# Weight Prediction : Flight Delays

Team Mint Squad

2011004534 윤덕진
2012004407 권규혁
2013012405 김현중
2013012676 이병곤

# CONTENTS

# Introduction

# 1)  Dataset of Flight Delays

## Dataset provided by <u>Bureau of Transportation Statistics</u>

# 1) Dataset of Flight Delays

## Dataset of Training - trainingFinal.csv

30000 Objects, including Total Delay column

| ID | Month | DayofMon | DayOfWe | FlightNun | ActualElap | CRSElapse | AirTime | ArrDelay | TotalDela | DepDelay | Origin | Dest | Distance | TaxiIn | TaxiOut | CarrierDel | WeatherD | NASDelay | SecurityD | LateAircraftDelay |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 3 | 4 | 335 | 128 | 150 | 116 | -14 | -6 | 8 | IAD | TPA | 810 | 4 | 8 | | | | | |
| 1 | 1 | 3 | 4 | 3231 | 128 | 145 | 113 | 2 | 21 | 19 | IAD | TPA | 810 | 5 | 10 | | | | | |
| 2 | 1 | 3 | 4 | 448 | 96 | 90 | 76 | 14 | 22 | 8 | IND | BWI | 515 | 3 | 17 | | | | | |
| 4 | 1 | 3 | 4 | 3920 | 90 | 90 | 77 | 34 | 68 | 34 | IND | BWI | 515 | 3 | 10 | 2 | 0 | 0 | 0 | 32 |
| 5 | 1 | 3 | 4 | 378 | 101 | 115 | 87 | 11 | 36 | 25 | IND | JAX | 688 | 4 | 10 | | | | | |
| 6 | 1 | 3 | 4 | 509 | 240 | 250 | 230 | 57 | 124 | 67 | IND | LAS | 1591 | 3 | 7 | 10 | 0 | 0 | 0 | 47 |
| 10 | 1 | 3 | 4 | 100 | 130 | 135 | 106 | 1 | 7 | 6 | IND | MCO | 828 | 5 | 19 | | | | | |
| 11 | 1 | 3 | 4 | 1333 | 121 | 135 | 107 | 80 | 174 | 94 | IND | MCO | 828 | 6 | 8 | 8 | 0 | 0 | 0 | 72 |
| 15 | 1 | 3 | 4 | 2272 | 52 | 50 | 37 | 11 | 20 | 9 | IND | MDW | 162 | 6 | 9 | | | | | |
| 16 | 1 | 3 | 4 | 675 | 228 | 240 | 213 | 15 | 42 | 27 | IND | PHX | 1489 | 7 | 8 | 3 | 0 | 0 | 0 | 12 |
| 17 | 1 | 3 | 4 | 1144 | 226 | 250 | 205 | -15 | -6 | 9 | IND | PHX | 1489 | 5 | 16 | | | | | |
| 18 | 1 | 3 | 4 | 4 | 123 | 135 | 110 | 16 | 44 | 28 | IND | TPA | 838 | 4 | 9 | 0 | 0 | 0 | 0 | 16 |
| 19 | 1 | 3 | 4 | 54 | 56 | 70 | 49 | 37 | 88 | 51 | ISP | BWI | 220 | 2 | 5 | 12 | 0 | 0 | 0 | 25 |
| 21 | 1 | 3 | 4 | 623 | 57 | 70 | 47 | 19 | 51 | 32 | ISP | BWI | 220 | 5 | 5 | 7 | 0 | 0 | 0 | 12 |
| 22 | 1 | 3 | 4 | 717 | 56 | 70 | 49 | 6 | 26 | 20 | ISP | BWI | 220 | 2 | 5 | | | | | |
| 23 | 1 | 3 | 4 | 1244 | 54 | 70 | 47 | -7 | 2 | 9 | ISP | BWI | 220 | 2 | 5 | | | | | |
| 25 | 1 | 3 | 4 | 2553 | 59 | 70 | 50 | 14 | 39 | 25 | ISP | BWI | 220 | 2 | 7 | | | | | |
| 26 | 1 | 3 | 4 | 188 | 155 | 195 | 143 | 47 | 134 | 87 | ISP | FLL | 1093 | 6 | 6 | 40 | 0 | 0 | 0 | 7 |
| 27 | 1 | 3 | 4 | 1754 | 165 | 190 | 155 | 4 | 33 | 29 | ISP | FLL | 1093 | 3 | 7 | | | | | |
| 30 | 1 | 3 | 4 | 362 | 147 | 165 | 134 | 64 | 146 | 82 | ISP | MCO | 972 | 6 | 7 | 5 | 0 | 0 | 0 | 59 |
| 33 | 1 | 3 | 4 | 1397 | 154 | 170 | 140 | -4 | 8 | 12 | ISP | MCO | 972 | 7 | 7 | | | | | |
| 34 | 1 | 3 | 4 | 3398 | 146 | 170 | 134 | -5 | 14 | 19 | ISP | MCO | 972 | 6 | 6 | | | | | |
| 35 | 1 | 3 | 4 | 3480 | 145 | 170 | 134 | 14 | 53 | 39 | ISP | MCO | 972 | 5 | 6 | | | | | |
| 37 | 1 | 3 | 4 | 422 | 135 | 145 | 118 | 72 | 154 | 82 | ISP | MDW | 765 | 6 | 11 | 3 | 0 | 0 | 0 | 69 |
| 38 | 1 | 3 | 4 | 1837 | 128 | 145 | 114 | 5 | 27 | 22 | ISP | MDW | 765 | 9 | 5 | | | | | |
| 39 | 1 | 3 | 4 | 2871 | 127 | 145 | 113 | 11 | 40 | 29 | ISP | MDW | 765 | 8 | 6 | | | | | |
| 40 | 1 | 3 | 4 | 1056 | 153 | 180 | 143 | 29 | 85 | 56 | ISP | PBI | 1052 | 5 | 5 | 0 | 0 | 0 | 0 | 29 |

