

# Spark Assignment - Week 8

Nick Sun

This Spark assignment was pretty rough. The way that I ended up doing it was using **PySpark** all the way through without exporting any data into BigQuery to tackle in SQL.

The key **Dataprep** step was rerunning the recipe from last week, but in addition to ICAO, LAT, and LONG, we also need PosTime which is the epoch time of each ping received. I then loaded this new table into BigQuery so that I could access it from the Spark compute cluster. I used this information later to calculate the total distance travelled for each ICAO.

The relevant **PySpark** code is included below. The basic logic is that I have a function written which can calculate the Haversine distance for an entire flight:

```
def getHaversineDistance(pings):
    distance = 0
    for first_ping, second_ping in zip(pings, pings[1:]):
        distance += haversine(
            float(first_ping[1]['Long']), float(first_ping[1]['Lat']),
            float(second_ping[1]['Long']), float(second_ping[1]['Lat'])
        )
    return distance
```

I can then use this function on each flight and then sum over all flights using **PySpark**.

```
vals = table_data.values()
vals = vals.map(lambda line: json.loads(line))

key_pings = vals.map(lambda x: (x['Icao'], x))
sorted_key_pings = key_pings.sortBy(lambda x: x[1]['PosTime'])

grouped_key_pings = sorted_key_pings.groupBy(lambda x: x[0])
grouped_key_pings_lst = grouped_key_pings.map(lambda (x, y): (x, list(y)))

distances = grouped_key_pings_lst.map(lambda (x,y): (x, getHaversineDistance(y)))

totalDistance = distances.map(lambda (x,y): y).reduce(add)
```

My final answer was 170081905.217 km. Below are relevant screenshots:

Equivalent **command line**

cs512-week7 ▼



[← Job details](#)

 REFRESH CLONE

✔ job-68204b3d

Start time: Feb 26, 2020, 9:27:36 PM Elapsed time: 10 min 26 sec Status:

Output Configuration

☐ Line wrapping

This account is managed by [oregonstate.edu](#).  
[Learn more](#)



Nicholas Sun  
sunn@oregonstate.edu  
[Privacy](#)

Google Account



Nicholas Sun  
njnicksun@gmail.com

Add account

Sign out

$$(1 + 1) / 2]$$

```
'The total distance flown is: 170081905.217 km.'
```