

HW4 SML

Jie Li j15246

February 28, 2019

Problem 1

Bayes Classifier:

$$\hat{y}_0 = \text{Max } P(Y = y|X = x_0)$$

Under 0-1 loss, Minimizing Bayes Error Rate:

$$\text{Min } E\left[\frac{1}{m} \sum_i^m I(y_i \neq \hat{y}_i)\right] = 1 - E[\text{Max } P(Y = y|X = x_0)]$$

Bayes Decision Boundary:

$$P(Y = y|X = x_0) = 1 - P(Y = y|X = x_0) = 1/2$$

Since the Bayes classifier will always choose the class for which is the largest expectation averages of probability over all possible values of X

Problem 2

a)

```
library(ggplot2)
library(purrr)
library(ggforce)
library(MASS)
library(factoextra)
library(glmnet)

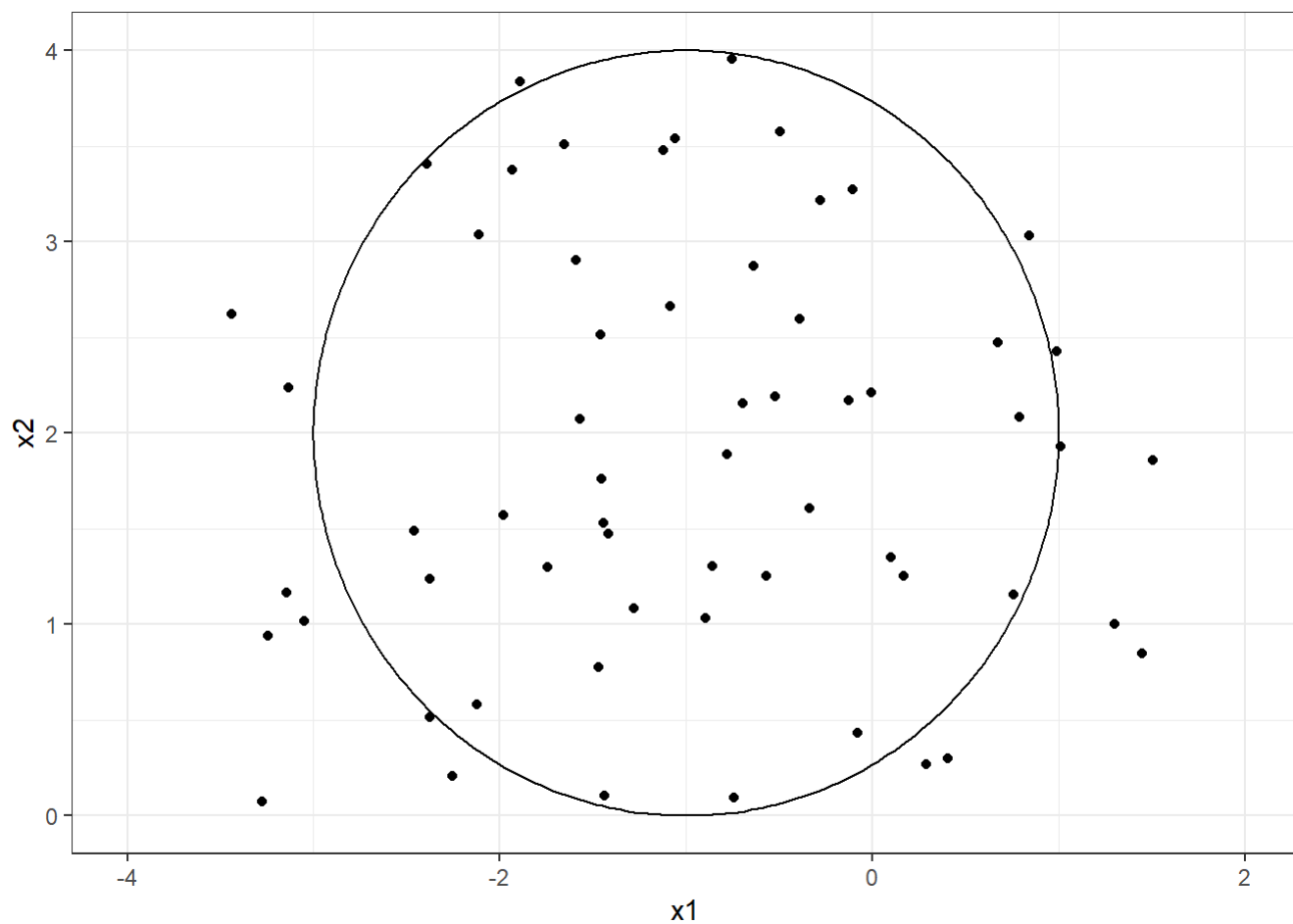
set.seed(123)
setwd("C:/Users/jay48/OneDrive/Documents/work/Statistical ML/HW3")

x1 = rnorm(100, -1, 2)
x2 = rnorm(100, 2, 2)

data = data.frame(x1, x2)

ggplot(data)+
  geom_point(aes(x1, x2), pch = 19)+
  geom_circle(aes(x0=-1, y0=2, r=2))+
  xlim(-4, 2)+
  ylim(0, 4)+
  theme_bw()
```

