

Problem 1.

Take the following training/test split: 13000 training and 6020 test. Perform all necessary tests, plots, and Monte Carlo simulations to determine your final choice of classifier. You can use the built-in function of your favorite program. Can you beat the 86.6% mean accuracy on test data (based on 100 runs)?

Solution

Comments are provided in all sample R code near print statements detailing the results of runs of the script.

KNN

K-Nearest Neighbors was my best performing method, performing just shy of 81%, due to the magnitude of the data set, finding the best k value time consuming, even with some parallel processing included to speed up the process.

```

1  #####
2  # Setup #
3  #####
4
5  library(foreach)
6  library(doParallel)
7
8  data <- read.csv("magic04.data", header=F, sep=",")
9  train.size <- 13000 # Given by homework specification
10 start.time <- proc.time()
11
12 #####
13 # KNN #
14 #####
15
16 # Try a bunch of different K-values,
17 cl <- makeCluster(8)
18 registerDoParallel(cl)
19
20 err <- foreach (K=1:1000, .combine = c) %dopar% {
21
22     # Need to load the library for knn on each thread.
23     library(class)
24
25     # Run KNN 100 times for each K value.
26     # Each run is independent, so we can speed things up a little
27     # bit by running it in parallel.
```

```
28     k.err <- foreach (i=1:1, .combine = c) %do% {
29
30         data <- data[sample(nrow(data)),] # Randomize the data set
31
32         train <- data[1:train.size, 1:10]
33         test <- data[(train.size+1):nrow(data), 1:10]
34         train.cl <- factor(data[1:train.size, 11])
35         test.cl <- factor(data[(train.size+1):nrow(data), 11]);
36
37         predict.cl <- knn(train, test, train.cl, k=K)
38         sum(test.cl != predict.cl) / nrow(test)
39     }
40
41     mean(k.err)
42 }
43
44 stopCluster(cl)
45
46 plot(1:1000, err)
47
48 # This was our best performing k value.
49 k <- which.min(err)
50 min.err <- min(err)
51 acc <- 1.0 - min.err
52
53 # About 80.975%
54 print(paste("Min K: ", k))
55 print(paste("KNN - Accuracy: ", acc))
56 print(proc.time() - start.time)
```

LDA

I had issues running hzTest on the entire data set with size of 19200, but I was able to run it on a random subset of size 15000. So instead I resorted to running it on a random subset of the data.

Both the hzTest and uniPlot indicated that the data was not multi-variate norm. From the plot we can see that some features look almost normal, but overall the data does not follow a normal distribution. Because of this, we can expect LDA and QDA to not perform particularly well.

