

Latam Airlines : Data Scientist Challenge

By: Aakruti Ambasana

Sections

- [Question 1: Data distribution](#)
- [Question 2: Generating Synthetic features](#)
- [Question 3: Data Visualization](#)
- [Question 4: Estimated likelihood of the flight delay](#)
- [Question 5: Predictive model and evaluation of model performance](#)

```
In [1]: import pandas as pd
import numpy as np
import seaborn as sb
from matplotlib import pyplot as plt

from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor
from sklearn.ensemble import RandomForestClassifier, RandomForestRegressor

from xgboost import XGBClassifier, XGBRegressor

from sklearn.model_selection import GridSearchCV, train_test_split
from sklearn import import tree, metrics

from sklearn.metrics import confusion_matrix, accuracy_score, classification_report,
mean_squared_error
```

Loading data

```
In [2]: types = {'Vlo-I': 'string', 'Ori-I': 'category', 'Des-I': 'category',
'Emp-I': 'category', 'Vlo-O': 'string', 'Ori-O': 'category',
'Des-O': 'category', 'Emp-O': 'category', 'DIA': 'int', 'MES': 'int', 'AÑO':
'int',
'DIANOM': 'category', 'TIPOVUELO': 'string', 'OPERA': 'category', 'SIGLAORI':
'category',
'SIGLADES': 'category'}
path = "data/dataset_SCL.csv"
data = pd.read_csv(path, dtype=types, parse_dates=['Fecha-I', 'Fecha-O'])
data.head()
```

```
Out[2]:
```

	Fecha-I	Vlo-I	Ori-I	Des-I	Emp-I	Fecha-O	Vlo-O	Ori-O	Des-O	Emp-O	DIA	MES	AÑO	DIANOM	TIPOVUELO	OPERA	SIGLAORI	SIGLADES
0	2017-01-01	226	SECL	KMIA	AAL	2017-01-01	226	SECL	KMIA	AAL	1	1	2017	Domingo				American Airlines
1	2017-01-02	226	SECL	KMIA	AAL	2017-01-02	226	SECL	KMIA	AAL	2	1	2017	Lunes				American Airlines
2	2017-01-03	226	SECL	KMIA	AAL	2017-01-03	226	SECL	KMIA	AAL	3	1	2017	Martes				American Airlines
3	2017-01-04	226	SECL	KMIA	AAL	2017-01-04	226	SECL	KMIA	AAL	4	1	2017	Miercoles				American Airlines
4	2017-01-05	226	SECL	KMIA	AAL	2017-01-05	226	SECL	KMIA	AAL	5	1	2017	Jueves				American Airlines

• Vlo-I contains alpha numeric value.

• Vlo-O contains NaN null values.

• default is object type. object contains any kind of data like strings, integers, etc.

• string datatype contains string.

Converting Vlo-I and Vlo-O into numeric values, as it contains flight number

- There are 16 data records in Vlo-O and 5 data records in Vlo-I which contains alpha numeric values.

```
In [3]: print ("Vlo-O:",data["Vlo-O"].str.contains("[A-Za-z]").sum())
print ("Vlo-I:",data["Vlo-I"].str.contains("[A-Za-z]").sum())

Vlo-O: 16
Vlo-I: 5
```

```
In [4]: # Replacing string with empty string
data["Vlo-I"] = data["Vlo-I"].str.replace("[A-Za-z]", "", regex=True)
data["Vlo-I"] = data["Vlo-I"].str.replace("[A-Za-z]", "", regex=True)
```

```
In [5]: # Checking that Vlo-I, Vlo-O value does not contain any alphabet now.
print ("After Vlo-O:",data["Vlo-O"].str.contains("[A-Za-z]").sum())
print ("After Vlo-I:",data["Vlo-I"].str.contains("[A-Za-z]").sum())

After Vlo-O: 0
After Vlo-I: 0
```

Dealing with missing data Vlo-O

```
In [6]: print(data.isna().sum())

Fecha-I  0
Vlo-I    0
Ori-I    0
Des-I    0
Emp-I    0
Fecha-O  0
Vlo-O    1
Ori-O    0
Des-O    0
Emp-O    0
DIA      0
MES      0
AÑO      0
DIANOM   0
TIPOVUELO 0
OPERA     0
SIGLAORI 0
SIGLADES 0
dtype: int64
```

```
In [7]: # Vlo-O contains one row containing null value
print(data.loc[data['Vlo-O'].isna()])

      Fecha-I  Vlo-I  Ori-I  Des-I  Emp-I  Fecha-O  Vlo-O  \
6968 2017-01-19 11:00:00 200  SECL  SPJC  LW  2017-01-19 11:03:00 <NA>
      Ori-O  Des-O  Emp-O  DIA  MES  AÑO  DIANOM  TIPOVUELO  \
6968  SECL  SPJC  56R  19  1  2017  Jueves  I

      OPERA  SIGLAORI  SIGLADES
6968  Latin American Kings  Santiago  Lima
```

Assigning Vlo-I to Vlo-O because Ori-I is same as Ori-O, Des-I is same as Des-O, and if DIANOM is Jueves, there is very high chances that Vlo-I is same as Vlo-O, and there is no delay also.

```
In [8]: data["Vlo-O"] = data["Vlo-O"].fillna(data["Vlo-I"])
data["Vlo-I"] = data["Vlo-I"].astype("int")
data["Vlo-O"] = data["Vlo-O"].astype("float")
data["Vlo-O"] = data["Vlo-O"].astype("int")
```

Some Vlo-O values are specified in **decimal** form but is integer. So, first converting into float and then int.

Converting TIPOVUELO feature into integer

International flights are assigned with 1 and national flights are assigned with 0 values.

```
In [9]: data["TIPOVUELO"] = data["TIPOVUELO"].astype('object')
data["TIPOVUELO"] = data["TIPOVUELO"].replace({'I', 'N'},[1,0]).astype("int")
```

Data Analysis

```
In [10]: data.describe().T

Out[10]:
```

	count	unique	top	Freq
Vlo-I	68206.0	969.827288	2029.024762	1.0
Vlo-O	68206.0	967.421092	2026.193621	1.0
DIA	68206.0	15.714790	8.782886	1.0
MES	68206.0	6.622585	3.523321	1.0
AÑO	68206.0	2017.000029	0.005415	2017.0
TIPOVUELO	68206.0	0.458024	0.498239	0.0

Here, Vlo-I and Vlo-O values are skewed, because mean value is greater than 50%. There is big gap between 75% and max values of Vlo-I and Vlo-O features.

```
In [11]: print(data.describe(include="category").T)

count unique top Freq
Ori-I 68206 1 SCFL 68206
Des-I 68206 64 SCFA 5787
Emp-I 68206 30 LAN 37611
Ori-O 68206 1 SCFL 68206
Des-O 68206 63 SCFA 5786
Emp-O 68206 32 LAN 29988
Vlo-O 68206 7 Viernes 18202
OPERA 68206 23 Grupo LATAM 48892
SIGLAORI 68206 1 Santiago 68206
SIGLADES 68206 62 Buenos Aires 6335
```

- Ori-I, Ori-O, and SIGLAORI features contains only one value. So we will not include this features into input while training the model.
- Analyze Des-I, Des-I, and SIGLADES which contains almost same number of categories. So some features can be duplicating.
- Compare Emp-I, Emp-O and OPERA features. Keep the features which are useful, and remove the duplicate columns.

```
In [12]: pd.crosstab(data["OPERA"],data["DIANOM"])

Out[12]:
```

	DIANOM	Domingo	Jueves	Lunes	Martes	Miercoles	Sabado	Viernes
OPERA								
Aerolineas Argentinas								
Aeromexico								
Air Canada								
Alitalia								
American Airlines								
Austral								
Avianca								
British Airways								
Copa Air								
Delta Air								
Go Trans								
Grupo LATAM								
Iberia								
K.L.M.								
Latin American Wings								
Qantas Airways								
Sky Airline								
United Airlines								
JetSmart SCA								
Laca								
Oceanair Linhas Aereas								
Plus Ultra Lineas Aereas								

Maximum data contains of flights operated by LATAM Airlines.

Some Airlines only fly on some days of the week.

- Alitalia has only one flight on Miercoles, but on other days of week approx. 50 flights are operated.
- Austral has no flights on Lunes, Miercoles, and Viernes.
- K.L.M. has less than 10 flights on Domingo and Viernes, while on other days of week it has approx 50 flights.
- Plus Ultra Lineas Aereas has flights on Domingo, Miercoles, and Viernes, while on other days of week it only have couple of flights.
- Qantas Airways has combined less than 10 flights on Jueves, Lunes, and Sabado, while on other days of week approx. 47 flights are operated.

```
In [13]: data.hist(figsize=(10,10))

Out[13]:
```

```
In [14]: print(data[["AÑO","TIPOVUELO"]].value_counts())

AÑO TIPOVUELO
2017 0 36966
2018 1 31298
2019 1 2
dtype: int64
```

Insights from histogram

- In MES, January and December month contains the highest number of flights operated. It is vacation time, so people travel a lot at that time.
- 54% flights are national, and 45% flights are international.
- AÑO is highly skewed and data is distributed more vertically. Because whole data is about 2017, only two data records are of 2018.

Q 1: Data Distribution

Calculating Skew

- Skewness means data distribution is not uniform means it has less symmetry. The shape of curve represents the data distribution.
- If curve is positively skewed then most of the values are less than median value. If curve is negatively skewed then most of the data is greater than median value.
- The value zero mean data distribution is symmetric.

```
In [15]: print(data.skew())

Vlo-I 3.893216
Vlo-O 3.183728
DIA 0.988439
MES -4.865328
AÑO 184.665914
TIPOVUELO 0.168582
dtype: float64
```

Observations

- Vlo-I and Vlo-O are positively skewed.
- AÑO is highly positively skewed, because all data is from 2017 but only 2 data records from 2018.
- Apparently, MES is slightly negatively skewed. DIA is closest to symmetric distribution.

Calculating Kurtosis

Kurtosis measures peak point of curve of data. There are 3 types of curve:

1. **Leptokurtic Curve:** This curve is **taller** than normal distribution curve. Its value is greater than 0.
2. **Mesokurtic Curve:** This curve is closest to normal distribution curve. Its value is 0.
3. **Platykurtic Curve:** The peak of this curve is **flat**. It is flatter than other 2 curves. The value is less than 0.

```
In [16]: print(data.kurtosis())

Vlo-I 8.893281
Vlo-O 8.162987
DIA -1.192468
MES -1.249986
AÑO 34106.499883
TIPOVUELO -1.971665
dtype: float64
```

Observations

- Vlo-I and Vlo-O will be represented by leptokurtic curve means data distribution is more vertical.
- DIA and MES represents platykurtic curve means data is distributed more horizontally but it is closest to normal data distribution.
- **Data distribution of AÑO is highly vertical represented by leptokurtic curve.**

Data distribution for categorical features

To view data distribution of categorical features we need to plot the features.

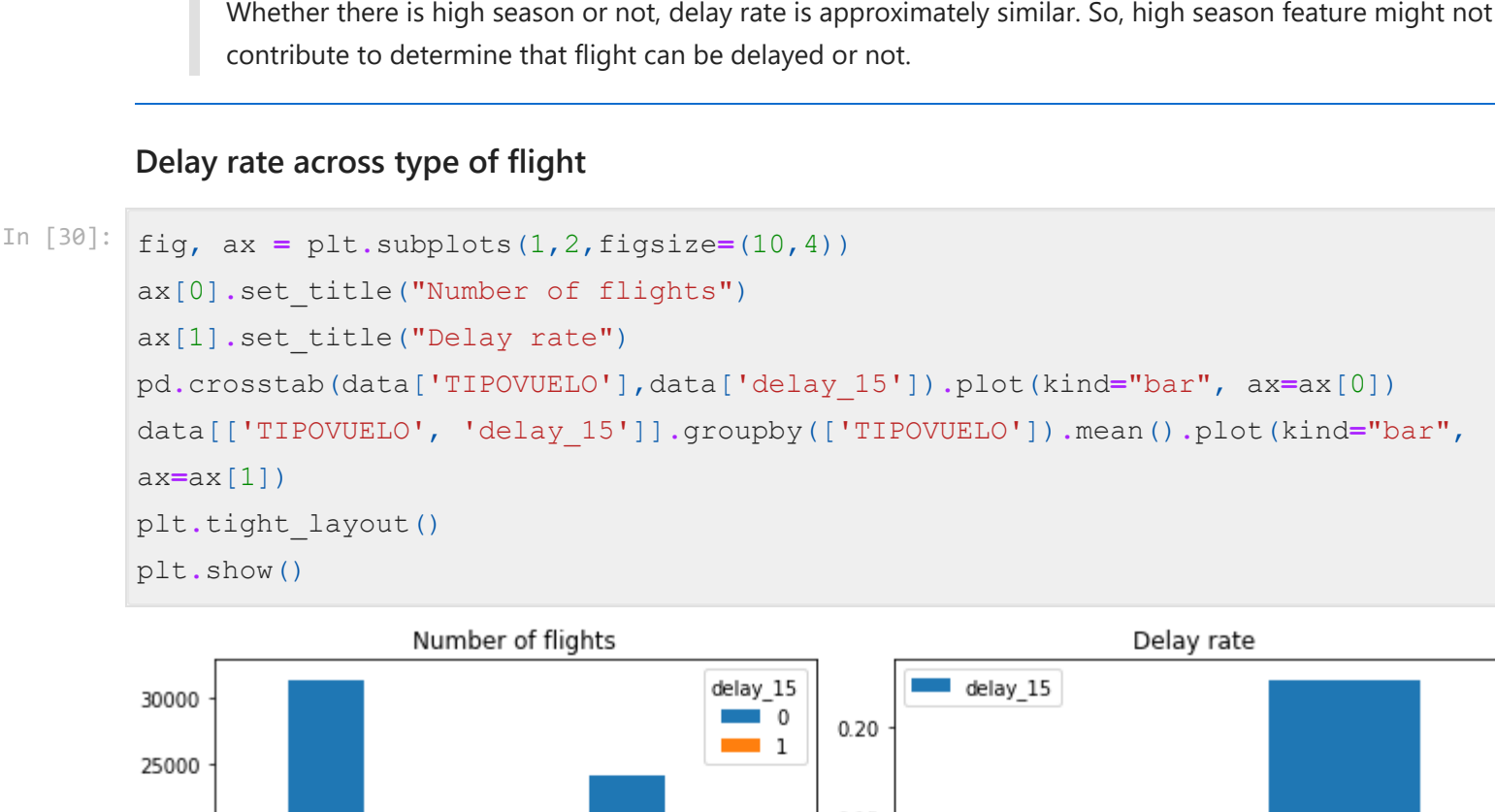
```
In [17]: fig, axes = plt.subplots(4,2,figsize=(20,25))
data["Des-I"].value_counts().plot(kind="bar", xlabel="Des-I", ax=axes[0][0])
data["Emp-I"].value_counts().plot(kind="bar", xlabel="Emp-I", ax=axes[0][1])
data["Des-O"].value_counts().plot(kind="bar", xlabel="Des-O", ax=axes[1][0])
data["Emp-O"].value_counts().plot(kind="bar", xlabel="Emp-O", ax=axes[1][1])
data["DIANOM"].value_counts().plot(kind="bar", xlabel="DIANOM", ax=axes[2][0])
data["SIGLADES"].value_counts().plot(kind="bar", xlabel="OPERA", ax=axes[2][1])
data["SIGLAORI"].value_counts().plot(kind="bar", xlabel="SIGLADES", ax=axes[3][0])
plt.tight_layout()
plt.show()
```




What is the behavior of the delay rate across week of the day?

All days of the week contains almost similar delay rates. So, week of the day feature might not contribute to determine that flight can be delayed or not.

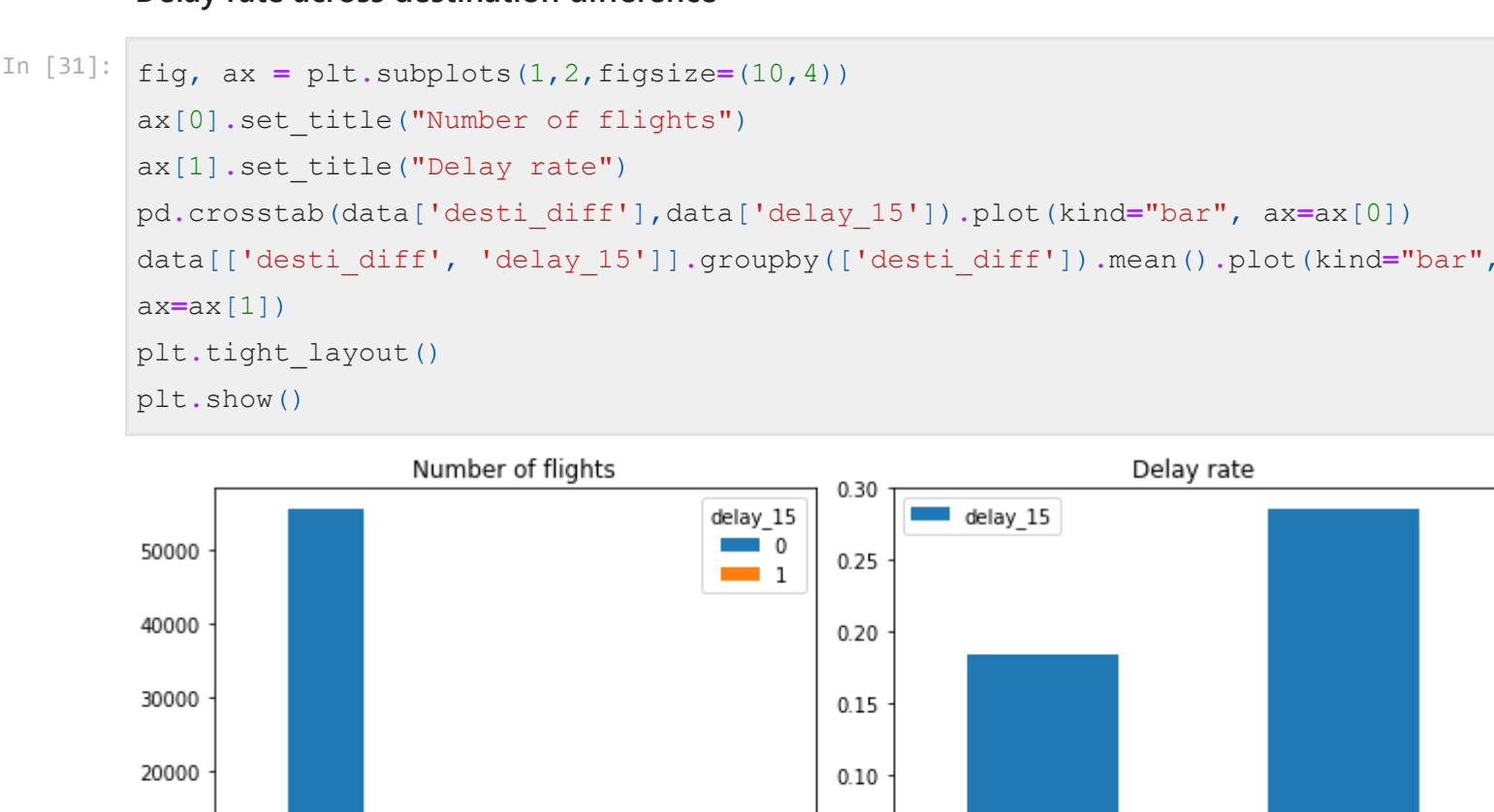
Delay rate across high season



What is the behavior of the delay rate across season?

Whether there is high season or not, delay rate is approximately similar. So, high season feature might not contribute to determine that flight can be delayed or not.

