

Computer Assignment 6 - Classification

Machine Learning, Spring 2020

YOUR NAME

Thinking further about initial cluster centers

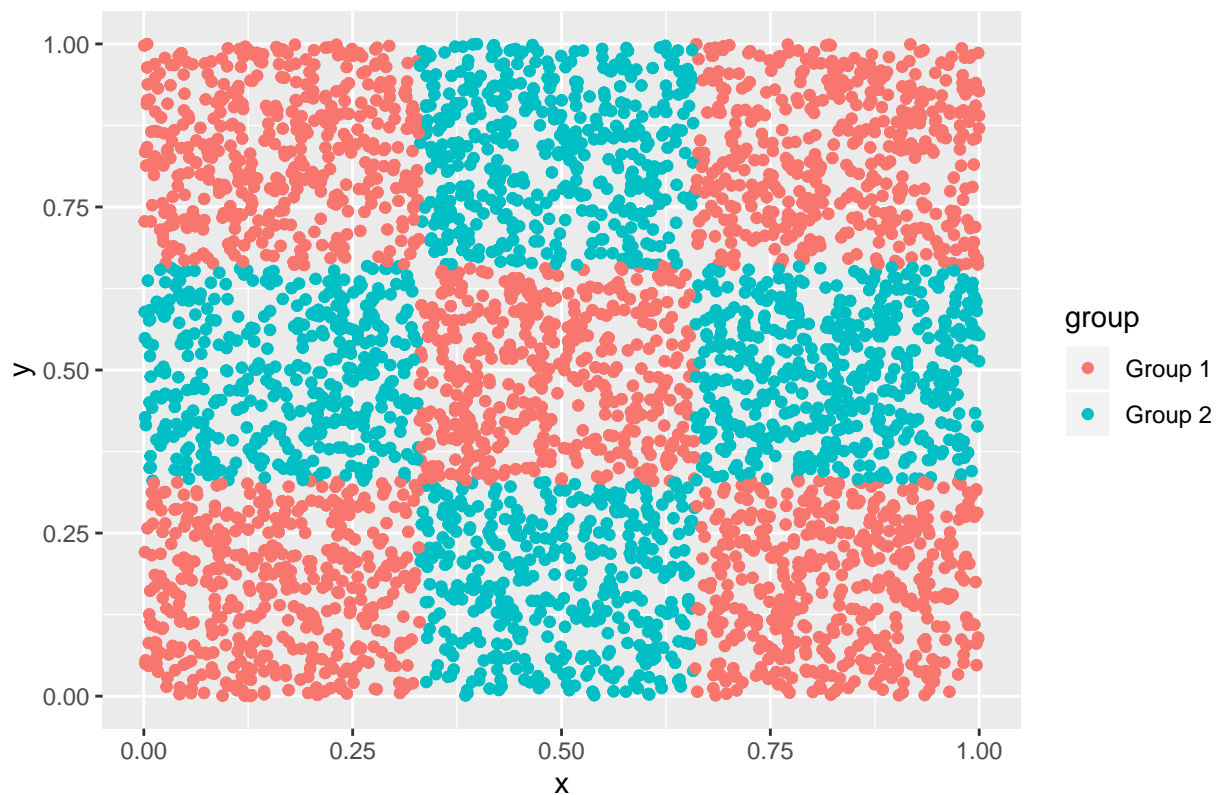
Greetings! The last problem on CA04 was met with some difficulty – many of you were unable to see the total within-cluster sum of squares change with different cluster centers. While it worked for some of you (and, for me, when I wrote the problem), ultimately it seems the kmeans algorithm worked too well even though we restricted it to a single updating step, and it found the optimal clustering regardless. The important punchline is that exhaustively finding the optimal cluster centers is a HARD problem and that we instead find the local maximum. The following is a repeat of the last problem on CA04, except with a new data set (and a `set.seed()` command for extra security) I created to truly show a change in total within-cluster sum of squares. -Murph

