# Forecasting Bike-Share Usage via Conditional Inference Trees and Random Forests

Andrew Trice

April 29, 2015

## 1 Introduction

A bike-share is a system in which bike are rented on a on-demand basis. They can be returned when and where they are no longer needed. They are paid for through online kiosks that record arrival location, trip length, trip duration, and departure location. This data is important because it can give researchers an indication of the transportation patterns of a city. These transportation patterns can be extrapolated to the city's vitality, naunces, and improvements that are needed

## 2 Materials

This is a purely research based competition hosted on the competitive Data Science platform, Kaggle. The dataset was provided by Hadi Fanaee Tork using data from Capital One Bikeshare and is hosted by the UCI machine learning library.

Capital One Bikeshare in Washington DC presents bikeshare rental data from January 1st 2011 to December 31st 2012. The season, whether it was a holiday, whether it was a workingday, severity of the weather, temperature, adjusted temperature, humidity, windspeed, casual rentals, and registered rentals are all included as features. The goal of the competition is to predict bicycle rental count for each individual hour.

We are using the R open-source statistical programming language and Rstudio software to conduct analysis. The data is already partitioned into training and testing sets in csv or comma separated value format. We import the "party" and "randomForests" packages from R.

## 3 Methods

### 3.1 Feature Engineering

All data given was numerical and different regression methods could suffice to forecast bike rentals. But we want to apply novel methods of forecasting demand by first using a decision tree to establish as a benchmark and then by creating a random forests ensemble.

First though we must feature engineer our data into more useful variables. The datetime feature contains the date in y-m-d form followed by the time. We substring the time from datetime and then an integer variable hour for the particular hour, we convert hour back to a factor. We use the weekdays and as.date functions on the datetime variable to create a day of the week factor variable. We factor the weather, season, and working day variables. We notice that there are minimal instances of weather = 4 so we reassign those variables to weather = 3 to increase model robustness.

With these newly engineered variables we can view average rental times by hour and day of the week and understand their effect on bike rentals.

## 3.2 Initial Plot Analysis

We notice that Saturday and Sunday have significantly different times than the weekdays. The weekdays have the most pronounced rental periods at 8AM and 5PM where the weekend has a gradual period at 10AM through 5PM. We proceed to build our our initial decision tree.
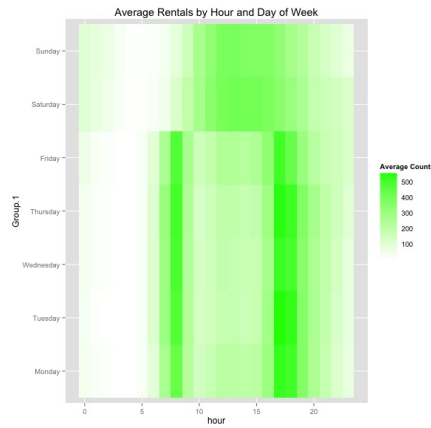


Figure 1: Bike Rental Averages by Day of Week and Hour

## 3.3 Conditional Inference Trees

### 3.3.1 Decision Tree Overview

Decision trees are like flowcharts in which each node represents the class label and each branch represents the decision made. Extrapolated to our case we may see a decision tree node of weather. If the weather is 3 or 4 then we will see a reduction in the number of bikes rented. A tree may have multiple nodes and branches to make a conclusion.

Decision trees are commonly used in operations research to determine optimal strategy but have applications to many fields. They are especially useful for determining variable importance in our forecasting methods.

### 3.3.2 Conditional Inference Tree Overview

We choose to implement a conditional inference tree, or ctree, as a benchmark over the traditional inference tree. The ctree avoids some of the variable selection bias of the traditional decision tree by using a significance test procedure to select variables as opposed to selecting the variables that maximize information gain.

### 3.3.3 Conditional Inference Tree Implementation

Our ctree is built with the variables: season, holiday, workingday, weather, temp, atemp, humidity, and hour. This builds an extensive tree with 393 individual nodes.

## 3.4 Random Forests

### 3.4.1 Random Forests Overview

Random forests construct many decision trees and output the mean prediction of the decision trees it builds. Each decision tree is constructed optimally from a random subset of the features in what is called "feature bagging" or "bagging". This helps correct for overfitting or when a model simply memorizes the training data as opposed to detecting the dependent variable relationships.

$$\hat{f} = \frac{1}{B} \sum_{b=1}^{B} \hat{f}_b(x')$$

Figure 2: Random Forests Unseen Predictions

$$\mathcal{D}_n = \{(X_i, Y_i)\}_{i=1}^{n}$$

Figure 3: Random Forests Variable Importance

### 3.4.2 Random Forests Implementation

We use the same variables as before except we are going to create two formulas: one for the casual rentals and one for the registered rentals. We hypothesize that the random forests can pick up on the nuances of the two rental types. We set the seed to be 192 or the mean rental of the training data, number of trees in the Forests to be 2000 each, the number of variables per level to be 2, which is the default, and for importance = true meaning high variable importance.

# 4  Results

## 4.1  Conditional Inference Tree

The submission of the conditional inference tree received a rank of 950 out of 2219 distinct teams at the time of submission. This makes it better than 57% of teams. The average bike rental value is 189.9 which differs only marginally from the training data which is 191.6. This is an appropriate benchmark and sets a reasonable standard for the random forest implementation.

## 4.2  Random Forests

At the time of this publication the random forests receives a placement of 479 out of 2778. This makes it better than 82% of teams. We can see that our random forests receives a significantly improved score over the conditional inference tree benchmark and we can observe the variables' importance.
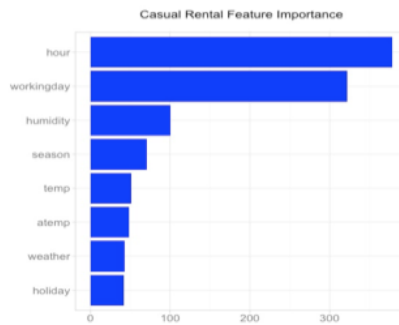
## 4.3  Variable Importance
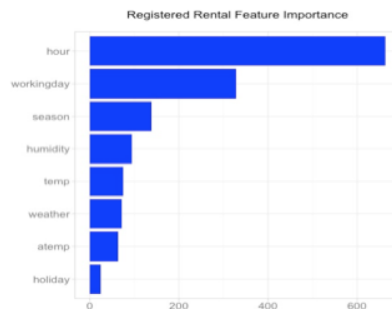


Figure 4: Casual Rental Variable Importance



Figure 5: Registered Rental Variable Importance

Notice that in order of importance the variables for casual are: hour, workingday, humidity, season, temp, atemp, weather, and holiday. For registered the

variables are: hour, workingday, season, humidity, temp, weather, atemp, and holiday.

The registered has a far greater importance on hour than the casual riders which primarily comes at the expense of the workingday feature.

# 5 References

Côme Etienne and Oukhellou Latifa. 2014. Model-Based Count Series Clustering for Bike Sharing System Usage Mining: A Case Study with the Vélib' System of Paris. ACM Trans. Intell. Syst. Technol. 5, 3, Article 39 (July 2014), 21 pages. DOI=10.1145/2560188 http://doi.acm.org/10.1145/2560188

Fanaee-T, Hadi, and Gama, Joao, Event labeling combining ensemble detectors and background knowledge, Progress in Artificial Intelligence (2013): pp. 1-15, Springer Berlin Heidelberg. Rixey, R. Alexander. "Station-Level Forecasting of Bikesharing Ridership."Transportation Research Record: Journal of the Transportation Research Board2387.1 (2013): 46-55.

# 6 Appendix

## 6.1 Acknowledgments

## 6.2 Conditional Inference Tree R Code

```
1  #package party is used to implement conditional inference tree
2  #package ggplot2 is used for visualizations
3  library('party')
4  library('ggplot2')
5
6  #set working directory
7  setwd("/users/andrewtrice/Desktop/Data Science Capstone")
8
9  #read in train/test
10 train <- read.csv("train.csv")
11 test <- read.csv("test.csv")
12
13 str(train)
14
15 #factorize training set
16 train_factor <- train
17 train_factor$weather <- factor(train$weather)
18 train_factor$holiday <- factor(train$holiday)
```

```r
19 train_factor$workingday <- factor(train$workingday)
20 train_factor$season <- factor(train$season)
21
22 #factorize test set
23 test_factor <- test
24 test_factor$weather <- factor(test$weather)
25 test_factor$holiday <- factor(test$holiday)
26 test_factor$workingday <- factor(test$workingday)
27 test_factor$season <- factor(test$season)
28
29 #create time column by stripping out timestamp
30 train_factor$time <- substring(train$datetime,12,20)
31 test_factor$time <- substring(test$datetime,12,20)
32
33 #factorize new timestamp column
34 train_factor$time <- factor(train_factor$time)
35 test_factor$time <- factor(test_factor$time)
36
37 #create day of week column
38 train_factor$day <- weekdays(as.Date(train_factor$datetime))
39 train_factor$day <- as.factor(train_factor$day)
40 test_factor$day <- weekdays(as.Date(test_factor$datetime))
41 test_factor$day <- as.factor(test_factor$day)
42
43 aggregate(train_factor[,"count"],list(train_factor$day),mean)
44
45 #create Sunday variable
46 train_factor$sunday[train_factor$day == "Sunday"] <- "1"
47 train_factor$sunday[train_factor$day != "1"] <- "0"
48
49 test_factor$sunday[test_factor$day == "Sunday"] <- "1"
50 test_factor$sunday[test_factor$day != "1"] <- "0"
51
52 #convert to factor
53 train_factor$sunday <- as.factor(train_factor$sunday)
54 test_factor$sunday <- as.factor(test_factor$sunday)
55
56 #convert time and create $hour as integer to evaluate for daypart
57 train_factor$hour<- as.numeric(substr(train_factor$time,1,2))
58 test_factor$hour<- as.numeric(substr(test_factor$time,1,2))
59
60 #create daypart column, default to 4 to make things easier for
      ourselves
61 train_factor$daypart <- "4"
62 test_factor$daypart <- "4"
63
64 #4AM - 10AM = 1
65 train_factor$daypart[(train_factor$hour < 10) & (train_factor$hour
      > 3)] <- 1
66 test_factor$daypart[(test_factor$hour < 10) & (test_factor$hour >
      3)] <- 1
67
68 #11AM - 3PM = 2
69 train_factor$daypart[(train_factor$hour < 16) & (train_factor$hour
      > 9)] <- 2
70 test_factor$daypart[(test_factor$hour < 16) & (test_factor$hour >
      9)] <- 2
71
72 #4PM - 9PM = 3
73 train_factor$daypart[(train_factor$hour < 22) & (train_factor$hour
      > 15)] <- 3
```

```
74  test_factor$daypart[(test_factor$hour < 22) & (test_factor$hour >
        15)] <- 3
75
76  #convert daypart to factor
77  train_factor$daypart <- as.factor(train_factor$daypart)
78  test_factor$daypart <- as.factor(test_factor$daypart)
79
80  #convert hour back to factor
81  train_factor$hour <- as.factor(train_factor$hour)
82  test_factor$hour <- as.factor(test_factor$hour)
83
84  #create formulas for casual and registered riders
85  formula <- count ~ season + holiday + workingday + weather + temp +
        atemp + humidity + hour + daypart + sunday
86
87  #conditional inference tree
88  fit.ctree <- ctree(formula, data=train_factor)
89
90  #examine model for variable importance
91  fit.ctree
92
93  #run model against test data set
94  predict.ctree <- predict(fit.ctree, test_factor)
95
96  #label importance
97  imp <- importance(fit.ctree, type=1)
98  featureImportance <- data.frame(Feature=row.names(imp), Importance=
        imp[,1])
99
100 #plot variable importance
101 p <- ggplot(featureImportance, aes(x=reorder(Feature, Importance),
        y=Importance)) +
102   geom_bar(stat="identity", fill="#12cfff") +
103   coord_flip() +
104   theme_light(base_size=20) +
105   xlab("Importance") +
106   ylab("") +
107   ggtitle("Random Forest Feature Importance\n") +
108   theme(plot.title=element_text(size=18))
109
110 #build a dataframe with our results
111 submit.ctree <- data.frame(datetime = test$datetime, count=ctree.
        count)
112
113 #write results to .csv for submission
114 write.table(submit.ctree, file="submit_ctree_v2.csv", sep=",",row.
        names=FALSE, col.names=v)
```

## 6.3 Random Forests R Code

```
1  #import necessary packages
2  library('randomForest')
3  library('ggplot2')
4
5  #set working directory
6  setwd("/users/andrewtrice/Desktop/Data Science Capstone")
7
8  #import datasets from working directory
9  train <- read.csv("train.csv") #use nrows=1000 rows for speed
        during feature engineering
10 test <- read.csv("test.csv") #use nrows=1000 rows for speed during
        feature engineering
```

7

```
11
12  str(train)
13
14  #factorize training set
15  train_factor <- train
16  train_factor$weather <- factor(train$weather)
17  train_factor$holiday <- factor(train$holiday)
18  train_factor$workingday <- factor(train$workingday)
19  train_factor$season <- factor(train$season)
20
21  #factorize test set
22  test_factor <- test
23  test_factor$weather <- factor(test$weather)
24  test_factor$holiday <- factor(test$holiday)
25  test_factor$workingday <- factor(test$workingday)
26  test_factor$season <- factor(test$season)
27
28  #create time column by stripping out timestamp
29  train_factor$time <- substring(train$datetime,12,20)
30  test_factor$time <- substring(test$datetime,12,20)
31
32  #factorize new timestamp column
33  train_factor$time <- factor(train_factor$time)
34  test_factor$time <- factor(test_factor$time)
35
36  #create day of week column
37  train_factor$day <- weekdays(as.Date(train_factor$datetime))
38  train_factor$day <- as.factor(train_factor$day)
39  test_factor$day <- weekdays(as.Date(test_factor$datetime))
40  test_factor$day <- as.factor(test_factor$day)
41
42  aggregate(train_factor[,"count"],list(train_factor$day),mean)
43
44  #create Sunday variable
45  train_factor$sunday[train_factor$day == "Sunday"] <- "1"
46  train_factor$sunday[train_factor$day != "1"] <- "0"
47
48  test_factor$sunday[test_factor$day == "Sunday"] <- "1"
49  test_factor$sunday[test_factor$day != "1"] <- "0"
50
51  #convert to factor
52  train_factor$sunday <- as.factor(train_factor$sunday)
53  test_factor$sunday <- as.factor(test_factor$sunday)
54
55  #convert time and create $hour as integer to evaluate for daypart
56  train_factor$hour<- as.numeric(substr(train_factor$time,1,2))
57  test_factor$hour<- as.numeric(substr(test_factor$time,1,2))
58
59  #create daypart column, default to 4 to make things easier for
        ourselves
60  train_factor$daypart <- "4"
61  test_factor$daypart <- "4"
62
63  #4AM - 10AM = 1
64  train_factor$daypart[(train_factor$hour < 10) & (train_factor$hour
      > 3)] <- 1
65  test_factor$daypart[(test_factor$hour < 10) & (test_factor$hour >
      3)] <- 1
66
67  #11AM - 3PM = 2
68  train_factor$daypart[(train_factor$hour < 16) & (train_factor$hour
      > 9)] <- 2
```

```r
69  test_factor$daypart[(test_factor$hour < 16) & (test_factor$hour >
        9)] <- 2

70
71  #4PM - 9PM = 3
72  train_factor$daypart[(train_factor$hour < 22) & (train_factor$hour
        > 15)] <- 3
73  test_factor$daypart[(test_factor$hour < 22) & (test_factor$hour >
        15)] <- 3

74
75  #convert daypart to factor
76  train_factor$daypart <- as.factor(train_factor$daypart)
77  test_factor$daypart <- as.factor(test_factor$daypart)

78
79  #convert hour back to factor
80  train_factor$hour <- as.factor(train_factor$hour)
81  test_factor$hour <- as.factor(test_factor$hour)

82
83  #get rid of weather 4
84  train$weather[train$weather==4] <- 3
85  test$weather[test$weather==4] <- 3

86
87  #variables
88  myNtree = 500
89  myMtry = 7
90  myImportance = TRUE

91
92  #set the seed to mean for model reproducibility
93  set.seed(192)

94
95  #fit and predict casual
96  casualFit <- randomForest(casual ~ season + holiday + workingday +
        weather + temp + atemp + humidity + hour, data=train_factor,
        ntree=myNtree, mtry=myMtry, importance=myImportance)
97  predictCasual <- predict(casualFit, test_factor)

98
99  #fit and predict registered
100 registeredFit <- randomForest(registered ~ season + holiday +
        workingday + weather + temp + atemp + humidity + hour, data=
        train_factor, ntree=myNtree, mtry=myMtry, importance=
        myImportance)
101 predictRegistered <- predict(registeredFit, test_factor)

102
103 #add both columns into final count, round to whole number
104 test$count <- casualFit + test$registered

105
106 #testplot
107 plot(train$count)
108 plot(test$count)

109
110 #write output to csv for submission
111 submit <- data.frame (datetime = test$datetime, count = test$count)
112 write.csv(submit, file = "randomForest_Prediction.csv", row.names=
        FALSE)

113
114 #label variable importances for both casual and registered
115 imp1 <- importance(casualFit, type=1)
116 imp2 <- importance(registeredFit, type=1)
117 featureImportance1 <- data.frame(Feature=row.names(imp1),
        Importance=imp1[,1])
118 featureImportance2 <- data.frame(Feature=row.names(imp2),
        Importance=imp2[,1])
119
```

```r
120 #plot variable importances
121 casualImportance <- ggplot(featureImportance1, aes(x=reorder(
        Feature, Importance), y=Importance)) +
122   geom_bar(stat="identity", fill="blue") +
123   coord_flip() +
124   theme_light(base_size=20) +
125   xlab("") +
126   ylab("") +
127   ggtitle("Casual Rental Feature Importance\n") +
128   theme(plot.title=element_text(size=18))
129
130 registeredImportance <- ggplot(featureImportance2, aes(x=reorder(
        Feature, Importance), y=Importance)) +
131   geom_bar(stat="identity", fill="blue") +
132   coord_flip() +
133   theme_light(base_size=20) +
134   xlab("") +
135   ylab("") +
136   ggtitle("Registered Rental Feature Importance\n") +
137   theme(plot.title=element_text(size=18))
```