

# CIS662 - INTRO TO MACHINE LEARNING

## PROJECT REPORT

Vaama Shrikrishna Nikam  
Samyuktha Chaparla  
Vinay Anjaneya Desiraju  
Raviteja Reddy Bayapu

## Table of Contents

<b>INTRODUCTION.....</b>	<b>3</b>
<b>DATA COLLECTION .....</b>	<b>4</b>
<b>DATA PREPROCESSING .....</b>	<b>4</b>
FEATURE SELECTION.....	4
DATA CLEANING .....	5
DATA MERGING.....	8
<b>MODEL TRAINING .....</b>	<b>12</b>
<b>MODEL TESTING.....</b>	<b>15</b>
<b>PREDICTIONS: .....</b>	<b>16</b>
<b>CONCLUSION .....</b>	<b>21</b>
<b>REFERENCES: .....</b>	<b>22</b>

## INTRODUCTION

In this project we are focusing on the real-world difficulty of predicting airline arrival times at the Syracuse (SYR) airport. Our primary goal is to understand the major factors that could be responsible for the flight delays including weather conditions and then creating a Machine Learning prediction model to predict whether the flights from major airports such as Chicago (ORD), New York (JFK), and Orlando (MCO) will arrive early, on time, or delayed at Syracuse (SYR) airport.

We all know how flight delays affect passengers' schedules and have experienced this at least once during travel. Accurate arrival and delay estimates are crucial. In our project, we are dealing with a classification problem where we are categorizing the arrival status of the flights as 'EARLY', 'ON-TIME', and 'DELAYED'. Furthermore, for each origin airport, we consider two flights flying into Syracuse, a 'previous flight' and a 'latter flight'. Our goal is to estimate the arrival status of the previous flight and whether the latter flight will arrive at SYR early, on time, or late, based on the condition of the previous flight.

Our approach involves training a machine learning model using historical flight and weather data. We recognized the importance of considering weather conditions, which have a major effect on flight schedules, in our analysis. Through thorough analysis and model refinement, we aim to accurately predict the arrival status of both early and later flights mentioned above.

This study will go from data collection state to the model development process, and finally predicting the arrival status of the flights. Our goal is to build the model in such a way we can make accurate predictions of the status of the flights. We will also follow the ground truth criteria, which classifies flights as early, on-time, or late based on their actual arrival times compared to scheduled times, with a 5-minute leeway for discrepancies.

## DATA COLLECTION

Collecting data is the first and foremost step in building any machine learning model. The quality and relevance of the data to our model plays an integral part in how accurate the predictions will be. We believed that to accurately predict the arrival statuses of flights from different airports to SYR, we need to get the historical flight and weather data.

We obtained flight data from 2019 to 2023 from the Bureau of Transportation Statistics (BTS) website. The BTS website provides comprehensive and customizable datasets on flight operations, consisting of details such as flight numbers, carrier codes, arrival and departure times, arrival delays, origin and destination airports and many other features. The focus was on flights flying from origin airports ORD, JFK and MCO to SYR, which we filtered. We then had a tailored dataset to work with, which still had some features which we felt were unnecessary and were removed in preprocessing.

Our next hurdle was to gather the historical weather data. We used Weatherbit API which gave us access to accurate historical weather data and let us customize the parameters that we considered were fit to help us train our model better. This API does have restrictions regarding how much data we can retrieve in a day (each key allows the retrieval of 2 years of data and has a daily restriction of up to 4 years of data), which we resolved by using multiple keys and running python scripts to fetch the hourly weather data for all our origin locations (ORD, JFK and MCO) and destination (SYR). A lot of effort was put into merging the weather datasets for multiple locations across multiple years.

We also had to collect the weather forecast to predict the flight status during the dates 04/19/2024 to 04/23/2024. We got this data using the same site, Weatherbit, because we needed the data with the same features which we originally considered while training the model.

All this data is of prime importance and laid the groundwork on which we developed the rest of our model.

## DATA PREPROCESSING

### FEATURE SELECTION:

Data preprocessing is the initial and crucial step for any machine learning project. All the features in our dataset might not be closely related to our target variable, which is the ‘arrival status’ in our case. For any machine learning model, their performance depends on the clean, consistent and relevant data. In our project, we are carefully undertaking the data preprocessing stage to prepare both the weather and airline data for analysis.

We considered same features for predicting both the former and latter flights statuses, with the exception of adding the ‘arrival statuses’ of the previous flight to the latter flight. Our research from internet sources helped to choose relevant factors for the weather data, such as ‘Temperature’, ‘Apparent Temperature’, ‘Clouds’, ‘Precipitation’, ‘Relative Humidity’, ‘Wind Speed’, ‘wind\_gust\_spd’, ‘Wind Direction’, ‘Pressure’, ‘Sea Level Pressure’, ‘Visibility’, ‘Snow’, ‘Date’ and ‘Hour’. We can better understand the underlying patterns and causes impacting flight arrival timings by looking at these variables, which are known to have a substantial impact on flight operations.

```

temp                float64
app_temp            float64
clouds              float64
precip              float64
rh                  float64
wind_spd            float64
wind_gust_spd       float64
wind_dir            float64
weather             object
pres                float64
slp                 float64
vis                 float64
snow                float64
Date                datetime64[ns]
Hour                int32
dtype: object

```

Figure: Features of weather data

Similarly, for the flights data, we chose all the relevant features including 'Scheduled Departure Hour', 'Scheduled Departure Minute', 'Scheduled Arrival Hour', 'Scheduled Arrival Minute', 'Date', 'Origin Airport', 'Day Name', and 'Arrival Status' which is calculated by using the 'Arrival Delay in minutes' feature. All these features are common for both earlier flights and the latter flight. Whereas for the later flight we include an extra feature which is the previous flight's status which is basically the arrival status of the previous flight from that origin airport. These features will enable us to assess the relationship between different factors and the likelihood of flights arriving early, on-time, or delayed.

```

Carrier Code        object
Date (MM/DD/YYYY)   datetime64[ns]
Flight Number        int64
Tail Number          object
Origin Airport        object
Destination Airport   object
Scheduled departure time object
Actual departure time object
Departure delay (Minutes) int64
Scheduled Arrival Time object
Actual Arrival Time   object
Arrival Delay (Minutes) int64
dtype: object

```

Figure: Features of Airlines data

## DATA CLEANING:

### Missing Values Handling:

- To begin data preprocessing, we checked for missing values in the dataset and calculated the cumulative count for each feature.
- In the airline dataset, we discovered that the 'Tail Number' field has missing values. We handled this by dropping the associated rows with the dropna() function.

```

Carrier Code        0
Date (MM/DD/YYYY)   0
Flight Number        0
Tail Number          53
Origin Airport        0
Destination Airport   0
Scheduled departure time 0
Actual departure time 0
Departure delay (Minutes) 0
Scheduled Arrival Time 0
Actual Arrival Time   0
Arrival Delay (Minutes) 0
dtype: int64