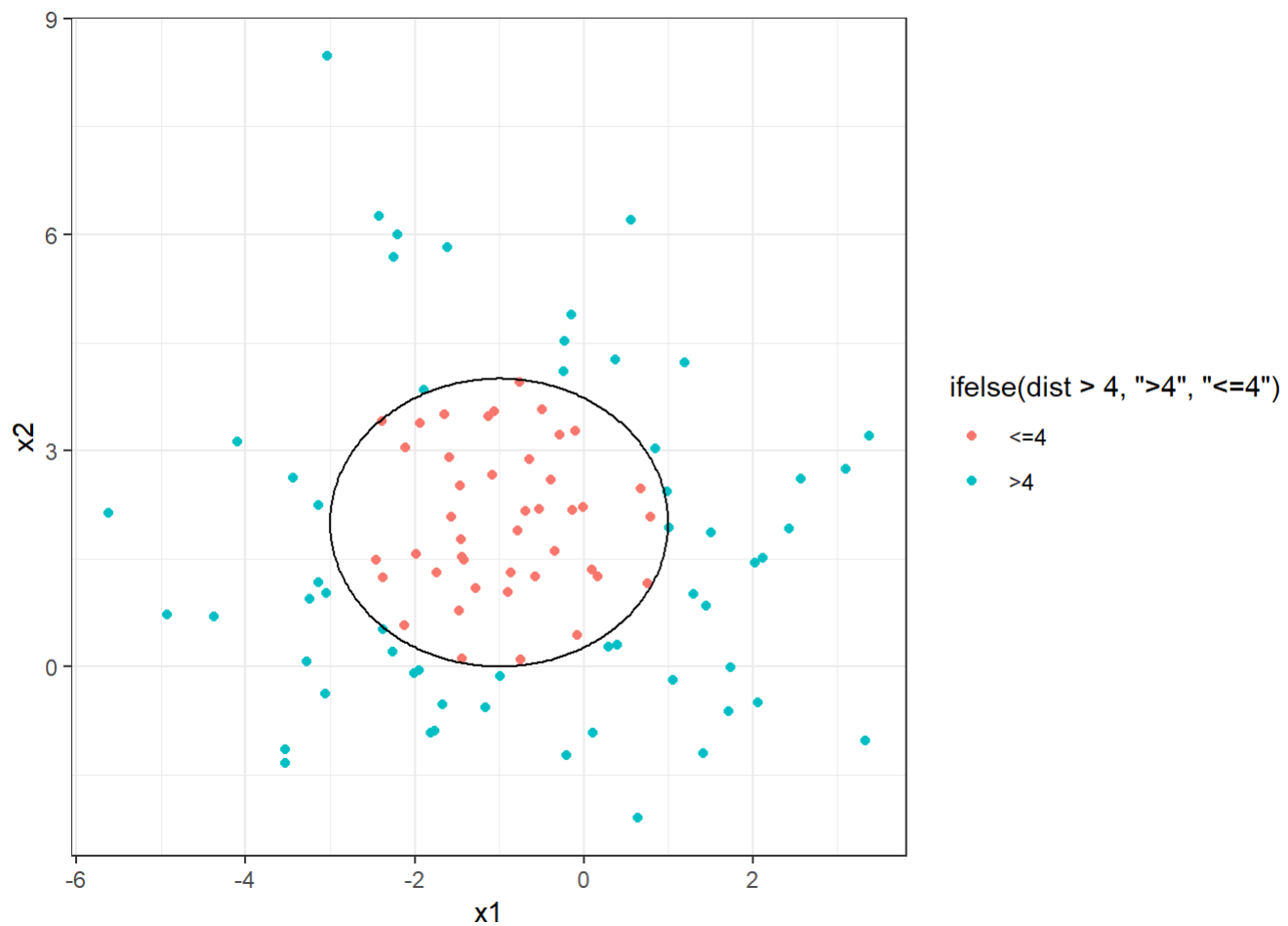
```
## Warning: Removed 40 rows containing missing values (geom_point).
```



b)

```
data$dist = (data$x1+1)^2 + (2-data$x2)^2

ggplot(data)+
  geom_point(aes(x1, x2, col = ifelse(dist > 4,'>4','<=4')), pch = 19)+
  geom_circle(aes(x0=-1, y0=2, r=2))+
  theme_bw()
```



c)

observation (0, 0): $f(0, 0) = 5 > 4 \rightarrow \text{Blue}$

observation (-1, 1): $f(-1, 1) = 1 < 4 \rightarrow \text{Red}$

observation (2, 2): $f(2, 2) = 9 > 4 \rightarrow \text{Blue}$

observation (3, 8): $f(3, 8) = 52 > 4 \rightarrow \text{Blue}$

d)

Decision Boundary: $1 + 2x_1 - 4x_2 + x_1^2 + x_2^2 = 0$. If only in terms of X_1 and X_2 , then the decision boundary is non-linear because of second-order terms. However, if in terms of X_1, X_2, X_1^2, X_2^2 , then the decision boundary is linear because all coef are treated as linear.