When passing a number x as the argument for `centers` to the kmeans algorithm, x cluster centers are chosen randomly, points are assigned to a cluster based on these centers, and then the cluster centers are iteratively updated in an attempt to find the centers that minimize total within-cluster sum of squares. The number of times `kmeans` updates is controlled by the argument `iter.max`. To examine how the choice of initial centers affects the algorithm, we will disable the updating steps via `iter.max` and choose some initial cluster centers of our own.

```
library(ggplot2)
# First, we create a dummy data set
set.seed(10)
random_data = data.frame(x = runif(4000),
                          y = runif(4000),
                          group = rep(NA, times = 8000))
random_data[which(random_data$x < 0.33 & random_data$y < 0.33 |
                  random_data$x > 0.66 & random_data$y < 0.33 |
                  random_data$x < 0.33 & random_data$y > 0.66 |
                  random_data$x > 0.66 & random_data$y > 0.66 |
                  (random_data$x > 0.33 & random_data$x < 0.66 &
                   random_data$y > 0.33 & random_data$y < 0.66)),]$group = "Group 1"
random_data[which(is.na(random_data$group)),]$group = "Group 2"

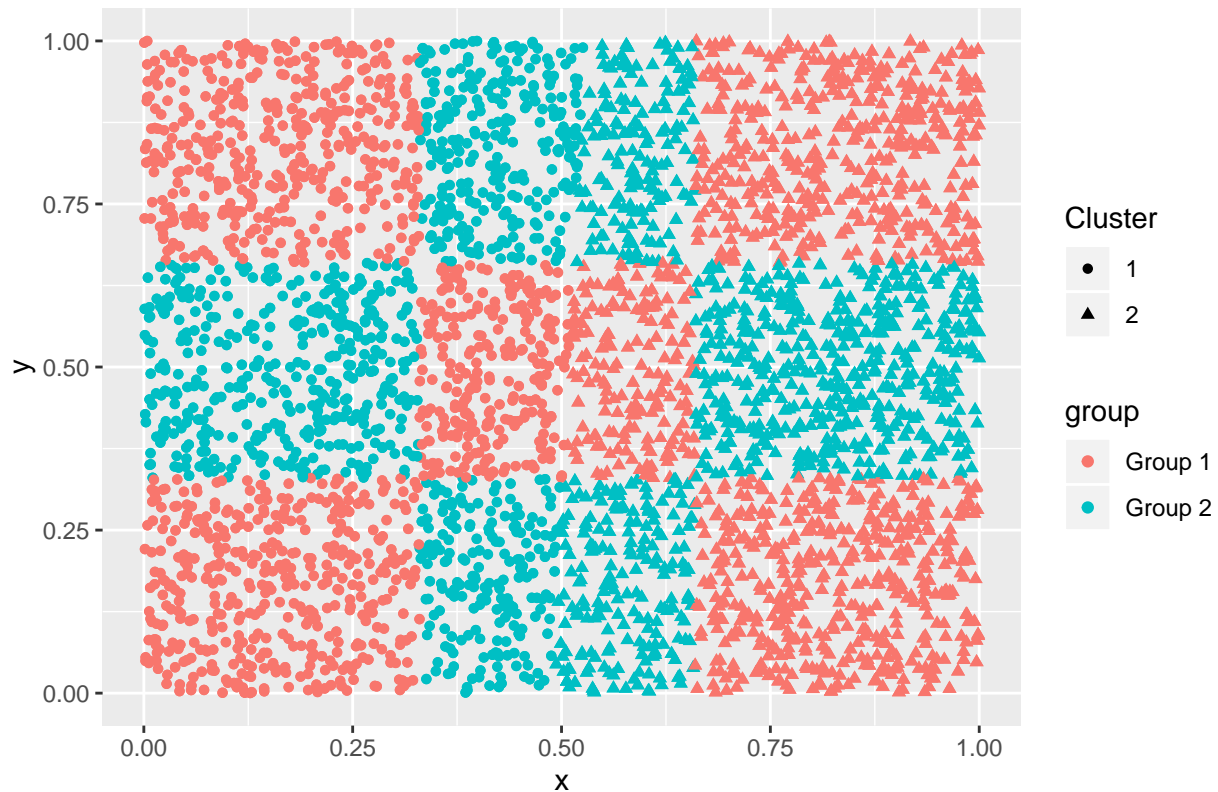
# This data set looks like:
ggplot(random_data, aes(x = x, y = y, col = group)) + geom_point() +
  ggtitle("Plot of Dummy Data")
```

