

Connor Brown

CS 348 - Homework 7: Data Warehousing (HIVE and BigQuery).
(40 Points)

Fall 2020

Due on: **12/1/2020 at 11:59 pm**

This assignment is to be completed by individuals. You should only talk to the instructor, and the TAs about this assignment. You may also post questions (and not answers) to Campuswire.

There will be a 10% penalty if the homework is submitted 24 hours after the due date, a 15% penalty if the homework is submitted 48 hours after the due date, or a 20% penalty if the homework is submitted 72 hours after the due date. The homework will not be accepted after 72 hours, as a solution will be posted by then.

Submission Instruction: Write your answers in this document. Create a pdf and upload to Gradescope.

For questions 1 and 2, consider the `us_cities.csv` file included with this homework (this file does not include column headers). `us_cities_sample.csv` includes the column headers and only a few rows. The create-table statement can be found in `createTable.sql`.

For more information about HIVE and BigQuery, including the supported SQL clauses, you can visit the HIVE and BigQuery documentation.

<https://cwiki.apache.org/confluence/display/Hive/LanguageManual>

<https://cloud.google.com/bigquery/docs/reference/standard-sql/query-syntax>

Question 1) (10 points)

Using HIVE, list each state along with its population. The population of a state is the total population of its cities. Sort your result by population in a descending order.

Write your answer (query) and insert a screenshot of your result (partial result is sufficient). **Your screenshot should show the HIVE command prompt.**

Expected Result:

state_name	totalPop
California	59608877
New York	37847876
Texas	33020438
Florida	30243610
Illinois	18266059
Pennsylvania	15641201
Ohio	13764673
Georgia	10750970
Michigan	10429997
Washington	10134988
Arizona	9581518
North Carolina	9203228
Massachusetts	8972816
Virginia	8197045
Colorado	7499152
Missouri	7463356
Minnesota	7275566
Maryland	6880970
Indiana	6161106

.....
OUTPUT is TRIMMED

```
SELECT state_name, SUM(population) AS totalPop
FROM cities2
GROUP BY state_name
ORDER BY totalPop DESC;
```

```
0: jdbc:hive2://localhost:10000/default> SELECT state_name, SUM(population) AS totalPop
. . . . .> FROM cities2
. . . . .> GROUP BY state_name
. . . . .> ORDER BY totalPop DESC;
```

state_name	totalpop
California	59608877
New York	37847876
Texas	33020438
Florida	30243610
Illinois	18266059
Pennsylvania	15641201
Ohio	13764673
Georgia	10750970
Michigan	10429997
Washington	10134988
Arizona	9581518
North Carolina	9203228
Massachusetts	8972816
Virginia	8197045
Colorado	7499152
Missouri	7463356
Minnesota	7275566
Maryland	6880970
Indiana	6161106
Tennessee	6065103
Wisconsin	5936025
New Jersey	5697038
District of Columbia	5289420
Oregon	5212252
Louisiana	4894945
Utah	4854868
South Carolina	4530198
Nevada	4434080
Alabama	4299420
Connecticut	4072787
Oklahoma	3657380
Kentucky	3287922
Iowa	3167941
Arkansas	2608972
Kansas	2586907
New Mexico	2102495
Mississippi	2076480
Puerto Rico	1940742
Nebraska	1885071
Hawaii	1839050
Rhode Island	1653103
Idaho	1632041
West Virginia	1375886
Montana	816367
New Hampshire	799613
North Dakota	727244
Alaska	695791
South Dakota	678953
Maine	616831
Delaware	557808
Wyoming	487683
Vermont	258863
FL	NULL
state_name	NULL

```
54 rows selected (10.97 seconds)
```

Question 2, HIVE) (15 points)

The explode command in HIVE can be used to flatten complex data types, such as maps and arrays. See the following example of Explode in the HIVE documentation:

