

(https://databricks.com)

```
from pyspark.sql.functions import isnan, when, count, col, split, trim, lit, avg, sum
from pyspark.sql.functions import col, rand
from pyspark.sql.functions import col, concat, lit, to_timestamp, format_string
from pyspark.sql.functions import concat, col, lpad, lit, to_date
from pyspark.sql.functions import *
import pyspark.sql.functions as F
from pyspark.sql.functions import isnan, when, count, col, split, trim, lit, avg, sum, length, regexp_replace
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from pyspark.sql.types import DoubleType, FloatType
from pyspark.ml.stat import Correlation
from pyspark.ml.feature import VectorAssembler
```

```
print("Welcome to the W261 final project!")
```

Welcome to the W261 final project!

Set Checkpoint directory

```
spark.sparkContext.setCheckpointDir("/ml_41/df_otpw_checkpoint")
```

Locate Raw Data

```
data_BASE_DIR = "dbfs:/mnt/mids-w261/"
display(dbutils.fs.ls(f"{data_BASE_DIR}"))
```

Table					
	path	name	size	modificationTime	
1	dbfs:/mnt/mids-w261/HW5/	HW5/	0	0	
2	dbfs:/mnt/mids-w261/OTPW_12M/	OTPW_12M/	0	0	
3	dbfs:/mnt/mids-w261/OTPW_1D_CSV/	OTPW_1D_CSV/	0	0	
4	dbfs:/mnt/mids-w261/OTPW_36M/	OTPW_36M/	0	0	
5	dbfs:/mnt/mids-w261/OTPW_3M/	OTPW_3M/	0	0	
6	dbfs:/mnt/mids-w261/OTPW_3M_2015.csv	OTPW_3M_2015.csv	1500620247	1679772070000	
10 rows					

Load Raw Data

```
# OTPW
df_otpw = spark.read.format("csv").option("header", "true").load(f"dbfs:/mnt/mids-w261/OTPW_60M/")
```

```
# Select only columns needed
df_otpw = df_otpw.select('DEP_DEL15', 'DEP_DELAY_NEW', 'ORIGIN_AIRPORT_ID', 'DEST_AIRPORT_ID', 'FL_DATE',
'CRS_DEP_TIME', 'DAY_OF_MONTH', 'DAY_OF_WEEK', 'OP_UNIQUE_CARRIER', 'DEP_TIME_BLK', 'TAIL_NUM', 'MONTH', 'YEAR',
'HourlyDewPointTemperature', 'HourlyDryBulbTemperature', 'HourlyWetBulbTemperature', 'HourlyRelativeHumidity',
'HourlyWindDirection', 'HourlyWindSpeed')

df_otpw.printSchema()
```

```
root
|-- DEP_DEL15: string (nullable = true)
|-- DEP_DELAY_NEW: string (nullable = true)
|-- ORIGIN_AIRPORT_ID: string (nullable = true)
|-- DEST_AIRPORT_ID: string (nullable = true)
|-- FL_DATE: string (nullable = true)
|-- CRS_DEP_TIME: string (nullable = true)
|-- DAY_OF_MONTH: string (nullable = true)
|-- DAY_OF_WEEK: string (nullable = true)
|-- OP_UNIQUE_CARRIER: string (nullable = true)
|-- DEP_TIME_BLK: string (nullable = true)
|-- TAIL_NUM: string (nullable = true)
|-- MONTH: string (nullable = true)
|-- YEAR: string (nullable = true)
|-- HourlyDewPointTemperature: string (nullable = true)
|-- HourlyDryBulbTemperature: string (nullable = true)
|-- HourlyWetBulbTemperature: string (nullable = true)
|-- HourlyRelativeHumidity: string (nullable = true)
|-- HourlyWindDirection: string (nullable = true)
|-- HourlyWindSpeed: string (nullable = true)
```

Drop all observations with null for target variable

```
df_otpw = df_otpw.dropna(subset=['DEP_DEL15'])
```

Downsample

```
# Show distribution of year on raw data
df_otpw.groupBy(['YEAR']).count().show()
```

```
+----+-----+
|YEAR|  count|
+----+-----+
|2016|5546861|
|2017|5572592|
|2015|5726181|
|2019|7263966|
|2018|7087730|
+----+-----+
```

```
# Count occurrences of 0 and 1 in DEP_DEL15
count_0 = df_otpw.filter(col("DEP_DEL15") == 0).count()
count_1 = df_otpw.filter(col("DEP_DEL15") == 1).count()
total_count = df_otpw.count()

# Randomly sample and filter to balance counts
df_0 = df_otpw.filter(col("DEP_DEL15") == 0).sample(False, count_1/count_0, seed=42)
df_1 = df_otpw.filter(col("DEP_DEL15") == 1)

# Now both DataFrames have the same number of rows for DEP_DEL15 being 0 and 1
df_otpw = df_0.union(df_1)

# Check the counts after balancing
df_otpw.groupBy("DEP_DEL15").count().show()
```

```
+-----+-----+
|DEP_DEL15| count|
+-----+-----+
|      0.0|5680557|
|      1.0|5676029|
+-----+-----+
```

```
df_otpw.groupBy(['YEAR']).count().show()
```

```
+---+-----+
|YEAR| count|
+---+-----+
|2019|2671488|
|2018|2591451|
|2016|1973748|
|2017|2024605|
|2015|2095294|
+---+-----+
```

Convert Data to Delta Lake Format

```
# Configure Path
DELTALAKE_GOLD_PATH = "/ml_41/flights_5Y.delta"

# Remove table if it exists
dbutils.fs.rm(DELTALAKE_GOLD_PATH, recurse=True)

# Save table as Delta Lake
df_otpw.write.format("delta").mode("overwrite").save(DELTALAKE_GOLD_PATH)
```

