

Finding Flight Delay Trends

DAT500 Project Group 17

Bhakti Prabhu & Stephan F. W. Brandasu

Faculty of Science and Technology

Why & What

the use case

As a traveller:

1. Not miss any important meetings/events/functions.
2. Pre-plan journey
3. Have idea of buffer time while flight booking

As an airline:

1. know when to increase workforce.
2. Opportunity to improve over competitors.

Objective:

1. visualise delay trends over time
2. how well do the airlines catch up when departing late
3. how often are flights cancelled
4. are longer flights delayed more often

The Dataset

From 20GB of unstructured .txt files to 6GB of .csv

```
ubuntu@namenode:~$ hadoop fs -du -h /txt
382.2 M 764.4 M /txt/2018-01.txt
349.4 M 698.9 M /txt/2018-02.txt
410.1 M 820.2 M /txt/2018-03.txt
400.0 M 799.9 M /txt/2018-04.txt
413.9 M 827.9 M /txt/2018-05.txt
420.6 M 841.2 M /txt/2018-06.txt
433.6 M 867.2 M /txt/2018-07.txt
428.1 M 856.2 M /txt/2018-08.txt
392.9 M 785.7 M /txt/2018-09.txt
414.5 M 828.9 M /txt/2018-10.txt
394.6 M 789.3 M /txt/2018-11.txt
399.7 M 799.4 M /txt/2018-12.txt
391.6 M 783.3 M /txt/2019-01.txt
358.1 M 716.1 M /txt/2019-02.txt
423.9 M 847.7 M /txt/2019-03.txt
410.5 M 821.0 M /txt/2019-04.txt
427.2 M 854.4 M /txt/2019-05.txt
428.0 M 856.0 M /txt/2019-06.txt
```

```
ubuntu@namenode:~$ hadoop fs -du -h /csv
1.4 G 2.8 G /csv/2018.csv
1.4 G 2.9 G /csv/2019.csv
936.6 M 1.8 G /csv/2020.csv
1.1 G 2.3 G /csv/2021.csv
1.2 G 2.5 G /csv/2022.csv
5.9 G 11.8 G /csv/all-years.csv
```

The Dataset

```
year
quarter
month
day_of_month
day_of_week
fl_date
op_unique_carrier
tail_num
op_carrier_fl_num
origin_airport_id
origin_airport_seq_id
origin_city_market_id
origin
origin_city_name
origin_state_nm
dest_airport_id
dest_airport_seq_id
dest_city_market_id
dest
dest_city_name
dest_state_nm
dep_delay
dep_delay_new
dep_del15
arr_delay
arr_delay_new
arr_del15
cancelled
cancellation_code
diverted
air_time
distance
distance_group
carrier_delay
weather_delay
nas_delay
security_delay
late_aircraft_delay
```

```
year
month
fl_date
op_unique_carrier
origin_airport_id
dest_airport_id
dep_delay_new
arr_delay_new
cancelled
diverted
air_time
```

```
year
month
op_unique_carrier
origin_airport_id
dest_airport_id
max_arr_delay
max_arr_delay_fl_date
avg_arr_delay
med_arr_delay
avg_time_recovered
nr_diverted
avg_airtime
flight_count
nr_cancelled
```

Mapping

```
1 import sys
2
3 row = []
4
5 for line in sys.stdin:
6     line = line.strip().replace(',', '').split()
7
8     if line[0] == "LATE_AIRCRAFT_DELAY":
9         try:
10             data = line[1]
11         except IndexError:
12             data = ' '
13         row.append(data)
14         print(','.join(row))
15         row = []
16     else:
17         try:
18             data = ' '.join(line[1:])
19         except IndexError:
20             data = ' '
21         row.append(data)
```

- create array of each row
- remove any unwanted commas from cell
- check for last column of row
- catch error in case of no data in current column
- print row

Reading the data

```
1 flight_data = spark.read.csv('hdfs://namenode
   :9000/csv/'+sys.argv[1]+''.csv', schema=
   flightSchema)\n
2 .withColumn('FL_DATE', to_date(to_timestamp('
   FL_DATE', 'M/d/yyyy h:mm:ss a')))\n
3
4 flight_data = flight_data.select( 'year'
5                                   , 'month'
6                                   , 'fl_date'
7                                   , 'op_unique_carrier'
8                                   , 'origin_airport_id'
9                                   , 'dest_airport_id'
10                                  , 'dep_delay_new'
11                                  , 'arr_delay_new'
12                                  , 'cancelled'
13                                  , 'diverted'
14                                  , 'air_time')\n
15
16 flight_data = flight_data.na.drop(subset=['year',
17   'origin_airport_id', 'dest_airport_id', '
18   fl_date'])\n
19 flight_data = flight_data.fillna({'arr_delay_new'
20   : 0.0})
```

