

CIS - 600 Intro To Machine Learning

Final Project

Anurima Swarup (853430232) Nirmit Patel (992323441) Manoj Kumar Gundapaneni (239719638) Purushothama Rajanna (770079902)

CIS-600 Intro To Machine Learning

Table of Contents

Introduction:	3
Initial Prediction:	3
Final Prediction:	8
Conclusion:	19
Reference:	20

Introduction:

- As part of the final project in our Introduction to Machine Learning class, we were provided with flight data for United Airlines departing from four different airports: Newark (EWR), Washington (IAD), Denver (DEN), and Chicago (ORD), all of which were destined for Syracuse (SYR) Airport.
- Our objective was to predict the flight status for each flight, which we classified
 as "early" if it arrived more than 10 minutes ahead of schedule, "on-time" if it
 arrived within plus or minus 10 minutes of the scheduled time, "late" if it arrived
 more than 10 but less than 30 minutes late, and "severely late" if it was more
 than 30 minutes late.
- The project consisted of two stages: initial prediction and final prediction. To
 accomplish these, we were provided with two different datasets containing flight
 data for the periods of April 12-15 and April 21-24, respectively. We had no
 restrictions on the data we could use to train our model or on the type of machine
 learning model we could employ.

Initial Prediction:

 For our initial prediction we were given a .csv file with the following data: Date, Day, Origin Airport, Flight Number, Arrival Time and a blank column Status. Our goal is to predict the status of these flights.

Approach:

- For our initial prediction we went for a very simple approach of just using historical flight data for our train data without adding any additional data to it.
- Our thought process behind this approach was to find just how accurate of a prediction we can get just by using the historical flight data without introducing any additional factors to our train data.
- For this approach we decided to use a machine learning model to predict the delay based on the historical data, whatever delay we get out of it we then categorized them into late, severely late, on-time and early.

Data Gathering and Analysis:

- We obtained information from the Bureau of Transportation Statistics (BTS), a federal body that gathers and disseminates statistics on airline performance and operations.
- We began by obtaining data for departures from four places (Washington, Denver, Newark, and Chicago) and arrival statistics of Syracuse.
- Aside from the data suggested by the professor, we gathered a few more metrics (which are useful in making accurate predictions).
- We integrated arrivals and departures based on date, origin airport, and flight number after cleaning the data. This process assisted us in removing any inconsistencies in the data.
- To train the model, we hot encoded the origin airports and scaled the data.
- We used a linear regression model for training and got a score of 0.89.
- Test predictions are moderately off. Mean absolute error is 13.9 and Error ratio is 7.22
- Finally, we exported the predictions into a CSV file.

Output:

 For our train data we first read arrival and departure csv data into our dataframe and then merged it together to get our final flightData

arrivalData - pd.read_csv('/Users/puru/Desktop/Syracuse Courses/Sem 3/Intro to ML/Project/Dataset/Flight_data/Detailed_Statistics_Arrivals -final.csv',parse_dates-['Date (MW/DD/YYYY)', 'Scheduled Ard departureData - pd.read_csv('/Users/puru/Desktop/Syracuse Courses/Sem 3/Intro to ML/Project/Dataset/Flight_data/final departures-details.csv',parse_dates-['Date (MW/DD/YYYY)', 'Scheduled departure to the control of the control

flightData = pd.merge(arrivalData, departureData, on = ['Date (NW/DD/YYYY)','Flight Number','Tail Number','Carrier Code'])

flightData.drop(columns=['Tail Number','Carrier Code','Source Airport','Destination Airport'],inplace=True)
flightData.dtypes
flightData

Python Python

	Date (MM/DD/YYYY)	Flight Number	Origin Airport	Scheduled Arrival Time	Actual Arrival Time	Arrival Delay (Minutes)	Scheduled departure time	Actual departure time	Scheduled elapsed time (Minutes)	Actual elapsed time (Minutes)	Departure delay (Minutes)	Wheels- off time	Taxi-Out time (Minutes)	Delay Carrier (Minutes)	Delay Weather (Minutes)	Delay National Aviation System (Minutes)	Delay Security (Minutes)	Delay Late Aircraft Arrival (Minutes)
0	2022-01-10	604	DEN	2023-04- 11 15:09:00	2023- 04-11 15:04:00		2023-04- 11 09:49:00	2023-04- 11 09:46:00	200	198		2023- 04-11 10:03:00						
1	2022-01-10	1917	ORD	2023-04- 11 16:57:00	2023- 04-11 16:54:00		2023-04- 11 14:10:00	2023-04- 11 14:06:00		108		2023- 04-11 14:32:00						
2	2022-01-10	1998	ORD	2023-04- 11 21:18:00	2023- 04-11 21:01:00		2023-04- 11 18:21:00	2023-04- 11 18:15:00		106		2023- 04-11 18:35:00						
3	2022-01-10	2198	IAD	2023-04- 11 23:31:00	2023- 04-11 23:52:00		2023-04- 11 22:20:00	2023-04- 11 22:47:00				2023- 04-11 23:01:00						
4	2022-02-10	604	DEN	2023-04- 11 15:09:00	2023- 04-11 14:55:00		2023-04- 11 09:49:00	2023-04- 11 09:44:00	200	191		2023- 04-11 09:58:00						
266	2022-12-30	1998	ORD	2023-04- 11 21:07:00	2023- 04-11 20:56:00		2023-04- 11 18:14:00	2023-04- 11 18:11:00		105		2023- 04-11 18:33:00						
267	2022-12-30	2488	EWR	2023-04- 11 23:14:00	2023- 04-11 23:07:00		2023-04- 11 21:59:00	2023-04- 11 22:17:00				2023- 04-11 22:31:00						
268	2022-12-31	604	DEN	2023-04- 11 14:58:00	2023- 04-11 14:46:00		2023-04- 11 09:45:00	2023-04- 11 09:48:00	193			2023- 04-11 10:00:00						
269	2022-12-31	1998	ORD	2023-04- 11 21:08:00	2023- 04-11 20:44:00		2023-04- 11 18:15:00	2023-04- 11 18:06:00		98		2023- 04-11 18:22:00						
270	2022-12-31	2488	EWR	2023-04- 11 23:14:00	2023- 04-11 00:46:00		2023-04- 11 21:59:00	2023-04- 11 23:33:00			94	2023- 04-11 00:03:00						
271 ro	ws × 18 columns													Do you mind	d taking a quick	feedback surv		×

