

Problem: Predicting Airplane Delays

The goals of this notebook are:

- Process and create a dataset from downloaded .zip files
- Perform exploratory data analysis (EDA)
- Establish a baseline model
- Move from a simple model to an ensemble model
- Perform hyperparameter optimization
- Check feature importance

Introduction to business scenario

You work for a travel booking website that wants to improve the customer experience for flights that were delayed. The company wants to create a feature to let customers know if the flight will be delayed because of weather when they book a flight to or from the busiest airports for domestic travel in the US.

You are tasked with solving part of this problem by using machine learning (ML) to identify whether the flight will be delayed because of weather. You have been given access to the a dataset about the on-time performance of domestic flights that were operated by large air carriers. You can use this data to train an ML model to predict if the flight is going to be delayed for the busiest airports.

About this dataset

This dataset contains scheduled and actual departure and arrival times reported by certified US air carriers that account for at least 1 percent of domestic scheduled passenger revenues. The data was collected by the U.S. Office of Airline Information, Bureau of Transportation Statistics (BTS). The dataset contains date, time, origin, destination, airline, distance, and delay status of flights for flights between 2013 and 2018.

Features

For more information about features in the dataset, see [On-time delay dataset features](https://www.transtats.bts.gov/Fields.asp) (<https://www.transtats.bts.gov/Fields.asp>).

Dataset attributions

Website: <https://www.transtats.bts.gov/> (<https://www.transtats.bts.gov/>)

Dataset(s) used in this lab were compiled by the U.S. Office of Airline Information, Bureau of Transportation Statistics (BTS), Airline On-Time Performance Data, available at https://www.transtats.bts.gov/DatabaseInfo.asp?DB_ID=120&DB_URL=Mode_ID=1&Mode_Desc=Aviation&Subject_ID2=0 (https://www.transtats.bts.gov/DatabaseInfo.asp?DB_ID=120&DB_URL=Mode_ID=1&Mode_Desc=Aviation&Subject_ID2=0).

Step 1: Problem formulation and data collection

Start this project by writing a few sentences that summarize the business problem and the business goal that you want to achieve in this scenario. You can write down your ideas in the following sections. Include a business metric that you would like your team to aspire toward. After you define that information, write the ML problem statement. Finally, add a comment or two about the type of ML this activity represents.

Project presentation: Include a summary of these details in your project presentation.

1. Determine if and why ML is an appropriate solution to deploy for this scenario.

In [2]:

```
# Write your answer here  
#Machine Learning is the scientific study of algorithms and statistical models to perform a  
#instead of instructions. Look at this scenario, the dataset is large and it has large numb  
#also lacks clear procedures to obtain the solution. At that time, using machine Learning i
```

2. Formulate the business problem, success metrics, and desired ML output.

In [3]:

```
# Write your answer here  
#I want to predict if the flight will be delayed because of weather.  
#This problem is a classification problem. We want to predict the flight will be  
#delayed or not.  
  
#why:let customers know if the flight will be delayed because of weather when they book a f  
#increase customer trust in the airline and choose this airline more often  
  
#success metrics: we are using ROC graph to evaluate models.  
#desired ML output:The result is either delayed or not delayed.
```

3. Identify the type of ML problem that you're working with.

In [4]:

```
# Write your answer here  
#Binary Classification Problem
```

4. Analyze the appropriateness of the data that you're working with.

In [5]:

```
# Write your answer here  
#first of all, we should extract, transform and load the data, making sure it's the right f  
#as we are faced with a massive amount of data. It has many features but many features are  
#useful features.  
#we should shuffle the data
```

Setup

Now that you have decided where you want to focus your attention, you will set up this lab so that you can start solving the problem.

Note: This notebook was created and tested on an `ml.m4.xlarge` notebook instance with 25 GB storage.

In [1]:

```
import os
from pathlib2 import Path
from zipfile import ZipFile
import time

import pandas as pd
import numpy as np
import subprocess

import matplotlib.pyplot as plt
import seaborn as sns

sns.set()
instance_type='ml.m4.xlarge'

import warnings
warnings.filterwarnings('ignore')

%matplotlib inline
```

Step 2: Data preprocessing and visualization

In this data preprocessing phase, you explore and visualize your data to better understand it. First, import the necessary libraries and read the data into a pandas DataFrame. After you import the data, explore the dataset. Look for the shape of the dataset and explore your columns and the types of columns that you will work with (numerical, categorical). Consider performing basic statistics on the features to get a sense of feature means and ranges. Examine your target column closely, and determine its distribution.

Specific questions to consider

Throughout this section of the lab, consider the following questions:

1. What can you deduce from the basic statistics that you ran on the features?
2. What can you deduce from the distributions of the target classes?
3. Is there anything else you can deduce by exploring the data?

Project presentation: Include a summary of your answers to these questions (and other similar questions) in your project presentation.

Start by bringing in the dataset from a public Amazon Simple Storage Service (Amazon S3) bucket to this notebook environment.

In [2]:

```
# download the files

zip_path = '/home/ec2-user/SageMaker/project/data/FlightDelays/'
base_path = '/home/ec2-user/SageMaker/project/data/FlightDelays/'
csv_base_path = '/home/ec2-user/SageMaker/project/data/csvFlightDelays/'

!mkdir -p {zip_path}
!mkdir -p {csv_base_path}
!aws s3 cp s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/ {zip_path}
```

```
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_proje
ct/data/On_Time_Reportin
Carrier_On_Time_Performance_1987_present_2014_1
1.zip to project/data/FlightDelays/On_Time_Reportin
Carrier_On_Time_Perf
ormance_1987_present_2014_11.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_proje
ct/data/On_Time_Reportin
Carrier_On_Time_Performance_1987_present_2014_
2.zip to project/data/FlightDelays/On_Time_Reportin
Carrier_On_Time_Perf
ormance_1987_present_2014_2.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_proje
ct/data/On_Time_Reportin
Carrier_On_Time_Performance_1987_present_2014_1
0.zip to project/data/FlightDelays/On_Time_Reportin
Carrier_On_Time_Perf
ormance_1987_present_2014_10.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_proje
ct/data/On_Time_Reportin
Carrier_On_Time_Performance_1987_present_2014_
1.zip to project/data/FlightDelays/On_Time_Reportin
Carrier_On_Time_Perf
ormance_1987_present_2014_1.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_proje
ct/data/On_Time_Reportin
Carrier_On_Time_Performance_1987_present_2014_1
2.zip to project/data/FlightDelays/On_Time_Reportin
Carrier_On_Time_Perf
ormance_1987_present_2014_12.zip
```

In [3]:

```
zip_files = [str(file) for file in list(Path(base_path).iterdir()) if '.zip' in str(file)]
len(zip_files)
```

Out[3]:

60

Extract comma-separated values (CSV) files from the .zip files.

In [4]:

```
def zip2csv(zipFile_name , file_path):
    """
    Extract csv from zip files
    zipFile_name: name of the zip file
    file_path : name of the folder to store csv
    """

    try:
        with ZipFile(zipFile_name, 'r') as z:
            print(f'Extracting {zipFile_name} ')
            z.extractall(path=file_path)
    except:
        print(f'zip2csv failed for {zipFile_name}')

    for file in zip_files:
        zip2csv(file, csv_base_path)

    print("Files Extracted")
```

```
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Report_Carrier_On_Time_Performance_1987_present_2018_2.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Report_Carrier_On_Time_Performance_1987_present_2015_11.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Report_Carrier_On_Time_Performance_1987_present_2014_3.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Report_Carrier_On_Time_Performance_1987_present_2017_12.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Report_Carrier_On_Time_Performance_1987_present_2014_12.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Report_Carrier_On_Time_Performance_1987_present_2018_8.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Report_Carrier_On_Time_Performance_1987_present_2016_3.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Report_Carrier_On_Time_Performance_1987_present_2014_9.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Report_Carrier_On_Time_Performance_1987_present_2015_5.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Report_Carrier_On_Time_Performance_1987_present_2010_10.zip
```

In [5]:

```
csv_files = [str(file) for file in list(Path(csv_base_path).iterdir()) if '.csv' in str(file)]
len(csv_files)
```

Out[5]:

60

Before you load the CSV file, read the HTML file from the extracted folder. This HTML file includes the background and more information about the features that are included in the dataset.

In [6]:

```
from IPython.display import IFrame  
  
IFrame(src=os.path.relpath(f"{csv_base_path}readme.html"), width=1000, height=600)
```

Out[6]:

Load sample CSV file

Before you combine all the CSV files, examine the data from a single CSV file. By using pandas, read the `On_Time_Reported_Carrier_On_Time_Performance_(1987_present)_2018_9.csv` file first. You can use the built-in `read_csv` function in Python ([pandas.read_csv documentation \(https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.read_csv.html\)](https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.read_csv.html)).

In [7]:

```
df_temp = pd.read_csv(f"{csv_base_path}On_Time_Reported_Carrier_On_Time_Performance_(1987_
```

Question: Print the row and column length in the dataset, and print the column names.

Hint: To view the rows and columns of a DataFrame, use the `<DataFrame>.shape` function. To view the column names, use the `<DataFrame>.columns` function.

In [8]:

```
df_shape = df_temp.shape  
print(f'Rows and columns in one CSV file is {df_shape}')
```

Rows and columns in one CSV file is (585749, 110)

Question: Print the first 10 rows of the dataset.

Hint: To print x number of rows, use the built-in `head(x)` function in pandas.

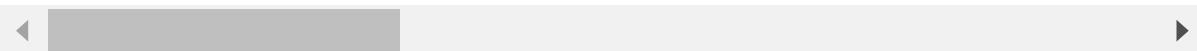
In [9]:

```
df_temp.head(10)
```

Out[9]:

	Year	Quarter	Month	DayofMonth	DayOfWeek	FlightDate	Reporting_Airline	DOT_ID_Report
0	2018	3	9	3	1	2018-09-03	9E	
1	2018	3	9	9	7	2018-09-09	9E	
2	2018	3	9	10	1	2018-09-10	9E	
3	2018	3	9	13	4	2018-09-13	9E	
4	2018	3	9	14	5	2018-09-14	9E	
5	2018	3	9	16	7	2018-09-16	9E	
6	2018	3	9	17	1	2018-09-17	9E	
7	2018	3	9	20	4	2018-09-20	9E	
8	2018	3	9	21	5	2018-09-21	9E	
9	2018	3	9	23	7	2018-09-23	9E	

10 rows × 110 columns

**Question:** Print all the columns in the dataset. To view the column names, use `<DataFrame>.columns`.

In [10]:

```
print(f'The column names are :')
print('#####')
for col in df_temp.columns:
    print(col)
```

```
The column names are :
#####
Year
Quarter
Month
DayofMonth
DayOfWeek
FlightDate
Reporting_Airline
DOT_ID_Reported_Airline
IATA_CODE_Reported_Airline
Tail_Number
Flight_Number_Reported_Airline
OriginAirportID
OriginAirportSeqID
OriginCityMarketID
Origin
OriginCityName
OriginState
: : : : : :
```

Question: Print all the columns in the dataset that contain the word *Del*. This will help you see how many columns have *delay data* in them.

Hint: To include values that pass certain `if` statement criteria, you can use a Python list comprehension.

For example: `[x for x in [1,2,3,4,5] if x > 2]`

Hint: To check if the value is in a list, you can use the `in` keyword ([Python in Keyword documentation](#) (https://www.w3schools.com/python/ref_keyword_in.asp)).

For example: `5 in [1,2,3,4,5]`

In [11]:

```
Del_cols = [col for col in df_temp.columns if 'Del' in col]
print(Del_cols)
```

```
['DepDelay', 'DepDelayMinutes', 'DepDel15', 'DepartureDelayGroups', 'ArrDelay',
 'ArrDelayMinutes', 'ArrDel15', 'ArrivalDelayGroups', 'CarrierDelay', 'WeatherDelay',
 'NASDelay', 'SecurityDelay', 'LateAircraftDelay', 'DivArrDelay']
```

Here are some more questions to help you learn more about your dataset.

Questions

1. How many rows and columns does the dataset have?
2. How many years are included in the dataset?
3. What is the date range for the dataset?
4. Which airlines are included in the dataset?

5. Which origin and destination airports are covered?

Hints

- To show the dimensions of the DataFrame, use `df_temp.shape` .
- To refer to a specific column, use `df_temp.columnName` (for example, `df_temp.CarrierDelay`).
- To get unique values for a column, use `df_temp.column.unique()` (for, example `df_temp.Year.unique()`).

