

Flight Delay Analysis

```
1 import pandas as pd
2 # Mount google drive
3 from google.colab import drive
4 drive.mount('/content/drive')
5 #Import Libraries that I will use
6 import pandas as pd
7 import numpy as np
8 import matplotlib.pyplot as plt
9 from sklearn.linear_model import LogisticRegression
10 from sklearn.tree import DecisionTreeClassifier
11 from sklearn.ensemble import GradientBoostingClassifier
12 from sklearn.model_selection import train_test_split, StratifiedKFold
13 from sklearn.preprocessing import OneHotEncoder, OrdinalEncoder, StandardScaler
14 from sklearn.pipeline import Pipeline
15 from sklearn.metrics import accuracy_score
16 import matplotlib.pyplot as plt
17 import seaborn as sns # For additional visualizations
18 import gc # for garbage collection the colab ran over RAM in my runtime gc.collect()
19 # Import excel files into df from google drive
20 file_path = '/content/drive/My Drive/Dataset/'
21
22 # **1. Importing Data**
23 df=None
24 airlines_df = pd.read_excel(file_path+'Airlines.xlsx')
25 airports_df = pd.read_excel(file_path+'airports.xlsx')
26 runways_df = pd.read_excel(file_path+'runways.xlsx')
27 # Check import is correct
28 # Investigate data structures (modify column names if needed)
29 print(airlines_df.info())
30 print(airlines_df.shape)
31 print(airlines_df.head())
32 print(airports_df.info())
33 print(airports_df.shape)
34 print(airports_df.head())
35 print(runways_df.info())
36 print(runways_df.shape)
37 print(runways_df.head())
38
39 # dropping unnecesary columns
40 runways_df = runways_df.drop(['le_ident', 'le_latitude_deg','le_longitude_deg', 'le_elevation_ft', 'le_heading_degT',
41 'le_displaced_threshold_ft', 'he_ident', 'he_latitude_deg','he_longitude_deg', 'he_elevation_ft', 'he_heading_degT',
42 'he_displaced_threshold_ft'], axis = 1)
43
44 airports_df =airports_df.drop(['continent', 'iso_country', 'iso_region','municipality', 'gps_code','local_code', 'home_link',
45 'wikipedia_link', 'keywords'], axis=1)
46
47 # Join 1: merge the runways and airport data.
48 df = pd.merge(airports_df, runways_df,how='left', left_on = "ident", right_on = "airport_ident")
49 df.head()
50 # Derive 'number_of_runways' Feature Engineering
51 df['number_of_runways'] = df.groupby('ident')['id_y'].transform('count')
52 df.isna().sum()
53 df.shape
54 df.drop(['id_x','id_y','surface'], axis=1, inplace=True)
55 # 2. Delete rows where 'Type' == 'closed'
56 df = df[df['type'] != 'closed']
57
58 # 3. Delete rows with null 'iata_code'
59 df = df[df['iata_code'].notnull()]
60 # drop nulls in certain columns:
61 df['elevation_ft'] = df['elevation_ft'].fillna(0) # Replace nulls in 'elevation_ft' with 0
62 df['length_ft'] = df['length_ft'].fillna(df['length_ft'].mean()) # Impute with mean
63 df['width_ft'] = df['width_ft'].fillna(df['length_ft'].mean()) # Impute with mean
64 # Impute nulls in 'lighted' column with the most frequent value
65 df['lighted'] = df['lighted'].fillna(df['lighted'].mode()[0])
66 # Now drop the nulls
67 df.dropna(inplace=True)
68 df.shape
```

```

69 df.info()
70 df['iata_code'].nunique()
71 df.drop_duplicates(subset='iata_code', inplace=True)
72
73 #Join 2: join Airlines on df airportfrm<>iata code
74 df.reset_index(inplace=True) # rest indices
75 airlines_df.reset_index(inplace=True)
76 airlines_df.info()
77 df_airrun=df.copy()
78 df = airlines_df.merge(df, how='left', left_on='AirportFrom', right_on='iata_code')
79 df.shape
80 df.columns
81
82 #merged_df.rename(columns={'id': 'airport_id', 'type': 'airport_type'}, inplace=True) # Rename for clarity
83 df=df.drop(['index_x', 'id', 'Flight',
84             'index_y', 'ident', 'name', 'latitude_deg', 'longitude_deg',
85             'scheduled_service', 'airport_ref',
86             'airport_ident' ],axis=1)
87 # make a temp df
88 df_airrun.columns
89 df_airrun=df_airrun.drop(['index', 'ident', 'type', 'name', 'latitude_deg', 'longitude_deg',
90             'scheduled_service', 'airport_ref',
91             'closed',
92             ],axis=1)
93 df_airrun.drop_duplicates(subset='iata_code', inplace=True)
94 df_airrun.isna().sum()
95 df_airrun.info()
96 # Join 3: airports_df <-> airlines_df (on arrival airport)
97 df = df.merge(df_airrun, how='left', left_on='AirportTo', right_on='iata_code', suffixes=('_from', '_to'))
98 df.isna().sum()
99 # Now drop the nulls
100 df.dropna(inplace=True)
101 df.info()
102 df.columns
103 df=df.drop(['closed'],axis=1)
104 # Upload to excel as pgm crashing
105 df.to_csv(file_path + 'df_output_air.csv')
106
107
108 # Columns to convert to integers
109 float_to_int_cols = ['DayOfWeek', 'Time', 'Length',
110                     'Delay', 'elevation_ft_from',
111                     'length_ft_from', 'width_ft_from', 'lighted_from',
112                     'number_of_runways_from', 'elevation_ft_to',
113                     'length_ft_to', 'width_ft_to', 'lighted_to',
114                     'number_of_runways_to']
115
116 # Columns to convert to strings
117 obj_to_str_cols = ['Airline', 'AirportFrom', 'AirportTo', 'iata_code_from',
118                     'iata_code_to']
119
120 # Apply transformations
121 for col in float_to_int_cols:
122     df[col] = df[col].astype(int)
123
124 for col in obj_to_str_cols:
125     df[col] = df[col].astype(str)
126 # Check Data
127 df.isnull().sum()
128 df.info()
129 df.head()
130 df.shape
131 df['Airline'].unique()
132 # Assuming your DataFrame is named "df"
133 count_co_records = df[df['Airline'] == 'CO'].shape[0]
134 print(f"There are {count_co_records} records with the 'Airline' value as 'CO'.")
135