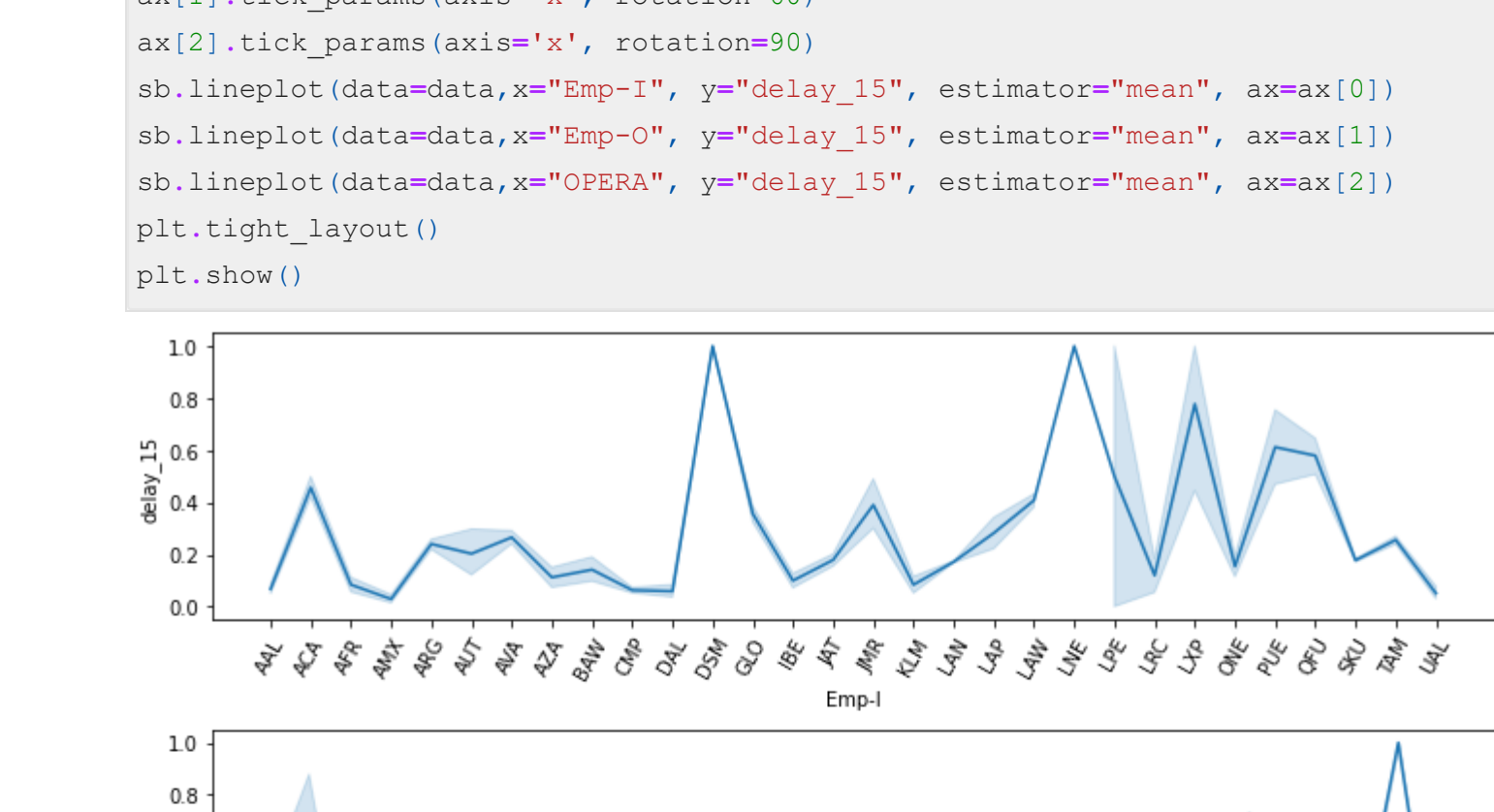
Delay rate across type of flight



What is the behavior of the delay rate across type of flight?

International flights have comparatively high delayrate than national flights. Number of international flights are less but still its delay rate is high. So, type of flight will contribute to determine that flight can be delayed or not.

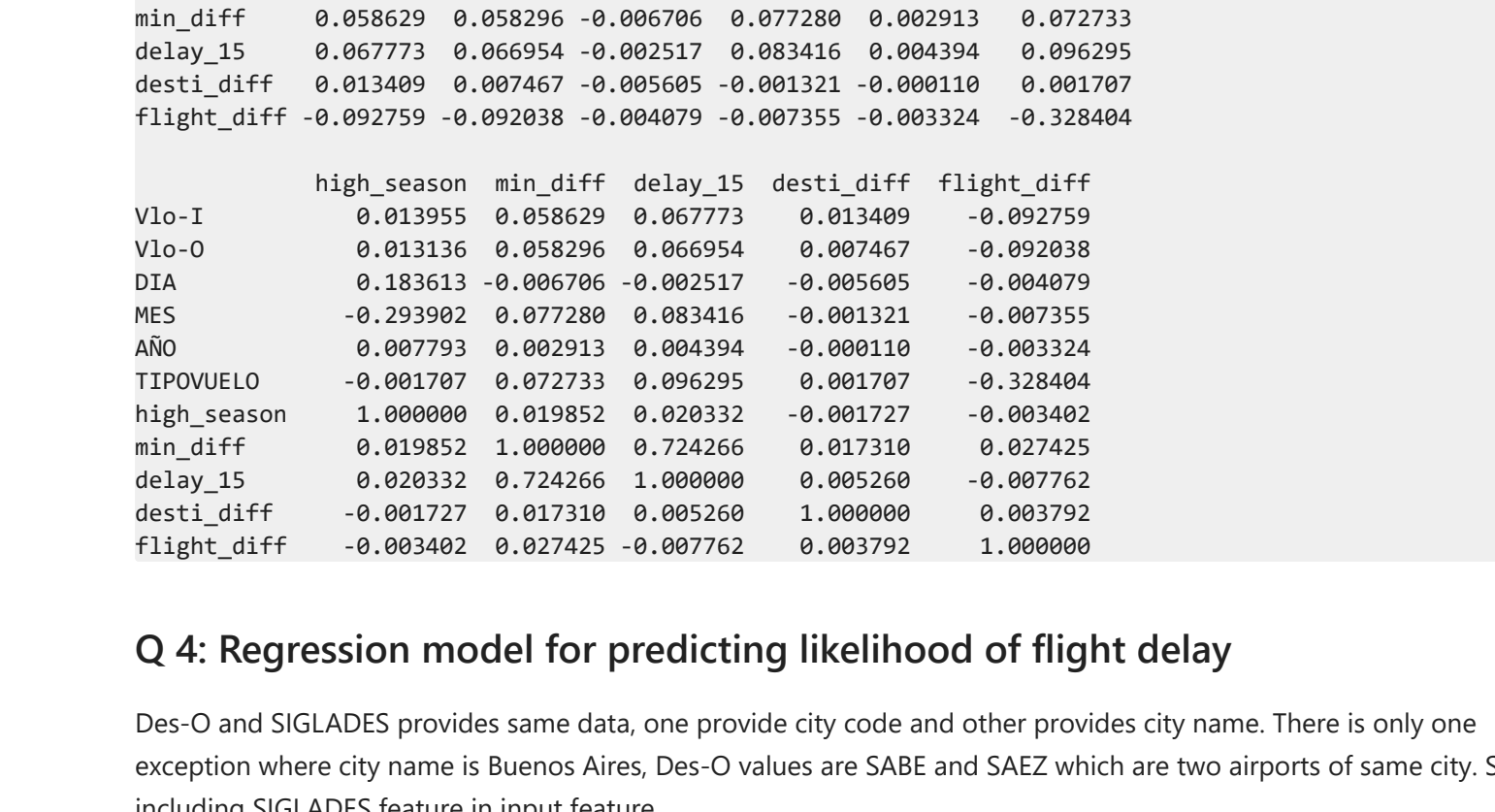
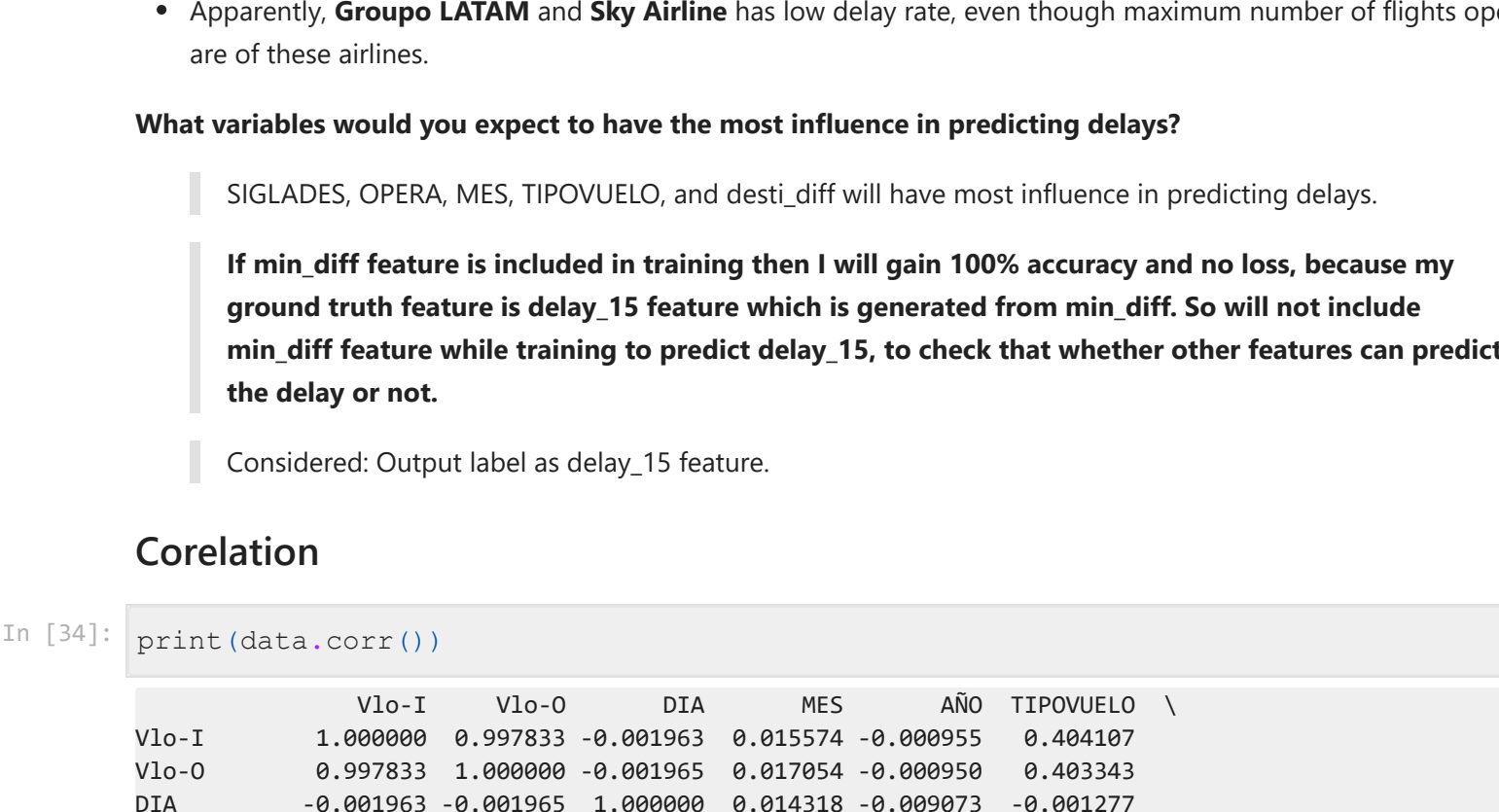
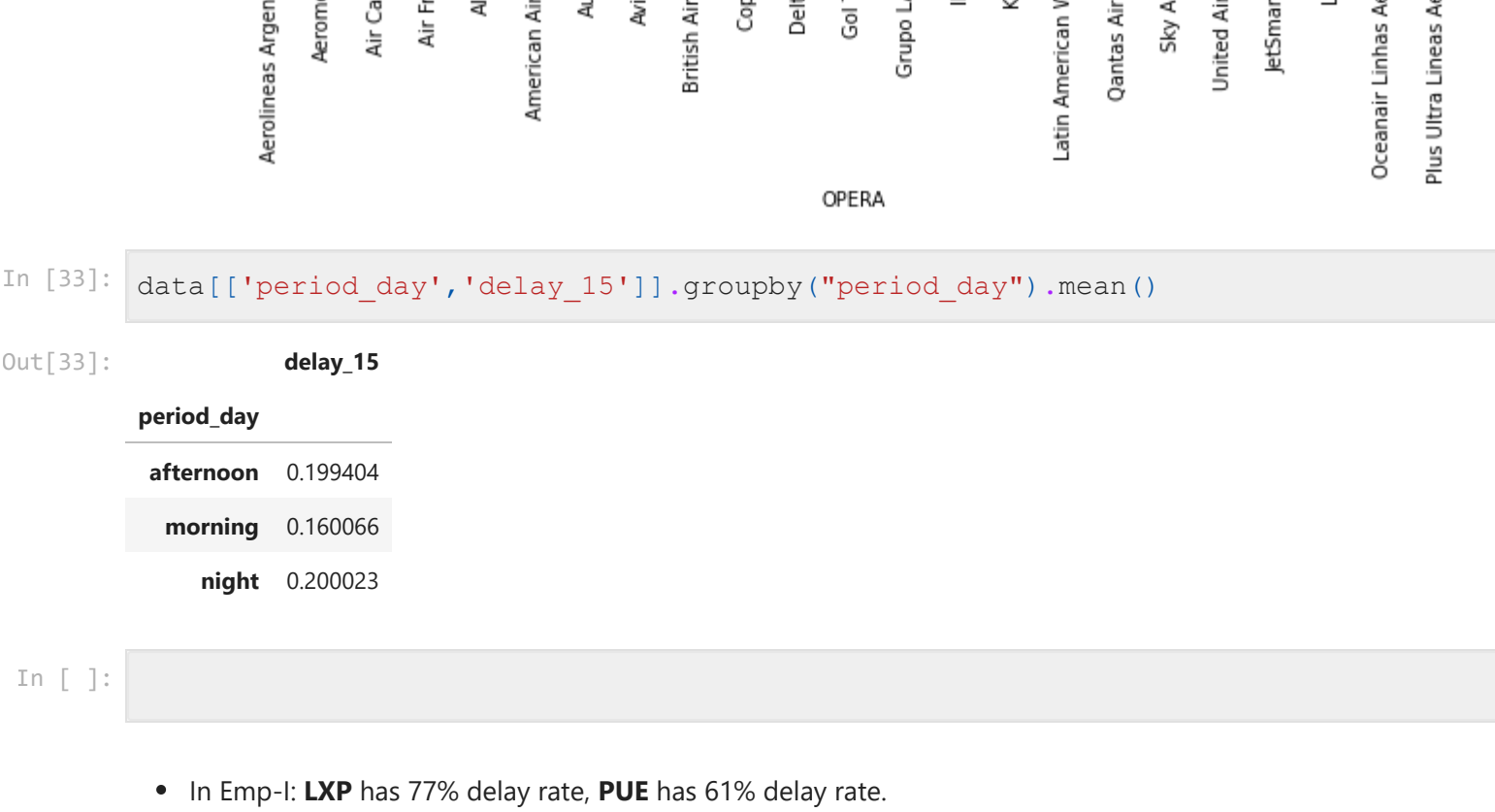
Delay rate across destination difference