Load Checkpoint Data

```
# Configure Path
DELTALAKE_GOLD_PATH = "/ml_41/flights_5Y.delta"

# Re-read as Delta Lake
df_otpw = spark.read.format("delta").load(DELTALAKE_GOLD_PATH)

# Sample if needed
# df_otpw = df_otpw.sample(False, 1000/df_otpw.count(), seed = 42)

# Review data
#display(df_otpw)
# df_otpw.show(2)
# df_otpw.count()
```

```
# Select only columns needed
df_otpw = df_otpw.select('DEP_DEL15', 'DEP_DELAY_NEW', 'ORIGIN_AIRPORT_ID', 'DEST_AIRPORT_ID', 'FL_DATE',
'CRS_DEP_TIME', 'DAY_OF_MONTH', 'DAY_OF_WEEK', 'OP_UNIQUE_CARRIER', 'DEP_TIME_BLK', 'TAIL_NUM', 'MONTH', 'YEAR',
'HourlyDewPointTemperature', 'HourlyDryBulbTemperature', 'HourlyWetBulbTemperature', 'HourlyRelativeHumidity',
'HourlyWindDirection', 'HourlyWindSpeed')

df_otpw.printSchema()
```

Cancelled

Drop all observations with null for target variable

```
df_otpw = df_otpw.dropna(subset=['DEP_DEL15'])
```

Cancelled

```
df_otpw.groupBy(['YEAR']).count().show()
```

Cancelled

```
# Count occurrences of 0 and 1 in DEP_DEL15
count_0 = df_otpw.filter(col("DEP_DEL15") == 0).count()
count_1 = df_otpw.filter(col("DEP_DEL15") == 1).count()
total_count = df_otpw.count()

# Randomly sample and filter to balance counts
df_0 = df_otpw.filter(col("DEP_DEL15") == 0).sample(False, count_1/count_0, seed=42)
df_1 = df_otpw.filter(col("DEP_DEL15") == 1)

# Now both DataFrames have the same number of rows for DEP_DEL15 being 0 and 1
df_otpw = df_0.union(df_1)

# Check the counts after balancing
df_otpw.groupBy("DEP_DEL15").count().show()
```

Cancelled

```
df_otpw.groupBy(['YEAR']).count().show()
```

Cancelled

Create DateTime var

```
df_otpw = df_otpw.withColumn("DAY_OF_MONTH", lpad(col("DAY_OF_MONTH"), 2, "0"))
df_otpw = df_otpw.withColumn("MONTH", lpad(col("MONTH"), 2, "0"))

date_string = concat(
    col("YEAR").cast("string"),
    lit("-"),
    col("MONTH"),
    lit("-"),
    col("DAY_OF_MONTH")
)
df_otpw=df_otpw.withColumn("DATE_VARIABLE", to_date(date_string, 'yyyy-MM-dd'))
#df_otpw=df_otpw.withColumn("DATE_VARIABLE", date_string)
# df_otpw.printSchema()

# Format CRS_DEP_TIME to HH:MM format
df_otpw = df_otpw.withColumn("CRS_DEP_TIME_INT", col("CRS_DEP_TIME").cast("integer"))
df_otpw = df_otpw.withColumn("CRS_DEP_TIME_STR", format_string("%04d", col("CRS_DEP_TIME_INT")))

# Create the HH:MM Format
df_otpw = df_otpw.withColumn("CRS_DEP_TIME_FMT", concat(col("CRS_DEP_TIME_STR").substr(1, 2), lit(":"),
col("CRS_DEP_TIME_STR").substr(3, 2)))

# Combine 'DATE_VARIABLE' and 'CRS_DEP_TIME_FMT' into a single timestamp
df_otpw = df_otpw.withColumn(
    "DateTime",
    to_timestamp(concat(col("DATE_VARIABLE"), lit(" "), col("CRS_DEP_TIME_FMT")), "yyyy-MM-dd HH:mm")
)

#
# df_otpw.select("CRS_DEP_TIME", "CRS_DEP_TIME_FMT", "DateTime").show(100, False)
```

Cancelled

Filter Data

```
# Drop all observations with null of target variable
df_otpw = df_otpw.dropna(subset=['DEP_DEL15'])

# Cast as numeric variables into numeric format from string
df_otpw = df_otpw.withColumn("DEP_DEL15", col("DEP_DEL15").cast('int')) \
    .withColumn("DEP_DELAY_NEW", regexp_replace("DEP_DELAY_NEW", "s", "").cast('int')) \
    .withColumn("YEAR", regexp_replace("YEAR", "s", "").cast('int')) \
    .withColumn("MONTH", regexp_replace("MONTH", "s", "").cast('int')) \
    .withColumn("DAY_OF_MONTH", regexp_replace("DAY_OF_MONTH", "s", "").cast('int')) \
    .withColumn("DAY_OF_WEEK", regexp_replace("DAY_OF_WEEK", "s", "").cast('int')) \
    .withColumn("HourlyDewPointTemperature", regexp_replace("HourlyDewPointTemperature", "s", "").cast('int')) \
    .withColumn("HourlyDryBulbTemperature", regexp_replace("HourlyDryBulbTemperature", "s", "").cast('int')) \
    .withColumn("HourlyWetBulbTemperature", regexp_replace("HourlyWetBulbTemperature", "s", "").cast('int')) \
    .withColumn("HourlyRelativeHumidity", regexp_replace("HourlyRelativeHumidity", "s", "").cast('int')) \
    .withColumn("HourlyWindDirection", regexp_replace("HourlyWindDirection", "s", "").cast('int')) \
    .withColumn("HourlyWindSpeed", regexp_replace("HourlyWindSpeed", "s", "").cast('int')) \

# # Format CRS_DEP_TIME to HH:MM format
# df_otpw = df_otpw.withColumn("CRS_DEP_TIME_STR", format_string("%04d", col("CRS_DEP_TIME")))
# df_otpw = df_otpw.withColumn("CRS_DEP_TIME_FMT", concat(col("CRS_DEP_TIME_STR").substr(1, 2), lit(":"),
# col("CRS_DEP_TIME_STR").substr(3, 2)))

# # Combine 'FL_DATE' and 'CRS_DEP_TIME' into a single timestamp
# df_otpw = df_otpw.withColumn(
#     "DateTime",
#     to_timestamp(concat(col("FL_DATE"), lit(" "), col("CRS_DEP_TIME_FMT")), "yyyy-MM-dd HH:mm")
# )
df_otpw.printSchema()
```

```
root
|-- DEP_DEL15: integer (nullable = true)
|-- DEP_DELAY_NEW: integer (nullable = true)
|-- ORIGIN_AIRPORT_ID: string (nullable = true)
|-- DEST_AIRPORT_ID: string (nullable = true)
|-- FL_DATE: string (nullable = true)
|-- CRS_DEP_TIME: string (nullable = true)
|-- DAY_OF_MONTH: integer (nullable = true)
|-- DAY_OF_WEEK: integer (nullable = true)
|-- OP_UNIQUE_CARRIER: string (nullable = true)
|-- DEP_TIME_BLK: string (nullable = true)
|-- TAIL_NUM: string (nullable = true)
|-- MONTH: integer (nullable = true)
|-- YEAR: integer (nullable = true)
|-- HourlyDewPointTemperature: integer (nullable = true)
|-- HourlyDryBulbTemperature: integer (nullable = true)
|-- HourlyWetBulbTemperature: integer (nullable = true)
|-- HourlyRelativeHumidity: integer (nullable = true)
|-- HourlyWindDirection: integer (nullable = true)
|-- HourlyWindSpeed: integer (nullable = true)
|-- DATE_VARIABLE: date (nullable = true)
```

