

Homework Week 2

2022-08-29

Question 3.1

Using the same data set (*credit_card_data.txt* or *credit_card_data-headers.txt*) as in Question 2.2, use the *ksvm* or *kknn* function to find a good classifier:

- using cross-validation (do this for the *k*-nearest-neighbors model; *SVM* is optional); and
- splitting the data into training, validation, and test data sets (pick either *KNN* or *SVM*; the other is optional).

KNN Cross Validation

First, lets grab the *kknn* package and bring in our data.

```
library(kknn)

loans <- read.csv("credit_card_data-headers.txt", sep = '\t')

str(loans)

## 'data.frame':   654 obs. of  11 variables:
## $ A1 : int  1 0 0 1 1 1 1 0 1 1 ...
## $ A2 : num  30.8 58.7 24.5 27.8 20.2 ...
## $ A3 : num  0 4.46 0.5 1.54 5.62 ...
## $ A8 : num  1.25 3.04 1.5 3.75 1.71 ...
## $ A9 : int  1 1 1 1 1 1 1 1 1 1 ...
## $ A10: int  0 0 1 0 1 1 1 1 1 1 ...
## $ A11: int  1 6 0 5 0 0 0 0 0 0 ...
## $ A12: int  1 1 1 0 1 0 0 1 1 0 ...
## $ A14: int  202 43 280 100 120 360 164 80 180 52 ...
## $ A15: int  0 560 824 3 0 0 31285 1349 314 1442 ...
## $ R1 : int  1 1 1 1 1 1 1 1 1 1 ...
```

Now we need to separate our data. Since we are using cross validation, we will want 3 sets of data; training, validation, and test. HOWEVER, the *cv.kknn* package takes in only one argument, training data. Why is that? Well it handles the splitting out of training/validation data for us given *kcv* number of folds.

What about testing data? Well, the *cv.kknn* model does not actually create a model, rather, it returns the fitted (*yhat*) values of a given value of *k*. Remember, the whole point of cross validation is not to build a model for predicting, it is to determine the best hyper-parameters so that we can *then* go build a model with those hyper parameters. Thus, we will:

- Split the data in two sets, training and test
 - Let *cv.kknn* handle the training/validating
- after finding the best hyper parameters, we will then use test data to report on the accuracy.

```
sample_size <- floor(.8 * nrow(loans))
set.seed(123)
train_ind <- sample(seq_len(nrow(loans)), size = sample_size)

k_train <-loans[train_ind,]
k_test <- loans[-train_ind,]
```