What is the behavior of the delay rate across destination difference?

Number of flights whose destination planned is different from destination operated, means (Des-I and Des-O are different) are less, but its delay rate is high. So, destination difference feature will contribute to determine that flight will be delayed or not.

Compare Emp-I, Emp-O, and OPERA contains almost same data



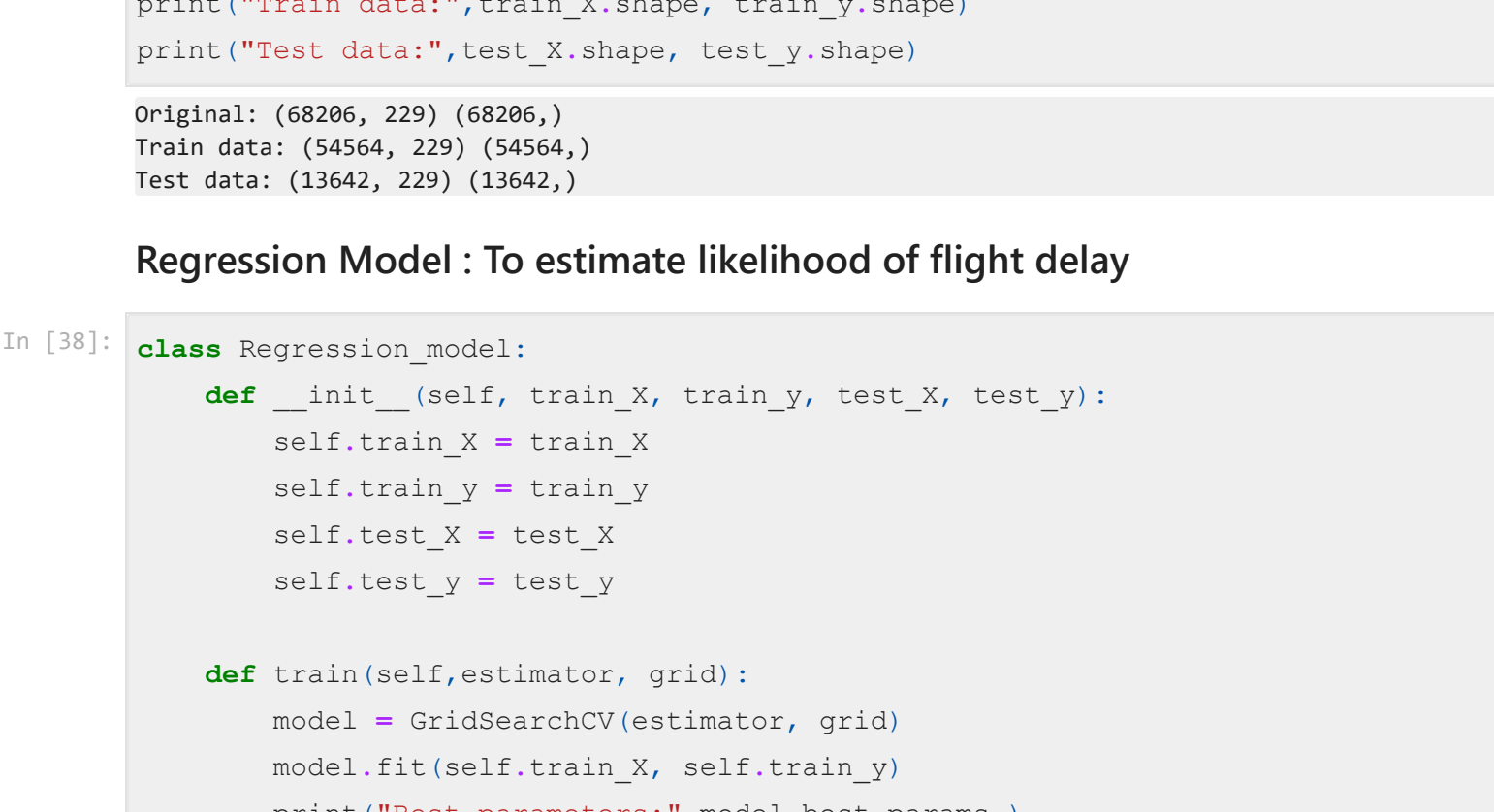
What variables would you expect to have the most influence in predicting delays?

SIGLADES, OPERA, MES, TIPOVUELO, and desti_diff will have most influence in predicting delays.

If min_diff feature is included in training then i will gain 100% accuracy and no loss, because my ground truth feature is delay_15 feature which is generated from min_diff. So will not include min_diff feature while training to predict delay_15, to check that whether other features can predict the delay or not.

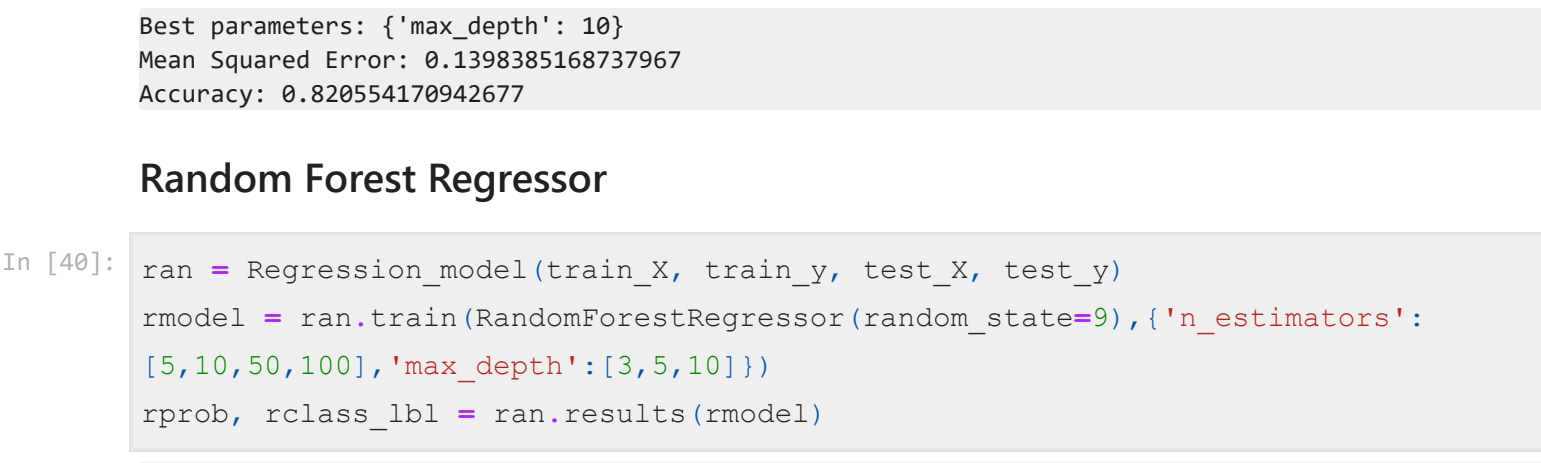
Considered: Output label as delay_15 feature.

Correlation



Q 4: Regression model for predicting likelihood of flight delay

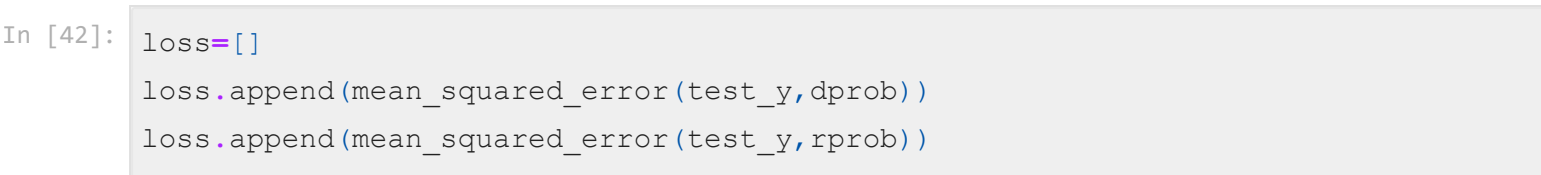
Des-O and SIGLADES provides same data, one provide city code and other provides city name. There is only one exception where city name is Buenos Aires, Des-O values are SAE and SAEZ which are two airports of same city. So only including SIGLADES feature in input feature.



