

Predicting Train Arrival Status - On Time or Late

Introduction

Adanna Alutu

June 6, 2017

At the beginning of the project, it was hard to come up with a good data to analyze and predict the outcome. Initially I wanted to work on data from my job but after we couldn't see much dependence among the fields that made sense, my mentor Dr. Shmuel Naaman advised me to scout for data from other internet sites he recommended.

Since I take the train most of the time and experienced delay issues many times that has ranged from 10 mins to 2 hours, I became interested in working on transportation data for trains. This is because I want to experience the process of predicting outcomes which is made possible through Data Science. I want to focus on the steps that will make it possible for me and my mentor Dr Shmuel Naaman to predict the arrival times of the train. The possibility of cutting down the delays experienced in waiting for the train no longer seems to be far fetched. My mentor agreed with me and the Septa Train data from Kaggle website was a good option to work on. There were 3 different datasets available to work on but I chose the "on time performance" which I felt has more relevant features, variables and observations and also has sufficient data for the analysis, tests involved.

The variables in the dataset include: 1. train_id 2. status 3. origin 4. direction 5. next_station 6. timeStamp 7. date

subtitle:

Data Exploration

Several steps were taken to ensure elaborate data analysis and wrangling. Every bit of the data was maximized. We went beyond using the provided variables by creating new ones, removing unnecessary data and testing with reliable tools to get quality, reliable results that can be tested with any dataset.

It was necessary to take the following steps to ensure that all the combinations, dicing and testing would yield a meaningful interpretation and prediction that will help tell us with high confidence when the train will be late:

I. We first tried to plot charts with the entire data but the plots were too crowded and blurry to make any sense. The scales were distorted with big units affected by the outliers.

II. GGPlot bar charts were used to plot and observe the trends and statistics summary but the dataset was too huge for the charts.

III. My mentor suggested shuffling the data and taking the first 20percent as sample to work on. Using the formula below, the row-wise shuffling was done first before the column was then shuffled:

IV. We used the data to fit in several models which include:

- GGPlot with different combination of the variables.
- Linear regression model which was used different ways to get the best statistical summary. Including using some of the observations as variables.
- CART model with focus on the classification method because most of the variables in the data are categorical and the prediction is binary with 0 as "on Time" and 1 as "Late"
- Random Forest which created it's own model that highlighted the top more meaningful variables that contributed majorly in predicting the outcome.

Each of these models were implemented because the train dataset contains a mixture of numerical and categorical variables. Converting their types to either numeric or factors wasn't sufficient. To get the benefit of all the variables, it was essential to test these models.

subtitle:

Data Wrangling

Some data manipulations were done which include: + splitting some of the original variables into separate variables. For example, time stamp variable was split into six variables: year, month, day, hour, min, seconds. + Irrelevant variables were removed or set to null so they would not appear in the dataframe used for the predictions. + Some of observations from the wkday and day of month variables were converted to variables and they significantly improved the statistics of the models. The additions however increased the number of variables from 11 to 58. + Units attached to the dependent variable observations were removed to enable conversions to different types and allow plotting with only the observations of the same type. + The dependent variable "status" observations of "on Time" were replaced with "0" using gsub so that all the observations for the variable will match and easier to manipulate. "On time" meant the train arrived as scheduled so it made sense to use "0" to represent no delay.

subtitle:

A Peek into some new variables

This section shows the summary of the SEPTA train data and the first few records using the head().

```
## Warning: Too many values at 150009 locations: 1, 2, 3, 4, 5, 6, 7, 8, 9,
## 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, ...
```

```
##      X.1          X          next_station  direction
## Min.   :      1  Min.   :      6  None           : 7019  N:74989
## 1st Qu.: 47150  1st Qu.: 472113  30th Street Station: 6029  S:75020
## Median : 94032  Median : 942842  Suburban Station : 5883
## Mean   : 94116  Mean   : 941520  Jefferson Station : 5809
## 3rd Qu.:141177  3rd Qu.:1411702  Temple U          : 5596
## Max.   :188210  Max.   :1882001  University City   : 2865
##                                     (Other)           :116808
##      status          origin      train_id
## Length:150009      Doylestown :10961 5368 : 692
## Class :character   Airport Terminal E-F: 9573 3542 : 646
## Mode  :character   Frazer Yard      : 9495 3524 : 622
##                                     Elwyn      : 7439 222 : 547
##                                     None        : 6994 216 : 541
##                                     Roberts Yard : 6991 5315 : 520
##                                     (Other)      :98556 (Other):146441
##      monthday      hour      minute
## Length:150009      Length:150009      Length:150009
## Class :character   Class :character   Class :character
## Mode  :character   Mode  :character   Mode  :character
##
##
##
##      wkday      month
## Length:150009      Min.   : 3.000
## Class :character   1st Qu.: 5.000
## Mode  :character   Median : 7.000
##                                     Mean   : 6.843
```

```
##          3rd Qu.: 9.000
##          Max.    :11.000
##
```

