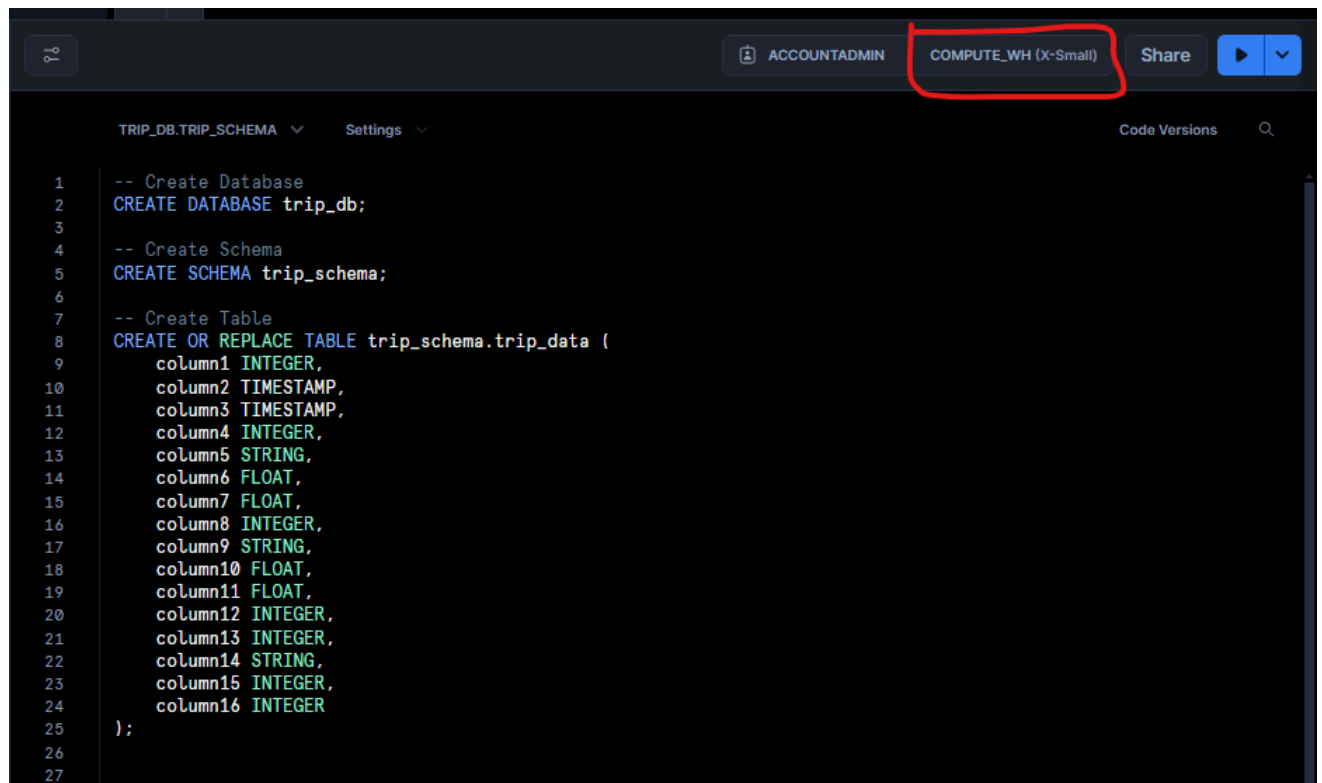


Homework 2

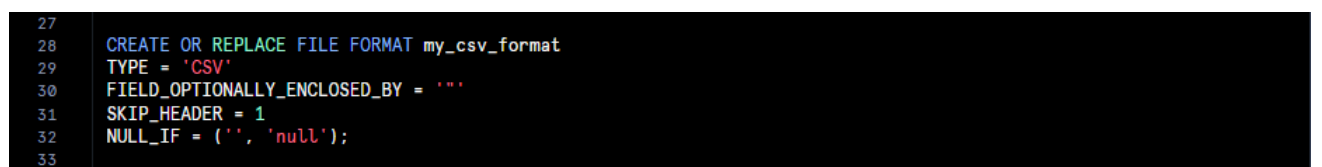
1. (+2) Create Database 'trip_db', Schema 'trip_schema' and Table 'trip_data' in Snowflake. Make sure the schema of trip_data should matches the structure of the file. Use the warehouse of size **XSMALL**



