

# Boston Marathon Finish Time Predictions

Harvard Stats E139 Fall 2015

*Nathaniel Burbank*

*Pooja Singh*

*David Wihl*

*December 21, 2015*

## Abstract

We wanted to build a reasonably accurate model for predicting Boston Marathon finish times based on a runner's gender, age, and the number of minutes it took for the runner to complete the first five kilometers of the race (5k split time). We first built a baseline model using untransformed linear regression. Then, we analyzed the data and performed appropriate transformations. After transforming the data, we clustered the race results into several different subgroups using an unsupervised learning algorithms. Using these different clusters of race results, we generated a new regression model for each subset. Finally, we compared the accuracy of predictions made via our clustered regression model with our baseline model.

## Objectives and Motivation

Runners who participate in the Boston Marathon are very concerned about their finish times. The reason is simple. The Boston Athletic Association (BAA), which organizes the race each April, limits the field of participants to runners who are capable of completing a marathon in a very fast time. Starting in the 1970s, when the BAA was concerned that the race was becoming unwieldy as it had grown to include more than 1000 runners for the first time ([source](#)), the BAA established a set of qualifying times that runners had to meet([source](#)) in a previously timed marathon in order to be granted acceptance to the Boston race. With marathoning having grown tremendously in popularity over the following four decades, today race organizers are physically constrained by the 30' width of the road in Hopkinton, MA ([source](#)) where the race begins. While other races like New York and Chicago have increased their participation sizes significantly in recent years, Boston is currently, and will be for the foreseeable future, restricted to a size of approximately 27,000 runners. Thus, while the Boston Marathon course is considered one of the harder marathon courses, the runners who participate are faster than average. The mean marathon finish time for Boston runners is significantly faster than global averages:

Mean Finish Time(min)	Men	Women
Worldwide(*)	253	282
Boston(**)	216	235

(\*) [The Guardian](#)

(\*\*) See Data Sources, below

Boston is also known as a hilly course. Race participants who have qualified for the race in marathons featuring flat courses are sometimes challenged by New England's rolling hills, especially those between miles 16-21, which are among the most difficult of the course. The peak of these hills is affectionately known as Heartbreak Hill.

For all of these reasons plus the long history of the race, the Boston Marathon is a very prestigious event in the running community. Forecasting an accurate finish time in the race is a key motivator for runners to improve their time during the hundreds of miles of training a marathon requires.

## Data Sources

We obtained our data from two sources. First, we downloaded a [previously assembled dataset](#) from a group of researchers at the University of North Carolina-Chapel Hill who [analyzed](#) historic runner performance with the goal of predicting the anticipated finish times for the approximately 5000 runners who were unable to complete the 2013 Boston Marathon due to the bombings at the finish line. This data set contained the split and finish times for 64106 runners in the 2010, 2011, and 2013 Boston Marathons. (The 2012 Boston marathon was intentionally excluded from their dataset because it was unusually hot that year.) We also obtained finish and split times for 78042 runners who completed the Chicago Marathon in 2014 and 2015. We obtained these results directly from the Chicago Marathon's [results website](#) using a custom-built scraper.

## Methodology

We are attempting to build the most accurate prediction model possible on the smallest amount of interesting data. We are not attempting to infer or analyze the data extensively. Our method consists of the following steps:

1. Data cleansing
2. Create a baseline regression using ordinary least squares (OLS) regression
3. Analyze the data and perform appropriate transformations
4. Cluster the data into subgroups using unsupervised learning, such as  $k$ -means
5. Verify the improved accuracy of our predictions.