```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 518556 entries, 0 to 518555
Data columns (total 9 columns):
 # Column Non-Null Count Dtype
--- -- ----- ----- -----

```

0   id          518556 non-null  int64
1   Airline      518556 non-null  object
2   Flight       518556 non-null  int64
3   AirportFrom  518556 non-null  object
4   AirportTo    518556 non-null  object
5   DayOfWeek    518556 non-null  int64
6   Time         518556 non-null  int64
7   Length       518556 non-null  int64
8   Delay        518556 non-null  int64
dtypes: int64(6), object(3)
memory usage: 35.6+ MB
None
(518556, 9)
   id Airline  Flight AirportFrom AirportTo DayOfWeek  Time  Length  Delay
0   1     CO     269       SFO      IAH        3    15    205     1
1   2     US    1558       PHX      CLT        3    15    222     1
2   3     AA    2400       LAX      DFW        3    20    165     1
3   4     AA    2466       SFO      DFW        3    20    195     1
4   5     AS    108       ANC      SEA        3    30    202     0
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 73805 entries, 0 to 73804
Data columns (total 18 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   id               73805 non-null   int64  
 1   ident             73805 non-null   object  
 2   type              73805 non-null   object  
 3   name              73805 non-null   object  
 4   latitude_deg     73805 non-null   float64
 5   longitude_deg    73805 non-null   float64
 6   elevation_ft     59683 non-null   float64
 7   continent         38086 non-null   object  
 8   iso_country       73546 non-null   object  
 9   iso_region        73805 non-null   object  
 10  municipality      68739 non-null   object  
 11  scheduled_service 73805 non-null   object  
 12  gps_code          42996 non-null   object  
 13  iata_code          9160 non-null   object  
 14  local_code         32975 non-null   object  
 15  home_link          3492 non-null   object  
 16  wikipedia_link    10705 non-null   object  
 17  keywords            13951 non-null   object  
dtypes: float64(3), int64(1), object(14)
memory usage: 10.1+ MB
None
(73805, 18)
   id ident      type          name  latitude_deg  \
0   4983  NTGA  medium_airport  Anaa Airport  -17.352600
1   27342 YARY  small_airport  Arrabury Airport  -26.696390
2   3176  HEAR  medium_airport  El Arish International Airport  31.078565
3   324756 AAD   small_airport  Adado Airport   6.095802
4   2060  DABB  medium airprot  Annaba Rabah Bitat Airport  36.826781

```

Double-click (or enter) to edit

▼ Web scraping

Airline experience feature engg and extraction using beautiful Soup fir web scrapping.

```

1 import requests
2 from bs4 import BeautifulSoup
3
4 def get_airline_experience(airline_name, url):
5     response = requests.get(url)
6     soup = BeautifulSoup(response.content, 'html.parser')
7
8     # Find all relevant tables
9     info_tables = soup.find_all('table', class_='wikitable')
10
11    if info_tables:
12        for table in info_tables:
13            # Search for the airline row within each table
14            for row in table.find_all('tr'):
15                cells = row.find_all('td')
16                if cells and airline_name in cells[2].text: #index 2 holds IATA code
17                    founded_year_str = cells[6].text.strip() # this index holds Founded open and check url
18
19                    try:
20                        founded_year = int(founded_year_str)
21                        return 2024 - founded_year # Return as soon as a match is found
22                    except ValueError:
23                        print('Founding year format error for:', airline_name)
24
25                # If no match is found in any of the tables:
26                print('Error: Airline not found', airline_name)
27                return None
28        else:
29            print('Error: No info tables found')
30            return None
31
32 # Set URL of website
33 airline_list_url = "https://en.wikipedia.org/wiki/List_of_airlines_of_the_United_States"
34 airline_experience = {}
35
36 # Assuming you have a DataFrame called 'airlines_df' with an 'airline_name' column
37 for airline in airlines_df['Airline'].unique():
38     experience = get_airline_experience(airline, airline_list_url)
39     airline_experience[airline] = experience
40
41 # Add airline experience to the merged dataset
42 df['airline_experience'] = df['Airline'].map(airline_experience)
43 # Checking and testing
44 df.tail()
45 df['airline_experience'].isna().sum()
46 # Find out airlines for which data na in website
47 grouped_df = df[df['airline_experience'].isna()].groupby('Airline').size().reset_index(name='Count')
48
49 # Print the websites where data NA
50 print(grouped_df)
51
52 # Assigning values manually as deleting these will be huge loss of data
53 df.loc[df['Airline'] == 'CO', 'airline_experience'] = 87
54 df.loc[df['Airline'] == 'US', 'airline_experience'] = 87
55 df.loc[df['Airline'] == 'EV', 'airline_experience'] = 37
56 df.isna().sum()
57 df.info()
58 df.shape
59 # df=merged_df.copy() #This copies df the original one is unchanged
60 df.dropna(inplace=True)
61 df.isnull().sum()
62 df.shape
63 df.info()

```

```

Error: Airline not found      CO
Error: Airline not found      US
Error: Airline not found      EV
   Airline  Count
0      CO  21118
1      EV  27983
2      US  34500
<class 'pandas.core.frame.DataFrame'>
Int64Index: 518494 entries, 0 to 518493
Data columns (total 25 columns):
 #   Column           Non-Null Count  Dtype

```

```
--- -----  
0 Airline      518494 non-null object  
1 AirportFrom  518494 non-null object  
2 AirportTo    518494 non-null object  
3 DayOfWeek   518494 non-null int64  
4 Time        518494 non-null int64  
5 Length      518494 non-null int64  
6 Delay       518494 non-null int64  
7 type_from   518494 non-null object  
8 elevation_ft_from 518494 non-null int64  
9 iata_code_from 518494 non-null object  
10 length_ft_from 518494 non-null int64  
11 width_ft_from 518494 non-null int64  
12 lighted_from 518494 non-null int64  
13 number_of_runways_from 518494 non-null int64  
14 elevation_ft_to 518494 non-null int64  
15 iata_code_to 518494 non-null object  
16 airport_ident 518494 non-null object  
17 length_ft_to 518494 non-null int64  
18 width_ft_to 518494 non-null int64  
19 lighted_to   518494 non-null int64  
20 number_of_runways_to 518494 non-null int64  
21 type_to     518494 non-null object  
22 elevation_ft 518494 non-null float64  
23 iata_code   518494 non-null object  
24 airline_experience 518494 non-null float64  
dtypes: float64(2), int64(14), object(9)  
memory usage: 102.9+ MB  
<class 'pandas.core.frame.DataFrame'>  
Int64Index: 518494 entries, 0 to 518493  
Data columns (total 25 columns):  
 #  Column          Non-Null Count  Dtype    
---  --  
0   Airline         518494 non-null  object  
1   AirportFrom    518494 non-null  object  
2   AirportTo      518494 non-null  object  
3   DayOfWeek      518494 non-null  int64  
4   Time           518494 non-null  int64  
5   Length          518494 non-null  int64  
6   Delay           518494 non-null  int64  
7   type_from      518494 non-null  object  
8   elevation_ft_from 518494 non-null  int64  
9   iata_code_from 518494 non-null  object  
10  length_ft_from 518494 non-null  int64  
11  width_ft_from 518494 non-null  int64  
12  lighted_from   518494 non-null  int64  
13  number_of_runways_from 518494 non-null  int64
```

