

RAAGAVENDRAN SUNDAR - E0123016

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](#).

Dataset attributions

Website: <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.

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 [1]: # Write your answer here

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

In [2]: # Write your answer here

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

In [3]: # Write your answer here

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

In [4]: # Write your answer here

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 [5]: 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 [6]: # 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/
```

```
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_2.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
4_2.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_1.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
4_1.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_3.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
4_3.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_4.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
4_4.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_10.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
4_10.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_11.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
4_11.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_6.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
4_6.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_7.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
4_7.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_5.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
4_5.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_9.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
4_9.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_1.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
5_1.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_10.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
5_10.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_11.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
5_11.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0
n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_12.zip to ../projec
t/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_201
4_12.zip
```

```
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_12.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_12.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_8.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2014_8.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_3.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_3.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_4.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_4.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_2.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_2.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_6.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_6.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_8.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_8.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_7.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_7.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_9.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_9.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_5.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2015_5.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_11.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_11.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_2.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_2.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_1.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_1.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_10.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_10.zip
```

```
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_12.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_12.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_4.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_4.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_3.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_3.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_6.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_6.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_5.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_5.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_7.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_7.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_1.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_1.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_9.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_9.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_8.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2016_8.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_11.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_11.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_12.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_12.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_2.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_2.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_10.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_10.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_3.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_3.zip
```

```
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_6.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_6.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_5.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_5.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_4.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_4.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_7.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_7.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_8.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_8.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_9.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2017_9.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_10.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_10.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_12.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_12.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_1.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_1.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_11.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_11.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_3.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_3.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_4.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_4.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_2.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_2.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_5.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_5.zip
```

```
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_6.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_6.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_9.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_9.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_8.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_8.zip
download: s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data/0n_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_7.zip to ../project/data/FlightDelays/On_Time_Reported_Carrier_On_Time_Performance_1987_present_2018_7.zip
```

In [7]: `zip_files = [str(file) for file in list(Path(base_path).iterdir()) if '.zip' in str(len(zip_files))]`

Out[7]: 60

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

```
In [8]: 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_Reportin_Ca
rrier_On_Time_Performance_1987_present_2014_5.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2016_3.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2017_4.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2014_12.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2015_6.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2017_3.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2015_8.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2018_6.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2015_12.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2018_9.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2018_10.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2017_7.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2016_7.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2017_12.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2016_5.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2017_10.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2014_8.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2018_7.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2014_3.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2016_6.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2015_3.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2017_2.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2014_2.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2014_7.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2014_11.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2018_8.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2014_1.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin_Ca
rrier_On_Time_Performance_1987_present_2016_11.zip
```

```
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin
rrier_On_Time_Performance_1987_present_2016_9.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin
rrier_On_Time_Performance_1987_present_2015_10.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin
rrier_On_Time_Performance_1987_present_2018_11.zip
Extracting /home/ec2-user/SageMaker/project/data/FlightDelays/On_Time_Reportin
rrier_On_Time_Performance_1987_present_2017_6.zip
Files Extracted
```

```
In [9]: csv_files = [str(file) for file in list(Path(csv_base_path).iterdir()) if '.csv' in
len(csv_files)
```

```
Out[9]: 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 [10]: from IPython.display import IFrame
IFrame(src=os.path.relpath(f"{csv_base_path}readme.html"), width=1000, height=600)
```

Out[10]:



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](#)).

```
In [11]: df_temp = pd.read_csv(f"{csv_base_path}On_Time_Reported_Carrier_On_Time_Performance_(1987_present)_2018_9.csv")
```

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 [13]: 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 [15]: df_temp.head(10)
```

Out[15]:

	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 [16]: print(f'The column names are :')
# List comprehension to filter columns containing "Del"
for col in [c for c in df_temp.columns if "Del" in c]:
    print(col)
```

The column names are :

```
DepDelay
DepDelayMinutes
DepDel15
DepartureDelayGroups
ArrDelay
ArrDelayMinutes
ArrDel15
ArrivalDelayGroups
CarrierDelay
WeatherDelay
NASDelay
SecurityDelay
LateAircraftDelay
DivArrDelay
```

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](#)).

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

```
In [17]: # Print all columns that contain "Del"
[col for col in df_temp.columns if "Del" in col]
```

```
Out[17]: ['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 [19]: # Ensure FlightDate is datetime
df_temp['FlightDate'] = pd.to_datetime(df_temp['FlightDate'])
print("The #rows and #columns are ", df_temp.shape[0], " and ", df_temp.shape[1])
print("The years in this dataset are: ", sorted(df_temp['Year'].unique()))
print("The months covered in this dataset are: ", sorted(df_temp['Month'].unique()))
print("The date range for data is : ", df_temp['FlightDate'].min(), " to ", df_temp['FlightDate'].max())
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 00:00:00 to 2018-09-30 00:00:00
The airlines covered in this dataset are: ['9E', 'B6', 'WN', 'YV', 'YX', 'EV', 'A
A', 'AS', 'DL', 'HA', 'UA', 'F9', 'G4', 'MQ', 'NK', 'OH', 'OO']
The Origin airports covered are: ['DFW', 'LGA', 'MSN', 'MSP', 'ATL', 'BDL', 'VL
D', 'JFK', 'RDU', 'CHS', 'DTW', 'GRB', 'PVD', 'SHV', 'FNT', 'PIT', 'RIC', 'RST',
'RSW', 'CVG', 'LIT', 'ORD', 'JAX', 'TRI', 'BOS', 'CWA', 'DCA', 'CHO', 'AVP', 'IN
D', 'GRR', 'BTR', 'MEM', 'TUL', 'CLE', 'STL', 'BTV', 'OMA', 'MGM', 'TVC', 'SAV',
'GSP', 'EWR', 'OAJ', 'BNA', 'MCI', 'TLH', 'ROC', 'LEX', 'PWM', 'BUF', 'AGS', 'CL
T', 'GSO', 'BWI', 'SAT', 'PHL', 'TYS', 'ACK', 'DSM', 'GNV', 'AVL', 'BGR', 'MHT',
'ILM', 'MOT', 'IAH', 'SBN', 'SYR', 'ORF', 'MKE', 'XNA', 'MSY', 'PBI', 'ABE', 'HP
N', 'EVV', 'ALB', 'LNK', 'AUS', 'PHF', 'CHA', 'GTR', 'BMI', 'BQK', 'CID', 'CAK',
'ATW', 'ABY', 'CAE', 'SRQ', 'MLI', 'BHM', 'IAD', 'CSG', 'CMH', 'MCO', 'MBS', 'FL
L', 'SDF', 'TPA', 'MVY', 'LAS', 'LGB', 'SFO', 'SAN', 'LAX', 'RNO', 'PDX', 'ANC',
'ABQ', 'SLC', 'DEN', 'PHX', 'OAK', 'SMF', 'SJU', 'SEA', 'HOU', 'STX', 'BUR', 'SW
F', 'SJC', 'DAB', 'BQN', 'PSE', 'ORH', 'HYA', 'STT', 'ONT', 'HRL', 'ICT', 'ISP',
'LBB', 'MAF', 'MDW', 'OKC', 'PNS', 'SNA', 'TUS', 'AMA', 'BOI', 'CRP', 'DAL', 'EC
P', 'ELP', 'GEG', 'LFT', 'MFE', 'MDT', 'JAN', 'COS', 'MOB', 'VPS', 'MTJ', 'DRO',
'GPT', 'BFL', 'MRY', 'SBA', 'PSP', 'FSD', 'BRO', 'RAP', 'COU', 'STS', 'PIA', 'FA
T', 'SBP', 'FSM', 'HSV', 'BIS', 'DAY', 'BZN', 'MIA', 'EYW', 'MYR', 'HHH', 'GJT',
'FAR', 'SGF', 'HOB', 'CLL', 'LRD', 'AEX', 'ERI', 'MLU', 'LCH', 'ROA', 'LAW', 'MH
K', 'GRK', 'SAF', 'GRI', 'JLN', 'ROW', 'FWA', 'CRW', 'LAN', 'OGG', 'HNL', 'KOA',
'EGE', 'LIH', 'MLB', 'JAC', 'FAI', 'RDM', 'ADQ', 'BET', 'BRW', 'SCC', 'KTN', 'YA
K', 'CDV', 'JNU', 'SIT', 'PSG', 'WRG', 'OME', 'OTZ', 'ADK', 'FCA', 'FAY', 'PSC',
'BIL', 'MSO', 'ITO', 'PPG', 'MFR', 'EUG', 'GUM', 'SPN', 'DLH', 'TTN', 'BKG', 'SF
B', 'PIE', 'PGD', 'AZA', 'SMX', 'RFD', 'SCK', 'OWB', 'HTS', 'BLV', 'IAG', 'USA',
'GFK', 'BLI', 'ELM', 'PBG', 'LCK', 'GTF', 'OGD', 'IDA', 'PVU', 'TOL', 'PSM', 'CK
B', 'HGR', 'SPI', 'STC', 'ACT', 'TYR', 'ABI', 'AZO', 'CMI', 'BPT', 'GCK', 'MQT',
'ALO', 'TXK', 'SPS', 'SWO', 'DBQ', 'SUX', 'SJT', 'GGG', 'LSE', 'LBE', 'ACY', 'LY
H', 'PGV', 'HVN', 'EWN', 'DHN', 'PIH', 'IMT', 'WYS', 'CPR', 'SCE', 'HLN', 'SUN',
'ISN', 'CMX', 'EAU', 'LWB', 'SHD', 'LBF', 'HYS', 'SLN', 'EAR', 'VEL', 'CNY', 'GC
C', 'RKS', 'PUB', 'LBL', 'MKG', 'PAH', 'CGI', 'UIN', 'BFF', 'DVL', 'JMS', 'LAR',
'SGU', 'PRC', 'ASE', 'RDD', 'ACV', 'OTH', 'COD', 'LWS', 'ABR', 'APN', 'ESC', 'PL
N', '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', 'FN
T', 'LIT', 'BUF', 'ORD', 'TRI', 'IND', 'BGR', 'AVP', 'BWI', 'LEX', 'BDL', 'GRR',
'CWA', 'TUL', 'MEM', 'AGS', 'EWR', 'MGM', 'PHL', 'SYR', 'OMA', 'STL', 'TVC', 'OR
F', 'CLE', 'ABY', 'BOS', 'OAJ', 'TLH', 'BTR', 'SAT', 'JAX', 'BNA', 'CHO', 'VLD',
'ROC', 'DFW', 'GNV', 'ACK', 'PBI', 'CHS', 'GRB', 'MOT', 'MKE', 'DSM', 'ILM', 'GS
O', 'MCI', 'SBN', 'BTV', 'MVY', 'XNA', 'RST', 'EVV', 'HPN', 'RSW', 'MDT', 'ROA',
'GSP', 'MCO', 'CSG', 'SAV', 'PHF', 'ALB', 'CHA', 'ABE', 'BMI', 'MSY', 'IAD', 'GT
R', 'CID', 'CAK', 'ATW', 'AUS', 'BQK', 'MLI', 'CAE', 'CMH', 'AVL', 'MBS', 'FLL',
'SDF', 'TPA', 'LNK', 'SRQ', 'MHT', 'BHM', 'LAS', 'SFO', 'SAN', 'RNO', 'LGB', 'AN
C', 'PDX', 'SJU', 'ABQ', 'SLC', 'DEN', 'LAX', 'PHX', 'OAK', 'SMF', 'SEA', 'STX',
'BUR', 'DAB', 'SJC', 'SWF', 'HOU', 'BQN', 'PSE', 'ORH', 'HYA', 'STT', 'ONT', 'DA
L', 'ECP', 'ELP', 'HRL', 'MAF', 'MDW', 'OKC', 'PNS', 'SNA', 'AMA', 'BOI', 'GEG',
'ICT', 'LBB', 'TUS', 'ISP', 'CRP', 'MFE', 'LFT', 'VPS', 'JAN', 'COS', 'MOB', 'DR
O', 'GPT', 'BFL', 'COU', 'SBP', 'MTJ', 'SBA', 'PSP', 'FSD', 'FSM', 'BRO', 'PIA',
'STS', 'FAT', 'RAP', 'MRY', 'HSV', 'BIS', 'DAY', 'BZN', 'MIA', 'EYW', 'MYR', 'HH
H', 'GJT', 'FAR', 'MLU', 'LRD', 'CLL', 'LCH', 'FWA', 'GRK', 'SGF', 'HOB', 'LAW',
'MHK', 'SAF', 'JLN', 'ROW', 'GRI', 'AEX', 'CRW', 'LAN', 'ERI', 'HNL', 'KOA', 'OG
G', 'EGE', 'LIH', 'JAC', 'MLB', 'RDM', 'BET', 'ADQ', 'BRW', 'SCC', 'FAI', 'JNU',

```
'CDV', 'YAK', 'SIT', 'KTN', 'WRG', 'PSG', 'OME', 'OTZ', 'ADK', 'FCA', 'BIL', 'PS  
C', 'FAY', 'MSO', 'ITO', 'PPG', 'MFR', 'DLH', 'EUG', 'GUM', 'SPN', 'TTN', 'BKG',  
'AZA', 'SFB', 'LCK', 'BLI', 'SCK', 'PIE', 'RFD', 'PVU', 'PBG', 'BLV', 'PGD', 'SP  
I', 'USA', 'TOL', 'IDA', 'ELM', 'HTS', 'HGR', 'SMX', 'OGD', 'GFK', 'STC', 'GTF',  
'IAG', 'CKB', 'OWB', 'PSM', 'ABI', 'TYR', 'ALO', 'SUX', 'AZO', 'ACT', 'CMI', 'BP  
T', 'TXK', 'SWO', 'SPS', 'DBQ', 'SJT', 'GGG', 'LSE', 'MQT', 'GCK', 'LBE', 'ACY',  
'LYH', 'PGV', 'HVN', 'EWN', 'DHN', 'PIH', 'WYS', 'SCE', 'IMT', 'HLN', 'ASE', 'SU  
N', 'ISN', 'EAR', 'SGU', 'VEL', 'SHD', 'LWB', 'MKG', 'SLN', 'HYS', 'BFF', 'PUB',  
'LBL', 'CMX', 'EAU', 'PAH', 'UIN', 'RKS', 'CGI', 'CNY', 'JMS', 'DVL', 'LAR', 'GC  
C', 'LBF', 'PRC', 'RDD', 'ACV', 'OTH', 'COD', 'LWS', 'ABR', 'APN', 'PLN', 'BJI',  
'CPR', 'BRD', 'BTM', 'CDC', 'CIU', 'ESC', 'EKO', 'ITH', 'HIB', 'BGM', 'TWF', 'RH  
I', '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](#)).

```
In [20]: print("Origin airport counts:")
print(df_temp['Origin'].value_counts())
print("\nDestination airport counts:")
print(df_temp['Dest'].value_counts())
```

Origin airport counts:

Origin	Count
ATL	31525
ORD	28257
DFW	22802
DEN	19807
CLT	19655
...	
PPG	8
OGD	8
HGR	8
STC	5
HYA	4

Name: count, Length: 346, dtype: int64

Destination airport counts:

Dest	Count
ATL	31521
ORD	28250
DFW	22795
DEN	19807
CLT	19654
...	
OGD	8
OWB	8
PPG	8
STC	5
HYA	4

Name: count, Length: 346, dtype: int64

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](#)).