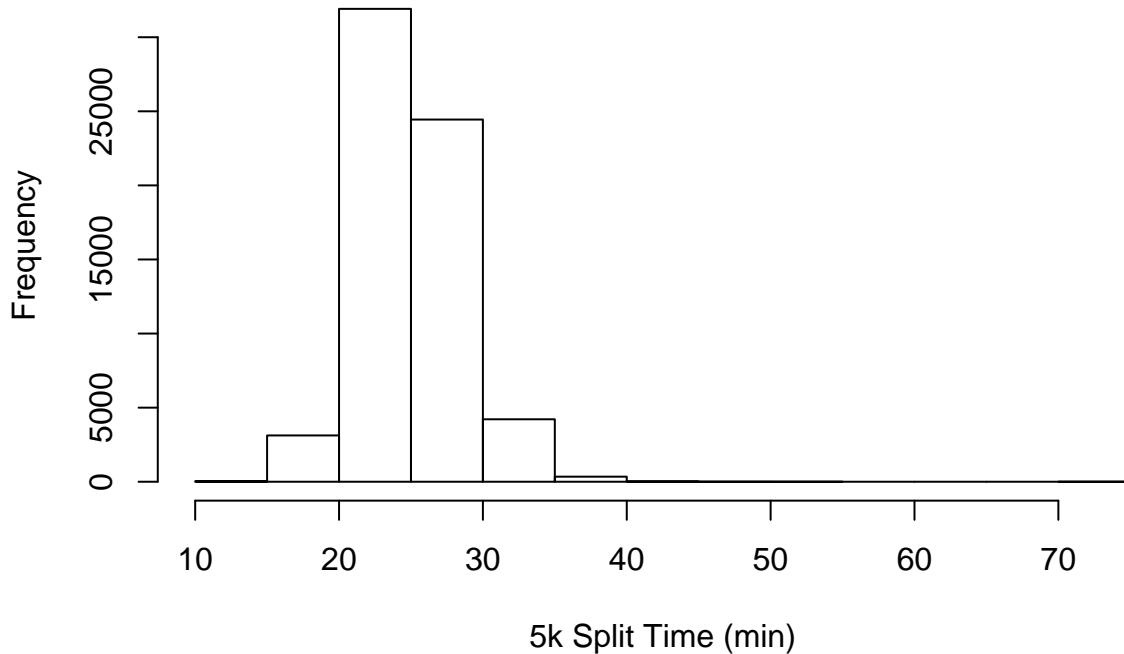
Throughout we will be using 10-fold cross validation testing with sum of squares errors as the key quality metric.

## Data Cleansing

Even though we downloaded our primary data from another study where the researchers had already performed some scrubbing, we needed to perform additional data cleansing to format the data for our needs. As our response variable (finish times) was calculated as a sum of a runners individual split times, we needed to eliminate, rather than impute, any rows that were missing any split times as their inclusion might have tainted the predictive model. Since we were using  $k$ -Fold cross-validation, we randomized the order of the data to ensure that each fold had a reasonable sample size. Finally, gender had a value of {1, 2} which we had to convert into a factor in R.

## Removing Outliers

### Distribution of 5k Split Times



As the histogram reveals, there were several outliers in the 5k split times. These could be due to timing errors, or an injury or accident at the start of the race. Due to the sensitivity of OLS regression, we elected to remove the outliers that were more than three standard deviations ( $\sigma$ ) greater than the mean 5k split time. In total there were 416 rows like this out of 64106 (0.6%).

Age was in the range [18,81]. (The minimum age to run the Boston Marathon is 18). No outliers were removed. Males represented 58.7% of total runners, which seemed reasonable.

## Baseline Regression - 5K Split

We ran a simple OLS regression model using all the Boston data. The predictor variables were age, gender (as a factor), and 5k split time. The response variable was the finish time in minutes.