Plot of Dummy Data



```
# We'll perform the kmeans with two random centers:
random.km = kmeans(random_data[,1:2], centers = 2, iter.max = 1)
ggplot(random_data, aes(x = x, y = y, col = group, pch = as.factor(random.km$cluster))) +
  geom_point() +
  scale_shape_discrete(name = "Cluster") +
  ggtitle("Plot of Dummy Data; Clustered")
```

Plot of Dummy Data; Clustered



```
# Which gives the following within-cluster sum of squares:  
random.km$tot.withinss
```

```
## [1] 829.424
```

```
my_centers = random_data[1:2, 1:2]  
random.km = kmeans(random_data[,1:2], centers = my_centers, iter.max = 1)  
ggplot(random_data, aes(x = x, y = y, col = group, pch = as.factor(random.km$cluster))) +  
  geom_point() +  
  scale_shape_discrete(name = "Cluster") +  
  ggtitle("Plot of Dummy Data; Clustered")
```



```
random.km$tot.withinss
```

```
## [1] 847.1953
```

Don't just look at the number!! Can you see how the clusters themselves changed? Now, choose your own cluster centers! These *do not* need to be actual observed points from our data. Try something weird, and report the total within-cluster sum of squares. Did it change? Was it better or worse than the random choice? Comment as to why you think that is.

YOUR CODE HERE

YOUR ANSWER HERE

Bayes Rule and Univariate Normal Simulations

Notice, in the above problem, how we created the `random_data` variable using both the dataframe the `runif` functions. This latter function (`runif`) draws random values from a $\text{uniform}(0, 1)$ distribution. As mentioned in our first few Computing Assignments, R can simulate a number of distributions, including the normal distributions

$$\mathbb{W} \sim \mathcal{N}(-2, 1) \quad \& \quad \mathbb{V} \sim \mathcal{N}(2, 1).$$

For this next exercise, we are going to simulate 300 observations from \mathbb{W} and 200 observations from \mathbb{V} , and create another variable \mathbb{Y} that classifies from which normal distribution each observation came. Indeed,

```
W_obs = rnorm(300, mean = -1)
V_obs = rnorm(200, mean = 1)
```

```
Y_class = c(rep(0, times = length(W_obs)), rep(1, times = length(V_obs)))
train_data = data.frame(X = c(W_obs, V_obs), Y = Y_class)
```

Question: Why did we not specify the value for standard deviation in the `rnorm` function?

YOUR ANSWER HERE

Since we have specified ourselves the model for our data (here: two different normals), we can assess the performance of any classification technique we use on this data. We will illustrate this by calculating the Bayes rule for our `train_data`. To do so, fill in the following quantities

```
pi_0 = YOUR CODE HERE (Probability of W_obs)
pi_1 = YOUR CODE HERE (Probability of V_obs)
```

To start, let's just calculate the Bayes' rule for the first observation of `train_data`. To find the conditional probability, you may use the `density` function in R to estimate the PDF and get the appropriate probability mass for x . That is, let $P(x|Y = 0)$ be the probability mass of x in the sub-population for which $Y = 0$. The following is an example of how to get a probability mass from an estimated density function using the full `training_data`. For your purposes, you will need to use a subset of the `training_data` instead of the full.

```
# Example of PDF estimation
x_obs = train_data[1,1]
full_density = density(train_data$X)
index_of_density = sum(full_density$x <= x_obs)
pdf_value_of_x_obs = full_density$y[index_of_density]
pdf_value_of_x_obs
```

```
x_obs = train_data[1,1]
prob_x_given_0 = YOUR CODE HERE
prob_x_given_1 = YOUR CODE HERE
Bayes_rule_for_x = prob_x_given_1 * pi_1 / (prob_x_given_1 * pi_1 + prob_x_given_0 * pi_0)
```

What hypothesis does this Bayes Rule test? Based on our calculation, in which distribution should we classify `x_obs`?

YOUR ANSWER HERE

Now, calculate the Bayes Rule for every value, and use them to compute a classifier for every observation in `train_data`. (DO NOT do this exhaustively. You should be using built-in features in R and/or a for loop.) Compare these classifiers with the true classifiers. Calculate, and report, the Bayes' Risk.