One-Hot Encoding

As majority of data contains categorical features, tree and ensemble algorithms will perform well. The categorical data needs to be converted into numerical data.

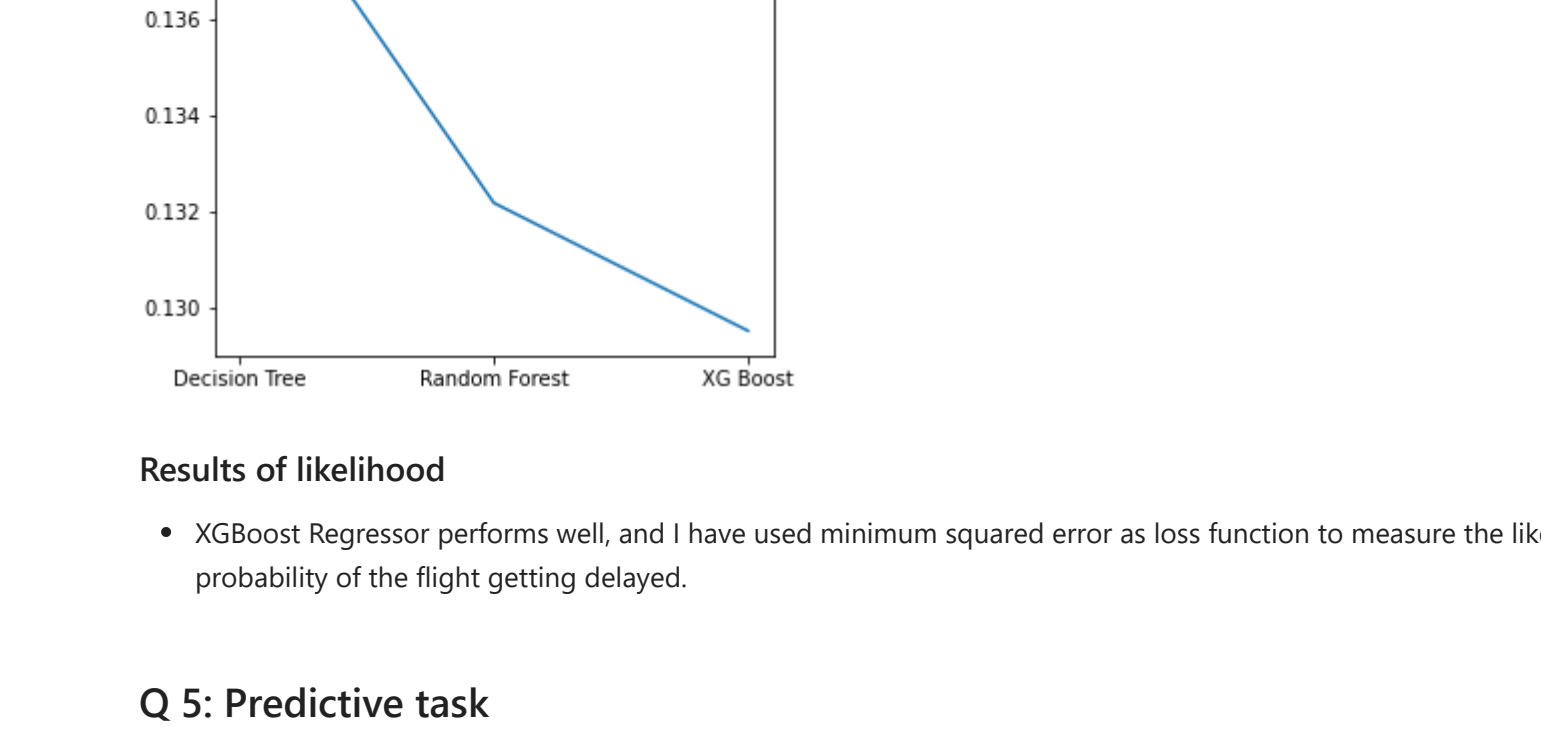
If I assign numbers to categorical, decision tree algorithm will treat them as ordinal values. So, categorical data should be converted to one-hot encoding.



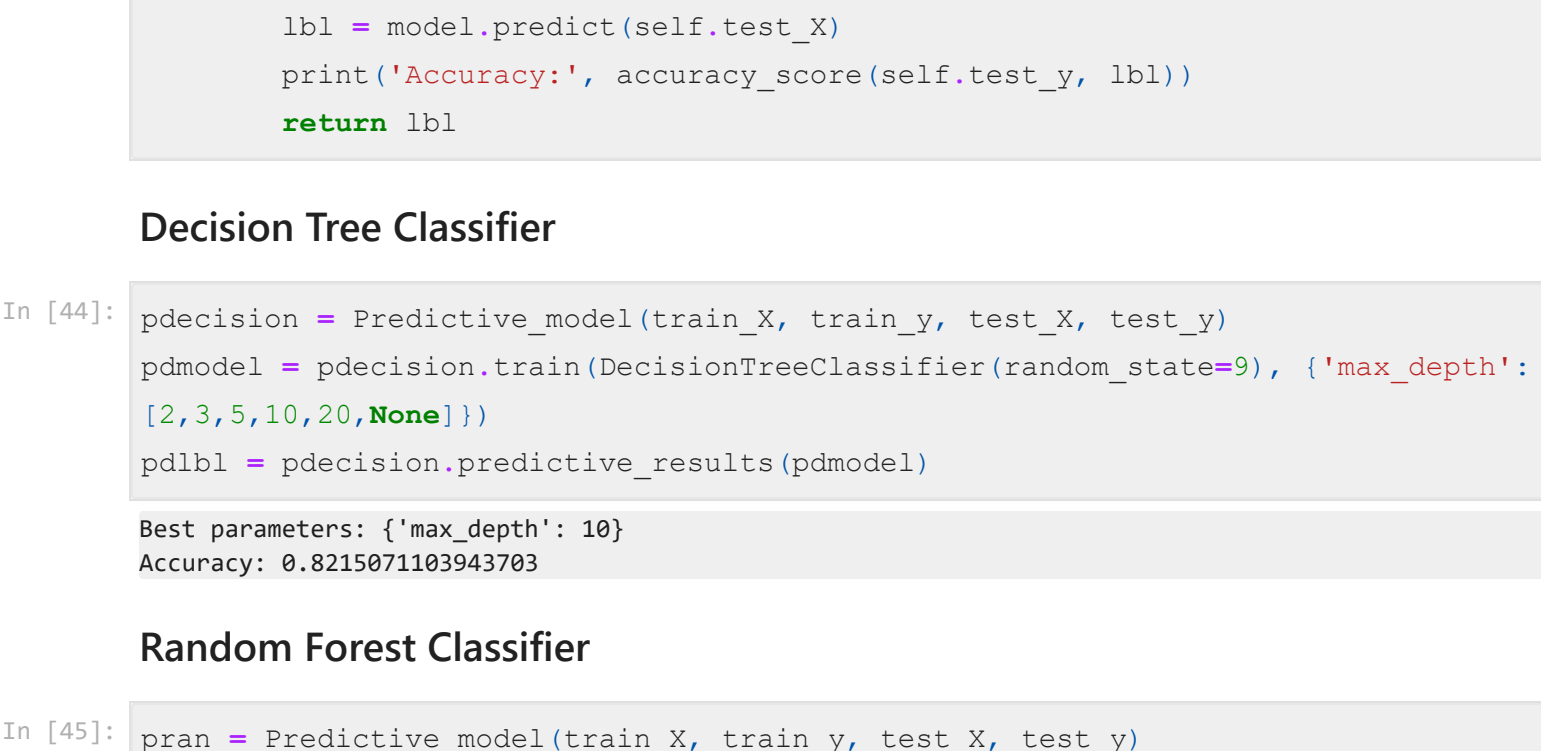
Split dataset

- Training dataset: 80%
- Test dataset: 20%

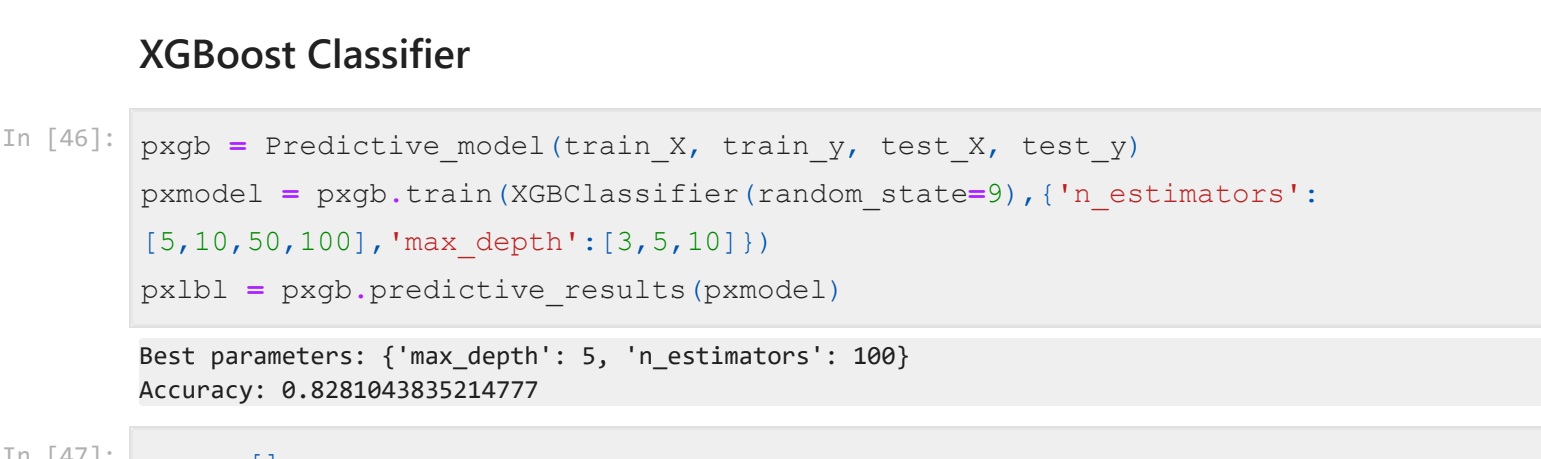
I have used 10-fold cross validation approach. So, it will generate validation sets while training the model.