The baseline model was quite good, but we felt we could do better.

```
##  
## Call:  
## lm(formula = totaltime ~ . - HalfMar, data = dfm)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max  
## -79.265  -9.239  -2.611   6.443 187.552  
##  
## Coefficients:  
##             Estimate Std. Error t value Pr(>|t|)  
## (Intercept) -24.32412    0.54076 -44.981  <2e-16 ***  
## Age         0.06106    0.00619   9.864  <2e-16 ***
```

```

## Gender1F2M2    3.16388    0.13832   22.874   <2e-16 ***
## K0.5          9.82057    0.02154   455.856   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 15.17 on 63686 degrees of freedom
## Multiple R-squared:  0.8036, Adjusted R-squared:  0.8036
## F-statistic: 8.688e+04 on 3 and 63686 DF,  p-value: < 2.2e-16

```

Using  $k$ -Fold cross validation, and taking the sum of squares error for each fold resulted in a baseline mean score of 230 or approximately 15 minutes. While  $\pm$  15 minutes over the course of a four hour race may not seem like a large error,  $\pm$  15 minutes could be the difference between qualifying and not qualifying for a competitive race such as the Boston Marathon.

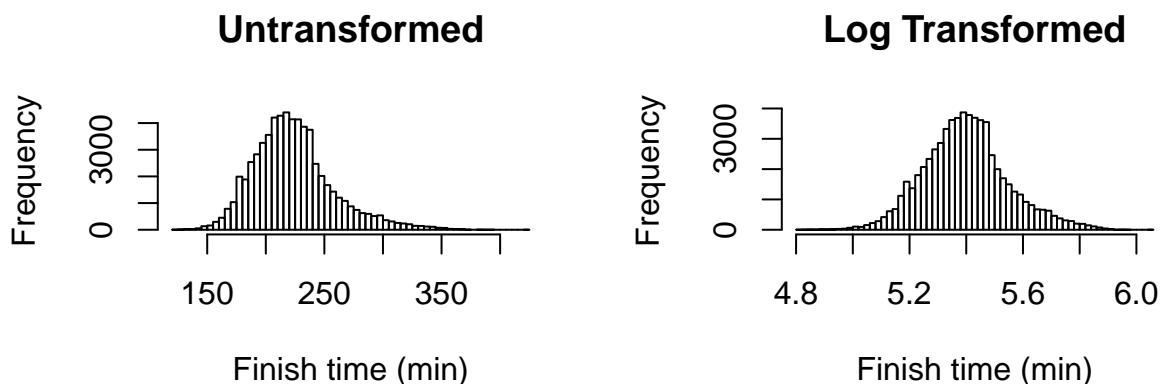
## Baseline Regression - Half Marathon Split

As a comparison, we also ran a regression using half marathon instead of 5k split time. The sum of squares error for each fold resulted in a baseline mean score of 113 or approximately 10 minutes, a significant improvement.  $R^2$  went from 0.8036325 to 0.9032388. Adjusted  $R^2$  improved from 0.8036232 to 0.9032342.

This was an expected result as closer we are to finish line, the better our predictive analytics will be, hence half marathon time was a better predictor than first split time.

