# Project Statement for Milestone 1

## Group 6

Group member: Jinyang Ruan, Rusu Wu, Brian Chan, Junqiao Mou, Yi Yao

## Part 1: Project description (4%).

**1. (1%) Problem statement: Answer the following questions:**
**a. Give a brief description of the candidate project and its tasks.**

As a startup company, we need to develop a program to aid client analytics and visualize the massive scale of the datasets. The functions of the program include but are not limited to:

Airport and airline search:
- Find list of airports operating in the Country X
- Find the list of Airlines having X stops
- List of airlines operating with code share
- Find the list of active airlines in the United States - Airline aggregation:
- Which country (or) territory has the highest number of Airports.
- Top K cities with most incoming/outgoing airlines

Geospatial analytics: You must use Apache Sedona (https://sedona.apache.org/)
- Find the closest airport to a city X's geospatial coordinate
- Find the airport in each US state's geospatial boundary

**b. Explain why do you want to choose this candidate project?**

We chose the Airline Search Engine project by voting. 4 of us voted this topic because we think it is doable and relatively interesting.

**c. Your group members and a description about the roles of each member (consider the four basic roles listed in Project Overview).**

Jinyang Ruan: Infrastructure manager, visualization expert
Rusu Wu: Data analyst, ETL (Extract-Transform-Load) programmer
Brian Chan: Visualization expert, Data analyst
Junqiao Mou: ETL (Extract-Transform-Load) programmer, visualization expert
Yi Yao: Data analyst, Infrastructure manager

**2. (1%) Datasets:**
**a. Give the link and description of the dataset.**

Airports, Counties(2.66MB): http://openflights.org/data.html

Airport Codes(35MB): https://datahub.io/core/airport-codes

Airport ID, Location, Region(13.8MB): https://ourairports.com/data/

Routes(2.38MB): https://www.kaggle.com/open-flights/flight-route-database

Total: 53.84MB

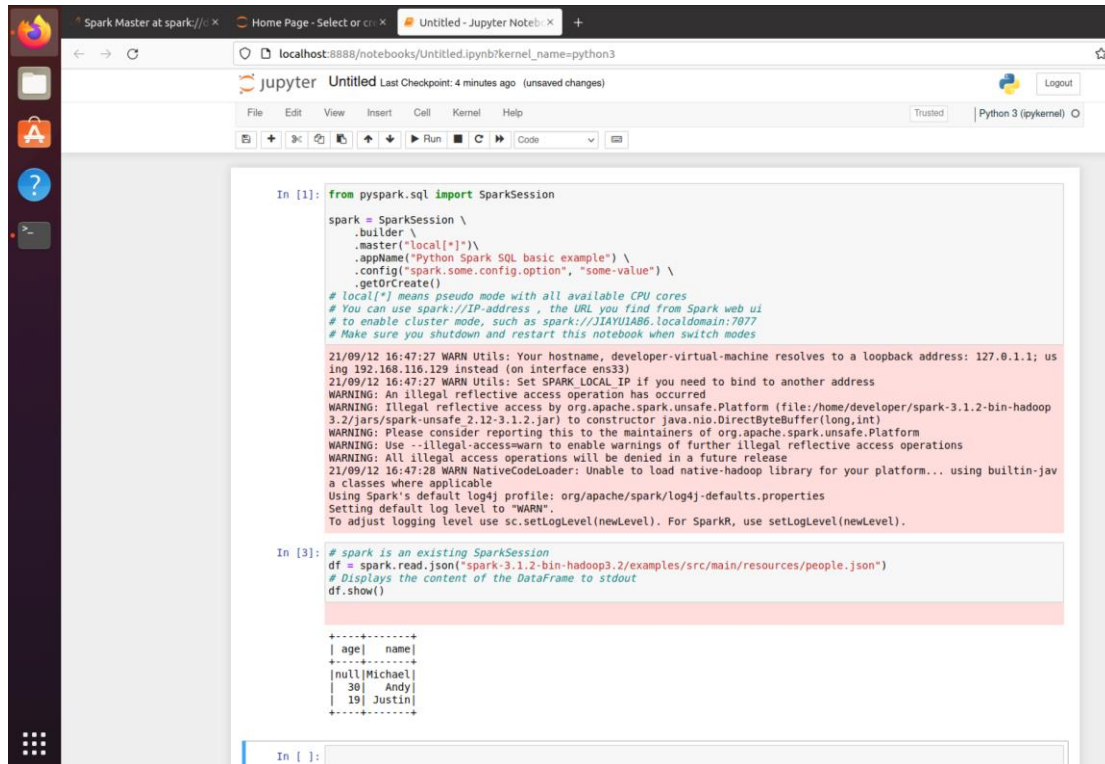**b. Assume we use relational models, what are the data tables you plan to create based on the input datasets? At this stage, a rough estimation is fine.**

Airport ID, Name of the airport, Airline, Airline ID, Main city served by airport, Country or territory where airport is located, local code, continent, Code of Airport, Hours offset from UTC, Timezone, etc.

**3. (1%) A timetable for your Milestones according to the roles of each group members. At this stage, a rough estimation is fine.**

| Phase | Beginning Date | Deadline |
|---|---|---|
| *establish the project environment* | 20th September | 26th September |
| *prepare the large-scale datasets* | 20th September | 10th October |
| *data analyst* | 4th October | 17th October |
| *visualization expert* | 18th October | 7th November |
| *programming* | 18th October | 7th November |
| *debugging* | 8th November | 21th November |
| *finishing* | 22th November | 28th November |

## Part 2: Set up PySpark environment (7%)

• Screenshot 1 (3%): A Jupyter notebook in a web browser that runs all commands (with outputs) listed in the given example notebook

• Screenshot 2 (4%): Spark web UI in a web browser that shows your Jupyter notebook application connects to the cluster. The cluster only needs to have 1 worker which is the master machine itself.

**Spark Master at spark://developer-virtual-machine:7077**

**URL:** spark://developer-virtual-machine:7077
**Alive Workers:** 1
**Cores in use:** 4 Total, 4 Used
**Memory in use:** 6.7 GiB Total, 1024.0 MiB Used
**Resources in use:**
**Applications:** 1 Running, 0 Completed
**Drivers:** 0 Running, 0 Completed
**Status:** ALIVE

**▾ Workers (1)**

| Worker Id | Address | State | Cores | Memory | Resources |
|-----------|---------|-------|-------|--------|-----------|
| worker-20210909080134-192.168.116.129-36839 | 192.168.116.129:36839 | ALIVE | 4 (4 Used) | 6.7 GiB (1024.0 MiB Used) | |

**▾ Running Applications (1)**

| Application ID | | Name | Cores | Memory per Executor | Resources Per Executor | Submitted Time | User | State | Duration |
|----------------|--|------|-------|---------------------|------------------------|----------------|------|-------|----------|
| app-20210912165945-0000 | (kill) | Python Spark SQL basic example | 4 | 1024.0 MiB | | 2021/09/12 16:59:45 | developer | RUNNING | 22 s |

**▾ Completed Applications (0)**

| Application ID | Name | Cores | Memory per Executor | Resources Per Executor | Submitted Time | User | State | Duration |
|----------------|------|-------|---------------------|------------------------|----------------|------|-------|----------|