EDA: Select features for correlation

```
corr_features = df_otpw.select(['DEP_DELAY_NEW', 'HourlyDewPointTemperature', 'HourlyDryBulbTemperature',
'HourlyRelativeHumidity', 'HourlyWetBulbTemperature', 'HourlyWindDirection', 'HourlyWindSpeed'])

corr_features_df = corr_features.toPandas()
corr_features_df.dtypes
```

Cancelled

EDA: Import Packages

```
##
```

Cancelled

EDA: Correlation matrix

```
# Spearman Correlation Matrix
# Spearman's correlation evaluates the monotonic relationship between two ranked variables. Instead of using raw
data, it uses the rank order of the data.
# Preferred this method since we have non-normal distributions and ordinal variables
NUM_FEATURES = ['DEP_DELAY_NEW', 'HourlyDewPointTemperature', 'HourlyDryBulbTemperature', 'HourlyRelativeHumidity',
'HourlyWetBulbTemperature', 'HourlyWindDirection', 'HourlyWindSpeed']
df_otpw_features = df_otpw.select(NUM_FEATURES).dropna()

# Assemble the columns into a single feature column
assembler = VectorAssembler(inputCols=df_otpw_features.columns, outputCol="features")
df_assembled = assembler.transform(df_otpw_features)

# Calculate the correlation matrix
corr_matrix = Correlation.corr(df_assembled, "features", method = 'spearman').head()
corr_matrix = corr_matrix[0].toArray()
print(corr_matrix)
```

Cancelled

EDA: Show heat map

```
# Plot correlation
fig, ax = plt.subplots()
im = ax.imshow(corr_matrix, cmap='coolwarm')

# Add colorbar
cbar = ax.figure.colorbar(im, ax=ax)

# Set ticks and tick labels
ax.set_xticks(np.arange(corr_matrix.shape[1]))
ax.set_yticks(np.arange(corr_matrix.shape[0]))
ax.set_xticklabels(NUM_FEATURES)
ax.set_yticklabels(NUM_FEATURES)

# Rotate the tick labels and set their alignment
plt.setp(ax.get_xticklabels(), rotation=45, ha="right", rotation_mode="anchor")

# Set title
ax.set_title("Spearman Correlation Matrix")

# Show the plot
plt.show()
```

Cancelled

Split data into train and test

```

from pyspark.sql import functions as F
from pyspark.sql.window import Window

## add checkpoint
# spark.sparkContext.setCheckpointDir("/ml_41/df_otpw_checkpoint")

total_records = df_otpw.count()
# df_otpw2.checkpoint()

# Specify the percentage for the training set (e.g., 80%)
training_percentage = 0.8

# Calculate the number of records for the training set
training_records = int(total_records * training_percentage)

# Determine the split_date based on the calculated training_records
window_spec = Window.orderBy("DATE_VARIABLE")
df_otpw_with_rank = df_otpw.withColumn("rank", F.row_number().over(window_spec))
split_date_row = df_otpw_with_rank.filter(col("rank") == training_records).select("DATE_VARIABLE").collect()[0]
split_date = split_date_row["DATE_VARIABLE"]

# drop rank
df_otpw = df_otpw.drop("rank")

# # Filter the DataFrame based on the dynamically determined split_date (Turn on for 5Y)
train_df = df_otpw.filter(col("DATE_VARIABLE") < ('2018-01-01'))
test_intermediate_df = df_otpw.filter((col("DATE_VARIABLE") >= '2018-01-01') & (col("DATE_VARIABLE") < '2019-01-01'))
test_df = df_otpw.filter(col("DATE_VARIABLE") >= ('2019-01-01'))

# Filter the DataFrame based on the dynamically determined split_date (Turn on for 3M or 1Y)
# full_train_df = df_otpw.filter(col("DATE_VARIABLE") < split_date)

# Create intermediate DF
# total_train_records = full_train_df.count()
# int_training_percentage = 0.8
# # Calculate the number of records for the training set
# int_training_records = int(total_train_records * int_training_percentage)
# # Determine the split_date based on the calculated training_records
# window_spec = Window.orderBy("DATE_VARIABLE")
# train_df_with_rank = full_train_df.withColumn("rank", F.row_number().over(window_spec))

# train_df = train_df_with_rank.filter((train_df_with_rank['rank'] < int_training_records))
# test_intermediate_df = train_df_with_rank.filter((train_df_with_rank['rank'] >= int_training_records))
# test_df = df_otpw.filter(col("DATE_VARIABLE") >= split_date)

print("df_otpw:", df_otpw.count())
#print("full_train_df:", full_train_df.count())
print("train_df:", train_df.count())
print("test_intermediate_df:", test_intermediate_df.count())
print("test_df:", test_df.count())

# print("test")
# test_df.show(1)
# test_df.select("DATE_VARIABLE").distinct().show(365)
# cache dataframes for performance
train_df.cache()
test_df.cache()

```

```

df_otpw: 31197330
train_df: 16845634
test_intermediate_df: 7087730
test_df: 7263966