In [12]:

```
print("The #rows and #columns are ", df_temp.shape[0] , " and ", df_temp.shape[1])
print("The years in this dataset are: ", df_temp['Year'].unique())
print("The months covered in this dataset are: ", df_temp['Month'].unique())
print("The date range for data is :" , min(df_temp['FlightDate']), " to ", max(df_temp['FlightDate']))
print("The airlines covered in this dataset are: ", list(df_temp['Reporting_Airline'].unique()))
print("The Origin airports covered are: ", list(df_temp['Origin'].unique()))
print("The Destination airports covered are: ", list(df_temp['Dest'].unique()))
```

The #rows and #columns are 585749 and 110
The years in this dataset are: [2018]
The months covered in this dataset are: [9]
The date range for data is : 2018-09-01 to 2018-09-30
The airlines covered in this dataset are: ['9E', 'B6', 'WN', 'YV', 'YX', 'EV', 'AA', 'AS', 'DL', 'HA', 'UA', 'F9', 'G4', 'MQ', 'NK', 'OH', 'OO']
The Origin airports covered are: ['DFW', 'LGA', 'MSN', 'MSP', 'ATL', 'BDL', 'VLD', 'JFK', 'RDU', 'CHS', 'DTW', 'GRB', 'PVD', 'SHV', 'FNT', 'PIT', 'RIC', 'RST', 'RSW', 'CVG', 'LIT', 'ORD', 'JAX', 'TRI', 'BOS', 'CWA', 'DCA', 'CHO', 'AVP', 'IND', 'GRR', 'BTR', 'MEM', 'TUL', 'CLE', 'STL', 'BTW', 'OMA', 'MGM', 'TVC', 'SAV', 'GSP', 'EWR', 'OAJ', 'BNA', 'MCI', 'TLH', 'ROC', 'LEX', 'PWM', 'BUF', 'AGS', 'CLT', 'GSO', 'BWI', 'SAT', 'PHL', 'TYS', 'ACK', 'DSM', 'GNV', 'AVL', 'BGR', 'MHT', 'ILM', 'MOT', 'IAH', 'SBN', 'SYR', 'ORF', 'MKE', 'XNA', 'MSY', 'PBI', 'ABE', 'HPN', 'EVV', 'ALB', 'LNK', 'AUS', 'PHF', 'CHA', 'GTR', 'BMI', 'BQK', 'CID', 'CAK', 'ATW', 'ABY', 'CAE', 'SRQ', 'MLI', 'BHM', 'IAD', 'CSG', 'CMH', 'MCO', 'MBS', 'FLL', 'SDF', 'TPA', 'MVY', 'LAS', 'LGB', 'SFO', 'SAN', 'LAX', 'RNO', 'PDX', 'ANC', 'ABQ', 'SLC', 'DEN', 'PHX', 'OAK', 'SMF', 'SJU', 'SEA', 'HOU', 'STX', 'BUR', 'SWF', 'SJC', 'DAB', 'BQN', 'PSE', 'ORH', 'HYA', 'STT', 'ONT', 'HRL', 'ICT', 'ISP', 'LBB', 'MAF', 'MDW', 'OKC', 'PNS', 'SNA', 'TUS', 'AMA', 'BOI', 'CRP', 'DAL', 'ECP', 'ELP', 'GEG', 'LFT', 'MFE', 'MDT', 'JAN', 'COS', 'MOB', 'VPS', 'MTJ', 'DRO', 'GPT', 'BFL', 'MRY', 'SBA', 'PSP', 'FSD', 'BRO', 'RAP', 'COU', 'STS', 'PIA', 'FAT', 'SBP', 'FSM', 'HSV', 'BIS', 'DAY', 'BZN', 'MIA', 'EYW', 'MYR', 'HHH', 'GJT', 'FAR', 'SGF', 'HOB', 'CLL', 'LRD', 'AEX', 'ERI', 'MLU', 'LCH', 'ROA', 'LAW', 'MHK', 'GRK', 'SAF', 'GRI', 'JLN', 'ROW', 'FWA', 'CRW', 'LAN', 'OGG', 'HNL', 'KOA', 'EGE', 'LIH', 'MLB', 'JAC', 'FAI', 'RDM', 'ADQ', 'BET', 'BRW', 'SCC', 'KTN', 'YAK', 'CDV', 'JNU', 'SIT', 'PSG', 'WRG', 'OME', 'OTZ', 'ADK', 'FCA', 'FAY', 'PSC', 'BIL', 'MSO', 'ITO', 'PPG', 'MFR', 'EUG', 'GUM', 'SPN', 'DLH', 'TTN', 'BKG', 'SFB', 'PIE', 'PGD', 'AZA', 'SMX', 'RFD', 'SCK', 'OWB', 'HTS', 'BLV', 'IAG', 'USA', 'GFK', 'BLI', 'ELM', 'PBG', 'LCK', 'GTF', 'OGD', 'IDA', 'PVU', 'TOL', 'PSM', 'CKB', 'HGR', 'SPI', 'STC', 'ACT', 'TYR', 'ABI', 'AZO', 'CMI', 'BPT', 'GCK', 'MQT', 'ALO', 'TXK', 'SPS', 'SWO', 'DBQ', 'SUX', 'SJT', 'GGG', 'LSE', 'LBE', 'ACY', 'LYH', 'PGV', 'HVN', 'EWN', 'DHN', 'PIH', 'IMT', 'WYS', 'CPR', 'SCE', 'HLN', 'SUN', 'ISN', 'CMX', 'EAU', 'LWB', 'SHD', 'LBF', 'HYS', 'SLN', 'EAR', 'VEL', 'CNY', 'GCC', 'RKS', 'PUB', 'LBL', 'MKG', 'PAH', 'CGI', 'UIN', 'BFF', 'DVL', 'JMS', 'LAR', 'SGU', 'PRC', 'ASE', 'RDD', 'ACV', 'OTH', 'COD', 'LWS', 'ABR', 'APN', 'ESC', 'PLN', 'BJI', 'BRD', 'BTM', 'CDC', 'CIU', 'EKO', 'TWF', 'HIB', 'BGM', 'RHI', 'ITH', 'INL', 'FLG', 'YUM', 'MEI', 'PIB', 'HDN']
The Destination airports covered are: ['CVG', 'PWM', 'RDU', 'MSP', 'MSN', 'SHV', 'CLT', 'PIT', 'RIC', 'IAH', 'ATL', 'JFK', 'DCA', 'DTW', 'LGA', 'TYS', 'PVD', 'FNT', 'LIT', 'BUF', 'ORD', 'TRI', 'IND', 'BGR', 'AVP', 'BWI', 'LEX', 'BDL', 'GRR', 'CWA', 'TUL', 'MEM', 'AGS', 'EWR', 'MGM', 'PHL', 'SYR', 'OMA', 'STL', 'TVC', 'ORF', 'CLE', 'ABY', 'BOS', 'OAJ', 'TLH', 'BTR', 'SAT', 'JAX', 'BNA', 'CHO', 'VLD', 'ROC', 'DFW', 'GNV', 'ACK', 'PBI', 'CHS', 'GRB', 'MOT', 'MKE', 'DSM', 'ILM', 'GSO', 'MCI', 'SBN', 'BTW', 'MVY', 'XNA', 'RST', 'EVV', 'HPN', 'RSW', 'MDT', 'ROA', 'GSP', 'MCO', 'CSG', 'SAV', 'PHF', 'ALB', 'CHA', 'ABE', 'BMI', 'MSY', 'IAD', 'GTR', 'CID', 'CAK', 'ATW', 'AUS', 'BQK', 'MLI', 'CAE', 'CMH', 'AVL', 'MBS', 'FLL', 'SDF', 'TPA', 'LNK', 'SRQ', 'MHT', 'BHM', 'LAS', 'SFO', 'SAN', 'RNO', 'LGB', 'ANC', 'PDX', 'SJI', 'ABQ', 'SLC', 'DEN', 'LAX', 'PHX', 'OAK', 'SMF', 'SEA', 'STX', 'BUR', 'DAB', 'SJC', 'SWF', 'HOU']

```
'BQN', 'PSE', 'ORH', 'HYA', 'STT', 'ONT', 'DAL', 'ECP', 'ELP', 'HRL', 'MAF',
'MDW', 'OKC', 'PNS', 'SNA', 'AMA', 'BOI', 'GEG', 'ICT', 'LBB', 'TUS', 'ISP',
'CRP', 'MFE', 'LFT', 'VPS', 'JAN', 'COS', 'MOB', 'DRO', 'GPT', 'BFL', 'COU',
'SBP', 'MTJ', 'SBA', 'PSP', 'FSD', 'FSM', 'BRO', 'PIA', 'STS', 'FAT', 'RAP',
'MRY', 'HSV', 'BIS', 'DAY', 'BNZ', 'MIA', 'EYW', 'MYR', 'HHH', 'GJT', 'FAR',
'MLU', 'LRD', 'CLL', 'LCH', 'FWA', 'GRK', 'SGF', 'HOB', 'LAW', 'MHK', 'SAF',
'JLN', 'ROW', 'GRI', 'AEX', 'CRW', 'LAN', 'ERI', 'HNL', 'KOA', 'OGG', 'EGE',
'LIH', 'JAC', 'MLB', 'RDM', 'BET', 'ADQ', 'BRW', 'SCC', 'FAI', 'JNU', 'CDV',
'YAK', 'SIT', 'KTN', 'WRG', 'PSG', 'OME', 'OTZ', 'ADK', 'FCA', 'BIL', 'PSC',
'FAY', 'MSO', 'ITO', 'PPG', 'MFR', 'DLH', 'EUG', 'GUM', 'SPN', 'TTN', 'BKG',
'AZA', 'SFB', 'LCK', 'BLI', 'SCK', 'PIE', 'RFD', 'PVU', 'PBG', 'BLV', 'PGD',
'SPI', 'USA', 'TOL', 'IDA', 'ELM', 'HTS', 'HGR', 'SMX', 'OGD', 'GFK', 'STC',
'GTF', 'IAG', 'CKB', 'OWB', 'PSM', 'ABI', 'TYR', 'ALO', 'SUX', 'AZO', 'ACT',
'CMI', 'BPT', 'TXK', 'SWO', 'SPS', 'DBQ', 'SJT', 'GGG', 'LSE', 'MQT', 'GCK',
'LBE', 'ACY', 'LYH', 'PGV', 'HVN', 'EWN', 'DHN', 'PIH', 'WYS', 'SCE', 'IMT',
'HLN', 'ASE', 'SUN', 'ISN', 'EAR', 'SGU', 'VEL', 'SHD', 'LWB', 'MKG', 'SLN',
'HYS', 'BFF', 'PUB', 'LBL', 'CMX', 'EAU', 'PAH', 'UIN', 'RKS', 'CGI', 'CNY',
'JMS', 'DVL', 'LAR', 'GCC', 'LBF', 'PRC', 'RDD', 'ACV', 'OTH', 'COD', 'LWS',
'ABR', 'APN', 'PLN', 'BJI', 'CPR', 'BRD', 'BTM', 'CDC', 'CIU', 'ESC', 'EKO',
'ITH', 'HIB', 'BGM', 'TWF', 'RHI', 'INL', 'FLG', 'YUM', 'MEI', 'PIB', 'HDN']
```

Question: What is the count of all the origin and destination airports?

Hint: To find the values for each airport by using the **Origin** and **Dest** columns, you can use the `values_count` function in pandas ([pandas.Series.value_counts documentation](#) (https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.Series.value_counts.html)).

In [13]:

```
counts = pd.DataFrame({'Origin':df_temp['Origin'].value_counts(), 'Destination':df_temp['De
counts
```

Out[13]:

Origin	Destination
ABE	303
ABI	169
ABQ	2077
ABR	60
ABY	79
...	...
WRG	60
WYS	52
XNA	1004
YAK	60
YUM	96

346 rows × 2 columns

Question: Print the top 15 origin and destination airports based on number of flights in the dataset.

Hint: You can use the `sort_values` function in pandas ([pandas.DataFrame.sort_values documentation](#) (https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.sort_values.html)).

In [14]:

```
counts.sort_values(by=['Origin', 'Destination'], ascending=False).head(15)
```

Out[14]:

	Origin	Destination
ATL	31525	31521
ORD	28257	28250
DFW	22802	22795
DEN	19807	19807
CLT	19655	19654
LAX	17875	17873
SFO	14332	14348
IAH	14210	14203
LGA	13850	13850
MSP	13349	13347
LAS	13318	13322
PHX	13126	13128
DTW	12725	12724
BOS	12223	12227
SEA	11872	11877

Given all the information about a flight trip, can you predict if it would be delayed?

The **ArrDel15** column is an indicator variable that takes the value *1** when the delay is more than 15 minutes. Otherwise, it takes a value of *0*.

You could use this as a target column for the classification problem.

Now, assume that you are traveling from San Francisco to Los Angeles on a work trip. You want to better manage your reservations in Los Angeles. Thus, want to have an idea of whether your flight will be delayed, given a set of features. How many features from this dataset would you need to know before your flight?

Columns such as **DepDelay**, **ArrDelay**, **CarrierDelay**, **WeatherDelay**, **NASDelay**, **SecurityDelay**, **LateAircraftDelay**, and **DivArrDelay** contain information about a delay. But this delay could have occurred at the origin or the destination. If there were a sudden weather delay 10 minutes before landing, this data wouldn't be helpful to managing your Los Angeles reservations.

So to simplify the problem statement, consider the following columns to predict an arrival delay:

Year, **Quarter**, **Month**, **DayofMonth**, **DayOfWeek**, **FlightDate**, **Reporting_Airline**, **Origin**, **OriginState**, **Dest**, **DestState**, **CRSDepTime**, **DepDelayMinutes**, **DepartureDelayGroups**, **Cancelled**, **Diverted**, **Distance**, **DistanceGroup**, **ArrDelay**, **ArrDelayMinutes**, **ArrDel15**, **AirTime**

You will also filter the source and destination airports to be:

- Top airports: ATL, ORD, DFW, DEN, CLT, LAX, IAH, PHX, SFO
- Top five airlines: UA, OO, WN, AA, DL

This information should help reduce the size of data across the CSV files that will be combined.