```

Figure: Handling Missing Values

**FEATURE ENGINEERING:**

In weather datasets and flight datasets, we formatted the 'timestamp\_local' column as datetime and then extracted the date and hour out of the timestamp with the perspective of analyzing the weather data on a daily and hourly basis with more precision. This extraction also helps us while merging with the flight data to get the most accurate weather information. To ensure consistency, we transformed the "Date" field to datetime format and deleted duplicate columns (datetime and timestamp\_local).

```
All_weather['timestamp_local'] = pd.to_datetime(All_weather['timestamp_local'])
syn_weather['timestamp_local'] = pd.to_datetime(syn_weather['timestamp_local'])
syn_weather.head()
```

	datetime	timestamp_local	temp	app_temp	clouds	precip	rh	wind_spd	wind_gust_spd	wind_dir	weather	pres	slp	vis	snow
0	2019-01-01:05	2019-01-01 00:00:00	7.20	4.30	87.00	1.50	85.00	4.59	9.80	150.00	Light rain	982.00	997.00	16.00	0.00
1	2019-01-01:06	2019-01-01 01:00:00	8.90	8.90	100.00	0.50	79.00	5.09	11.80	190.00	Overcast clouds	984.00	999.00	16.00	0.00
2	2019-01-01:07	2019-01-01 02:00:00	9.40	9.40	100.00	0.00	83.00	4.59	12.90	180.00	Overcast clouds	982.00	997.00	16.00	0.00
3	2019-01-01:08	2019-01-01 03:00:00	11.10	11.10	87.00	0.00	80.00	8.19	13.90	200.00	Overcast clouds	980.00	996.00	16.00	0.00
4	2019-01-01:09	2019-01-01 04:00:00	10.60	10.60	100.00	0.50	85.00	6.20	16.80	230.00	Overcast clouds	981.00	996.00	16.00	0.00

```
All_weather['Date'] = All_weather['timestamp_local'].dt.date
All_weather['Hour'] = All_weather['timestamp_local'].dt.hour

syn_weather['Date'] = syn_weather['timestamp_local'].dt.date
syn_weather['Hour'] = syn_weather['timestamp_local'].dt.hour

syn_weather.head()
```

	datetime	timestamp_local	temp	app_temp	clouds	precip	rh	wind_spd	wind_gust_spd	wind_dir	weather	pres	slp	vis	snow	Date	Hour
0	2019-01-01:05	2019-01-01 00:00:00	7.20	4.30	87.00	1.50	85.00	4.59	9.80	150.00	Light rain	982.00	997.00	16.00	0.00	2019-01-01	0
1	2019-01-01:06	2019-01-01 01:00:00	8.90	8.90	100.00	0.50	79.00	5.09	11.80	190.00	Overcast clouds	984.00	999.00	16.00	0.00	2019-01-01	1
2	2019-01-01:07	2019-01-01 02:00:00	9.40	9.40	100.00	0.00	83.00	4.59	12.90	180.00	Overcast clouds	982.00	997.00	16.00	0.00	2019-01-01	2
3	2019-01-01:08	2019-01-01 03:00:00	11.10	11.10	87.00	0.00	80.00	8.19	13.90	200.00	Overcast clouds	980.00	996.00	16.00	0.00	2019-01-01	3
4	2019-01-01:09	2019-01-01 04:00:00	10.60	10.60	100.00	0.50	85.00	6.20	16.80	230.00	Overcast clouds	981.00	996.00	16.00	0.00	2019-01-01	4

```
All_weather['Date'] = pd.to_datetime(All_weather['Date'])
syn_weather['Date'] = pd.to_datetime(syn_weather['Date'])

syn_weather.drop(['datetime', 'timestamp_local'], axis=1, inplace=True)
All_weather.drop(['datetime', 'timestamp_local'], axis=1, inplace=True)
```

Figure: Data Type Conversions

New features were added to the flight data to gain deeper insights from the existing data. The airline dataset's date and time columns were used to infer attributes like 'day of the week'. The 'Day Name' column adds context and interpretability to the dataset by displaying the day of the week in a human-readable fashion and helps to analyze how the day of the week affects arrival times. Analysts can readily discover and analyze patterns and trends in flight operations, such as fluctuations in volume or delays on various days of the week.

```
#even_rows['Date (MM/DD/YYYY)'] = pd.to_datetime(even_rows['Date (MM/DD/YYYY)'])

first_df['Day of Week'] = first_df['Date (MM/DD/YYYY)'].dt.dayofweek

# Map the day of the week to day names
day_names = ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday', 'Sunday']
first_df['Day Name'] = first_df['Day of Week'].map(lambda x: day_names[x])

first_df['Day Name'] = pd.Categorical(first_df['Day Name'], categories=day_names)

print(first_df['Day Name'].cat.categories)

Index(['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday',
       'Sunday'],
      dtype='object')

# Convert 'Scheduled departure time' to datetime
first_df['Scheduled departure time'] = pd.to_datetime(first_df['Scheduled departure time'], format='%H:%M:%S')

first_df['Scheduled Arrival Time'] = pd.to_datetime(first_df['Scheduled Arrival Time'], format='%H:%M:%S')

# Extract hour, minute, and second components
first_df['Scheduled departure hour'] = first_df['Scheduled departure time'].dt.hour
first_df['Scheduled departure minute'] = first_df['Scheduled departure time'].dt.minute

first_df['Scheduled Arrival hour'] = first_df['Scheduled Arrival Time'].dt.hour
first_df['Scheduled Arrival minute'] = first_df['Scheduled Arrival Time'].dt.minute
```

Figure: Adding new columns

Like weather data, scheduled hour and minute components for departure and arrival were extracted for the flight data which improves the analysis by providing granularity. This detail allows for a more detailed analysis of departure and arrival patterns, identifying peak departure times and common delays at specific times of day. By breaking down times into smaller parts, we may better understand the hidden patterns.

## Arrival Status Classification:

We added an "Arrival Status" feature to the airline dataset, which categorizes flights depending on arrival delays. This is our dependent variable (goal), and we need to predict this value.

Flights were classified as 'EARLY', 'ON-TIME', or 'DELAYED' based on predetermined criteria for 'arrival delay in minutes'.

```
# Define a function to determine the arrival status
def determine_arrival_status(delay):
    if delay > 5:
        return 2
    elif delay < -5:
        return 1
    else:
        return 0

# Apply the function to create the "Arrival Status" column
first_df['Arrival Status'] = first_df['Arrival Delay (Minutes)'].apply(determine_arrival_status)

first_df.drop('Arrival Delay (Minutes)', axis=1, inplace=True)
```

```
# Define a dictionary to map numerical labels to their meanings
label_mapping = {
    2: "late",
    1: "early",
    0: "ontime"
}
```

Figure: determining the arrival status

By effectively handling the missing values, converting data types, finding important features, adding new features, and classifying arrival statuses, we made sure that our datasets were well-preprocessed and structured for further analysis and modeling. These preprocessing approaches establish the fundamentals for creating reliable predictive models and gaining useful insights to help enhance flight operations.

