

Exercise 1: The idea of credit scoring is that you try to classify people into a group that will pay the money back and a group that won't. We'll use a version of a well-known dataset called Kredit, which contained data from 1000 Germans. This version we'll use is from kaggle.

- a. Load the csv `german_credit_data`.

Solution:

Reading in the data we find that the checking and saving account variables contain all the missing values. Simply removing all the rows with missing values reduces our number of observations to 522, which is about half. I am not necessarily sure we should be throwing out so much of our data, for two variables which are not entirely indicative of an individual's credit responsibility. There are a few actions we can take to try and save this data. We could impute the medium value for these missing values, use a discriminant method that is robust to missing values like decision trees, or we could simply drop these categorical predictors. I ended up deciding to remove the variables.

Code:

```
## Reading in the data.
data <- read.csv('german_credit_data.csv', header = TRUE)

## Note that checking and saving account variable contains a lot of NAs.
length(na.omit(data$Saving.accounts))
# [1] 817
length(na.omit(data$Checking.account))
# [1] 606
length(na.omit(data)[,1])

## Dropping these two columns removes all NAs. Perhaps we
## can explore a method more robust to NA, like decision trees
## and bring that data in.
data <- data[, -c(1,7,8)]

## Turning categorical variables into factors.
data$Sex <- as.numeric(as.factor(data$Sex))
data$Housing <- as.numeric(as.factor(data$Housing))
data$Purpose <- as.numeric(as.factor(data$Purpose))
```

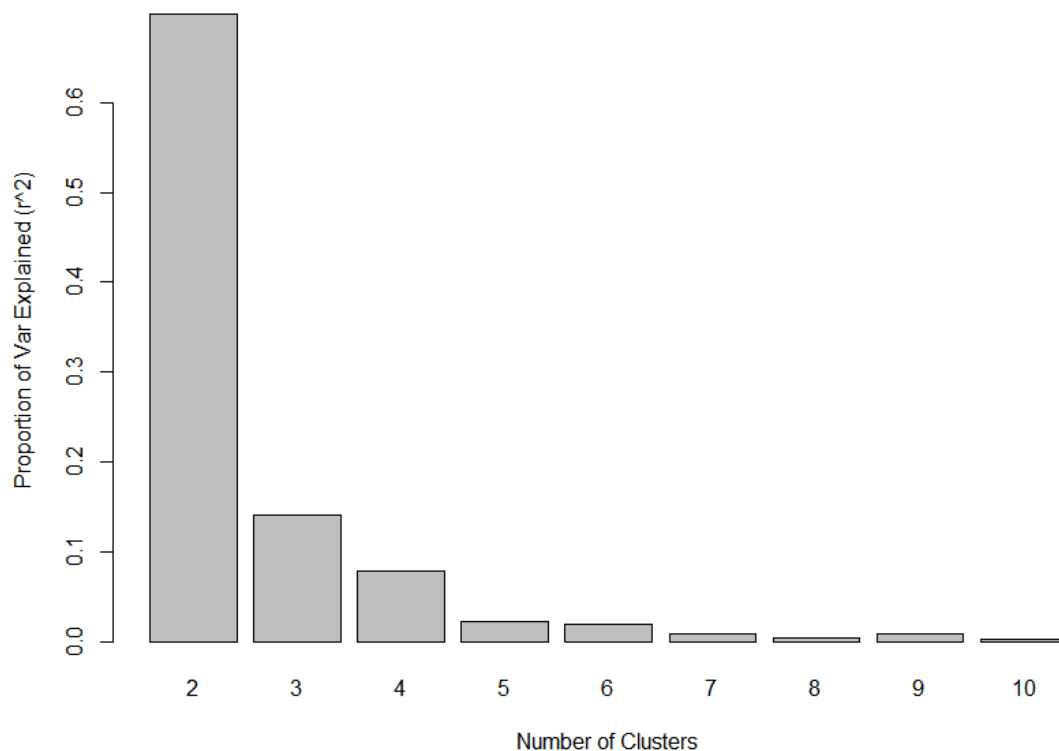
- b. Try a k-nearest neighbors classification. Select a value of k that you think is reasonable, then split the data into a training set and a testing set, then look at how well it classified the test data.

Solution:

From the data set we know that there are likely only two groups represented in the

data. We can verify this by simulating kmeans clusters and computing the rsquared between them. Doing so we found that $k = 2$ is explaining a majority of the variance. Using a 80/20 split for training testing data we found that the kmeans clustering performed rather poorly, misclassifying 86 of the 200 testing data achieving an accuracy of 57%.

Figure 1: Proportion of rSquared for each K.



Code:

```
## Simulating Clusters and Computing rSquared.
rsquared <- rep(NA, 10)
for (i in 1:10){
  tmp <- kmeans(data[, -1], centers = i)
  rsquared[i] <- tmp$betweenss/tmp$totss
}
plot(rsquared, type="l", lwd=2, xlab = '# of Clusters')
barplot((rsquared[2:10] - rsquared[1:9]), names.arg = seq(2,10),
        ylab = 'Proportion of Var Explained (r^2)',
        xlab = 'Number of Clusters')
```

```
## Split the data and cluster using kmeans
split <- sample(c(rep(0, 0.8 * nrow(data)), rep(1, 0.2 * nrow(data))))
test_data <- data[split == 1,]
train_data <- data[split == 0,]

kmeansClustering <- knn(train = train_data[, -1], test = test_data[, -1], cl = train_c
table(test_data[, 1], kmeansClustering)

## kmeansClustering
##    0    1
## 0 24  41
## 1 45  90
```

- c. Now repeat (b), but this time perform a cross validation. Which of these confusion matrices do you think gives you a better idea as to how well KNN will classify new observations.

Solution:

Applying one-out cross validation we get that the KNN misclassified 509 of the 1000 observations achieving an accuracy of 49%. Generally I will tend to trust the more skeptical model validation method so in this case I think the $\sim 50\%$ accuracy is true to what the KNN model is doing. Ideally we would be training the model on all the data, and gather more data to validate the model, rather than doing a split or cross validation.

Code:

```
## Performing one-out cross validation
CV_classifications <- rep(NA, 1000)
for(i in 1:1000){
  tmp_rule <- D2.disc(data=data[-i, 2:8],
                     grouping=data[-i, 1])
  CV_classifications[i] <- predict(tmp_rule,
                                 newdata=data[i, 2:8])$class
}
## Generating confusion matrix
table(data[, 1], CV_classifications - 1)
#      0    1
# 0 137 163
# 1 328 372
```

- d. Try several values of k and look at the test/trina confusion matrix from (b) for each value of k . What is your preferred value of k .

Solution: