

Peer-Graded Assignment: Analyzing Big Data with SQL

Name: Jessica Chen

Date: SEP 27 2020

(Include your name and today's date above.)

Assignment

Recommend which pair of United States airports should be connected with a high-speed passenger rail tunnel. To do this, write and run a SELECT statement to return pairs of airports that are between **300** and **400** miles apart and that had at least **5,000** (five thousand) flights per year on average *in each direction* between them. Arrange the rows to identify which one of these pairs of airports has largest total number of seats on the planes that flew between them. Your SELECT statement must return all the information required to fill in the table below.

Recommendation

I recommend the following tunnel route:

	First Direction	Second Direction
Three-letter airport code for origin	SFO	LAX
Three-letter airport code for destination	LAX	SFO
Average flight distance in miles	337	337
Average number of flights per year	14712	14540
Average annual passenger capacity	1996597	1981059
Average arrival delay in minutes	10	14

(Replace AAA and BBB with the actual airport codes, and fill in all the cells of the table.)

Method

I identified this route by running the following SELECT statement using Impala on the VM:

```
Select
    f.origin as origin,
    f.dest as destination,
    avg(f.distance) as avg_dist,
    round(count(f.flight)/10) as avg_annual_num_flights,
    round(sum(p.seats)/10) as avg_annual_seat_capacity,
    round(avg(arr_delay)) as avg_delay
from    flights f
left outer join planes p
on f.tailnum=p.tailnum
where f.distance >=300 and f.distance <=400 -- two airports between 300 and 400 miles apart
group by origin, destination
having avg_annual_num_flights >5000
order by avg_annual_seat_capacity desc
limit 10;
```

(Fill in the blank to indicate whether you used Hive or Impala, and fill in the SQL query.)

Notes

(This section is optional. You may use it to describe your process, add details or caveats, explain your interpretations, or describe any further analysis that you performed.)