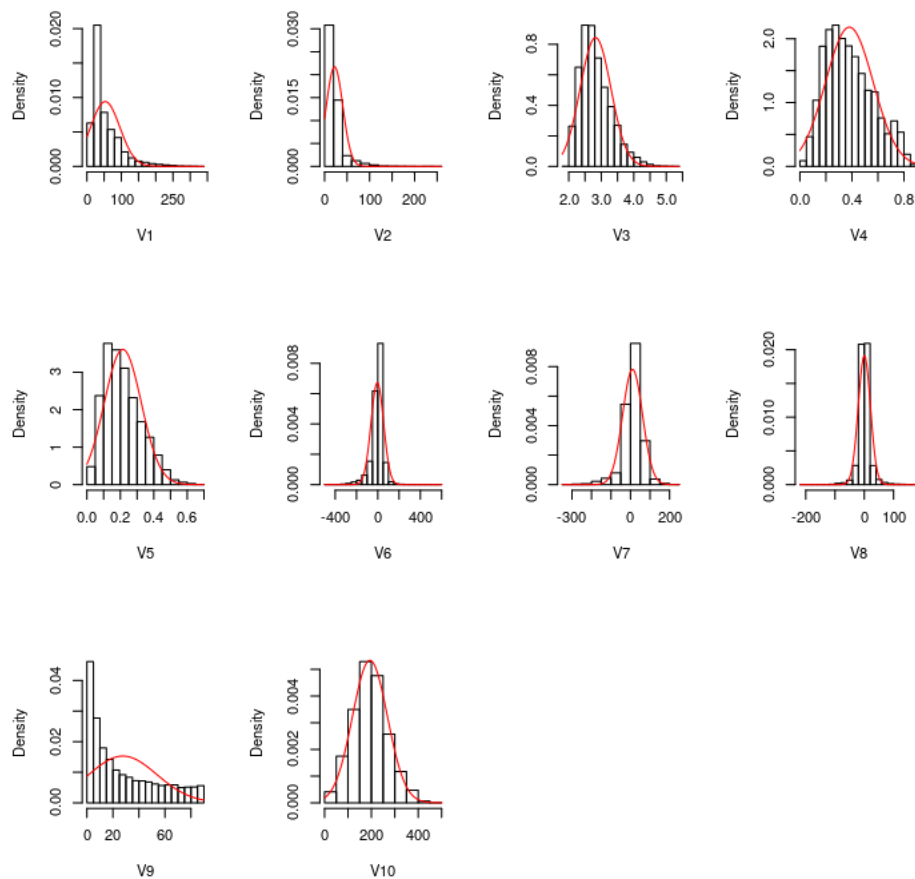
Running these methods anyway produced the expected results both performed at about 78% accuracy on test, which is less than I was able to achieve with KNN.

```
data : data[1:15000, 1:10]
```

```
HZ      : 39.14529
```

```
p-value : 0
```

```
Result  : Data are not multivariate normal.
```



```
1 #####
2 # Setup #
3 #####
4
5 library(foreach)
6 library(doParallel)
7
8 data <- read.csv("magic04.data", header=F, sep=",")
9 train.size <- 13000 # Given by homework specification
```

```
10 start.time <- proc.time()
11
12 #####
13 # Test for Normality #
14 #####
15
16 data <- data[sample(nrow(data)),] # Randomize the data set
17 hz = hzTest(data[1:15000,1:10], cov=TRUE, qqplot=FALSE)
18 print(hz)
19 uniPlot(data, type="histogram")
20
21 #####
22 # LDA #
23 #####
24
25 cl <- makeCluster(4)
26 registerDoParallel(cl)
27
28 err <- foreach (i=1:100, .combine = c) %dopar% {
29   library(MASS)
30
31   data <- data[sample(nrow(data)),] # Randomize the data set
32
33   train <- data[1:train.size, 1:10]
34   test <- data[(train.size+1):nrow(data), 1:10]
35   train.cl <- factor(data[1:train.size, 11])
36   test.cl <- factor(data[(train.size+1):nrow(data), 11]);
37
38   model <- lda(x = train, grouping = train.cl)
39   predict.cl <- predict(model, test)$class
40   sum(test.cl != predict.cl) / nrow(test)
41 }
42
43 stopCluster(cl)
44 acc <- 1.0 - mean(err)
45
46 # About 78.429%
47 print(paste("LDA - Accuracy: ", acc))
48 print(proc.time() - start.time)
```

QDA

Due to the poor performance of LDA, and the non-linear nature of the input data set, we do not expect QDA to perform any better. It may have one assumption less than LDA, the

assumption it keeps still fails. Running the algorithm produces a mean accuracy roughly equivalent to that produced by LDA.

```
1 #####
2 # Setup #
3 #####
4
5 library(foreach)
6 library(doParallel)
7
8 data <- read.csv("magic04.data", header=F, sep=",")
9 train.size <- 13000 # Given by homework specification
10 start.time <- proc.time()
11
12 #####
13 # QDA #
14 #####
15
16 cl <- makeCluster(4)
17 registerDoParallel(cl)
18
19 err <- foreach (i=1:100, .combine = c) %dopar% {
20   library(MASS)
21
22   data <- data[sample(nrow(data)),] # Randomize the data set
23
24   train <- data[1:train.size, 1:10]
25   test <- data[(train.size+1):nrow(data), 1:10]
26   train.cl <- factor(data[1:train.size, 11])
27   test.cl <- factor(data[(train.size+1):nrow(data), 11]);
28
29   model <- qda(x = train, grouping = train.cl)
30   predict.cl <- predict(model, test)$class
31   sum(test.cl != predict.cl) / nrow(test)
32 }
33
34 stopCluster(cl)
35 acc <- 1.0 - mean(err)
36
37 # About 78.4276%
38 print(paste("QDA - Accuracy: ", acc))
39 print((proc.time() - start.time))
```