```



```
DataFrame[DEP_DEL15: int, DEP_DELAY_NEW: int, ORIGIN_AIRPORT_ID: string, DEST_AIRPORT_ID: string, FL_DATE: string, CRS_DEP_TIME: string, DAY_OF_MONTH: int, DAY_OF_WEEK: int, OP_UNIQUE_CARRIER: string, DEP_TIME_BLK: string, TAIL_NUM: string, MONTH: int, YEAR: int, HourlyDewPointTemperature: int, HourlyDryBulbTemperature: int, HourlyWetBulbTemperature: int, HourlyRelativeHumidity: int, HourlyWindDirection: int, HourlyWindSpeed: int, DATE_VARIABLE: date, CRS_DEP_TIME_INT: int, CRS_DEP_TIME_STR: string, CRS_DEP_TIME_FMT: string, DateTime: timestamp]
```

Define function to impute missing values with moving averages

```
from pyspark.sql import Window, DataFrame
from pyspark.sql.types import TimestampType
from datetime import timedelta

def calculate_moving_averages(df: DataFrame, range_low: int, range_high: int) -> DataFrame:
    """
    Applies moving average to specified columns in the DataFrame.

    Args:
    df: The input DataFrame.
    range_low: The lower bound of the window for the moving average.
    range_high: The upper bound of the window for the moving average.

    Returns:
    DataFrame with moving averages applied.
    """
    # Define window specification
    windowSpec = Window.partitionBy("ORIGIN_AIRPORT_ID").orderBy("FL_DATE", "CRS_DEP_TIME").rowsBetween(range_low, range_high)

    # Columns to calculate moving average
    weather_columns = ['HourlyDewPointTemperature', 'HourlyDryBulbTemperature', 'HourlyRelativeHumidity', 'HourlyWindDirection', 'HourlyWindSpeed']

    # Calculate moving average and impute missing values
    for column in weather_columns:
        moving_avg_col = f"{column}_MovingAvg"
        df = df.withColumn(moving_avg_col, avg(col(column)).over(windowSpec))

    # Replace null values in the original column with the moving average
    df = df.withColumn(column, when(col(column).isNull(), col(moving_avg_col)).otherwise(col(column)))

    # Optionally, drop the moving average column if it's no longer needed
    df = df.drop(moving_avg_col)

    return df
```

Time-based features: Seasonality, Part of Day, Recency, Past Flight Delays

```
def add_seasonality_feature(df: DataFrame, date_column: str) -> DataFrame:
    """
    Adds a 'Season' column to the DataFrame based on the month of the year from the specified date column.

    Args:
    df: The input DataFrame.
    date_column: The name of the column containing date information.

    Returns:
    DataFrame with an additional 'Season' column.
    """
    # Add a seasonality column
    df_with_season = df.withColumn("Season",
                                    when((month(date_column) >= 3) & (month(date_column) <= 5), lit("Spring"))
                                    .when((month(date_column) >= 6) & (month(date_column) <= 8), lit("Summer"))
                                    .when((month(date_column) >= 9) & (month(date_column) <= 11), lit("Autumn"))
                                    .otherwise(lit("Winter")))

    return df_with_season
```

```
def add_part_of_day_feature(df: DataFrame, timestamp_column: str) -> DataFrame:
    """
    Adds a 'PartOfDay' column to the DataFrame based on the time of day from the specified timestamp column.

    Args:
    df: The input DataFrame.
    timestamp_column: The name of the column containing timestamp information.

    Returns:
    DataFrame with an additional 'PartOfDay' column.
    """
    # Function to categorize part of day
    def part_of_day(hour_col):
        return (when((hour_col >= 6) & (hour_col < 12), lit("Morning"))
                .when((hour_col >= 12) & (hour_col < 18), lit("Afternoon"))
                .when((hour_col >= 18) & (hour_col < 24), lit("Evening"))
                .otherwise(lit("Night")))

    # Add a part of day column
    df_with_part_of_day = df.withColumn("PartOfDay", part_of_day(hour(timestamp_column)))

    return df_with_part_of_day
```

```
def add_recency_feature(df: DataFrame, timestamp_column: str, partition_column: str) -> DataFrame:
    """
    Adds a 'Recency' column to the DataFrame, indicating the time difference in hours from the last delay
    at the same partition (like an airport).

    Args:
    df: The input DataFrame.
    timestamp_column: The name of the column containing timestamp information.
    partition_column: The column name to partition by (e.g., airport ID).

    Returns:
    DataFrame with an additional 'Recency' column.
    """
    # Define the window specification
    windowSpec = Window.partitionBy(partition_column).orderBy(timestamp_column)

    # Calculate the time difference from the last event
    df_with_recency = df.withColumn("LastEventTime", lag(timestamp_column).over(windowSpec))
    df_with_recency = df_with_recency.withColumn("Recency",
                                                (col(timestamp_column).cast("long") -
                                                 col("LastEventTime").cast("long")) / 3600) # in hours

    return df_with_recency
```

```
def add_prev_flight_delay_feature(df, tail_num_col, delay_flag_col, datetime_col):
    """
    Adds a feature to indicate whether the previous flight for the same plane was delayed.

    Args:
    df (DataFrame): Spark DataFrame containing flight data.
    tail_num_col (str): Column name for the tail number identifying each plane.
    delay_flag_col (str): Column name indicating whether a flight is delayed.
    datetime_col (str): Column name for the datetime of each flight.

    Returns:
    DataFrame: Modified DataFrame with an additional 'Prev_Flight_Delayed' column.
    """
    # Define window specification: partition by plane and order by datetime
    windowSpec = Window.partitionBy(tail_num_col).orderBy(datetime_col)

    # Create 'Prev_Flight_Delayed' column using lag function
    df_with_feature = df.withColumn("Prev_Flight_Delayed", lag(delay_flag_col, 1).over(windowSpec))

    # Replace nulls in 'Prev_Flight_Delayed' for first flights of each plane
    df_with_feature = df_with_feature.fillna({"Prev_Flight_Delayed": 0})

    return df_with_feature
```

```

# Impute null values by calling the calculate_moving_averages for train_df
train_df = calculate_moving_averages(train_df, -10, -1)
train_df = calculate_moving_averages(train_df, -16, -1)
train_df = calculate_moving_averages(train_df, -30, -1)
train_df = calculate_moving_averages(train_df, -30, -1)
train_df = calculate_moving_averages(train_df, -50, -1)

# Call the add_seasonality_feature function on train_df
train_df = add_seasonality_feature(train_df, "FL_DATE")

# Call the add_seasonality_feature function on test_df
test_intermediate_df = add_seasonality_feature(test_intermediate_df, "FL_DATE")

# Call the add_part_of_day_feature function on train_df
train_df = add_part_of_day_feature(train_df, "DateTime")

# Call the add_part_of_day_feature function on test_df
test_intermediate_df = add_part_of_day_feature(test_intermediate_df, "DateTime")

# Call the add_recency_feature function on train_df
train_df = add_recency_feature(train_df, "DateTime", "ORIGIN_AIRPORT_ID")

# Call the add_recency_feature function on test_df
test_intermediate_df = add_recency_feature(test_intermediate_df, "DateTime", "ORIGIN_AIRPORT_ID")

# Call the add_prev_flight_delay_feature function on train_df
train_df = add_prev_flight_delay_feature(train_df, "TAIL_NUM", "DEP_DEL15", "DateTime")

# Call the add_prev_flight_delay_feature function on test_df
test_intermediate_df = add_prev_flight_delay_feature(test_intermediate_df, "TAIL_NUM", "DEP_DEL15", "DateTime")