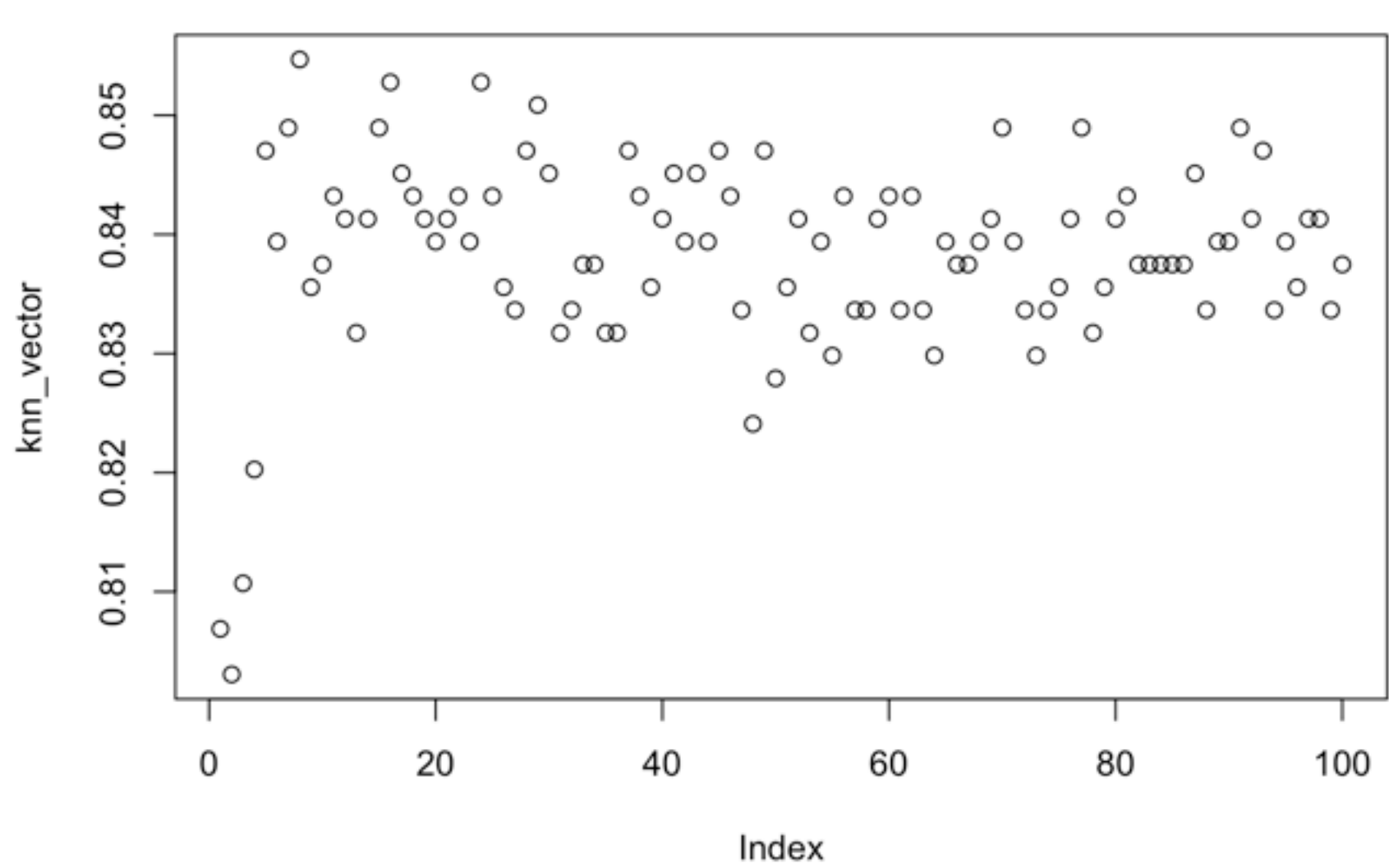
Now, lets loop over values of *K* to find the one that performs best.

```
k_list <- as.list(c(1:100))
k_accuracy <- c()
pred <- rep(0, nrow(k_train)) # create a vector of numbers for as many rows in the data set
knn_vector <- vector("numeric") #create an empty vector to store the accuracy of the model
set.seed(123)
for (K in seq_along(k_list)){ #loop over k to find the accuracy

  k_model <- cv.kknn(R1~.,
                    k_train,
                    kcv = 10,
                    k = K,
                    scale = TRUE)

  pred <- as.data.frame(round(k_model[[1]][,2]))
  actual <-as.data.frame(k_model[[1]][,1])
  k_accuracy <- sum(pred == actual) / nrow(pred) #calc how many predictions were correct

  knn_vector <- c(knn_vector, k_accuracy) #store the models accuracy for each value of k
}
plot(knn_vector)
```



Now find the most accurate value of *K*

```
c(which.max(knn_vector)),max(knn_vector))

## [1] 8.0000000 0.8546845
```

As mentioned above, *cv.kknn* doesn't actually return a model we can reuse to predict new points (it returns a list *y* (actual feature value) and *yhat* (predicted feature value) so that we can calculate the accuracy of each different values of *k*. Therefore, we have to now create a model object using the best value of *K* as determined above using the WHOLE training set.

```
best_k_model <- kknn(k_train$R1~.,
                    k_train,
                    k_test,
                    k = 8,
                    scale = TRUE)
```

Now, we use the test set to measure the accuracy of our model!

```
best_pred <- round(fitted(best_k_model))
sum(best_pred == k_test$R1)/nrow(k_test)

## [1] 0.8244275
```

KNN Cross Validation - Leave one out

Lets also use *train.kknn* to do the leave one out method of cross validation. This is essentially the same thing, just using a value of *kcv* (number of folds) equal to *n* (number of rows in the training data).

