

CSE 414 Homework 7: Parallel Data Processing and Spark

Objectives: To write distributed queries. To learn about Spark and running distributed data processing in the cloud using AWS.

Assigned date: Saturday, Dec 5, 2020.

Due date: Sunday, Dec 13, 2020, at 11 pm.

What to turn in:

Your Spark code in the SparkApp.java file with createLocalSession() enabled.

Resources:

- [Spark programming guide](#)
- [Spark Javadoc](#)
- [Amazon web services EMR \(Elastic MapReduce\) documentation](#)
- [Amazon S3 documentation](#)
- [Starter code files](#)

Spark Programming Assignment (75 points)

In this homework, you will be writing Spark and Spark SQL code, to be executed both locally on your machine and also using Amazon Web Services.

We will be using a similar flight dataset used in previous homework. This time, however, we will be using the *entire* data dump from the [US Bureau of Transportation Statistics](#), which consists of information about all domestic US flights from 1987 to 2011 or so. The data is in [Parquet](#) format. Your local runs/tests will use a subset of the data (in the flights_small directory) and your cloud jobs will use the full data (stored on Amazon [S3](#)).

Your main task is to complete the Spark programming locally. Once your code works locally, you can follow our instructions to get it running on AWS, but this is optional.

Here is a rough outline on how to do the HW:

A. **(Optional)** Sign up for AWS and apply for credit for AWS

B. Complete the HW locally

C. **(Optional)** Run your solutions on AWS Elastic MapReduce one at a time when you are fairly confident with your solutions

A. (Optional) Sign up on Amazon Web Services

Follow these steps to set up your Amazon Web Services account.

1. If you do not already have an Amazon Web Services account, go to [AWS Educate](#) and sign up with your @uw email.
2. Click "Join AWS Educate", and choose the student account. Then complete the necessary forms, and wait until your account is approved. **This may take a while, so do this early.**
3. To get \$\$\$ to use Amazon AWS, sign in to your AWS Educate account, click "AWS Account", and create a starter account. **This will only work if you have not applied for a starter account in the past.**
4. Follow the steps to finalize your credits. If everything works, you should have \$100 credit in your account.

IMPORTANT: If you run AWS in any other way rather than how we instruct you to do so below, you must remember to manually terminate the AWS clusters when you are done. While the credit that you receive should be more than enough for this homework assignment, you will be responsible for paying the extra bill should your credits be exhausted.

Now you are ready to run applications using the Amazon cloud. But before you do that let's write some code and run it locally.

B. Get Code Working Locally

We have created empty method bodies for each of the questions below (QA, QB, and QC). *Do not change any of the method signatures.* You are free to define extra methods and classes if you need to. We have also provided a warmup method that shows fully-functional examples of three ways that the same query could be solved using Spark's different APIs.

There are many ways to write the code for this assignment. Here are some documentation links that we think would get you started up about what is available in the Spark functional APIs:

- [Spark 2.4.5 Manual](#)
- [Spark 2.4.5 Javadocs](#)
- [Dataset](#)
- [Row](#) (see also RowFactory)
- [JavaRDD](#) (see also JavaPairRDD)

- [Tuple2](#)

The quickstart documentation also has more depth and examples of using [RDDs](#) and [Datasets](#).

For questions a, b, and c, you will get the points for writing a correct query.

(a) (15 points) Complete the method QA in SparkApp.java. Use the Spark functional APIs or SparkSQL. Select all flights that leave from 'Seattle, WA', and return the destination city names. Only return each destination city name once. Return the results in an RDD where the Row is a single column for the destination city name.

(b) (30 points) Complete the method QB in SparkApp.java. Only use the Spark functional APIs. Find the number of non-canceled ($\neq 1$) flights per month-origin city pair. Return the results in an RDD where the row has three columns that are the origin city name, month, and count, in that order.

(c) (30 points) Complete the method QC in SparkApp.java. Only use the Spark functional APIs. Compute the average delay from all departing flights for each city. Flights with NULL delay values should not be counted, and canceled flights should not be counted. Return the results in an RDD where the row has two columns that are the origin city name and average, in that order.

Testing Locally

We provide cardinality testing when you run

```
$ mvn test
```

You are responsible for verifying you have the correct format and contents in your results.

Running Local Jobs

To actually execute the main method, toggle the SparkSession initialization on lines 147 and 148 of SparkApp.java to allow it to run locally (local SparkSession, not cluster). Run from the top-level directory (with pom.xml in it):

```
$ mvn clean compile assembly:single
$ java -jar target/sparkapp-1.0-jar-with-dependencies.jar \
    flights_small output
```

Note that, on Windows, this should be executed in the root directory of this repo, so that the program can find bin/winutils.exe. (The code uses the directory of execution .; you can change

this in the createLocalSession() method if you must.) For reference, winutils.exe was obtained for this Hadoop version from a [Github repo](#).

For this quarter, we added code to add compatibility with Java 9+. It should work fine. In case there is a problem, you can force a Java 8 execution by downloading a Java 8 JRE and setting your JAVA_HOME variable to your Java 8 runtime.

C. (Optional) Run Code on Elastic Map Reduce (EMR)

Run your jobs on Elastic Map Reduce (EMR) as described below, and copy the resulting output from EMR to QA.txt, QB.txt, and QC.txt, respectively.

Running all jobs at the same time with the provided configuration took less than 30 min for the solutions.

We will use Amazon's [Elastic Map Reduce](#) (EMR) to deploy our code on AWS. Follow these steps to do so after you have set up your account, received credits as mentioned above, and have tested your solution locally. **Read this carefully!**

1. Toggle the SparkSession initialization of SparkApp.java to allow it to run on AWS (cluster SparkSession; comment out local). Then create a jar file from the top level directory that packages everything needed to run the Spark application. The following command creates the jar file in the targets folder:

```
$ mvn clean compile assembly:single
```

2. Login to [S3](#) and create a "bucket". S3 is Amazon's cloud storage service, and a bucket is similar to a folder. Give your bucket a meaningful name and select **US East (N. Virginia)**, and leave the other settings as default. Upload the jar file that you created in Step 1 to that bucket by selecting that file from your local drive and click "Upload" once you have selected the file.
3. Login to [EMR](#). Make sure you select US East (N. Virginia) on the upper right. **This is the only region supported by the starter account.**
4. We will first configure cluster software. Click on the "Create Cluster" link, then click "Go to advanced options" in the Amazon EMR console.
 - **Check the boxes for Spark and Hadoop** so that your screen looks like this:

Software Configuration

Release: emr-5.29.0

<input checked="" type="checkbox"/> Hadoop 2.8.5	<input type="checkbox"/> Zeppelin 0.8.2	<input type="checkbox"/> Livy 0.6.0
<input type="checkbox"/> JupyterHub 1.0.0	<input type="checkbox"/> Tez 0.9.2	<input type="checkbox"/> Flink 1.9.1
<input type="checkbox"/> Ganglia 3.7.2	<input type="checkbox"/> HBase 1.4.10	<input type="checkbox"/> Pig 0.17.0
<input type="checkbox"/> Hive 2.3.6	<input type="checkbox"/> Presto 0.227	<input type="checkbox"/> ZooKeeper 3.4.14
<input type="checkbox"/> MXNet 1.5.1	<input type="checkbox"/> Sqoop 1.4.7	<input type="checkbox"/> Mahout 0.13.0
<input type="checkbox"/> Hue 4.4.0	<input type="checkbox"/> Phoenix 4.14.3	<input type="checkbox"/> Oozie 5.1.0
<input checked="" type="checkbox"/> Spark 2.4.4	<input type="checkbox"/> HCatalog 2.3.6	<input type="checkbox"/> TensorFlow 1.14.0

Multiple master nodes (optional)

☐ Use multiple master nodes to improve cluster availability. [Learn more](#)

AWS Glue Data Catalog settings (optional)

☐ Use for Spark table metadata

Edit software settings

☒ Enter configuration ☐ Load JSON from S3

```
classification=config-file-name,properties=[myKey1=myValue1,myKey2=myValue2]
```