Combine all CSV files

First, create an empty DataFrame that you will use to copy your individual DataFrames from each file. Then, for each file in the `csv_files` list:

1. Read the CSV file into a dataframe
2. Filter the columns based on the `filter_cols` variable

```
columns = ['col1', 'col2']
df_filter = df[columns]
```

3. Keep only the `subset_vals` in each of the `subset_cols`. To check if the `val` is in the DataFrame column, use the `isin` function in pandas ([pandas.DataFrame.isin documentation](#) (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.isin.html>)). Then, choose the rows that include it.

```
df_eg[df_eg['col1'].isin('5')]
```

4. Concatenate the DataFrame with the empty DataFrame

In [15]:

```
def combine_csv(csv_files, filter_cols, subset_cols, subset_vals, file_name):
    """
    Combine csv files into one Data Frame
    csv_files: list of csv file paths
    filter_cols: list of columns to filter
    subset_cols: list of columns to subset rows
    subset_vals: list of list of values to subset rows
    """

    df = pd.DataFrame()

    for file in csv_files:
        df_temp = pd.read_csv(file)
        df_temp = df_temp[filter_cols]
        for col, val in zip(subset_cols, subset_vals):
            df_temp = df_temp[df_temp[col].isin(val)] 

        df = pd.concat([df, df_temp], axis=0)

    df.to_csv(file_name, index=False)
    print(f'Combined csv stored at {file_name}')
```

In [16]:

```
#cols is the list of columns to predict Arrival Delay
cols = ['Year', 'Quarter', 'Month', 'DayofMonth', 'DayOfWeek', 'FlightDate',
         'Reporting_Airline', 'Origin', 'OriginState', 'Dest', 'DestState',
         'CRSDepTime', 'Cancelled', 'Diverted', 'Distance', 'DistanceGroup',
         'ArrDelay', 'ArrDelayMinutes', 'ArrDel15', 'AirTime']

subset_cols = ['Origin', 'Dest', 'Reporting_Airline']

# subset_vals is a list collection of the top origin and destination airports and top 5 air
subset_vals = [[['ATL', 'ORD', 'DFW', 'DEN', 'CLT', 'LAX', 'IAH', 'PHX', 'SFO'],
                ['ATL', 'ORD', 'DFW', 'DEN', 'CLT', 'LAX', 'IAH', 'PHX', 'SFO'],
                ['UA', 'OO', 'WN', 'AA', 'DL']]
```

Use the previous function to merge all the different files into a single file that you can read easily.

Note: This process will take 5-7 minutes to complete.

In [17]:

```
start = time.time()
combined_csv_filename = f"{base_path}combined_files.csv"
combine_csv(csv_files, cols, subset_cols, subset_vals, combined_csv_filename)
print(f'CSVs merged in {round((time.time() - start)/60,2)} minutes')
```

Combined csv stored at /home/ec2-user/SageMaker/project/data/FlightDelays/co
mbined_files.csv
CSVs merged in 4.57 minutes

Load the dataset

Load the combined dataset.

In [18]:

```
data = pd.read_csv(combined_csv_filename)
```

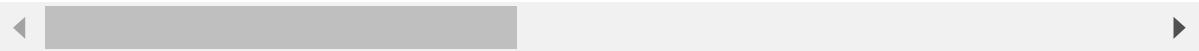
Print the first five records.

In [19]:

data.head(5)

Out[19]:

	Year	Quarter	Month	DayofMonth	DayOfWeek	FlightDate	Reporting_Airline	Origin	Origin
0	2016	3	7	1	5	2016-07-01	AA	SFO	
1	2016	3	7	2	6	2016-07-02	AA	SFO	
2	2016	3	7	3	7	2016-07-03	AA	SFO	
3	2016	3	7	4	1	2016-07-04	AA	SFO	
4	2016	3	7	5	2	2016-07-05	AA	SFO	



Here are some more questions to help you learn more about your dataset.

Questions

1. How many rows and columns does the dataset have?
2. How many years are included in the dataset?
3. What is the date range for the dataset?
4. Which airlines are included in the dataset?
5. Which origin and destination airports are covered?

In [20]:

```
print("The #rows and #columns are ", data.shape[0], " and ", data.shape[1])
print("The years in this dataset are: ", list(data['Year'].unique()))
print("The months covered in this dataset are: ", sorted(list(data['Month'].unique())))
print("The date range for data is : ", min(data['FlightDate']), " to ", max(data['FlightDate']))
print("The airlines covered in this dataset are: ", list(data['Reporting_Airline'].unique()))
print("The Origin airports covered are: ", list(data['Origin'].unique()))
print("The Destination airports covered are: ", list(data['Dest'].unique()))
```

The #rows and #columns are 1658130 and 20
The years in this dataset are: [2016, 2015, 2018, 2014, 2017]
The months covered in this dataset are: [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]
The date range for data is : 2014-01-01 to 2018-12-31
The airlines covered in this dataset are: ['AA', 'DL', 'UA', 'WN', 'OO']
The Origin airports covered are: ['SFO', 'DFW', 'LAX', 'DEN', 'ORD', 'IAH', 'PHX', 'CLT', 'ATL']
The Destination airports covered are: ['DFW', 'SFO', 'ORD', 'DEN', 'CLT', 'LAX', 'PHX', 'IAH', 'ATL']

Define your target column: **is_delay** (1* means that the arrival time delayed more than 15 minutes, and *0

means all other cases). To rename the column from **ArrDel15** to *is_delay*, use the `rename` method.

Hint: You can use the `rename` function in pandas ([pandas.DataFrame.rename documentation](https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.rename.html) (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.rename.html>)).

For example:

```
data.rename(columns={'col1':'column1'}, inplace=True)
```

In [21]:

```
data.rename(columns={'ArrDel15':'is_delay'}, inplace=True)
```

Look for nulls across columns. You can use the `isnull()` function ([pandas.isnull documentation](https://pandas.pydata.org/pandas-docs/version/0.17.0/generated/pandas.isnull.html) (<https://pandas.pydata.org/pandas-docs/version/0.17.0/generated/pandas.isnull.html>)).

Hint: `isnull()` detects whether the particular value is null or not. It returns a boolean (*True* or *False*) in its place. To sum the number of columns, use the `sum(axis=0)` function (for example, `df.isnull().sum(axis=0)`).

In [22]:

```
data.isnull().sum(axis=0)
```

Out[22]:

Year	0
Quarter	0
Month	0
DayofMonth	0
DayOfWeek	0
FlightDate	0
Reporting_Airline	0
Origin	0
OriginState	0
Dest	0
DestState	0
CRSDepTime	0
Cancelled	0
Diverted	0
Distance	0
DistanceGroup	0
ArrDelay	22540
ArrDelayMinutes	22540
is_delay	22540
AirTime	22540
dtype: int64	

The arrival delay details and airtime are missing for 22,540 out of 1,658,130 rows, which is 1.3 percent. You can either remove or impute these rows. The documentation doesn't mention any information about missing rows.

In [23]:

```
### Remove null columns
data = data[~data.is_delay.isnull()]
data.isnull().sum(axis = 0)
```

Out[23]:

Year	0
Quarter	0
Month	0
DayofMonth	0
DayOfWeek	0
FlightDate	0
Reporting_Airline	0
Origin	0
OriginState	0
Dest	0
DestState	0
CRSDepTime	0
Cancelled	0
Diverted	0
Distance	0
DistanceGroup	0
ArrDelay	0
ArrDelayMinutes	0
is_delay	0
AirTime	0
dtype:	int64

Get the hour of the day in 24-hour-time format from CRSDepTime.

In [24]:

```
data['DepHourOfDay'] = (data['CRSDepTime']//100)
```

The ML problem statement

- Given a set of features, can you predict if a flight is going to be delayed more than 15 minutes?
- Because the target variable takes only a value of 0* or *1, you could use a classification algorithm.

Before you start modeling, it's a good practice to look at feature distribution, correlations, and others.

- This will give you an idea of any non-linearity or patterns in the data
 - Linear models: Add power, exponential, or interaction features
 - Try a non-linear model
- Data imbalance
 - Choose metrics that won't give biased model performance (accuracy versus the area under the curve, or AUC)
 - Use weighted or custom loss functions
- Missing data
 - Do imputation based on simple statistics -- mean, median, mode (numerical variables), frequent class (categorical variables)
 - Clustering-based imputation (k-nearest neighbors, or KNNs, to predict column value)

- Drop column

Data exploration

Check the classes *delay* versus *no delay*.

In [25]:

```
(data.groupby('is_delay').size()/len(data) ).plot(kind='bar')# Enter your code here  
plt.ylabel('Frequency')  
plt.title('Distribution of classes')  
plt.show()
```



Question: What can you deduce from the bar plot about the ratio of *delay* versus *no delay*?

In [26]:

```
# Enter your answer here  
#About 80% of flight were not delayed and only 20% of flights were delayed.
```

Run the following two cells and answer the questions.

In [32]:

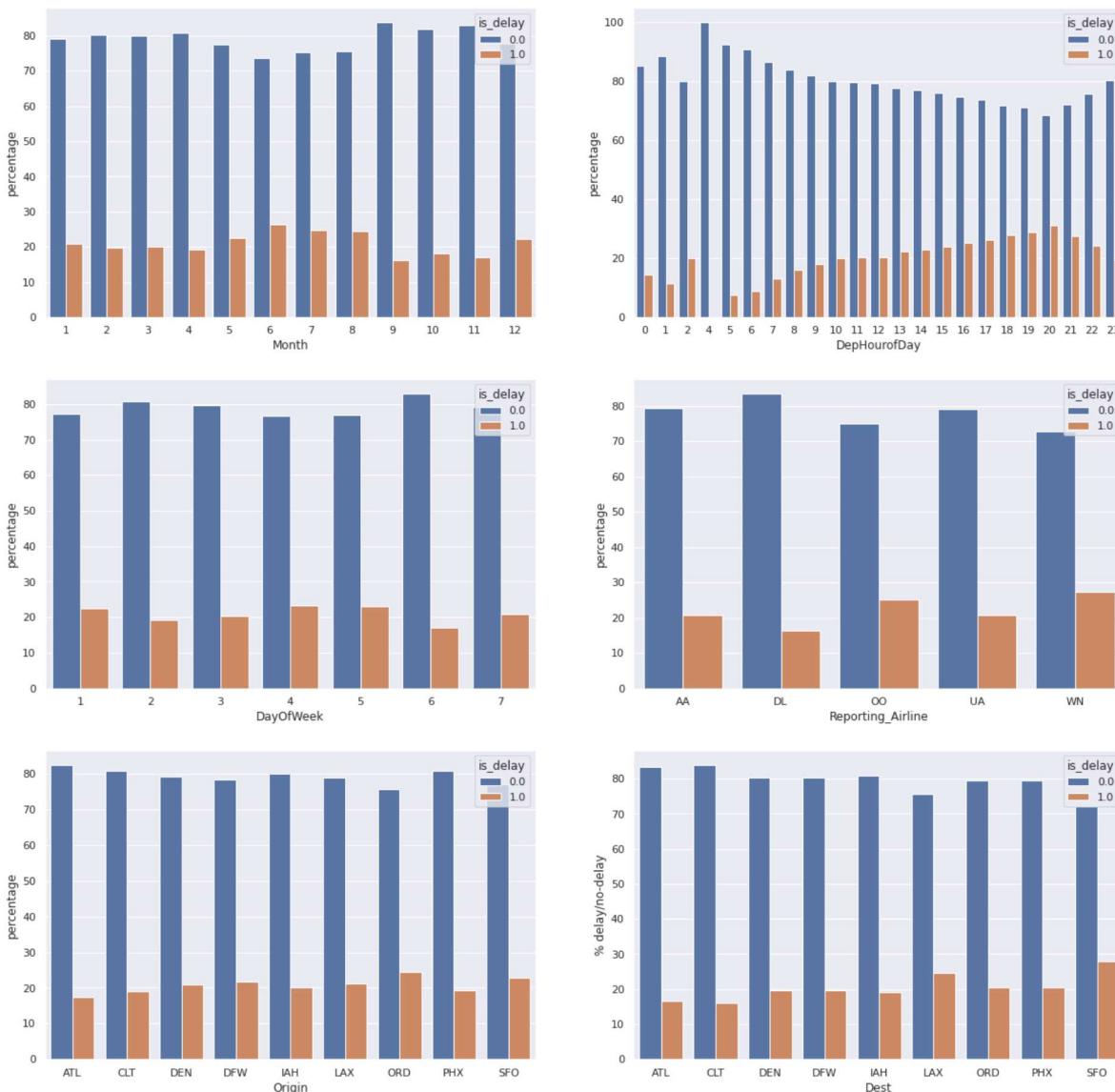
```

viz_columns = ['Month', 'DepHourOfDay', 'DayOfWeek', 'Reporting_Airline', 'Origin', 'Dest']
fig, axes = plt.subplots(3, 2, figsize=(20,20), squeeze=False)
# fig.autofmt_xdate(rotation=90)

for idx, column in enumerate(viz_columns):
    ax = axes[idx//2, idx%2]
    temp = data.groupby(column)['is_delay'].value_counts(normalize=True).rename('percentage')
    mul(100).reset_index().sort_values(column)
    sns.barplot(x=column, y="percentage", hue="is_delay", data=temp, ax=ax)
    plt.ylabel('% delay/no-delay')

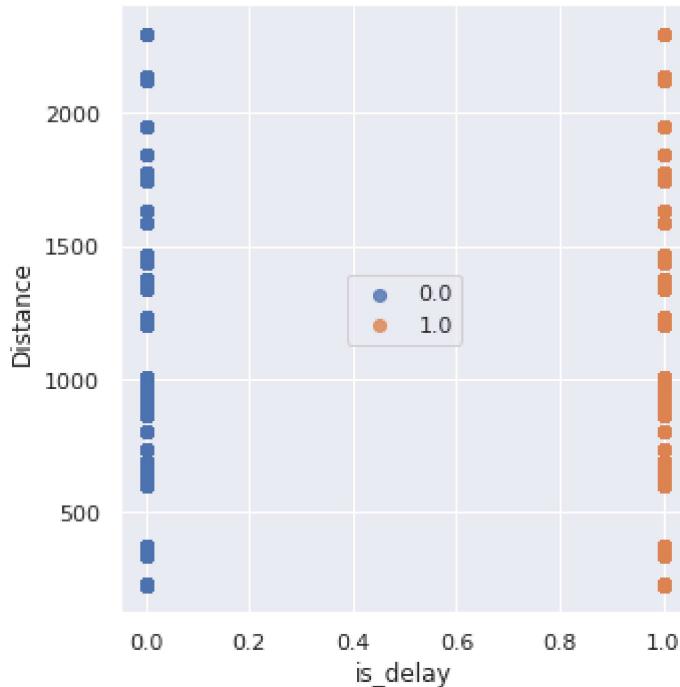
plt.show()

```



In [33]:

```
sns.lmplot( x="is_delay", y="Distance", data=data, fit_reg=False, hue='is_delay', legend=False)
plt.legend(loc='center')
plt.xlabel('is_delay')
plt.ylabel('Distance')
plt.show()
```



Questions

Using the data from the previous charts, answer these questions:

- Which months have the most delays?
- What time of the day has the most delays?
- What day of the week has the most delays?
- Which airline has the most delays?
- Which origin and destination airports have the most delays?
- Is flight distance a factor in the delays?

