

# Lab 3

Bangda Sun, bs2996

October 11, 2016

## Instructions

Before you leave lab today make sure that you upload a knitted HTML file to the canvas page (this should have a .html extension). No need to upload the .Rmd file. Include output for each question in its own individual code chunk and don't print out any vector that has more than 20 elements.

Objectives: KNN Classification and Cross-Validation

## Background

Today we'll be using the *Weekly* dataset from the *ISLR* package. This data is similar to the *Smarket* data from class. The dataset contains 1089 weekly returns from the beginning of 1990 to the end of 2010. Make sure that you have the *ISLR* package installed and loaded by running (without the code commented out) the following:

```
# install.packages("ISLR")
library(ISLR)
```

We'd like to see if we can accurately predict the direction of a week's return based on the returns over the last five weeks. *Today* gives the percentage return for the week considered and *Year* provides the year that the observation was recorded. *Lag1* - *Lag5* give the percentage return for 1 - 5 weeks previous and *Direction* is a factor variable indicating the direction ('UP' or 'DOWN') of the return for the week considered.

## Part 1: Visualizing the relationship between this week's returns and the previous week's returns.

1. Explore the relationship between a week's return and the previous week's return. You should plot more graphs for yourself, but include in the lab write-up a scatterplot of the returns for the weeks considered (*Today*) vs the return from two weeks previous (*Lag2*), and side-by-side boxplots for the lag one week previous (*Lag1*) divided by the direction of this week's return (*Direction*).