Problem 3

```

### all images corresponding to digit "3"
zip.3 = read.table("./train_3.txt", header = FALSE, sep = ",")
zip.3 = cbind(zip.3, y=rep(3, 658)) %>% data.frame()

### all images corresponding to digit "5"
zip.5 = read.table("./train_5.txt", header = FALSE, sep = ",")
zip.5 = as.matrix(cbind(zip.5, y=rep(5, 556)))

### all images corresponding to digit "8"
zip.8 = read.table("./train_8.txt", header = FALSE, sep = ",")
zip.8 = as.matrix(cbind(zip.8, y=rep(8, 542)))

train = data.frame(rbind(zip.3, zip.5, zip.8))
test = read.table("./zip_test.txt", header = F, sep = " ")
test$y = test[,1]
names(test)[2:257] = names(test)[1:256]
test = test[test$y == c(3, 5, 8), 2:258]

# output.image = function(vector)
# {
#   digit = as.matrix(vector, nrow = 16, ncol = 16)
#   index = seq(from = 16, to = 1, by = -1)
#   sym_digit = digit[, index]
#   image(sym_digit, col = gray((8:0)/8), axes = FALSE)
# }
#
# par(mfrow = c(5, 5), mai = c(0.1, 0.1, 0.1, 0.1))
# for(i in 1:25)
# {
#   output.image(zip.3[i, ])
# }

```

1)

```

score = function(model, train, test)
{
  fit = predict(model, train)
  pred = predict(model, test)

  cat("Training Confusion Matrix: \n")
  print(table(y=train$y, fit$class))

  cat("Testing Confusion Matrix: \n")
  print(table(y=test$y, pred$class))

  cat("Misclassification Error: \n")
  accuracy = c(1-mean(fit$class == train$y), 1-mean(pred$class == test$y))
  names(accuracy) = c("Training", "Testing")
  return(accuracy)
}

# scaled.train = scale(train, center = T, scale = T) %>% data.frame()
# scaled.test = scale(test, center = T, scale = T) %>% data.frame()

model = lda(y ~ ., data=train)
score(model, train, test)

```

```

## Training Confusion Matrix:
##
## y      3   5   8
## 3 644   5   9
## 5   6 549   1
## 8    2   5 535
## Testing Confusion Matrix:
##
## y      3   5   8
## 3 46   5   2
## 5   4 52   0
## 8   1   1 59
## Misclassification Error:

```

```

## Training    Testing
## 0.01594533 0.07647059

```

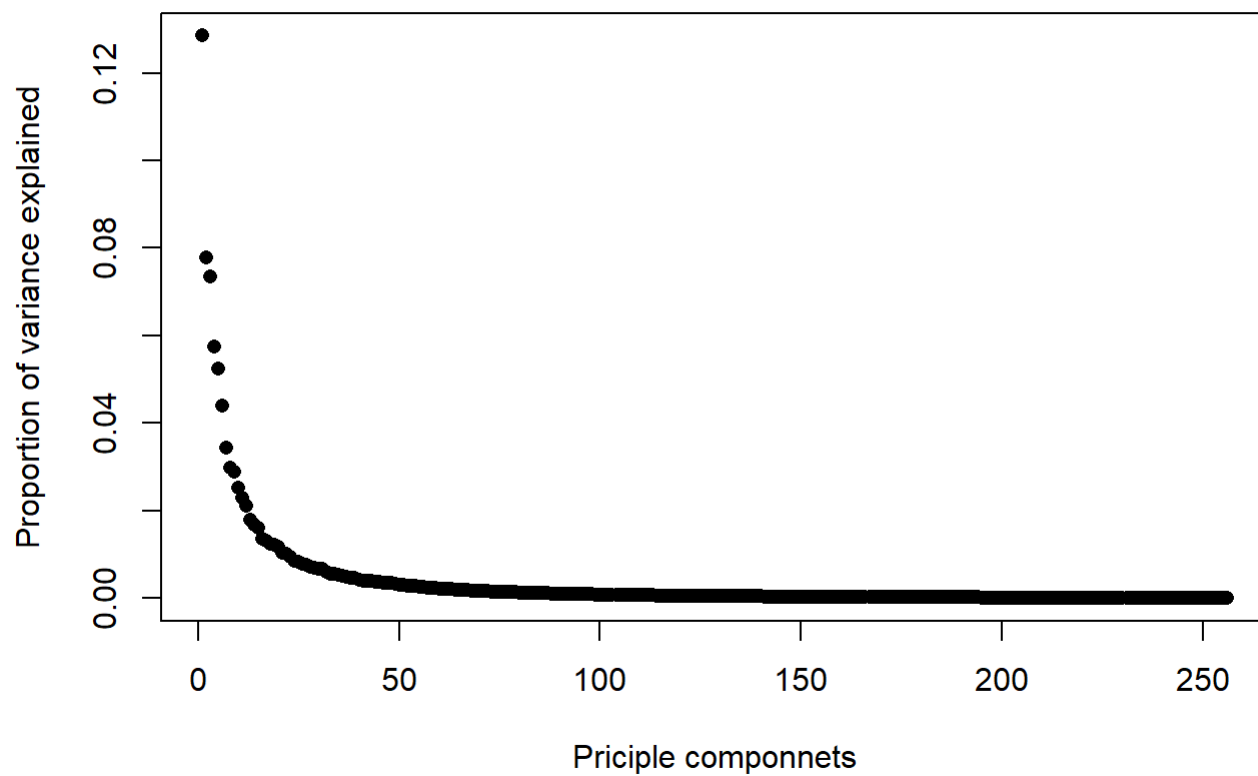
2)

```

scaled.train = scale(train[, -257], center = T, scale = F)
scaled.test = scale(test[, -257], center = T, scale = F)
pca = svd(scaled.train)

plot(seq(from = 1, to = 256, by = 1), (pca$d)^2/sum((pca$d)^2), xlab = "Principle componnets", ylab = "Proportion of variance explained", pch = 16)

```



```
pca.train = scaled.train %>% pca$v[, 1:49]
pca.train = cbind(pca.train, y=train$y) %>% data.frame()
pca.test = cbind(scaled.test[, -257] %>% pca$v[, 1:49], y=test$y) %>% data.frame()

model2 = lda(y ~ ., data=pca.train)
score(model2, pca.train, pca.test)
```

```
## Training Confusion Matrix:
##
## y      3   5   8
## 3 631  16  11
## 5  19 529   8
## 8  12  11 519
## Testing Confusion Matrix:
##
## y      3   5   8
## 3  47   4   2
## 5   5  49   2
## 8   4   1  56
## Misclassification Error:
```

```
## Training Testing
## 0.04384966 0.10588235
```

3)

```

filter = function(data)
{
  n = dim(data)[1]
  data.filter = matrix(NA, n, 64)
  temp = kronecker(diag(8), cbind(c(0.5, 0.5)))
  trans.mat = kronecker(temp, temp)

  for(i in 1:n)
  {
    data.filter[i, ] = as.numeric(data[i, 1:256]) %*% trans.mat
  }

  data.filter = data.frame(cbind(data.filter, y=data[, 257]))
  return(data.filter)
}

train.filter = filter(train)
test.filter = filter(test)

model3 = lda(y ~ ., data=train.filter)
score(model3, train.filter, test.filter)

```

```

## Training Confusion Matrix:
##
## y      3   5   8
## 3 631  14  13
## 5  12 540   4
## 8   9   7 526
## Testing Confusion Matrix:
##
## y      3   5   8
## 3  47   4   2
## 5   2  51   3
## 8   3   1  57
## Misclassification Error:

```

```

## Training Testing
## 0.03359909 0.08823529

```

4)

```
x = as.matrix(train.filter[, -65], 1756, 64)
x.test = as.matrix(test.filter[, -65], 170, 64)
y = train.filter$y

model4 = glmnet(x, y, family = "multinomial", alpha = 0)
fit = predict(model4, x, type = "class", s=0)
pred = predict(model4, x.test, type = "class", s=0)

cat("Training Confusion Matrix: \n")
```

```
## Training Confusion Matrix:
```

```
table(y, fit)
```

```
##      fit
## y      3   5   8
## 3 635  14   9
## 5  11 536   9
## 8   6   7 529
```

```
cat("Testing Confusion Matrix: \n")
```

```
## Testing Confusion Matrix:
```

```
table(y=test.filter$y, pred)
```

```
##      pred
## y      3   5   8
## 3 49   2   2
## 5   1 52   3
## 8   3   1 57
```

```
cat("Misclassification Error: \n")
```

```
## Misclassification Error:
```

```
accuracy = c(1-mean(fit == train.filter$y), 1-mean(pred == test.filter$y))
names(accuracy) = c("Training", "Testing")
accuracy
```

```
## Training Testing
## 0.03189066 0.07058824
```