In [34]:

```
#1.June has the most delays
#2.20:00 has the most delays
#3.Thursday and Friday have the most delays
#4.Airline WN has the most delays
#5.Origin ORD and destination SFO have the most delays
#6.No, distance is not a factor in the delays.
```

Features

Look at all the columns and what their specific types are.

In [27]:

```
data.columns
```

Out[27]:

```
Index(['Year', 'Quarter', 'Month', 'DayofMonth', 'DayOfWeek', 'FlightDate',
       'Reporting_Airline', 'Origin', 'OriginState', 'Dest', 'DestState',
       'CRSDepTime', 'Cancelled', 'Diverted', 'Distance', 'DistanceGroup',
       'ArrDelay', 'ArrDelayMinutes', 'is_delay', 'AirTime', 'DepHourOfDay'],
      dtype='object')
```

In [28]:

```
data.dtypes
```

Out[28]:

Year	int64
Quarter	int64
Month	int64
DayofMonth	int64
DayOfWeek	int64
FlightDate	object
Reporting_Airline	object
Origin	object
OriginState	object
Dest	object
DestState	object
CRSDepTime	int64
Cancelled	float64
Diverted	float64
Distance	float64
DistanceGroup	int64
ArrDelay	float64
ArrDelayMinutes	float64
is_delay	float64
AirTime	float64
DepHourOfDay	int64
dtype: object	

Filtering the required columns:

- Date is redundant, because you have Year, Quarter, Month, DayofMonth, and DayOfWeek to describe the date.
- Use Origin and Dest codes instead of OriginState and DestState.
- Because you are only classifying whether the flight is delayed or not, you don't need TotalDelayMinutes, DepDelayMinutes, and ArrDelayMinutes.

Treat DepHourOfDay as a categorical variable because it doesn't have any quantitative relation with the target.

- If you needed to do a one-hot encoding of this variable, it would result in 23 more columns.

- Other alternatives to handling categorical variables include hash encoding, regularized mean encoding, and bucketizing the values, among others.
- In this case, you only need to split into buckets.

To change a column type to category, use the `astype` function ([pandas.DataFrame.astype documentation](https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.astype.html) (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.astype.html>)).

In [29]:

```
data_orig = data.copy()
data = data[[ 'is_delay', 'Quarter', 'Month', 'DayofMonth', 'DayOfWeek',
    'Reporting_Airline', 'Origin', 'Dest', 'Distance', 'DepHourOfDay']]
categorical_columns = ['Quarter', 'Month', 'DayofMonth', 'DayOfWeek',
    'Reporting_Airline', 'Origin', 'Dest', 'DepHourOfDay']
for c in categorical_columns:
    data[c] = data[c].astype('category')
```

To use one-hot encoding, use the `get_dummies` function in pandas for the categorical columns that you selected. Then, you can concatenate those generated features to your original dataset by using the `concat` function in pandas. For encoding categorical variables, you can also use *dummy encoding* by using a keyword `drop_first=True`. For more information about dummy encoding, see [Dummy variable \(statistics\)](https://en.wikiversity.org/wiki/Dummy_variable_(statistics)) ([https://en.wikiversity.org/wiki/Dummy_variable_\(statistics\)](https://en.wikiversity.org/wiki/Dummy_variable_(statistics))).

For example:

```
pd.get_dummies(df[['column1','column2']], drop_first=True)
```

In [30]:

```
data_dummies = pd.get_dummies(data[categorical_columns], drop_first=True)
data = pd.concat([data, data_dummies], axis = 1)
data.drop(categorical_columns, axis=1, inplace=True)
```

Check the length of the dataset and the new columns.

Hint: Use the `shape` and `columns` properties.

In [31]:

```
data.shape
```

Out[31]:

```
(1635590, 94)
```

In [32]:

```
data.columns
```

Out[32]:

```
Index(['is_delay', 'Distance', 'Quarter_2', 'Quarter_3', 'Quarter_4',
       'Month_2', 'Month_3', 'Month_4', 'Month_5', 'Month_6', 'Month_7',
       'Month_8', 'Month_9', 'Month_10', 'Month_11', 'Month_12',
       'DayofMonth_2', 'DayofMonth_3', 'DayofMonth_4', 'DayofMonth_5',
       'DayofMonth_6', 'DayofMonth_7', 'DayofMonth_8', 'DayofMonth_9',
       'DayofMonth_10', 'DayofMonth_11', 'DayofMonth_12', 'DayofMonth_13',
       'DayofMonth_14', 'DayofMonth_15', 'DayofMonth_16', 'DayofMonth_17',
       'DayofMonth_18', 'DayofMonth_19', 'DayofMonth_20', 'DayofMonth_21',
       'DayofMonth_22', 'DayofMonth_23', 'DayofMonth_24', 'DayofMonth_25',
       'DayofMonth_26', 'DayofMonth_27', 'DayofMonth_28', 'DayofMonth_29',
       'DayofMonth_30', 'DayofMonth_31', 'DayOfWeek_2', 'DayOfWeek_3',
       'DayOfWeek_4', 'DayOfWeek_5', 'DayOfWeek_6', 'DayOfWeek_7',
       'Reporting_Airline_DL', 'Reporting_Airline_00', 'Reporting_Airline_U
A',
       'Reporting_Airline_WN', 'Origin_CLT', 'Origin_DEN', 'Origin_DFW',
       'Origin_IAH', 'Origin_LAX', 'Origin_ORD', 'Origin_PHX', 'Origin_SFO',
       'Dest_CLT', 'Dest_DEN', 'Dest_DFW', 'Dest_IAH', 'Dest_LAX', 'Dest_OR
D',
       'Dest_PHX', 'Dest_SFO', 'DepHourOfDay_1', 'DepHourOfDay_2',
       'DepHourOfDay_4', 'DepHourOfDay_5', 'DepHourOfDay_6', 'DepHourOfDay_
7',
       'DepHourOfDay_8', 'DepHourOfDay_9', 'DepHourOfDay_10',
       'DepHourOfDay_11', 'DepHourOfDay_12', 'DepHourOfDay_13',
       'DepHourOfDay_14', 'DepHourOfDay_15', 'DepHourOfDay_16',
       'DepHourOfDay_17', 'DepHourOfDay_18', 'DepHourOfDay_19',
       'DepHourOfDay_20', 'DepHourOfDay_21', 'DepHourOfDay_22',
       'DepHourOfDay_23'],
      dtype='object')
```

You are now ready to train the model. Before you split the data, rename the **is_delay** column to *target*.

Hint: You can use the `rename` function in pandas ([pandas.DataFrame.rename documentation](https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.rename.html) (<https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.rename.html>)).

In [33]:

```
data.rename(columns = {'is_delay':'target'}, inplace=True )
```

End of Step 2

Save the project file to your local computer. Follow these steps:

1. In the file explorer on the left, right-click the notebook that you're working on.
2. Choose **Download**, and save the file locally.

This action downloads the current notebook to the default download folder on your computer.

Step 3: Model training and evaluation

You must include some preliminary steps when you convert the dataset from a DataFrame to a format that a machine learning algorithm can use. For Amazon SageMaker, you must perform these steps:

1. Split the data into `train_data`, `validation_data`, and `test_data` by using `sklearn.model_selection.train_test_split`.
2. Convert the dataset to an appropriate file format that the Amazon SageMaker training job can use. This can be either a CSV file or record protobuf. For more information, see [Common Data Formats for Training](https://docs.aws.amazon.com/sagemaker/latest/dg/cdf-training.html) (<https://docs.aws.amazon.com/sagemaker/latest/dg/cdf-training.html>).
3. Upload the data to your S3 bucket. If you haven't created one before, see [Create a Bucket](https://docs.aws.amazon.com/AmazonS3/latest/gsg/CreatingABucket.html) (<https://docs.aws.amazon.com/AmazonS3/latest/gsg/CreatingABucket.html>).

Use the following cells to complete these steps. Insert and delete cells where needed.

Project presentation: In your project presentation, write down the key decisions that you made in this phase.

Train-test split

In [34]:

```
from sklearn.model_selection import train_test_split
def split_data(data):
    train, test_and_validate = train_test_split(data, test_size=0.2, random_state=42, stratify=data['target'])
    test, validate = train_test_split(test_and_validate, test_size=0.5, random_state=42, stratify=test_and_validate['target'])
    return train, validate, test
```

In [35]:

```
train, validate, test = split_data(data)
print(train['target'].value_counts())
print(test['target'].value_counts())
print(validate['target'].value_counts())
```

```
0.0    1033806
1.0    274666
Name: target, dtype: int64
0.0    129226
1.0    34333
Name: target, dtype: int64
0.0    129226
1.0    34333
Name: target, dtype: int64
```

Sample answer

```
0.0    1033570
1.0    274902
Name: target, dtype: int64
0.0    129076
1.0    34483
Name: target, dtype: int64
0.0    129612
1.0    33947
Name: target, dtype: int64
```

Baseline classification model

In [36]:

```
import sagemaker
from sagemaker.serializers import CSVSerializer
from sagemaker.amazon.amazon_estimator import RecordSet
import boto3

# Instantiate the LinearLearner estimator object with 1 ml.m4.xlarge
classifier_estimator = sagemaker.LinearLearner(role=sagemaker.get_execution_role(),
                                                instance_count=1,
                                                instance_type='ml.m4.xlarge',
                                                predictor_type='binary_classifier',
                                                binary_classifier_model_selection_criteria='
```

Sample code

```
num_classes = len(pd.unique(train_labels))
classifier_estimator = sagemaker.LinearLearner(role=sagemaker.get_execution_role(),
(),                                                 instance_count=1,
                                                 instance_type='ml.m4.xlarge',
                                                 predictor_type='binary_classifier',
                                                 binary_classifier_model_selection_cr
iteria = 'cross_entropy_loss')
```

Linear learner accepts training data in protobuf or CSV content types. It also accepts inference requests in protobuf, CSV, or JavaScript Object Notation (JSON) content types. Training data has features and ground-truth labels, but the data in an inference request has only features.

In a production pipeline, AWS recommends converting the data to the Amazon SageMaker protobuf format and storing it in Amazon S3. To get up and running quickly, AWS provides the `record_set` operation for converting and uploading the dataset when it's small enough to fit in local memory. It accepts NumPy arrays like the ones you already have, so you will use it for this step. The `RecordSet` object will track the temporary Amazon S3 location of your data. Create train, validation, and test records by using the `estimator.record_set` function. Then, start your training job by using the `estimator.fit` function.

In [37]:

```
### Create train, validate, and test records
train_records = classifier_estimator.record_set(train.values[:, 1:]).astype(np.float32), tra
val_records = classifier_estimator.record_set(validate.values[:, 1:]).astype(np.float32), va
test_records = classifier_estimator.record_set(test.values[:, 1:]).astype(np.float32), test.
```

Now, train your model on the dataset that you just uploaded.

Sample code

```
linear.fit([train_records, val_records, test_records])
```

In [39]:

```
### Fit the classifier
classifier_estimator.fit([train_records, val_records, test_records])
```

Defaulting to the only supported framework/algorithm version: 1. Ignoring
framework/algorithm version: 1.
Defaulting to the only supported framework/algorithm version: 1. Ignoring
framework/algorithm version: 1.

Model evaluation

In this section, you will evaluate your trained model.

First, examine the metrics for the training job:

In [40]:

```
sagemaker.analytics.TrainingJobAnalytics(classifier_estimator._current_job_name,
                                         metric_names = ['test:objective_loss',
                                                          'test:binary_f_beta',
                                                          'test:precision',
                                                          'test:recall'])
                                         ).dataframe()
```

Out[40]:

	timestamp	metric_name	value
0	0.0	test:objective_loss	0.491039
1	0.0	test:binary_f_beta	0.008497
2	0.0	test:precision	0.548507
3	0.0	test:recall	0.004282

Next, set up some functions that will help load the test data into Amazon S3 and perform a prediction by using the batch prediction function. Using batch prediction will help reduce costs because the instances will only run when predictions are performed on the supplied test data.

Note: Replace <LabBucketName> with the name of the lab bucket that was created during the lab setup.