# 1) Dataset of Flight Delays

## Dataset of Test - testFinal.csv

300 Objects, excluding Total Delay column

# 1) Dataset of Flight Delays

## Variables :

| | | | |
|---|---|---|---|
| DayOf Month | Flight Number | CRS Elapsed Time | Actual Elapsed Time |
| Distance | TaxiIn | TaxiOut | Carrier Delay |
| Airtime | NAS Delay | Security Delay | Late AirCraft Delay |

Weather Delay

# 2) Our Goal

```
Change
Algorithm / Tool        01100
                        10110
                        11110
```

Analysis
Train Dataset

Predict
Test Dataset

Find the most
influential Weight

**Our Goal**

**Predict total delay of flights without weather conditions.**

PART 2

—

# Data Analysis

# 1) Dataset of Flight Delays

## Variables :



DayOf Month

Flight Number

CRS Elapsed Time

Actual Elapsed Time

Distance

TaxiIn

TaxiOut

Carrier Delay

Airtime

NAS Delay

Security Delay

Late AirCraft Delay

Weather Delay

# 2) Variables Description

**DayOfMonth –** the impacts by the day of month.

**Flight Number –** The influence of different types of airplanes.

**CRS/Actual Elapsed Time –** Difference between scheduled time and actual departure time.

**Distance –** Flight distance.

**TaxiIn/Out –** Time to get on and off the plane.

**Carrier Time –** Time to find luggage.

**Air Time –** Flight time.

# 2) Variables Description

**NAS Delay –** Time it takes in National Air System delay

**Security Delay –** Time it takes to perform the security check.

**Late Aircraft Delay –** time it take to maintain an airplane.

**Weather Delay –** Time delayed due to weather.

# Pre-Processing

# 1) Variables Format

| DayOfMonth | 1-12 | CarrierDelay | In minutes |
|---|---|---|---|
| FlightNumber | Factor | AirTime | In minutes |
| CRSElapsedTime | In minutes | NASDelay | In minutes |
| ActualElapsedTime | In minutes | SecurityDelay | In minutes |
| Distance | In miles | LateAirCraftDelay | In minutes |
| TaxiIn/Out | Taxi in time, in minutes | WeatherDelay | In minutes |

# 1) Variables Format

```
 4   # Train Attributes 20
 5   train_data$ID <- as.character(train_data$ID)
 6   train_data$FlightNum <- as.factor (train_data$FlightNum)
 7   train_data$DayofMonth <- as.factor (train_data$DayofMonth)
 8
 9   # Test Attributes 19
10   test_data$ID <- as.character(test_data$ID)
11   test_data$FlightNum <- as.factor (test_data$FlightNum)
12   test_data$DayofMonth <- as.factor (test_data$DayofMonth)
```

Description

**We set non-numeric values in categories.**

# 2) Incomplete dataset

```
> sum(is.na(train_data$DayOfWeek))
[1] 0
> sum(is.na(train_data$FlightNum))
[1] 0
> sum(is.na(train_data$CRSElapsedTime))
[1] 0
> sum(is.na(train_data$ActualElapsedTime))
[1] 78
> sum(is.na(train_data$Distance))
[1] 0
> sum(is.na(train_data$TaxiIn))
[1] 78
> sum(is.na(train_data$TaxiOut))
[1] 0
> sum(is.na(train_data$CarrierDelay))
[1] 13554
> sum(is.na(train_data$AirTime))
[1] 78
> sum(is.na(train_data$NASDelay))
[1] 13554
> sum(is.na(train_data$SecurityDelay))
[1] 13554
> sum(is.na(train_data$LateAircraftDelay))
```