<https://cwiki.apache.org/confluence/display/Hive/LanguageManual+LateralView>

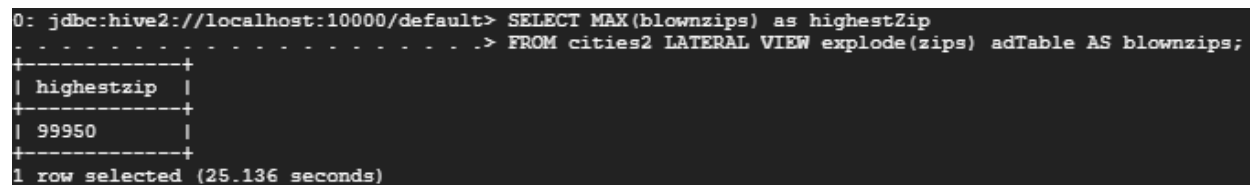
Using explode to flatten the zips column in the cities data set, find the highest zip code in the US. Write your query and insert a screenshot of your result. **Your screenshot should show the HIVE command prompt.**

Expected Result:

```
+-----+
| highestzip |
+-----+
| 99950      |
+-----+
```

SELECT MAX(blownzips) as highestZip

FROM cities2 LATERAL VIEW explode(zips) adTable AS blownzips;



```
0: jdbc:hive2://localhost:10000/default> SELECT MAX(blownzips) as highestZip
- - - - -> FROM cities2 LATERAL VIEW explode(zips) adTable AS blownzips;
+-----+
| highestzip |
+-----+
| 99950      |
+-----+
1 row selected (25.136 seconds)
```

Question 3) (5 points)

What is the MongoDB operator that is similar to Explode in HIVE?

The mongoDB function similar to explode in Hive is the \$unwind operator.

Question 3, BigQuery) (10 points)

Choose a public data set in BigQuery and write two queries. Insert screenshots of your queries and results (partial results are sufficient).

Google Cloud PlatformCS348-HW7Search products and resources

BigQueryFEATURES & INFOSHORTCUT

Query history

Saved queries

Job history

Transfers

Scheduled queries

Reservations

BI Engine

Resources

+ ADD DATA

Search for your tables and datasets

san_francisco_sffd_service_ca...

san_francisco_sfpd_incidents

san_francisco_transit_muni

san_francisco_trees

street_trees

sdoh_bea_cainc30

sdoh_cdc_wonder_natality

sdoh_cms_dual_eligible_enroll...

sdoh_hrsa_shortage_areas

Query editor

+ COMPOSE NEW QUERYHIDE EDITORFULL SCREEN

```
1 SELECT
2   COUNT(tree_id) count,
3   tree_id,
4   STRING_AGG(DISTINCT title) movies
5 FROM
6   `bigquery-public-data.san_francisco_trees.street_trees` t
7 JOIN
8   `bigquery-public-data.san_francisco_film_locations.film_locations` f
9 ON
10  f.locations = t.address
11 GROUP BY
12  tree_id, f.locations
13 ORDER BY
14  tree_id ASC
```

RunSave querySave viewSchedule queryMore

This query will process 4.7 MB when run.

Query results

SAVE RESULTSEXPLORE DATA

Query complete (0.7 sec elapsed, 4.7 MB processed)

Job informationResultsJSONExecution details

Row	count	tree_id	movies
1	1	29050	Bitter Melon
2	1	29051	Bitter Melon
3	2	31413	Quitters
4	2	31414	Quitters

Rows per page: 1001 - 46 of 46First pageLast page

Google Cloud PlatformCS348-HW7Search products and resources

BigQueryFEATURES & INFOSHORTCUT

Query history

Saved queries

Job history

Transfers

Scheduled queries

Reservations

BI Engine

Resources

+ ADD DATA

Search for your tables and datasets

san_francisco_bikeshare

san_francisco_film_locations

film_locations

san_francisco_sffd_service_ca...

san_francisco_sfpd_incidents

sfpd_incidents

san_francisco_transit_muni

san_francisco_trees

street_trees

sdoh_bea_cainc30

Query editor

+ COMPOSE NEW QUERYHIDE EDITORFULL SCREEN

```
1 SELECT
2   COUNT(locations) count,
3   title,
4 FROM
5   `bigquery-public-data.san_francisco_film_locations.film_locations` f
6 GROUP BY
7   title, f.locations
8 ORDER BY
9   count DESC
```

RunSave querySave viewSchedule queryMore

This query will process 162.4 KB when run.

Query results

SAVE RESULTSEXPLORE DATA

Query complete (0.4 sec elapsed, 162.4 KB processed)

Job informationResultsJSONExecution details

Row	count	title
1	6	Groove
2	4	A View to a Kill
3	4	San Andreas
4	2	A Jitney Elopement

Rows per page: 1001 - 100 of 1819First pageLast page