In [41]:

```
import io
bucket='sagemaker-us-east-1-743068201958'
prefix='flight-linear'
train_file='flight_train.csv'
test_file='flight_test.csv'
validate_file='flight_validate.csv'
whole_file='flight.csv'
s3_resource = boto3.Session().resource('s3')

def upload_s3_csv(filename, folder, dataframe):
    csv_buffer = io.StringIO()
    dataframe.to_csv(csv_buffer, header=False, index=False)
    s3_resource.Bucket(bucket).Object(os.path.join(prefix, folder, filename)).put(Body=csv_
```

In [42]:

```
def batch_linear_predict(test_data, estimator):
    batch_X = test_data.iloc[:,1:];
    batch_X_file='batch-in.csv'
    upload_s3_csv(batch_X_file, 'batch-in', batch_X)

    batch_output = "s3://{}/{}/batch-out/".format(bucket,prefix)
    batch_input = "s3://{}/{}/batch-in/{}".format(bucket,prefix,batch_X_file)

    classifier_transformer = estimator.transformer(instance_count=1,
                                                    instance_type='ml.m4.xlarge',
                                                    strategy='MultiRecord',
                                                    assemble_with='Line',
                                                    output_path=batch_output)

    classifier_transformer.transform(data=batch_input,
                                    data_type='S3Prefix',
                                    content_type='text/csv',
                                    split_type='Line')

    classifier_transformer.wait()

    s3 = boto3.client('s3')
    obj = s3.get_object(Bucket=bucket, Key="{}{}".format(prefix, 'batch-in.csv'))
    target_predicted_df = pd.read_json(io.BytesIO(obj['Body'].read()),orient="records",lines=True)
    return test_data.iloc[:,0], target_predicted_df.iloc[:,0]
```

To run the predictions on the test dataset, run the `batch_linear_predict` function (which was defined previously) on your test dataset.

In [50]:

```
test_labels, target_predicted = batch_linear_predict(test, classifier_estimator)
```

Defaulting to the only supported framework/algorithm version: 1. Ignoring framework/algorithm version: 1.

To view a plot of the confusion matrix, and various scoring metrics, create a couple of functions:

In [82]:

```
from sklearn.metrics import confusion_matrix

def plot_confusion_matrix(test_labels, target_predicted):
    matrix = confusion_matrix(test_labels, target_predicted)
    df_confusion = pd.DataFrame(matrix)
    colormap = sns.color_palette("BrBG", 10)
    sns.heatmap(df_confusion, annot=True, fmt=' .2f', cbar=None, cmap=colormap)
    plt.title("Confusion Matrix")
    plt.tight_layout()
    plt.ylabel("True Class")
    plt.xlabel("Predicted Class")
    plt.show()
```

In [85]:

```

from sklearn import metrics

def plot_roc(test_labels, target_predicted):
    TN, FP, FN, TP = confusion_matrix(test_labels, target_predicted).ravel()
    # Sensitivity, hit rate, recall, or true positive rate
    Sensitivity = float(TP)/(TP+FN)*100
    # Specificity or true negative rate
    Specificity = float(TN)/(TN+FP)*100
    # Precision or positive predictive value
    Precision = float(TP)/(TP+FP)*100
    # Negative predictive value
    NPV = float(TN)/(TN+FN)*100
    # Fall out or false positive rate
    FPR = float(FP)/(FP+TN)*100
    # False negative rate
    FNR = float(FN)/(TP+FN)*100
    # False discovery rate
    FDR = float(FP)/(TP+FP)*100
    # Overall accuracy
    ACC = float(TP+TN)/(TP+FP+FN+TN)*100

    print("Sensitivity or TPR: ", Sensitivity, "%")
    print("Specificity or TNR: ", Specificity, "%")
    print("Precision: ", Precision, "%")
    print("Negative Predictive Value: ", NPV, "%")
    print("False Positive Rate: ", FPR, "%")
    print("False Negative Rate: ", FNR, "%")
    print("False Discovery Rate: ", FDR, "%")
    print("Accuracy: ", ACC, "%")

    test_labels = test.iloc[:,0];
    print("Validation AUC", metrics.roc_auc_score(test_labels, target_predicted) )

    fpr, tpr, thresholds = metrics.roc_curve(test_labels, target_predicted)
    roc_auc = metrics.auc(fpr, tpr)

    plt.figure()
    plt.plot(fpr, tpr, label='ROC curve (area = %0.2f)' % (roc_auc))
    plt.plot([0, 1], [0, 1], 'k--')
    plt.xlim([0.0, 1.0])
    plt.ylim([0.0, 1.05])
    plt.xlabel('False Positive Rate')
    plt.ylabel('True Positive Rate')
    plt.title('Receiver operating characteristic')
    plt.legend(loc="lower right")

    # create the axis of thresholds (scores)
    ax2 = plt.gca().twinx()
    ax2.plot(fpr, thresholds, markeredgecolor='r', linestyle='dashed', color='r')
    ax2.set_ylabel('Threshold', color='r')
    ax2.set_ylim([thresholds[-1], thresholds[0]])
    ax2.set_xlim([fpr[0], fpr[-1]])

    print(plt.figure())

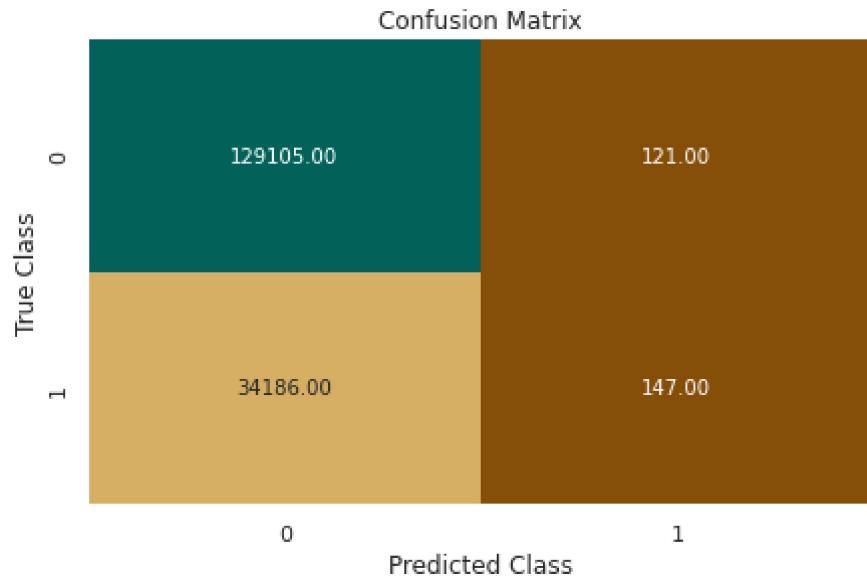
```

To plot the confusion matrix, call the `plot_confusion_matrix` function on the `test_labels` and the

target_predicted data from your batch job:

In [53]:

```
plot_confusion_matrix(test_labels, target_predicted)
```

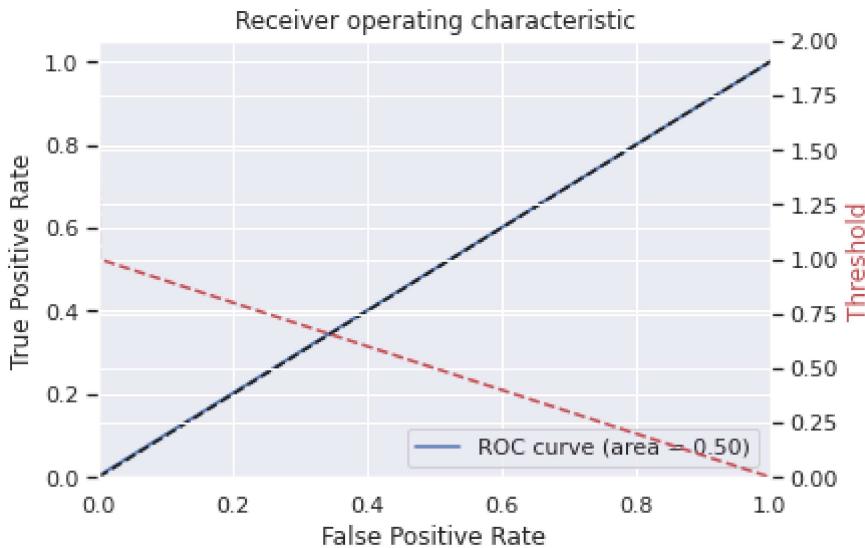


To print statistics and plot a receiver operating characteristic (ROC) curve, call the `plot_roc` function on the `test_labels` and `target_predicted` data from your batch job:

In [54]:

```
plot_roc(test_labels, target_predicted)
```

Sensitivity or TPR: 0.42815949669414266 %
 Specificity or TNR: 99.90636559206352 %
 Precision: 54.850746268656714 %
 Negative Predictive Value: 79.06436974481142 %
 False Positive Rate: 0.09363440793648337 %
 False Negative Rate: 99.57184050330585 %
 False Discovery Rate: 45.149253731343286 %
 Accuracy: 79.02469445276627 %
 Validation AUC 0.5016726254437883
 Figure(432x288)



<Figure size 432x288 with 0 Axes>

Key questions to consider:

1. How does your model's performance on the test set compare to its performance on the training set? What can you deduce from this comparison?
2. Are there obvious differences between the outcomes of metrics like accuracy, precision, and recall? If so, why might you be seeing those differences?
3. Given your business situation and goals, which metric (or metrics) is the most important for you to consider? Why?
4. From a business standpoint, is the outcome for the metric (or metrics) that you consider to be the most important sufficient for what you need? If not, what are some things you might change in your next iteration? (This will happen in the feature engineering section, which is next.)

Use the following cells to answer these (and other) questions. Insert and delete cells where needed.

Project presentation: In your project presentation, write down your answers to these questions -- and other similar questions that you might answer -- in this section. Record the key details and decisions that you made.

Question: What can you summarize from the confusion matrix?

In [55]:

```
# Enter your answer here
#The data is imbalanced. The amount of delayed flights data is small and it could cause some
#The accuracy of this model is only 79%, which means this model is not very good.
```

End of Step 3

Save the project file to your local computer. Follow these steps:

1. In the file explorer on the left, right-click the notebook that you're working on.
2. Select **Download**, and save the file locally.

This action downloads the current notebook to the default download folder on your computer.

Iteration II

Step 4: Feature engineering

You have now gone through one iteration of training and evaluating your model. Given that the first outcome that you reached for your model probably wasn't sufficient for solving your business problem, what could you change about your data to possibly improve model performance?

Key questions to consider:

1. How might the balance of your two main classes (*delay* and *no delay*) impact model performance?
2. Do you have any features that are correlated?
3. At this stage, could you perform any feature-reduction techniques that might have a positive impact on model performance?
4. Can you think of adding some more data or datasets?
5. After performing some feature engineering, how does the performance of your model compare to the first iteration?

Use the following cells to perform specific feature-engineering techniques that you think could improve your model performance (use the previous questions as a guide). Insert and delete cells where needed.

Project presentation: In your project presentation, record your key decisions and the methods that you use in this section. Also include any new performance metrics that you obtain after you evaluate your model again.

Before you start, think about why the precision and recall are around 80 percent, and the accuracy is at 99 percent.

Add more features:

1. Holidays
2. Weather

Because the list of holidays from 2014 to 2018 is known, you can create an indicator variable `is_holiday` to mark them.

The hypothesis is that airplane delays could be higher during holidays compared to the rest of the days. Add a boolean variable `is_holiday` that includes the holidays for the years 2014-2018.

In [43]:

```
# Source: http://www.calendarpedia.com/holidays/federal-holidays-2014.html

holidays_14 = ['2014-01-01', '2014-01-20', '2014-02-17', '2014-05-26', '2014-07-04', '2014-09-07', '2014-11-27']
holidays_15 = ['2015-01-01', '2015-01-19', '2015-02-16', '2015-05-25', '2015-06-03', '2015-07-04', '2015-09-07', '2015-11-26']
holidays_16 = ['2016-01-01', '2016-01-18', '2016-02-15', '2016-05-30', '2016-07-04', '2016-09-06', '2016-11-25']
holidays_17 = ['2017-01-02', '2017-01-16', '2017-02-20', '2017-05-29', '2017-07-04', '2017-09-05', '2017-11-24']
holidays_18 = ['2018-01-01', '2018-01-15', '2018-02-19', '2018-05-28', '2018-07-04', '2018-09-03', '2018-11-23']
holidays = holidays_14 + holidays_15 + holidays_16 + holidays_17 + holidays_18

### Add indicator variable for holidays
data_orig['is_holiday'] = [1 if i in holidays else 0 for i in data_orig['FlightDate']]
```

Weather data was fetched from <https://www.ncei.noaa.gov/access/services/data/v1?dataset=daily-summaries&stations=USW00023174,USW00012960,USW00003017,USW00094846,USW00013874,USW0002301-01&endDate=2018-12-31> (<https://www.ncei.noaa.gov/access/services/data/v1?dataset=daily-summaries&stations=USW00023174,USW00012960,USW00003017,USW00094846,USW00013874,USW0002301-01&endDate=2018-12-31>).

This dataset has information on wind speed, precipitation, snow, and temperature for cities by their airport codes.

Question: Could bad weather because of rain, heavy winds, or snow lead to airplane delays? You will now check.

In [44]:

```
!aws s3 cp s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data2/daily-summaries.csv .
#!wget 'https://www.ncei.noaa.gov/access/services/data/v1?dataset=daily-summaries&stations=USW00023174,USW00012960,USW00003017,USW00094846,USW00013874,USW0002301-01&endDate=2018-12-31'
```

download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data2/daily-summaries.csv to project/data/daily-summaries.csv

Import the weather data that was prepared for the airport codes in the dataset. Use the following stations and airports for the analysis. Create a new column called `airport` that maps the weather station to the airport name.

In [45]:

```
weather = pd.read_csv('/home/ec2-user/SageMaker/project/data/daily-summaries.csv')
station = ['USW00023174', 'USW00012960', 'USW0003017', 'USW00094846', 'USW00013874', 'USW000232
airports = ['LAX', 'IAH', 'DEN', 'ORD', 'ATL', 'SFO', 'DFW', 'PHX', 'CLT']

### Map weather stations to airport code
station_map = {s:a for s,a in zip(station, airports)}
weather['airport'] = weather['STATION'].map(station_map)
```

From the **DATE** column, create another column called **MONTH**.