The screenshot shows the Snowflake SQL Editor interface. At the top, the user is logged in as 'ACCOUNTADMIN' and the current warehouse is 'COMPUTE_WH (X-Small)', which is highlighted with a red rectangle. The editor displays the following SQL code:

```
1  -- Create Database
2  CREATE DATABASE trip_db;
3
4  -- Create Schema
5  CREATE SCHEMA trip_schema;
6
7  -- Create Table
8  CREATE OR REPLACE TABLE trip_schema.trip_data (
9      column1 INTEGER,
10     column2 TIMESTAMP,
11     column3 TIMESTAMP,
12     column4 INTEGER,
13     column5 STRING,
14     column6 FLOAT,
15     column7 FLOAT,
16     column8 INTEGER,
17     column9 STRING,
18     column10 FLOAT,
19     column11 FLOAT,
20     column12 INTEGER,
21     column13 INTEGER,
22     column14 STRING,
23     column15 INTEGER,
24     column16 INTEGER
25 );
26
27
```

- Creating database trip_db and schema trip_schema is pretty straight forward, I took a slightly different approach inorder to match the schema to the given S3 database. First using the link I printed out the first 10 rows of the database in snowflake, and using the first row, I concurred that there are 16 columns with different data types which I matched with the columns while creating the table as seen in the screenshot.
2. (+1) Create Stage 'trip_stage' in Snowflake and it can be any stage type of your choice (internal or external)



The screenshot shows the Snowflake SQL Editor with the following SQL code:

```
27
28 CREATE OR REPLACE FILE FORMAT my_csv_format
29 TYPE = 'CSV'
30 FIELD_OPTIONALLY_ENCLOSED_BY = ''
31 SKIP_HEADER = 1
32 NULL_IF = ('', 'null');
33
```

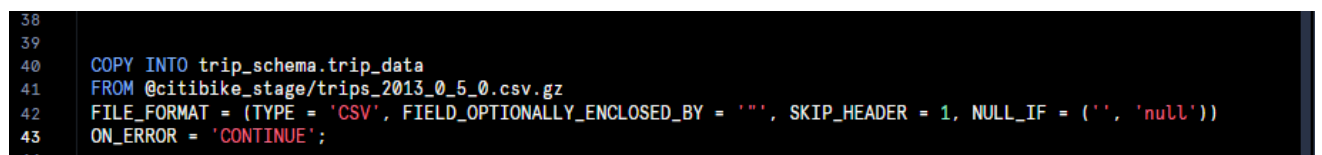
3. (+1) Create File Format 'trip_file_format' in Snowflake.



The screenshot shows the Snowflake SQL Editor with the following SQL code:

```
34
35 CREATE OR REPLACE STAGE citibike_stage
36 URL = 's3://snowflake-workshop-lab/citibike-trips-csv/'
37 FILE_FORMAT = my_csv_format;
```

4. (+1) Load Data into the 'trip_data' table (COPY INTO)



The screenshot shows the Snowflake SQL Editor with the following SQL code:

```
38
39
40 COPY INTO trip_schema.trip_data
41 FROM @citibike_stage/trips_2013_0_5_0.csv.gz
42 FILE_FORMAT = (TYPE = 'CSV', FIELD_OPTIONALLY_ENCLOSED_BY = '', SKIP_HEADER = 1, NULL_IF = ('', 'null'))
43 ON_ERROR = 'CONTINUE';
44
```

Results		Chart			
	COLUMN1	COLUMN2	COLUMN3	COLUMN4	COLUMN5
1	1717	2013-06-28 18:08:27.000	2013-06-28 18:37:04.000	349	Rivington St & Ridge St
2	838	2013-06-28 18:08:27.000	2013-06-28 18:22:25.000	2012	E 27 St & 1 Ave
3	972	2013-06-28 18:08:30.000	2013-06-28 18:24:42.000	358	Christopher St & Greenwich
4	602	2013-06-28 18:08:30.000	2013-06-28 18:18:32.000	540	Lexington Ave & E 29 St
5	542	2013-06-28 18:08:31.000	2013-06-28 18:17:33.000	521	8 Ave & W 31 St N
6	1344	2013-06-28 18:08:32.000	2013-06-28 18:30:56.000	477	W 41 St & 8 Ave
7	443	2013-06-28 18:08:34.000	2013-06-28 18:15:57.000	515	W 43 St & 10 Ave
8	903	2013-06-28 18:08:35.000	2013-06-28 18:23:38.000	400	Pitt St & Stanton St
9	525	2013-06-28 18:08:37.000	2013-06-28 18:17:22.000	484	W 44 St & 5 Ave
10	1176	2013-06-28 18:08:39.000	2013-06-28 18:28:15.000	454	E 51 St & 1 Ave

- As you can see the data has been loaded using COPY INTO

5. (+4) Write an SQL query to produce a report that shows, for each hour of the day, the following:
1. The total number of trips that started during that hour.
 2. The average duration of these trips in mins.
 3. The average distance traveled during these trips in kms.