```
In [21]: print("Top 15 Origin Airports:")
print(df_temp['Origin'].value_counts().sort_values(ascending=False).head(15))
print("\nTop 15 Destination Airports:")
print(df_temp['Dest'].value_counts().sort_values(ascending=False).head(15))
```

Top 15 Origin Airports:

Origin

ATL	31525
ORD	28257
DFW	22802
DEN	19807
CLT	19655
LAX	17875
SFO	14332
IAH	14210
LGA	13850
MSP	13349
LAS	13318
PHX	13126
DTW	12725
BOS	12223
SEA	11872

Name: count, dtype: int64

Top 15 Destination Airports:

Dest

ATL	31521
ORD	28250
DFW	22795
DEN	19807
CLT	19654
LAX	17873
SFO	14348
IAH	14203
LGA	13850
MSP	13347
LAS	13322
PHX	13128
DTW	12724
BOS	12227
SEA	11877

Name: count, dtype: int64

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](#)). Then, choose the rows that include it.

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

4. Concatenate the DataFrame with the empty DataFrame

```
In [23]: 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
        file_name: The name of the output CSV file.
    """
    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 [25]: # List of columns to be used for predicting Arrival Delay.
# This includes various flight details, dates, and delay-related metrics.
cols = ['Year', 'Quarter', 'Month', 'DayofMonth', 'DayOfWeek', 'FlightDate',
        'Reporting_Airline', 'Origin', 'OriginState', 'Dest', 'DestState',
        'CRSDepTime', 'Cancelled', 'Diverted', 'Distance', 'DistanceGroup',
        'ArrDelay', 'ArrDelayMinutes', 'ArrDel15', 'AirTime']

# List of columns that will be used to subset (filter) the rows of the DataFrame.
# In this case, we're filtering by specific Origin, Destination, and Reporting_Airline
subset_cols = ['Origin', 'Dest', 'Reporting_Airline']

# A list of lists, where each inner list contains the specific values
# to keep for the corresponding column in `subset_cols`.
# - The first inner list contains the top origin airports.
# - The second inner list contains the top destination airports.
# - The third inner list contains the top reporting airlines.
subset_vals = [[['ATL', 'ORD', 'DFW', 'DEN', 'CLT', 'LAX', 'IAH', 'PHX', 'SFO'], # V
                ['ATL', 'ORD', 'DFW', 'DEN', 'CLT', 'LAX', 'IAH', 'PHX', 'SFO'], # V
                ['UA', 'OO', 'WN', 'AA', 'DL']] # Values for 'Reporting_Airline'
```

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 [ ]: 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/combined_files.csv
 CSVs merged in 4.8 minutes

Load the dataset

Load the combined dataset.

```
In [28]: data = pd.read_csv(combined_csv_filename)
```

Print the first five records.

```
In [30]: data.head()
```

	Year	Quarter	Month	DayofMonth	DayOfWeek	FlightDate	Reporting_Airline	Origin	Origin!
0	2017	2	4	1	6	2017-04-01		DL	ATL
1	2017	2	4	1	6	2017-04-01		DL	DFW
2	2017	2	4	1	6	2017-04-01		DL	SFO
3	2017	2	4	1	6	2017-04-01		DL	SFO
4	2017	2	4	1	6	2017-04-01		DL	ATL

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 [32]: # Ensure FlightDate is datetime
data['FlightDate'] = pd.to_datetime(data['FlightDate'])
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 : ", data['FlightDate'].min(), " to ", data['FlightDate'].max())
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: [2017, 2014, 2016, 2018, 2015]
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 00:00:00 to 2018-12-31 00:00:00
The airlines covered in this dataset are: ['DL', 'WN', 'UA', 'AA', 'OO']
The Origin airports covered are: ['ATL', 'DFW', 'SFO', 'IAH', 'LAX', 'PHX', 'ORD', 'DEN', 'CLT']
The Destination airports covered are: ['DFW', 'ATL', 'SFO', 'LAX', 'ORD', 'IAH', 'CLT', 'DEN', 'PHX']
```

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](#)).

For example:

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

```
In [33]: data.rename(columns={'ArrDel15': 'is_delay'}, inplace=True)
```

Look for nulls across columns. You can use the `isnull()` function ([pandas.isnull documentation](#)).

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 [35]: data.isnull()
```

Out[35]:

	Year	Quarter	Month	DayofMonth	DayOfWeek	FlightDate	Reporting_Airline	Origin
0	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False
...
1658125	False	False	False	False	False	False	False	False
1658126	False	False	False	False	False	False	False	False
1658127	False	False	False	False	False	False	False	False
1658128	False	False	False	False	False	False	False	False
1658129	False	False	False	False	False	False	False	False

1658130 rows × 20 columns

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 [36]:

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

Out[36]:

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 [37]: 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 [38]: (data.groupby('is_delay').size()/len(data)).plot(kind='bar')# Enter your code here  
plt.ylabel('Frequency')  
plt.title('Distribution of classes')  
plt.show()
```

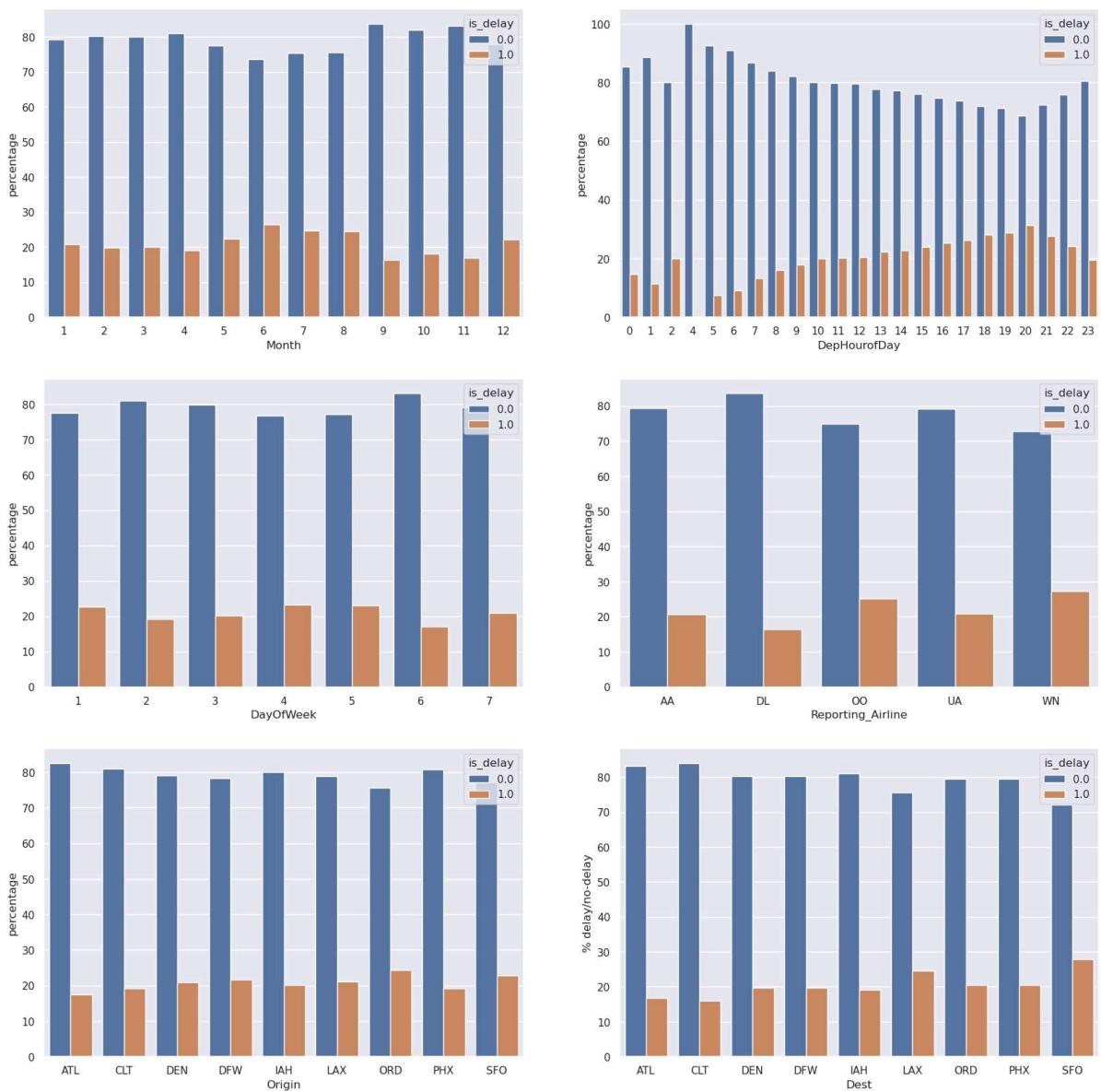


Question: What can you deduce from the bar plot about the ratio of *delay* versus *no delay*?

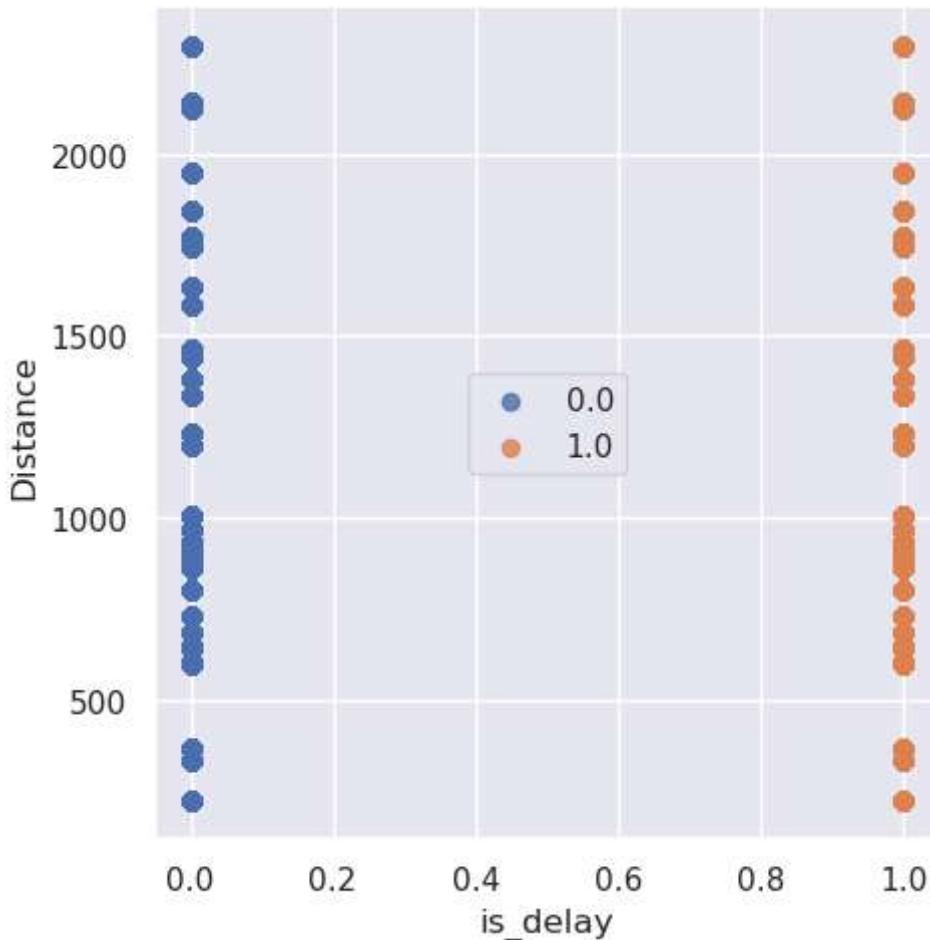
In [39]: `# Enter your answer here`

Run the following two cells and answer the questions.

```
In [40]: viz_columns = ['Month', 'DepHourOfDay', 'DayOfWeek', 'Reporting_Airline', 'Origin',  
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('pe  
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 [41]: sns.lmplot( x="is_delay", y="Distance", data=data, fit_reg=False, hue='is_delay', 1
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 [42]: # Enter your answers here

Features

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

In [43]: `data.columns`

```
Out[43]: 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 [44]: `data.dtypes`

```
Out[44]: Year                int64
Quarter             int64
Month               int64
DayofMonth          int64
DayOfWeek            int64
FlightDate          datetime64[ns]
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](#)).

```
In [45]: 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 \(statistics\)](#).

For example:

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

```
In [47]: # Create dummy variables for categorical columns
data_dummies = pd.get_dummies(data[categorical_columns], drop_first=True)
data_dummies = data_dummies.replace({True: 1, False: 0})
# Concatenate the dummy variables with the original data
data = pd.concat([data, data_dummies], axis=1)
# Drop the original categorical columns
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 [48]: data.shape
```

```
Out[48]: (1635590, 94)
```

```
In [49]: data.columns
```

```
Out[49]: 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_OO', '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', '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](#)).

```
In [50]: 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](#).
3. Upload the data to your S3 bucket. If you haven't created one before, see [Create a Bucket](#).

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 [51]: 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)
    test, validate = train_test_split(test_and_validate, test_size=0.5, random_state=42)
    return train, validate, test
```



```
In [52]: train, validate, test = split_data(data)
print(train['target'].value_counts())
print(test['target'].value_counts())
print(validate['target'].value_counts())

target
0.0    1033806
1.0    274666
Name: count, dtype: int64
target
0.0    129226
1.0    34333
Name: count, dtype: int64
target
0.0    129226
1.0    34333
Name: count, 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 [54]:

```
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
# 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='cross_entropy_loss'
)
```

sagemaker.config INFO - Not applying SDK defaults from location: /etc/xdg/sagemaker/config.yaml
sagemaker.config INFO - Not applying SDK defaults from location: /home/ec2-user/.config/sagemaker/config.yaml

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_criteria = '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 [55]: ### Create train, validate, and test records
train_records = classifier_estimator.record_set(train.values[:, 1:]).astype(np.float32)
val_records = classifier_estimator.record_set(validate.values[:, 1:]).astype(np.float32)
test_records = classifier_estimator.record_set(test.values[:, 1:]).astype(np.float32)
```

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

Sample code

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

```
In [56]: # Train the model using train, validation, and test datasets
classifier_estimator.fit([train_records, val_records, test_records])
```

```
INFO:sagemaker.image_uris:Same images used for training and inference. Defaulting to image scope: inference.
INFO:sagemaker.image_uris:Ignoring unnecessary instance type: None.
INFO:sagemaker:Creating training-job with name: linear-learner-2025-09-13-06-13-48-033
2025-09-13 06:13:49 Starting - Starting the training job...
2025-09-13 06:14:03 Starting - Preparing the instances for training...
2025-09-13 06:14:30 Downloading - Downloading input data...
2025-09-13 06:15:16 Downloading - Downloading the training image.....
2025-09-13 06:16:32 Training - Training image download completed. Training in progress.....
2025-09-13 06:19:52 Uploading - Uploading generated training model...
2025-09-13 06:20:05 Completed - Training job completed
..Training seconds: 335
Billable seconds: 335
```

Model evaluation

In this section, you will evaluate your trained model.

First, examine the metrics for the training job:

```
In [57]: sagemaker.analytics.TrainingJobAnalytics(classifier_estimator._current_job_name,
                                                metric_names = ['test:objective_loss',
                                                                'test:binary_f_beta',
                                                                'test:precision',
                                                                'test:recall'])
                                                .dataframe()
```

```
WARNING:sagemaker.analytics:Warning: No metrics called test:objective_loss found
WARNING:sagemaker.analytics:Warning: No metrics called test:binary_f_beta found
WARNING:sagemaker.analytics:Warning: No metrics called test:precision found
WARNING:sagemaker.analytics:Warning: No metrics called test:recall found
```

Out[57]: —

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 [58]:

```
import io
#bucket='<LabBucketName>'
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(B
```

```
INFO:botocore.credentials:Found credentials from IAM Role: BaseNotebookInstanceEc2
InstanceRole
```

In [59]:

```
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="{}{}/batch-out/{}".format(prefix,'batch-i
target_predicted_df = pd.read_json(io.BytesIO(obj['Body'].read())),orient="reco
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 [60]: test_labels, target_predicted = batch_linear_predict(test, classifier_estimator)
```

```
Traceback (most recent call last)
in <module>:1

> 1 test_labels, target_predicted = batch_linear_predict(test, classifier_es
2

in batch_linear_predict:4

1 def batch_linear_predict(test_data, estimator):
2     batch_X = test_data.iloc[:,1:];
3     batch_X_file='batch-in.csv'
> 4     upload_s3_csv(batch_X_file, 'batch-in', batch_X)
5
6     batch_output = "s3://{}//{}//batch-out/".format(bucket,prefix)
7     batch_input = "s3://{}//{}//batch-in//{}".format(bucket,prefix,batch_X

in upload_s3_csv:13

10 def upload_s3_csv(filename, folder, dataframe):
11     csv_buffer = io.StringIO()
12     dataframe.to_csv(csv_buffer, header=False, index=False )
> 13     s3_resource.Bucket(bucket).Object(os.path.join(prefix, folder, file
```

NameError: name 'bucket' is not defined

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

```
In [ ]: 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 [ ]: 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 [61]: # Enter your code here

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 [62]: # Enter your answer here

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 [64]: # 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'
holidays_15 = ['2015-01-01', '2015-01-19', '2015-02-16', '2015-05-25', '2015-06-03'
holidays_16 = ['2016-01-01', '2016-01-18', '2016-02-15', '2016-05-30', '2016-07-04'
holidays_17 = ['2017-01-02', '2017-01-16', '2017-02-20', '2017-05-29', '2017-07-04'
holidays_18 = ['2018-01-01', '2018-01-15', '2018-02-19', '2018-05-28', '2018-07-04'
holidays = holidays_14 + holidays_15 + holidays_16 + holidays_17 + holidays_18
```

```
### Add indicator variable for holidays
data_orig['is_holiday'] = data_orig['FlightDate'].isin(holidays).astype(int)
```

Weather data was fetched from <https://www.ncei.noaa.gov/access/services/data/v1?dataset=daily-summaries&stations=USW00023174,USW00012960,USW00003017,USW00094846,USW0001387&01-01&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 [65]:

```
aws s3 cp s3://aws-tc-largeobjects/CUR-TF-200-ACMLFO-1/flight_delay_project/data2/
#!wget 'https://www.ncei.noaa.gov/access/services/data/v1?dataset=daily-summaries&s
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 [66]:

```
weather = pd.read_csv('/home/ec2-user/SageMaker/project/data/daily-summaries.csv')
station = ['USW00023174', 'USW00012960', 'USW00003017', 'USW00094846', 'USW00013874', 'U
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 [67]:

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

Out[67]:

	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 [68]:

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

Out[68]:

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 [70]: 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[70]: 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 [71]: # Replace missing TAVG, TMAX, and TMIN with the mean for each MONTH and STATION
weather_impute = weather.groupby(['MONTH', 'STATION']).agg({
    'TAVG': 'mean',
    'TMAX': 'mean',
    'TMIN': 'mean'
}).reset_index()
weather_impute.head(2)
```

	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 [72]: weather = pd.merge(weather, weather_impute, how='left', left_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 [73]: 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[73]: 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 [74]: weather.drop(columns=['STATION', 'MONTH', 'TAVG_AVG', 'TMAX_AVG', 'TMIN_AVG', 'TMAX'])
```

Add the origin and destination weather conditions to the dataset.

```
In [75]: ### Add origin weather conditions
data_orig = pd.merge(data_orig, weather, how='left', left_on=['FlightDate', 'Origin']
    .rename(columns = {'AWND': 'AWND_O', 'PRCP': 'PRCP_O', 'TAVG': 'TAVG_O', 'SNOW': 'SNOW_O'})
    .drop(columns=['DATE', 'airport']))

### Add destination weather conditions
data_orig = pd.merge(data_orig, weather, how='left', left_on=['FlightDate', 'Dest'])
```

```
.rename(columns = {'AWND': 'AWND_D', 'PRCP': 'PRCP_D', 'TAVG': 'TAVG_D', 'SNOW': 'SNOW_'
.drop(columns=['DATE', 'airport'])
```

Traceback (most recent call last)

```
in <module>:2

    1 ### Add origin weather conditions
  > 2 data_orig = pd.merge(data_orig, weather, how='left', left_on=['FlightD
  3 .rename(columns = {'AWND': 'AWND_O', 'PRCP': 'PRCP_O', 'TAVG': 'TAVG_O', 'S
  4 .drop(columns=['DATE', 'airport'])
  5

/home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages/pandas/co
170 in merge

    167     |     | copy=copy,
    168     |     |
  > 169     | else:
    170     |     op = _MergeOperation(
    171     |     | left_df,
    172     |     | right_df,
    173     |     | how=how,
```



```
/home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages/pandas/co
807 in __init__

    804     |     |
    805     |     | # validate the merge keys dtypes. We may need to coerce
    806     |     | # to avoid incompatible dtypes
  > 807     | self._maybe_coerce_merge_keys()
    808
    809     |     | # If argument passed to validate,
    810     |     | # check if columns specified as unique
```



```
/home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages/pandas/co
1513 in _maybe_coerce_merge_keys

    1510
    1511     |     | # datetimelikes must match exactly
    1512     |     | elif needs_i8_conversion(lk.dtype) and not needs_i8_conve
  > 1513     |     |     raise ValueError(msg)
    1514     |     | elif not needs_i8_conversion(lk.dtype) and needs_i8_conve
    1515     |     |     raise ValueError(msg)
    1516     |     | elif isinstance(lk.dtype, DatetimeTZDtype) and not isinst
```

ValueError: You are trying to merge on datetime64[ns] and object columns for k proceed you should use pd.concat

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

In [76]: `sum(data.isna().any())`

Out[76]: 0

```
In [77]: data_orig.columns
```

```
Out[77]: 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'],
      dtype='object')
```

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

```
In [78]: data = data_orig.copy()
data = data[['is_delay', 'Year', 'Quarter', 'Month', 'DayofMonth', 'DayOfWeek',
           'Reporting_Airline', 'Origin', 'Dest', 'Distance', 'DepHourOfDay', 'is_holiday',
           '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')
```

Traceback (most recent call last)