Passenger data departure airport from url

```

1 # Define function to process passenger info
2 def get_airport_tfc(airport_name, url):
3     response = requests.get(url)
4     soup = BeautifulSoup(response.content, 'html.parser')
5
6     # Find all relevant tables
7     info_tables = soup.find_all('table', class_='wikitable')
8
9     if info_tables:
10         for table in info_tables:
11             # Search for the airline row within each table
12             for row in table.find_all('tr'):
13                 cells = row.find_all('td')
14                 if cells and airport_name in cells[2].text: #index 2 holds IATA code
15                     passenger_tfc_str = cells[6].text.strip() # this index holds Founded open and check url
16
17                 try:
18                     passenger_tfc = int(passenger_tfc_str.replace(',', '')) # Convert to integer
19                     return passenger_tfc # Return as soon as a match is found
20                 except ValueError:
21                     print('Passenger tfc format error for:', airport_name)
22
23     # If no match is found in any of the tables:
24     print('Error: Passenger data not found', ' ', airport_name)
25     return None
26 else:
27     print('Error: No info tables found')
28     return None
29
30 # Set URL of website
31 airline_list_url = "https://en.wikipedia.org/wiki/List_of_the_busiest_airports_in_the_United_States"
32 airport_tfc = {}
33
34 # Assuming you have a DataFrame called 'airlines_df' with an 'airline_name' column
35 for airport in df['AirportFrom'].unique():
36     tfc = get_airport_tfc(airport, airline_list_url)
37     airport_tfc[airport] = tfc
38
39 notnull_count = sum(1 for value in airport_tfc.values() if value is not None)
40 list_notnull = [key for key, value in airport_tfc.items() if value is not None]
41 # Add airline passenger to the merged dataset departure
42 df['passenger_tfc_from'] = df['AirportFrom'].map(airport_tfc)
43 df.tail()
44 df['passenger_tfc_from'].isnull().sum()
45 df['passenger_tfc_from'].unique()
46 # Find out airlines for which data na in website
47 grouped_passenger_df = df[df['passenger_tfc_from'].isna()].groupby('AirportFrom').size().reset_index(name='Count')
48 #df['passenger_tfc_from'].fillna(df['passenger_tfc_from'].min(), inplace=True)
49 # Print the websites where data NA
50 print(grouped_passenger_df)
51
52 # Consider dropping- data too large so will make cat variable
53 #df.dropna(inplace=True)
54 df.isnull().sum()
55 # I have 2 lakh nulls in the passenger field. As no data from website
56 # Will make a categorical variable hub to classify airport to large, medium, small
57 # Categorize airports as large/medium hubs based on passenger traffic
58 df['hub_type_from'] = np.where(df['passenger_tfc_from'] >= 8000000, 'large',
59                                 np.where(df['passenger_tfc_from'] >= 2000000, 'medium', 'small'))
60
61 # Trouble shoot immute later
62 #df['passenger_tfc_from'].fillna(df['passenger_tfc_from'].min(), inplace=True)
63 # Find outhub summary
64 grouped_hub_df = df[df['hub_type_from'].isna()].groupby('AirportFrom').size().reset_index(name='Count')
65
66 # Print the websites where data NA
67 print(grouped_hub_df)
68 # Count the occurrences of each hub type
69 hub_type_counts_from = df['hub_type_from'].value_counts()
70
71 # Print the counts
72 print(hub_type_counts_from)
73 df['hub_type_from'].unique()

```

```
74 df.shape
75 df.info()
76 df['Airline'].dtype
77 df['Delay'].dtype
78 df.head()
79 print(df['Delay'].unique())
80 df['Delay'] = df['Delay'].astype(bool)
```

| | |
|---------------------------------|-----|
| Error: Passenger data not found | FAI |
| Error: Passenger data not found | BQN |
| Error: Passenger data not found | PSE |
| Error: Passenger data not found | BIS |
| Error: Passenger data not found | IYK |
| Error: Passenger data not found | GFK |
| Error: Passenger data not found | GSO |
| Error: Passenger data not found | LMT |
| Error: Passenger data not found | DLH |
| Error: Passenger data not found | FAR |
| Error: Passenger data not found | MFE |
| Error: Passenger data not found | VPS |
| Error: Passenger data not found | MAF |
| Error: Passenger data not found | LWS |
| Error: Passenger data not found | RST |
| Error: Passenger data not found | ALB |
| Error: Passenger data not found | DSM |
| Error: Passenger data not found | MSN |
| Error: Passenger data not found | PNS |
| Error: Passenger data not found | BHM |
| Error: Passenger data not found | LIT |
| Error: Passenger data not found | SAV |
| Error: Passenger data not found | ICT |
| Error: Passenger data not found | ECP |
| Error: Passenger data not found | DHN |
| Error: Passenger data not found | MGM |
| Error: Passenger data not found | CAE |
| Error: Passenger data not found | PWM |
| Error: Passenger data not found | ACV |
| Error: Passenger data not found | EKO |
| Error: Passenger data not found | RIC |
| Error: Passenger data not found | BTR |
| Error: Passenger data not found | HRL |
| Error: Passenger data not found | MYR |
| Error: Passenger data not found | TUS |
| Error: Passenger data not found | SBN |
| Error: Passenger data not found | CAK |
| Error: Passenger data not found | TVC |
| Error: Passenger data not found | DAY |
| Error: Passenger data not found | MFR |
| Error: Passenger data not found | BTV |
| Error: Passenger data not found | TLH |
| Error: Passenger data not found | TYS |
| Error: Passenger data not found | CHA |
| Error: Passenger data not found | LRD |
| Error: Passenger data not found | BRO |
| Error: Passenger data not found | CRP |
| Error: Passenger data not found | LAN |
| Error: Passenger data not found | PVD |
| Error: Passenger data not found | FWA |
| Error: Passenger data not found | OKC |
| Error: Passenger data not found | ORF |
| Error: Passenger data not found | AEX |
| Error: Passenger data not found | SYR |
| Error: Passenger data not found | SHV |
| Error: Passenger data not found | VLD |
| Error: Passenger data not found | FAT |
| Error: Passenger data not found | BZN |

passenger data arr airport

```

1 # Add airline passenger to the merged dataset departure
2 df['passenger_tfc_to'] = df['AirportTo'].map(airport_tfc)
3 df.tail()
4 df['passenger_tfc_to'].isnull().sum()
5 df['passenger_tfc_to'].unique()
6 # Find out airlines for which data na in website
7 grouped_passenger_to_df = df[df['passenger_tfc_to'].isna()].groupby('AirportTo').size().reset_index(name='Count')
8
9 # Print the websites where data NA
10 print(grouped_passenger_to_df)
11
12 # Consider dropping- data too large so will make cat variable
13 #df.dropna(inplace=True)
14 df.isnull().sum()
15 # I have 2 lakh nulls in the passenger field. As no data from website
16 # Will make a categorical variable hub to classify airport to large, medium, small
17 # Categorize airports as large/medium hubs based on passenger traffic
18 df['hub_type_to'] = np.where(df['passenger_tfc_to'] >= 8000000, 'large',
19                               np.where(df['passenger_tfc_to'] >= 2000000, 'medium', 'small'))
20 # Trouble shoot will immute later
21 #df['passenger_tfc_to'].fillna(df['passenger_tfc_to'].min(), inplace=True)
22 # Find outhub summary
23 grouped_hub_to_df = df[df['hub_type_to'].isna()].groupby('AirportTo').size().reset_index(name='Count')
24
25 # Print the websites where data NA
26 print(grouped_hub_to_df)
27 # Count the occurrences of each hub type
28 hub_type_counts_to = df['hub_type_to'].value_counts()
29
30 # Print the counts
31 print(hub_type_counts_to)

```

| AirportTo | Count | |
|-----------|-------|------|
| 0 | ABE | 257 |
| 1 | ABI | 211 |
| 2 | ABR | 2 |
| 3 | ABY | 87 |
| 4 | ACT | 59 |
| .. | ... | ... |
| 221 | VPS | 633 |
| 222 | WRG | 59 |
| 223 | XNA | 1123 |
| 224 | YAK | 57 |
| 225 | YUM | 338 |

[226 rows x 2 columns]
Empty DataFrame
Columns: [AirportTo, Count]
Index: []
large 334086
medium 99490
small 84918
Name: hub_type_to, dtype: int64