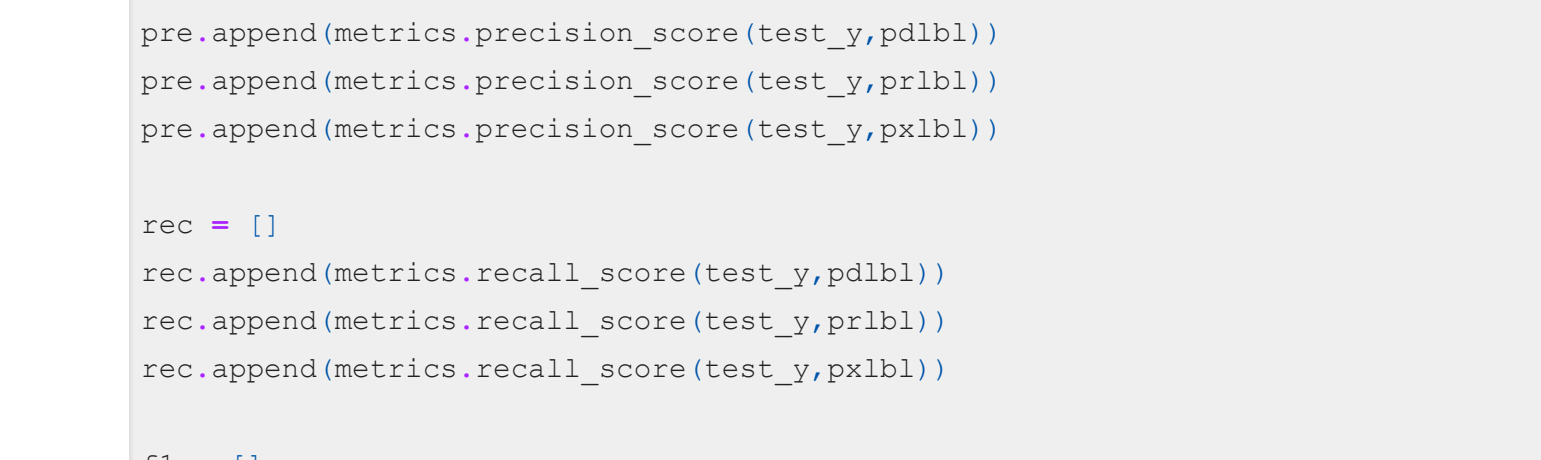
Regression Model : To estimate likelihood of flight delay



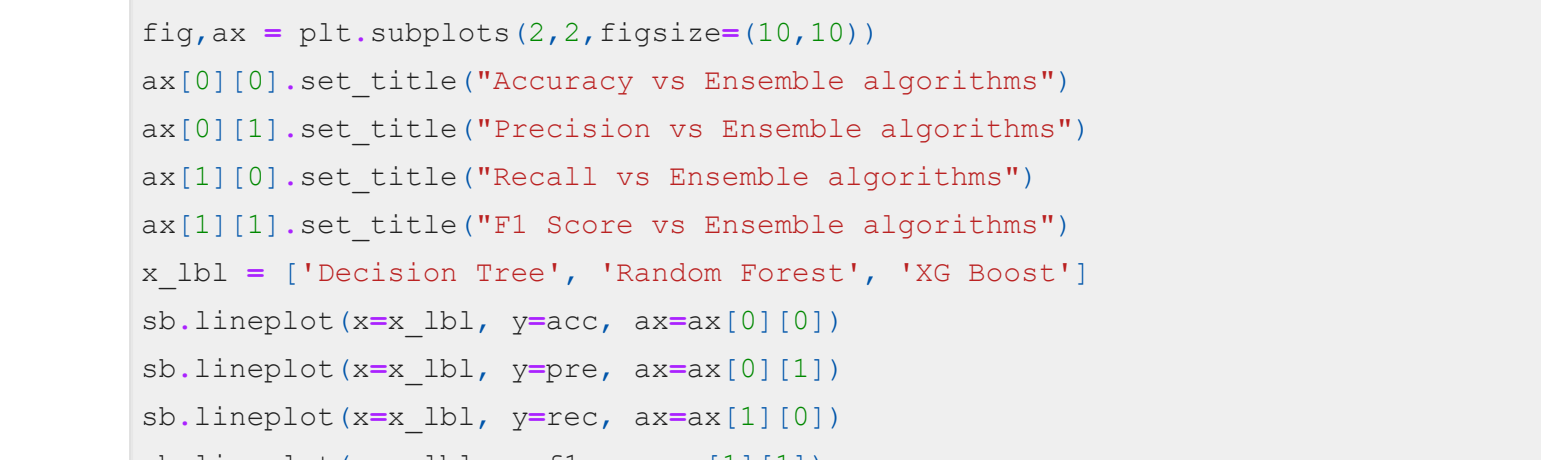
Decision Tree Regressor



Random Forest Regressor

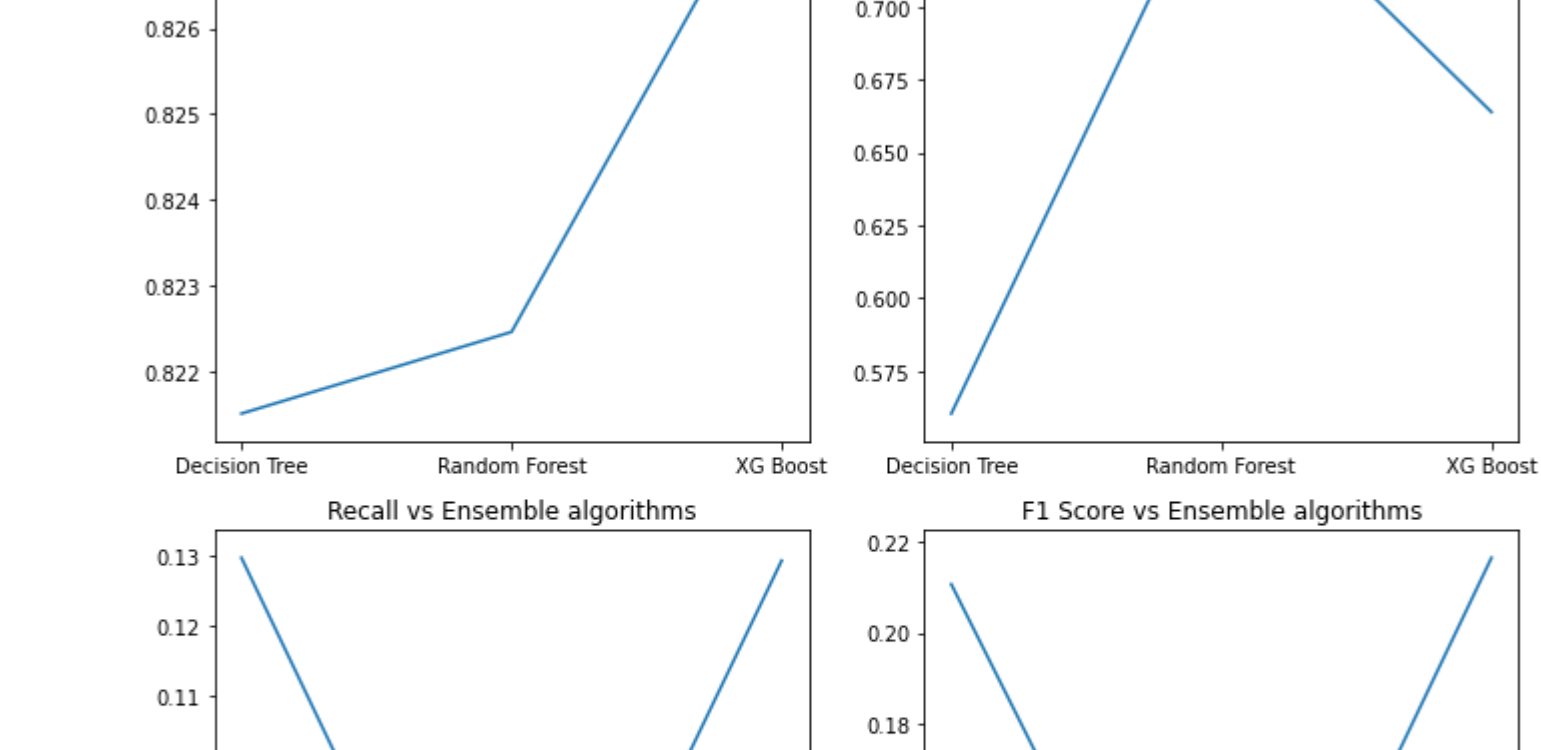


XGBoost Regressor



XG Boost Algorithm is go to algorithm for categorical data.