In [46]:

```
weather['MONTH'] = weather['DATE'].apply(lambda x: x.split('-')[1])
weather.head()
```

Out[46]:

	STATION	DATE	AWND	PRCP	SNOW	SNWD	TAVG	TMAX	TMIN	airport	MONTH
0	USW00023174	2014-01-01	16	0	NaN	NaN	131.0	178.0	78.0	LAX	01
1	USW00023174	2014-01-02	22	0	NaN	NaN	159.0	256.0	100.0	LAX	01
2	USW00023174	2014-01-03	17	0	NaN	NaN	140.0	178.0	83.0	LAX	01
3	USW00023174	2014-01-04	18	0	NaN	NaN	136.0	183.0	100.0	LAX	01
4	USW00023174	2014-01-05	18	0	NaN	NaN	151.0	244.0	83.0	LAX	01

Sample output

	STATION	DATE	AWND	PRCP	SNOW	SNWD	TAVG	TMAX	TMIN	airport	MONTH
0	USW00023174	2014-01-01	16	0	NaN	NaN	131.0	178.0	78.0	LAX	01
1	USW00023174	2014-01-02	22	0	NaN	NaN	159.0	256.0	100.0	LAX	01
2	USW00023174	2014-01-03	17	0	NaN	NaN	140.0	178.0	83.0	LAX	01
3	USW00023174	2014-01-04	18	0	NaN	NaN	136.0	183.0	100.0	LAX	01
4	USW00023174	2014-01-05	18	0	NaN	NaN	151.0	244.0	83.0	LAX	01

Analyze and handle the **SNOW** and **SNWD** columns for missing values by using `fillna()`. To check the missing values for all the columns, use the `isna()` function.

In [47]:

```
weather.SNOW.fillna(0, inplace=True)
weather.SNWD.fillna(0, inplace=True)
weather.isna().sum()
```

Out[47]:

STATION	0
DATE	0
AWND	0
PRCP	0
SNOW	0
SNWD	0
TAVG	62
TMAX	20
TMIN	20
airport	0
MONTH	0
	dtype: int64

Question: Print the index of the rows that have missing values for *TAVG*, *TMAX*, *TMIN*.

Hint: To find the rows that are missing, use the `isna()` function. Then, to get the index, use the list on the `idx` variable.

In [48]:

```
idx = np.array([i for i in range(len(weather))])
TAVG_idx = idx[weather.TAVG.isna()]
TMAX_idx = idx[weather.TMAX.isna()]
TMIN_idx = idx[weather.TMIN.isna()]
TAVG_idx
```

Out[48]:

```
array([ 3956,  3957,  3958,  3959,  3960,  3961,  3962,  3963,  3964,
       3965,  3966,  3967,  3968,  3969,  3970,  3971,  3972,  3973,
       3974,  3975,  3976,  3977,  3978,  3979,  3980,  3981,  3982,
       3983,  3984,  3985,  4017,  4018,  4019,  4020,  4021,  4022,
      4023,  4024,  4025,  4026,  4027,  4028,  4029,  4030,  4031,
      4032,  4033,  4034,  4035,  4036,  4037,  4038,  4039,  4040,
      4041,  4042,  4043,  4044,  4045,  4046,  4047, 13420])
```

Sample output

```
array([ 3956,  3957,  3958,  3959,  3960,  3961,  3962,  3963,  3964,
       3965,  3966,  3967,  3968,  3969,  3970,  3971,  3972,  3973,
       3974,  3975,  3976,  3977,  3978,  3979,  3980,  3981,  3982,
       3983,  3984,  3985,  4017,  4018,  4019,  4020,  4021,  4022,
      4023,  4024,  4025,  4026,  4027,  4028,  4029,  4030,  4031,
      4032,  4033,  4034,  4035,  4036,  4037,  4038,  4039,  4040,
      4041,  4042,  4043,  4044,  4045,  4046,  4047, 13420])
```

You can replace the missing *TAVG*, *TMAX*, and *TMIN* values with the average value for a particular station or airport. Because consecutive rows of *TAVG_idx* are missing, replacing them with a previous value would not be possible. Instead, replace them with the mean. Use the `groupby` function to aggregate the variables with a mean value.

Hint: Group by `MONTH` and `STATION`.

In [49]:

```
weather_impute = weather.groupby(['MONTH', 'STATION']).agg({'TAVG':'mean', 'TMAX':'mean', 'TMIN':'mean'})  
weather_impute.head(2)
```

Out[49]:

	MONTH	STATION	TAVG	TMAX	TMIN
0	01	USW00003017	-2.741935	74.000000	-69.858065
1	01	USW00003927	79.529032	143.767742	20.696774

Merge the mean data with the weather data.

In [50]:

```
weather = pd.merge(weather, weather_impute, how='left', left_on=['MONTH', 'STATION'], right_on=['MONTH', 'STATION'])  
.rename(columns = {'TAVG_y':'TAVG_AVG',  
                  'TMAX_y':'TMAX_AVG',  
                  'TMIN_y':'TMIN_AVG',  
                  'TAVG_x':'TAVG',  
                  'TMAX_x':'TMAX',  
                  'TMIN_x':'TMIN'})
```

Check for missing values again.

In [51]:

```
weather.TAVG[TAVG_idx] = weather.TAVG_AVG[TAVG_idx]
weather.TMAX[TMAX_idx] = weather.TMAX_AVG[TMAX_idx]
weather.TMIN[TMIN_idx] = weather.TMIN_AVG[TMIN_idx]
weather.isna().sum()
```

Out[51]:

STATION	0
DATE	0
AWND	0
PRCP	0
SNOW	0
SNWD	0
TAVG	0
TMAX	0
TMIN	0
airport	0
MONTH	0
TAVG_AVG	0
TMAX_AVG	0
TMIN_AVG	0
dtype: int64	

Drop STATION, MONTH, TAVG_AVG, TMAX_AVG, TMIN_AVG, TMAX, TMIN, SNWD from the dataset.

In [52]:

```
weather.drop(columns=['STATION', 'MONTH', 'TAVG_AVG', 'TMAX_AVG', 'TMIN_AVG', 'TMAX', 'TMIN', 'SNWD'])
```

Add the origin and destination weather conditions to the dataset.

In [53]:

```
### Add origin weather conditions
data_orig = pd.merge(data_orig, weather, how='left', left_on=['FlightDate', 'Origin'], right_rename(columns = {'AWND': 'AWND_O', 'PRCP': 'PRCP_O', 'TAVG': 'TAVG_O', 'SNOW': 'SNOW_O'})\n.drop(columns=['DATE', 'airport'])

### Add destination weather conditions
data_orig = pd.merge(data_orig, weather, how='left', left_on=['FlightDate', 'Dest'], right_rename(columns = {'AWND': 'AWND_D', 'PRCP': 'PRCP_D', 'TAVG': 'TAVG_D', 'SNOW': 'SNOW_D'})\n.drop(columns=['DATE', 'airport'])
```

Note: It's always a good practice to check for nulls or NAs after joins.

In [54]:

```
sum(data.isna().any())
```

Out[54]:

0

In [55]:

```
data_orig.columns
```

Out[55]:

```
Index(['Year', 'Quarter', 'Month', 'DayofMonth', 'DayOfWeek', 'FlightDate',
       'Reporting_Airline', 'Origin', 'OriginState', 'Dest', 'DestState',
       'CRSDepTime', 'Cancelled', 'Diverted', 'Distance', 'DistanceGroup',
       'ArrDelay', 'ArrDelayMinutes', 'is_delay', 'AirTime', 'DepHourOfDay',
       'is_holiday', 'AWND_O', 'PRCP_O', 'SNOW_O', 'TAVG_O', 'AWND_D',
       'PRCP_D', 'SNOW_D', 'TAVG_D'],
      dtype='object')
```

Convert the categorical data into numerical data by using one-hot encoding.

In [56]:

```
data = data_orig.copy()
data = data[['is_delay', 'Year', 'Quarter', 'Month', 'DayofMonth', 'DayOfWeek',
            'Reporting_Airline', 'Origin', 'Dest', 'Distance', 'DepHourOfDay', 'is_holiday', 'AWND_O',
            'TAVG_O', 'AWND_D', 'PRCP_D', 'TAVG_D', 'SNOW_O', 'SNOW_D']]

categorical_columns = ['Year', 'Quarter', 'Month', 'DayofMonth', 'DayOfWeek',
                       'Reporting_Airline', 'Origin', 'Dest', 'is_holiday']
for c in categorical_columns:
    data[c] = data[c].astype('category')
```

In [57]:

```
data_dummies = pd.get_dummies(data[['Year', 'Quarter', 'Month', 'DayofMonth', 'DayOfWeek']],
data = pd.concat([data, data_dummies], axis = 1)
data.drop(categorical_columns, axis=1, inplace=True)
```

Check the new columns.

In [58]:

```
data.shape
```

Out[58]:

```
(1635590, 86)
```

In [59]:

```
data.columns
```

Out[59]:

```
Index(['is_delay', 'Distance', 'DepHourOfDay', 'AWND_O', 'PRCP_O', 'TAVG_O',
       'AWND_D', 'PRCP_D', 'TAVG_D', 'SNOW_O', 'SNOW_D', 'Year_2015',
       'Year_2016', 'Year_2017', 'Year_2018', 'Quarter_2', 'Quarter_3',
       'Quarter_4', 'Month_2', 'Month_3', 'Month_4', 'Month_5', 'Month_6',
       'Month_7', 'Month_8', 'Month_9', 'Month_10', 'Month_11', 'Month_12',
       'DayofMonth_2', 'DayofMonth_3', 'DayofMonth_4', 'DayofMonth_5',
       'DayofMonth_6', 'DayofMonth_7', 'DayofMonth_8', 'DayofMonth_9',
       'DayofMonth_10', 'DayofMonth_11', 'DayofMonth_12', 'DayofMonth_13',
       'DayofMonth_14', 'DayofMonth_15', 'DayofMonth_16', 'DayofMonth_17',
       'DayofMonth_18', 'DayofMonth_19', 'DayofMonth_20', 'DayofMonth_21',
       'DayofMonth_22', 'DayofMonth_23', 'DayofMonth_24', 'DayofMonth_25',
       'DayofMonth_26', 'DayofMonth_27', 'DayofMonth_28', 'DayofMonth_29',
       'DayofMonth_30', 'DayofMonth_31', 'DayOfWeek_2', 'DayOfWeek_3',
       'DayOfWeek_4', 'DayOfWeek_5', 'DayOfWeek_6', 'DayOfWeek_7',
       'Reporting_Airline_DL', 'Reporting_Airline_OO', 'Reporting_Airline_U
A',
       'Reporting_Airline_WN', 'Origin_CLT', 'Origin_DEN', 'Origin_DFW',
       'Origin_IAH', 'Origin_LAX', 'Origin_ORD', 'Origin_PHX', 'Origin_SFO',
       'Dest_CLT', 'Dest_DEN', 'Dest_DFW', 'Dest_IAH', 'Dest_LAX', 'Dest_OR
D',
       'Dest_PHX', 'Dest_SFO', 'is_holiday_1'],
      dtype='object')
```

Sample output

```
Index(['Distance', 'DepHourofDay', 'is_delay', 'AWND_O', 'PRCP_O', 'TAVG_O',
       'AWND_D', 'PRCP_D', 'TAVG_D', 'SNOW_O', 'SNOW_D', 'Year_2015',
       'Year_2016', 'Year_2017', 'Year_2018', 'Quarter_2', 'Quarter_3',
       'Quarter_4', 'Month_2', 'Month_3', 'Month_4', 'Month_5', 'Month_6',
       'Month_7', 'Month_8', 'Month_9', 'Month_10', 'Month_11', 'Month_12',
       'DayofMonth_2', 'DayofMonth_3', 'DayofMonth_4', 'DayofMonth_5',
       'DayofMonth_6', 'DayofMonth_7', 'DayofMonth_8', 'DayofMonth_9',
       'DayofMonth_10', 'DayofMonth_11', 'DayofMonth_12', 'DayofMonth_13',
       'DayofMonth_14', 'DayofMonth_15', 'DayofMonth_16', 'DayofMonth_17',
       'DayofMonth_18', 'DayofMonth_19', 'DayofMonth_20', 'DayofMonth_21',
       'DayofMonth_22', 'DayofMonth_23', 'DayofMonth_24', 'DayofMonth_25',
       'DayofMonth_26', 'DayofMonth_27', 'DayofMonth_28', 'DayofMonth_29',
       'DayofMonth_30', 'DayofMonth_31', 'DayOfWeek_2', 'DayOfWeek_3',
       'DayOfWeek_4', 'DayOfWeek_5', 'DayOfWeek_6', 'DayOfWeek_7',
       'Reporting_Airline_DL', 'Reporting_Airline_00', 'Reporting_Airline_UA',
       'Reporting_Airline_WN', 'Origin_CLT', 'Origin_DEN', 'Origin_DFW',
       'Origin_IAH', 'Origin_LAX', 'Origin_ORD', 'Origin_PHX', 'Origin_SFO',
       'Dest_CLT', 'Dest_DEN', 'Dest_DFW', 'Dest_IAH', 'Dest_LAX', 'Dest_ORD',
       'Dest_PHX', 'Dest_SFO', 'is_holiday_1'],
      dtype='object')
```

Rename the `is_delay` column to `target` again. Use the same code that you used previously.

In [60]:

```
data.rename(columns = {'is_delay':'target'}, inplace=True )
```

Create the training sets again.

Hint: Use the `split_data` function that you defined (and used) earlier.

In [61]:

```
train, validate, test = split_data(data)
```

New baseline classifier

Now, see if these new features add any predictive power to the model.