✓ Upload webscrapped file

This is to keeo copy of webscraped df

```

1 # Upload to csv for continuing proj from here
2 # Only passenger from and to field has na values
3 df.to_csv(file_path + 'df_wescraped_air.csv')
4 print(df.info())
5 df.shape
6 # The shape of the df is similar to Airline table. If majorly different then the joins are incorrect.

<class 'pandas.core.frame.DataFrame'>
Int64Index: 518494 entries, 0 to 518493
Data columns (total 29 columns):
 #   Column            Non-Null Count  Dtype  
--- 
 0   Airline           518494 non-null   object 
 1   AirportFrom       518494 non-null   object 
 2   AirportTo         518494 non-null   object 
 3   DayOfWeek         518494 non-null   int64  
 4   Time              518494 non-null   int64  

```

```

5 Length 518494 non-null int64
6 Delay 518494 non-null bool
7 type_from 518494 non-null object
8 elevation_ft_from 518494 non-null int64
9 iata_code_from 518494 non-null object
10 length_ft_from 518494 non-null int64
11 width_ft_from 518494 non-null int64
12 lighted_from 518494 non-null int64
13 number_of_runways_from 518494 non-null int64
14 elevation_ft_to 518494 non-null int64
15 iata_code_to 518494 non-null object
16 airport_ident 518494 non-null object
17 length_ft_to 518494 non-null int64
18 width_ft_to 518494 non-null int64
19 lighted_to 518494 non-null int64
20 number_of_runways_to 518494 non-null int64
21 type_to 518494 non-null object
22 elevation_ft 518494 non-null float64
23 iata_code 518494 non-null object
24 airline_experience 518494 non-null float64
25 passenger_tfc_from 433524 non-null float64
26 hub_type_from 518494 non-null object
27 passenger_tfc_to 433576 non-null float64
28 hub_type_to 518494 non-null object
dtypes: bool(1), float64(4), int64(13), object(11)
memory usage: 115.2+ MB
None

```

Load Webscraped file from Gdrive. Midcourse start as Web scraping is time consuming.

```

1 df= None
2
3 # read file into df
4 df=pd.read_csv(file_path+'df_wescraped_air.csv')
5 print(df.info())

```

| # | Column | Non-Null Count | Dtype |
|----|------------------------|----------------|------------------|
| 0 | Unnamed: 0 | 518494 | int64 |
| 1 | Airline | 518494 | object |
| 2 | AirportFrom | 518494 | non-null object |
| 3 | AirportTo | 518494 | non-null object |
| 4 | DayOfWeek | 518494 | non-null int64 |
| 5 | Time | 518494 | non-null int64 |
| 6 | Length | 518494 | non-null int64 |
| 7 | Delay | 518494 | non-null bool |
| 8 | type_from | 518494 | non-null object |
| 9 | elevation_ft_from | 518494 | non-null int64 |
| 10 | iata_code_from | 518494 | non-null object |
| 11 | length_ft_from | 518494 | non-null int64 |
| 12 | width_ft_from | 518494 | non-null int64 |
| 13 | lighted_from | 518494 | non-null int64 |
| 14 | number_of_runways_from | 518494 | non-null int64 |
| 15 | elevation_ft_to | 518494 | non-null int64 |
| 16 | iata_code_to | 518494 | non-null object |
| 17 | airport_ident | 518494 | non-null object |
| 18 | length_ft_to | 518494 | non-null int64 |
| 19 | width_ft_to | 518494 | non-null int64 |
| 20 | lighted_to | 518494 | non-null int64 |
| 21 | number_of_runways_to | 518494 | non-null int64 |
| 22 | type_to | 518494 | non-null object |
| 23 | elevation_ft | 518494 | non-null float64 |
| 24 | iata_code | 518494 | non-null object |
| 25 | airline_experience | 518494 | non-null float64 |
| 26 | passenger_tfc_from | 433524 | non-null float64 |
| 27 | hub_type_from | 518494 | non-null object |
| 28 | passenger_tfc_to | 433576 | non-null float64 |
| 29 | hub_type_to | 518494 | non-null object |

```

dtypes: bool(1), float64(4), int64(14), object(11)
memory usage: 115.2+ MB
None

```

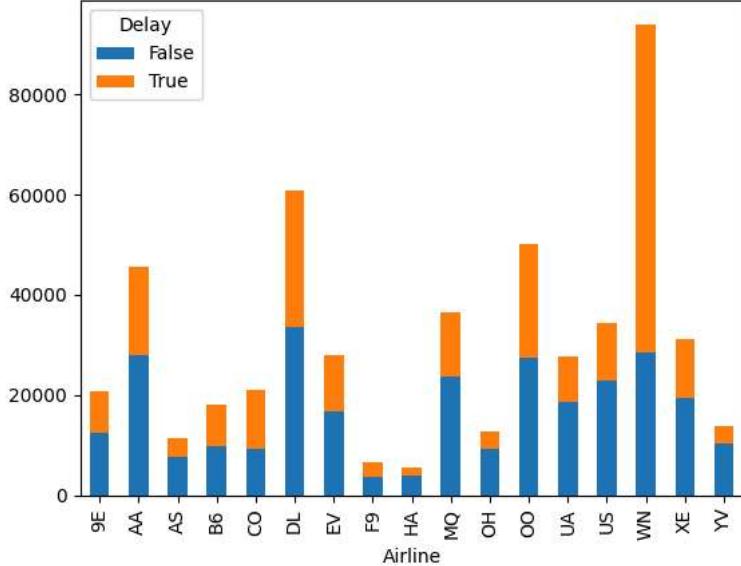
Visualisation

```

1
2 # Reset the index to ensure it starts from 0 and is consecutive
3 df.reset_index(drop=True, inplace=True)
4
5 # Add a new column named 'ID'
6 df['ID'] = df.index + 1 # Index starts from 0, so we add 1
7
8 df_pivot = pd.pivot_table(df, index='Airline', columns='Delay', values='ID', aggfunc='count').fillna(0)
9 plt.figure(figsize=(10,7))
10 df_pivot.plot.bar(stacked=True)
11 plt.show()
12

