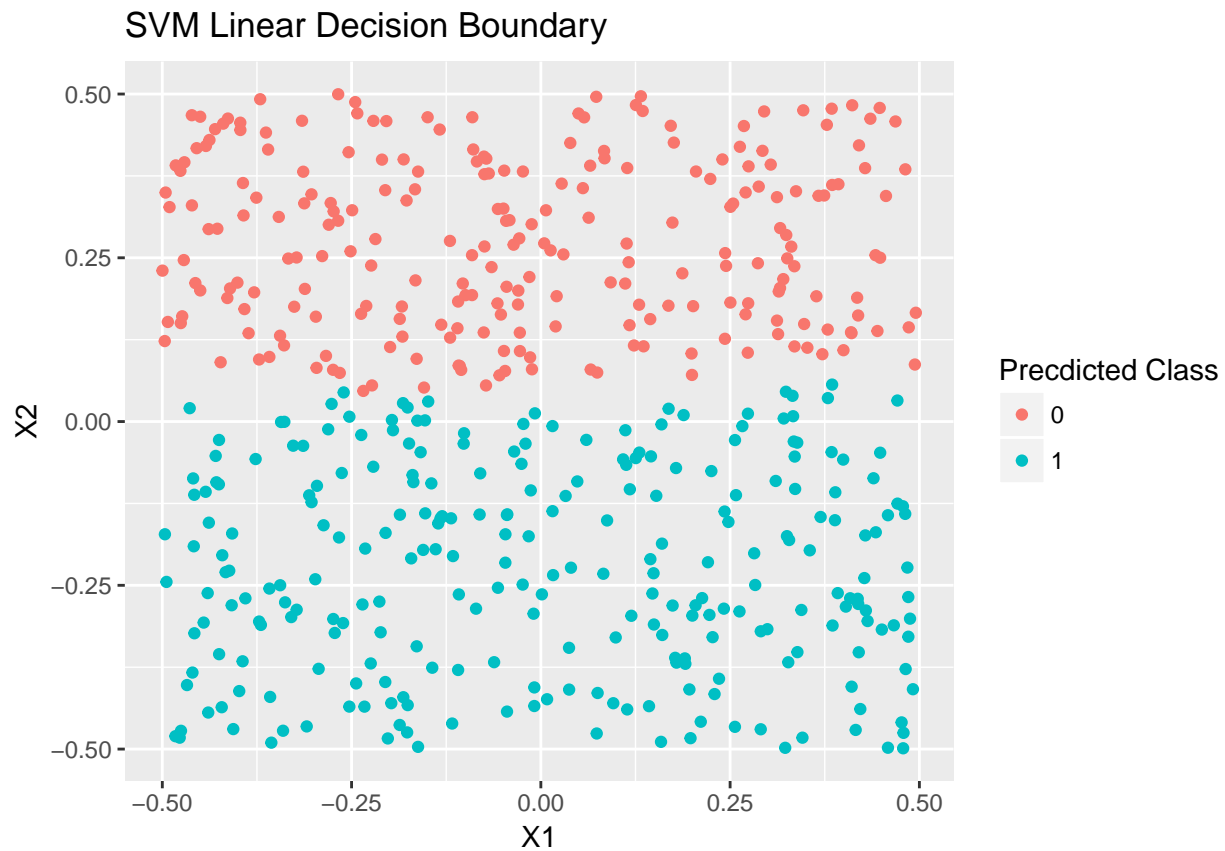


- G.

```
svm.mod <- svm(y ~ x1 + x2,
  data = df,
  kernel = 'linear')
svm.mod.y.hat <- predict(svm.mod, newdata = df)

df$svm_y_hat <- svm.mod.y.hat

ggplot(df, aes(x1, x2, color = svm_y_hat)) +
  geom_point() +
  ggtitle("SVM Linear Decision Boundary") +
  xlab("X1") +
  ylab("X2") +
  labs(color = "Predicted Class")
```