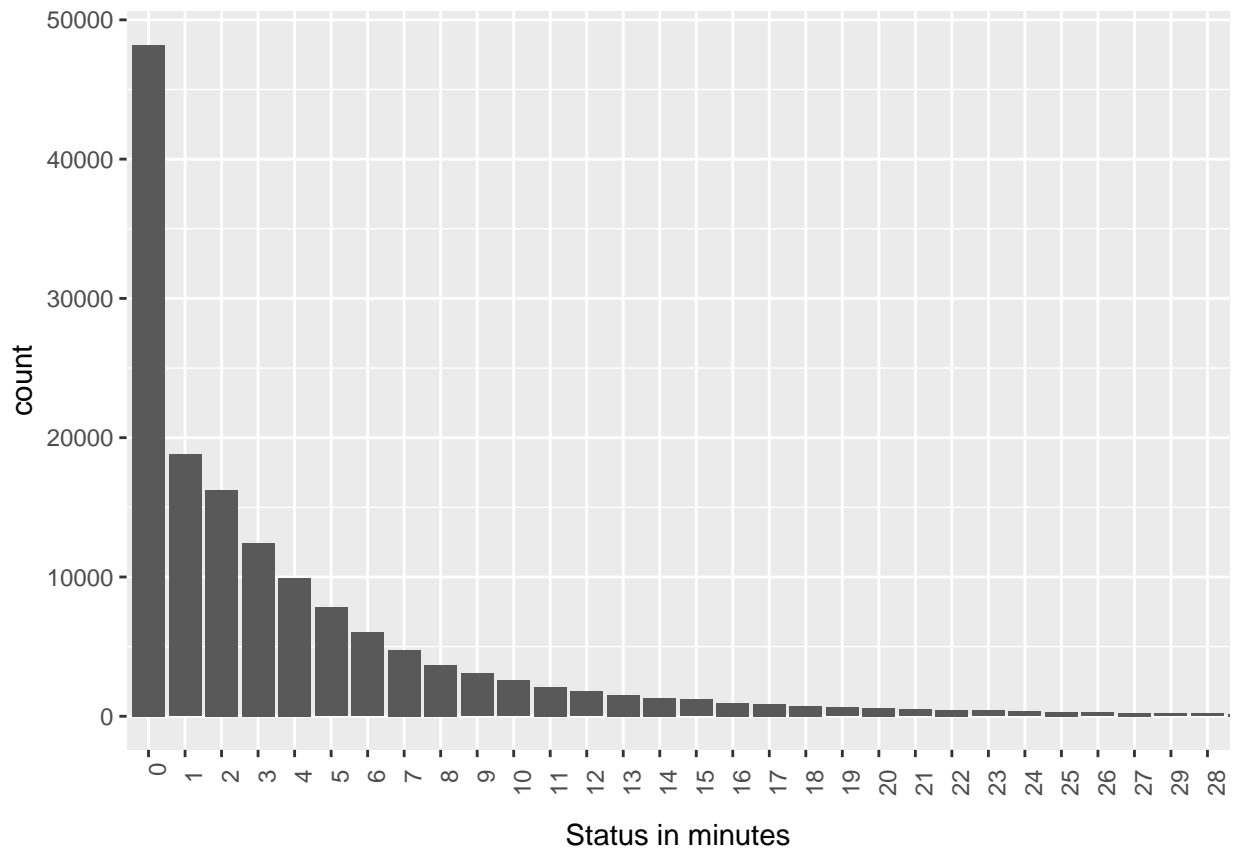
title:

Some Initial plots

GGplot graphs used initially to see trends and relationships within the datasets.

subtitle:

Status variable chart

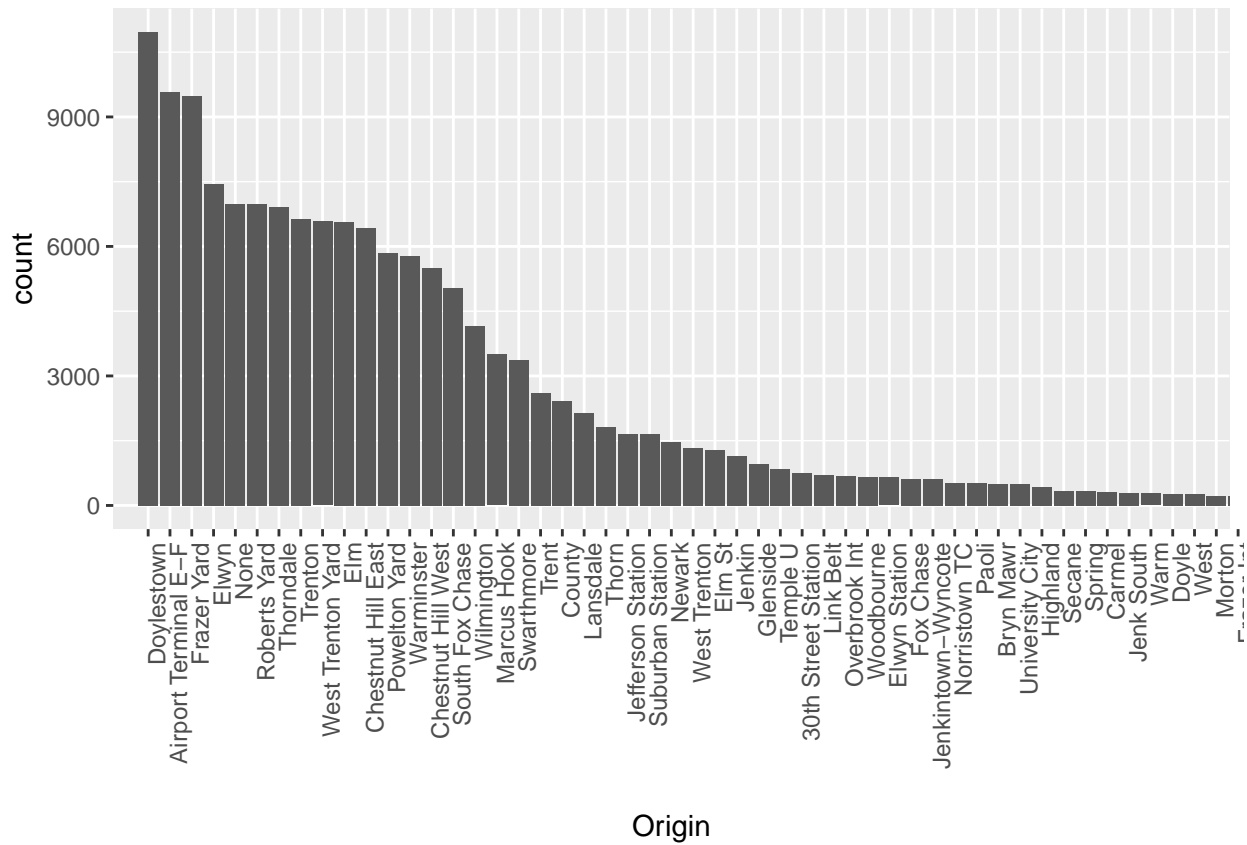


Status is the name of the dependent variable being predicted in this project. The bar chart shows the frequency of the delays experienced by passengers at the train station when the train is late.

From the chart, we can tell that the trains are on time ~50% of the time and late 50% of the time. In this project, we want to predict when to expect the train to be late and when it will be early to avoid waste of time when possible.

subtitle:

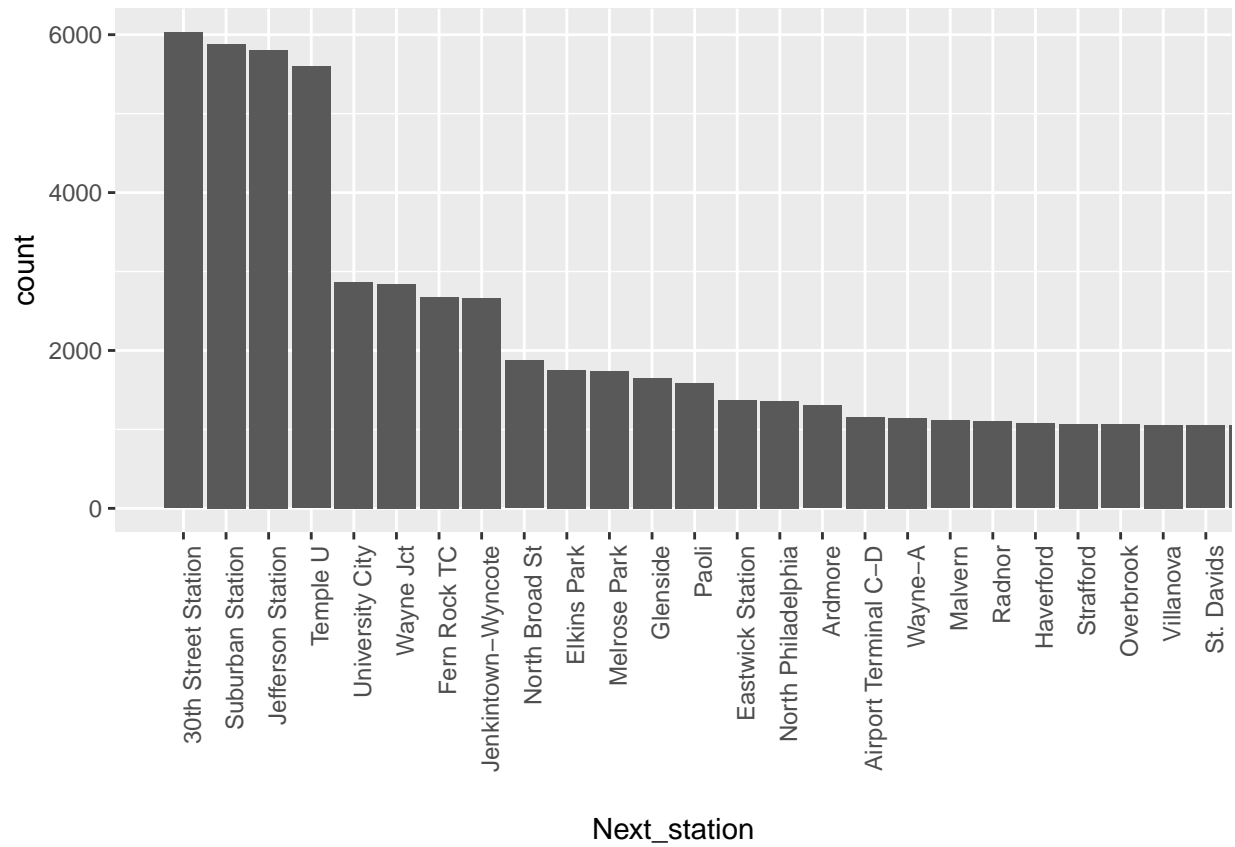
Origin variable bar chart



This chart shows the origin of each passenger's trip. This also represents the train station where the passenger's journey begins.