```

<Figure size 1000x700 with 0 Axes>

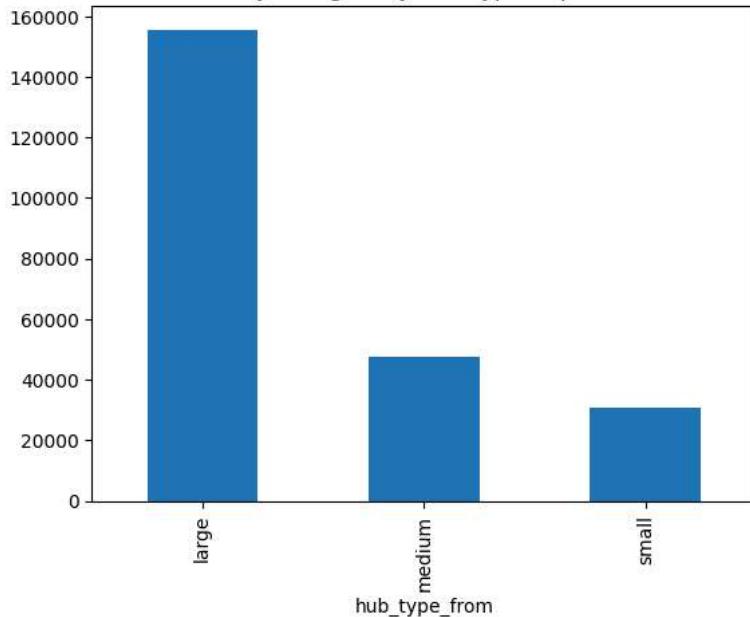


```

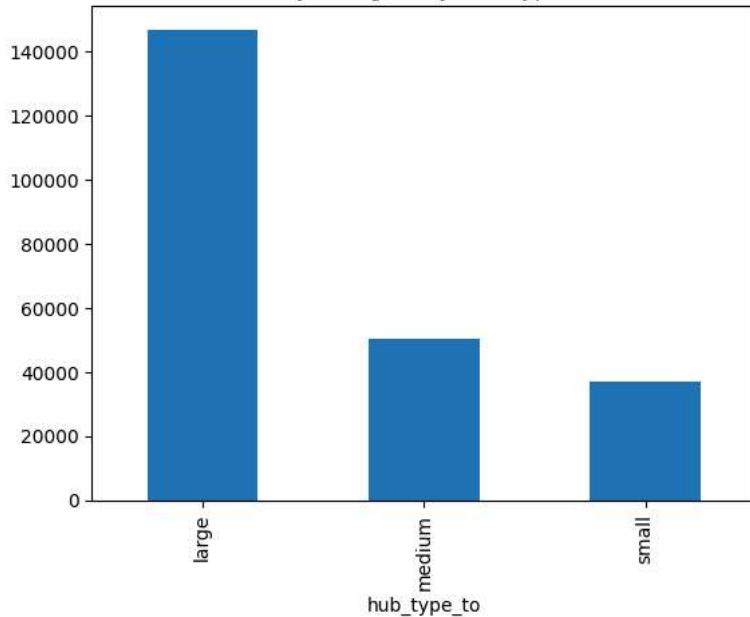
1 # Visualization: Delayed flights by hub type departure
2 delayed_by_hub_from = df[df['Delay'] == True].groupby('hub_type_from').size().plot(kind='bar')
3 plt.title('Delayed Flights by Hub Type Departure')
4 plt.show()
5 #Plot dealy vs Arr airfd hub type
6 delayed_by_hub_to = df[df['Delay'] == True].groupby('hub_type_to').size().plot(kind='bar')
7 plt.title('Delayed Flights by Hub Type Arr')
8 plt.show()

```

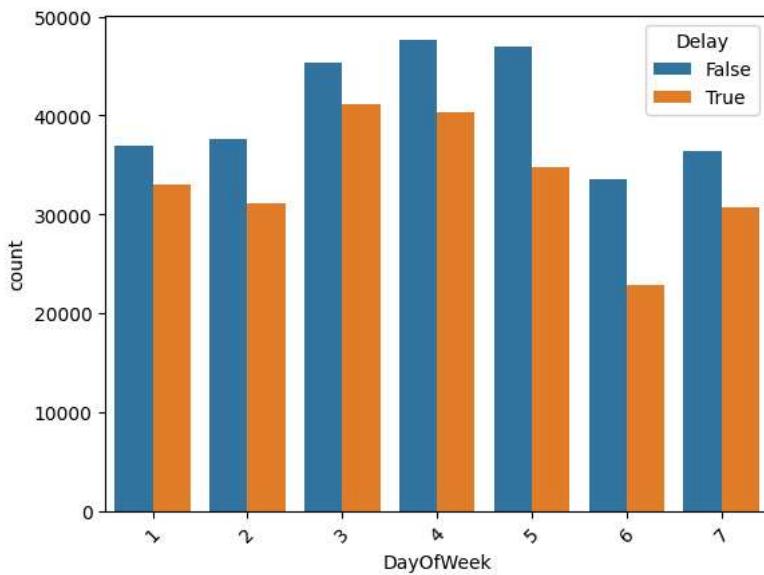
Delayed Flights by Hub Type Departure



Delayed Flights by Hub Type Arr



```
1 # Create a count plot of DayOfWeek vs Delay
2 sns.countplot(
3     x = 'DayOfWeek',
4     hue = 'Delay',
5     data = df
6 )
7
8 # Rotate x labels to prevent overlapping
9 plt.xticks(rotation=45)
10 plt.show()
```