## DATA MERGING

### First Flight:

We started with merging the Syracuse weather data (destination weather) with the earlier flight data which in our case is the ‘first flight data’ by mapping the datasets using common columns such as 'Date' and 'Arrival Hour'. Next, we merge this with all the origin locations weather data similarly by using the common columns like 'Date', 'Departure Hour', and 'Origin Airport'.

```
merged_df = pd.merge(first_df, syr_weather, left_on=["Date (MM/DD/YYYY)", "Scheduled Arrival hour"], right_on=["Date", "Hour"], how="left")

merged_df.head()
```

	Carrier Code	Date (MM/DD/YYYY)	Origin Airport	Day Name	Scheduled departure hour	Scheduled departure minute	Scheduled Arrival hour	Scheduled Arrival minute	Arrival Status	temp	...	wind_spd	wind_gust_spd	wind_dir	weather	pres	slp	vis	snow	Date	Hour
0	UA	2023-01-01	ORD	Sunday	21	10	23	57	1	3.90	...	1.60	3.60	280.00	Overcast clouds	1,000.00	1,015.00	11.00	0.00	2023-01-01	23
1	UA	2023-01-02	ORD	Monday	18	14	21	7	1	4.40	...	1.20	2.40	55.00	Overcast clouds	1,006.00	1,021.00	16.00	0.00	2023-01-02	21

```
merged_df = pd.merge(merged_df, All_weather, how="left", left_on=["Date (MM/DD/YYYY)", "Scheduled departure hour", "Origin Airport"], right_on=["Date", "Hour", "Location"])

merged_df.drop(['Date_x', 'Hour_x', 'Date_y', 'Hour_y'], axis=1, inplace=True)
merged_df.head()
```

	Carrier Code	Date (MM/DD/YYYY)	Origin Airport	Day Name	Scheduled departure hour	Scheduled departure minute	Scheduled Arrival hour	Scheduled Arrival minute	Arrival Status	temp_x	...	rh_y	wind_spd_y	wind_gust_spd_y	wind_dir_y	weather_description	pres_y	slp_y	vis_y	snow_y	Locat
0	UA	2023-01-01	ORD	Sunday	21	10	23	57	1	3.90	...	87	1.60	3.60	250	Overcast clouds	992	1017	6	0.00	C
1	UA	2023-01-02	ORD	Monday	18	14	21	7	1	4.40	...	85	4.59	6.40	80	Overcast clouds	991	1015	11	0.00	C
2	UA	2023-01-03	ORD	Tuesday	18	14	21	7	0	5.00	...	88	2.10	2.80	20	Haze	977	1001	0	0.00	C
3	UA	2023-01-04	ORD	Wednesday	18	24	21	17	2	7.20	...	88	3.60	8.40	220	Overcast clouds	984	1009	6	0.00	C

Figure: Merging earlier flight data with origin and destination locations weather data

Columns in our merged data frame such as 'Scheduled departure hour', 'Scheduled departure minute', 'Scheduled Arrival hour', and 'Scheduled Arrival minute' are converted to object type (string) for further processing.

Next in the process we are converting the scheduled departure and arrival hours into categorical variables with specified categories ranging from 0 to 23. There are 24 hours in a day, and the intuition was to consider ‘hour’ as a categorical variable with 24 categories, rather than interpreting it as a numeric value with continuous quantity. This adds relevance to the feature in the model.



Similarly, the scheduled departure and arrival minutes are converted into categorical variables with specified categories ranging from 0 to 59 (there are 60 minutes in an hour).

```
all_hours = pd.Series(range(24))

# Convert to categorical variable with specified categories
merged_df['Scheduled departure hour'] = pd.Categorical(merged_df['Scheduled departure hour'], categories=all_hours, ordered=True)
merged_df['Scheduled Arrival hour'] = pd.Categorical(merged_df['Scheduled Arrival hour'], categories=all_hours, ordered=True)

all_minutes = pd.Series(range(60))

# Convert to categorical variable with specified categories
merged_df['Scheduled departure minute'] = pd.Categorical(merged_df['Scheduled departure minute'], categories=all_minutes, ordered=True)
merged_df['Scheduled Arrival minute'] = pd.Categorical(merged_df['Scheduled Arrival minute'], categories=all_minutes, ordered=True)

merged_df['weather'] = pd.Categorical(merged_df['weather'], ordered=True)
merged_df['weather_description'] = pd.Categorical(merged_df['weather_description'], ordered=True)
```

Figure: Conversion to categorical variables

As we know not all features truly contribute or influence flight arrival and departure status so after considering all the relevant features, we dropped columns like 'weather', 'weather\_description', and 'Carrier Code' from the resultant data frame. We did try out these features, but they did not improve the test accuracies by a great deal and showed signs of overfitting. We dropped the Date column before encoding. Then we used `get_dummies` on the remaining columns to convert categorical variables into a numerical format making them more suitable for model training. Finally, the 'Arrival Status' column is converted to a categorical variable for better representation and understanding during model training.

```
merged_df.drop(['weather', 'weather_description', 'Carrier Code'], axis=1, inplace=True)

date_column = first_df['Date (MM/DD/YYYY)']

# Drop the date column before encoding
first_df_encoded = pd.get_dummies(merged_df.drop(columns=['Date (MM/DD/YYYY)']), drop_first=True)

# Add the date column back to the encoded DataFrame
#first_df_encoded['Date (MM/DD/YYYY)'] = date_column
#first_df_encoded.head()

# Convert 'Arrival Status' to categorical
first_df_encoded['Arrival Status'] = first_df_encoded['Arrival Status'].astype('category')
```

Figure: Arrival Status feature conversion to categorical

## Second Flight:

After reading both earlier and later flights data we then used the function `determine_arrival_status` function (mentioned above) on the Arrival Delay minutes column to create columns like 'Arrival Status First' and 'Arrival Status Second' in the respective data frames before merging. We chose 'Carrier Code', 'Date (MM/DD/YYYY)', 'Origin Airport', 'Scheduled departure time', 'Arrival Status First' columns to keep in the first flight data frame to be merged with the later flight. We merged the later flight dataset with the earlier flight dataset based on common columns like 'Origin Airport' and 'Date (MM/DD/YYYY)' using the `pd.merge()` function. This merger allows us to combine information from both datasets into a single comprehensive dataset for analysis.

```
# Apply the function to create the "Arrival Status" column
second_df['Arrival Status Second'] = second_df['Arrival Delay (Minutes)'].apply(determine_arrival_status)
first_df_reload['Arrival Status First'] = first_df_reload['Arrival Delay (Minutes)'].apply(determine_arrival_status)

first_df_reload.head()
second_df.head()
```

	Carrier Code	Date (MM/DD/YYYY)	Flight Number	Tail Number	Origin Airport	Destination Airport	Scheduled departure time	Actual departure time	Departure delay (Minutes)	Scheduled Arrival Time	Actual Arrival Time	Arrival Delay (Minutes)	Arrival Status First
0	UA	2023-01-01	2645	N23721	ORD	SYR	21:10:00	21:07:00	-3	23:57:00	23:47:00	-10	1
1	UA	2023-01-02	1998	N802UA	ORD	SYR	18:14:00	18:17:00	3	21:07:00	20:46:00	-21	1
2	UA	2023-01-03	1998	N854UA	ORD	SYR	18:14:00	18:13:00	-1	21:07:00	21:04:00	-3	0
3	UA	2023-01-04	1998	N893UA	ORD	SYR	18:24:00	18:41:00	17	21:17:00	21:31:00	14	2
4	UA	2020-01-05	2012	N825UA	ORD	SYR	19:00:00	19:21:00	21	21:47:00	21:55:00	8	2