 After removing NaN, we removed unnecessary columns, changed the data type of a few columns into datetime and perform a one hot encoding on "Origin Airport" column

```
flightData['Scheduled Arrival Time(hour)'] = flightData['Scheduled Arrival Time'].dt.hour
flightData['Scheduled drival Time(minutes)'] = flightData['Scheduled departure time'].dt.hour
flightData['Scheduled departure time(minutes)'] = flightData['Scheduled departure time'].dt.hour
flightData['Actual Arrival Time(hour)'] = flightData['Scheduled departure time'].dt.minute

flightData['Actual Arrival Time(minutes)'] = flightData['Actual departure time'].dt.minute

flightData('Actual Arrival Time(minutes)'] = flightData['Actual departure time'].dt.minute

flightData('Actual departure time(minutes)'] = flightData['Actual departure time'].dt.minute

flightData('Mheels-off time(hour)'] = flightData['Actual departure time'].dt.minute

flightData('Mheels-off time(minutes)'] = flightData['Mheels-off time'].dt.minute

flightData('Mheels-off time(minutes)') = flightData('Mheels-off time').dt.minute

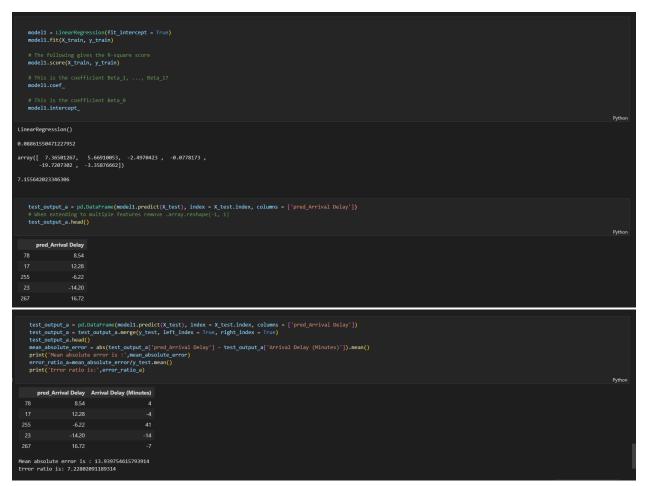
flightData('Mheels-off time').dt.minute
```

 After that we performed a train test split and used standard scalar on the data so we can fit it into our model.

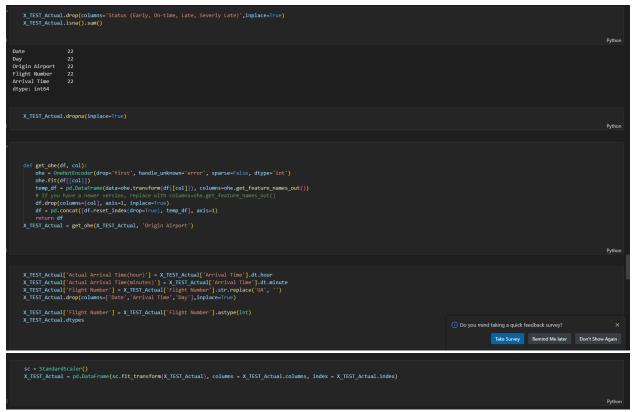
```
X_train, X_test, y_train, y_test = train_test_split(Final_flightData.drop(columns = ['Arrival Delay (Minutes)']), Final_flightData['Arrival Delay (Minutes)'], test_size=0.05, random_state=35)
X_train
X_test
y_train
y_test
Python
```

```
if True:
    from sklearn.preprocessing import StandardScaler
    sc = StandardScaler()
    X_train = pd.DataFrame(sc.fit_transform(X_train), columns = X_train.columns, index = X_train.index)
    X_test = pd.DataFrame(sc.transform(X_test), columns = X_test.columns, index = X_test.index)
    Fython
Python
```

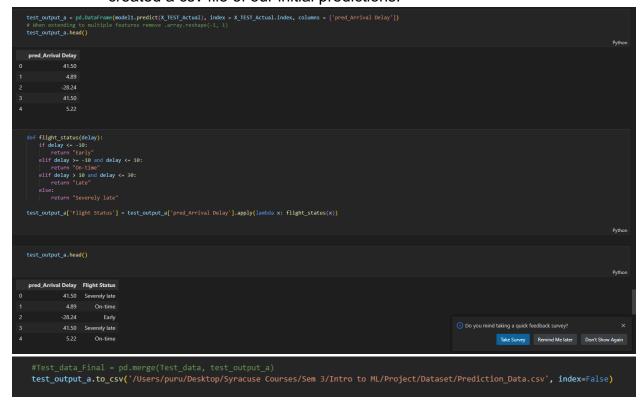
 We use Linear Regression, first we used test data to find out the score of our model.



 After checking the score of our model we moved on with our actual test data. We added that data into another dataframe and performed all the necessary scaling operations and one hot encoding to get the data in the same format as our test data that we used for our model.



• Finally, we used our test data to predict our initial prediction, converted that prediction into the status that was asked from the question and created a csv file of our initial predictions.



 According to the ground truth our model predicted 8 out of 32 correctly with the accuracy of 25%

Final Prediction:

After our initial prediction we can see that just by using the historical data in this
particular situation will not give us very accurate predictions. So, to increase the
accuracy of our prediction we have to incorporate multiple datasets.

Approach:

- For our final prediction(21-24) we have considered both flight data and weather data. Because weather is one of the most common causes of flight delays, taking into account both flight data and weather data is a wise approach. This may improve our understanding of the influence of weather conditions on flight schedules and generate more accurate predictions by including weather data.
- We planned to consider the weather on both the arrival and departure city.
 It is critical to collect information about weather conditions in both the
 departure and arrival cities when utilizing weather data for prediction. This
 data assisted us in identifying possible difficulties that may develop during
 the journey, such as severe rain or thunderstorms, which may cause flight
 delays or cancellations.
- Our thought about using weather data is getting the reasons behind the
 delay and interpolate with the timings. After gathering weather data, we
 evaluated it to see how it will affect the flight schedule. For example, if
 thunderstorms are expected at the destination airport during the
 scheduled arrival time, the flight may be delayed or even canceled. If the
 weather is clear and there are no serious complications, the flight may
 arrive on time, if not early.
- By merging weather and flight data, we can create a more accurate forecast model that accounts for the specific aspects that might affect the flight schedule. For example, based on current weather conditions, we may utilize meteorological data to interpolate timings and predict the most likely time of arrival or departure.

Data Gathering and Analysis:

- We obtained information from the Bureau of Transportation Statistics
 (BTS), a federal body that gathers and disseminates statistics on airline
 performance and operations. We acquired insights into different elements
 of airline operations, such as estimated flight delays, cancellations, ontime performance by using data from the BTS. This information has been
 used to discover trends, patterns, and concerns affecting airline
 performance, as well as to construct predictive models for anticipating
 future performance.
- For weather data we used weather bit api to get weather from the past and also forecast data. This information can be used to create predictive models that foresee probable weather-related concerns, allowing us to take proactive efforts to avoid delays and enhance overall flight operations.
- We have used date as a key to merge flight and weather data to create a complete dataset for predicting aircraft delays. Because flight schedules and weather conditions are both time-dependent, using date as a key to merge these datasets is our analysis.
- While using weather data we considered metrics like temp, pressure, wind speed, wind direction, precipitation and others which are key factors.
- Later the extra metrics have been divided into four categories : Cloudy, rain, snow and clear. Which gives us little clarity in data.
- Weather conditions such as cloudy, rain, snow, and clear may all have an
 affect on flight operations in different ways. Cloudy circumstances, for
 example, might produce reduced visibility, resulting in aircraft delays or
 cancellations. Rain and snow can have an effect on airport operations,
 such as runway conditions and aircraft de-icing needs. Clear weather, on
 the other hand, may have no effect on aircraft operations.
- We have used the same approach as initial prediction while training the model but this time we have trained multiple models to get the best predictions.
- LinearRegression: 0.12
 - Mean absolute error is 31.06
 - o Error ratio is 2.14
- GradientBoostingRegressor
 - Mean absolute error is 32.00
 - o Error ratio is 2.20
- Lasso Regression
 - o Score: 0.11

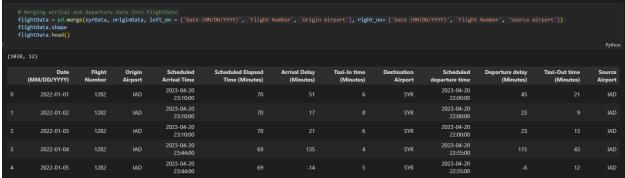
- Mean absolute error is 30.89
- DecisionTreeRegressor
 - We have used cross validation technique to get best prediction but MAE is same for both the results
 - Mean absolute error is 28.2
- RandomForestRegressor
 - o Mean absolute error is 30.89
 - 0 31.44

Output:

 We started by adding the flight data for our source airport and destination airport to a dataframe. In this data we had some unnecessary columns, we removed them, dropped NaN rows and merged two datasets on date, flight number and origin airport so we can get our final flight data.

CIS-600 Intro To Machine Learning

```
print('Syracuse flight NaN :')
   syrData.isna().sum()
   print('origin flight NaN :')
   originData.isna().sum()
Syracuse flight NaN:
Date (MM/DD/YYYY)
Flight Number
                                    0
Origin Airport
                                    0
Scheduled Arrival Time
Scheduled Elapsed Time (Minutes)
                                    0
Arrival Delay (Minutes)
Taxi-In time (Minutes)
dtype: int64
origin flight NaN :
Date (MM/DD/YYYY)
Flight Number
                             ø
Destination Airport
Scheduled departure time
Departure delay (Minutes)
                             0
Taxi-Out time (Minutes)
                             a
Source Airport
dtype: int64
```



After getting our flight we started working on getting weather data. We read weather data for Syracuse, Chicago, Newark, Washington D.C. and Denver into dataframes, dropped all unnecessary columns and removed NaNs to get our final weather data.

```
# Reading Weather data

syrWeather = pd.read_csv("C:/Users/nirmi/Desktop/SU_SEM3/IntroToML/Project/Data sets/Weather_2022_final/Weather_2022/syracuse_weather_2022_csv",parse_dates=['datetime', 'timestamp_utc', 'timestamp_tc', 'timestamp_utc', '
```

CIS-600 Intro To Machine Learning

```
# checking for NaN
syrWeather.isna().sum()
chiWeather.isna().sum()
wasWeather.isna().sum()
newWeather.isna().sum()
denWeather.isna().sum()

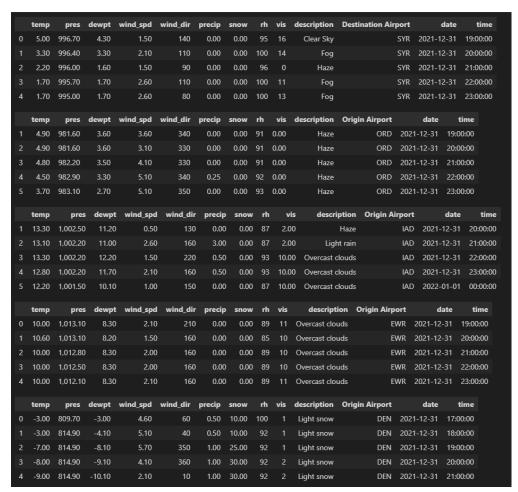
chiWeather.dropna(inplace= True)
wasWeather.dropna(inplace= True)
```

```
denkeather['Origin Airport'] = 'DEN'
chiweather['Origin Airport'] = 'ORD'
waskeather['Origin Airport'] = 'IAD'
newkeather['Origin Airport'] = 'IAD'
newkeather['Destination Airport'] = 'SYR'

syrWeather['Destination Airport'] = 'SYR'

syrWeather['date'] = syrWeather['timestamp_local'].dt.date
syrWeather['time'] = chiweather['timestamp_local'].dt.date
chiweather['date'] = chiweather['timestamp_local'].dt.date
chiweather['date'] = wasWeather['timestamp_local'].dt.date
wasWeather['date'] = wasWeather['timestamp_local'].dt.date
wasWeather['date'] = newWeather['timestamp_local'].dt.date
newWeather['time'] = newWeather['timestamp_local'].dt.date
denkweather['time'] = newWeather['timestamp_local'].dt.date
denkweather['date'] = denkweather['timestamp_local'].dt.date
denkweather['date'] = denkweather['timestamp_local'].dt.time

syrWeather.drop(columns=['datetime', 'timestamp_utc', 'uv', 'solar_rad', 'clouds', 'pod', 'timestamp_local'],inplace = True)
chiweather.drop(columns=['datetime', 'timestamp_utc', 'uv', 'solar_rad', 'clouds', 'pod', 'timestamp_local'],inplace = True)
mesWeather.drop(columns=['datetime', 'timestamp_utc', 'uv', 'solar_rad', 'clouds', 'pod', 'timestamp_local'],inplace = True)
denkweather.drop(columns=['datetime', 'timestamp_utc', 'uv', 'solar_rad', 'clouds', 'pod', 'timestamp_local'],inplace = True)
syrWeather.head()
chiweather.head()
denkweather.head()
denkweather.head()
denkweather.head()
denkweather.head()
denkweather.head()
```



 Finally to get our final data we merged all the data using different keys and performed some cleaning to get our final data.

```
# Merging the weather data with flight data
flightData['Scheduled Arrival Time'] = flightData['Scheduled Arrival Time'].dt.round('H')
flightData['Scheduled departure time'] = flightData['Scheduled departure time'].dt.round('H')
flightData['Scheduled Arrival Time'] = flightData['Scheduled Arrival Time'].dt.time
flightData['Scheduled departure time'] = flightData['Scheduled departure time'].dt.time
flightData.drop(columns = ['Source Airport'], inplace = True)
flightData
flightData.dtypes
```

```
syrWeather = syrWeather.rename(columns=destinationColumnNames)
     denWeather = denWeather.rename(columns=sourceColumnNames)
     wasWeather = wasWeather.rename(columns=sourceColumnNames)
     chiWeather = chiWeather.rename(columns=sourceColumnNames)
     newWeather = newWeather.rename(columns=sourceColumnNames)
    syrWeather['date_d'] = pd.to_datetime(syrWeather['date_d'])
denWeather['date_s'] = pd.to_datetime(denWeather['date_s'])
wasWeather['date_s'] = pd.to_datetime(wasWeather['date_s'])
chiWeather['date_s'] = pd.to_datetime(chiWeather['date_s'])
newWeather['date_s'] = pd.to_datetime(newWeather['date_s'])
     syrWeather.head()
     chiWeather.head()
     wasWeather.head()
     newWeather.head()
     originWeatherData = pd.concat([denWeather, wasWeather, chiWeather, newWeather], ignore_index=True)
     originWeatherData.shape
     originWeatherData
finalData = pd.merge(flightData, syrWeather, right_on=['Destination Airport', 'date_d', 'time_d'], left_on=['Destination Airport', 'Date (WM/DD/YYYY)', 'Scheduled Arrival Time'])
finalData = pd.merge(finalData, originWeatherData, right_on=['Origin Airport', 'date_s', 'time_s'], left_on=['Origin Airport', 'Date (WM/DD/YYYY)', 'Scheduled departure time'])
   finalData.drop(columns = ['Date (MM/DD/YYYY)', 'Destination Airport', 'date_d', 'time_d', 'date_s', 'time_s'], inplace = True)
Arrival T
Delay
Minutes) (Min
                                                                                                    ... temp_s
                                                                                                                pres_s dewpt_s wind_spd_s wind_dir_s precip_s snow_s rh_s vis_s description_s
                                                                                                                                                                    63 10.00
                        23:00:00
                                                                22:00:00
                                                                                             -11.10 ... -1.60 1,016.10
                                                                                                                                     2.80
                        00:00:00
                                                                23:00:00
                                                                                             -11.40 ...
                                                                                                        -1.20 1.016.40
                                                                                                                          -4.20
                                                                                                                                                                    80 10.00
                                                                23:00:00
                                                                                                          5.10 998.10
                                                                                                                          1.90
                                                                                                                                                                     80 10.00
                                                                                                                                                       0.00
 1028
                                                                                                         -1.00 818.00
                                                                                                                                                                    63 16.00
                        21:00:00
                                                                18:00:00
                                                                                                        13.30 988.60
                                                                                                                                                                     78 16.00
                        23:00:00
                 EWR
                                                                22:00:00
                                                                                                         4.40 1.025.00
                                                                                                                                                                    62 16.00
                                                                                        26 17.20 ... -3.00 820.00 -11.00
```

Before getting our test train split, we performed one hot encoding on the origin airport and weather description columns.

```
set(finalData['description_d'])

set(finalData['description_d'])

finalData['description_s'] = finalData['description_s'].replace(('Broken clouds': 'Cloudy', 'Few clouds': 'Cloudy', 'Fog': 'Cloudy', 'Naze': 'Cloudy', 'Naze': 'Cloudy', 'Naze': 'Cloudy', 'Statered clouds': 'Cloudy', 'Statered clouds': 'Cloudy', 'Fey: 'Cloudy', 'Naze': 'Cloudy', 'Overcast clouds': 'Statered clouds': 'Cloudy', 'Fey: 'Cloudy', 'Naze': 'Cloudy', 'Overcast clouds': 'Statered clouds': 'Cloudy', 'Fey: 'Cloudy', 'Naze': 'Cloudy', 'Overcast clouds': 'Statered clouds': 'Cloudy', 'Fey: 'Cloudy', 'Naze': 'Cloudy', 'Overcast clouds': 'Cloudy', 'Statered clouds': 'Cloudy', 'Fey: 'Cloudy', 'Naze': 'Cloudy', 'Overcast clouds': 'Cloudy', 'Statered clouds': 'Cloudy', 'Fey: 'Cloudy', 'Naze': 'Cloudy', 'Naze': 'Cloudy', 'Statered clouds': 'Cloudy', 'Fey: 'Cloudy', 'Naze': 'Cloudy', 'Naze': 'Cloudy', 'Naze': 'Cloudy', 'Naze': 'Cloudy', 'Statered clouds': 'Cloudy', 'Fey: 'Cloudy', 'Fey: 'Cloudy', 'Naze': 'Cloudy', 'Naze': 'Cloudy', 'Statered clouds': 'Cloudy', 'Fey: 'Cloudy', 'Fey: 'Cloudy', 'Naze': 'Cloudy', 'Naze': 'Cloudy', 'Naze': 'Cloudy', 'Fey: 'Cloudy', 'Fey:
```

To find the scores of the ML model we split the data into training and testing data and used a standard scaler to scale the data.

```
# Scaling the data

xTrain, xTest, yTrain, yTest = train_test_split(subsetData.drop(columns = ['Arrival Delay (Minutes)']), subsetData['Arrival Delay (Minutes)'], test_size=0.20, random_state=40)

xTrain

yTest

from sklearn.preprocessing import StandardScaler

sc = StandardScaler()

xTrain = pd.DataFrame(sc.fit_transform(xTrain), columns = xTrain.columns, index = xTrain.index)

xTest = pd.DataFrame(sc.transform(xTest), columns = xTest.columns, index = xTest.index)

xTrain

xTest

yTrain

yTest
```