This function works a little different than the above, instead of giving us the *yhat* values and having us calculate the accuracy ourself, loop through various values of *K* to determine the best, *train.kknn* actually does the lifting for us - we tell it the maximum values of *K* we want it to loop through for us (or specific values of *K* we want it to loop over) and it just return the best value of *k* it found (the most accurate). Then we will need to build a new model using that parameter if we want to test the model.

```
k_accuracy <- c()
#pred <- rep(0, nrow(k_train)) # create a vector of numbers for as many rows in the data set

#set.seed(3)

k_loo_model <- train.kknn(R1~.,
                        k_train,
                        kmax = 100,
                        kernel = "optimal",
                        scale = TRUE) #see "further analysis"

#actual <- as.data.frame(k_train$R1)
#pred <- round((predict(k_loo_model, k_train)))
#k_accuracy <- sum(pred == actual) / nrow(k_train)

k_loo_model

##
## Call:
## train.kknn(formula = R1 ~ ., data = k_train, kmax = 100, kernel = "optimal",      scale = TRUE)
##
## Type of response variable: continuous
## Minimal mean absolute error: 0.1873805
## Minimal mean squared error: 0.1035524
## Best kernel: optimal
## Best k: 48
```

```
#k_accuracy
```

The best value appears to be 48, so lets build a model using that value and run our test data through it.

```
best_k_model <- kknn(R1~.,
                    k_train,
                    k_test,
                    k = k_loo_model$best.parameters[[2]], #returns the best k from above
                    scale = TRUE)
best_pred <- round(fitted(best_k_model))
sum(best_pred == k_test$R1)/nrow(k_test)

## [1] 0.8320611
```

As it appears, our model is 83% accurate!

Question 3.1b

Now, instead of doing cross validation, lets do "normal validation".