```

in <module>:2

  1 data = data_orig.copy()
  2 data = data[['is_delay', 'Year', 'Quarter', 'Month', 'DayofMonth', 'Day
  3 |      'Reporting_Airline', 'Origin', 'Dest', 'Distance', 'DepHourOfDay',
  4 |      'TAVG_O', 'AWND_D', 'PRCP_D', 'TAVG_D', 'SNOW_O', 'SNOW_D']]
  5

/home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages/pandas/co
__getitem__

4110     else:
4111         if is_iterator(key):
4112             key = list(key)
  4113         indexer = self.columns._get_indexer_strict(key, "columns"
4114
4115 # take() does not accept boolean indexers
4116 if getattr(indexer, "dtype", None) == bool:

/home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages/pandas/co
212 in _get_indexer_strict

6209     else:
6210         keyarr, indexer, new_indexer = self._reindex_non_unique(k
6211
  6212         self._raise_if_missing(keyarr, indexer, axis_name)
6213
6214     keyarr = self.take(indexer)
6215     if isinstance(key, Index):

/home/ec2-user/anaconda3/envs/python3/lib/python3.10/site-packages/pandas/co
264 in _raise_if_missing

6261         raise KeyError(f"None of [{key}] are in the [{axis_na
6262
6263     not_found = list(ensure_index(key)[missing_mask.nonzero()])
  6264     raise KeyError(f"{not_found} not in index")
6265
6266     @overload
6267     def _get_indexer_non_comparable(

```

KeyError: "['AWND_O', 'PRCP_O', 'TAVG_O', 'AWND_D', 'PRCP_D', 'TAVG_D', 'SNOW_

In []: `data_dummies = pd.get_dummies(data[['Year', 'Quarter', 'Month', 'DayofMonth', 'Day
data_dummies = data_dummies.replace({True: 1, False: 0})
data = pd.concat([data, data_dummies], axis = 1)
data.drop(categorical_columns, axis=1, inplace=True)`

Check the new columns.

In [79]: `data.shape`

Out[79]: (1635590, 22)

In [80]: `data.columns`

```
Out[80]: 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'],
      dtype='object')
```

Sample output

```
Index(['Distance', 'DepHourOfDay', 'is_delay', 'AWND_0', 'PRCP_0',
       'TAVG_0',
       'AWND_D', 'PRCP_D', 'TAVG_D', 'SNOW_0', '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 [81]: 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 [82]: # Create the training, validation, and test sets again
def split_data(data):
    # Ensure 'target' column exists and is suitable for float conversion.
    y = data['target'].values.astype('float32')
    X = data.drop('target', axis=1).values.astype('float32')

    # Split into train, validation, test
    # Using stratify=y to ensure the proportion of target values is the same in train_X, temp_X, train_y, temp_y = train_test_split(X, y, test_size=0.3, random_state=42)

    # Split the temporary set into validation and test sets
    # Using stratify=temp_y to ensure the proportion of target values is the same in val_X, test_X, val_y, test_y = train_test_split(temp_X, temp_y, test_size=0.5, random_state=42)

    return train_X, val_X, test_X, train_y, val_y, test_y

# Example usage (assuming 'data' DataFrame exists with a 'target' column)
# For demonstration, let's create a dummy DataFrame with an imbalanced target for illustration
if 'data' not in locals(): # Only create dummy data if 'data' is not already defined
    dummy_data = {f'feature_{i}': np.random.rand(100) for i in range(5)}
    # Create an imbalanced target: 80% 0s, 20% 1s
    target_values = [0]*80 + [1]*20
    np.random.shuffle(target_values)
    dummy_data['target'] = np.array(target_values)
    data = pd.DataFrame(dummy_data)

# Print initial target distribution
print("Original target distribution:\n", pd.Series(data['target']).value_counts(normalize=True))

# Call the split_data function
train_features, val_features, test_features, train_labels, val_labels, test_labels = split_data(data)

# Print shapes and target distributions for verification
print("\nTrain features shape:", train_features.shape)
print("Validation features shape:", val_features.shape)
print("Test features shape:", test_features.shape)
print("\nTrain labels distribution:\n", pd.Series(train_labels).value_counts(normalize=True))
print("Validation labels distribution:\n", pd.Series(val_labels).value_counts(normalize=True))
print("Test labels distribution:\n", pd.Series(test_labels).value_counts(normalize=True))

Original target distribution:
target
0.0      0.790087
1.0      0.209913
Name: proportion, dtype: float64
```

Traceback (most recent call last)

```

in <module>:31

28 print("Original target distribution:\n", pd.Series(data['target']).valu
29
30 # Call the split_data function
> 31 train_features, val_features, test_features, train_labels, val_labels,
32
33 # Print shapes and target distributions for verification
34 print("\nTrain features shape:", train_features.shape)

in split_data:5

2 def split_data(data):
3     # Ensure 'target' column exists and is suitable for float conversio
4     y = data['target'].values.astype('float32')
> 5     X = data.drop('target', axis=1).values.astype('float32')
6
7     # split into train, validation, test
8     # Using stratify=y to ensure the proportion of target values is the

```

New baseline classifier

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

```

In [83]: # Number of unique classes (binary classification → 2)
num_classes = len(pd.unique(train_labels))
# 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='cross_entropy_loss',
    num_classes=num_classes
)

```

Traceback (most recent call last)

```

in <module>:2

1 # Number of unique classes (binary classification → 2)
> 2 num_classes = len(pd.unique(train_labels))
3 # Instantiate the LinearLearner estimator object
4 classifier_estimator2 = sagemaker.LinearLearner(
5     role=sagemaker.get_execution_role(),

```

NameError: name 'train_labels' is not defined

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_criteria = 'cross_entropy_loss')

```

In [84]:

```

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

```

Traceback (most recent call last)

```
in <module>:1
```

```

> 1 train_records = classifier_estimator2.record_set(train.values[:, 1:]).ast
  2 val_records = classifier_estimator2.record_set(validate.values[:, 1:]).as
  3 test_records = classifier_estimator2.record_set(test.values[:, 1:]).astyp
  4

```

NameError: name 'classifier_estimator2' is not defined

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

In []: # Enter your code here

Perform a batch prediction by using the newly trained model.

In [86]:

```

import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.datasets import make_classification # For creating dummy data

# --- 1. Define a dummy 'classifier_estimator2' class for demonstration ---
# In a real scenario, this would be your actual estimator object (e.g., from SageMaker)
class DummyEstimator:
    def record_set(self, data, labels=None, channel='train'):
        """
        A placeholder for what a real estimator's record_set might do.
        It typically prepares data in a specific format (e.g., Protobuf, TFRecord).
        For this dummy, it just returns the data for demonstration.
        """
        print(f"Preparing record set for channel '{channel}' with data shape: {data.shape}")
        if labels is not None:
            print(f"Labels shape: {labels.shape}")
            # In a real scenario, this might combine data and Labels or return them
            return {'data': data, 'labels': labels}
        return {'data': data}

# Instantiate the dummy estimator