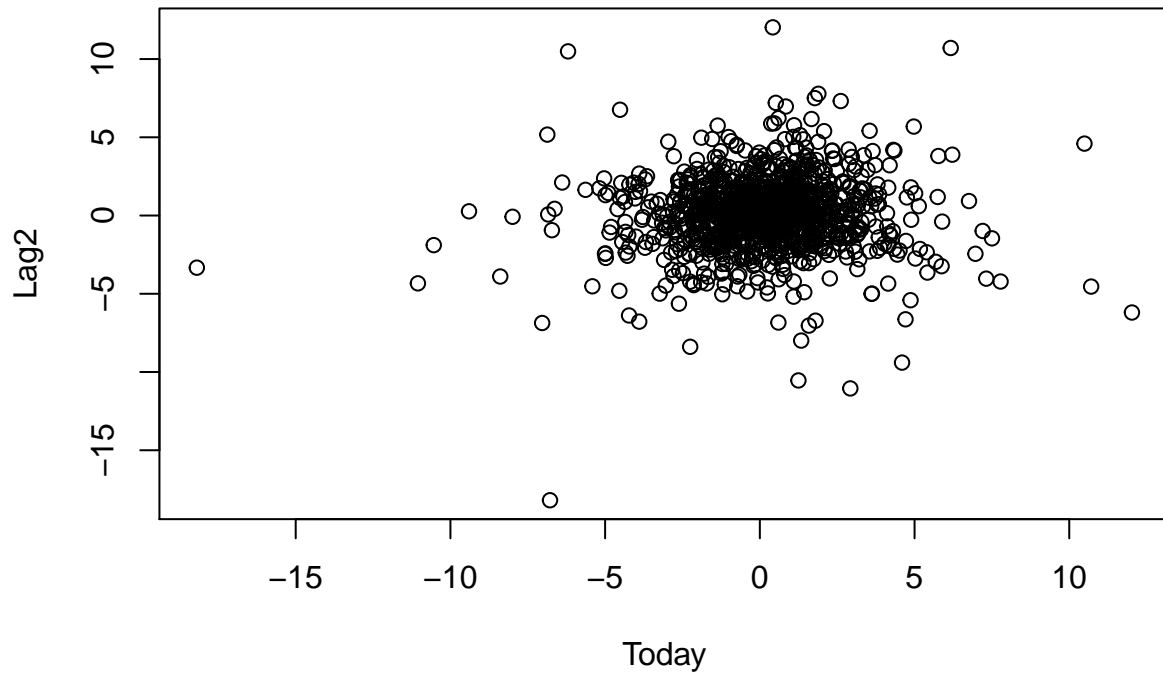
```
# view data
head(Weekly)
```

##	Year	Lag1	Lag2	Lag3	Lag4	Lag5	Volume	Today	Direction
## 1	1990	0.816	1.572	-3.936	-0.229	-3.484	0.1549760	-0.270	Down
## 2	1990	-0.270	0.816	1.572	-3.936	-0.229	0.1485740	-2.576	Down
## 3	1990	-2.576	-0.270	0.816	1.572	-3.936	0.1598375	3.514	Up
## 4	1990	3.514	-2.576	-0.270	0.816	1.572	0.1616300	0.712	Up
## 5	1990	0.712	3.514	-2.576	-0.270	0.816	0.1537280	1.178	Up
## 6	1990	1.178	0.712	3.514	-2.576	-0.270	0.1544440	-1.372	Down

```
# scatter plot
```

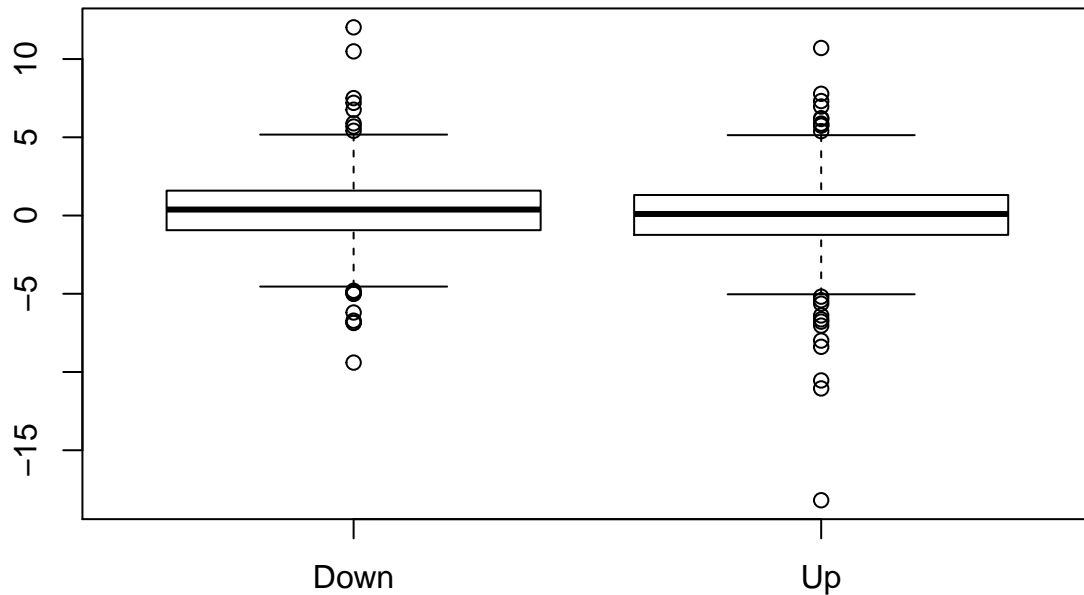
```
plot(Weekly$Today, Weekly$Lag2, xlab = "Today", ylab = "Lag2", main = "Relationship between the consider
```

### Relationship between the considered week's return and the previous two week's return



```
# side by side boxplot  
boxplot(Weekly$Lag1~Weekly$Direction, main = "Box plots of percentage return of one week previous")
```

## Box plots of percentage return of one week previous



## Part 2: Building a classifier

Recall the KNN procedure. We classify a new point with the following steps:

- Calculate the Euclidean distance between the new point and all other points.
- Create the set  $\mathcal{N}_{new}$  containing the  $K$  closest points (or, nearest neighbors) to the new point.
- Determine the number of ‘UPs’ and ‘DOWNS’ in  $\mathcal{N}_{new}$  and classify the new point according to the most frequent.