```

Create PageRank Features

```

# from functools import reduce
# from pyspark.sql.functions import col, lit, when
# from graphframes import *

```

```

# # Select unique airport IDs from the DataFrame
# vertices = train_df.select("ORIGIN_AIRPORT_ID").distinct()
# vertices = vertices.withColumnRenamed("ORIGIN_AIRPORT_ID", "id")

```

```

# edges = train_df.groupBy(['ORIGIN_AIRPORT_ID', 'DEST_AIRPORT_ID']).agg({'DEP_DEL15': 'sum'})
# edges = edges.withColumnRenamed("ORIGIN_AIRPORT_ID", "src")
# edges = edges.withColumnRenamed("DEST_AIRPORT_ID", "dst")
# edges = edges.withColumnRenamed("sum(DEP_DEL15)", "sum_delays")

```

```

# g = GraphFrame(vertices, edges)

```

```

# results = g.pageRank(resetProbability=0.15, tol=0.01)

```

```

# display(results.vertices)

```

```
# vertices_df = results.vertices.select("id", "pagerank")
# vertices_df = vertices_df.withColumnRenamed("id", "ORIGIN_AIRPORT_ID")
# train_df = train_df.join(vertices_df, on=['ORIGIN_AIRPORT_ID'], how='left')
```

Set Response and Predictor Variables

```
print("-----")
print("Setting variables to predict flight delays")
# Target variable
myY = "DEP_DEL15"
# Categorical predictor variables
categoricals = ['DAY_OF_MONTH', 'DAY_OF_WEEK', 'OP_UNIQUE_CARRIER', 'DEP_TIME_BLK', 'MONTH']
# Numeric predictor variables
numerics = ['YEAR', 'HourlyDewPointTemperature', 'HourlyDryBulbTemperature', 'HourlyRelativeHumidity',
'HourlyWindDirection', 'HourlyWindSpeed']
# All predictor variables in a single list
myX = categoricals + numerics

# Create new dataframe with predictor vars, target var, etc.
train_df_bl = train_df.select(myX + [myY, "DATE_VARIABLE"])
```

Setting variables to predict flight delays

Baseline Logistic Regression

```

from pyspark.ml import Pipeline
from pyspark.ml.feature import StringIndexer, VectorAssembler, OneHotEncoder
from pyspark.ml.feature import StandardScaler, Imputer
from pyspark.ml.classification import LogisticRegression
from pyspark.ml.evaluation import BinaryClassificationEvaluator
from pyspark.ml.tuning import CrossValidator, ParamGridBuilder

## Current possible ways to handle categoricals in string indexer is 'error', 'keep', and 'skip'
# For each column in categoricals, create "StringIndexer" object to index categorical features by assigning
# them numerical values. Then convert categorical indices created by 'StringIndexer' into a sparse vector.
# Imputer used for imputing missing values and replacing them with mean or median of column.
# indexers = map(lambda c: StringIndexer(inputCol=c, outputCol=c+"_idx", handleInvalid = 'keep'), categoricals)
# ohes = map(lambda c: OneHotEncoder(inputCol=c + "_idx", outputCol=c+"_class"),categoricals)
indexers = map(lambda c: StringIndexer(inputCol=c, outputCol=c+"_idx", handleInvalid='keep'), categoricals)
ohes = map(lambda c: OneHotEncoder(inputCol=c+"_idx", outputCol=c+"_class"), categoricals)
imputers = Imputer(inputCols = numerics, outputCols = numerics)

# Establish features columns
# Categorical variables with "_class" appended for each of the categoricals to use the sparse-encoded vars
# Numerics
featureCols = list(map(lambda c: c+"_class", categoricals)) + numerics

# Build the stage for the ML pipeline
# List contains all transformation stages necessary for preprocessing data:
# -indexing and encoding categorical variables
# -imputing missing values in numerical columns
# -assembling feature columns into a vector
# -indexing the target variable
model_matrix_stages = list(indexers) + list(ohes) + [imputers] + \
    [VectorAssembler(inputCols=featureCols, outputCol="features"),
StringIndexer(inputCol="DEP_DEL15", outputCol="label",handleInvalid='skip')]

# Apply StandardScaler to create scaledFeatures
# Specifies input column containing the features that need to be scaled, 'features'
# specifies output column where scaled features will be stored 'scaledFeatures'
# withSTD std will be 1
# withMean means the mean will be 0
scaler = StandardScaler(inputCol="features",
                        outputCol="scaledFeatures",
                        withStd=True,
                        withMean=True)

```

Define blocking Time Series Split

```
class BlockingTimeSeriesSplit:
    def __init__(self, n_splits=5):
        self.n_splits = n_splits

    def split(self, X, y=None, groups=None):

        n_samples = X.count()
        print("n_samples " + str(n_samples))
        k_fold_size = n_samples // self.n_splits
        print("k_fold_size " + str(k_fold_size))

        margin = 0

        for i in range(self.n_splits):
            start = i * k_fold_size

            print("start " + str(start))
            stop = start + k_fold_size
            print("stop " + str(stop))
            mid = int(0.75 * (stop - start)) + start
            print("mid " + str(mid))
            print("mid+margin " + str(mid+margin))
            print("X " + str(X.count()))

            X_with_rank = X.withColumn("rank", F.row_number().over(window_spec))
            train_train_indices = X_with_rank.filter((X_with_rank['rank'] >= start) & (X_with_rank['rank'] < mid))
            train_valid_indices = X_with_rank.filter((X_with_rank['rank'] >= mid + margin) & (X_with_rank['rank'] <
stop))

            print("train_train_indices " + str(train_train_indices.count()))
            print("train_valid_indices " + str(train_valid_indices.count()))
            print("train indices ")

            print("train valid indices ")

            yield train_train_indices, train_valid_indices
```