	Carrier Code	Date (MM/DD/YYYY)	Flight Number	Tail Number	Origin Airport	Destination Airport	Scheduled departure time	Actual departure time	Departure delay (Minutes)	Scheduled Arrival Time	Actual Arrival Time	Arrival Delay (Minutes)	Arrival Status Second
0	MQ	2019-01-01	3538	N523AE	ORD	SYR	11:45:00	11:43:00	-2	14:27:00	14:35:00	8	2
1	MQ	2019-01-01	3685	N538EG	ORD	SYR	18:50:00	18:50:00	0	21:39:00	21:59:00	20	2
2	MQ	2019-01-01	3946	N256NN	ORD	SYR	14:47:00	16:30:00	103	17:34:00	19:31:00	117	2

```
merged_df_new = pd.merge(second_df, first_df_reload, on=['Origin Airport', 'Date (MM/DD/YYYY)'], how='left')
merged_df_new.head()
```

	Carrier Code_x	Date (MM/DD/YYYY)	Flight Number	Tail Number	Origin Airport	Destination Airport	Scheduled departure time_x	Actual departure time	Departure delay (Minutes)	Scheduled Arrival Time	Actual Arrival Time	Arrival Delay (Minutes)	Arrival Status Second	Carrier Code_y	Scheduled departure time_y	Arrival Status First
0	MQ	2019-01-01	3538	N523AE	ORD	SYR	11:45:00	11:43:00	-2	14:27:00	14:35:00	8	2	NaN	NaN	NaN
1	MQ	2019-01-01	3685	N538EG	ORD	SYR	18:50:00	18:50:00	0	21:39:00	21:59:00	20	2	NaN	NaN	NaN

Figure: Merging earlier flight data with later flights data

To identify flights within a considerable departure gap, we calculated the time gap between the first and second flights' scheduled departure times. Later, we filtered away rows with time differences less than a particular threshold, as indicated in the code below, to ensure that only relevant data records were evaluated for further analysis. We wanted to consider data where the second flight's departure time was not more than 3 hours after the first flight's departure.

```
# Convert 'Scheduled departure time_x' and 'Scheduled departure time_y' columns to datetime format
merged_df_new['Scheduled departure time_x'] = pd.to_datetime(merged_df_new['Scheduled departure time_x'], format='%H:%M:%S')
merged_df_new['Scheduled departure time_y'] = pd.to_datetime(merged_df_new['Scheduled departure time_y'], format='%H:%M:%S')

# Calculate the time difference in hours
merged_df_new['time_diff'] = (merged_df_new['Scheduled departure time_x'] - merged_df_new['Scheduled departure time_y']).dt.total_seconds() / 3600

# Drop rows where the time difference is greater than or equal to 3 hours or less than 0
merged_df_new = merged_df_new.drop(merged_df_new[(merged_df_new['time_diff'] > 3) | (merged_df_new['time_diff'] < 0)].index)

# Drop the 'time_diff' column if not needed anymore
merged_df_new = merged_df_new.drop(columns=['time_diff'])

merged_df_new.head()
```

Figure: Extracting relevant data

After merging our combined flight dataset with Syracuse weather data (destination) and then merging all the origin locations', we preprocessed the data further by transforming categorical variables into dummy variables with get\_dummies. This process prepares the dataset for usage in machine learning models, allowing us to do analysis on airline arrivals.

```
merged_df_new = pd.merge(merged_df_new, syr_weather, left_on=["Date (MM/DD/YYYY)", "Scheduled Arrival hour"], right_on=["Date", "Hour"], how="left")
merged_df_new.head()
```

Figure: Merging combined flights data with Syracuse weather data

```
merged_df_new = pd.merge(merged_df_new, All_weather, how='left', left_on=['Date (MM/DD/YYYY)', 'Scheduled departure hour', 'Origin Airport'], right_on=['Date', 'Hour', 'Location'])

merged_df_new.drop(['Scheduled departure time_x', 'Actual departure time', 'Scheduled Arrival Time', 'Actual Arrival Time', 'Date (MM/DD/YYYY)', 'Day of Week'], axis=1, inplace=True)

merged_df_new.drop(['Date_x', 'Hour_x', 'Date_y', 'Hour_y'], axis=1, inplace=True)
```

Figure: Merging combined flights data with All origin locations weather data

Later we converted the hour and minute columns into categorical variables, considering them as discrete categories rather than continuous integers. This allows the model to understand the ordinal relationship between different hours and minutes, which could be important in the context of flight scheduling.

```
# Convert columns to object type
merged_df_new['Scheduled departure hour'] = merged_df_new['Scheduled departure hour'].astype(str)
merged_df_new['Scheduled departure minute'] = merged_df_new['Scheduled departure minute'].astype(str)
merged_df_new['Scheduled Arrival hour'] = merged_df_new['Scheduled Arrival hour'].astype(str)
merged_df_new['Scheduled Arrival minute'] = merged_df_new['Scheduled Arrival minute'].astype(str)

all_hours = pd.Series(range(24))

# Convert to categorical variable with specified categories
merged_df_new['Scheduled departure hour'] = pd.Categorical(merged_df_new['Scheduled departure hour'], categories=all_hours, ordered=True)
merged_df_new['Scheduled Arrival hour'] = pd.Categorical(merged_df_new['Scheduled departure hour'], categories=all_hours, ordered=True)

all_minutes = pd.Series(range(60))

# Convert to categorical variable with specified categories
merged_df_new['Scheduled departure minute'] = pd.Categorical(merged_df_new['Scheduled departure minute'], categories=all_minutes, ordered=True)
merged_df_new['Scheduled Arrival minute'] = pd.Categorical(merged_df_new['Scheduled Arrival minute'], categories=all_minutes, ordered=True)
```

Figure: Conversion of hour and minute to categorical variables

In the final stage of our preprocessing, we used `get_dummies` to extract categorical features from our dataset, such as 'Scheduled departure hour', 'Scheduled departure minute', 'Scheduled Arrival hour', and 'Scheduled Arrival minute', for example. Our goal is to find hidden patterns and relationships between specific times of day and airplane arrival status which in turn results in improved accuracy with the predictions.

```
Arr_status = merged_df_new['Arrival Status Second']

merged_df_new = pd.get_dummies(merged_df_new.drop(columns=['Arrival Status Second', 'Carrier Code_x', 'Carrier Code_y']), drop_first=True)

merged_df_new['Arrival Status Second'] = Arr_status

merged_df_new.head()
```

Figure: Getting dummies

Therefore, the process of merging and preparing data is crucial because it allows for the integration of multiple data sources, refines data quality, and prepares the dataset for future analysis. Finally, these procedures play an important role in enabling informed decision-making in airline operations.