- select only relevant columns
- drop any rows that are missing essential information
- fill 0.0 in rows to avoid issues during calculations

Manipulating the data - step 1

```
1 # grab the fl_date of the flight with the highest
   delay for a given group
2 windowSpec = Window.partitionBy(
3     'year'
4     , 'month'
5     , 'op_unique_carrier'
6     , 'origin_airport_id'
7     , 'dest_airport_id').orderBy(col('
   arr_delay_new').desc())
8
9 arr_delay_dates = flight_data.withColumn(
10     'rank'
11     , rank().over(windowSpec)
12 ).filter(
13     col('rank') == 1
14 ).groupBy(
15     'year'
16     , 'month'
17     , 'op_unique_carrier'
18     , 'origin_airport_id'
19     , 'dest_airport_id'
20 ).agg(
21     round(max('arr_delay_new'), 2).alias('
   max_arr_delay')
22 , first('fl_date').alias('max_arr_delay_fl_date')
23 )
```

- create a window sorted by the arrival delay
- select the top delay date by only selecting single row

Manipulating the data - step 2

```
1 flight_data = flight_data.groupby('year'
2 , 'month'
3 , 'op_unique_carrier'
4 , 'origin_airport_id'
5 , 'dest_airport_id').agg( round(avg('
    arr_delay_new'), 2).alias('avg_arr_delay')
6 , round(percentile_approx('arr_delay_new',
    0.5), 2).alias('med_arr_delay')
7 , round(avg(col('dep_delay_new') - col('
    arr_delay_new')), 2).alias('
    avg_time_recovered')
8 , sum('diverted').alias('nr_diverted')
9 , round(avg('air_time'), 2).alias('
    avg_airtime')
10 , count('*').alias('flight_count')
11 , sum('cancelled').alias('nr_cancelled'))
12
13 flight_data = arr_delay_dates.join( flight_data
14 , on=['year', 'month', 'op_unique_carrier', '
    origin_airport_id', 'dest_airport_id']
15 , how='left')
```

- do a groupby select for to grab delay statistics
- each result is rounded to 2 decimals to avoid ugly numbers
- join the result of this operation with the previous one

Upserting into the DeltaTable

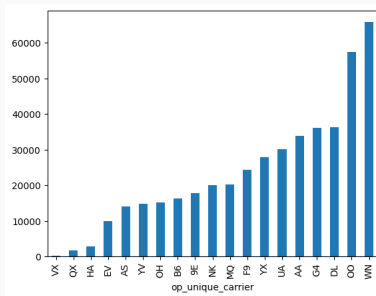
```
1 # Checking if Table exists
2 if DeltaTable.isDeltaTable(spark, "hdfs://namenode:9000/spark-warehouse/sample_flight_table"):
3     # Perform the upsert operation
4     deltaDF = DeltaTable.forPath(spark, "hdfs://namenode:9000/spark-warehouse/sample_flight_table")
5     merge_condition = "existing.year = upsert.year \
6     AND existing.month = upsert.month \
7     AND existing.op_unique_carrier = upsert.op_unique_carrier \
8     AND existing.origin_airport_id = upsert.origin_airport_id \
9     AND existing.dest_airport_id = upsert.dest_airport_id "
10
11     deltaDF.alias('existing') \
12         .merge(flight_data.alias('upsert'), merge_condition) \
13         .whenMatchedUpdateAll() \
14         .whenNotMatchedInsertAll() \
15         .execute()
16 else:
17     # Create new delta table
18     flight_data.write.format("delta").mode("overwrite").saveAsTable("sample_flight_table")
```

Skew & Spill

Optmization
(spark-defaults.conf)

Spill	Memory	Disk
Before	3.5 GiB	177.8 MiB
After	128.4 MiB	7.1 MiB

```
1 spark.executor.memory           6g
2 spark.executor.instances        3
3 spark.executor.cores            4
4 spark.sql.shuffle.partition     64
5 spark.sql.adaptive.skewedJoin.enabled true
6 spark.sql.adaptive.skewJoin.enabled true
```