**Train.csv**

```
> sum(is.na(test_data$DayOfWeek))
[1] 0
> sum(is.na(test_data$FlightNum))
[1] 0
> sum(is.na(test_data$CRSElapsedTime))
[1] 0
> sum(is.na(test_data$ActualElapsedTime))
[1] 3
> sum(is.na(test_data$Distance))
[1] 0
> sum(is.na(test_data$TaxiIn))
[1] 3
> sum(is.na(test_data$TaxiOut))
[1] 0
> sum(is.na(test_data$CarrierDelay))
[1] 62
> sum(is.na(test_data$AirTime))
[1] 3
> sum(is.na(test_data$NASDelay))
[1] 62
> sum(is.na(test_data$SecurityDelay))
[1] 62
> sum(is.na(test_data$LateAircraftDelay))
```

**Test.csv**

Problem

**We have found that NULL values exist.**

# 3) complete dataset

```r
14  # Null of ActualElapsedTime
15  train_data$ActualElapsedTime[is.na(train_data$ActualElapsedTime)] <- median(train_data$ActualElapsedTime, na.rm=TRUE)
16  test_data$ActualElapsedTime[is.na(test_data$ActualElapsedTime)] <- median(test_data$ActualElapsedTime, na.rm=TRUE)
17
18  # Null of AirTime
19  train_data$AirTime[is.na(train_data$AirTime)] <- median(train_data$AirTime, na.rm=TRUE)
20  test_data$AirTime[is.na(test_data$AirTime)] <- median(test_data$AirTime, na.rm=TRUE)
21
22  #Null of TaxiIn
23  train_data$TaxiIn[is.na(train_data$TaxiIn)] <- median(train_data$TaxiIn, na.rm=TRUE)
24  test_data$TaxiIn[is.na(test_data$TaxiIn)] <- median(test_data$TaxiIn, na.rm=TRUE)
```

```r
20  # Train 5가지 Delay 빈값 = 0
21  train_data$CarrierDelay[is.na(train_data$CarrierDelay)] <- 0
22  train_data$WeatherDelay[is.na(train_data$WeatherDelay)] <- 0
23  train_data$NASDelay[is.na(train_data$NASDelay)] <- 0
24  train_data$SecurityDelay[is.na(train_data$SecurityDelay)] <- 0
25  train_data$LateAircraftDelay[is.na(train_data$LateAircraftDelay)] <- 0
26
27  # Test 5가지 Delay 빈값 = 0
28  test_data$CarrierDelay[is.na(test_data$CarrierDelay)] <- 0
29  test_data$WeatherDelay[is.na(test_data$WeatherDelay)] <- 0
30  test_data$NASDelay[is.na(test_data$NASDelay)] <- 0
31  test_data$SecurityDelay[is.na(test_data$SecurityDelay)] <- 0
32  test_data$LateAircraftDelay[is.na(test_data$LateAircraftDelay)] <- 0
```
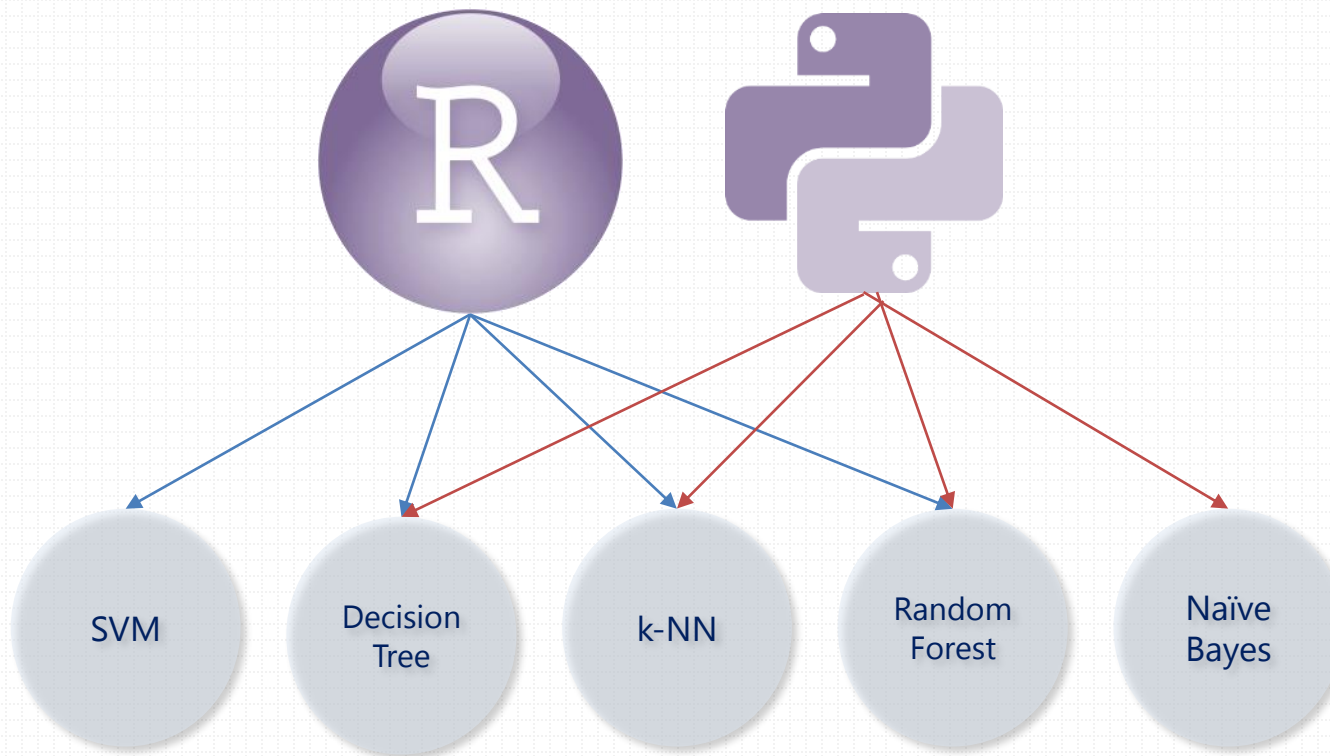
Solution