title:

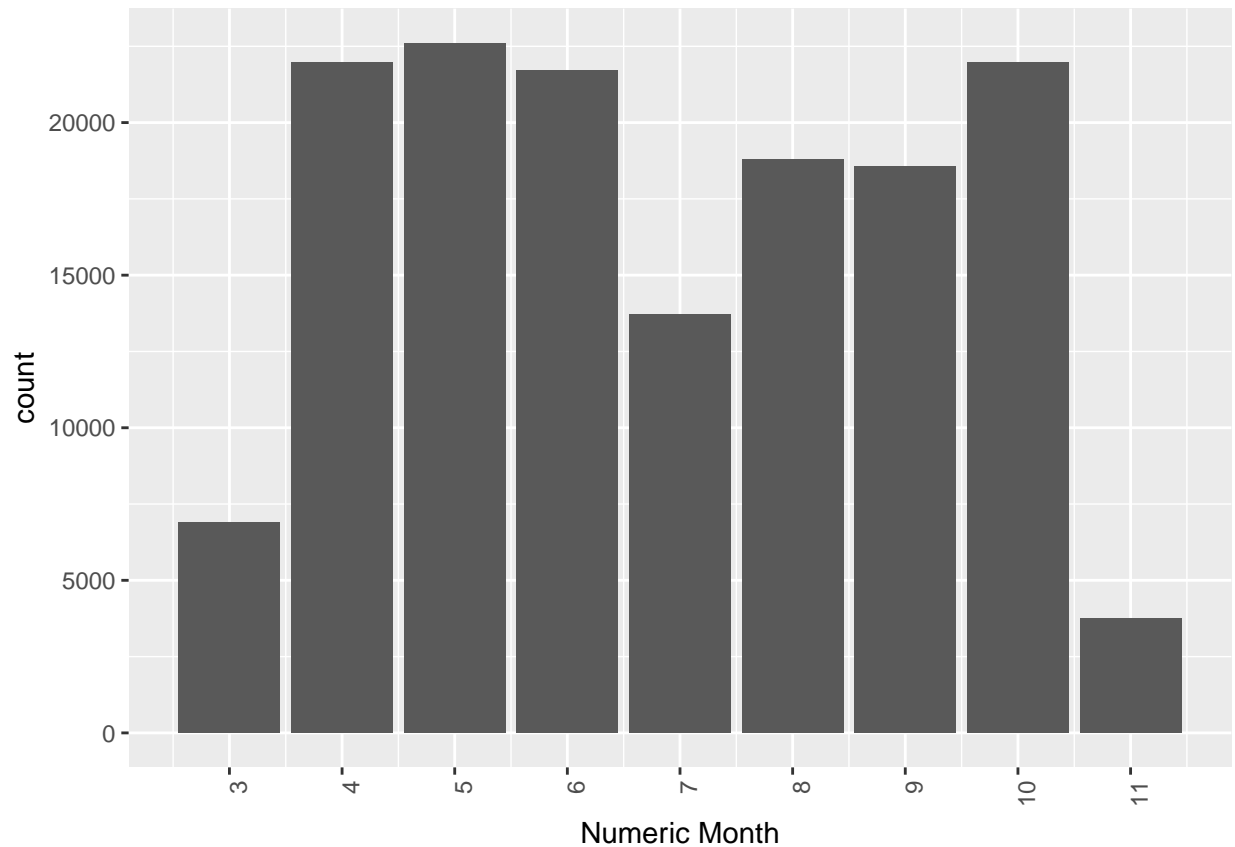
Next Station variable bar chart



NextStation variable represents the next station the train will stop.

title:

Month bar chart

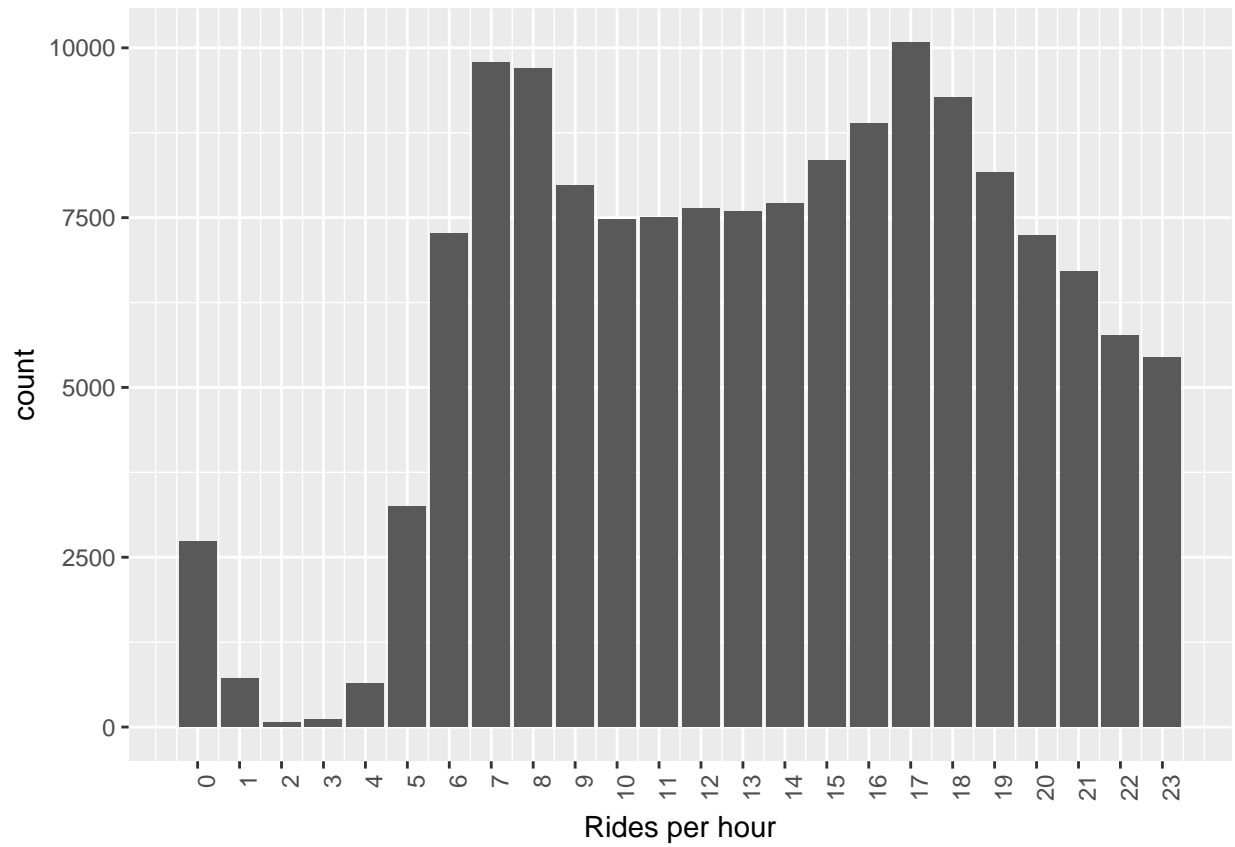


This chart is used to show the month with the most train rides and the least.

This is one of the new variables improvised by splitting up the timestamp variable. The numbers on the X-axis represent the numbers of real months. The y-axis represent the total number of train rides per month. The highest number of rides occurs in May, June and October.

title:

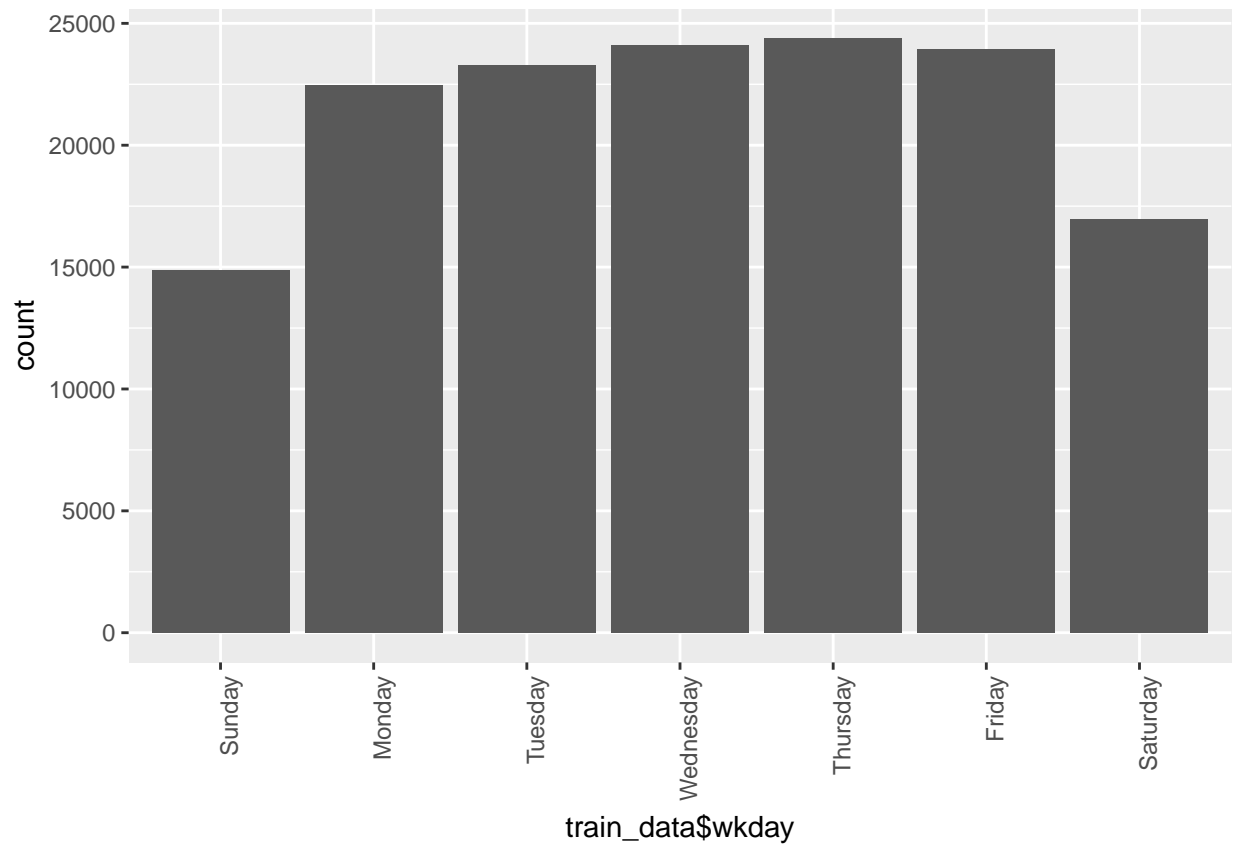
Hour of the day chart



The hour of the day chart shows the time that passengers ride the most.

title:

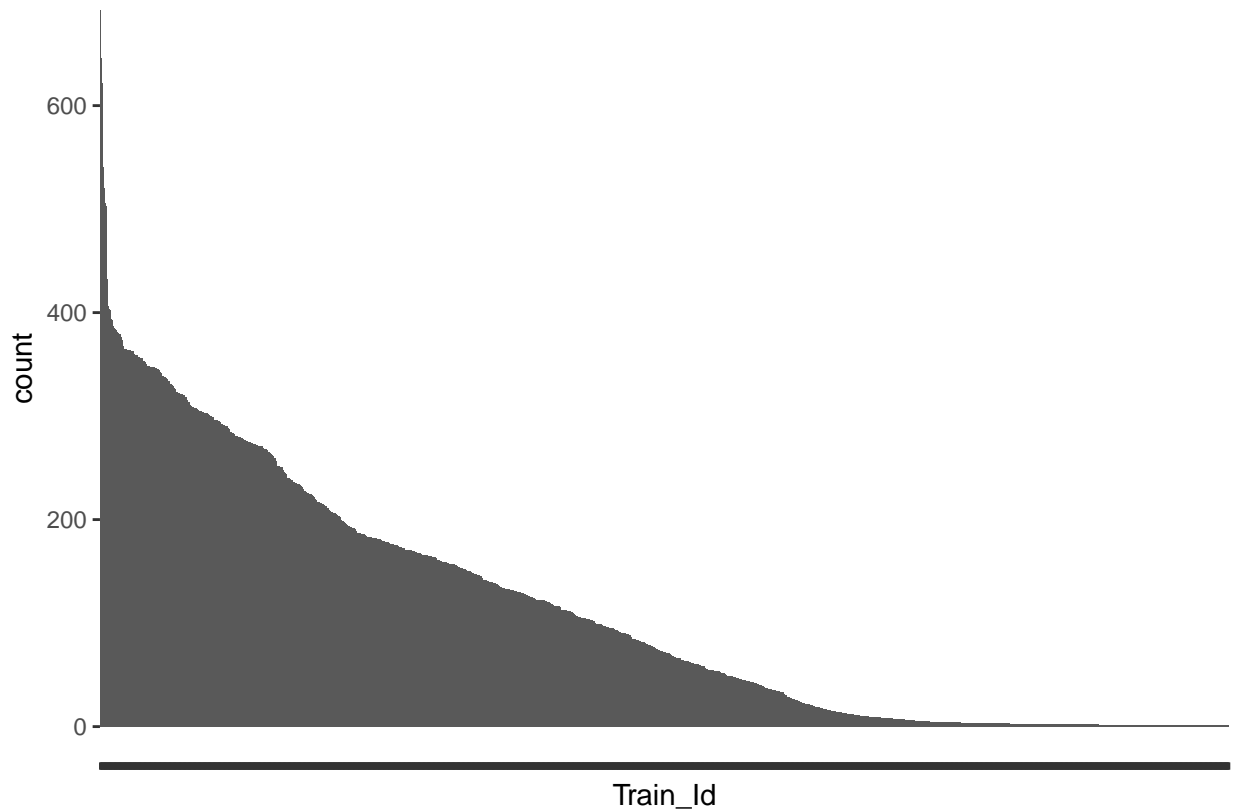
The Weekday chart



The weekday chart shows the number of trains that run different days of the week. More trains run during the week and fewer trains on the weekends. The busiest day is Thursday.

title:

Train id



This is a chart of all the train ids in descending order. These are all the trains that transported passengers during the period of one year in our dataset.

subtitle:

More conversion of observations into variables

Convert the weekday and monthday observations to columns. The purpose is to increase the number of variables that contribute to the status (delay and on time arrivals) of the train.

To achieve this, a matrix was used. In this case, the matrix translated the values of the new columns to “0” and “1”. “1” was printed when the train travelled in the specified day or month. “0” was used to fill the rest of the observations that the train did not ride during the week day or day of the month. for example, the new column “wkday1” represents “Monday”. The value “1” in that column represents the train rides that happened on Mondays.

subtitle:

CART Model /Decision Tree

:

In this section, the CART model is implemented. The two options considered are Classification and Regression CART models/trees but the regression model is preferred so that the results can be compared with the linear regression model used above. It’s like comparing apples to apples or oranges to oranges.