In [55]:

```
# Instantiate the LinearLearner estimator object
classifier_estimator2 = sagemaker.LinearLearner(role=sagemaker.get_execution_role(),
                                                 instance_count=1,
                                                 instance_type='ml.m4.xlarge',
                                                 predictor_type='binary_classifier',
                                                 binary_classifier_model_selection_criteria =
```

Sample code

```
num_classes = len(pd.unique(train_labels))
classifier_estimator2 = sagemaker.LinearLearner(role=sagemaker.get_execution_role()
(),  
                                                instance_count=1,  
                                                instance_type='ml.m4.xlarge',
                                                predictor_type='binary_classifier',
                                                binary_classifier_model_selection_c
riteria = 'cross_entropy_loss')
```

In [76]:

```
train_records = classifier_estimator2.record_set(train.values[:, 1:]).astype(np.float32), train
val_records = classifier_estimator2.record_set(validate.values[:, 1:]).astype(np.float32), validate
test_records = classifier_estimator2.record_set(test.values[:, 1:]).astype(np.float32), test
```

Train your model by using the three datasets that you just created.

In [77]:

```
classifier_estimator2.fit([train_records, val_records,test_records])
```

Defaulting to the only supported framework/algorithm version: 1. Ignoring framework/algorithm version: 1.

Defaulting to the only supported framework/algorithm version: 1. Ignoring framework/algorithm version: 1.

Perform a batch prediction by using the newly trained model.

In [78]:

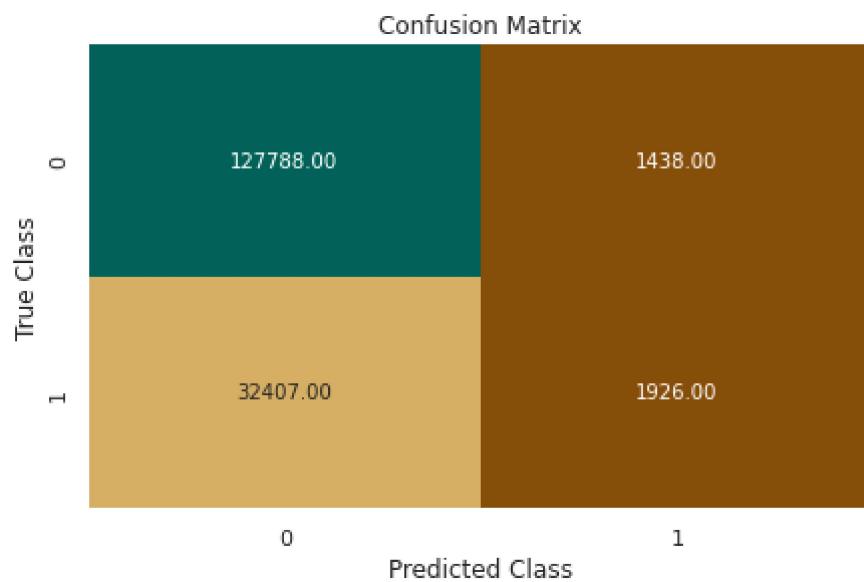
```
test_labels, target_predicted = batch_linear_predict(test, classifier_estimator2)
```

Defaulting to the only supported framework/algorithm version: 1. Ignoring framework/algorithm version: 1.

Plot a confusion matrix.

In [79]:

```
plot_confusion_matrix(test_labels, target_predicted)
```

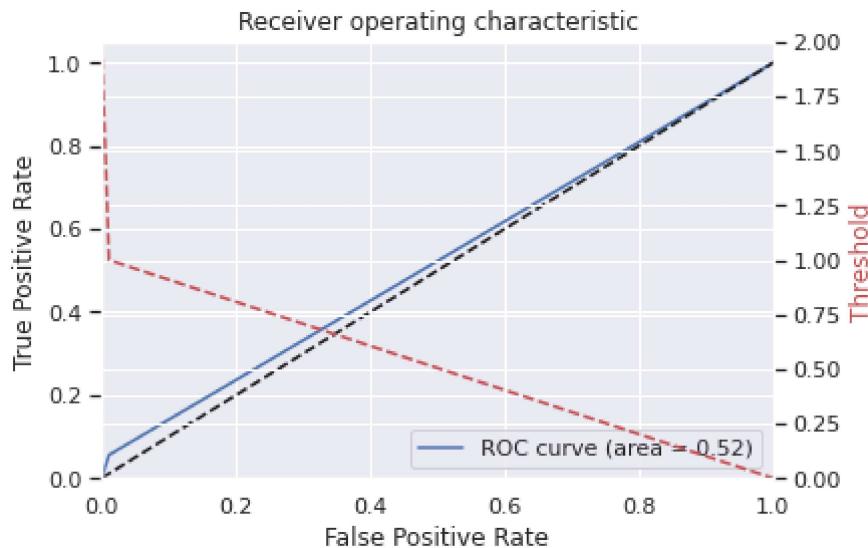


Plot the ROC curve.

In [80]:

```
plot_roc(test_labels, target_predicted)
```

```
Sensitivity or TPR:  5.609763201584482 %
Specificity or TNR:  98.88722083791187 %
Precision:  57.25326991676576 %
Negative Predictive Value:  79.77027997128499 %
False Positive Rate:  1.1127791620881247 %
False Negative Rate:  94.39023679841551 %
False Discovery Rate:  42.74673008323425 %
Accuracy:  79.3071613301622 %
Validation AUC 0.5224849201974817
Figure(432x288)
```



<Figure size 432x288 with 0 Axes>

The linear model shows only a little improvement in performance. Try a tree-based ensemble model, which is called *XGBoost*, with Amazon SageMaker.

Try the XGBoost model

Perform these steps:

1. Use the training set variables and save them as CSV files: train.csv, validation.csv and test.csv.
2. Store the bucket name in the variable. The Amazon S3 bucket name is provided to the left of the lab instructions.
 - a. `bucket = <LabBucketName>`
 - b. `prefix = 'flight-xgb'`
3. Use the AWS SDK for Python (Boto3) to upload the model to the bucket.

In [62]:

```
bucket='sagemaker-us-east-1-743068201958'
prefix='flight-xgb'
train_file='flight_train.csv'
test_file='flight_test.csv'
validate_file='flight_validate.csv'
whole_file='flight.csv'
s3_resource = boto3.Session().resource('s3')

def upload_s3_csv(filename, folder, dataframe):
    csv_buffer = io.StringIO()
    dataframe.to_csv(csv_buffer, header=False, index=False )
    s3_resource.Bucket(bucket).Object(os.path.join(prefix, folder, filename)).put(Body=csv_)

upload_s3_csv(train_file, 'train', train)
upload_s3_csv(test_file, 'test', test)
upload_s3_csv(validate_file, 'validate', validate)
```

Use the `sagemaker.inputs.TrainingInput` function to create a `record_set` for the training and validation datasets.

In [63]:

```
train_channel = sagemaker.inputs.TrainingInput(
    "s3://{}/{}/train/".format(bucket,prefix,train_file),
    content_type='text/csv')

validate_channel = sagemaker.inputs.TrainingInput(
    "s3://{}/{}/validate/".format(bucket,prefix,validate_file),
    content_type='text/csv')

data_channels = {'train': train_channel, 'validation': validate_channel}
```

In [64]:

```
from sagemaker.image_uris import retrieve
container = retrieve('xgboost',boto3.Session().region_name,'1.0-1')
```

In [65]:

```
sess = sagemaker.Session()
s3_output_location="s3://{} / {} / output/".format(bucket,prefix)

xgb = sagemaker.estimator.Estimator(container,
                                      role = sagemaker.get_execution_role(),
                                      instance_count=1,
                                      instance_type=instance_type,
                                      output_path=s3_output_location,
                                      sagemaker_session=sess)

xgb.set_hyperparameters(max_depth=5,
                       eta=0.2,
                       gamma=4,
                       min_child_weight=6,
                       subsample=0.8,
                       silent=0,
                       objective='binary:logistic',
                       eval_metric = "auc",
                       num_round=100)

xgb.fit(inputs=data_channels)
```

```
2022-08-18 14:43:26 Starting - Starting the training job...
2022-08-18 14:43:51 Starting - Preparing the instances for training
ProfilerReport-1660833805: InProgress
.....
2022-08-18 14:45:09 Downloading - Downloading input data...
2022-08-18 14:45:50 Training - Downloading the training image.....
2022-08-18 14:47:10 Training - Training image download completed. Training in progress.
INFO:sagemaker-containers:Imported framework sagemaker_xgboost_container.training
INFO:sagemaker-containers:Failed to parse hyperparameter eval_metric value auc to Json.
Returning the value itself
INFO:sagemaker-containers:Failed to parse hyperparameter objective value binary:logistic to Json.
Returning the value itself
INFO:sagemaker-containers:No GPUs detected (normal if no gpus installed)
INFO:sagemaker_xgboost_container.training:Running XGBoost Sagemaker in algorithm mode
INFO:root:Determined delimiter of CSV input is ','
INFO:root:Determined delimiter of CSV input is ','
```

Use the batch transformer for your new model, and evaluate the model on the test dataset.

In [66]:

```
batch_X = test.iloc[:,1:];
batch_X_file='batch-in.csv'
upload_s3_csv(batch_X_file, 'batch-in', batch_X)
```

In [67]:

```
batch_output = "s3://{}/{}/batch-out/".format(bucket,prefix)
batch_input = "s3://{}/{}/batch-in/{}".format(bucket,prefix,batch_X_file)

xgb_transformer = xgb.transformer(instance_count=1,
                                  instance_type=instance_type,
                                  strategy='MultiRecord',
                                  assemble_with='Line',
                                  output_path=batch_output)

xgb_transformer.transform(data=batch_input,
                         data_type='S3Prefix',
                         content_type='text/csv',
                         split_type='Line')
xgb_transformer.wait()
```

```
.....[2022-08-18:15:01:23:INFO] No GPUs detected (normal if no gpus installed)
[2022-08-18:15:01:23:INFO] No GPUs detected (normal if no gpus installed)
[2022-08-18:15:01:23:INFO] nginx config:
worker_processes auto;
daemon off;
pid /tmp/nginx.pid;
error_log /dev/stderr;
worker_rlimit_nofile 4096;
events {
    worker_connections 2048;
}
http {
    include /etc/nginx/mime.types;
    default_type application/octet-stream;
    access_log /dev/stdout combined;
    upstream gunicorn {
        server unix:/tmp/gunicorn.sock;
    }
}
```

Get the predicted target and test labels.

In [68]:

```
s3 = boto3.client('s3')
obj = s3.get_object(Bucket=bucket, Key="{}#/batch-out/{}".format(prefix, 'batch-in.csv.out'))
target_predicted = pd.read_csv(io.BytesIO(obj['Body'].read()), ',', names=['target'])
test_labels = test.iloc[:,0]
```

Calculate the predicted values based on the defined threshold.

Note: The predicted target will be a score, which must be converted to a binary class.

In [69]:

```
print(target_predicted.head())

def binary_convert(x):
    threshold = 0.55
    if x > threshold:
        return 1
    else:
        return 0

target_predicted['target'] = target_predicted['target'].apply(binary_convert)

test_labels = test.iloc[:,0]

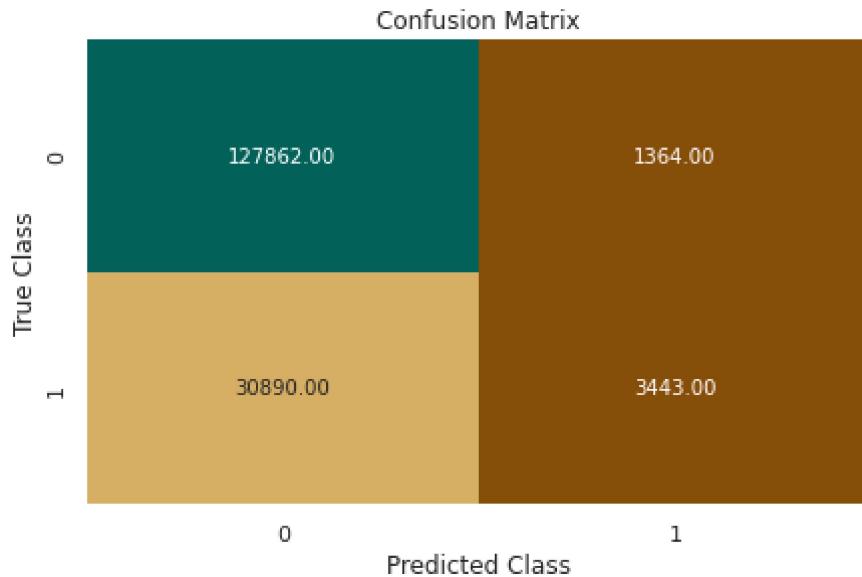
print(target_predicted.head())
```

```
target
0 0.064856
1 0.082722
2 0.719806
3 0.345698
4 0.249509
target
0      0
1      0
2      1
3      0
4      0
```

Plot a confusion matrix for your `target_predicted` and `test_labels`.