	Flight Number	Scheduled Elapsed Time (Minutes)	temp_d	pres_d	dewpt_d	wind_spd_d	wind_dir_d	precip_d	snow_d	rh_d	vis_s	description_d_Cloudy	description_d_Rainy	description_d_Snow	description_s_Cloudy	description_s
502	-1.07	1.41	1.30	-1.36	1.47	1.75	-0.14	0.34	-0.11	0.10	0.56	0.42	-0.29	-0.16	0.57	
957	0.86	-0.32	-0.96	0.86		0.90	0.78		-0.11	-0.35	0.56	0.42	-0.29	-0.16	0.57	
262	0.90	-1.15	0.33	-0.20	0.23	0.69	-0.38	-0.27	-0.11	-0.40	-1.14	0.42	-0.29	-0.16	0.57	
341	0.94	-0.34	1.25	-0.39		0.00	1.13		-0.11	-0.68	-1.14	0.42	-0.29	-0.16	0.57	
	0.90	-1.15	-0.86	-1.18	-0.74	1.11	0.78	-0.27	-0.11	0.21	-1.14	0.42	-0.29	-0.16	-1.74	
626	1.72	-1.09	0.69	-0.20	0.97	-0.72	0.20	-0.27	-0.11	0.43	0.56	0.42	-0.29	-0.16	0.57	
016		-1.05	-2.14	-1.73	-2.43	2.65	0.09	-0.27	-0.11	-0.68	0.56	0.42	-0.29	-0.16		
165	0.69	-1.17	-0.76	-1.40	-0.27	-0.21	-1.88	0.95	-0.11	1.27	-3.41	-2.39	3.48	-0.16	0.57	
	-0.13		-2.12	1.41	-1.80	0.69	0.96			1.05	-1.14	0.42	-0.29	-0.16	-1.74	
219	-0.46	1.39	0.27	0.35	-0.95	1.11	0.78	-0.27	-0.11	-2.08	0.56	0.42	-0.29	-0.16	0.57	
		Scheduled														
	Flight Number	Flancod	temp_d	pres_d	dewpt_d	wind_spd_d	wind_dir_d	precip_d	snow_d	rh_d	vis_s	description_d_Cloudy	description_d_Rainy	description_d_Snow	description_s_Cloudy	description_s
.66	Flight	Elapsed Time	temp_d 0.58	pres_d -0.06	dewpt_d 0.75	wind_spd_d -1.40	wind_dir_d	precip_d -0.27	snow_d -0.11			description_d_Cloudy	description_d_Rainy -0.29	description_d_Snow -0.16	description_s_Cloudy 0.57	description_s
666	Flight Number	Elapsed Time (Minutes)								0.16	0.56					description_s
43	Flight Number 1.72	Elapsed Time (Minutes) -1.09	0.58	-0.06	0.75	-1.40	0.72	-0.27	-0.11	0.16 -1.02	0.56 0.56	0.42	-0.29	-0.16	0.57	description_s
43 37	Flight Number 1.72 -1.07	Elapsed Time (Minutes) -1.09 1.39	0.58	-0.06 -0.06	0.75 0.97	-1.40 0.43	0.72 1.13	-0.27 -0.27	-0.11 -0.11	0.16 -1.02 0.43	0.56 0.56 0.56	0.42	-0.29 -0.29	-0.16 -0.16	0.57 0.57	description_s
43 37 98	Flight Number 1.72 -1.07 -1.07	Elapsed Time (Minutes) -1.09 1.39 1.31	0.58 1.40 -0.60	-0.06 -0.06 -0.71	0.75 0.97 -0.39	-1.40 0.43 1.33	0.72 1.13 0.78	-0.27 -0.27 -0.27	-0.11 -0.11 -0.11	0.16 -1.02 0.43 -0.85	0.56 0.56 0.56 0.56	0.42 0.42 0.42	-0.29 -0.29 -0.29	-0.16 -0.16 -0.16	0.57 0.57 -1.74	description_s
43 37 98	Flight Number 1.72 -1.07 -1.07	Elapsed Time (Minutes) -1.09 1.39 1.31	0.58 1.40 -0.60 1.62	-0.06 -0.06 -0.71 -0.20	0.75 0.97 -0.39 1.29	-1.40 0.43 1.33 -0.85	0.72 1.13 0.78 -2.11	-0.27 -0.27 -0.27 -0.27	-0.11 -0.11 -0.11 -0.11	0.16 -1.02 0.43 -0.85	0.56 0.56 0.56 0.56	0.42 0.42 0.42 0.42	-0.29 -0.29 -0.29 -0.29	-0.16 -0.16 -0.16 -0.16	0.57 0.57 -1.74 0.57	description_s
43 37 98 85	1.72 -1.07 -1.07 -1.07 -0.46	Elapsed Time (Minutes) -1.09 1.39 1.31 1.41	0.58 1.40 -0.60 1.62 -0.09	-0.06 -0.06 -0.71 -0.20 -0.25	0.75 0.97 -0.39 1.29 0.56	-1.40 0.43 1.33 -0.85	0.72 1.13 0.78 -2.11 1.59	-0.27 -0.27 -0.27 -0.27 -0.27	-0.11 -0.11 -0.11 -0.11	0.16 -1.02 0.43 -0.85 1.66	0.56 0.56 0.56 0.56 -2.56	0.42 0.42 0.42 0.42	-0.29 -0.29 -0.29 -0.29 -0.29	-0.16 -0.16 -0.16 -0.16 -0.16	0.57 0.57 -1.74 0.57	description_s
43 37 98 85 	1.72 -1.07 -1.07 -1.07 -0.46	Elapsed Time (Minutes) -1.09 1.39 1.31 1.41 1.39	0.58 1.40 -0.60 1.62 -0.09	-0.06 -0.06 -0.71 -0.20 -0.25	0.75 0.97 -0.39 1.29 0.56	-1.40 0.43 1.33 -0.85 -0.64	0.72 1.13 0.78 -2.11 1.59	-0.27 -0.27 -0.27 -0.27 -0.27	-0.11 -0.11 -0.11 -0.11 -0.11 	0.16 -1.02 0.43 -0.85 1.66	0.56 0.56 0.56 0.56 -2.56 	0.42 0.42 0.42 0.42 0.42	-0.29 -0.29 -0.29 -0.29 -0.29	-0.16 -0.16 -0.16 -0.16 -0.16	0.57 0.57 -1.74 0.57 0.57	description_s
