

STATS 769 - Lab 03 - bole001

Bernard O'Leary

18 August 2019

Read in data and Linux commands

Data format is CSV (Comma Separated Values), reading into R using the `read.csv` function. Because there are a large number of files an alternative approach has been taken, to list all files in the Lab02 directory according to the file structure we are looking for ("trips*.csv") and then pick up each file and add to a dataframe (called "trips.df").

The machine that was used to create the report is a Microsoft Windows 10 machine that has an Ubuntu Linux Bash shell built into it as part of the WSL (Windows Subsystem for Linux) feature that is an optional feature of the Windows 10 OS (<https://docs.microsoft.com/en-us/windows/wsl/install-win10>). Once access was gained to the STATS 769 Linux machines, the files were copied to the Windows 10 machine used to create the report using the following scp command:

```
bernardo@PKS10198:/mnt/d/Study/UoA-STATS-769$ scp bole001@sc-cer00014-04.its.auckland.ac.nz:/course/Lab02/
```

Whereas on the Windows 10 machine used to create the report, it was set to the local filesystem where the csv files for the lab were copied to. Code used to create the initial dataframe for the scooter Trips data follows.

```
# Files are already in our local directory
directory = '.'
scooter_files.ls <- intersect(list.files(path=directory, pattern="scooter"), list.files(path=directory,
# First apply read.csv, then rbind
scooter_trips.df = do.call(rbind, lapply(scooter_files.ls, function(x) read.csv(x, stringsAsFactors = F
```

Get rid of non-positive and non-finite values in duration and distance

Remove from the entire dataset.

```
scooter_trips.df <- subset(scooter_trips.df, scooter_trips.df$Trip.Duration > 0)
scooter_trips.df <- subset(scooter_trips.df, scooter_trips.df$Trip.Distance > 0)
scooter_trips.df <- subset(scooter_trips.df, is.finite(scooter_trips.df$Trip.Duration))
scooter_trips.df <- subset(scooter_trips.df, is.finite(scooter_trips.df$Trip.Distance))
```

Add logged duration and distance variables

Name the variables `Trip.Duration.Logged` and `Trip.Distance.Logged` and add to the end of the dataframe. Note that because the data are not cleansed by this point NaNs are produced here.

```
scooter_trips.df$Trip.Duration.Logged <- log(scooter_trips.df$Trip.Duration)
scooter_trips.df$Trip.Distance.Logged <- log(scooter_trips.df$Trip.Distance)
```

Perform 10-fold cross validation

```
# Define MSE function
MSE <- function(m, o) {
  mean((m - o)^2)
}

# Randomly shuffle the data
scooter_trips.df <- scooter_trips.df[sample(nrow(scooter_trips.df)),]

MSE_for_fold.df <- data.frame(fold=integer(),pred_mean=double(),pred_lm=double())

# Create 10 equally size folds
folds <- cut(seq(1,nrow(scooter_trips.df)),breaks=10,labels=FALSE)

# Perform 10 fold cross validation
for(i in 1:10){

  # Segment your data by fold using the which() function
  test_indexes <- which(folds==i,arr.ind=TRUE)
  test_set <- scooter_trips.df[test_indexes, ]
  training_set <- scooter_trips.df[-test_indexes, ]

  # Make the model
  y_train = training_set$Trip.Duration.Logged
  x_train = training_set$Trip.Distance.Logged
  fit_mean = mean(y_train)
  fit_lm <- lm(y ~ x, data.frame(y=y_train, x=x_train))

  # Test results
  y_test <- test_set$Trip.Duration.Logged
  x_test <- test_set$Trip.Distance.Logged

  # Make a vector that is populated with the fit_mean
  pred_mean <- rep(fit_mean, length(y_test))

  # Get a vector of predictions based on test data
  pred_lm <- predict(fit_lm, data.frame(x=x_test))

  # get the RMSE
  MSE_for_fold.df[i,] = c(i,MSE(pred_mean, y_test),MSE(pred_lm, y_test))
}

# Print the result
MSE_for_fold.df
```

```
##      fold pred_mean  pred_lm
## 1      1  0.9611076 0.3791317
## 2      2  0.9165046 0.3727245
## 3      3  0.9493668 0.3861235
## 4      4  0.9471102 0.3803816
## 5      5  0.9110766 0.3877613
## 6      6  0.8984726 0.3875488
```

```
## 7      7 0.9255221 0.3900825
## 8      8 0.8929401 0.3899385
## 9      9 0.9320337 0.3906315
## 10    10 0.8886390 0.3687420
```

```
# Print the average MSE across all folded data
```

```
mean(MSE_for_fold.df$pred_mean)
```

```
## [1] 0.9222773
```

```
mean(MSE_for_fold.df$pred_lm)
```

```
## [1] 0.3833066
```

Visualise the linear model

Plot `scooter_trips.df$Trip.Duration` vs `scooter_trips.df$Trip.Distance` and overlay the average with random error model and our derived logistic regression model. This is to help us visualise the data and the derived model as it applies to the test dataset.

Linear Model vs Average with Random Error (All Data)