In [89]:

```
plot_confusion_matrix(test_labels, target_predicted)
```

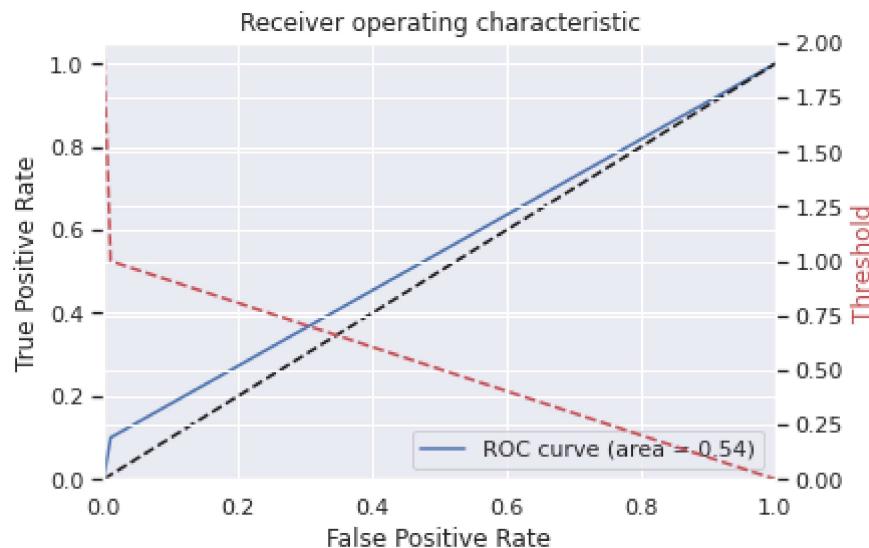


Plot the ROC chart:

In [90]:

```
plot_roc(test_labels, target_predicted)
```

Sensitivity or TPR: 10.028252701482538 %
 Specificity or TNR: 98.94448485598873 %
 Precision: 71.62471395881008 %
 Negative Predictive Value: 80.54197742390647 %
 False Positive Rate: 1.0555151440112671 %
 False Negative Rate: 89.97174729851746 %
 False Discovery Rate: 28.37528604118993 %
 Accuracy: 80.27989899669232 %
 Validation AUC 0.5448636877873564
 Figure(432x288)



<Figure size 432x288 with 0 Axes>

Try different thresholds

Question: Based on how well the model handled the test set, what can you conclude?

In [91]:

```
#Enter your answer here
#With the addition of both weather and holiday features, we can see that the accuracy has improved.
#But the XGBoost model also shows only a little improvement in performance.
#We should consider hyperparameter optimization to achieve better performance.
```

Hyperparameter optimization (HPO)

In [72]:

```
from sagemaker.tuner import IntegerParameter, CategoricalParameter, ContinuousParameter, HyperparameterTuner
### You can spin up multiple instances to do hyperparameter optimization in parallel

xgb = sagemaker.estimator.Estimator(container,
                                      role=sagemaker.get_execution_role(),
                                      instance_count=1, # make sure you have a limit set for
                                      instance_type=instance_type,
                                      output_path='s3://{}/{}/output'.format(bucket, prefix),
                                      sagemaker_session=sess)

xgb.set_hyperparameters(eval_metric='auc',
                        objective='binary:logistic',
                        num_round=100,
                        rate_drop=0.3,
                        tweedie_variance_power=1.4)

hyperparameter_ranges = {'alpha': ContinuousParameter(0, 1000, scaling_type='Linear'),
                        'eta': ContinuousParameter(0.1, 0.5, scaling_type='Linear'),
                        'min_child_weight': ContinuousParameter(3, 10, scaling_type='Linear'),
                        'subsample': ContinuousParameter(0.5, 1),
                        'num_round': IntegerParameter(10,150)}

objective_metric_name = 'validation:auc'

tuner = HyperparameterTuner(xgb,
                             objective_metric_name,
                             hyperparameter_ranges,
                             max_jobs=10, # Set this to 10 or above depending upon budget and time
                             max_parallel_jobs=1)
```

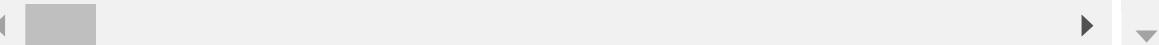
In [74]:

```
tuner.fit(inputs=data_channels)
tuner.wait()
```

No finished training job found associated with this estimator. Please make sure this estimator is only used for building workflow config
No finished training job found associated with this estimator. Please make sure this estimator is only used for building workflow config

.....

!



Wait until the training job is finished. It might take 25-30 minutes.

To monitor hyperparameter optimization jobs:

1. In the AWS Management Console, on the **Services** menu, choose **Amazon SageMaker**.
2. Choose **Training > Hyperparameter tuning jobs**.
3. You can check the status of each hyperparameter tuning job, its objective metric value, and its logs.

Check that the job completed successfully.

In [75]:

```
boto3.client('sagemaker').describe_hyper_parameter_tuning_job(  
    HyperParameterTuningJobName=tuner.latest_tuning_job.job_name)[ 'HyperParameterTuningJobs'
```

Out[75]:

```
'Completed'
```

The hyperparameter tuning job will have a model that worked the best. You can get the information about that model from the tuning job.

In [76]:

```
sage_client = boto3.Session().client('sagemaker')
tuning_job_name = tuner.latest_tuning_job.job_name
print(f'tuning job name:{tuning_job_name}')
tuning_job_result = sage_client.describe_hyper_parameter_tuning_job(HyperParameterTuningJob)
best_training_job = tuning_job_result['BestTrainingJob']
best_training_job_name = best_training_job['TrainingJobName']
print(f"best training job: {best_training_job_name}")

best_estimator = tuner.best_estimator()

tuner_df = sagemaker.HyperparameterTuningJobAnalytics(tuning_job_name).dataframe()
tuner_df.head()
```

tuning job name:sagemaker-xgboost-220818-1503
 best training job: sagemaker-xgboost-220818-1503-007-f362e223

2022-08-18 16:06:15 Starting - Preparing the instances for training
 2022-08-18 16:06:15 Downloading - Downloading input data
 2022-08-18 16:06:15 Training - Training image download completed. Training in progress.
 2022-08-18 16:06:15 Uploading - Uploading generated training model
 2022-08-18 16:06:15 Completed - Training job completed

Out[76]:

	alpha	eta	min_child_weight	num_round	subsample	TrainingJobName	TrainingJobStatus
0	979.302463	0.466431	3.616380	64.0	0.971856	sagemaker-xgboost-220818-1503-010-9b33c1ff	C
1	978.249131	0.235477	3.728918	140.0	0.877761	sagemaker-xgboost-220818-1503-009-d7b09904	C
2	407.627044	0.102146	4.171985	77.0	0.946617	sagemaker-xgboost-220818-1503-008-769e8c8b	C
3	773.666107	0.249259	9.293058	99.0	0.908308	sagemaker-xgboost-220818-1503-007-f362e223	C
4	695.736779	0.239366	8.543048	47.0	0.619368	sagemaker-xgboost-220818-1503-006-0798d82f	C

Use the estimator `best_estimator` and train it by using the data.

Tip: See the previous XGBoost estimator fit function.

In [77]:

```
best_estimator.fit(inputs = data_channels)
```

```
2022-08-18 16:35:00 Starting - Starting the training job...
2022-08-18 16:35:24 Starting - Preparing the instances for training
ProfilerReport-1660840500: InProgress
.....
2022-08-18 16:36:45 Downloading - Downloading input data...
2022-08-18 16:37:26 Training - Downloading the training image.....
2022-08-18 16:38:24 Training - Training image download completed. Training in progress.
INFO:sagemaker-containers:Imported framework sagemaker_xgboost_container.training
INFO:sagemaker-containers:Failed to parse hyperparameter _tuning_objective_metric value validation:auc to Json.
Returning the value itself
INFO:sagemaker-containers:Failed to parse hyperparameter eval_metric value auc to Json.
Returning the value itself
INFO:sagemaker-containers:Failed to parse hyperparameter objective value binary:logistic to Json.
Returning the value itself
INFO:sagemaker-containers:No GPUs detected (normal if no gpus installed)
```

Use the batch transformer for your new model, and evaluate the model on the test dataset.

In [78]:

```
batch_output = "s3://{}/{}/batch-out/".format(bucket,prefix)
batch_input = "s3://{}/{}/batch-in/{}".format(bucket,prefix,batch_X_file)

xgb_transformer = best_estimator.transformer(instance_count=1,
                                             instance_type=instance_type,
                                             strategy='MultiRecord',
                                             assemble_with='Line',
                                             output_path=batch_output)

xgb_transformer.transform(data=batch_input,
                         data_type='S3Prefix',
                         content_type='text/csv',
                         split_type='Line')
xgb_transformer.wait()
```

```
.....
[2022-08-18:16:50:35:INFO] No GPUs detected (normal if no gpus installed)
[2022-08-18:16:50:35:INFO] No GPUs detected (normal if no gpus installed)
[2022-08-18:16:50:35:INFO] nginx config:
worker_processes auto;
daemon off;
pid /tmp/nginx.pid;
error_log /dev/stderr;
[2022-08-18:16:50:35:INFO] No GPUs detected (normal if no gpus installed)
[2022-08-18:16:50:35:INFO] No GPUs detected (normal if no gpus installed)
[2022-08-18:16:50:35:INFO] nginx config:
worker_processes auto;
daemon off;
pid /tmp/nginx.pid;
error_log /dev/stderr;
worker_rlimit_nofile 4096;
events {
    worker_connections 2048;
}
....
```

In [79]:

```
s3 = boto3.client('s3')
obj = s3.get_object(Bucket=bucket, Key="{}{}/batch-out/{}".format(prefix,'batch-in.csv.out'))
target_predicted = pd.read_csv(io.BytesIO(obj['Body'].read()),',',names=['target'])
test_labels = test.iloc[:,0]
```

Get the predicted target and test labels.

In [80]:

```
print(target_predicted.head())

def binary_convert(x):
    threshold = 0.55
    if x > threshold:
        return 1
    else:
        return 0

target_predicted['target'] = target_predicted['target'].apply(binary_convert)

test_labels = test.iloc[:,0]

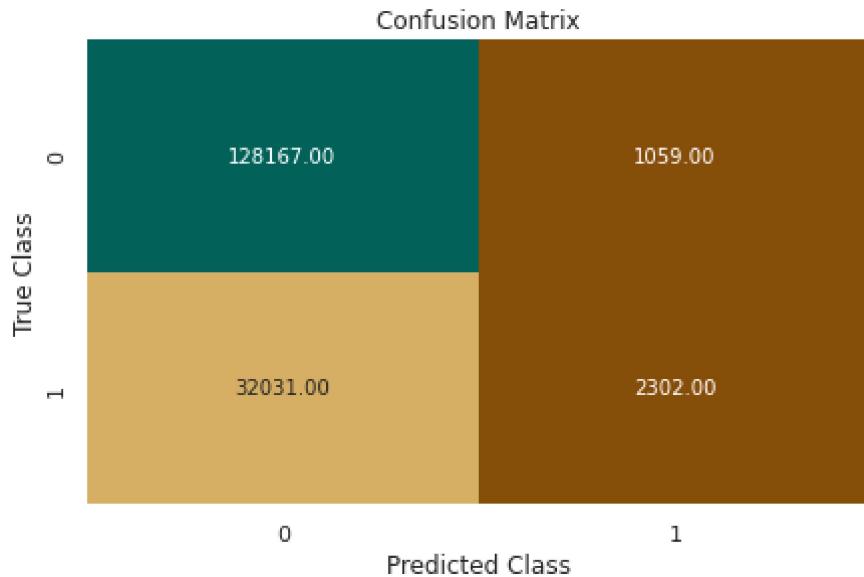
print(target_predicted.head())
```

```
target
0 0.065113
1 0.091417
2 0.476038
3 0.322747
4 0.222725
target
0      0
1      0
2      0
3      0
4      0
```

Plot a confusion matrix for your `target_predicted` and `test_labels`.

In [83]:

```
plot_confusion_matrix(test_labels, target_predicted)
```

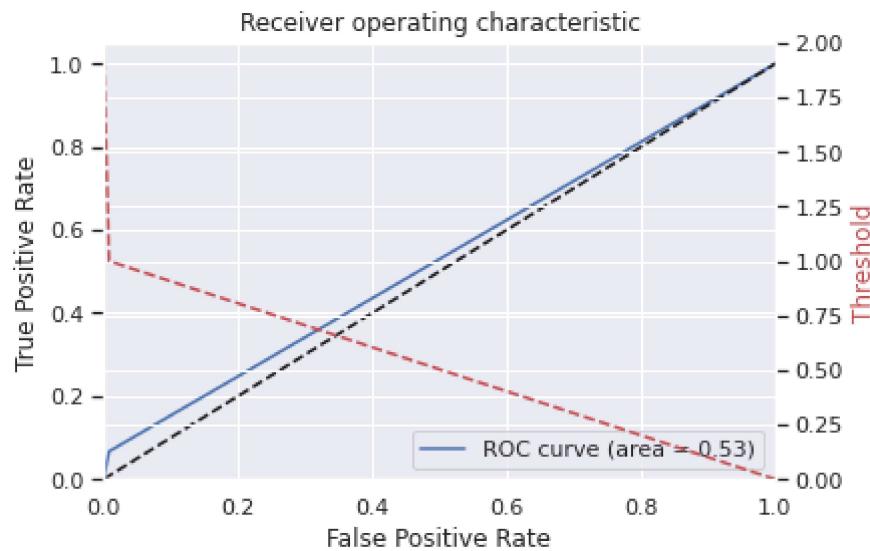


Plot the ROC chart:

In [86]:

```
plot_roc(test_labels, target_predicted)
```

Sensitivity or TPR: 6.704919465237527 %
Specificity or TNR: 99.18050547103525 %
Precision: 68.49152038083903 %
Negative Predictive Value: 80.00536835665864 %
False Positive Rate: 0.8194945289647595 %
False Negative Rate: 93.29508053476248 %
False Discovery Rate: 31.508479619160962 %
Accuracy: 79.76876845664256 %
Validation AUC 0.5294271246813638
Figure(432x288)



<Figure size 432x288 with 0 Axes>

Question: Try different hyperparameters and hyperparameter ranges. Do these changes improve the model?