Naive Bayes (Normal)

Naive Bayes was less interesting, again, as was the case with LDA, the normal assumption of the data set failed, but I went ahead and ran it anyway. The results were unexciting, serving as the worst performing classifier over 100 random test cases.

```

1 #####
2 # Setup #
3 #####
4
5 library(foreach)
6 library(doParallel)
7 library(MVN)
8
9 data <- read.csv("magic04.data", header=F, sep=",")
10 train.size <- 13000 # Given by homework specification
11 start.time <- proc.time()
12
13 #####
14 # Naive Bayes (Normal) #
15 #####
16
17 cl <- makeCluster(4)
18 registerDoParallel(cl)
19
20 err <- foreach (i=1:100, .combine = c) %dopar% {
21   library(klaR)
22   library(caret)
23
24   data <- data[sample(nrow(data)),] # Randomize the data set
25
26   train <- data[1:train.size, 1:10]
27   test <- data[(train.size+1):nrow(data), 1:10]
28   train.cl <- factor(data[1:train.size, 11])
29   test.cl <- factor(data[(train.size+1):nrow(data), 11]);
30
31   model <- NaiveBayes(x = train, grouping = train.cl, usekernel=FALSE)
32   predict.cl <- predict(model, test)$class
33   sum(test.cl != predict.cl) / nrow(test)
34 }
35
36 stopCluster(cl)
37 acc <- 1.0 - mean(err)
38
39 # About 72.6714%

```

```

40 print(paste("Naive Bayes (Normal) - Accuracy: ", acc))
41 print(proc.time() - start.time)

```

Naive Bayes (Kernel)

Removing the normal assumption of the input data set increased the processing time required to train a model using Naive Bayes with a Kernel density estimation. But this resulted in an improvement in accuracy over Naive Bayes with a Normal assumption. Despite this it still performed worse than LDA and QDA, much less KNN.

```

1  #####
2  # Setup #
3  #####
4
5  library(foreach)
6  library(doParallel)
7  library(MVN)
8
9  data <- read.csv("magic04.data", header=F, sep=",")
10 train.size <- 13000 # Given by homework specification
11 start.time <- proc.time()
12
13 #####
14 # Naive Bayes (Kernel) #
15 #####
16
17 cl <- makeCluster(4)
18 registerDoParallel(cl)
19
20 err <- foreach (i=1:100, .combine = c) %dopar% {
21   library(klaR)
22   library(caret)
23
24   data <- data[sample(nrow(data)),] # Randomize the data set
25
26   train <- data[1:train.size, 1:10]
27   test <- data[(train.size+1):nrow(data), 1:10]
28   train.cl <- factor(data[1:train.size, 11])
29   test.cl <- factor(data[(train.size+1):nrow(data), 11]);
30
31   model <- NaiveBayes(x = train, grouping = train.cl, usekernel=TRUE)
32   predict.cl <- predict(model, test)$class
33   sum(test.cl != predict.cl) / nrow(test)
34 }

```

```
35
36 stopCluster(cl)
37 acc <- 1.0 - mean(err)
38
39 # About 76.2375%
40 print(paste("Naive Bayes (Kernel) - Accuracy: ", acc))
41 print((proc.time() - start.time))
```

Problem 2.

- (a) Classify the test point $(0, 1)$ using QDA and calculate the posterior class probabilities. Do the calculations by hand.
- (b) Classify the test point $(0, 1)$ using naive Bayes assuming normality and calculate the posterior class probabilities. Do the calculations by hand.
- (c) Verify your results for both classifiers using Matlab, R, etc. You can use the built-in functions.