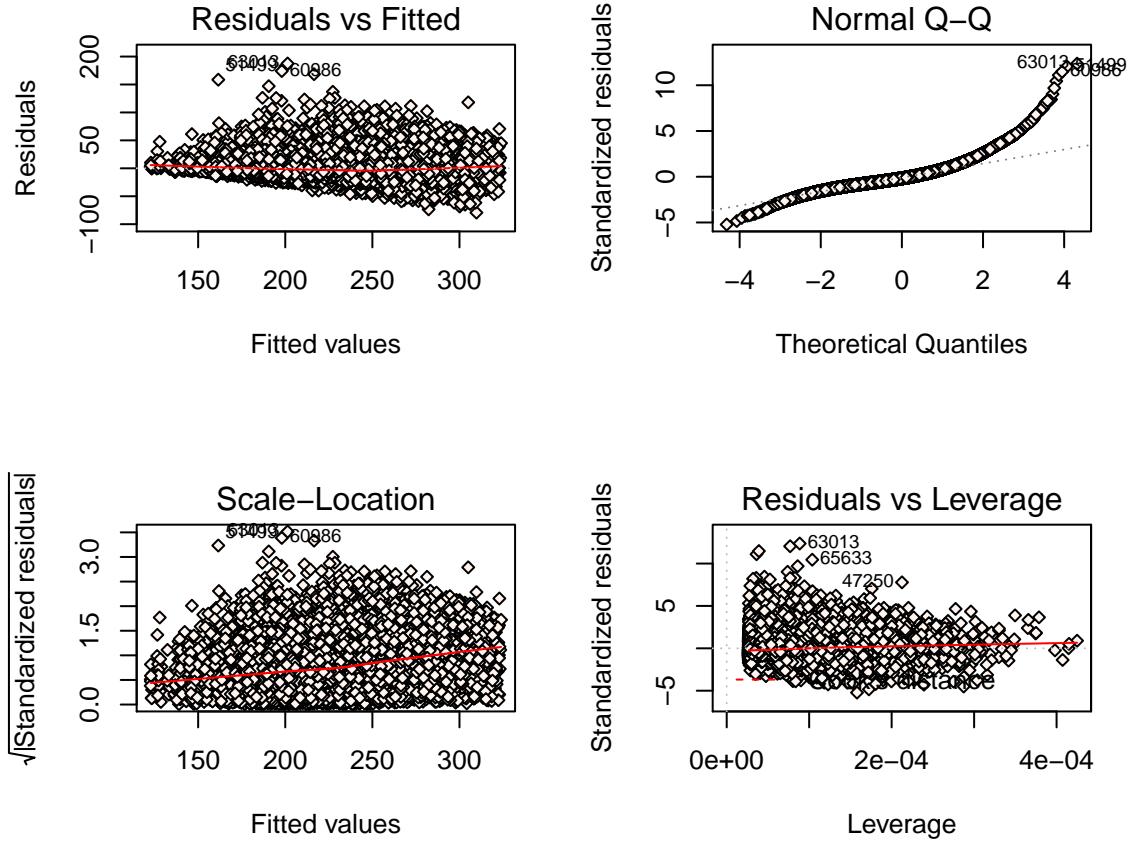
While this may be useful to someone planning to meet the runner at the finish, we felt it would be more challenging to minimize prediction errors using the 5k split time.

## Examination of the Data and Transformations



The untransformed data appeared to be slightly right skewed. A log transformation of the response variable would result in a more normal distribution.

We then examined the base model regression diagnostics such as residuals vs. predicted, Cook's Distance, and heteroskedasticity.



Given the right-tail distribution of the QQplot, the fan out of the residuals vs. predicted, and the slight right-skewness of the finish time histogram, we elected to perform a log transformation on the response variable.

Using  $k$ -Fold cross validation, and taking the sum of squares error for each fold resulted in a log-transformed mean score of 227 or approximately 15 minutes, a negligible improvement.

$R^2$  went from 0.8036325 to 0.8224044, a slight improvement. Adjusted  $R^2$  improved by an equally small amount.

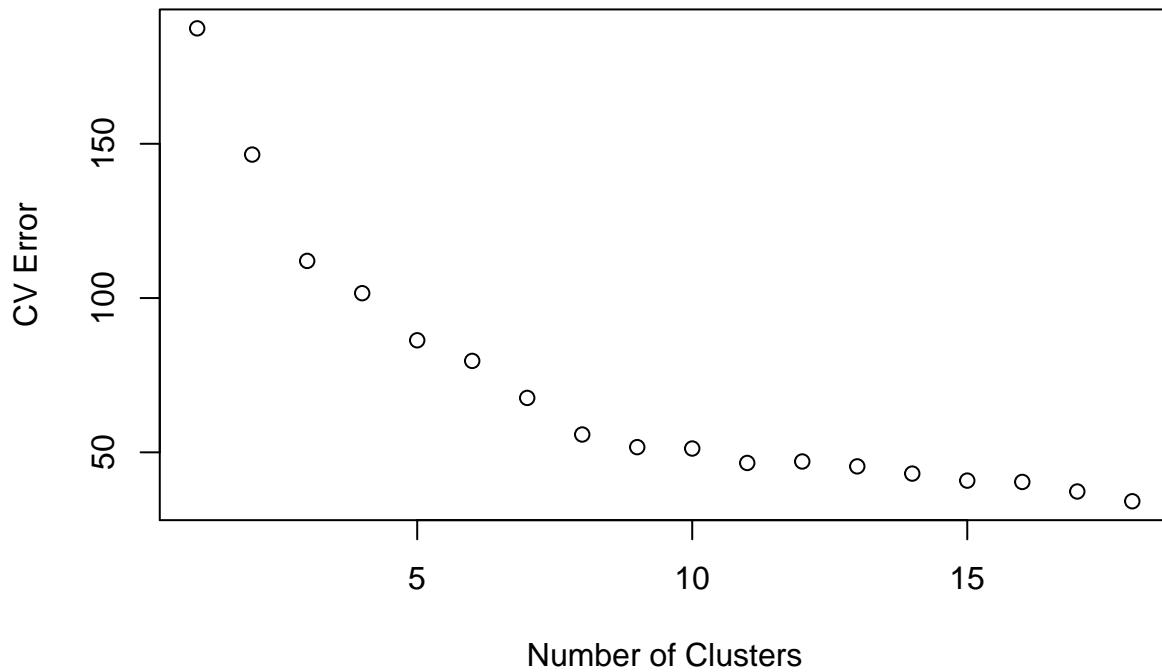
We also attempted a range of polynomial transformations on the 5K times using an exponent in the range [1.1,2.5]. However, this did not improve model's accuracy significantly. We have omitted these results in order to be concise within the report, but they are available on our project GitHub site (See Appendix B).

## Clustering the Data into Subgroups

Since the transformation yielded only a small improvement, we elected to try clustering the data into subgroups using  $k$ -Means. We attempted  $k = [3, 20]$ .

By clustering, the resulting error rates decreased significantly. Using  $k = 3$ , the mean error rate was 187 or  $\pm 13$  min. At  $k = 20$ , the mean error rate was 34 or  $\pm 5$  min, a significant improvement.

## CV Error vs Number of Clusters



Using the [elbow method](#), we determined that the optimal number of clusters is  $k = 8$ , which results in a cross validation error of 79 or  $\pm 8$  min, which is approximately half the error as compared to unclustered linear regression.

```
## Warning: Quick-TRANSfer stage steps exceeded maximum (= 3184500)
```



As you can see, the unsupervised learning found groupings, especially in the dense middle of the field that would have been difficult if not impossible to us to find on our own.

In a given race, there are multiple categories of runners. Elite runners are significantly faster of course. Even though the Boston Marathon requires fast qualification times, an allowance of approximately 10% of the field is made for runners who are running for a charity such as [Dana-Farber](#). The charity runners do not need to qualify and are much slower than an average qualified runner.

The number of runners two  $\sigma$  faster than the mean is 409. There are 2951 runners two  $\sigma$  slower than the mean. (This also confirms our earlier right skewness to the data). The difference in slope for fast vs. slow runners is indicative that a single regression model could be improved upon.

## Applying the Model to Chicago Data

To evaluate the robustness of our model, we attempted to verify the predictability of the model to a different marathon. The Chicago marathon does not require qualification and is a flatter course.

The SSE using clustering and individual regression models is 46. Using the Boston Marathon baseline linear regression is 24. So this model does not generalize well to a marathon with a different profile.

## Conclusion and Next Steps

The following table summarizes our results:

Method	Mean Error (min)
Baseline 5k	15
Baseline Half Marathon	10
Log(y) Transform 5k	15

Method	Mean Error (min)
8-Fold Cluster	8

By subsetting the data, we were able to bring down the cross validation error significantly and improve the accuracy of the model.