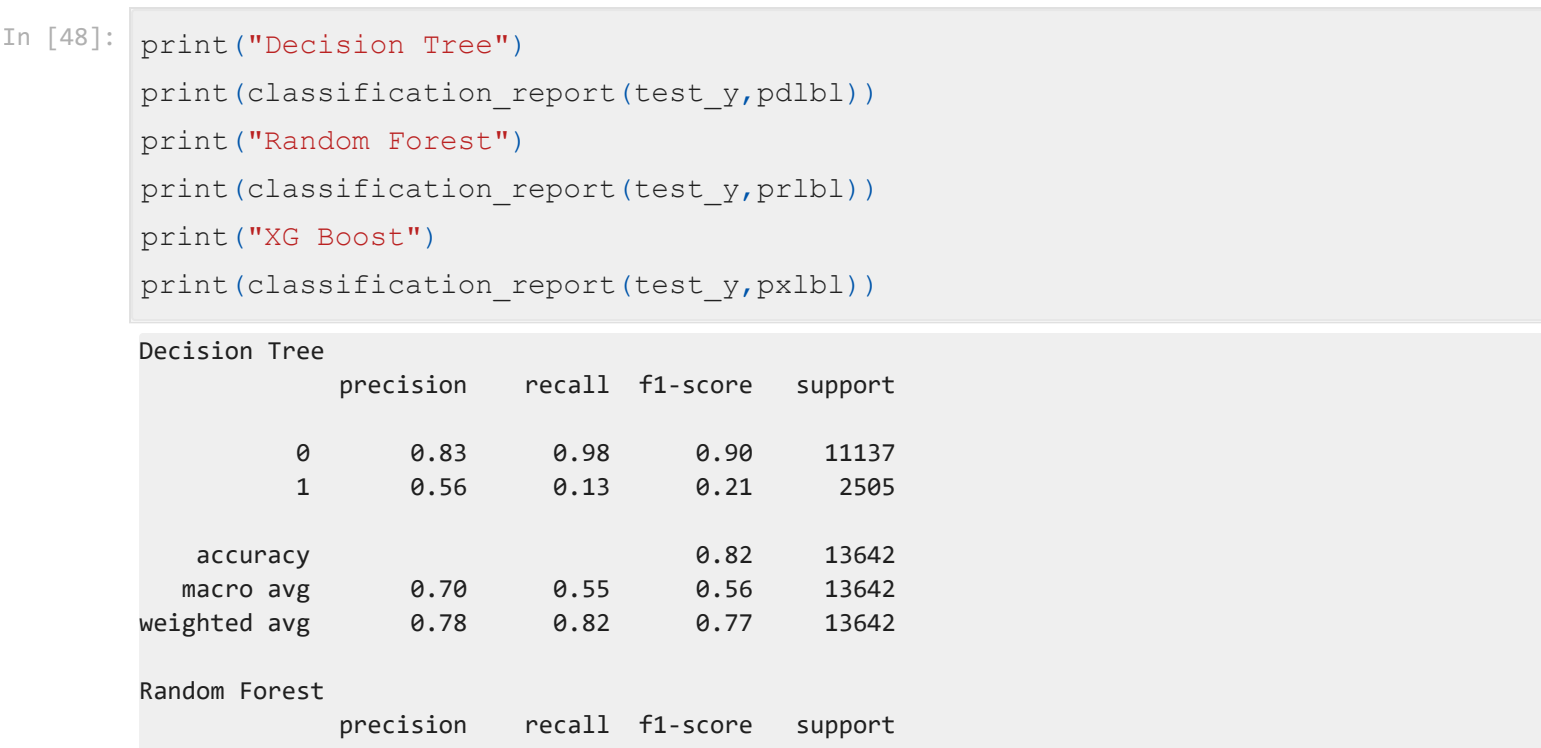
Loss across various ensemble algorithm



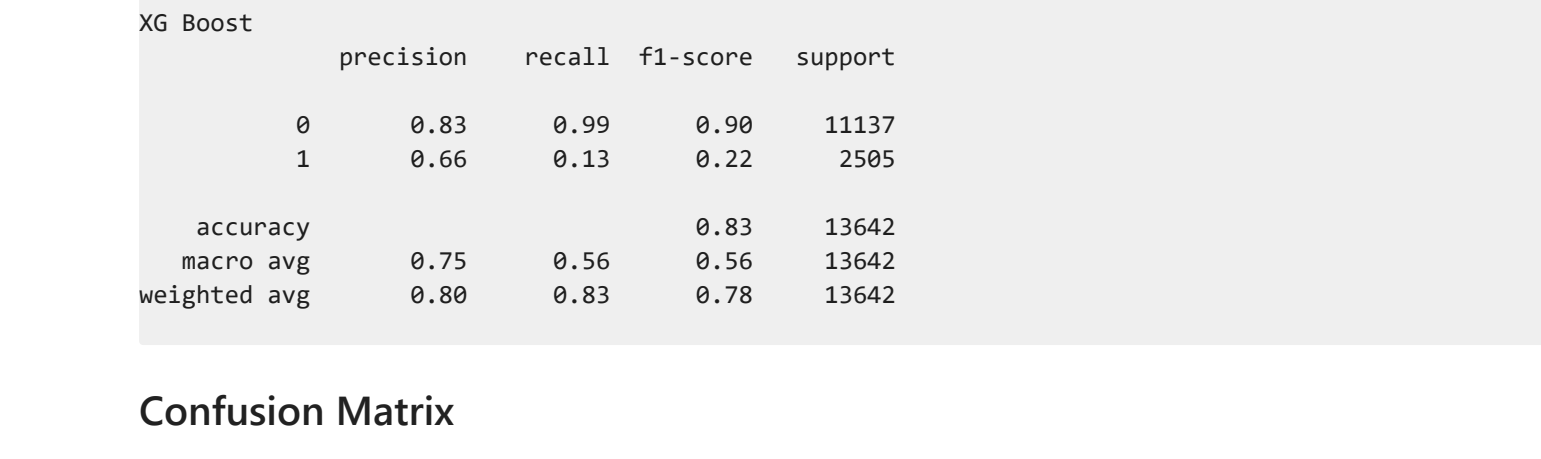
Results of likelihood

- XGBoost Regressor performs well, and I have used minimum squared error as loss function to measure the likelihood probability of the flight getting delayed.

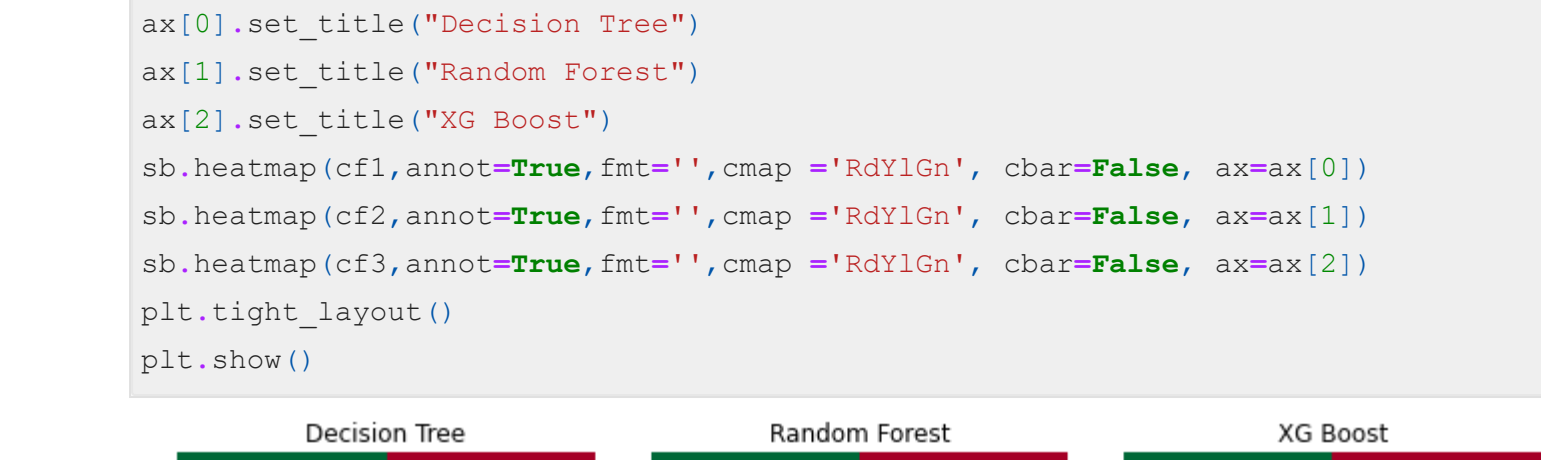
Q 5: Predictive task



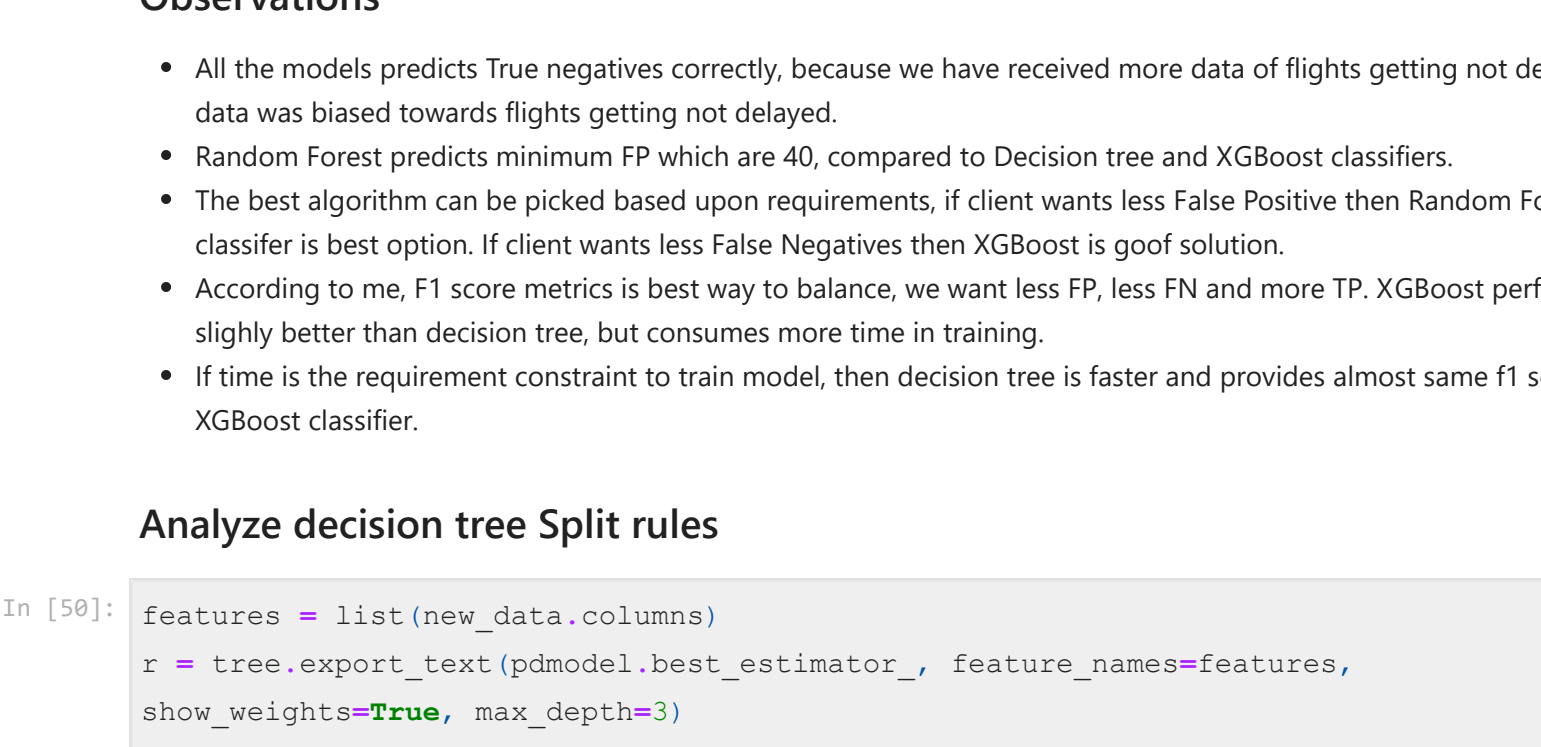
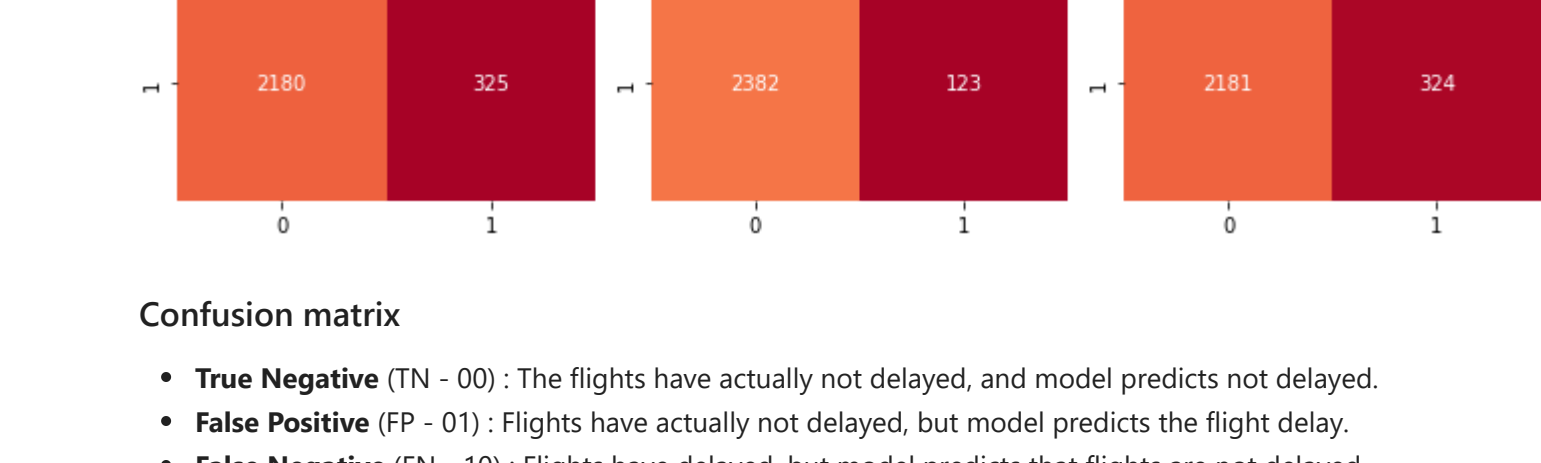
Decision Tree Classifier