Using UDF's to clean the data

```
1 def replace_null(value, default):
2     if value is None:
3         return default
4     return value
5
6 def drop_null(*cols):
7     for col in cols:
8         if col is None:
9             return False
10    return True
11
12 replace_null_udf = udf(lambda value, default:
13     replace_null(value, default), FloatType())
14 drop_null_udf = udf(lambda *cols: drop_null(*cols
15     ), BooleanType())
16
17 flight_data = flight_data.filter(drop_null_udf(*[
18     col(c) for c in ['year', '
19     origin_airport_id', 'dest_airport_id', '
20     fl_date']]))
21
22 flight_data = flight_data.withColumn('
23     arr_delay_new', replace_null_udf(col('
24     arr_delay_new'), lit(0.0)))
```

Performance difference between
each type of join

	normal	udf
minutes	2.8	7.1

Shuffle

sort merge join

```
1 arr_delay_dates = arr_delay_dates.sort(['year', 'month', 'op_unique_carrier', 'origin_airport_id', 'dest_airport_id'])
2 flight_data = flight_data.sort(['year', 'month', 'op_unique_carrier', 'origin_airport_id', 'dest_airport_id'])
3
4 # join the highest delay with the res of the group
5 flight_data = arr_delay_dates.join(flight_data
6 , on=['year', 'month', 'op_unique_carrier', 'origin_airport_id', 'dest_airport_id']
7 , how='left')
```

broadcast join

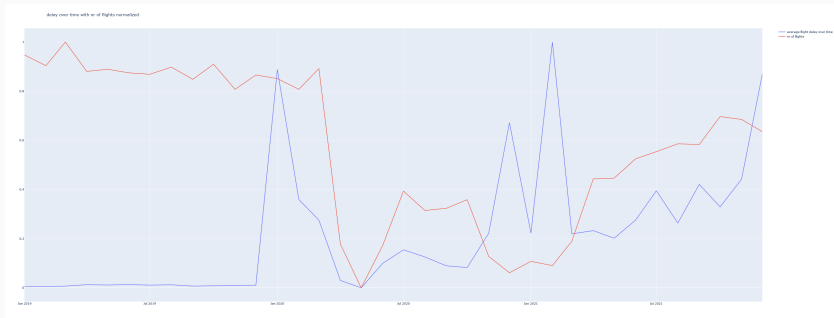
```
1 flight_data = arr_delay_dates.join(
2     broadcast(flight_data)
3     , on=['year', 'month', 'op_unique_carrier', 'origin_airport_id', 'dest_airport_id']
4     , how='left')
```

	Sort Merge	Broadcast
Minutes	4.1	3.2

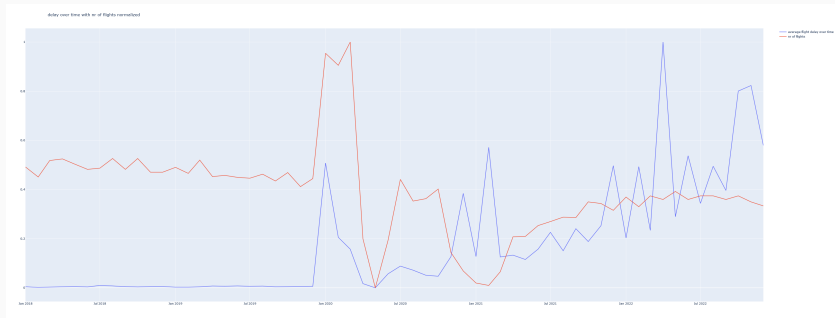
Graphs from the delta table

```
1 builder = SparkSession.builder.appName('flight_plot')
2 spark = configure_spark_with_delta_pip(builder).getOrCreate()
3
4 flight_data = spark.read.format('delta').load('hdfs://namenode:9000/spark-warehouse/flight_data_table'
5 )
6 flight_data = flight_data.filter((flight_data.op_unique_carrier == 'AA') & (flight_data.
7     origin_airport_id == 12892) & (flight_data.dest_airport_id == 12478))
8 flight_data = flight_data.orderBy('year', 'month')
9
10 flight_data = flight_data.withColumn('year_month', concat('year', lit('-'), 'month'))
11
12 flight_data = flight_data.toPandas()
13
14 avgdelay = normalize(flight_data, "avg_arr_delay")
15 flightcount = normalize(flight_data, "flight_count")
16
17 data = [
18     plt.Scatter(x=flight_data.year_month
19         , y=avgdelay
20         , name='average flight delay over time'
21         , text=flight_data.avg_arr_delay),
22     plt.Scatter(x=flight_data.year_month
23         , y=flightcount
24         , name='nr of flights'
25         , text=flight_data.flight_count),
26 ]
27
28 fig = plt.Figure(data, layout_title_text='delay over time with nr of flights normalized')
29 plot(fig, filename='plot.html')
```

A Graph!



A Graph! with a new year



A Graph! with bad information for 1 year