## Model Training:

### First Flight:

```

from sklearn.ensemble import StackingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import GradientBoostingClassifier
from sklearn.pipeline import make_pipeline
from sklearn.ensemble import RandomForestClassifier

# Define your base models
base_models = [
    ('logistic', LogisticRegression(fit_intercept=True, solver='newton-cg', multi_class='multinomial', penalty=None, max_iter=100000)),
    ('random_forest', RandomForestClassifier(
        n_estimators=200, # Number of trees in the forest
        max_depth=None, # Maximum depth of the tree
        min_samples_split=2, # Minimum number of samples required to split an internal node
        min_samples_leaf=1, # Minimum number of samples required to be at a leaf node
        max_features=10, # Number of features to consider when looking for the best split
        random_state=42 # Random state for reproducibility
    ))
]

# Define meta-model
meta_model = GradientBoostingClassifier(learning_rate=0.001, max_features='log2', n_estimators=500, min_samples_leaf=1)

# Create the stacking classifier
stacking_clf = StackingClassifier(estimators=base_models, final_estimator=meta_model)

# Create pipeline with standard scaler and stacking classifier
pipeline = make_pipeline(stacking_clf)

X_train_1, X_test_1, y_train_1, y_test_1 = train_test_split(first_df_encoded.drop(columns=['Arrival Status']), first_df_encoded['Arrival Status'], test_size=0.20, random_state=42)

# Train the stacking model
pipeline.fit(X_train_1, y_train_1)

# Make predictions
train_accuracy = pipeline.score(X_train_1, y_train_1)
test_accuracy = pipeline.score(X_test_1, y_test_1)

print("Train Accuracy:", train_accuracy)
print("Test Accuracy:", test_accuracy)

```

Figure: Model training for the earlier flights

For the first flight, we have first imported the necessary libraries such as `StackingClassifier`, `LogisticRegression`, `RandomForestClassifier` and `GradientBoostingClassifier` and `make_pipeline` from `sklearn`. We are using the stacking classifier. The `StackingClassifier` class in `scikit-learn` implements this stacking approach:

- **estimators:** This parameter takes a list of tuples, where each tuple consists of a string identifier and an instance of a machine learning model.
- **final\_estimator** option defines the meta-model that generates the final prediction based on base model outputs.

After importing the models, we first generated a collection of tuples called `base_models`, with each tuple representing a base model utilized in the stacking ensemble. The first element is the model's string identifier. 'logistic', and the second element is an instance of the model, which in this case is `LogisticRegression` with specified parameters.

- **fit\_intercept=True:** Adds a constant term to the model.
- **solver='newton-cg':** Specifies the optimization algorithm.
- **multi\_class='multinomial':** Enables support for multinomial classification tasks.
- **penalty=None:** Disables regularization.

- **max\_iter=100000**: Sets a high limit for iterations, allowing sufficient time for convergence.

Next in the `base_models` we have model's string identifier which is 'random\_forest', and the second element is an instance of the model, which in this case is `RandomForestClassifier` with specified parameters.

- **n\_estimators=200** : Number of trees to be in the forest
- **max\_depth=None** : Maximum depth of the tree
- **min\_samples\_split=2** : Minimum number of samples required to split an internal node
- **min\_samples\_leaf=1** : Minimum number of samples required to be at a leaf node
- **max\_features=10** : Number of features to consider when looking for the best split
- **random\_state=42**: Random state for reproducibility

We then have a 'meta\_model' which is an instance of `GradientBoostingClassifier` that acts as the final estimator of the stacking ensemble. This has hyperparameters such as `learning_rate=0.001`, `max_features='log2'`, `n_estimators=500`, `min_samples_leaf=1` which are used to control the model's complexity.

We then created a stacking classifier. Stacking, or stacked generalization, is a different form of ensemble learning. It uses multiple base models (sometimes referred to as level-1 models), which make individual predictions on the dataset. The predictions from these base models are then passed to a final model, known as a meta-model or level-2 model, which makes the final prediction.

In our code, `base_models` consist of two tuples, representing a `LogisticRegression` model, `RandomForestClassifier` model and the `meta_model` is a `GradientBoostingClassifier`, which acts as the final estimator. The base models make predictions on the training set, producing a new dataset of their outputs. The meta-model is then trained on this new dataset, learning to make the final predictions by leveraging the outputs of the base models.

A pipeline is created to streamline the model training process. The pipeline consists of the final model used to train and evaluate the ensemble model.

We then split the dataset into training and test sets using `train_test_split`. The `test_size=0.20` allocates 20% of the data to the test set while the `random_state=42` ensures reproducibility of the split. The pipeline's `fit` method trains the stacking model on the training set.

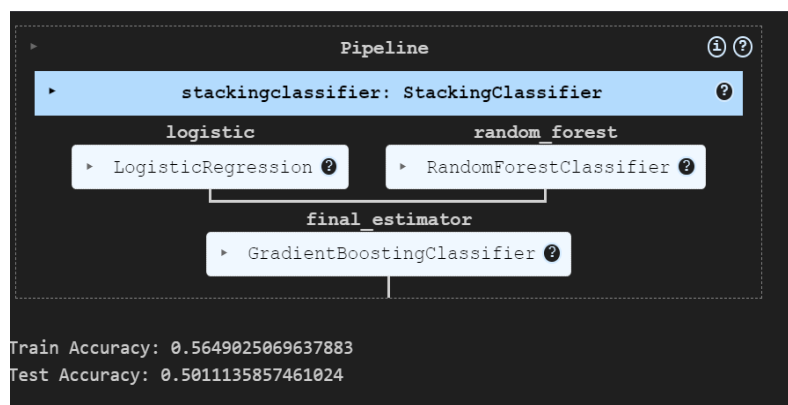


Figure: Prediction accuracy for the first model

The train accuracy for the first flight is 0.56.

## Second Flight:

```

from sklearn.ensemble import StackingClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import GradientBoostingClassifier

# Define your base models
base_models = [
    ('logistic', LogisticRegression(fit_intercept=True, solver = 'newton-cg', multi_class='multinomial', penalty = None, max_iter = 100000)),
    ('random_forest', RandomForestClassifier(
        n_estimators=100, # Number of trees in the forest
        max_depth=None, # Maximum depth of the tree
        min_samples_split=2, # Minimum number of samples required to split an internal node
        min_samples_leaf=1, # Minimum number of samples required to be at a leaf node
        max_features=10, # Number of features to consider when looking for the best split
        random_state=42 # Random state for reproducibility
    ))]

# Define your meta-model
meta_model = GradientBoostingClassifier(learning_rate = 0.01, max_features = 'sqrt', n_estimators = 300, min_samples_leaf = 2)

# Create the stacking classifier
stacking_clf = StackingClassifier(estimators=base_models, final_estimator=meta_model)

# Create pipeline with standard scaler and stacking classifier
pipeline_2 = make_pipeline(stacking_clf)

X_train_2, X_test_2, y_train_2, y_test_2 = train_test_split(merged_df_new.drop(columns = ['Arrival Status Second']), merged_df_new['Arrival Status Second'], test_size=0.20, random_state=42)

# Train the stacking model
pipeline_2.fit(X_train_2, y_train_2)

# Make predictions
train_accuracy = pipeline_2.score(X_train_2, y_train_2)
test_accuracy = pipeline_2.score(X_test_2, y_test_2)

print("Train Accuracy:", train_accuracy)
print("Test Accuracy:", test_accuracy)

```

Figure: Model training for the later flights

For the second flight, we have first imported the necessary libraries such as StackingClassifier, Logistic Regression and GradientBoostingClassifier from sklearn. We are using the stacking classifier again which uses the stacking approach where it first takes the estimators which is a list of tuples of base models and then the final estimator which is the meta model.