Solution**Part (a)**

First we need to calculate the class means $\hat{\mu}_0$ and $\hat{\mu}_1$

$$\begin{aligned}
 \hat{\mu}_0 &= \frac{1}{3}([0.6585, 0.2444] + [2.2460, 0.5281] + [-2.7665, -3.8303]) \\
 &= \frac{1}{3}[0.138, -3.8303] \\
 &= [0.046, -1.019267] \\
 \hat{\mu}_1 &= \frac{1}{3}([-1.2565, 3.4912] + [-0.7973, 1.2288] + [1.1170, 2.2637]) \\
 &= \frac{1}{3}[-0.9368, 6.9837] \\
 &= [-0.3122667, 2.3279000]
 \end{aligned}$$

Next calculate the covariance matrix for each class.

$$\begin{aligned}
 \hat{\Sigma}_0 &= \frac{1}{2}(([0.6585, 0.2444] - \hat{\mu}_0)([0.6585, 0.2444] - \hat{\mu}_0)^\top \\
 &\quad + ([2.2460, 0.5281] - \hat{\mu}_0)([2.2460, 0.5281] - \hat{\mu}_0)^\top \\
 &\quad + ([-2.7665, -3.8303] - \hat{\mu}_0)([-2.7665, -3.8303] - \hat{\mu}_0)^\top) \\
 &= \frac{1}{2} \left(\begin{bmatrix} 0.33516 & 0.77400 \\ 0.77400 & 1.59685 \end{bmatrix} + \begin{bmatrix} 4.8400 & 3.4042 \\ 3.4042 & 2.394 \end{bmatrix} + \begin{bmatrix} 7.9102 & 7.9060 \\ 7.9060 & 7.9019 \end{bmatrix} \right) \\
 &= \begin{bmatrix} 6.5627 & 6.0421 \\ 6.0421 & 5.9466 \end{bmatrix}
 \end{aligned}$$

$$\begin{aligned}
\hat{\Sigma}_1 &= \frac{1}{2}(([-1.2565, 3.4912] - \hat{\mu}_1)([-1.2565, 3.4912] - \hat{\mu}_1)^\top \\
&\quad + ([-0.7973, 1.2288] - \hat{\mu}_1)([-0.7973, 1.2288] - \hat{\mu}_1)^\top \\
&\quad + ([1.1170, 2.2637] - \hat{\mu}_1)([1.1170, 2.2637] - \hat{\mu}_1)^\top \\
&= \frac{1}{2} \begin{bmatrix} 0.89158 & -1.09843 \\ -1.09843 & 1.35327 \end{bmatrix} + \begin{bmatrix} 0.23526 & 0.53310 \\ 0.53310 & 1.20802 \end{bmatrix} + \begin{bmatrix} 2.0426033 & -0.0917589 \\ -0.0917589 & 2.0426033 \end{bmatrix} \\
&= \begin{bmatrix} 1.58482 & -0.32854 \\ -0.32854 & 1.28270 \end{bmatrix}
\end{aligned}$$

We would like to maximize $\hat{P}[Y = k] \hat{f}(X = x|Y = k)$ over $k \in \{0, 1\}$ for our input $x=(0,1)$.

$$\arg \max_{k \in \{0,1\}} \hat{P}[Y = k] \left(\frac{1}{(2\pi)^{\frac{d}{2}} |\hat{\Sigma}_k|^{\frac{1}{2}}} \right) \exp\left(\frac{1}{2}(x - \hat{\mu}_k)^\top \hat{\Sigma}_k^{-1}(x - \hat{\mu}_k)\right)$$