This subset model works well for the Boston Marathon, but does not appear to generalize well for different profile marathons.

Potential next steps of this project include:

- Elaborate the data set to build a more robust model that would apply well to multiple marathons.
- Attempt other clustering models besides *k*-Means ([example](#)).
- Attempt other algorithms of regression, such as Random Forests.
- Using other split times, such as half marathon, predict the probability of a runner not finishing the race.
- Making the prediction model available on a public website using R/Shiny. We fill that our prediction model is far more accurate than what is currently available to runners.

## Appendix A - Code

The following is the code used to prepare this report:

```
# Boston Marathon Finish Time Predictions
# Harvard Stats E139 Fall 2015
# Nathaniel Burbank, Pooja Singh, David Wihl
#
# December 21, 2015

library(ggplot2)

# Common Functions

percent <- function(x, digits = 1, format = "f", ...){
  paste0(formatC(100 * x, format = format, digits = digits, ...), "%")
}

create_folds = function (df, k){
  # Perform k-Fold Cross Validation and return average score
  # Inspired by Scikit-Learn's cross_val_score
  if (missing(k)) {
    k = 10
  }

  train1_start = train1_end = train2_start = train2_end = test_start = test_end = c()
  testsize = nrow(df) / k
  for (i in 1:k){
    if (i < k){
      train1_start[i] = 1
      train1_end[i] = floor((k-i) * testsize)
    } else {
      train1_start[i] = 0
      train1_end[i] = 0
    }
    test_start[i] = max(ceiling((k-i) * testsize),1)
    test_end[i] = floor((k-i+1) * testsize)
    train2_start[i] = 0
    train2_end[i] = 0
    if (i > 1){
      train2_start[i] = ceiling((k-i+1) * testsize)
      train2_end[i] = nrow(df)
    }
  }
  return (data.frame(train1_start, train1_end, train2_start, train2_end, test_start, test_end))
}

# we'll do 10-fold cross validation
k = 10

# Read in the data for multiple years
dfm <- read.csv("Previous Boston Marathon study/BAA data.txt", header=T, sep=" ")
```

```

dfm$Age2014 = NULL # remove unneeded column which is mostly NA
ok = complete.cases(dfm)
dfm = dfm[ok,] # remove rows that have any NA values
times = as.matrix(dfm[,7:15], ncol=9)
dfm$totaltime = rowSums(times)
dfm<- dfm[c("totaltime","Age","Gender1F2M","K0.5","HalfMar")] # keep only columns we need
dfm$Gender1F2M = as.factor(dfm$Gender1F2M) # make gender into a factor
dfm = dfm[!is.na(dfm$totaltime), ] # eliminate rows with no finish times
dfm = dfm[sample(nrow(dfm)),] # in case the data is sorted, randomize the order
# Find mean finish time by gender
agg = aggregate(dfm$totaltime, by=list(dfm$Gender1F2M), FUN=mean)[2]
men = as.integer(agg$x[2])
women = as.integer(agg$x[1])
# Remove outliers (5k > 3 sigma from mean)
hist (dfm$K0.5, breaks=15, main="Distribution of 5k Split Times", xlab="5k Split Time (min)")
mean5k = mean(dfm$K0.5)
sd5k = sd(dfm$K0.5)
outliers5k = dfm$K0.5>(mean5k + 3*sd5k)
outliers5ks = sum(outliers5k)
dfm.rows = nrow(dfm)
dfm = dfm[!outliers5k,]
malepct = table(dfm$Gender1F2M)[2] / (table(dfm$Gender1F2M)[2] + table(dfm$Gender1F2M)[1])
# Create a set of k-fold sets for cross-validation
folds = create_folds(dfm,k)
score = c()

# Baseline regression
base.mod = lm(totaltime~.-HalfMar,data=dfm)

# get baseline score
for (i in 1:k) {
  train = dfm[c(folds[i,"train1_start"] : folds[i,"train1_end"], folds[i,"train2_start"] : folds[i,"train2_end"])]
  test = dfm[folds[i,"test_start"] : folds[i,"test_end"],]

  cvmodel = lm(totaltime~.-HalfMar,data=train)
  score[i] = sum((test$totaltime - predict(cvmodel,new=test))^2) / as.numeric(nrow(test))
}
score.mean.5k = mean(score)
summary(base.mod)
# Half Marathon Baseline regression
base.mod.half = lm(totaltime~.-K0.5,data=dfm)

for (i in 1:k) {
  train = dfm[c(folds[i,"train1_start"] : folds[i,"train1_end"], folds[i,"train2_start"] : folds[i,"train2_end"])]
  test = dfm[folds[i,"test_start"] : folds[i,"test_end"],]

  cvmodel = lm(totaltime~.-K0.5,data=train)
  score[i] = sum((test$totaltime - predict(cvmodel,new=test))^2) / as.numeric(nrow(test))
}
score.mean.half = mean(score)
dfm$HalfMar = NULL # Remove half marathon column as it is no longer needed
# Display histograms of response variable untransformed and log transformed
par(mfrow=c(1,2))