As such, we will randomly split our data into three groups. See documentation therein for process:

```
set.seed(123)

train_size <- floor(.5 * nrow(loans)) #first, calculate how many rows we want. Here I am saying 50% of the rows f
or training
train_ind <- sample(seq_len(nrow(loans)), size = train_size) # now, get indecies for 505 of the data at random

test_val <- loans[-train_ind,] #remove the training data so we can split the remaining data into validation/test
data

val_size <- floor(.5 * nrow(test_val)) # we want 50% of the REMAINING data to be for validation, so calculate wha
t that number is

val_ind <- sample(seq_len(nrow(test_val)), size = val_size) # now get the indecies

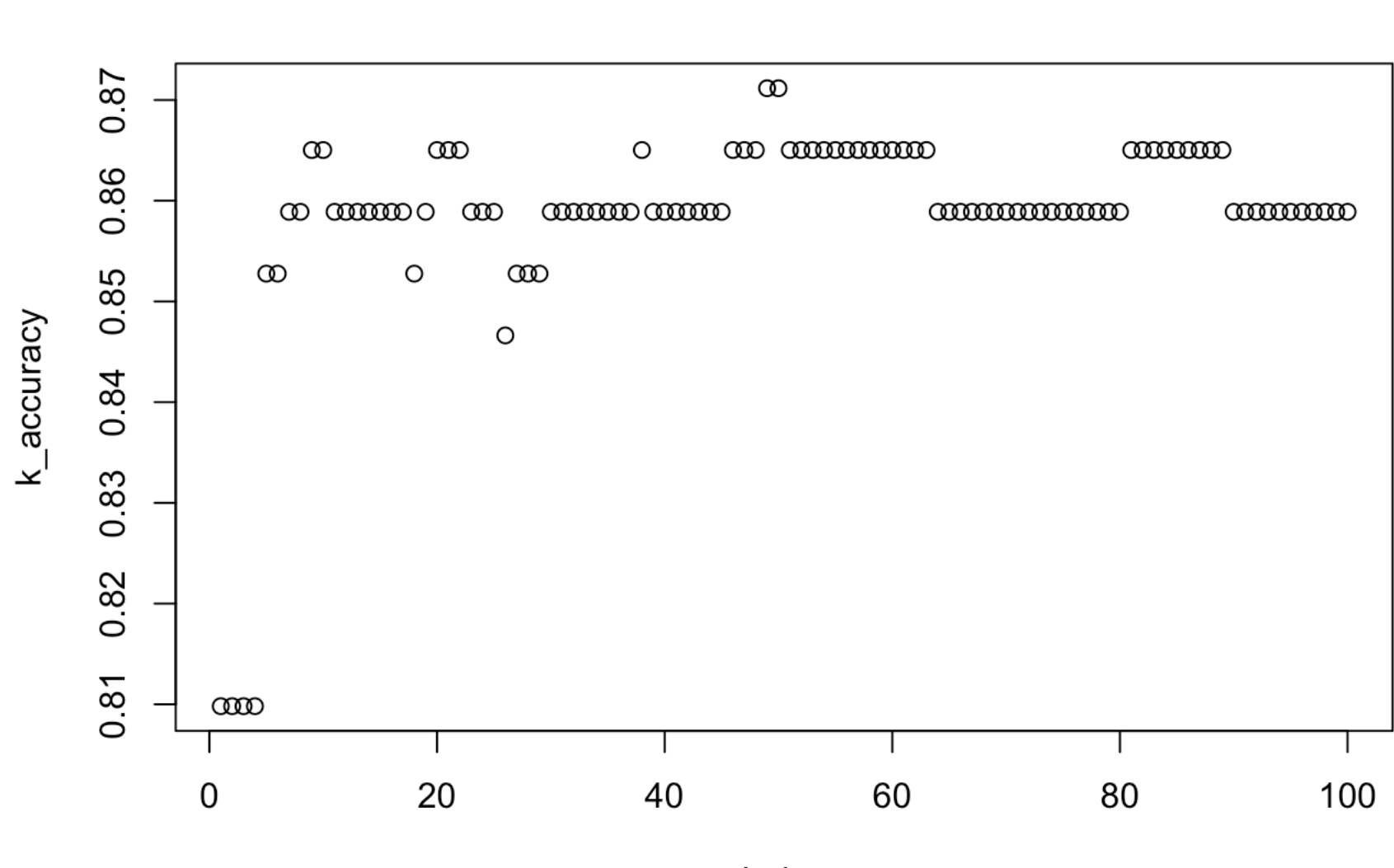
#now store those indecies into variables
k_train <-loans[train_ind,]
k_val <- test_val[val_ind,]
k_test <- test_val[-val_ind,]
```

Now, lets loop over 100 values of *K* and plot the accuracy.

```
k_list <- as.list(c(1:100))
k_accuracy <- c()
pred <- rep(0, nrow(k_train))

for (K in k_list) {
  k_model <- kknn(k_train$R1~.,
                k_train,
                k_val,
                k = K,
                scale = TRUE)
  pred <- round(fitted(k_model))
  k_accuracy[K] <- sum(pred == k_val$R1)/nrow(k_val)
}

plot(k_accuracy)
```



```
c(which.max(k_accuracy),max(k_accuracy))

## [1] 49.0000000 0.8711656
```

Now we see the best value of *k* is 49 with an accuracy of 87%. Now lets make a model with that info and see how accurate it is on our test data:

```
best_k_model <- kknn(k_train$R1~.,
                    k_train,
                    k_test,
                    k = which.max(k_accuracy),
                    scale = TRUE)
best_pred <- round(fitted(best_k_model))
sum(best_pred == k_test$R1)/nrow(k_test)

## [1] 0.8414634
```

It appears the model is 84% accurate.

Further Thoughts

It would appear we had a higher accuracy with normal validation vs cross validation, but that is not an accurate statement. That only happened by chance, as if we remove the seed values from the above code that may shift. The change an accuracy is really only a result in being trained on different data and used different values of *k*. However, I will note, they were extremely similar, both were within 15 accuracy of each other and *K* only differed by 1. Had our data set been much larger, there may have been even less of a difference.

