Spark Assignment 2

Brian Cervantes Alvarez March 10, 2025



Overview

I used Spark on Google Dataproc to find out how far each plane traveled between location recordings. I used the haversine formula to calculate the shortest distance on a sphere. Since the raw data had errors (like a plane suddenly showing up at the South Pole), we removed data points that didn't make sense before adding up the distances. This made the final numbers much more realistic.

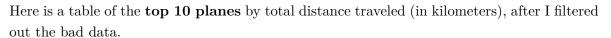
Process Description

I brought the plane data from BigQuery into Spark, fixed missing information, and removed rows with huge jumps that were likely mistakes. By focusing on reasonable coordinates and time sequences, I got better distance estimates. Additionally, I also created a new export folder each time to avoid conflicts in BigQuery. Below are screenshots of my Dataproc job runs and the Spark output:

25/83/11 66:25:53 IMFO FileInputFormat: Total input files to process: 24
Top 10 planes by total distance (km):
25/83/11 66:27:13 IMFO GoogleHadoopOutputStream: hflush(): No-op due to fate limit (RateLimiter[stableRate=0.2qps]): readers will *not* yet see flushed data for gs://
dataproc-temp-us-west1-762991648884-arm7uuc4/3834b4b7-d765-4918-8550-73348813f723/spark-job-history/application_17416787888967_0812.inprogress [CONTEXT ratelimit_period="1 MINUTES [skipped: 13]"]
[Row[Caoa 22C2C2C, total_dist_km=477865.1231855612], Row[Caoa* 282EBD*, total_dist_km=17992.8645493853], Row[Caoa* 4CC8DC*, total_dist_km=479865.231855612], Row[Caoa* 282EBD*, total_dist_km=69780.2739318836], Row[Caoa* 4053847*, total_dist_km=64366.588835942186),
Row[Caoa**462818*, total_dist_km=69727.331866231915], Row[Caoa**48187*, total_dist_km=69780.2739318836], Row[Caoa**405347*, total_dist_km=64366.588835942186),
Row[Caoa**462818*, total_dist_km=69727.331866231915], Row[Caoa**48187*, total_dist_km=69780.2739318836], Row[Caoa**405347*, total_dist_km=64366.588835942186),
Row[Caoa**462818*, total_dist_km=69727.331866231915], Row[Caoa**48187*, total_dist_km=69780.2739318836], Row[Caoa**405347*, total_dist_km=69780.473931836],
Row[Caoa**462818*, total_dist_km=69727.331866231915], Row[Caoa**462818*, total_dist_km=69780.27393183856], Row[Caoa**462849466419]]
Sum of all distances (km):
25/83/11 86:32-94 IMFO GhfsGlobalStorageStatistics: periodic connector metrics: {action_http_delete_request_duration=293, action_http_bost_request_max=194, action_http_bost_request_max=204, action_http_bost_request_max=204, action_http_bost_request_max=204, action_http_post_request_max=204, action_h

Output is complete





| Distance (km) |
|---------------|
| 477,805.12 |
| 179,792.06 |
| 139,628.47 |
| 139,000.06 |
| 72,371.39 |
| 69,870.27 |
| $64,\!366.59$ |
| 62,772.31 |
| $60,\!230.98$ |
| 59,411.56 |
| |

Total Distance

The total distance flown by all planes in the data is about 97,497,513 km.

Final Thoughts

By removing incorrect coordinates and making sure each plane's next point was valid, I stopped the data from showing unrealistic trips around the globe. We also fixed issues with Spark like export folder conflicts, missing data fields, and ClassNotFound errors (which I solved by adding the BigQuery connector). Working alone on this project helped me understand every step, from cleaning the data to performing the final analysis, so I now feel ready to move on to the final project.