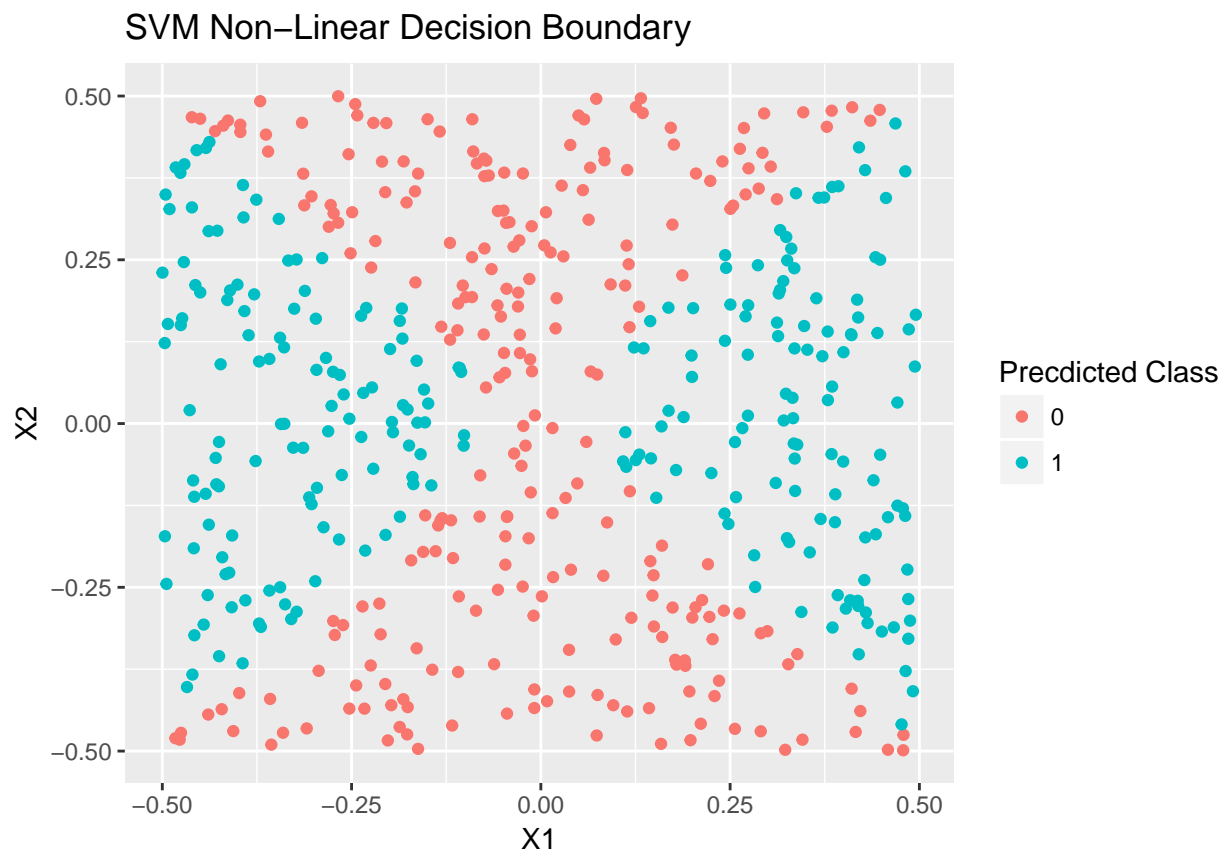


- H.

```
svm.mod.nonlin <- svm(y ~ x1 + x2,
  data = df,
  kernel = 'polynomial',
  degree = 2)
svm.mod.nonlin.y.hat <- predict(svm.mod.nonlin, newdata = df)

df$svm_y_hat_nonlin <- svm.mod.nonlin.y.hat

ggplot(df, aes(x1, x2, color = svm_y_hat_nonlin)) +
  geom_point() +
  ggtitle("SVM Non-Linear Decision Boundary") +
  xlab("X1") +
  ylab("X2") +
  labs(color = "Predicted Class")
```



- I. The comparison between SVM's and Logistic Regression above illustrates an interesting concept when seeking to model some natural phenomena; instead of toggling the different statistical learning methods that one would use to model said phenomena, *transforming the attributes that are fed into a statistical learning method can have a similar affect on modeling that phenomena.*

Typically, one is given a matrix \mathbf{X} and a vector y , and a variety of models/functions, $f_1(\mathbf{X}, y)$, $f_2(\mathbf{X}, y)$, $f_3(\mathbf{X}, y)$, $f_j(\mathbf{X}, y)$, are used to try and model the relationship between \mathbf{X} and y .