Results		Chart			
	HOUR_OF_DAY	TOTAL_TRIPS	AVG_DURATION_MINS	AVG_DISTANCE_KMS	
1	0	1900	422.284211	40.731788075	<div>Query Details</div> <div>Query duration 428ms</div> <div>Rows 24</div> <div>Query ID 01b70192-0002-e7a5-0...</div> <div>Show more</div> <div>HOUR_OF_DAY</div> <div>TOTAL_TRIPS</div> <div>AVG_DURATION_MINS</div> <div>AVG_DISTANCE_KMS</div>
2	1	1050	423.681905	40.731909883	
3	2	699	405.340486	40.732597666	
4	3	453	419.456954	40.732713679	
5	4	297	414.909091	40.735083739	
6	5	470	426.714894	40.734007014	
7	6	1671	441.868941	40.73604748	
8	7	3507	435.335044	40.73630456	
9	8	5332	424.010690	40.733586983	
10	9	4066	419.421545	40.732889609	
11	10	3797	432.680801	40.733275058	
12	11	4502	431.010440	40.732353958	
13	12	5673	431.122158	40.732453917	
14	13	6014	431.758730	40.73171905	
15	14	6424	437.827833	40.731967256	
16	15	5655	430.698674	40.732154404	
17	16	6270	434.710686	40.733400151	
18	17	8861	434.075612	40.734096079	
19	18	11363	423.358004	40.733277689	
20	19	9291	425.566139	40.732600832	
21	20	6036	422.352386	40.731374105	
22	21	4573	421.676580	40.732084337	
23	22	3714	421.341949	40.732306824	
24	23	2738	416.083272	40.732084247	

6. (+1) Create a virtual warehouse of XLARGE.

```

57
58
59 CREATE OR REPLACE WAREHOUSE trip_warehouse
60 WITH
61     WAREHOUSE_SIZE = 'XLARGE'
62     AUTO_SUSPEND = 600
63     AUTO_RESUME = TRUE;
64
65
66 USE WAREHOUSE trip_warehouse;
67

```

- (+3) Rerun the same query from step 4 & 5 using the XLARGE warehouse **after dropping trip_data table**, Analyze the performance of the node upon changing configurations from x-small to x-large in Snowflake.

x-small:

```

100
101
102 COPY INTO trip_schema.trip_data
103 FROM @bulkdata-stage/trips_2013_0_5_0.csv.gz
104 FILE_FORMAT = (TYPE = 'CSV', FIELD_OPTIONALLY_ENCLOSED_BY = '"', SKIP_HEADER = 1, NULL_IF = ('', 'null'))
105 ON_ERROR = CONTINUE;
106
107
108 SELECT * FROM trip_schema.trip_data LIMIT 10;
109
110
111 SELECT
112     EXTRACT(HOUR FROM column2) AS hour_of_day,
113     COUNT(*) AS total_trips,
114     AVG(column4) AS avg_duration_mins,
115     AVG(column6) AS avg_distance_kms
116 FROM trip_schema.trip_data
117 GROUP BY EXTRACT(HOUR FROM column2)
118 ORDER BY hour_of_day;
119
120 -- Drop the trip_data table
121 DROP TABLE IF EXISTS trip_schema.trip_data;

```

file	status	rows_parsed	rows_loaded	error_limit	errors_seen	first_error	first_error_line	first_error_character	first_error_column_name
s3://snowflake-workshop-lab/citibike-trips-csv/trips_2013_0_5_0.csv.gz	LOADED	104356	104356	104356	0	null	null	null	null

Query Details
 Query duration: 2.5s
 Rows: 1
 Query ID: 03b70245-0002-479d-8...