After importing the models, we have first created a list of tuples called base\_models. Here as well, the first element is a string identifier for the model ie. 'logistic' and the second element is an instance of the model itself, in this case, LogisticRegression with specific hyperparameters such as fit\_intercept, solver, multi\_class, penalty and max\_iter. In addition to the above base\_model we also have a model with string identifier which is 'random\_forest', and the second element is an instance of the model, which in this case is RandomForestClassifier with specified parameters. The parameters are same as it was for the first model, except for n\_estimators which is now 100.

Then we created a 'meta\_model' which is an instance of GradientBoostingClassifier acting as the final estimator for stacking ensemble. This has hyperparameters such as learning\_rate=0.01, max\_features='sqrt', n\_estimators=300, min\_samples\_leaf=2 which are used to control the model's complexity.

We then create a stacking classifier. A pipeline is then created to streamline the model training process. We then split the dataset into training and test sets using train\_test\_split. The test\_size=0.20 allocates 20% of the data to the test set while the random\_state=42 ensures reproducibility of the split. The pipeline's fit method trains the stacking model on the training set.

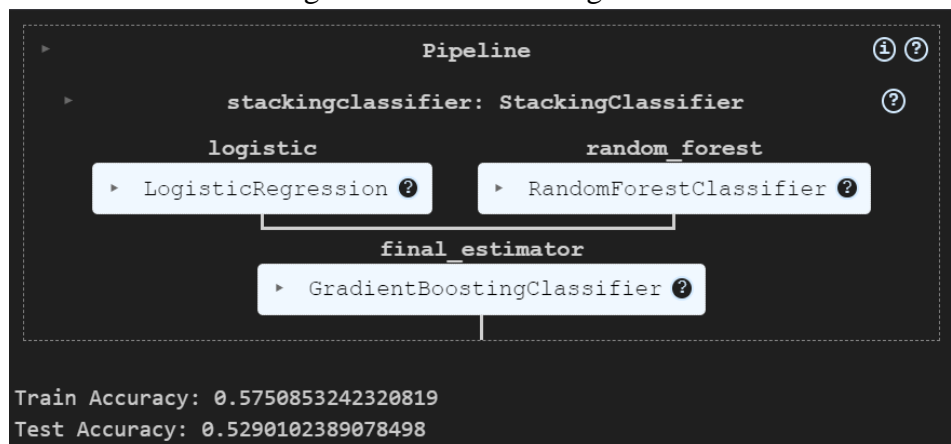


Figure: Prediction accuracy for the later flights model

The train accuracy for the second flight is 0.58.

### Model Testing:

Model testing is a crucial stage in any machine learning project. It involves evaluating how well the model performs on unseen data, providing a reliable measure of its generalization ability. In our project, after training our model on the training dataset, we assess its performance on a separate test dataset, which consists of 20% of the overall data.

The main purpose of testing is to see how well the model generalizes to new, unseen data. This reflects its real-world applicability, as a model needs to handle data beyond the training set.

The testing accuracy for the first flight is 0.50 and the testing accuracy for second flight is 0.54.



## PREDICTIONS:

First, we will be predicting for the first flights whether they are ‘early’/‘on time’/‘late’. Based on this information we will be predicting the second flights whether they are ‘early’/‘on time’/‘late’.



Figure: Extracting the earlier and later flights and then merging with weather data

We read the “Final\_Earlier\_flights.xlsx”, “Final\_Later\_flights.xlsx” from which the first and second flights for which the predictions should be done on “even\_rows”, “odd\_rows” respectively. Then, we read the merged weather forecast data of all the source and destination locations for the flights.

We have to preprocess these files in the exact manner as we did for the flight data because the model is trained on it and recognizes features of that format only.

Thus, we only keep the columns that are helpful in predictions such as ['datetime', 'temp', 'clouds', 'precip', 'rh', 'wind\_spd', 'wind\_gust\_spd', 'wind\_dir', 'pres', 'slp', 'vis', 'snow', 'Location', 'app\_temp', 'timestamp\_local'] in the merged weather forecast data.

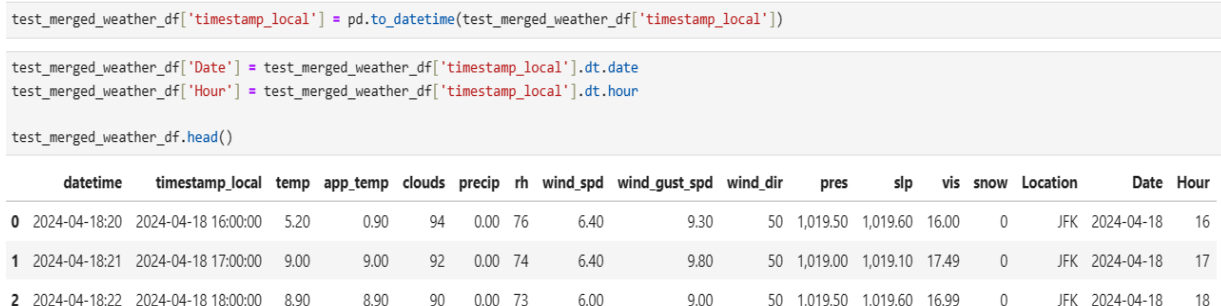


Figure: Extracting date and hour from timestamp\_local feature



The “timestamp\_local” column in merged weather forecast data is converted into of type “datetime”. Then, the date and hour from the “timestamp\_local” datetime column is extracted in two new columns “Date”, “Hour”.

```
even_rows['Date (MM/DD/YYYY)'] = pd.to_datetime(even_rows['Date (MM/DD/YYYY)'])

even_rows['Day of Week'] = even_rows['Date (MM/DD/YYYY)'].dt.dayofweek

# Map the day of the week to day names
day_names = ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday', 'Sunday']
even_rows['Day Name'] = even_rows['Day of Week'].map(lambda x: day_names[x])

even_rows['Day Name'] = pd.Categorical(even_rows['Day Name'], categories=day_names)

print(even_rows['Day Name'].cat.categories)

Index(['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday', 'Saturday',
       'Sunday'],
      dtype='object')

# Convert 'Scheduled departure time' to datetime
even_rows['Scheduled departure time'] = pd.to_datetime(even_rows['DEPARTURE TIME'], format='%H:%M:%S')

even_rows['Scheduled Arrival Time'] = pd.to_datetime(even_rows['ARRIVAL TIME'], format='%H:%M:%S')

# Extract hour, minute, and second components
even_rows['Scheduled departure hour'] = even_rows['Scheduled departure time'].dt.hour
even_rows['Scheduled departure minute'] = even_rows['Scheduled departure time'].dt.minute

even_rows['Scheduled Arrival hour'] = even_rows['Scheduled Arrival Time'].dt.hour
even_rows['Scheduled Arrival minute'] = even_rows['Scheduled Arrival Time'].dt.minute

even_rows.drop(['FLIGHT NUMBER', 'Scheduled departure time', 'Scheduled Arrival Time', 'Day of Week'], axis=1, inplace=True)
even_rows.drop(['DEPARTURE TIME', 'ARRIVAL TIME', 'ARRIVAL STATUS', 'ARRIVAL STATUS_Prev_flight_early', 'ARRIVAL STATUS_Prev_flight_ontime'], axis=1, inplace=True)

# Define the new column names and order
new_column_names = [ 'Date (MM/DD/YYYY)', 'Origin Airport', 'Day Name',
                     'Scheduled departure hour', 'Scheduled departure minute',
                     'Scheduled Arrival hour', 'Scheduled Arrival minute']

# Rename the columns
even_rows = even_rows.rename(columns={
    'ORIGIN': 'Origin Airport'
})
```

