

Background

The airline data that we've used for learning how to use map-reduce includes a large number of fields for things like day of the week (`DayOfWeek`), origin airport (`Origin`), and departure delay (`DepDelay`). The following problems requires you to process the airport data to answer questions about these times.

For each question, use the following data sets and approach:

Dataset: <https://cse.sc.edu/~rose/590B/airline/2008.csv>

Method 1: map-reduce

Method 2: Pig Latin

Method 3: HiveQL

If you are using `wget` to transfer the files, the command is of the form:

`wget https://cse.sc.edu/~rose/590B/airline/2008.csv`

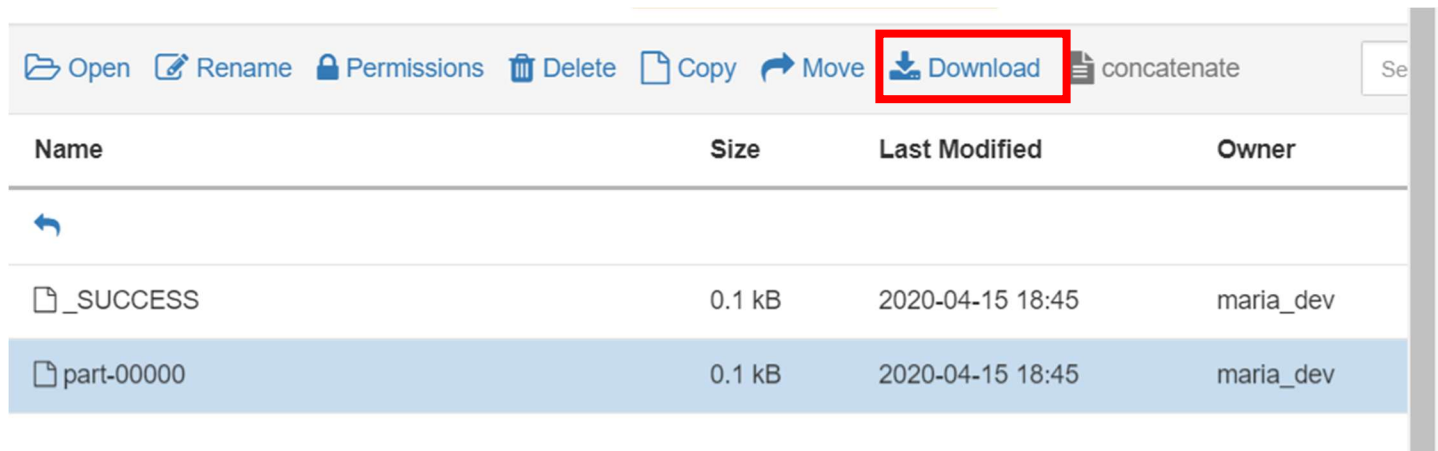
My suggestion is that you develop your code (map-reduce, pig, and hive) using the smallest data set of 12 records (`testNA.csv`), then try the data set containing ~25K flight records (`test_25K.csv`). When you get that to work, then move on to the large file from 2008.

Transferring the input data set to HDFS

The data set 2008.csv is too large to use the drag-and-drop approach with Ambari. You will have to use `wget` to transfer it to the linux filesystem of your vm and then “`hadoop fs -put`” it to HDFS.

Saving output from HDFS to your local machine:

- 1) Login to Ambari (`vm-hadoop-xx.cse.sc.edu:8080`)
- 2) click on the table icon (tick-tack-toe) on the menu bar
- 3) select “Files View” from the table drop down menu
- 4) navigate to user/`maria_dev` then to the directory containing your output
- 5) click on the output file to select it.
- 6) You should then see a menu bar with file operations in blue lettering. The 7th operation is “Download”. See the red box in the figure below:
- 7) Click on the download icon. This will download the file to the “Downloads” folder on your local machine.



I. MAPREDUCE Problem 1 Find the average departure delay time by day of week using map-reduce: aveDepDelayByDayOfWeek

This problem is an organic extension of our understanding of map-reduce from counting words, to counting flights, to calculating average arrival delay to now finding the average departure delay times by day of week. Our hypothesis is that some days of the week exhibit larger departure delays. As a first step we need to calculate the average departure delay to see if this hypothesis might have merit. Start by extracting the `DayOfWeek` and `DepDelay` columns. The goal of this problem is find the average departure delay time for each day of the week. Since the mapper and reducer communicate via key-value tuples, it would be easiest to select `DayOfWeek` and `DepDelay` in the mapper and output the key-value pair as `DayOfWeek`, `DepDelay`. As in the case of our mapreduce lab, the reducer must handle missing values (NA). you saw how to do this in the reducer.) Thus, the reducer must handle missing values (NA) and simply calculate the average of the `DepDelay` that it receives for each day of the week.

Review your results and write a brief report describing how your results either support or do not support the hypothesis that some days of the week exhibit larger departure delays. Independent of whether your results support this hypothesis or not, explain how you could definitively test this hypothesis **using ideas that you should have learned this semester. (No credit for methods NOT discussed in this course.)** Be sure to provide enough details as to how you would test this hypothesis.

Submit the python code for your mapper, reducer, and your output results.

- 1) Name your mapper file `aveDepDelayByDayOfWeekMapper.py`
- 2) Name your reducer file `aveDepDelayByDayOfWeekReducer.py`
- 3) Name your output file `aveDepDelayByDayOfWeekMapReduceResults`
- 4) Name your report `aveDepDelayByDayOfWeekReport`

Hints:

- 1) Read the file “TipsForMR.txt” This file is in the Class Materials folder on Teams. **Be smart! Test your mapper and reducer separately as explained in “TipsForMR.txt” BEFORE using hadoop in streaming mode.**
- 2) Start by refreshing your understanding of the `aveArrDelay` map-reduce example. The solution to that problem is very close to this problem.
- 3) Be sure to check for missing values (NA) in your reducer
- 4) Be sure to explicitly cast objects to be the type you want. If you want to interpret X as an integer, use `int(X)`. If you want to interpret the integer Y as a string use `str(Y)`.

Problem 2 Find the average departure delay time by day of week using Pig: aveDepDelayByDayOfWeek

Create a pig solution for finding the average departure delay by origin day of week. This is the same problem you solved in map-reduce, but now you are creating a solution in Pig Latin. When indexing columns, don't forget that column numbers start from 0 in Pig Latin.

Submit your Pig Latin script and your output results.

- 1) Name your Pig Latin script aveDepDelayByDayOfWeek.pig
- 2) Name your output file aveDepDelayByDayOfWeekPigResults

Hints:

- 1) Start by refreshing your understanding of the aveArrDelay Pig Latin example. The solution to that problem is very close to this problem.
- 2) In Pig Latin, missing values are referred to as nulls. In the case of arithmetic, if either subexpression is null, the resulting expression is null. Since the AVG() function handles null values, you don't need worry missing values for this problem.

Problem 3 Find the average departure delay time by day of week using Hive:**aveDepDelayByDayOfWeek**

Create a hive solution for finding the average departure delay by origin day of week. This is the same problem you solved in map-reduce, but now you are creating a solution in HiveQL. When indexing columns, don't forget that column numbers start from 1 in HiveQL.

Submit your HiveQL script and output results.

- 1) Name your HiveQL script aveDepDelayByDayOfWeek.hql
- 2) Name your output file aveDepDelayByDayOfWeekHiveResults

Hints:

- 1) Start by refreshing your understanding of aveArrDelay HiveQL. The solution to that problem is very close to this problem.
- 2) Be sure to check for missing values (NA) for DepDelay when using a SELECT statement