Define getMetrics function

```
def getMetrics(predictions):  
    # Calculate precision, recall, and AUC  
    binary_evaluator = BinaryClassificationEvaluator()  
    multi_evaluator = MulticlassClassificationEvaluator()  
  
    auc = binary_evaluator.evaluate(predictions)  
    tp = multi_evaluator.evaluate(predictions, {multi_evaluator.metricName: 'truePositiveRateByLabel'})  
    fp = multi_evaluator.evaluate(predictions, {multi_evaluator.metricName: 'falsePositiveRateByLabel'})  
  
    # Precision: TP/(TP+FP)  
    precision = multi_evaluator.evaluate(predictions, {multi_evaluator.metricName: 'precisionByLabel',  
    multi_evaluator.metricLabel: 1.0}) # should be list  
    # Recall: TP/(TP+FN)  
    recall = multi_evaluator.evaluate(predictions, {multi_evaluator.metricName: 'recallByLabel',  
    multi_evaluator.metricLabel: 1.0})  
    # F1: Harmonic mean of precision and recall  
    f1 = multi_evaluator.evaluate(predictions, {multi_evaluator.metricName: 'fMeasureByLabel',  
    multi_evaluator.metricLabel: 1.0})  
  
    # Display or use the precision, recall, and AUC values as needed  
    print(f"AUC: {auc}, Precision: {precision}, Recall: {recall}, F1 Score: {f1}")  
  
    return f1
```

Define Model Parameter for Baseline

```
# Use logistic regression  
# Set number of iterations for optimization algorithm, 10  
# Set input column to be used for training  
lr = LogisticRegression(maxIter=10, featuresCol = "scaledFeatures")  
  
# Build our ML pipeline  
pipeline_lr = Pipeline(stages=model_matrix_stages+[scaler]+[lr])
```


Run K-fold Cross validation on logistic regression (baseline) and evaluate test dat

```

from pyspark.ml.evaluation import MulticlassClassificationEvaluator
from pyspark.mllib.evaluation import MulticlassMetrics

# Define the time series split
btscv = BlockingTimeSeriesSplit(n_splits=5)

# cv_model_set_lr = []
# cv_model_set_lr_f1 = []

# for i, (train_index, valid_index) in enumerate(btscv.split(train_df_bl)):
#     print(f"Fold {i + 1}")

#     train_set = train_df_bl.filter(train_df_bl["index"].isin(train_index))
#     valid_set = train_df_bl.filter(train_df_bl["index"].isin(valid_index))

### Aditya's code ###
cv_model_set_lr = []
cv_model_set_lr_f1 = []

for i, (train_train_index, train_valid_index) in enumerate(btscv.split(train_df_bl)):
    print(f"Fold {i + 1}")
    print("prepping training...")

    train_train_index.show(1)
    print("train count "+ str(train_train_index.count()))
    #print(time.strftime('%H:%M%p %Z on %b %d, %Y'))

    train_valid_index.show(1)
    print("valid count "+str(train_valid_index.count()))
    #print(time.strftime('%H:%M%p %Z on %b %d, %Y'))

    train_train_index = train_train_index.checkpoint()
    train_valid_index = train_valid_index.checkpoint()

    col_to_drop = ["index", "rank", "DATE_VARIABLE"]
    train_train_index = train_train_index.drop(*col_to_drop)
    train_valid_index = train_valid_index.drop(*col_to_drop)

    spark.sparkContext.setCheckpointDir("/ml/5Y_checkpoint")
    #####
    ### fitting the model #####

    # ## lr pipeline
    # cv_model_lr = pipeline_lr.fit(train_train_index)
    # cv_model_set_lr.append(cv_model_lr)

    ## lr pipeline
    cv_model_lr = pipeline_lr.fit(train_train_index)
    model_name = f"cv_model_set_lr_{i}" # Creating unique model name
    cv_model_set_lr.append((model_name, cv_model_lr)) # Appending model with name to the list

    #print(time.strftime('%H:%M%p %Z on %b %d, %Y'))

    predictions_lr = cv_model_lr.transform(train_valid_index)
    cv_model_set_lr_f1.append(getMetrics(predictions_lr))

    #####
    ## best model LR EN #####
    # print("Running test data on best model for LR pipeline ....")

    # best_i_lr = cv_model_set_lr_f1.index(max(cv_model_set_lr_f1))
    # print(f"The best model is in fold {best_i_lr+1} ")

```

```
# print("Running test data on best model ....")
# cvModel_test = cv_model_set[best_i_lr]
# predictions_test = cvModel_test.transform(test_intermediate_df)
# test_f1 = getMetrics(predictions_test)
# print("Success pipeline for logistic regression finished")
```

```
n_samples 16845634
k_fold_size 3369126
start 0
stop 3369126
mid 2526844
mid+margin 2526844
X 16845634
train_train_indices 2526843
train_valid_indices 842282
train indicces
train valid indices
Fold 1
prepping training...
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|DAY_OF_MONTH|DAY_OF_WEEK|OP_UNIQUE_CARRIER|DEP_TIME_BLK|MONTH|YEAR|HourlyDewPointTemperature|HourlyDryBulbTemperature|Ho
urlyRelativeHumidity|HourlyWindDirection|HourlyWindSpeed|DEP_DEL15|DATE_VARIABLE|rank|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|            1|            4|            00| 1200-1259| 1|2015|                28.0|                49.0|
44.0|            260.0|            6.0|            0| 2015-01-01| 1|
Cancelled
```

```
# Get best performing model
print(cv_model_set_lr_f1)
```