Steps (optional)

A step is a unit of work you submit to the cluster. For instance, a step might contain one or more Hadoop or Spark jobs. You can also submit additional steps to a cluster after it is running. [Learn more](#)

Concurrency: ☐ Run multiple steps at the same time to improve cluster utilization

After last step completes: ☐ Clusters enters waiting state

☒ Cluster auto-terminates

- Next, scroll to the Steps section at the bottom of the page and **create a Spark application step**. A "step" is a single job to be executed. You could specify multiple Spark jobs to be executed one after another in a cluster. Fill out the Spark application step details by filling in the boxes so that your screen looks like this (with SparkApp instead of HW6):

Add step

Step type

Spark application

Run Spark application using spark-submit. [Learn more](#)

Name

Spark application

Deploy mode

Cluster

Run your driver on a slave node (cluster mode) or on the master node as an external client (client mode).

Spark-submit options

--class edu.uw.cs.SparkApp

Specify other options for spark-submit.

Application location*

s3://<your bucket containing the jar file>/HW6.jar

Path to a JAR with your application and dependencies (client deploy mode only supports a local path).

Arguments

s3://us-east-1.elasticmapreduce.samples/flightdata/input

s3://<your output bucket>

Specify optional arguments for your application.

Action on failure

Terminate cluster

What happens if the step fails

Cancel

Add

The `--class` option under "Spark-submit options" tells Spark where your main method lives. (`--class edu.uw.cs.SparkApp`)

The "Application location" should just point to where your uploaded jar file is. You can use the folder button to navigate.

The full flights data location is the first argument:
`s3://us-east-1.elasticmapreduce.samples/flightdata/input`

The output destination is the second argument. You can use the bucket that holds your jar. You can modify the "output" folder name prefix to be something different if you like.

Make sure you fill out the correct bucket names. There are two arguments listed (and separated by white space, as if you were running the program locally):

Change "Action on failure" to "Terminate cluster" (or else you will need to terminate the cluster manually).

- **Click Add.**
- Back to the main screen, now on the After last step completes: option at the bottom of the page, select Cluster auto-terminates so the cluster will shut down once your Spark application is finished. Click Next.
- On the next screen, we will now configure the hardware for a five-node EMR cluster to execute the code. We recommend using the "m4.large" "instance type", which is analogous to some set of allocated resources on a server (in AWS terminology, "m" stands for high memory, "4" represents the generation of servers, and "large" is the relative size of allocated resources). You get to choose how many machines you want in your cluster. For this assignment 1 master instance and 4 core (i.e., worker) instances of m4.large should be good. You are free to add more or pick other types, but make sure you think about the price tag first... Grabbing 100 machines at once will probably drain your credit in a snap :(If m4.large is not available, choose another instance with a similar name (m4.xlarge, m5.large, etc.). **Click Next.**

Network: [Create a VPC](#)

EC2 Subnet:

Root device EBS volume size: GiB

Choose the instance type, number of instances, and a purchasing option. You can choose to use On-Demand Instances, Spot Instances, or both. The instance type and purchasing option apply to all EC2 instances in each instance group, and you can only specify these options for an instance group when you create it. [Learn more about instance purchasing options](#)

Node type	Instance type	Instance count	Purchasing option	Auto Scaling
Master Master - 1	m4.large 4 vCore, 8 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	1 Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot <input type="text" value="Use on-demand as max price"/>	Not available for Master
Core Core - 2	m4.large 4 vCore, 8 GiB memory, EBS only storage EBS Storage: 32 GiB Add configuration settings	<input type="text" value="4"/> Instances	<input checked="" type="radio"/> On-demand <input type="radio"/> Spot <input type="text" value="Use on-demand as max price"/>	Not enabled

[+ Add task instance group](#)

[Cancel](#) [Previous](#) [Next](#)

- Under “General Options” *uncheck the “Termination protection” option*. We recommend that you allow the default logging information in case you need to debug a failure. **Click Next.**
- Click *Create cluster* once you are done and your cluster will start spinning up!