2. We’d like to perform KNN on the *Weekly* data, as we did with the *Smarket* data in class. In class we wrote the following function which takes as input a new point ( $Lag1_{new}, Lag2_{new}$ ) and provides the KNN decision using as defaults  $K = 5$ ,  $Lag1$  data given in *Smarket*  $Lag1$ , and  $Lag2$  data given in *Smarket*  $Lag2$ . Update the function to calculate the KNN decision for weekly market direction using the *Weekly* dataset with  $Lag1 - Lag5$  as predictors. Your function should have only three input values: (1) a new point which should be a vector of length 5, (2) a value for  $K$  (with default of 5), and (3) the  $Lag$  data which should be a data frame with five columns (and  $n$  rows) and have as default the *Weekly* data. test your function on a point  $c(-.5, .5, -.5, -.5, .5)$ . The result should be ‘UP’.

```
KNNclass <- function(Lag1.new, Lag2.new, K = 5, Lag1 = Smarket$Lag1, Lag2 = Smarket$Lag2) {
  n <- length(Lag1)

  stopifnot(length(Lag2) == n, length(Lag1.new) == 1, length(Lag2.new) == 1, K <= n)

  dists      <- sqrt((Lag1-Lag1.new)^2 + (Lag2-Lag2.new)^2)
```

```

neighbors <- order(dists)[1:K]
neighb.dir <- Smarket$Direction[neighbors]
choice <- names(which.max(table(neighb.dir)))
return(choice)
}

# new KNN.decision function
KNN.Decision <- function(predictor, K = 5, training.set){
  # this function has 3 parameters:
  # predictor: new lag1 - lag5, should be a length 5 vector
  # K: the number of closet points, default 5
  # training.set: existed Lag1 - Lag5, as the training data
  # this function will return the prediction direction (up or down)

  # numbers of observations
  n <- dim(training.set)[1]

  # set the requirement of parameters
  stopifnot(length(predictor) == 5, length(K) == 1, dim(training.set)[2] == 5, K <= n)

  # make sure the mode of predictor is numeric
  predictor <- as.numeric(predictor)

  # compute the euclidean distance
  dist <- sqrt((training.set[,1] - predictor[1])^2 +
               (training.set[,2] - predictor[2])^2 +
               (training.set[,3] - predictor[3])^2 +
               (training.set[,4] - predictor[4])^2 +
               (training.set[,5] - predictor[5])^2)

  # sort the distance and select top K nearest
  neighbor <- order(dist)[1:K]

  # find these K points direction
  neighbor.dir <- Weekly$Direction[neighbor]

  # predict the direction
  prediction.dir <- names(which.max(table(neighbor.dir)))
  return(prediction.dir)
}

# prediction
predictor <- c(-.5, .5, -.5, -.5, .5)
KNN.Decision(predictor, training.set = Weekly[,2:6])

```

```
## [1] "Up"
```

- Now train your model using data from 1990 - 2008 and use the data from 2009-2010 as test data. To do this, divide the data into two data frames, *test* and *train*. Then write a loop that iterates over the test points in the test dataset calculating a prediction for each based on the training data with  $K = 5$ . Save these predictions in a vector. Finally, calculate your test error, which you should store as a variable named *test.error*. The test error calculates the proportion of your predictions which are incorrect (don't match the actual directions).

```

# test data
test <- Weekly[Weekly$Year <= 2010 & Weekly$Year >= 2009, ]
# train data

```

```

train <- Weekly[Weekly$Year <= 2008 & Weekly$Year >= 1990, ]
# number of test data
n <- dim(test)[1]
prediction.dir <- rep(NA, n)
for (i in 1:n){
  prediction.dir[i] <- KNN.Decision(predictor = test[i, 2:6], K = 5, training.set = train[,2:6])
}
# prediction of direction
head(prediction.dir)

## [1] "Down" "Up"   "Up"    "Down" "Up"    "Down"

# test error
test.error <- mean(prediction.dir != test$Direction); test.error

## [1] 0.4519231

```

4. Do the same thing as in question 3, but instead use  $K = 3$ . Which has a lower test error?

```

n <- dim(test)[1]
prediction.dir <- rep(NA, n)
for (i in 1:n){
  prediction.dir[i] <- KNN.Decision(predictor = test[i, 2:6], K = 3, training.set = train[,2:6])
}
# prediction of direction
head(prediction.dir)

## [1] "Down" "Up"   "Up"    "Down" "Up"    "Down"

# test error
test.error <- mean(prediction.dir != test$Direction); test.error

## [1] 0.4423077

```

We can see that  $K = 3$  has a lower test error.

