# Spark Plane Distances Part 1

Brian Cervantes Alvarez

February 26, 2025

Working on this assignment was both challenging and satisfying. First, I set up a DataProc cluster on Google Cloud and configured my PySpark code so it could read data from BigQuery. I ran into a big problem called a "Conflict occurred creating export directory," which happened because an old folder in Google Cloud Storage was still there. This made Spark fail when trying to export data from BigQuery. After searching online and testing different fixes, I solved the issue by deleting the old folder before I tried to export again. That way, the folder was empty and ready to be used. Also, I started out trying to use the **FSeen** column in the plane data, but realized it wasn't needed for my main task. By removing it, my code ran smoothly without KeyError messages.

Once everything was set up, the DataProc job started up its worker nodes, accessed my BigQuery data, and showed me the first few records so I knew it worked. If you look at the screenshot, you'll see proof of my PySpark script running successfully on DataProc—no more errors this time. Along the way, I got valuable practice using Google Cloud Platform services—DataProc, GCS, and BigQuery—and I learned how important it is to carefully manage export folders so things don't crash or get stuck.