

CS236

Database Management System

Final Project Report

Sakib Md Bin Malek
862004540

The goal of the project was to leverage spark's parallel ability to divide up the SQL queries to consolidate the precipitation data of the 50 states over the real data collected for about 4 years from various weather stations.

Script to run the program

I used python as the programming language to complete the project. I also used pyspark package and used hadoop (version 2.4.5)

```
#Setting up environment
```

```
$ pip3 install pyspark
```

```
$ wget https://downloads.apache.org/spark/spark-2.4.5/  
spark-2.4.5-bin-hadoop2.7.tgz
```

```
$ tar xvf spark-2.4.5-bin-hadoop2.7.tgz
```

```
#After properly setting up the spark and hadoop
```

```
#Both input and output folder path are optional.
```

```
$ spark-submit project.py -o output_folder_path -i input_folder_path
```

```
$ less <output_folder_path>/output.txt
```

Project Approach

I used the following steps to complete the project.

- Read weather station metadata
- Read weather station data
- Remove invalid row and remove redundant columns

- Run queries to get aggregated data grouped by state and month
- Create a new table using the aggregated data in the earlier step
- Run queries on new table to get Max, Min and Difference data

Results

Following is the snapshot of the consolidated grouped by State and Month.

```

output.txt
~/Academics/CS236/Project/CSE236/Output

Time take to read data: 193.897000074 secs

[(u'PA', 12, 0.11152692569870487), (u'LA', 12, 0.12237618871186857), (u'CO', 4, 0.030123210123210117),
(u'AL', 11, 0.10613475177304964), (u'VA', 2, 0.028723879082978288), (u'SC', 8, 0.14045536357091432), (u'UT',
4, 0.0313863004740886), (u'KY', 9, 0.13152872444011682), (u'IA', 2, 0.023119444444444442), (u'WV', 8,
0.10961102904972919), (u'AK', 12, 0.053286931329945234), (u'ID', 10, 0.046768845189897816), (u'NC', 7,
0.10464413145539904), (u'DE', 7, 0.0707798960138648), (u'GA', 8, 0.11706457546713135), (u'NH', 8,
0.1294584031267448), (u'ME', 11, 0.13843258042436682), (u'ND', 7, 0.08175451620885921), (u'HI', 6,
0.02864836100468283), (u'AR', 2, 0.09512895535327265), (u'NC', 3, 0.05507599655864638), (u'NJ', 11,
0.10527668845315902), (u'RI', 6, 0.06500149476831091), (u'AZ', 2, 0.018955171196611366), (u'IA', 9,
0.060688394760889434), (u'AK', 6, 0.04086302936402538), (u'MI', 11, 0.035182254196642684), (u'MD', 3,
0.03289814066289409), (u'NH', 12, 0.10877350776778412), (u'WA', 11, 0.18171547768043392), (u'AR', 6,
0.13829818829818827), (u'NaN', 7, 0.0015092865300234408), (u'TX', 11, 0.0269833060368359), (u'TN', 1,
0.13709698137443802), (u'FL', 12, 0.06467918209060464), (u'MT', 4, 0.051100787198348166), (u'NE', 9,
0.07398487066861882), (u'PR', 7, 0.18088541666666667), (u'OH', 8, 0.11637287790993187), (u'OK', 2,
0.022275366086127714), (u'WY', 1, 0.01723443223443224), (u'RI', 4, 0.07020758122743681), (u'KY', 10,
0.16832728372655775), (u'OH', 9, 0.08832031656773923), (u'IA', 7, 0.06918501048218029), (u'AZ', 1,
0.01745526625950927), (u'FL', 5, 0.10000887223974761), (u'OK', 12, 0.039555010511562724), (u'VT', 12,
0.12191142191142189), (u'TN', 4, 0.16476049868766401), (u'PA', 6, 0.15678846827900453), (u'NaN', 5,
0.0035490223660149104), (u'NC', 5, 0.07308534821244059), (u'NM', 9, 0.050567535190176705), (u'DE', 12,
0.06238546603475512), (u'PR', 11, 0.15143589743589742), (u'VA', 7, 0.07606354334948842), (u'MO', 5,
0.1751697883356175), (u'IL', 5, 0.08553565081396174), (u'MO', 7, 0.16298430286241922), (u'TN', 9,
0.11095299145299145), (u'VA', 6, 0.09173496718896479), (u'PR', 5, 0.6332307692307693), (u'WY', 12,
0.02028676021964612), (u'AZ', 11, 0.006272269821841424), (u'KS', 11, 0.0235119825708061), (u'AL', 2,
0.148646481404673), (u'TX', 1, 0.03375084090144635), (u'WY', 6, 0.04584437338206183), (u'NJ', 1,
0.09548258706467659), (u'WV', 3, 0.09613403880070544), (u'SD', 4, 0.07332629504504502), (u'MN', 3,