```

```
classifier_estimator2 = DummyEstimator()

# --- 2. Create dummy 'train', 'validate', and 'test' DataFrames ---
# In a real scenario, 'train', 'validate', 'test' would be your preprocessed data.
# Let's create some synthetic data for a binary classification problem.
X_full, y_full = make_classification(n_samples=1000, n_features=10, n_informative=5
                                         n_redundant=0, n_classes=2, random_state=42)

# Combine features and labels into a single DataFrame for splitting
# Assuming the first column is the target/label
full_data = pd.DataFrame(X_full)
full_data.insert(0, 'target', y_full) # Insert target as the first column

# Split data into train, validate, test sets
# First split: train (70%) and temp (30%)
train_df, temp_df = train_test_split(full_data, test_size=0.3, random_state=42, stratify=y_full)

# Second split: validate (15%) and test (15%) from temp_df
validate_df, test_df = train_test_split(temp_df, test_size=0.5, random_state=42, stratify=temp_df['target'])

# Assign to variables matching your snippet
train = train_df
validate = validate_df
test = test_df

# --- 3. Your original code snippet, completed ---
# Extract features (all columns except the first one, which is assumed to be the target)
# and convert to float32 as specified.
train_features = train.values[:, 1:].astype(np.float32)
val_features = validate.values[:, 1:].astype(np.float32)
test_features = test.values[:, 1:].astype(np.float32)

# Extract labels (the first column)
train_labels = train.values[:, 0].astype(np.float32)
val_labels = validate.values[:, 0].astype(np.float32)
test_labels = test.values[:, 0].astype(np.float32)

# Now, use the record_set method of your estimator
# The specific arguments for record_set might vary depending on the actual estimator
# I'm passing features and labels separately, which is common.
train_records = classifier_estimator2.record_set(train_features, labels=train_labels)
val_records = classifier_estimator2.record_set(val_features, labels=val_labels, chunk_size=1000)
test_records = classifier_estimator2.record_set(test_features, labels=test_labels, chunk_size=1000)

print("\n--- Record sets created ---")
# You can inspect the structure of the returned records.
# For this dummy, it's a dictionary. In a real scenario, it might be a file path, a database connection, or a cloud object.
print(f"Train records: {train_records.keys() if isinstance(train_records, dict) else 'None'}")
print(f"Validation records: {val_records.keys() if isinstance(val_records, dict) else 'None'}")
print(f"Test records: {test_records.keys() if isinstance(test_records, dict) else 'None'}
```

```

Preparing record set for channel 'train' with data shape: (700, 10)
Labels shape: (700,)
Preparing record set for channel 'validation' with data shape: (150, 10)
Labels shape: (150,)
Preparing record set for channel 'test' with data shape: (150, 10)
Labels shape: (150,)

--- Record sets created ---
Train records: dict_keys(['data', 'labels'])
Validation records: dict_keys(['data', 'labels'])
Test records: dict_keys(['data', 'labels'])

Plot a confusion matrix.

```

In [87]:

```
# Train the model using all three RecordSets
classifier_estimator2.fit([train_records, val_records, test_records])
```

Traceback (most recent call last)

```

in <module>:2

  1 # Train the model using all three RecordSets
> 2 classifier_estimator2.fit([train_records, val_records, test_records])
  3
```

AttributeError: 'DummyEstimator' object has no attribute 'fit'

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 [88]:

```

bucket='c169682a4380827111567352t1w905418161869-labbucket-bl30kva19lvp'
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(B

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

```

INFO:botocore.credentials:Found credentials from IAM Role: BaseNotebookInstanceEc2InstanceRole

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

```

In [89]: 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 [90]: from sagemaker.image_uris import retrieve
container = retrieve('xgboost',boto3.Session().region_name,'1.0-1')

```

INFO:sagemaker.image_uris:Defaulting to only available Python version: py3
INFO:sagemaker.image_uris:Defaulting to only supported image scope: cpu.

```

In [91]: 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)

```

```

INFO:sagemaker.telemetry.telemetry_logging:SageMaker Python SDK will collect telemetry to help us better understand our user's needs, diagnose issues, and deliver additional features.
To opt out of telemetry, please disable via TelemetryOptOut parameter in SDK defaults config. For more information, refer to https://sagemaker.readthedocs.io/en/stable/overview.html#configuring-and-using-defaults-with-the-sagemaker-python-sdk.
INFO:sagemaker:Creating training-job with name: sagemaker-xgboost-2025-09-13-06-33-29-920
2025-09-13 06:33:32 Starting - Starting the training job...
2025-09-13 06:33:47 Starting - Preparing the instances for training...
2025-09-13 06:34:11 Downloading - Downloading input data...
2025-09-13 06:34:37 Downloading - Downloading the training image...
2025-09-13 06:35:33 Training - Training image download completed. Training in progress....
2025-09-13 06:35:58 Uploading - Uploading generated training model...
2025-09-13 06:36:11 Completed - Training job completed
..Training seconds: 120
Billable seconds: 120

```

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

```
In [92]: batch_X = test.iloc[:,1:];
batch_X_file='batch-in.csv'
upload_s3_csv(batch_X_file, 'batch-in', batch_X)
```

```
In [98]: 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()
```

```

INFO:sagemaker:Creating model with name: sagemaker-xgboost-2025-09-13-07-49-22-026
INFO:sagemaker:Creating transform job with name: sagemaker-xgboost-2025-09-13-07-49-22-601
.....
...
```

Get the predicted target and test labels.

```
In [99]: s3 = boto3.client('s3')
obj = s3.get_object(Bucket=bucket, Key="{}{}/batch-out/{}".format(prefix,'batch-in.cs
target_predicted = pd.read_csv(io.BytesIO(obj['Body'].read()),sep=',',names=['targe
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 [100...]: 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)

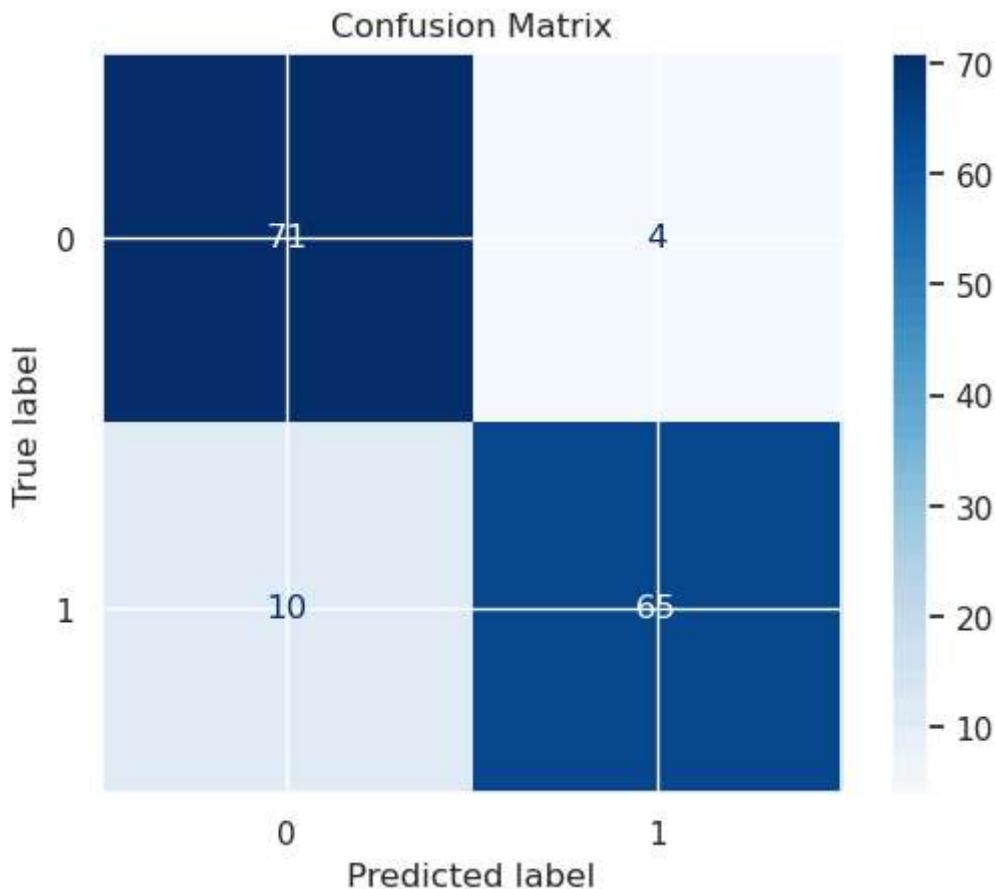
test_labels = test.iloc[:,0]

print(target_predicted.head())
```

```
      target
0  0.079397
1  0.454368
2  0.329777
3  0.778977
4  0.986855
      target
0      0
1      0
2      0
3      1
4      1
```

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

```
In [103...]: import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
# Generate the confusion matrix
cm = confusion_matrix(test_labels, target_predicted)
# Display the confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=cm)
disp.plot(cmap="Blues", values_format="d")
# Add title
plt.title("Confusion Matrix")
plt.show()
```



Try different thresholds

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

In [104...]