**We have replaced Null values with average values.**

# Algorithm and Models

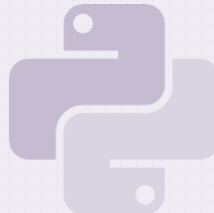# 1) 2 Tools & 5 Algorithms

# 2) Decision Tree

[Decision Tree Classifier]

Decision tree analysis is the most popular decision support technique in data mining field. From the root, input training set is divided recursively at split point of input features(attributes) making a tree. The benefit of using decision tree is simple understanding of classification process.

| Python | R |
|---|---|
| Library : sklearn.tree.DecisionTreeClassifier | Library : library(rpart) |
| Accuracy Result : **79.57%** | Accuracy Result : **73.81%** |

# 3) k-NN (k-Nearest Neighbors)

[k-NN(k-nearest neighbors)]

kNN is a prediction method for classification as well as regression type prediction problems. Its one of the simplest machine-learning algorithm, that test case is simply assigned to the class that most k nearest training cases are found. The result get different depending on k value you set.



| Python | R |
|---|---|
| Library : sklearn.neighbors.KNeighborsClassifier | Library : library(kknn) |
| Accuracy Result : **55.89%** | Accuracy Result : **34.16%** |

# 4) Random Forest

[Random Forest Classifier]

A random forest is a meta classifier that has a number of decision tree classifiers on various sub-samples of the dataset. It uses mode of the classes for classification and mean prediction for regression to improve the predictive accuracy and control over-fitting.



| Python | R |
|---|---|
| Library : sklearn.ensemble.RandomForestClassifier | Library : library(randomForest) |
| Accuracy Result : **73.99%** | Accuracy Result : **86.86%** |

# 5) Naïve Bayes

[Naive Bayes Net Classifier]

Naive Bayes Classification is a simple probabilistic classification based on applying Bayes theorem using the naive independence assumptions. In this algorithm, there is an assumption that every features are independent between each other.



### Python

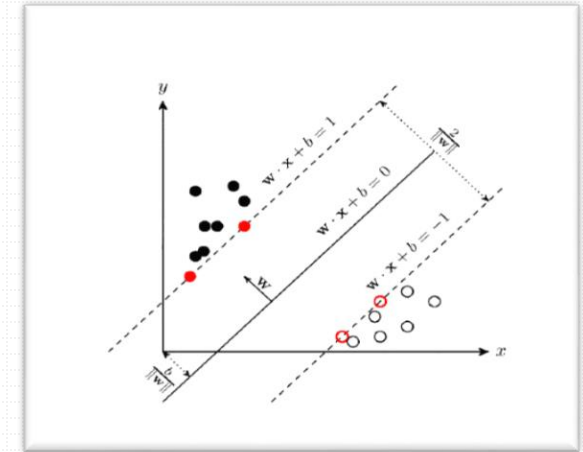Library : sklearn.naive.bayes.GaussianNB

Accuray Result : **45.94%**

# 6) SVM (Support Vector Machine)

[Support Vector Machine]