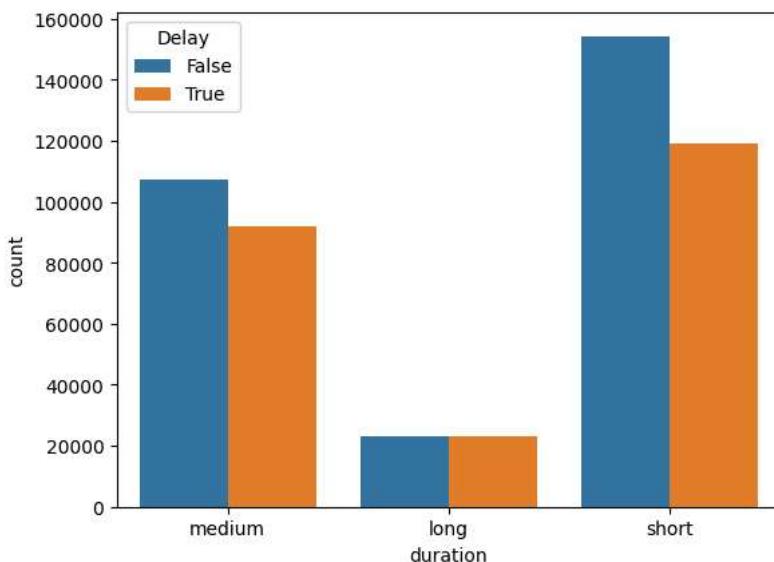


Feature engineering

```

1 # Create the 'duration' field
2 def categorize_duration(row):
3     if row['Length'] < 120:
4         return 'short'
5     elif 120 <= row['Length'] < 240:
6         return 'medium'
7     else:
8         return 'long'
9
10 df['duration'] = df.apply(categorize_duration, axis=1)
11
12 # Plot duration vs. delay
13 sns.countplot(x='duration', hue='Delay', data=df)
14 plt.show()
15
16

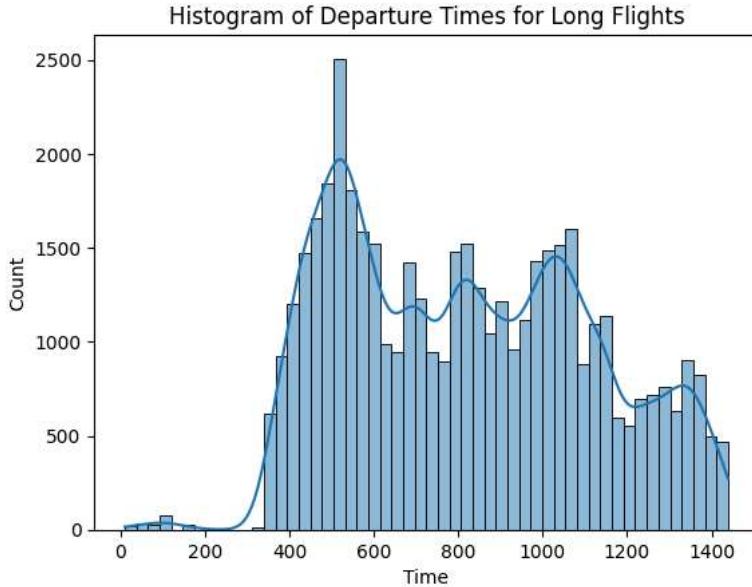
```



```

1 # Filter for long duration flights
2 long_flights = df[df['duration'] == 'long']
3
4 # Create the histogram
5 sns.histplot(long_flights['Time'], kde=True) # kde=True adds a density curve
6 plt.xlabel('Time')
7 plt.ylabel('Count')
8 plt.title('Histogram of Departure Times for Long Flights')
9 plt.show()

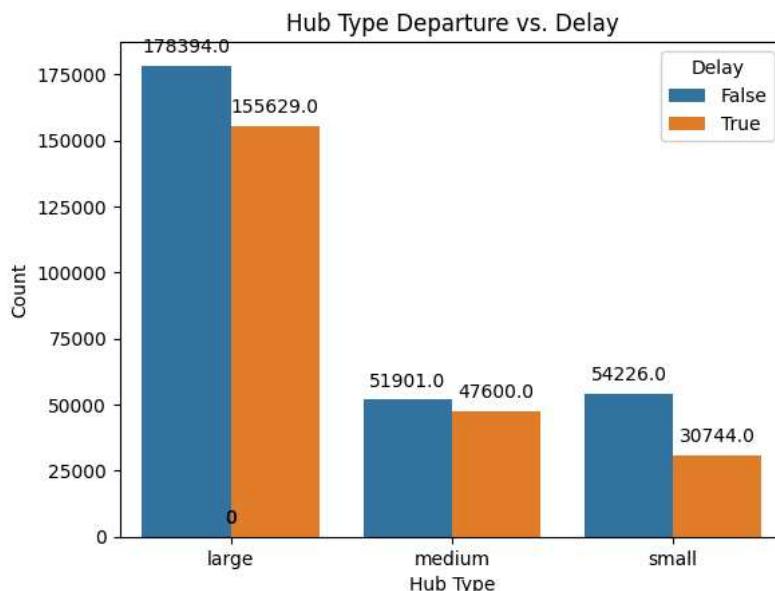
```



```

1 # Plot hub vs. delay
2 sns.countplot(x='hub_type_from', hue='Delay', data=df)
3 # Get the current axes object
4 ax = plt.gca()
5 # Add value labels
6 for p in ax.patches:
7     height = p.get_height()
8     ax.annotate(f'{height}', xy=(p.get_x() + p.get_width() / 2, height),
9                 xytext=(0, 5), # Offset slightly above the bar
10                textcoords='offset points',
11                ha='center', va='bottom')
12
13 # Customize the plot
14 plt.xlabel('Hub Type')
15 plt.ylabel('Count')
16 plt.title('Hub Type Departure vs. Delay')
17 plt.legend(title='Delay', loc='upper right')
18 plt.show()
19

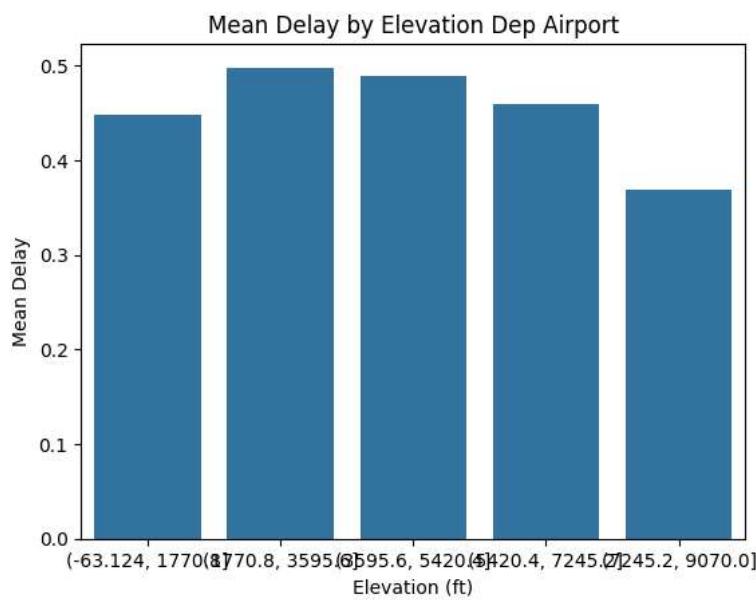
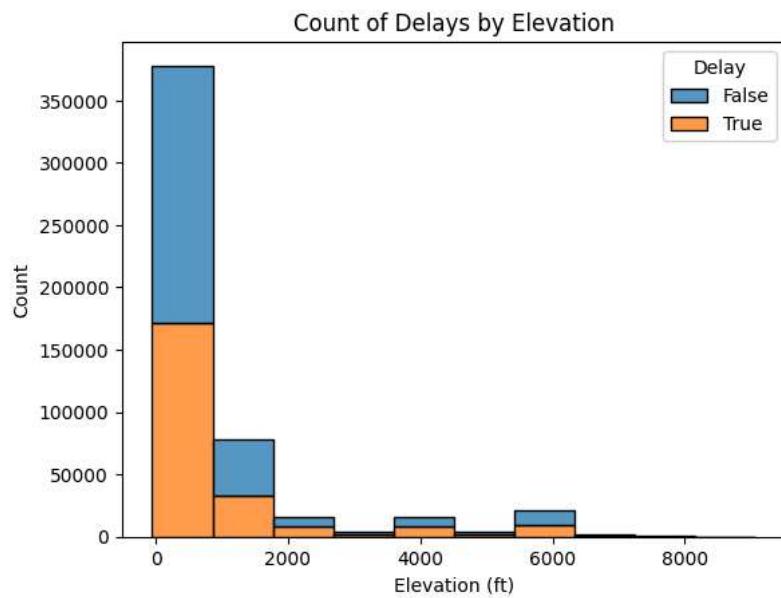
```



```

1 # Define the number of bins (adjust as needed)
2 num_bins = 10
3
4 # Create the histogram
5 sns.histplot(
6     df,
7     x = 'elevation_ft_from',
8     hue = 'Delay', # Separate histograms for delays vs non-delays
9     bins = num_bins,
10    multiple = 'stack' # Stack the histograms on top of each other
11 )
12 plt.xlabel('Elevation (ft)')
13 plt.ylabel('Count')
14 plt.title('Count of Delays by Elevation')
15 plt.show()
16
17 # Calculate grouped means
18 # Define the number of bins (adjust as needed)
19 num_bins = 5
20
21 df_grouped = df.groupby(pd.cut(df['elevation_ft_from'], bins=num_bins))['Delay'].mean()
22
23 # Create the bar plot
24 sns.barplot(
25     x = df_grouped.index.astype(str), # Convert bin labels to string for plotting
26     y = df_grouped.values
27 )
28 plt.xlabel('Elevation (ft)')
29 plt.ylabel('Mean Delay')
30 plt.title('Mean Delay by Elevation Dep Airport')
31 plt.figure(figsize=(12, 20))
32 plt.show()
33 # Hypothesis testing (example: effect of airport elevation on delays)
34 from scipy.stats import ttest_ind
35
36 high_elev = df[df['elevation_ft_from'] > 5000]
37 low_elev = df[df['elevation_ft_from'] <= 5000]
38
39 t_stat, p_value = ttest_ind(high_elev['Delay'], low_elev['Delay'], equal_var=False)
40
41 print("t-statistic:", t_stat)
42 print("p-value:", p_value)
43