I found this site very helpful because they explained in detail the conditions for the variables before a succesful model can be achieved - https://rstudio-pubs-static.s3.amazonaws.com/27179_e64f0de316fc4f169d6ca300f18ee2aa.html.

Prior to finding this site, only the root or just one circle with a number (4.6) was drawn.

```
## [1] "The performance/r-squared value of the Testing data  0.207282993650837"
## [1] "The performance/r-squared value of the Training data  0.208606693807776"
```

title:

Check Performance of CART Regression model

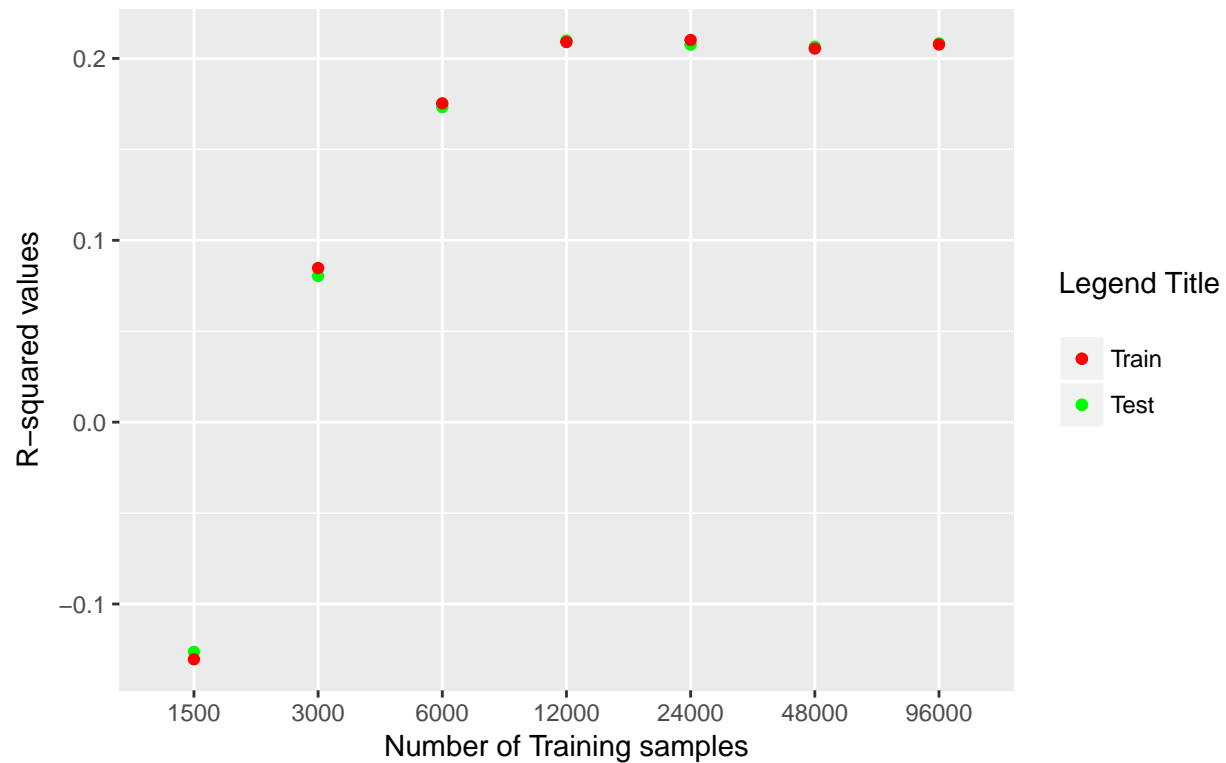
```
## [1] 1500
## [1] 3000
## [1] 6000
## [1] 12000
## [1] 24000
## [1] 48000
## [1] 96000
```

```
## [1] -0.12627769  0.08038021  0.17325549  0.20986960  0.20760042  0.20648902
## [7]  0.20828764
```

```
## [1] -0.13042049  0.08474824  0.17532162  0.20901182  0.21014193  0.20543047
## [7]  0.20764966
```

```
##      jj      R80trn      R80tst
## 1  1500 -0.12627769 -0.13042049
## 2  3000  0.08038021  0.08474824
## 3  6000  0.17325549  0.17532162
## 4 12000  0.20986960  0.20901182
## 5 24000  0.20760042  0.21014193
## 6 48000  0.20648902  0.20543047
## 7 96000  0.20828764  0.20764966
```

Performance CART (Decision Tree) model



Further test to check the performance of the Training and Test variables. The Trainset data will be trained and performance measured at 5%, 10%, 30%, 50%, 70%, 80% of the observations. A for loop is used to accomplish this in the code.