SVMs are supervised learning algorithms that are mostly used for classification and regression. SVMs produce linear classifiers called hyperplane that separate the data into multiple subsections.
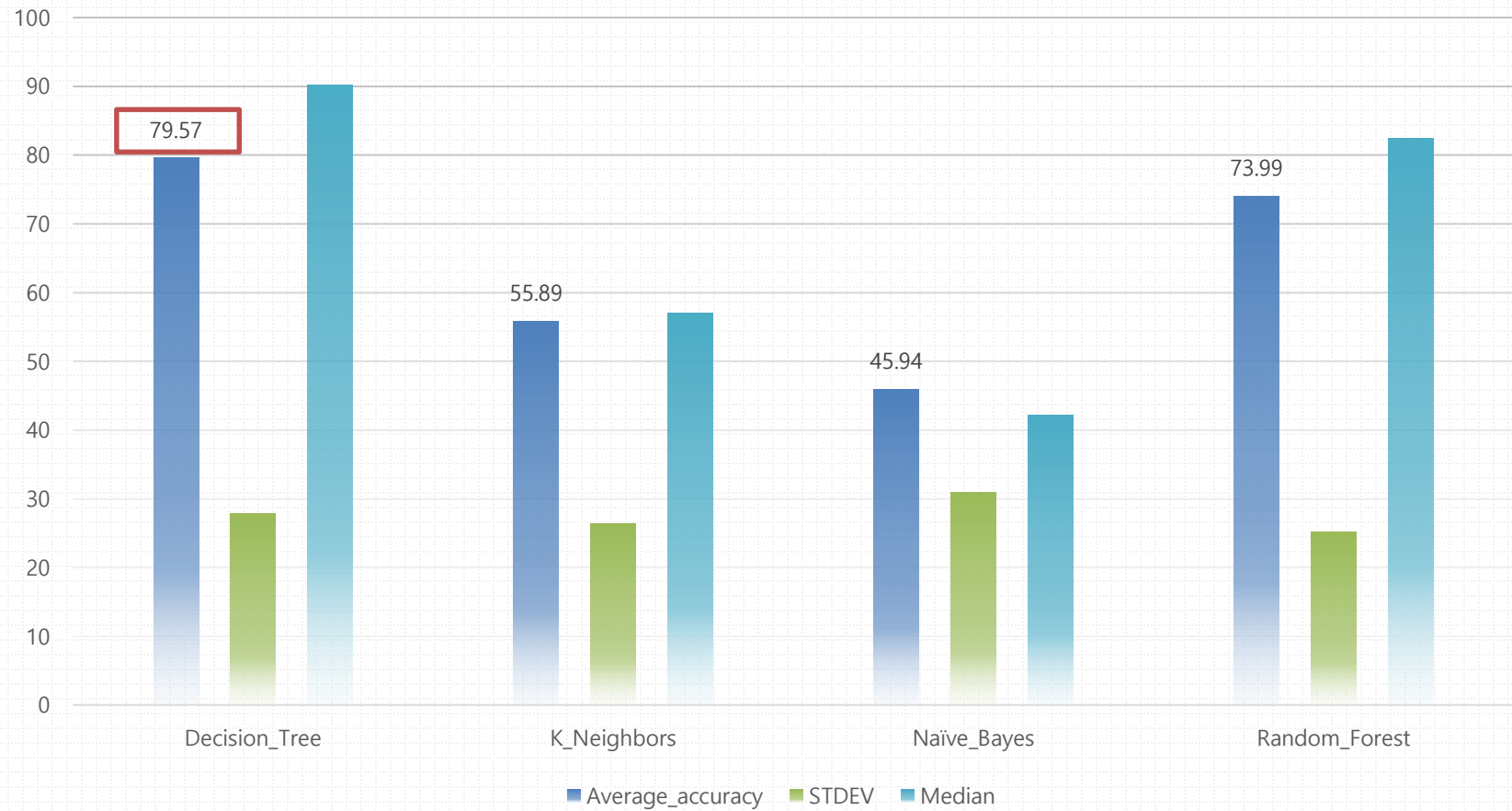