```



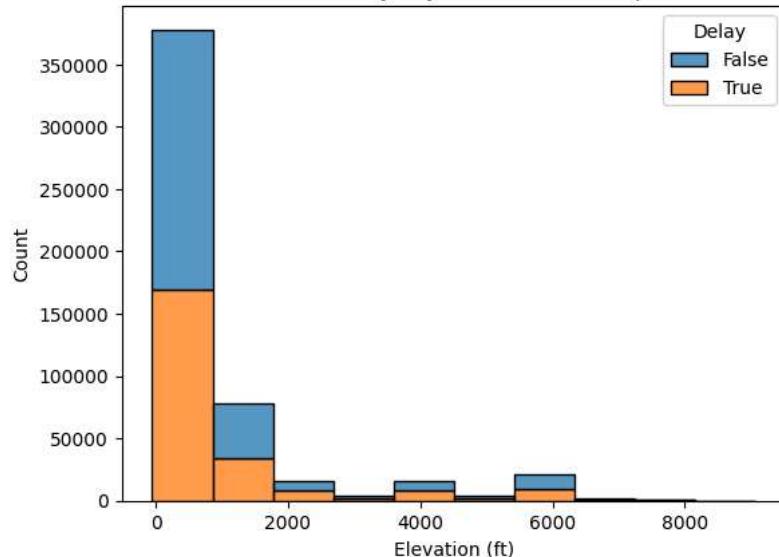
<Figure size 1200x2000 with 0 Axes>

t-statistic: 2.415435752905638

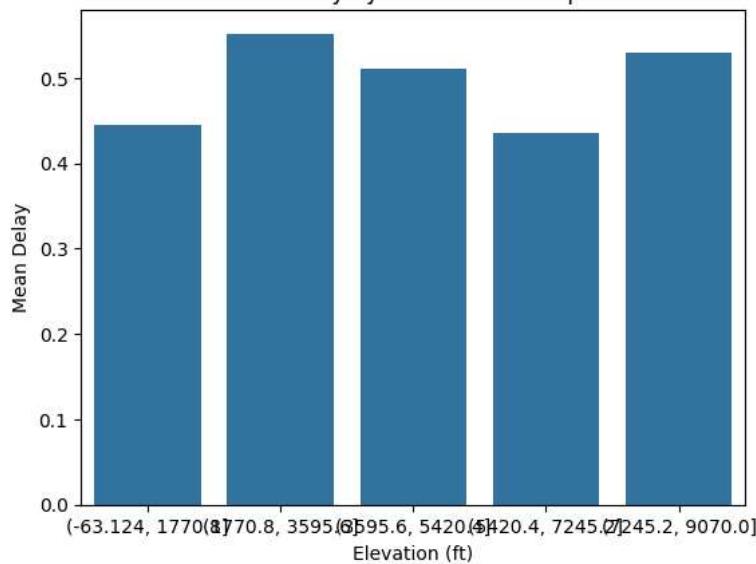
p-value: 0.015722462914667622

```
1
2 # Define the number of bins (adjust as needed)
3 num_bins = 10
4
5 # Create the histogram
6 sns.histplot(
7     df,
8     x = 'elevation_ft_to',
9     hue = 'Delay', # Separate histograms for delays vs non-delays
10    bins = num_bins,
11    multiple = 'stack' # Stack the histograms on top of each other
12 )
13 plt.xlabel('Elevation (ft)')
14 plt.ylabel('Count')
15 plt.title('Count of Delays by Elevation: Arr Airport')
16 plt.show()
17
18 # Calculate grouped means
19 # Define the number of bins (adjust as needed)
20 num_bins = 5
21
22 df_grouped = df.groupby(pd.cut(df['elevation_ft_to'], bins=num_bins))['Delay'].mean()
23
24 # Create the bar plot (using barplot for flexibility)
25 sns.barplot(
26     x = df_grouped.index.astype(str), # Convert bin labels to string for plotting
27     y = df_grouped.values
28 )
29 plt.xlabel('Elevation (ft)')
30 plt.ylabel('Mean Delay')
31 plt.title('Mean Delay by Elevation Arr Airport')
32 plt.figure(figsize=(12, 20))
33 plt.show()
34 # Hypothesis testing (example: effect of airport elevation on delays)
35 from scipy.stats import ttest_ind
36
37 high_elev = df[df['elevation_ft_to'] > 5000]
38 low_elev = df[df['elevation_ft_to'] <= 5000]
39
40 t_stat, p_value = ttest_ind(high_elev['Delay'], low_elev['Delay'], equal_var=False)
41
42 print("t-statistic:", t_stat)
43 print("p-value:", p_value)
44
```

Count of Delays by Elevation: Arr Airport



Mean Delay by Elevation Arr Airport



<Figure size 1200x2000 with 0 Axes>

t-statistic: 0.7498156315532386

p-value: 0.4533717588728631

Correlation

```
1 df.info()
2 df.columns
3 # Select relevant numerical features for correlation analysis
4 predictors = ['DayOfWeek', 'Time', 'Length',
5               'Delay', 'elevation_ft_from',
6               'length_ft_from', 'width_ft_from',
7               'number_of_runways_from', 'elevation_ft_to',
8               'length_ft_to', 'width_ft_to',
9               'number_of_runways_to',
10              'airline_experience', 'passenger_tfc_from',
11              'passenger_tfc_to' ]
12 corr_matrix = df[predictors].corr()
13
14 # Heatmap visualization
15 import seaborn as sns
16 # Increase size
17 plt.figure(figsize=(12, 8)) # Adjust width and height as desired
18
19 # Select a colormap:
20 cmap = sns.color_palette("coolwarm", as_cmap=True) # Try other palettes!
21
22 # Create the heatmap
23 sns.heatmap(corr_matrix, annot=True, cmap=cmap)
24 plt.title('Correlation Heatmap Enhanced')
25 plt.show()
26 # Convert the correlation matrix to a DataFrame (if it isn't already)
27 if not isinstance(corr_matrix, pd.DataFrame):
28     corr_matrix = pd.DataFrame(corr_matrix)
29
30 # Export to Excel
31 corr_matrix.to_excel(file_path + 'correlation_matrix.xlsx', index=True)
32 # 'index=True' will include the row and column labels
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 518494 entries, 0 to 518493
Data columns (total 32 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Unnamed: 0        518494 non-null   int64  
 1   Airline          518494 non-null   object  
 2   AirportFrom      518494 non-null   object  
 3   AirportTo        518494 non-null   object  
 4   DayOfWeek        518494 non-null   int64  
 5   Time             518494 non-null   int64  
 6   Length           518494 non-null   int64  
 7   Delay            518494 non-null   bool    
 8   type_from        518494 non-null   object  
 9   elevation_ft_from 518494 non-null   int64  
 10  iata_code_from   518494 non-null   object  
 11  length_ft_from   518494 non-null   int64  
 12  width_ft_from    518494 non-null   int64  
 13  lighted_from     518494 non-null   int64  
 14  number_of_runways_from 518494 non-null   int64  
 15  elevation_ft_to   518494 non-null   int64  
 16  iata_code_to     518494 non-null   object  
 17  airport_ident    518494 non-null   object  
 18  length_ft_to     518494 non-null   int64  
 19  width_ft_to      518494 non-null   int64  
 20  lighted_to       518494 non-null   int64  
 21  number_of_runways_to 518494 non-null   int64  
 22  type_to          518494 non-null   object  
 23  elevation_ft     518494 non-null   float64 
 24  iata_code        518494 non-null   object  
 25  airline_experience 518494 non-null   float64 
 26  passenger_tfc_from 433524 non-null   float64 
 27  hub_type_from    518494 non-null   object  
 28  passenger_tfc_to 433576 non-null   float64 
 29  hub_type_to      518494 non-null   object  
 30  ID               518494 non-null   int64  
 31  duration         518494 non-null   object  
dtypes: bool(1), float64(4), int64(15), object(12)
memory usage: 123.1+ MB