NameError: name 'cv_model_set_lr_f1' is not defined

```
# predictions_test = cv_model_set_lr[0][1].transform(test_intermediate_df)
# test_f1 = getMetrics(predictions_test)
# print("Success pipeline for logistic regression finished")
```

```
# df = spark.createDataFrame([(x,) for x in cv_model_set_lr_f1], ["value"])
# ##best_i_lr = cv_model_set_lr_f1.index(max(cv_model_set_lr_f1))
# best_i_lr = int(df.agg(spark_max("value")).collect()[0][0])
# print(f"The best model is in fold {best_i_lr+1} ")
# print("Running test data on best model ....")
# cvModel_test = cv_model_set_lr[best_i_lr]
# predictions_test = cvModel_test.transform(test_intermediate_df)
# test_f1 = getMetrics(predictions_test)
# print("Success pipeline for logistic regression finished")
```

```
train_df.printSchema()
```

```
root
|-- ORIGIN_AIRPORT_ID: string (nullable = true)
|-- DEP_DEL15: integer (nullable = true)
|-- DEP_DELAY_NEW: integer (nullable = true)
|-- DEST_AIRPORT_ID: string (nullable = true)
|-- FL_DATE: string (nullable = true)
|-- CRS_DEP_TIME: string (nullable = true)
|-- DAY_OF_MONTH: integer (nullable = true)
|-- DAY_OF_WEEK: integer (nullable = true)
|-- OP_UNIQUE_CARRIER: string (nullable = true)
|-- DEP_TIME_BLK: string (nullable = true)
|-- TAIL_NUM: string (nullable = true)
|-- MONTH: integer (nullable = true)
```

```
|-- YEAR: integer (nullable = true)
|-- HourlyDewPointTemperature: double (nullable = true)
|-- HourlyDryBulbTemperature: double (nullable = true)
|-- HourlyWetBulbTemperature: integer (nullable = true)
|-- HourlyRelativeHumidity: double (nullable = true)
|-- HourlyWindDirection: double (nullable = true)
|-- HourlyWindSpeed: double (nullable = true)
```

Feature Selection of Feature-Engineered Vars

```
print("-----")
print("Setting variables to predict flight delays")
# Target variable
myY = "DEP_DEL15"
# Categorical predictor variables
categoricals = ['DAY_OF_MONTH', 'DAY_OF_WEEK', 'OP_UNIQUE_CARRIER', 'DEP_TIME_BLK', 'MONTH', 'Season', 'PartOfDay']
# Numeric predictor variables
numerics = ['YEAR', 'HourlyDewPointTemperature', 'HourlyDryBulbTemperature', 'HourlyRelativeHumidity',
            'HourlyWindDirection', 'HourlyWindSpeed', 'Recency', 'Prev_Flight_Delayed']
# All predictor variables in a single list
myX = categoricals + numerics

# Create new dataframe with predictor vars, target var, etc.
train_df_adv = train_df.select(myX + [myY, "DATE_VARIABLE"])
```

Setting variables to predict flight delays

Create Model Classifier and Pipelines

```

from pyspark.ml import Pipeline
from pyspark.ml.feature import StringIndexer, VectorAssembler, OneHotEncoder
from pyspark.ml.feature import StandardScaler, Imputer
from pyspark.ml.classification import LogisticRegression
from pyspark.ml.evaluation import BinaryClassificationEvaluator
from pyspark.ml.tuning import CrossValidator, ParamGridBuilder

from xgboost.spark import SparkXGBClassifier
from pyspark.ml.classification import DecisionTreeClassifier
from pyspark.ml.classification import MultilayerPerceptronClassifier

lr_en = LogisticRegression(maxIter=10, elasticNetParam=0.5, featuresCol = "scaledFeatures")
pipeline_lr_en = Pipeline(stages=model_matrix_stages+[scaler]+[lr_en])

# XGB pipeline
xgb_classifier = SparkXGBClassifier(
    features_col="features",
    label_col="label",
    prediction_col="prediction",
    num_workers=2,
    device="cuda",
)
pipeline_xgb = Pipeline(stages=model_matrix_stages+[scaler]+[xgb_classifier])

# Define DecisionTreeClassifier
dt_classifier = DecisionTreeClassifier(
    featuresCol="features",
    labelCol="label",
    predictionCol="prediction")

# DT Pipeline
pipeline_dt = Pipeline(stages=model_matrix_stages + [scaler, dt_classifier])

## MLP
# specify layers for the neural network:
# input layer of size 4 (features), two intermediate of size 5 and 4
# and output of size 2 (classes)
# layers = [4, 5, 4, 2]

# # create the trainer and set its parameters
# trainer = MultilayerPerceptronClassifier(maxIter=100, layers=layers, blockSize=128, seed=1234)
# # model_matrix_stages + [scaler]+[mean_encoder] +
# pipeline_mlp = Pipeline(stages= [imputers] + [trainer])

# mlp_train_df = train_df
# mlp_test_df = test_df

```

Run model pipelines for K-folds and choose best model and evaluate on test data

```

from pyspark.ml.evaluation import MulticlassClassificationEvaluator
from pyspark.mllib.evaluation import MulticlassMetrics

# Define the time series split
btscv = BlockingTimeSeriesSplit(n_splits=5)

cv_model_set_lr_en = []
cv_model_set_lr_en_f1 = []

cv_model_set_xgb = []
cv_model_set_xgb_f1 = []

cv_model_set_dt = []
cv_model_set_dt_f1 = []

# cv_model_set_mlp = []
# cv_model_set_mlp_f1 = []

for i, (train_train_index, train_valid_index) in enumerate(btscv.split(train_df_adv)):
    print(f"Fold {i + 1}")
    print("prepping training...")

    train_train_index.show(1)
    print("train count "+ str(train_train_index.count()))
    #print(time.strftime('%H:%M%p %Z on %b %d, %Y'))

    train_valid_index.show(1)
    print("valid count "+str(train_valid_index.count()))
    #print(time.strftime('%H:%M%p %Z on %b %d, %Y'))

    train_train_index = train_train_index.cache()
    train_valid_index = train_valid_index.cache()

    col_to_drop = ["index", "rank", "DATE_VARIABLE"]
    train_train_index = train_train_index.drop(*col_to_drop)
    train_valid_index = train_valid_index.drop(*col_to_drop)

    spark.sparkContext.setCheckpointDir("/ml/5Y_checkpoint")
    #####
    ### fitting the model #####

    ## lr pipeline
    cv_model_lr_en = pipeline_lr_en.fit(train_train_index)
    model_name = f"cv_model_lr_en_{i}" # Creating unique model name
    cv_model_set_lr_en.append((model_name, cv_model_lr_en)) # Appending model with name to the list
    #print(time.strftime('%H:%M%p %Z on %b %d, %Y'))

    predictions_lr_en = cv_model_lr_en.transform(train_valid_index)
    print(f"LR EN Fold {i + 1} Metrics .....")
    cv_model_set_lr_en_f1.append(getMetrics(predictions_lr_en))

    ## xgb pipeline
    cv_model_xgb = pipeline_xgb.fit(train_train_index)
    model_name = f"cv_model_xgb_{i}" # Creating unique model name
    cv_model_set_xgb.append((model_name, cv_model_xgb)) # Appending model with name to the list
    #print(time.strftime('%H:%M%p %Z on %b %d, %Y'))

    predictions_xgb = cv_model_xgb.transform(train_valid_index)
    print(f"XGB Fold {i + 1} Metrics .....")
    cv_model_set_xgb_f1.append(getMetrics(predictions_xgb))

    ## dt pipeline