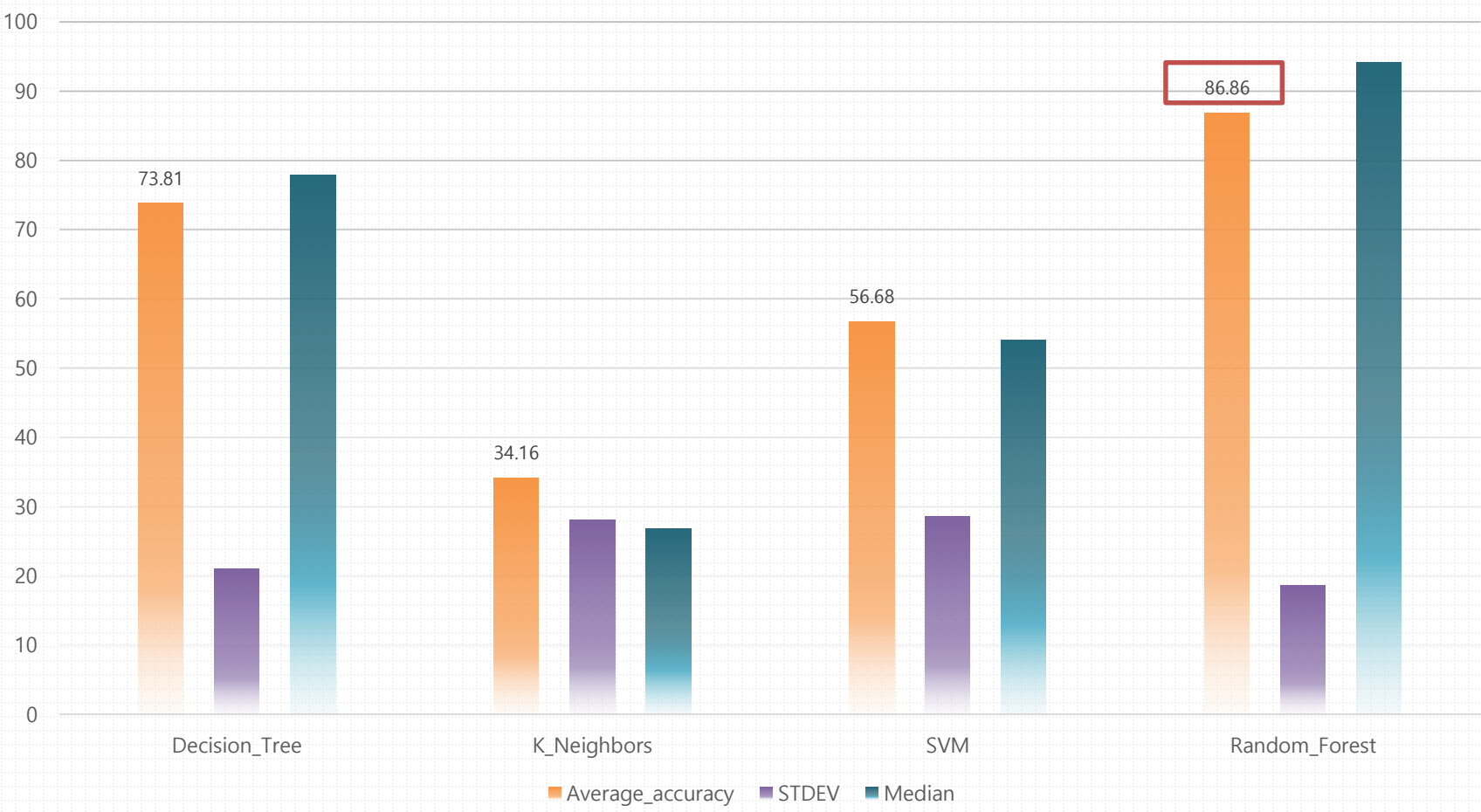
R

Library : library(kernlab)

Accuray Result : **56.68%**

# Algorithm Results

## 1) Python

## 2) R

# 3) Compare Accuracy Python & R

| R | Algorithm | Python |
|---|---|---|
| 73.81% | Decision Tree | 79.57% |
| 86.86% | Random Forest | 73.99% |
| 34.16% | K-NN | 55.89% |
| | Naïve Bayes | 45.94% |
| 56.58% | SVM | |