```

```

hist(dfm$totaltime, breaks=50, main="Untransformed", xlab="Finish time (min)")
hist(log(dfm$totaltime), breaks=50, main="Log Transformed", xlab="Finish time (min)")
par(mfrow=c(2,2))
plot(base.mod, pch=23 ,bg="seashell",cex=.8)
# Try regression on log(response)
dfm$logtotaltime = log(dfm$totaltime)
folds = create_folds(dfm,k)
score = c()

tx.mod = lm(logtotaltime~.-totaltime,data=dfm)

for (i in 1:k) {
  train = dfm[c(folds[i],"train1_start"):folds[i,"train1_end"],folds[i,"train2_start"]:folds[i,"train2_end"]]
  test = dfm[folds[i,"test_start"]:folds[i,"test_end"],]

  cv.tx.mod = lm(totaltime~.-logtotaltime,data=train)
  score[i] = sum((test$totaltime - predict(cv.tx.mod,new=test))^2) / as.numeric(nrow(test))
}
score.mean.tx = mean(score)

# Clustering the Data into Subgroups
# Let's try subsetting the data set in subgroups using an unsupervised learning algorithm.

minclusters = 3
maxclusters = 20
overallclustererror = c()
for (num_clusters in minclusters:maxclusters) {
  fit.km <- kmeans(dfm, num_clusters)
  meanclustererror = c()
  # Find cv score for each cluster
  for (i in 1:num_clusters){
    df = dfm[ fit.km$cluster == i, ]
    folds = create_folds(df,k)
    foldscore = c()
    for (j in 1:k) {
      train = df[c(folds[j],"train1_start"):folds[j,"train1_end"],folds[j,"train2_start"]:folds[j,"train2_end"]]
      test = df[folds[j,"test_start"]:folds[j,"test_end"],]

      model = lm(totaltime~.-logtotaltime,data=train)
      foldscore[j] = sum((test$totaltime - predict(model,new=test))^2) / as.numeric(nrow(test))
    }
    meanclustererror[i] = mean(foldscore)
  }
  # store mean error for a given k clusters
  overallclustererror[num_clusters-2] = mean(meanclustererror)
}

# TODO fix x axis of plot because it is off by -2
plot(overallclustererror, xlab="Number of Clusters", ylab="CV Error", main="CV Error vs Number of Clusters")
# Show the resulting clusters
fit.km = kmeans(dfm, 8)
dfm$cluster = fit.km$cluster
# Create indicator variables for fast and slow runners to plot different slopes
mean.total = mean(dfm$totaltime)