case: k = 0

$$\begin{aligned}
\hat{P}[Y = 0] &= \frac{3}{6} = \frac{1}{2} \\
|\hat{\Sigma}_0| &= 2.5180 \Rightarrow |\hat{\Sigma}_0|^{\frac{1}{2}} = 1.5868 \\
\left(\frac{1}{(2\pi)^{\frac{2}{2}} |\hat{\Sigma}_0|^{\frac{1}{2}}} \right) &= \left(\frac{1}{9.9702} \right) = 0.10030 \\
\exp\left(\frac{1}{2}(x - \hat{\mu}_0)^\top \hat{\Sigma}_0^{-1}(x - \hat{\mu}_0)\right) &= \\
\exp\left(-\frac{1}{2} * 11.078\right) &= \exp(-5.5389) = 0.0039308 \\
\hat{P}[Y = 0] \hat{f}(X = x|Y = 0) &= .00019713
\end{aligned}$$

case: k = 1

$$\begin{aligned}
\hat{P}[Y = 1] &= \frac{3}{6} = \frac{1}{2} \\
|\hat{\Sigma}_1| &= 1.9249 \Rightarrow |\hat{\Sigma}_1|^{\frac{1}{2}} = 1.3874 \\
\left(\frac{1}{(2\pi)^{\frac{2}{2}} |\hat{\Sigma}_1|^{\frac{1}{2}}} \right) &= \frac{1}{8.7174} = 0.11471 \\
\exp\left(-\frac{1}{2}(x - \hat{\mu}_1)^\top \hat{\Sigma}_1^{-1}(x - \hat{\mu}_1)\right) &= \\
\exp\left(-\frac{1}{2} * 1.3752\right) &= \exp(-0.6876) = 0.50278 \\
\hat{P}[Y = 1] \hat{f}(X = x|Y = 1) &= 0.028838
\end{aligned}$$

Noting that $0.00019713 < 0.028838$, $k=1$ maximizes our function, thus we predict a class label of 1 for $(0,1)$. The posterior class probability is given by the following function.

$$\begin{aligned} P[Y = k|X = x] &= \frac{P[Y = k]f_{1\dots d}(X = x|Y = k)}{\sum_{i=0}^k P[Y = i]f_{Y\dots d}(X = x|Y = i)} \\ &= \frac{0.028838}{0.028838 + 0.00019713} \\ &= 0.99321 \end{aligned}$$

Thus we find that we have a posterior class probability of 99.321% for class $k=1$.

Part (b)

The naive Bayes classification uses the same class norms $\hat{\mu}_0$ and $\hat{\mu}_1$. However we will need to compute $\hat{\sigma}^2$ values for each feature of each class.

case: $k = 0$

First compute the variance of each feature.

$$\begin{aligned} \hat{\sigma}_0^2 &= \frac{1}{3} * \sum_{j \in C_0} (x_j - \hat{\mu}_j)^2 \\ &= \frac{1}{3} (((0.6585, 0.2444) - \hat{\mu}_0)^2 + ((2.2460, 0.5281) - \hat{\mu}_0)^2 + ((-2.7665, -3.8303) - \hat{\mu}_0)^2) \\ &= \frac{1}{3} (13.125, 11.893) \\ &= (4.3751, 3.9644) \end{aligned}$$

Now we can compute $f_i((0, 1)|Y = 0)$ for use in our classifier.

$$\begin{aligned} f((0, 1)|Y = 0) &= \frac{1}{\sqrt{2\pi\hat{\sigma}_0^2}} e^{-\frac{(\langle 0, 1 \rangle - \hat{\mu}_0)^2}{2\hat{\sigma}_0^2}} \\ &= \frac{1}{(5.2431, 4.9909)} (0.99976, 0.59794) \\ &= (0.19068, 0.11981) \end{aligned}$$

We would like to maximize the following equation over all k $\hat{P}[Y = 0] \prod_{i=1}^d f_i(X_i|Y = k)$.

$$\begin{aligned} \hat{P}[Y = 0] &= \frac{3}{6} \\ f_0(X_0|Y = 0) &= 0.19068 \\ f_1(X_1|Y = 0) &= 0.11981 \\ \frac{1}{2} * 0.19068 * 0.11981 &= 0.011423 \end{aligned}$$

case: k = 1

First compute the variance of each feature.