**Result**

**Algorithms with highest accuracy that fit our dataset are Decision Tree & RandomForest**
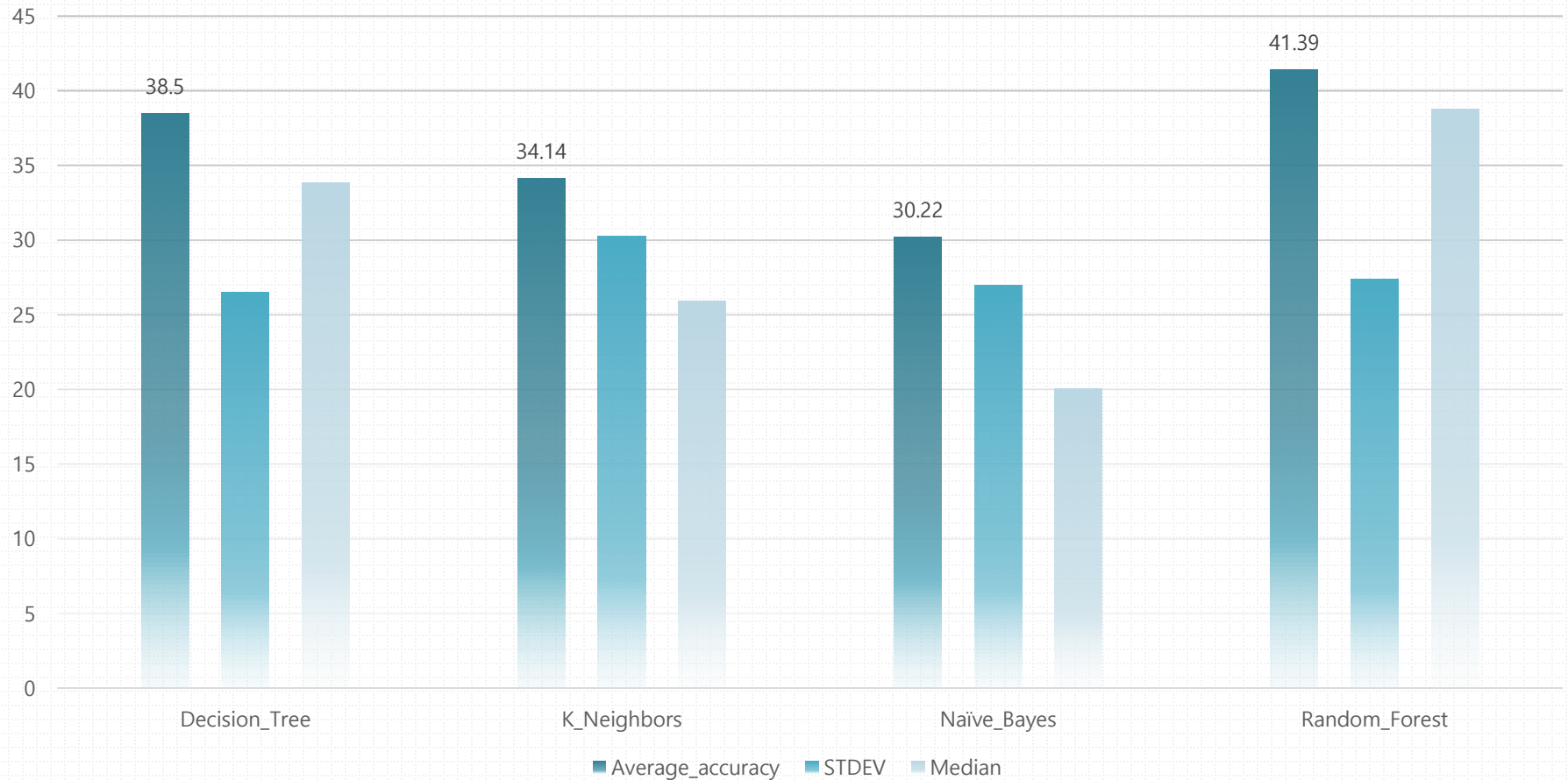
# 4) Select the most 6 influential weight

## Python

```
  ('Month', 0.0)
2 ('DayofMonth', 0.11247138377205146)
  ('DayOfWeek', 0.0766193284139051161)
1 ('FlightNum', 0.13885770294697292)
  ('ActualElapsedTime', 0.088081910076580242)
  ('CRSElapsedTime', 0.056354987631322184)
5 ('AirTime', 0.090983138460288035)
3 ('Distance', 0.10704811910243951)
  ('TaxiIn', 0.085552483282536118)
6 ('TaxiOut', 0.092105597532971623)
4 ('CarrierDelay', 0.047890236735033974)
  ('WeatherDelay', 0.0070799993406601893)
  ('NASDelay', 0.044285889182208056)
  ('SecurityDelay', 0.0015940150062221661)
  ('LateAircraftDelay', 0.051075208516808308)
  1.0
```