```

Following is the states sorted by the difference in maximum and minimum precipitation based on month in ascending order.

```

Row(STATE='DC', DIFF_PRCP=0.0, MAX_PRCP=0.0, MIN_PRCP=0.0)
Row(STATE='NaN', DIFF_PRCP=0.002393434166724912, MAX_PRCP=0.0035490223660149104,
MIN_PRCP=0.0011555881992899985)
Row(STATE='NV', DIFF_PRCP=0.016430704822198952, MAX_PRCP=0.025384193194291978, MIN_PRCP=0.008953488372093026)
Row(STATE='RI', DIFF_PRCP=0.02807150010215058, MAX_PRCP=0.07020758122743681, MIN_PRCP=0.04213608112528623)
Row(STATE='UT', DIFF_PRCP=0.028588202826524446, MAX_PRCP=0.05096310193871169, MIN_PRCP=0.022374899112187242)
Row(STATE='MI', DIFF_PRCP=0.0395590884907908, MAX_PRCP=0.07095705521472392, MIN_PRCP=0.031397956365644845)
Row(STATE='ID', DIFF_PRCP=0.045042798976555776, MAX_PRCP=0.05972575115950796, MIN_PRCP=0.014682952182952184)
Row(STATE='IL', DIFF_PRCP=0.05178858515945441, MAX_PRCP=0.08596879788271349, MIN_PRCP=0.03418021272325907)
Row(STATE='WY', DIFF_PRCP=0.05182458498220462, MAX_PRCP=0.06414096916299558, MIN_PRCP=0.01231638418079096)
Row(STATE='MN', DIFF_PRCP=0.0537049132599366, MAX_PRCP=0.05808902657753138, MIN_PRCP=0.004384113317594779)
Row(STATE='AK', DIFF_PRCP=0.05499419962607067, MAX_PRCP=0.08550112826753889, MIN_PRCP=0.03050692864146822)
Row(STATE='MD', DIFF_PRCP=0.0550405872492519, MAX_PRCP=0.08597362296353761, MIN_PRCP=0.03093303571428571)
Row(STATE='WI', DIFF_PRCP=0.056384101983851034, MAX_PRCP=0.07517805839513596, MIN_PRCP=0.018793956411284927)
Row(STATE='IN', DIFF_PRCP=0.05903946510110894, MAX_PRCP=0.12301587301587301, MIN_PRCP=0.06397640791476407)
Row(STATE='NC', DIFF_PRCP=0.06100515622699178, MAX_PRCP=0.10464413145539904, MIN_PRCP=0.04363897522840726)
Row(STATE='CO', DIFF_PRCP=0.06136640976688339, MAX_PRCP=0.07250497434286311, MIN_PRCP=0.011138564575979724)
Row(STATE='AZ', DIFF_PRCP=0.062294371064985074, MAX_PRCP=0.0685666408868265, MIN_PRCP=0.006272269821841424)

```

Runtime

Total runtime of my program is **212.44 sec**. However, it only took **18.54 sec** to run both of the queries and create intermediate table. Most of the time (**193.90 sec**) is spent reading the data to the memory and clean the data.

Comment

There are number of issues that I faced while completing the project.

One of the issues that I encountered was that my program was very slow. I was looping through all the state and making queries to the SQL with a state filter. This made my project extremely slow, taking more than 10 mins. This issue was easily resolved by adding a **Group By** in the SQL query. That brought down the runtime to less than 4 mins, about **3X faster**.

Another issue was due to the configuration of the hadoop. I was getting **OutOfMemoryError Exception in Java** from hadoop. This was due to the fact that default configuration does not allocate enough memory for spark units. I used Python's garbage collector to removed already used table from the memory. That did not resolve the issue, thus I had to make changes to the config file. It solved the problem.