$$\begin{aligned}
 \hat{\sigma}_1^2 &= \frac{1}{3} * \sum_{j \in C_1} (x_i - \hat{\mu}_i)^2 \\
 &= \frac{1}{3} (((-1.2565, 3.4912) - \hat{\mu}_1)^2 + ((-0.7973, 1.2288) - \hat{\mu}_1)^2 + (1.1170, 2.2637) - \hat{\mu}_1)^2) \\
 &= \frac{1}{3} (3.1696, 2.5654) \\
 &= (1.05655, 0.85514)
 \end{aligned}$$

Now we can compute $f_i((0, 1)|Y = 1)$ for use in our classifier.

$$\begin{aligned}
 f((0, 1)|Y = 1) &= \frac{1}{\sqrt{2\pi\hat{\sigma}_1^2}} e^{-\frac{(\langle 0, 1 \rangle - \hat{\mu}_1)^2}{2\hat{\sigma}_1^2}} \\
 &= \frac{1}{(2.5765, 2.3180)} (0.95490, 0.35664) \\
 &= (0.37062, 0.15386)
 \end{aligned}$$

We would like to maximize the following equation over all k $\hat{P}[Y = 1] \prod_{i=1}^d f_i(X_i|Y = k)$.

$$\begin{aligned}
 \hat{P}[Y = 1] &= \frac{3}{6} \\
 f_0(X_0|Y = 0) &= 0.37062 \\
 f_1(X_1|Y = 0) &= 0.15386 \\
 \frac{1}{2} * 0.37062 * 0.15386 &= 0.028512
 \end{aligned}$$

Noting that $0.011423 < 0.028512$, k=1 maximizes our function, thus we predict a class label of 1 for (0,1). The posterior class probability is given by the following function.

$$\begin{aligned}
 P[Y = k|X = x] &= \frac{P[Y = k]f_{1\dots d}(X = x|Y = k)}{\sum_{i=0}^k P[Y = i]f_{Y\dots d}(X = x|Y = i)} \\
 &= \frac{0.028512}{0.028512 + 0.011423} \\
 &= 0.71396
 \end{aligned}$$

Thus we find that we have a posterior class probability of 71.396% for class k=1.

Part (c)

I have included the results of the R script that follows as comments directly below the corresponding print statements. The result of `qda(...)` in R produces the exact same posterior class probability as I have computed above. Currently my results for `NaiveBayes(...)` are off by about 4%, so I must have an error in my work somewhere.

```

1 #####
2 # Setup #
3 #####
4
5 library(MASS)
6 library(klaR)
7
8 data <- data.frame(
9     c(0.6585, 2.2460, -2.7665, -1.2565, -0.7973, 1.1170),
10    c(0.2444, 0.5281, -3.8303, 3.4912, 1.2288, 2.2637),
11    c(0, 0, 0, 1, 1, 1)
12 )
13
14 test <- data.frame(c(0), c(1))
15
16 names(data) <- c("F1", "F2", "CLASS")
17 names(test) <- c("F1", "F2")
18
19 train.size <- nrow(data)
20 train <- data[1:train.size, 1:2]
21 train.cl <- factor(data[1:train.size, 3])
22
23 #####
24 # QDA #
25 #####
26
27 model <- qda(x = train, grouping = train.cl)
28 predict <- predict(model, test)
29 print(predict)
30 print(paste("QDA: ", predict$class))
31
32 # Posterior Class Probabilities
33 # 0: 0.0067895
34 # 1: 0.9932105
35
36 #####
37 # Naive Bayes #
38 #####

```

```
39
40 model <- NaiveBayes(x=train, grouping=train.cl, usekernel=FALSE, fL=0)
41 predict <- predict(model, test)
42 print(predict)
43 print(paste("Naive Bayes: ", predict$class))
44
45 # Posterior Class Probabilities
46 # 0: 0.2493135
47 # 1: 0.7506865
```