YOUR CODE, AND ANALYSIS, HERE

k-nearest Neighbors and LDA

Using the same `train_data` from the last exercise, fit a k -nearest neighbors model for $k \in \{1, 3, 10\}$. The code for $k = 1$ is provided.

YOUR CODE, AND ANALYSIS, HERE.

```
# Take special care to separate the classifier from the rest of the data when fitting a
# knn model. Consider the following code to give you an idea as to how one does this.
# Note, however, that this code is based off of my naming practices, and may require editing
# depending on your previous code.
train_data_classifiers = as.factor(train_data$Y)
train_data_observations = data.frame(train_data$X)
knn.1 <- knn(train_data_observations, train_data_observations, cl = train_data_classifiers, k=1)
```

```
R_knn_1 = 100 * sum(train_data_classifiers == knn.1)/length(knn.1)
R_knn_1
```

Comment on the performance for the different values of k . Why does $k = 1$ do so well? What is it doing that gives it such great performance?

YOUR ANSWER HERE

Now let's do the same thing with Fisher's Linear Discriminate Analysis (LDA). We have provided the follow code as an example. Assess the risk using the derived predictions (you will need to grab the `class` attribute from this variable).

```
library(MASS)
library(dplyr)
# Fit the model
model <- lda(Y~X, data = train_data)
# Make predictions
predictions <- model %>% predict(train_data)
```

YOUR CODE, AND ANALYSIS, HERE

Method Evaluation

Now that we have explored Bayes' Rule, k-nearest neighbors, and LDA, we will see how each method performs on data that was NOT used to originally set them up. Consider the following new data set drawn from the same random variables \mathbb{W} & \mathbb{V} .

```
W_obs = rnorm(150, mean = -2)
V_obs = rnorm(50, mean = 2)
Y_class_test = c(rep(0, times = length(W_obs)), rep(1, times = length(V_obs)))
test_data = data.frame(X = c(W_obs, V_obs), Y = Y_class_test)
```

Use the information from `train_data` to classify values in the `test_data`, then compare these calculated classes with the true classes found in `Y_class_test`. Report the Bayes' Risk, and compare it to the same metric calculated from k-nearest neighbors and LDA.

Hints:

1. In the case of Bayes rule, you will use the same π_0 & π_1 and calculate the conditional probabilities in the same way, except this time your `x_obs` will be from `test_data`. DO NOT use the class labels from `test_data` ANYWHERE in this calculation (until you evaluate at the end).
2. Use the manual page `?knn` to see how one inputs a different dataset for the "test" parameter. You may have to separate the observations X from the class labels Y .
3. In a similar fashion, you are expected to read the manual pages for the functions used in the LDA process.

YOUR CODE, AND ANALYSIS, HERE

Now, let's do this 1000 more times! During each iteration of the following for loop, use the models you have created to calculate classifiers for the `test_data` and calculate the risk for each.

```
set.seed(13)
all_bayes_risks = c()
all_knn_risks = c()
all_lda_risks = c()
```

```

for(iteration in 1:1000){
  W_obs = rnorm(150, mean = -2)
  V_obs = rnorm(50, mean = 2)
  Y_class_test = c(rep(0, times = length(W_obs)), rep(1, times = length(V_obs)))
  test_data = data.frame(X = c(W_obs, V_obs), Y = Y_class_test)

  YOUR MODEL CODE HERE

  bayes_risk = YOUR CODE HERE
  knn_risk = YOUR CODE HERE
  lda_risk = YOUR CODE HERE

  all_bayes_risks = c(all_bayes_risks, bayes_risk)
  all_knn_risks = c(all_knn_risks, knn_risk)
  all_lda_risks = c(all_lda_risks, lda_risk)
}

```

If done correctly, you should have three vectors of risk values, each from a different classification method. Create an intuitive plot that compares these three values. Make sure this plot compares the values AT THE SAME ITERATION. Our suggestion would be a line plot with the x -axis as the iteration number and the y -axis as the risk value, colored by classification method.

YOUR CODE HERE

Comment on your model

YOUR ANSWER HERE