Figure: Feature Engineering

The “DATE” column in both the first and second flights is renamed as “Date (MM/DD/YYYY)”. Then, the column is converted to a datetime column. Then, two new columns are added “Day of Week” and “Day Name” from this column. The columns “Scheduled departure time”, “Scheduled Arrival Time” are converted into datetime columns and four new columns “Scheduled departure hour”, “Scheduled departure minute”, “Scheduled Arrival hour”, “Scheduled Arrival minute” are created to store the hour and minute respectively. As shown in the image above, few of the columns in first flights data file are removed, renamed and rearranged. Similar processing was done for the second flights data file as well.

```
# Filter test_merged_weather_df to include only rows where Location is SYR
filtered_weather_df = test_merged_weather_df[test_merged_weather_df['Location'] == 'SYR']
filtered_weather_df.head()

# Merge even_rows with the filtered_weather_df, which has the destination airport weather info.
merged_df_test_even = pd.merge(even_rows, filtered_weather_df, left_on=["Date (MM/DD/YYYY)", "Scheduled Arrival hour"], right_on=["Date", "Hour"], how="left")

merged_df_test_even.head()
```

	temp	app_temp	clouds	precip	rh	wind_spd	wind_gust_spd	wind_dir	pres	slp	vis	snow	Location	Date	Hour
324	4.80	2.10	87	0.00	76	3.20	4.90	90	1,002.00	1,017.40	12.40	0	SYR	2024-04-18	16
325	12.60	12.60	91	0.76	77	3.20	4.20	70	1,002.50	1,017.90	12.50	0	SYR	2024-04-18	17

Then, we filtered the Syracuse weather data from the merged weather forecast data of all sources and destinations of the first and second flights. The filtered Syracuse weather data is merged with the first flights data with `left_on=["Date (MM/DD/YYYY)", "Scheduled Arrival hour"]`, `right_on=["Date", "Hour"]`, `how="left"`) so, that it can act as a feature while predicting the first flight arrival status.

```
# Merge merged_df_test_even with the test_merged_weather_df, which has all airports weather info, and map it with the matching origin airport.
merged_df_test_even = pd.merge(merged_df_test_even, test_merged_weather_df, how='left', left_on=['Date (MM/DD/YYYY)', 'Scheduled departure hour', 'Origin Airport'],
                                right_on=['Date', 'Hour', 'Location'])
```

Figure: Merging the flights data with weather data

Then, the merged weather forecast data of all the source and destination locations is merged with the first flights data with `how='left'`, `left_on=['Date (MM/DD/YYYY)', 'Scheduled departure hour', 'Origin Airport']`, `right_on=['Date', 'Hour', 'Location']`) so, that it can act as a feature while predicting the first flight arrival status.

```
merged_df_test_even.drop(['Date_x', 'Hour_x', 'Date_y', 'Hour_y', 'Location_y'], axis=1, inplace=True)

merged_df_test_odd.drop(['Date_x', 'Hour_x', 'Date_y', 'Hour_y', 'Location_y'], axis=1, inplace=True)

#Even rows
all_hours = pd.Series(range(24))

# Convert to categorical variable with specified categories
merged_df_test_even['Scheduled departure hour'] = pd.Categorical(merged_df_test_even['Scheduled departure hour'], categories=all_hours, ordered=True)
merged_df_test_even['Scheduled Arrival hour'] = pd.Categorical(merged_df_test_even['Scheduled Arrival hour'], categories=all_hours, ordered=True)

all_minutes = pd.Series(range(60))

# Convert to categorical variable with specified categories
merged_df_test_even['Scheduled departure minute'] = pd.Categorical(merged_df_test_even['Scheduled departure minute'], categories=all_minutes, ordered=True)
merged_df_test_even['Scheduled Arrival minute'] = pd.Categorical(merged_df_test_even['Scheduled Arrival minute'], categories=all_minutes, ordered=True)

date_column = merged_df_test_even['Date (MM/DD/YYYY)']

# Drop the date column before encoding
merged_df_test_even = pd.get_dummies(merged_df_test_even.drop(columns=['Date (MM/DD/YYYY)']), drop_first=True)

merged_df_test_even.head()
```

Figure: Features extraction and Data conversion

Few of the redundant columns are removed and the hour and minute columns such as “Scheduled departure hour”, “Scheduled departure minute”, “Scheduled arrival hour”, “Scheduled arrival minute” are categorized in ranges 24 and 60 respectively. Then, the column “Date (MM/DD/YYYY)” is dropped, and we create dummy columns with the function `pd.get_dummies`.

```

# Define a dictionary to map numerical labels to their meanings
label_mapping = {
    2: "late",
    1: "early",
    0: "ontime"
}

# Define the possible combinations of values for two columns
column_combinations = [(True, False), (False, True), (False, False)]

# Initialize an empty dictionary to store predictions
predictions_dict = {}

# Iterate over each combination of column values
for comb in column_combinations:
    # Create a copy of the DataFrame
    df_copy = merged_df_test_odd.copy()

    # Set the values of the columns based on the current combination
    df_copy['Arrival Status First_1.0'] = comb[0]
    df_copy['Arrival Status First_2.0'] = comb[1]

```

Figure: Adding earlier flights arrivals status feature to the later flights dataframe.

In the second flights data, a new column “Arrival Status First” is created, to store the arrival status of the first flight 0(on-time), 1(early), 2(late). And the column is categorized into three columns so that the predictions for second flights can be done for all the scenarios where the first flights are on-time/early/late.

```

# Merge merged_df_test_odd with the test_merged_weather_df, which has all airports weather info, and map it with the matching origin airport.
merged_df_test_odd = pd.merge(merged_df_test_odd, test_merged_weather_df, how='left', left_on=['Date (MM/DD/YYYY)', 'Scheduled departure hour', 'Origin Airport'],
                              right_on=['Date', 'Hour', 'Location'])

```

The merged weather forecast data of all the source and destination locations is merged with the second flights data with (how='left', left\_on=['Date (MM/DD/YYYY)', 'Scheduled departure hour', 'Origin Airport'], right\_on=['Date', 'Hour', 'Location']). So, that this can be an additional feature while predicting the arrival status of second flights.

```
merged_df_test_odd.drop(['Date_x', 'Hour_x', 'Date_y', 'Hour_y', 'Location_y'], axis=1, inplace=True)
```

```

# odd rows
all_hours = pd.Series(range(24))

# Convert to categorical variable with specified categories
merged_df_test_odd['Scheduled departure hour'] = pd.Categorical(merged_df_test_odd['Scheduled departure hour'], categories=all_hours, ordered=True)
merged_df_test_odd['Scheduled Arrival hour'] = pd.Categorical(merged_df_test_odd['Scheduled Arrival hour'], categories=all_hours, ordered=True)

all_minutes = pd.Series(range(60))

# Convert to categorical variable with specified categories
merged_df_test_odd['Scheduled departure minute'] = pd.Categorical(merged_df_test_odd['Scheduled departure minute'], categories=all_minutes, ordered=True)
merged_df_test_odd['Scheduled Arrival minute'] = pd.Categorical(merged_df_test_odd['Scheduled Arrival minute'], categories=all_minutes, ordered=True)

date_column = merged_df_test_odd['Date (MM/DD/YYYY)']

# Drop the date column before encoding
merged_df_test_odd = pd.get_dummies(merged_df_test_odd.drop(columns=['Date (MM/DD/YYYY)'], drop_first=True))

merged_df_test_odd.head()