It will take a bit for AWS to both provision the machines and run your Spark job. As a reference, it took about 10 mins to run the warmup job on EMR. You can monitor the status of your cluster on the EMR homepage.

To rerun a similar job (maybe you want to try a different jar), use the "Clone" cluster button to copy the settings into a new job when you run your actual HW problems.

Make sure your cluster is terminated! It should do so if you selected the options above. You should check this each time you look at the HW, just to make sure you don't get charged for leaving a cluster running. It's fine if you see warning (or even occasional error) messages in the logs. If your EMR job finishes successfully, you should see something similar to the below in the main EMR console screen:

j-2BIXJO3MBV2B7	Terminated All steps completed	2017-11-08 17:51 (UTC-8)	17 minutes	24
-----------------	-----------------------------------	--------------------------	------------	----

Debugging AWS jobs

Debugging AWS jobs is not easy for beginners. Besides making sure your program works locally before running on AWS, here are some general tips:

- Make sure that you set ALL the job details (i.e., options, arguments, bucket names, etc) correctly!
- Make sure you switched the two lines of SparkSession code mentioned to run your job on AWS instead of locally.

- Make sure you freshly compile your solution and replace your jar to test a new version of code!
- **99% of cluster failures or errors are due to the first three points!**
- The easiest way to debug is to look at the output/logging files. Spark generates a lot of log files, the most useful ones are probably the stderr.gz files listed under containers/application.../container/stderr.gz. You will have one container folder per machine instance. So make sure you check all folders. You should also check out the log files that you have specified when you created the job in Step 8 above. You can also see the names of those files listed as "Log File" under "Steps":

Cluster: My cluster Terminated Steps completed with errors

Summary Application history Monitoring Hardware Events Steps Configurations Bootstrap actions

Add step Clone step Cancel step No job available

Steps

Filter: All steps Filter steps ... 2 steps (all loaded)

ID	Name	Status	Start time (UTC-7)	Elapsed time	Log files
s-WRF8SWXWIS96	Spark application	Failed	2018-04-28 15:14 (UTC-7)	1 minute	

Status : FAILED
Reason :
Log File : s3://aws-logs-488156842632-us-east-1/elasticmapreduce/j-3NW5TJH8IMCL/steps/s-WRF8SWXWIS96/stderr.gz
Details : Exception in thread "main" org.apache.spark.SparkException: Application application_1524953570520_0001 finished with failed status
JAR location : command-runner.jar

- It is rare that your HW solution is fine but the cluster fails. This is usually due to AWS not being able to grab your machines due to the demand for the instance type saturating the supply available. If you can't find available instances in a region, try changing to a different **EC2 subnet**, like so:

Hardware Configuration

If you need more than 20 EC2 instances, [see this topic](#).

Instance group configuration ☒ **Uniform instance groups**
Specify a single instance type and purchasing option for each node type.

☐ **Instance fleets**
Specify target capacity and how Amazon EMR fulfills it for each node type. Mix instance types and purchasing options. [Learn more](#)

Network vpc-869e18fe (172.31.0.0/16) (default) [Create a VPC](#)

EC2 Subnet subnet-76d06149 | Default in us-east-1e

Root device EBS volume size 10 GiB

- Spark has a web UI that you can set up to check on job progress etc. You can check out [their webpage](#) for details. But these are more involved so you are probably better off to first try examining the logs. Specifically, try the "Application History" tab and the dropdown.

IMPORTANT: Cleanup after completing the HW

Double-check that the clusters you have created are all terminated.

S3 charges by [downloading/uploading data from/to the buckets](#). So once you are done with the assignment you might want to delete all the buckets that you have created (in addition to shutting down any EMR clusters that you have created).

The amount you are allocated from Amazon should be more than enough to complete the assignment. And every year we have students forgetting to shut down their cluster/clean up their buckets and that can result in substantial charges that they need to pay out of pocket. So be warned!

Submission Instructions

Turn in your SparkApp.java by submitting it to Gradescope.

Note: if you changed your code to run on AWS with the `createClusterSession()` enabled, make sure to change it back to `createLocalSession()` before submitting