However, as introduced in chapter 7, *basis functions* can be used to transform the attributes of \mathbf{X} into a new matrix, $b(\mathbf{X})$, where $b_j(x_j)$ is a mathematical transformation of the j^{th} column of \mathbf{X} (illustrated below).

$$\mathbf{X} = \begin{pmatrix} x_{1,1} & x_{1,2} & x_{1,3} & \cdots & x_{1,p} \\ x_{2,1} & x_{2,2} & x_{2,3} & \cdots & x_{2,p} \\ x_{3,1} & x_{3,2} & x_{3,3} & \cdots & x_{3,p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & x_{n,3} & \cdots & x_{n,p} \end{pmatrix}$$

$$b(\mathbf{X}) = \begin{pmatrix} b_1(x_{1,1}) & b_2(x_{1,2}) & b_3(x_{1,3}) & \cdots & b_\infty(x_{1,\infty}) \\ b_1(x_{2,1}) & b_2(x_{2,2}) & b_3(x_{2,3}) & \cdots & b_\infty(x_{2,\infty}) \\ b_1(x_{3,1}) & b_2(x_{3,2}) & b_3(x_{3,3}) & \cdots & b_\infty(x_{3,\infty}) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ b_1(x_{n,1}) & b_2(x_{n,2}) & b_3(x_{n,3}) & \cdots & b_\infty(x_{n,\infty}) \end{pmatrix}$$

This opens up up a whole new world of possibilities, where the challenge becomes finding the proper basis functions as opposed to finding the proper model framework. Once these basis functions are found, the model framework can vary and still produce similar results, as shown by pitting Logistic Regression against SVM's with the same basis functions, and obtaining similar decision boundaries.

6

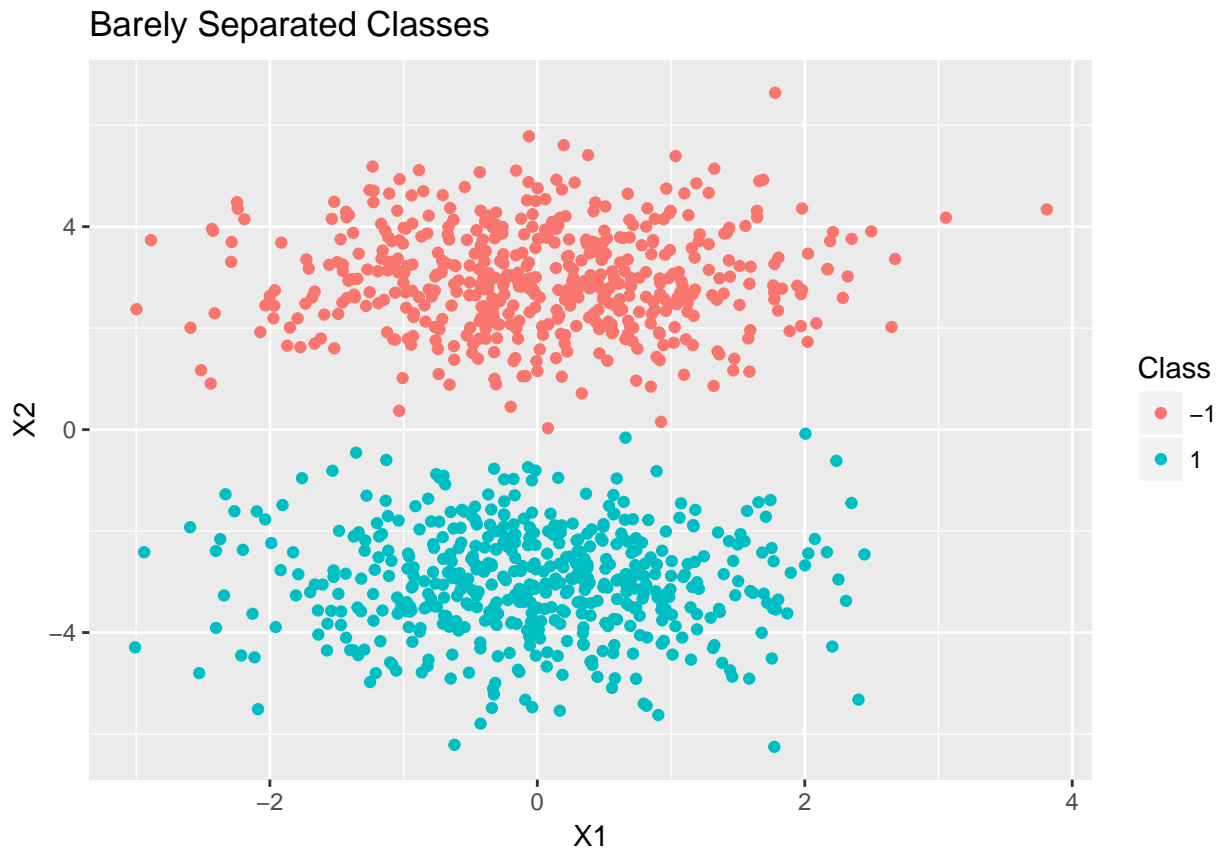
- A.

```
# set seed for reproducibility
set.seed(1)
# randomly generate data from a Gaussian Distribution
x_1 <- rnorm(1000)
x_2 <- rnorm(1000)
# create an index that will be used to separate the classes
idx <- sample(1000, 500)
x_2[idx] <- x_2[idx] + 3
x_2[-idx] <- x_2[-idx] - 3

# create separating hyperplane at 0.5
hyperplane <- 0
y <- ifelse(x_2 <= hyperplane, 1, -1)

train.df <- data.frame(x1 = x_1,
                       x2 = x_2,
                       y = factor(y, levels = c(-1, 1)))
```

```
# plot to illustrate barely separated classes
ggplot(train.df, aes(x1, x2, color = y)) +
  geom_point() +
  ggtitle("Barely Separated Classes") +
  xlab("X1") +
  ylab("X2") +
  labs(color = 'Class')
```