Random Forest Classifier

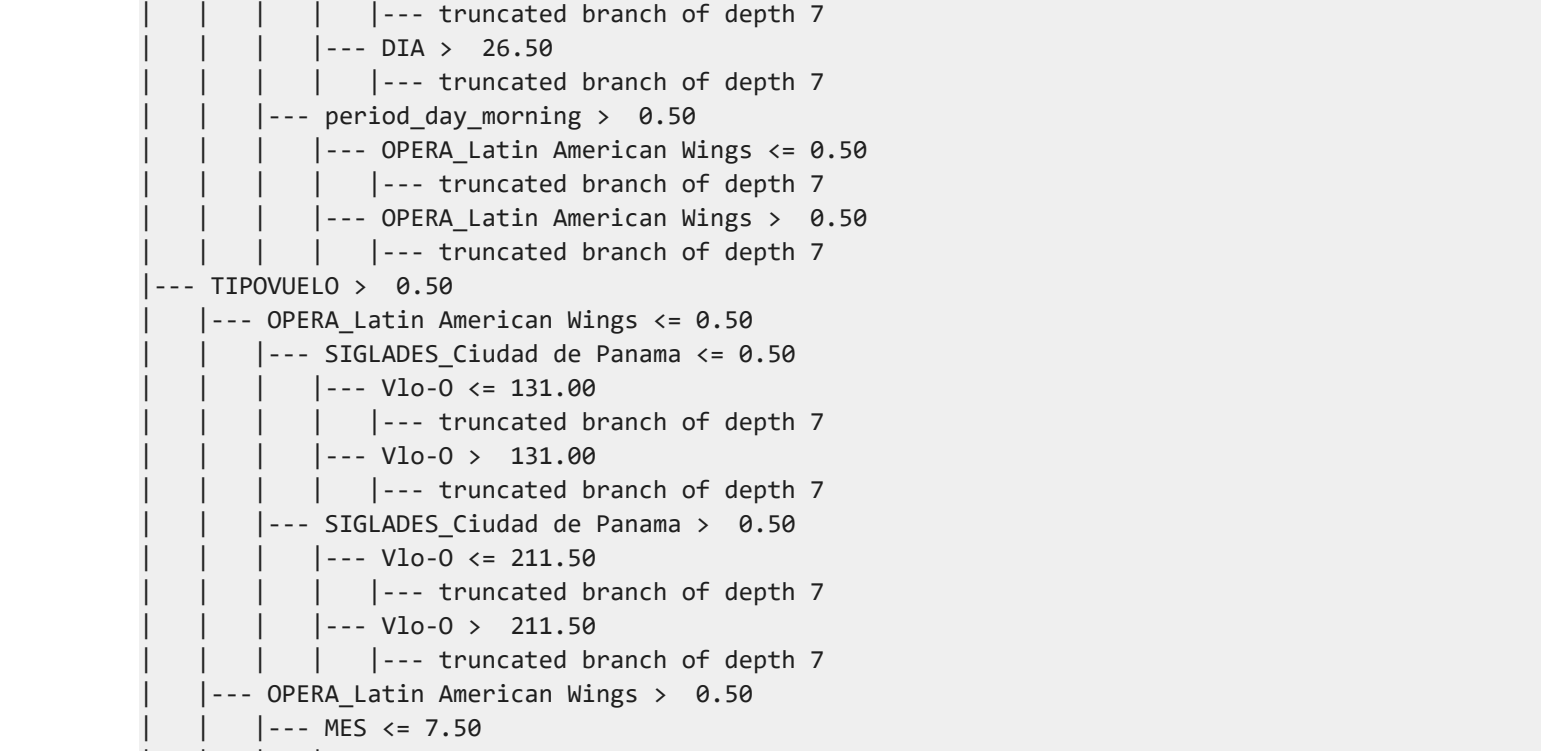


XGBoost Classifier

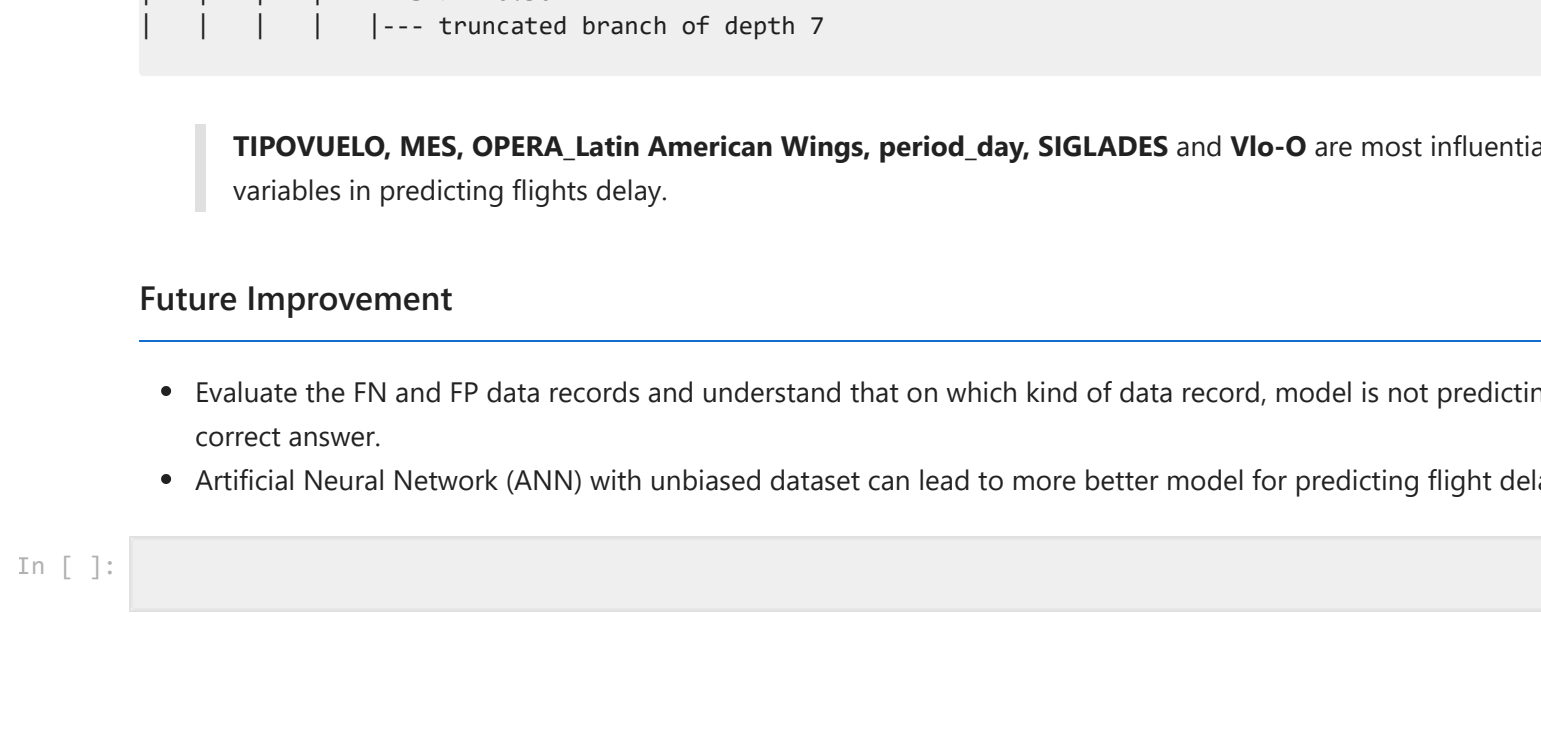


The best performance is of XGBoost Classifier to predict that whether flight is delayed or not. Initially, used accuracy to measure the performance of model.

Classification Report



Confusion Matrix



Confusion matrix

- True Negative (TN - 00)**: The flights have actually not delayed, and model predicts not delayed.
- False Positive (FP - 01)**: Flights have actually not delayed, but model predicts the flight delayed.
- False Negative (FN - 10)**: Flights have delayed, but model predicts that flights are not delayed.
- True Positive (TP - 11)**: Flights have delays, and model predicts the flights are delayed.

Observations

- All the models predicts True negatives correctly, because we have received more data of flights getting not delays. So data was biased towards flights getting not delayed.
- Random Forest predicts minimum FP which are 40, compared to Decision tree and XGBoost classifiers.
- The best algorithm can be picked based upon requirements, if client wants less False Positive then Random Forest classifier is best option. If client wants less False Negatives then XGBoost is good solution.
- According to me, F1 score metrics is best way to balance, we want less FP, less FN and more TP. XGBoost performs slightly better than decision tree, but consumes more time in training.
- If time is the requirement constraint to train model, then decision tree is faster and provides almost same f1 score as XGBoost classifier.

Analyze decision tree Split rules



TIPOVUELO, MES, OPERA, Latin American Wings, period_day, SIGLADES and Vlo-O are most influential variables in predicting flights delay.

Future Improvement

- Evaluate the FN and FP data records and understand that on which kind of data record, model is not predicting correct answer.
- Artificial Neural Network (ANN) with unbiased dataset can lead to more better model for predicting flight delays.