```

```
n_samples 16845634
k_fold_size 3369126
start 0
stop 3369126
mid 2526844
mid+margin 2526844
X 16845634
train_train_indices 2526843
train_valid_indices 842282
train indicces
train valid indices
Fold 1
prepping training...
+-----+
+-----+
+-----+
|DAY_OF_MONTH|DAY_OF_WEEK|OP_UNIQUE_CARRIER|DEP_TIME_BLK|MONTH|Season|PartOfDay|YEAR|HourlyDewPointTemperature|HourlyDryBulbTemperature|HourlyRelativeHumidity|HourlyWindDirection|HourlyWindSpeed|Recency|Prev_Flight_Delayed|DEP_DEL15|DATE_VARIABLE|rank|
+-----+
+-----+
/databricks/python/lib/python3.10/site-packages/xgboost/sklearn.py:782: UserWarning: Loading a native XGBoost model with Scikit-Learn interface.
  warnings.warn("Loading a native XGBoost model with Scikit-Learn interface.")
```

```
XGB Fold 1 Metrics .....
AUC: 0.6820255192275666, Precision: 0.46565127716960175, Recall: 0.015252933656234735, F1 Score: 0.029538306705245942
DT Fold 1 Metrics .....
AUC: 0.5300075059992718, Precision: 0.0, Recall: 0.0, F1 Score: 0.0
start 3369126
stop 6738252
mid 5895970
mid+margin 5895970
X 16845634
train_train_indices 2526844
train_valid_indices 842282
train indices
train valid indices
Fold 2
prepping training...
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|DAY_OF_MONTH|DAY_OF_WEEK|OP_UNIQUE_CARRIER|DEP_TIME_BLK|MONTH|Season|PartOfDay|YEAR|HourlyDewPointTemperature|HourlyDryB
ulbTemperature|HourlyRelativeHumidity|HourlyWindDirection|HourlyWindSpeed|Recency|Prev_Flight_Delayed|DEP_DEL1
5|DATE_VARIABLE| rank|
```

Table							
	DAY_OF_MONTH ▲	DAY_OF_WEEK ▲	OP_UNIQUE_CARRIER ▲	DEP_TIME_BLK ▲	MONTH ▲	Season ▲	PartOfDay ▲
1	1	4	AA	0700-0759	1	Winter	Morning
2	1	4	AA	2000-2059	1	Winter	Evening
3	2	5	AA	0800-0859	1	Winter	Morning
4	6	2	AA	0900-0959	1	Winter	Morning
5	6	2	AA	1100-1159	1	Winter	Morning

6	6	2	AA	1700-1750	1	Winter	Afternoon
---	---	---	----	-----------	---	--------	-----------

10,000 rows | Truncated data

```
# Get best performing LR EN model
print(cv_model_set_lr_en_f1)
```

```
[0.0003636684798138017, 0.001282203490442835, 0.027000883907873453, 0.002132350283711012, 0.003693772361618769]
```

```
predictions_test = cv_model_set_lr_en[2][1].transform(test_intermediate_df)
test_f1 = getMetrics(predictions_test)
print("Success pipeline for logistic regression finished")
```

```
AUC: 0.655669707182829, Precision: 0.45839388973066575, Recall: 0.011368002263784585, F1 Score: 0.022185804617053915
Success pipeline for logistic regression finished
```

```
# Get best performing XGB model
print(cv_model_set_xgb_f1)
```

```
[0.029538306705245942, 0.015577718767368616, 0.06661190705609174, 0.029263151731608067, 0.02002813870727468]
```

```
# predictions_test = cv_model_set_xgb[4][1].transform(test_intermediate_df)
# test_f1 = getMetrics(predictions_test)
# print("Success pipeline for logistic regression finished")
```

```
# Get best performing DT model
print(cv_model_set_dt_f1)
```

```
[0.16890825757006536, 0.0, 0.27898878752002226, 0.0008866034222892101, 0.0009792720744246776]
```

```
# predictions_test = cv_model_set_dt_f1[4][1].transform(test_intermediate_df)
# test_f1 = getMetrics(predictions_test)
# print("Success pipeline for logistic regression finished")
```

```
for i, model in enumerate(cv_model_set_lr):
    # Example: Checkpointing each model
    checkpoint_dir = f"/ml_41/df_otpw_checkpoint/cv_model_set_lr_{i}"
    model.write().overwrite().save(checkpoint_dir)
```

```
AttributeError: 'tuple' object has no attribute 'write'
```

```
for i, model in enumerate(cv_model_set_lr_en):
    # Example: Checkpointing each model
    checkpoint_dir = f"/ml_41/df_otpw_checkpoint/cv_model_set_lr_en_{i}"
    model.write().overwrite().save(checkpoint_dir)
```

```
AttributeError: 'tuple' object has no attribute 'write'
```

```
for i, model in enumerate(cv_model_set_xgb):
    # Example: Checkpointing each model
    checkpoint_dir = f"/ml_41/df_otpw_checkpoint/cv_model_set_xgb_{i}"
    model.write().overwrite().save(checkpoint_dir)
```

```
AttributeError: 'tuple' object has no attribute 'write'
```

```
AttributeError: 'tuple' object has no attribute 'write'
```