```

Figure: Features extraction and Data conversion

Few of the redundant columns are removed and the hour and minute columns such as “Scheduled departure hour”, “Scheduled departure minute”, “Scheduled arrival hour”, “Scheduled arrival minute” are categorized in ranges 24 and 60 respectively. Then, the column “Date (MM/DD/YYYY)” is dropped and dummy columns are created again with the function `pd.get_dummies()`.

Now that the data has been preprocessed, we pass it to the model to make predictions. For the previous flight predictions, we predict using the first model trained.

```
# Convert the numerical labels into their corresponding meanings
gb_y_train_1_pred_meanings = [label_mapping[label] for label in gb_y_train_1_pred]

# Print the predicted labels with their corresponding meanings
print(gb_y_train_1_pred_meanings)

['early', 'early', 'early', 'early', 'early', 'early', 'early', 'early', 'early', 'early', 'early', 'early', 'early']
```

Figure: Final predictions for the earlier flights

For the second flight, we use the second model trained, which also uses all possible combinations of the previous flights’ arrival statuses.

```
# Make predictions using the model pipeline
predictions = pipeline_2.predict(df_copy)

# Map numerical labels to their corresponding meanings
predictions_meanings = [label_mapping[label] for label in predictions]

# Store the predictions in the dictionary with the combination as the key
predictions_dict[comb] = predictions_meanings

# Access predictions for each combination from the dictionary
for comb, predictions in predictions_dict.items():
    print("Predictions for combination:", comb, ":", predictions)

Predictions for combination: (True, False) : ['early', 'early', 'late', 'early', 'ontime', 'early', 'early', 'early', 'early', 'early', 'late']
Predictions for combination: (False, True) : ['early', 'early', 'ontime', 'early', 'early', 'early', 'early', 'late', 'early', 'early', 'late']
Predictions for combination: (False, False) : ['early', 'early', 'late', 'early', 'early', 'early', 'early', 'late', 'early', 'early', 'late']
```

Combination (True, False) is for when the previous flight is 'Early'.

Combination (False, True) is for when the previous flight is 'Late'.

Combination (False, False) is for when the previous flight is 'On Time'.

Figure: Final predictions for the later flights

The predictions done by the above models for the given flights is as follows:

DATE	DAY	FLIGHT NU	ORIGIN	DEPARTUR	ARRIVAL TI	ARRIVAL STATUS	ARRIVAL STATUS_Prev_flight	ARRIVAL STATUS_Prev_flight_c	ARRIVAL STATUS_Prev_flight_Late		
4/19/24	FRIDAY	UA 1400	ORD	6:52 PM	9:47 PM	Early	NA	NA	NA		Early
4/19/24	FRIDAY	AA 3402	ORD	7:59 PM	10:52 PM	NA	Early	Early	Early		On time
4/19/24	FRIDAY	B6 116	JFK	1:34 PM	2:51 PM	Early	NA	NA	NA		Late
4/19/24	FRIDAY	DL 5182	JFK	2:55 PM	4:21 PM	NA	Early	Early	Early		
4/19/24	FRIDAY	WN 5285	MCO	11:35 AM	2:20 PM	Early	NA	NA	NA		
4/19/24	FRIDAY	B6 656	MCO	1:35 PM	4:25 PM	NA	Late	Late	On time		
4/20/24	SATURDAY	UA 1400	ORD	6:52 PM	9:47 PM	Early	NA	NA	NA		
4/20/24	SATURDAY	AA 3402	ORD	7:59 PM	10:52 PM	NA	Early	Early	Early		
4/20/24	SATURDAY	B6 116	JFK	1:25 PM	2:41 PM	Early	NA	NA	NA		
4/20/24	SATURDAY	DL 5182	JFK	2:55 PM	4:21 PM	NA	On time	Early	Early		
4/20/24	SATURDAY	B6 656	MCO	1:35 PM	4:25 PM	Early	NA	NA	NA		
4/21/24	SUNDAY	UA 1400	ORD	6:52 PM	9:47 PM	Early	NA	NA	NA		
4/21/24	SUNDAY	AA 3402	ORD	7:59 PM	10:52 PM	NA	Early	Early	Early		
4/21/24	SUNDAY	B6 116	JFK	1:35 PM	2:51 PM	Early	NA	NA	NA		
4/21/24	SUNDAY	DL 5182	JFK	2:55 PM	4:21 PM	NA	Early	Early	Early		
4/21/24	SUNDAY	WN 5285	MCO	11:05 AM	1:50 PM	Early	NA	NA	NA		
4/21/24	SUNDAY	B6 656	MCO	1:35 PM	4:25 PM	NA	Early	Late	Late		
4/22/24	MONDAY	UA 1400	ORD	6:52 PM	9:47 PM	Early	NA	NA	NA		
4/22/24	MONDAY	AA 3402	ORD	7:59 PM	10:52 PM	NA	Early	Early	Early		
4/22/24	MONDAY	B6 116	JFK	1:35 PM	2:51 PM	Early	NA	NA	NA		
4/22/24	MONDAY	DL 5182	JFK	2:55 PM	4:21 PM	NA	Early	Early	Early		
4/22/24	MONDAY	WN 5285	MCO	11:35 AM	2:20 PM	Early	NA	NA	NA		
4/22/24	MONDAY	B6 656	MCO	1:34 PM	4:25 PM	NA	Late	Late	Late		

Based on the ground truth provided by the professor, our models were able to predict 14 arrival statuses correctly.

## Conclusion:

In conclusion, the pipeline models trained has proven to be a useful tool for predicting the arrival time of flights into Syracuse (SYR) from four major airports - Chicago (ORD), New York (JFK), and Orlando (MCO)- up to four days in advance. The results of the analysis show that the model was able to accurately predict whether flights would arrive early, on-time, delayed, based on a variety of factors such as weather conditions, flight schedules.

The pipeline models were particularly effective in identifying the most important variables for predicting flight arrival times, allowing for a more accurate and efficient prediction model. The model was able to provide valuable insights to airlines, airports, and travelers, allowing them to plan and adjust their schedules accordingly.

Overall, the use of pipeline models trained using decision trees has the potential to significantly improve the efficiency and reliability of the airline industry, ultimately benefiting both businesses and consumers. With further research and development, this approach could potentially be applied to other areas of transportation and logistics as well.

**REFERENCES:**

- <https://www.transtats.bts.gov/ontime/arrivals.aspx>
- <https://www.weatherbit.io/>
-