43 37 98 85 39	1.72 -1.07 -1.07 -1.07 -0.46 	Elapsed Time (Minutes) -1.09 1.39 1.31 1.41 1.39 	0.58 1.40 -0.60 1.62 -0.09 	-0.06 -0.06 -0.71 -0.20 -0.25 	0.75 0.97 -0.39 1.29 0.56 	-1.40 0.43 1.33 -0.85 -0.64 	0.72 1.13 0.78 -2.11 1.59 	-0.27 -0.27 -0.27 -0.27 -0.27 0.04 	-0.11 -0.11 -0.11 -0.11 -0.11 -0.11	0.16 -1.02 0.43 -0.85 1.66 	0.56 0.56 0.56 0.56 -2.56 0.56 -1.14	0.42 0.42 0.42 0.42 0.42 0.42	-0.29 -0.29 -0.29 -0.29 -0.29 	-0.16 -0.16 -0.16 -0.16 -0.16 -0.16	0.57 0.57 -1.74 0.57 0.57 	description <u>.</u>
43 37 98 85 39 1 15 94	1.72 -1.07 -1.07 -1.07 -1.07 -0.46 -0.43 -0.13 -0.13	Elapsed Time (Minutes) -1.09 -1.31 -1.41 -1.39 -1.15 -1.17 -1.43	0.58 1.40 -0.60 1.62 -0.09 -0.91 -1.89 -2.94 0.84	-0.06 -0.06 -0.71 -0.20 -0.25 -1.54 0.63 2.11 0.12	0.75 0.97 -0.39 1.29 0.56 -0.62 -1.69 -2.85	-1.40 0.43 1.33 -0.85 -0.64 1.33 -0.42 -1.19	0.72 1.13 0.78 -2.11 1.59 0.32 1.48 -0.95	-0.27 -0.27 -0.27 -0.27 -0.27 -0.27 -0.27 -0.27	-0.11 -0.11 -0.11 -0.11 -0.11 -0.11 -0.11 -0.11	0.16 -1.02 0.43 -0.85 1.66 0.71 0.60 0.49	0.56 0.56 0.56 -2.56 0.56 -1.14 -3.41	0.42 0.42 0.42 0.42 0.42 0.42 0.42	-0.29 -0.29 -0.29 -0.29 -0.29 -0.29 -0.29	-0.16 -0.16 -0.16 -0.16 -0.16 -0.16 -0.16	0.57 0.57 -1.74 0.57 0.57 0.57 -1.74	description <u>:</u>
43 37 98 85 39 1 15	1.72 -1.07 -1.07 -1.07 -0.46 -0.46 -0.13	Elapsed Time (Minutes) -1.09 -1.39 -1.31 -1.41 -1.391.39 -1.15 -1.17	0.58 1.40 -0.60 1.62 -0.09 -0.91 -1.89 -2.94	-0.06 -0.06 -0.71 -0.20 -0.25 -1.54 0.63 2.11	0.75 0.97 -0.39 1.29 0.56 -0.62 -1.69 -2.85	-1.40 0.43 1.33 -0.85 -0.64 1.33 -0.42 -1.19	0.72 1.13 0.78 -2.11 1.59 0.32 1.48 -0.95	-0.27 -0.27 -0.27 -0.27 -0.04 -0.27 -0.27	-0.11 -0.11 -0.11 -0.11 -0.11 -0.11 -0.11	0.16 -1.02 0.43 -0.85 1.66 0.71 0.60 0.49	0.56 0.56 0.56 -2.56 0.56 -1.14 -3.41	0.42 0.42 0.42 0.42 0.42 0.42 0.42	-0.29 -0.29 -0.29 -0.29 -0.29 -0.29	-0.16 -0.16 -0.16 -0.16 -0.16 	0.57 0.57 -1.74 0.57 0.57 0.57 0.57	description_
43 37 98 85 39 1 15 94	1.72 -1.07 -1.07 -1.07 -1.07 -0.46 -0.43 -0.13 -0.13	Elapsed Time (Minutes) -1.09 1.39 1.31 1.41 1.39 1.39 -1.15 -1.17 1.43 -0.28	0.58 1.40 -0.60 1.62 -0.09 -0.91 -1.89 -2.94 0.84	-0.06 -0.06 -0.71 -0.20 -0.25 -1.54 0.63 2.11 0.12	0.75 0.97 -0.39 1.29 0.56 -0.62 -1.69 -2.85	-1.40 0.43 1.33 -0.85 -0.64 1.33 -0.42 -1.19	0.72 1.13 0.78 -2.11 1.59 0.32 1.48 -0.95	-0.27 -0.27 -0.27 -0.27 -0.27 -0.27 -0.27 -0.27	-0.11 -0.11 -0.11 -0.11 -0.11 -0.11 -0.11 -0.11	0.16 -1.02 0.43 -0.85 1.66 0.71 0.60 0.49	0.56 0.56 0.56 -2.56 0.56 -1.14 -3.41	0.42 0.42 0.42 0.42 0.42 0.42 0.42	-0.29 -0.29 -0.29 -0.29 -0.29 -0.29 -0.29	-0.16 -0.16 -0.16 -0.16 -0.16 -0.16 -0.16	0.57 0.57 -1.74 0.57 0.57 0.57 -1.74	description
13 37 98 35 15 94	1.72 -1.07 -1.07 -0.46 -0.13 -0.13 -1.07 0.86 ws × 29 colums	Elapsed Time (Minutes) -1.09 1.39 1.31 1.41 1.39 1.39 -1.15 -1.17 1.43 -0.28	0.58 1.40 -0.60 1.62 -0.09 -0.91 -1.89 -2.94 0.84	-0.06 -0.06 -0.71 -0.20 -0.25 -1.54 0.63 2.11 0.12	0.75 0.97 -0.39 1.29 0.56 -0.62 -1.69 -2.85	-1.40 0.43 1.33 -0.85 -0.64 1.33 -0.42 -1.19	0.72 1.13 0.78 -2.11 1.59 0.32 1.48 -0.95	-0.27 -0.27 -0.27 -0.27 -0.27 -0.27 -0.27 -0.27	-0.11 -0.11 -0.11 -0.11 -0.11 -0.11 -0.11 -0.11	0.16 -1.02 0.43 -0.85 1.66 0.71 0.60 0.49	0.56 0.56 0.56 -2.56 0.56 -1.14 -3.41	0.42 0.42 0.42 0.42 0.42 0.42 0.42	-0.29 -0.29 -0.29 -0.29 -0.29 -0.29 -0.29	-0.16 -0.16 -0.16 -0.16 -0.16 -0.16 -0.16	0.57 0.57 -1.74 0.57 0.57 0.57 -1.74	description_

```
56
626
        380
165
          6
         -6
219
Name: Arrival Delay (Minutes), Length: 826, dtype: int64
466
643
498
      -19
285
239
       17
15
694
Name: Arrival Delay (Minutes), Length: 207, dtype: int64
```