File: 100% filled
 Status: 100% filled
 Ask Copilot

x-large:

```

100
101
102 COPY INTO trip_schema.trip_data
103 FROM @bulkdata-stage/trips_2013_0_5_0.csv.gz
104 FILE_FORMAT = (TYPE = 'CSV', FIELD_OPTIONALLY_ENCLOSED_BY = '"', SKIP_HEADER = 1, NULL_IF = ('', 'null'))
105 ON_ERROR = CONTINUE;
106
107
108 SELECT * FROM trip_schema.trip_data LIMIT 10;
109
110
111 SELECT
112     EXTRACT(HOUR FROM column2) AS hour_of_day,
113     COUNT(*) AS total_trips,
114     AVG(column4) AS avg_duration_mins,
115     AVG(column6) AS avg_distance_kms
116 FROM trip_schema.trip_data
117 GROUP BY EXTRACT(HOUR FROM column2)
118 ORDER BY hour_of_day;
119
120 -- Drop the trip_data table
121 DROP TABLE IF EXISTS trip_schema.trip_data;

```

file	status	rows_parsed	rows_loaded	error_limit	errors_seen	first_error	first_error_line	first_error_character	first_error_column_name
s3://snowflake-workshop-lab/citibike-trips-csv/trips_2013_0_5_0.csv.gz	LOADED	104356	104356	104356	0	null	null	null	null

Query Details
 Query duration: 2.5s
 Rows: 1
 Query ID: 03b70245-0002-479d-8...

File: 100% filled
 Status: 100% filled
 Ask Copilot

ACCOUNTADMIN TRIP_WAREHOUSE (X-Large) Share

TRIP_DB.TRIP_SCHEMA Settings Code Versions

```
35 CREATE OR REPLACE STAGE citibike_stage
36 URL = 's3://snowflake-workshop-lab/citibike-trips-csv/'
37 FILE_FORMAT = my_csv_format;
38
39
40 COPY INTO trip_schema.trip_data
41 FROM @citibike_stage/trips_2013_0_5_0.csv.gz
42 FILE_FORMAT = (TYPE = 'CSV', FIELD_OPTIONALLY_ENCLOSED_BY = '', SKIP_HEADER = 1, NULL_IF = ('', 'null'))
43 ON_ERROR = 'CONTINUE';
44
45
46 SELECT * FROM trip_schema.trip_data LIMIT 10;
47
48 SELECT
49     EXTRACT(HOUR FROM column2) AS hour_of_day,
50     COUNT(*) AS total_trips,
51     AVG(column4) AS avg_duration_mins,
52     AVG(column6) AS avg_distance_kms
53 FROM trip_schema.trip_data
54 GROUP BY EXTRACT(HOUR FROM column2)
55 ORDER BY hour_of_day;
56
57
58 CREATE OR REPLACE WAREHOUSE trip_warehouse
59 WITH
60     WAREHOUSE_SIZE = 'XLARGE'
61     AUTO_SUSPEND = 600
62     AUTO_RESUME = TRUE;
63
64
65 USE WAREHOUSE trip_warehouse;
66
67
68 SELECT
69     EXTRACT(HOUR FROM column2) AS hour_of_day,
70     COUNT(*) AS total_trips,
71     AVG(column4) AS avg_duration_mins,
72     AVG(column6) AS avg_distance_kms
73 FROM trip_schema.trip_data
74 GROUP BY EXTRACT(HOUR FROM column2)
75 ORDER BY hour_of_day;
```

Results Chart

	HOUR_OF_DAY	TOTAL_TRIPS	AVG_DURATION_MINS	AVG_DISTANCE_KMS
1	0	1900	422.284211	40.731788075
2	1	1050	423.681905	40.731909883
3	2	699	405.340486	40.732597666
4	3	453	419.456954	40.732713679
5	4	297	414.909091	40.735083739
6	5	470	426.714894	40.734007014
7	6	1671	441.868941	40.73604748
8	7	3507	435.335044	40.73630456
9	8	5332	424.010690	40.733586983
10	9	4066	419.421545	40.732889609
11	10	3797	432.680801	40.733275058

Query Details

Query duration 81ms

Rows 24

Query ID 01b701c2-0002-e7a9-0...

Show more

HOUR_OF_DAY

TOTAL_TRIPS

Ask Copilot

We can see that x-small took 214ms to load the data into the schema whereas after dropping the table and loading the data in x-large took half as much time.

We can see that on the x-small warehouse the query took 428ms to run whereas in the x-large warehouse the same query took just 81ms which is 5x faster.