Appendix

```
Oregon State
University
```

```
import pyspark
from pyspark.sql import SparkSession
import pprint
import json
from pyspark.sql.types import StructType, FloatType, LongType,
    StringType, StructField
from pyspark.sql import Window
from math import radians, cos, sin, asin, sqrt
from pyspark.sql.functions import lead, udf, col
import datetime
# Haversine: compute distance in km between two lat/lon points
def haversine(lon1, lat1, lon2, lat2):
    lon1, lat1, lon2, lat2 = map(radians, [lon1, lat1, lon2, lat2]) #
   convert to radians
   dlon = lon2 - lon1
   dlat = lat2 - lat1
   a = \sin(dlat / 2)**2 + \cos(lat1) * \cos(lat2) * \sin(dlon / 2)**2
   c = 2 * asin(sqrt(a))
   r = 6371 # Earth radius in km
   return float(c * r)
# Convert numeric fields safely (PosTime, Lat, Long)
def To_numb(x):
   x['PosTime'] = int(x['PosTime']) # ensure integer
   x['Lat'] = float(x['Lat']) # ensure float
               = float(x['Long']) # ensure float
   x['Long']
   return x
# Start Spark
sc = pyspark.SparkContext() # create a SparkContext
spark = SparkSession.builder \
    .appName("flights") \
   .master("yarn") \
    .getOrCreate()
# Get references for your Dataproc bucket/project
bucket = sc._jsc.hadoopConfiguration().get('fs.gs.system.bucket')
project = sc._jsc.hadoopConfiguration().get('fs.gs.project.id')
```



```
# Create unique directory to avoid path conflicts
                 = datetime.datetime.now().strftime('%Y%m%d %H%M%S')
timestamp
input_directory =
    f'gs://{bucket}/hadoop/tmp/bigquerry/pyspark_input_{timestamp}'
output_directory = f'gs://{bucket}/pyspark_demo_output_{timestamp}'
# Clean existing path
input_path = sc._jvm.org.apache.hadoop.fs.Path(input_directory)
input_path.getFileSystem(sc._jsc.hadoopConfiguration()).delete(input_path,
    True)
# Set BigQuery config
conf = {
    "mapred.bq.project.id":
                                   project,
    "mapred.bq.gcs.bucket":
                                   bucket,
    "mapred.bq.temp.gcs.path":
                                   input_directory,
    # Update to match your actual table in BigQuery
    "mapred.bq.input.project.id": "cs512-projects",
    "mapred.bq.input.dataset.id": "aircraft_data",
    "mapred.bq.input.table.id":
                                  "plane_data",
}
# Read from BigQuery using the connector
table_data = sc.newAPIHadoopRDD(
    "com.google.cloud.hadoop.io.bigquery.JsonTextBigQueryInputFormat",
    "org.apache.hadoop.io.LongWritable",
    "com.google.gson.JsonObject",
    conf=conf
# Convert JSON lines to Python dict
vals = table_data.values().map(json.loads)
vals = vals.filter(lambda row: "PosTime" in row and "Lat" in row and
    "Long" in row)
# Convert data to numeric
vals = vals.map(To_numb)
# Define schema for DataFrame
schema = StructType([
    StructField("Icao", StringType(), True),
```

```
Oregon State
University
```

```
FloatType(), True),
    StructField("Lat",
                         FloatType(), True),
    StructField("Long",
    StructField("PosTime", LongType(),
                                         True),
])
# Create DataFrame
df1 = spark.createDataFrame(vals, schema=schema)
# Repartition if desired
df1 = df1.repartition(6)
# Define a window to get the next Lat/Long in time order
window = Window.partitionBy("Icao").orderBy("PosTime").rowsBetween(1, 1)
       = df1.withColumn("Lat2", lead("Lat").over(window))
df1
      = df1.withColumn("Long2", lead("Long").over(window))
df1
# Remove rows with null Lat2/Long2
df1
      = df1.na.drop()
# UDF for distance
haver_udf = udf(haversine, FloatType())
# Calculate distance for each consecutive segment
df1 = df1.withColumn("dist", haver_udf(col("Long"), col("Lat"),
    col("Long2"), col("Lat2")))
# Create a temporary view for SQL queries
df1.createOrReplaceTempView("planes")
top10 = spark.sql("""
    SELECT Icao, SUM(dist) AS dist
   FROM planes
   GROUP BY Icao
   ORDER BY dist DESC
   LIMIT 10
""")
print("Top 10 planes by distance:")
print(top10.collect())
miles = spark.sql("SELECT SUM(dist) FROM planes")
```

```
Oregon State
University
```