- After getting the data we fitted this data into multiple models to find out the best score and accuracy to get our final predictions.
- According to our data and scores we got the decision tree regressor that gave us the best result in terms of accuracy and score vise.

 To get our final prediction we added test data into a dataframe, cleaned the data up and got it into the same format as our test data. We used the decision tree regressor on our train data to get our final predictions and put it into a csv file.

```
finalTesttData = pd.read_csv("C:/Users/nirmi/Desktop/SU_SEM3/IntroToML/Project/Data sets/project_test_data.csv")
     finalTesttData.head()
        finalTesttData['description_d'] = finalTesttData['description_d'].replace({'Light shower rain' : 'Rainy'})
        finalTesttData
       testSubsetData = getOhe(finalTesttData, 'description_d')
       testSubsetData = getOhe(testSubsetData, 'description_s')
testSubsetData = getOhe(testSubsetData, 'Origin Airport')
       testSubsetData.drop(columns=['Scheduled Arrival Time', 'Scheduled departure time'],inplace=True)
       testSubsetData
       from sklearn.preprocessing import StandardScaler
       sc = StandardScaler()
       testData = pd.DataFrame(sc.fit_transform(testSubsetData), columns = testSubsetData.columns,
                                                         index = testSubsetData.index)
       testData
   from sklearn.tree import DecisionTreeRegressor from sklearn.model selection import cross val score
   modelDT = DecisionTreeRegressor(max_depth=5, min_samples_split=6, min_samples_leaf=2, random_state=10)
modelDT.fit(xTrain, yTrain)
modelDT.score(xTrain, yTrain)
   testOutputDI = pd.DataFrame(modelDI.predict(xTest), index-xTest.index, columns=['predArrivalDelay'])
testOutputDI = testOutputDI.merge(yTest, left_index-True, right_index=True)
mean_absolute_error = abs(testOutputDI['predArrivalDelay'] - testOutputDT['Arrival Delay (Minutes)']).mean()
print('Decision Tree Regression - Mean absolute error is ')
0.44291032914796336
Decision Tree Regression - Mean absolute error is 28.41618340194153
   test_output_d.head()
C:\Users\nirmi\anaconda\\lib\site-packages\sklearn\base.py:493: FutureWarning: The feature names should match those that were passed during fit. Starting version 1.2, an error will be raised. Feature names unseen at fit time:
- Scheduled Elapsed Time
Feature names seen at fit time, yet now missing:
- Scheduled Elapsed Time (Minutes)
 warnings.warn(message, FutureWarning)
   pred_Arrival Delay
            5.88
```

```
def flight_status(delay):
    if delay < -10:
        return "Early"
    elif delay >= -10 and delay < 10:
        return "On-time"
    elif delay > 10 and delay <= 30:
        return "Late"
    else:
        return "Severely late"

test_output_d['Flight Status'] = test_output_d['pred_Arrival Delay'].apply(lambda x: flight_status(x))
test_output_d</pre>
```

• For our final prediction we got 11 out of 32 predictions correct, that is 34.48% accuracy, that is a clear improvement over our initial prediction that was just 20%.

Conclusion:

This project gave us an insight into the power of predictive analysis using Machine Learning algorithms like regression. During the course of this project, we employed all our knowledge from our coursework including scaling the data so that they are in similar ranges, handling categorical data as a feature, using one-hot encoding for such features and employing different models and fine tuning their parameters to get the desired results. We were also able to witness the advantages and limitations of using certain models and their effects on accuracy.

The future scope for this project is to analyze other features apart from weather that affect the flight delays like Air Traffic Control delays, technical issues, Security concerns, etc. and predict these flight delays more accurately.

Reference:

- Flight data: https://www.transtats.bts.gov/ontime/
- Weather data: https://www.weatherbit.io/account/create
- SKlearn: https://scikit-learn.org/stable/