```

```

sd.total = sd(dfm$totaltime)
dfm$fastrunner = (dfm$totaltime < (mean.total - (2*sd.total)))
dfm$slowrunner = (dfm$totaltime > (mean.total + (2*sd.total)))
slope.model = lm(totaltime~.,data=dfm)
ggplot(dfm, aes(x = K0.5, y = totaltime, color=factor(cluster))) + geom_point(shape=1) +
  labs(list(title = "Clustered Total time Vs. First split time", x = "K0-5", y = "Total Time", colour=""))
  theme(legend.title = element_text(size=6, face="bold") , title = element_text(size=8, face="bold")) +
  geom_abline(intercept = 37, slope = 6,color="red") +
  geom_abline(intercept = 37, slope = 3,color="blue")
# Save the different regression models per cluster
clust.mod = c()
for (i in 1:length(fit.km$size)) {
  df = dfm[which(dfm$cluster==i), ]
  df = df[c("Age","Gender1F2M","K0.5","totaltime")]
  clust.mod[[i]] = lm(totaltime~.,data=df)
}

# Validating model against Chicago Marathon
dfChi14 <- read.csv("ChicagoScraper/Chicago2014Formated.csv",header=T)
dfChi15 <- read.csv("ChicagoScraper/Chicago2015Formated.csv",header=T)
dfChi <- rbind(dfChi14,dfChi15)
### dfChi$Gender1F2M <- as.factor(dfChi$Gender1F2M) # needs to stay numeric for cluster assignment
Chitimes = as.matrix(dfChi[,7:15], ncol=9)
dfChi$totaltime = rowSums(Chitimes)
dfChi$logtotaltime = log(dfChi$totaltime)
dfChi = dfChi[!is.na(dfm$totaltime), ]
#Draw a random sample
dfChi = dfChi[sample(nrow(dfChi), 1000), ][c("Age","Gender1F2M","K0.5","totaltime","HalfMar")]
# Remove outliers
mean5k = mean(dfChi$K0.5)
sd5k = sd(dfChi$K0.5)
outliers5k = dfChi$K0.5>(mean5k + 3*sd5k)
dfChi = dfChi[!outliers5k,]
assign.cluster <- function(df, centers) {
  # compute squared euclidean distance from each sample to each cluster center
  # There is probably a vectorized way of doing this.
  cols = c("Age","Gender1F2M","K0.5")
  clusters = c()
  for (i in 1:nrow(df)){
    diff = c()
    for (j in 1:nrow(centers) ){
      y = centers[j,cols]
      x = as.matrix(df[i,cols])
      diff[j] = sum((x-y)^2)
    }
    clusters[i] = which.min(diff)
  }
  return (clusters)
}

# find which cluster would be most appropriate
dfChi$cluster = assign.cluster(dfChi, fit.km[["centers"]])

```

```

# Run appropriate regression model for that cluster
dfChi$Gender1F2M = as.factor(dfChi$Gender1F2M) # need to convert to factor to match model
newtot = c()
for (n in 1:nrow(dfChi)) {
  newtot[n] = predict(clust.mod[[dfChi[n, "cluster"]]], newdata=dfChi[n,])
}

SSE = as.integer(sqrt(sum((newtot - dfChi$totaltime)^2) / nrow(dfChi)))

y.hat.chi = predict(base.mod,new=dfChi)
SSE2 = as.integer(sqrt(sum((y.hat.chi - dfChi$totaltime)^2) / nrow(dfChi)))
##
```

## Appendix B - References

This entire project code including source data can be found on GitHub (<https://github.com/wihl/statse139-project>)

[K-Means Clustering in R Sample Code](#)

### Solutions Found on StackOverflow

[Place R Code in Appendix](#)

[Takes mean by subgroup in R](#)

[Formatting R Values as Percent](#)

[R Markdown Page Break](#)

[Remove Rows with NA in DataFrame \(use of complete.cases\)](#)

[Assign a new datapoint to existing clusters](#)