- **B.** As one can see from the table below, cost values of 0.1, 1 and 2 all corresponded to 0 training errors.

```
cost.values <- c(0.1, 1, 5, 10, 50, 100, 1000, 10000)

tuned.svm <- tune(svm,
  y ~ .,
  data = train.df,
  kernel = 'linear',
  ranges = list(cost = cost.values))

error.df <- data.frame(Cost = tuned.svm$performances$cost,
  Train_MisClass = tuned.svm$performances$error * 1000)
error.df

##      Cost Train_MisClass
## 1 1e-01                2
```

## 2	1e+00	1
## 3	5e+00	2
## 4	1e+01	3
## 5	5e+01	4
## 6	1e+02	3
## 7	1e+03	0
## 8	1e+04	0

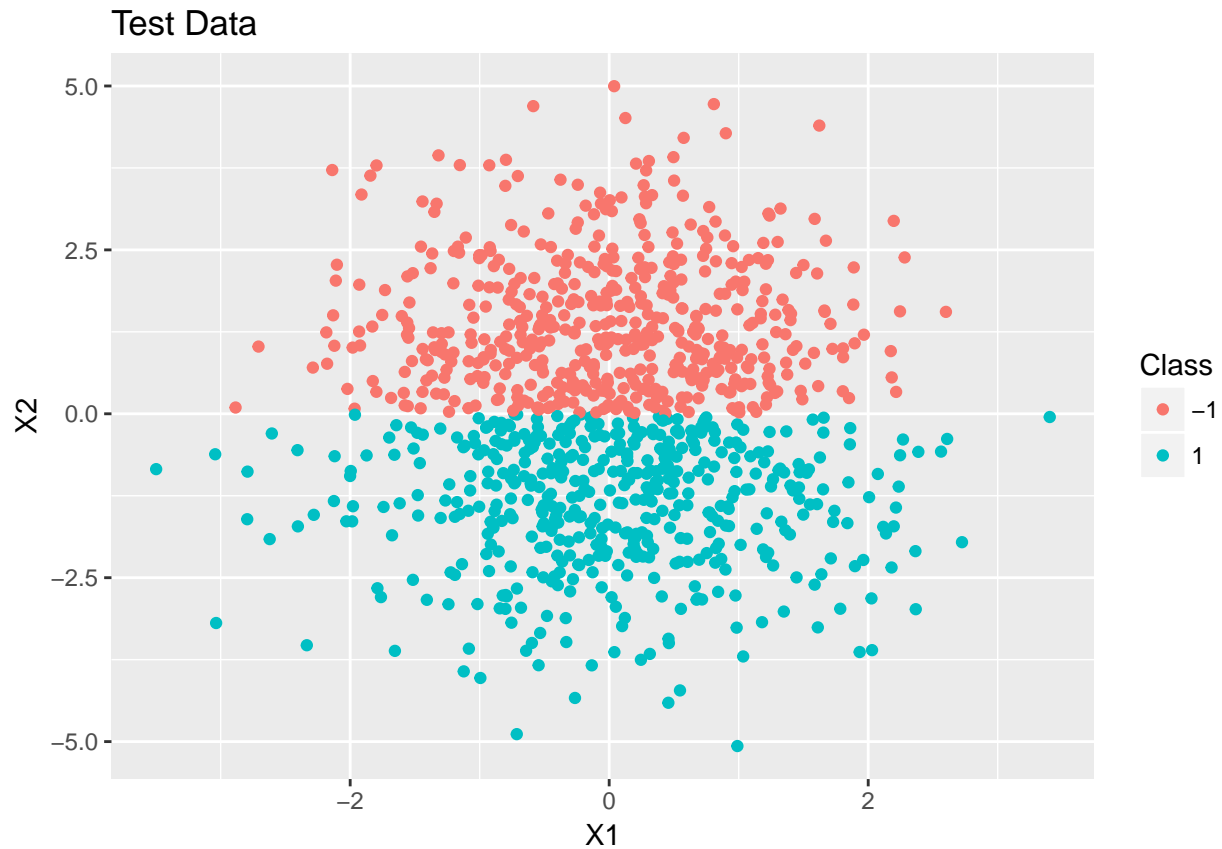
- C.

```
# generate test data by changing seed
set.seed(5)
x_1 <- rnorm(1000)
x_2 <- rnorm(1000) + rnorm(1000)
idx <- sample(1000, 500)
x_2[idx] <- x_2[idx] + 1
x_2[-idx] <- x_2[-idx] - 1

hyperplane <- 0
y <- ifelse(x_2 <= hyperplane, 1, -1)

test.df <- data.frame(x1 = x_1,
                      x2 = x_2,
                      y = factor(y, levels = c(-1, 1)))

# plot to illustrate barely separated classes
ggplot(test.df, aes(x1, x2, color = y)) +
  geom_point() +
  ggtitle("Test Data") +
  xlab("X1") +
  ylab("X2") +
  labs(color = 'Class')
```



```
# prediction
test.errors <- rep(0, length(cost.values))
for (i in 1:length(cost.values)) {
  svm.mod <- svm(y ~ .,
    data = train.df,
    kernel = 'linear',
    cost = cost.values[i])
  y.hat <- predict(svm.mod, newdata = test.df)
  test.errors[i] <- sum(y.hat != test.df$y)
}
```

```
error.df$Test_MisClass <- test.errors
error.df
```

```
##      Cost Train_MisClass Test_MisClass
## 1 1e-01           2         12
## 2 1e+00           1         11
## 3 5e+00           2          8
## 4 1e+01           3         28
## 5 5e+01           4         16
## 6 1e+02           3         24
## 7 1e+03           0         28
## 8 1e+04           0         28
```


- **D.** The point that is trying to be driven home here is one of the Bias-Variance Trade off; even if a Support Vector Classifier is able to correctly classify all *training* observations due to a high cost of violations to the margin/hyperplane, a *different* set of data (aka **test** data) might not be as cleanly separated as the training data, leading to overfitting if a high cost model is chosen. This is shown comparing the training misclassifications to the testing misclassifications, where the high cost SVM clearly overfits the data.