## R

```
> fit$variable.importance
 FlightNum DayofMonth    TaxiOut    Distance    AirTime     TaxiIn
42369268.7 16412143.4  1544937.4  1051584.0   737314.7   735131.1
```
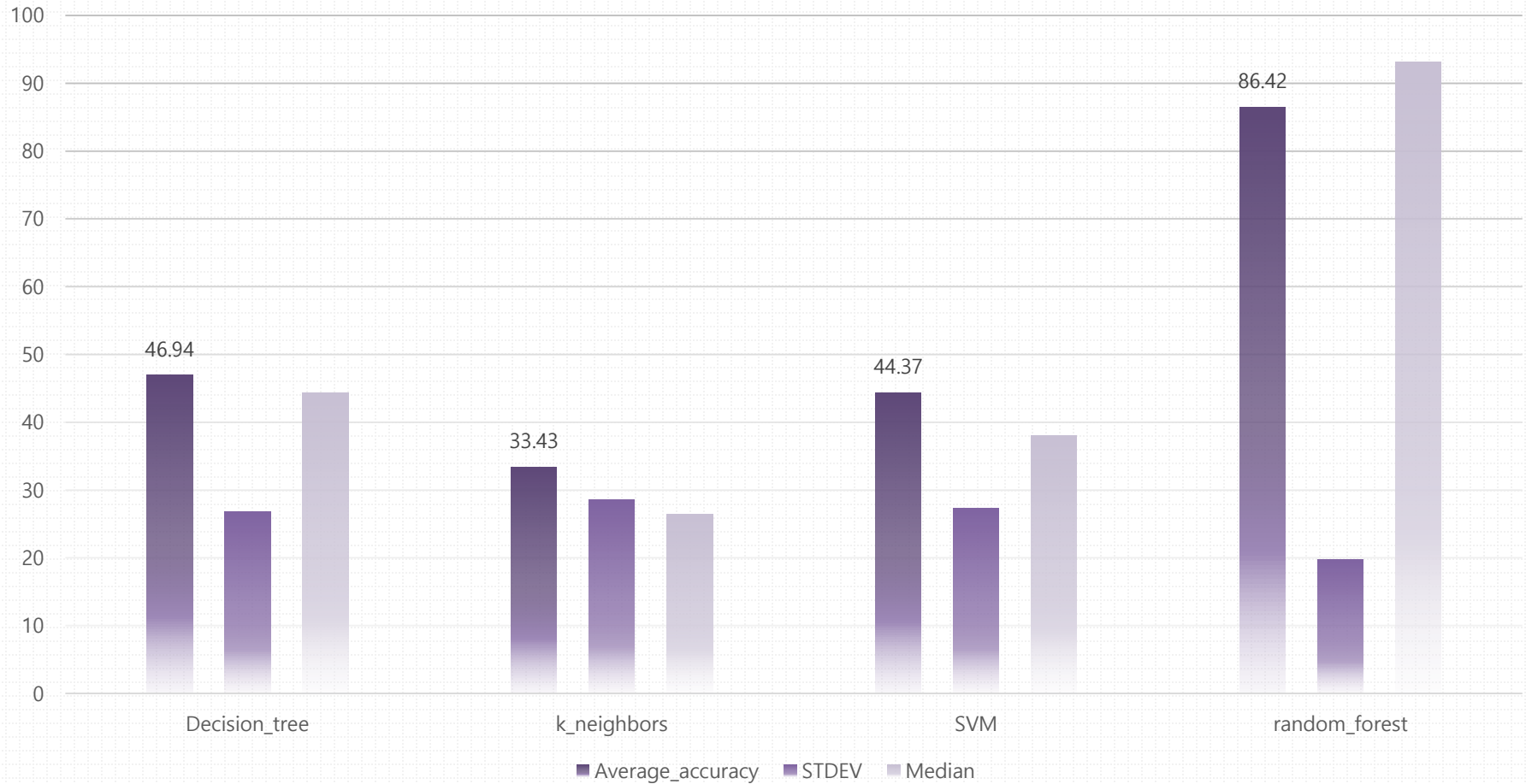
# 4) Python with Top6 features

# 5) Compare python with Top 6 & with all features

| With all features | Algorithm | With Top 6 features |
|---|---|---|
| 79.57% | Decision Tree | 38.5% |
| 73.99% | Random Forest | 34.14% |
| 55.89% | K-NN | 30.22% |
| 45.94% | Naïve Bayes | 41.49% |

# 6) R with Top6 features

# 7) Compare R with Top 6 & with all features

| With all features | Algorithm | With Top 6 features |
|:---:|:---:|:---:|
| 73.81% | Decision Tree | 46.94% |
| 86.86% | Random Forest | 86.42% |
| 34.16% | K-NN | 33.43% |
| 56.58% | SVM | 44.37% |

**Result**

**As a result of selecting and running the top 6 variables, the accuracy was lower than all the other variables.**

# 8) Conclusion

1. After comparing the accuracy using several algorithms, **Decision Tree and RandomForest fit our dataset the best.**

2. As a result of checking the weights of all variables through the above two algorithms, **we confirmed that the top 6 variables (FlightNumber, DayOfMonths, Distance, TaxiOut, AirTime, TaxiIn) affect 65%  of the total delay.**

3. As a result of estimating the delay with only the top 6 variables, **we confirmed that considering less variables leads to a considerable loss of accuracy.**

## Our Conclusion

After having reviewed the results of our project, we have concluded that considering more variables for the machine increases accuracy. Furthermore, the dataset itself is not appropriate for machine learning. So the more concrete data is needed to our study.

감사합니다.