```

Correlation Heatmap Enhanced



Machine learning

```

1 df.columns
2 df['type_to'].unique()
3 # Drop unnecessary columns from df
4 df=df.drop(['passenger_tfc_from','passenger_tfc_to','elevation_ft','iata_code_from','iata_code_to','airport_ident'],axis=1)
5 df=df.drop(['iata_code'],axis=1)
6
7 # Assuming 'AirportFrom' is already in the categorical_features list
8 encoder = OrdinalEncoder()
9 df['AirportFrom_encode'] = encoder.fit_transform(df[['AirportFrom']])
10 df['AirportFrom'].nunique()
11 df=df.drop(['AirportFrom_encode'],axis=1)
12
13 df.columns
14 df.shape
15

```

(518494, 25)

✓ Label encode

I am encoding it out of the pipeline due to repeated errors

```
1 from sklearn.preprocessing import LabelEncoder
2 le = LabelEncoder()
3
4 df['Airline'] = le.fit_transform(df['Airline'])
5 df['AirportFrom'] = le.fit_transform(df['AirportFrom'])
6 df['AirportTo'] = le.fit_transform(df['AirportTo'])
7 df.head()
```

| | Unnamed: 0 | Airline | AirportFrom | AirportTo | DayOfWeek | Time | Length | Delay | type_from | elevation_ft_from | ... | length_ft_to | width_ft_ |
|---|---------------|---------|-------------|-----------|-----------|------|--------|-------|---------------|-------------------|-----|--------------|-----------|
| 0 | 0 | 4 | 251 | 133 | 3 | 15 | 205 | True | large_airport | 13 | ... | 9000 | 1 |
| 1 | 1 | 13 | 215 | 59 | 3 | 15 | 222 | True | large_airport | 1135 | ... | 7502 | 1 |
| 2 | 2 | 1 | 152 | 78 | 3 | 20 | 165 | True | large_airport | 125 | ... | 9000 | 2 |
| 3 | 3 | 1 | 251 | 78 | 3 | 20 | 195 | True | large_airport | 13 | ... | 9000 | 2 |
| 4 | 4 | 2 | 13 | 250 | 3 | 30 | 202 | False | large_airport | 152 | ... | 9426 | 1 |

5 rows × 25 columns

✓ Pipeline Model Fitting and testing

This is the main ML codeblock where ppl has been used to encode categorical data and model fitting is done in a loop. This is a very useful and reusable code 😊

```

1
2 from sklearn.compose import ColumnTransformer
3 from sklearn.preprocessing import OrdinalEncoder, StandardScaler
4 from sklearn.pipeline import Pipeline
5 from sklearn.metrics import accuracy_score
6 from sklearn.model_selection import train_test_split, StratifiedKFold
7 from sklearn.linear_model import LogisticRegression
8 from sklearn.tree import DecisionTreeClassifier
9 from sklearn.ensemble import GradientBoostingClassifier
10 from sklearn.preprocessing import LabelEncoder
11 from sklearn.feature_extraction import FeatureHasher
12 from sklearn.ensemble import RandomForestClassifier
13
14 # Assuming 'delay' is your target variable (binary: delayed or not)
15 X = df[['Airline', 'AirportFrom', 'AirportTo', 'DayOfWeek',
16          'Time', 'Length', 'type_from', 'elevation_ft_from',
17          'length_ft_from', 'width_ft_from', 'lighted_from',
18          'number_of_runways_from', 'elevation_ft_to', 'length_ft_to',
19          'width_ft_to', 'lighted_to', 'number_of_runways_to', 'type_to',
20          'airline_experience', 'hub_type_from', 'hub_type_to']]
21 y = df['Delay']
22
23
24
25 # Identify Feature Types
26 ordinal_features = ['type_from','type_to','hub_type_from', 'hub_type_to'] # Replace with your actual ordinal feature
27 numerical_features = ['DayOfWeek','Time', 'Length', 'elevation_ft_from','length_ft_from', 'width_ft_from',
28                        'number_of_runways_from', 'elevation_ft_to', 'length_ft_to',
29                        'width_ft_to', 'number_of_runways_to',
30                        'airline_experience',] # Your numerical features
31 airport_type_order = ['small_airport', 'medium_airport', 'large_airport'] # Order for 'type_from', 'type_to'
32 hub_type_order=['large', 'medium', 'small']
33 ordinal_order=[airport_type_order,airport_type_order,hub_type_order,hub_type_order]
34
35
36 # Build the Pipeline
37 preprocessor = ColumnTransformer(
38     transformers=[
39
40         ('oe_type', OrdinalEncoder(categories='auto'), ordinal_features),
41         ('scaler', StandardScaler(), numerical_features)
42     ],
43     remainder='passthrough'
44 )
45
46
47
48 # Split into train/test
49 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25, random_state=42)
50 X_train_transformed = preprocessor.fit_transform(X_train)
51 X_test_transformed = preprocessor.transform(X_test)
52
53 # Model Building
54 models = {
55     'Logistic Regression': LogisticRegression(solver='liblinear'),
56     'Decision Tree': DecisionTreeClassifier(max_depth=5), # Example depth control
57     # Hyper tuning parameter set
58     'Gradient Boosting': GradientBoostingClassifier(max_depth=8, n_estimators=100),
59     # tuning hyper parameters based on grid search
60     'Random Forest Classifier' : RandomForestClassifier(max_depth= 8, min_samples_split= 10, n_estimators=200)
61 }
62
63 # Stratified K-Fold for voting
64 skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
65
66 for name, model in models.items():
67     accuracy_scores = []
68     for train_index, test_index in skf.split(X_train_transformed, y_train):
69         X_train_fold, y_train_fold = X_train_transformed[train_index], y_train.iloc[train_index]
70         model.fit(X_train_fold, y_train_fold)
71         y_pred = model.predict(X_train_transformed[test_index])
72         accuracy = accuracy_score(y_train.iloc[test_index], y_pred)
73         accuracy_scores.append(accuracy)

```

```

74     print(f"\nname) Average Accuracy: {np.mean(accuracy_scores):.3f}")
Plot the gradient booster
    Decision Tree Average Accuracy: 0.634
1 import matplotlib.pyplot as plt
2
3 model = models['Gradient Boosting']
4 feature_importances = model.feature_importances_
5 feature_names = X.columns # Or update with transformed column names
6 sorted_indices = np.argsort(feature_importances)[::-1] # Sort importances
7
8 plt.barh(feature_names[sorted_indices], feature_importances[sorted_indices])
9 plt.xlabel("Feature Importance")
10 plt.title("Gradient Boosting Feature Importance")
11 plt.show()
12
13 #Print top features
14 print('Top 10 Features (Gradient Boosting):')
15 for i in sorted_indices[:10]: # Print the top 10
16     print(f'{feature_names[i]}: {feature_importances[i]:.3f}')

```