The peak of the performance was reached at the fourth loop with 20,000 records. There is no need to use more than 20,000 records since more data won't add more value to the performance.

title:

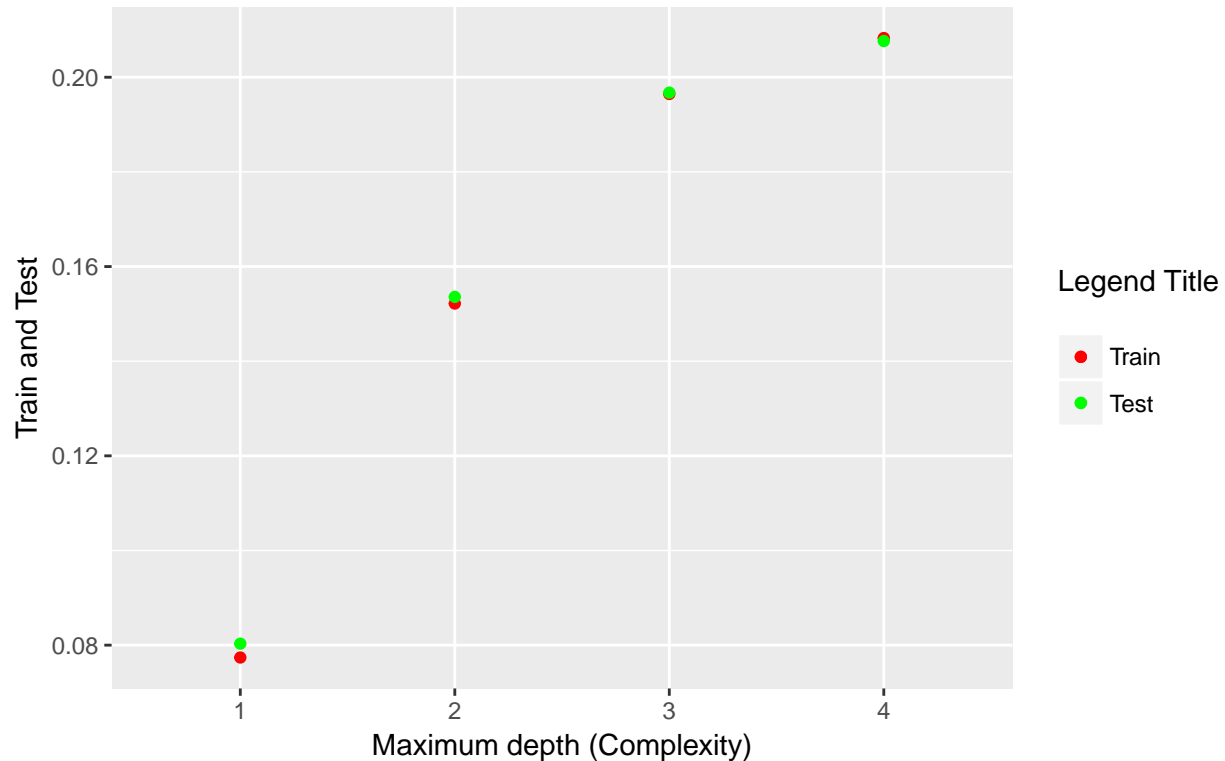
Checking max depth on the Test set for the CART models

```
## [1] 0.08031145 0.15358830 0.19674656 0.20764966
```

```
## [1] 0.07738756 0.15218426 0.19646121 0.20828764
```

```
##   jj      R2ts      R2tn
## 1  1 0.08031145 0.07738756
## 2  2 0.15358830 0.15218426
## 3  3 0.19674656 0.19646121
## 4  4 0.20764966 0.20828764
```

Decision Tree Learning curve Performance



As the depth increases, the complexity of the model increases.

The depths start leveling around the fourth level for both the Training and testing data. There is no need to increase the complexity further since it doesn't improve the performance of the data.

subtitle:

Random Forest model

Random Forest model implementation. This is the fourth model used in this project to try to come up with a reasonable prediction for the time delay prediction for the train dataset. The result at the end is similar to the Linear regression Cart model.

subtitle:

Check performance on the RandomForest Training and Testing data

The plots show that the Training set performed better than the Test set.

title:

Linear Regression model test

Linear model statistical summary for Status based on the independent variable $x = \text{origin}$. The linear regression model is one of the models implemented in this project to predict the delays of the train.

```
## [1] "The training data performance is 0.0822650944610929"
```

```
## [1] "The training data performance is 0.0807062477736474"
```

subtitle:

Conclusion from the models

The r-squared value from the Linear Regression model is 0.33506 ~34%. The Training set performance (r-squared) is .308 ~31%. The Testing set r-squared is 0.

The r-squared value from CART Regression model/ Decision Tree is .318 ~ 32%. The Training set best performance is 18% while the Test performance is 0.

The r-squared value from the Random Forest Training model at ntree=22 is .214 ~ 21% and as expected, the performance on the Training set is better than the performance on the Testing set. The r-squared value from the Random Forest Testing model is .101 ~ 10%. The word document named "finalproject_randomforest.docx" lists the different r-squared values at ntree values = 15, 10, 5 and 1. At the ntree value of 5, the r-squared or performance is 21% which is slightly better than the performance at ntree = 22 which is .202 ~ 20%.

The model that gave us the best performance is the Linear Regression model because its R-squared is slightly higher than the CART regression and Random Forest models. The Testing set performance is 31%, which is higher than the training performance of the other models.