## Part 3: Cross-validation

Ideally we'd like to use our model to predict future returns, but how do we know which value of  $K$  to choose? We could choose the best value of  $K$  by training with data from 1990 - 2008, testing with the 2009 - 2010 data, and selecting the model with the lowest test error as in the previous section. However, in order to build the best model, we'd like to use ALL the data we have to train the model. In this case, we could use all of the *Weekly* data and choose the best model by comparing the training error, but unfortunately this isn't usually a good predictor of the test error.

In this section, we instead consider a class of methods that estimate the test error rate by holding out a (random) subset of the data to use as a test set, which is called  $k$ -fold cross validation. (Note this lower case  $k$  is different than the upper case  $K$  in KNN. They have nothing to do with each other, it just happens that the standard is to use the same letter in both.) This approach involves randomly dividing the set of observations into  $k$  groups, or folds, of equal size. The first fold is treated as a test set, and the model is fit on the remaining  $k - 1$  folds. The error rate,  $ERR_1$ , is then computed on the observations in the held-out fold. This procedure is repeated  $k$  times; each time, a different group of observations is treated as a test set. This process results in  $k$  estimates of the test error:  $ERR_1, ERR_2, \dots, ERR_k$ . The  $k$ -fold CV estimate of

the test error is computed by averaging these values,

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k ERR_k.$$

We'll run a 9-fold cross-validation in the following. Note that we have 1089 rows in the dataset, so each fold will have exactly 121 members.

5. Create a vector *fold* which has *n* elements, where *n* is the number of rows in *Weekly*. We'd like for the *fold* vector to take values in 1-9 which assign each corresponding row of the *Weekly* dataset to a fold. Do this in two steps: (1) create a vector using *rep()* with the values 1-9 *each* repeated 121 times (note  $1089 = 121 \cdot 9$ ), and (2) use *sample()* to randomly reorder the vector you created in (1). After you've created the *fold* vector, run *table(fold)* to make sure your code worked.

```
fold <- rep(c(1:9), 121)
fold <- sample(fold, replace = FALSE)
table(fold)
```

```
## fold
##  1  2  3  4  5  6  7  8  9
## 121 121 121 121 121 121 121 121 121
```

6. Iterate over the 9 folds, treating a different fold as the test set and all others the training set in each iteration. Using a KNN classifier with  $K = 5$  calculate the test error for each fold. Then calculate the cross-validation approximation to the test error which is the average of ERR1, ERR2, ..., ERR9.

```
prediction.dir <- matrix(rep(NA, 1089), nrow = 9, ncol = 121)
test.error      <- rep(NA, 9)
Weekly$group    <- fold
for (i in 1:9){
  test  <- Weekly[Weekly$group == i,]
  train <- Weekly[Weekly$group != i,]
  for (j in 1:121){
    prediction.dir[i, j] <- KNN.Decision(predictor = test[j, 2:6], K = 5, training.set = train[,2:6])
  }
  test.error[i] <- mean(prediction.dir[i,] != test$Direction)
}
cv <- mean(test.error); cv
```

```
## [1] 0.4701561
```

7. Repeat step (6) for  $K = 1$ ,  $K = 3$ , and  $K = 7$ . For which set is the cross-validation approximation to the test error the lowest?

```
ptm <- proc.time()
prediction.dir <- matrix(rep(NA, 1089), nrow = 9, ncol = 121)
test.error1    <- rep(NA, 9)
Weekly$group    <- fold
cv              <- rep(NA, 3)
K               <- c(1,3,7)
for (k in 1:3){
  for (i in 1:9){
    test  <- Weekly[Weekly$group == i,]
    train <- Weekly[Weekly$group != i,]
    for (j in 1:121){
      prediction.dir[i, j] <- KNN.Decision(predictor = test[j, 2:6], K = K[k], training.set = train[,2:6])
    }
  }
}
```

```
    test.error1[i] <- mean(prediction.dir[i,] != test$Direction)
  }
  cv[k] <- mean(test.error1)
}
cv
```

```
## [1] 0.4848485 0.4719927 0.4784206
```

```
proc.time() - ptm
```

```
##      user  system elapsed
##      1.44    0.00    1.43
```

We can see that  $K = 3$  has the lowest test error.