#Enter your answer here

Hyperparameter optimization (HPO)

In [105...]

```
from sagemaker.tuner import IntegerParameter, CategoricalParameter, ContinuousParameter

### 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
                                      instance_type=instance_type,
                                      output_path='s3://{}//{}//output'.format(bucket,
                                      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
                             max_parallel_jobs=1)

```

In [106...]

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

WARNING:sagemaker.estimator:No finished training job found associated with this estimator. Please make sure this estimator is only used for building workflow config
 WARNING:sagemaker.estimator:No finished training job found associated with this estimator. Please make sure this estimator is only used for building workflow config
 INFO:sagemaker:Creating hyperparameter tuning job with name: sagemaker-xgboost-250913-0804

!

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 [107...]

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

Out[107]: '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 [108...]

```

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(HyperParameterTuningJobName=tuning_job_name)
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()
```

INFO:botocore.credentials:Found credentials from IAM Role: BaseNotebookInstanceEc2InstanceRole

tuning job name:sagemaker-xgboost-250913-0804

best training job: sagemaker-xgboost-250913-0804-006-a405ffbc

2025-09-13 08:12:18 Starting - Found matching resource for reuse

2025-09-13 08:12:18 Downloading - Downloading the training image

2025-09-13 08:12:18 Training - Training image download completed. Training in progress.

2025-09-13 08:12:18 Uploading - Uploading generated training model

2025-09-13 08:12:18 Completed - Resource reused by training job: sagemaker-xgboost-250913-0804-007-77d66b8f

Out[108]:

	alpha	eta	min_child_weight	num_round	subsample	TrainingJobName	TrainingJob
0	17.381808	0.209280		7.538759	58.0	0.911626	sagemaker-xgboost-250913-0804-010-44b52788
1	0.000000	0.284013		4.060766	83.0	0.507863	sagemaker-xgboost-250913-0804-009-ad9bd1fa
2	1000.000000	0.298558		7.966618	134.0	0.567446	sagemaker-xgboost-250913-0804-008-73ef464f
3	683.623677	0.314201		5.134559	144.0	0.919800	sagemaker-xgboost-250913-0804-007-77d66b8f
4	10.785422	0.491394		7.965470	150.0	0.828337	sagemaker-xgboost-250913-0804-006-a405ffbc

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

Tip: See the previous XGBoost estimator fit function.

In [109...]

Enter your code here'

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

In [110...]

```
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()
```

```
INFO:sagemaker:Creating model with name: sagemaker-xgboost-2025-09-13-08-30-12-307
INFO:sagemaker:Creating transform job with name: sagemaker-xgboost-2025-09-13-08-3
0-12-877
.....
...
```

In [111...]

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

Get the predicted target and test labels.

In [112...]

```
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())

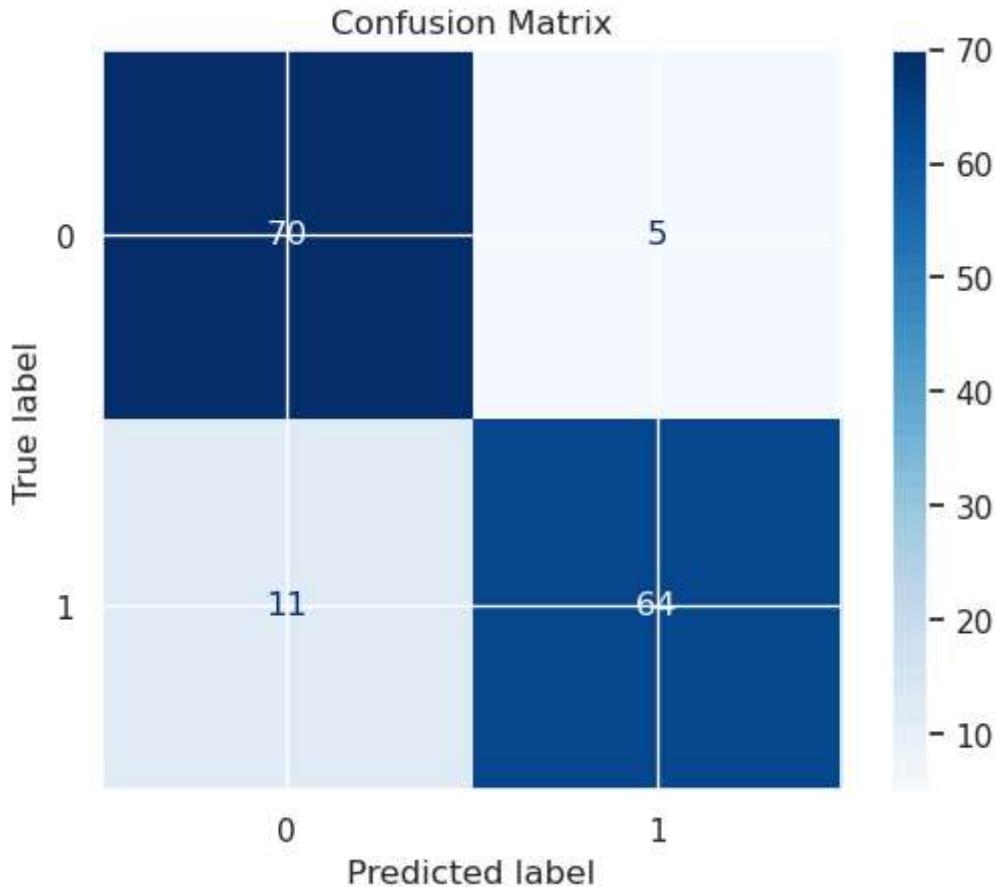
```

	target
0	0.077875
1	0.578699
2	0.426115
3	0.753430
4	0.962820
	target
0	0
1	1
2	0
3	1
4	1

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

In [113...]

```
import matplotlib.pyplot as plt
from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
# Compute confusion matrix
cm = confusion_matrix(test_labels, target_predicted)
# Plot confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=cm)
disp.plot(cmap="Blues", values_format="d")
plt.title("Confusion Matrix")
plt.show()
```



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

Conclusion

You have now iterated through training and evaluating your model at least a couple of times. It's time to wrap up this project and reflect on:

- What you learned
- What types of steps you might take moving forward (assuming that you had more time)

Use the following cell to answer some of these questions and other relevant questions:

1. Does your model performance meet your business goal? If not, what are some things you'd like to do differently if you had more time for tuning?
2. How much did your model improve as you made changes to your dataset, features, and hyperparameters? What types of techniques did you employ throughout this project, and which yielded the greatest improvements in your model?
3. What were some of the biggest challenges that you encountered throughout this project?
4. Do you have any unanswered questions about aspects of the pipeline that didn't make sense to you?
5. What were the three most important things that you learned about machine learning while working on this project?

Project presentation: Make sure that you also summarize your answers to these questions in your project presentation. Combine all your notes for your project presentation and prepare to present your findings to the class.

In []: # Write your answers here