**TEAM – 4**

**ADTA 5240 PROJECT**

**DOCUMENT WITH SCREENSHOTS**

**PRESENTED BY:**

Abhiram Karthik Meduri

Abhiram Nallamothu

Pavan Kola

Jacob Larkin

**Two Datasets:**

**Static dataset: https://data.cityofdenton.com/dataset/traffic-closed-cases/resource/b12cba00-24e4-4f45-89f8-ffd1c0bf6b95**

**Dynamic dataset: https://data.cityofdenton.com/dataset/denton-crime-data**

**WORKING WITH DYNAMIC DATASET**

**1. Fetching Data from the City of Denton API**

**Type of API Request**: HTTP GET Request

**Explanation**:

So, what we're doing here is pulling some data from the City of Denton Open Data portal using a public API. Instead of grabbing everything in one go (which would be a lot and could easily overload the system), we're using SQL queries to get smaller, more manageable chunks. The API lets us use SQL right at the endpoint (https://data.cityofdenton.com/api/3/action/datastore_search_sql), which makes interacting with the dataset super straightforward.

Steps:

1. Build the SQL Query:

   o First, we put together an SQL query that pulls a set number of rows—10,000 at a time—and uses an OFFSET to make sure we're paging through the data step-by-step.

Sql:

SELECT * FROM "RESOURCE_ID" LIMIT 10000 OFFSET

   o Here, RESOURCE_ID represents the unique ID for the dataset we're working with in the City of Denton's data portal.

   o The LIMIT tells it to grab just 10,000 rows in each request. The OFFSET helps in fetching these rows in batches, so we don't end up asking for too much data all at once.
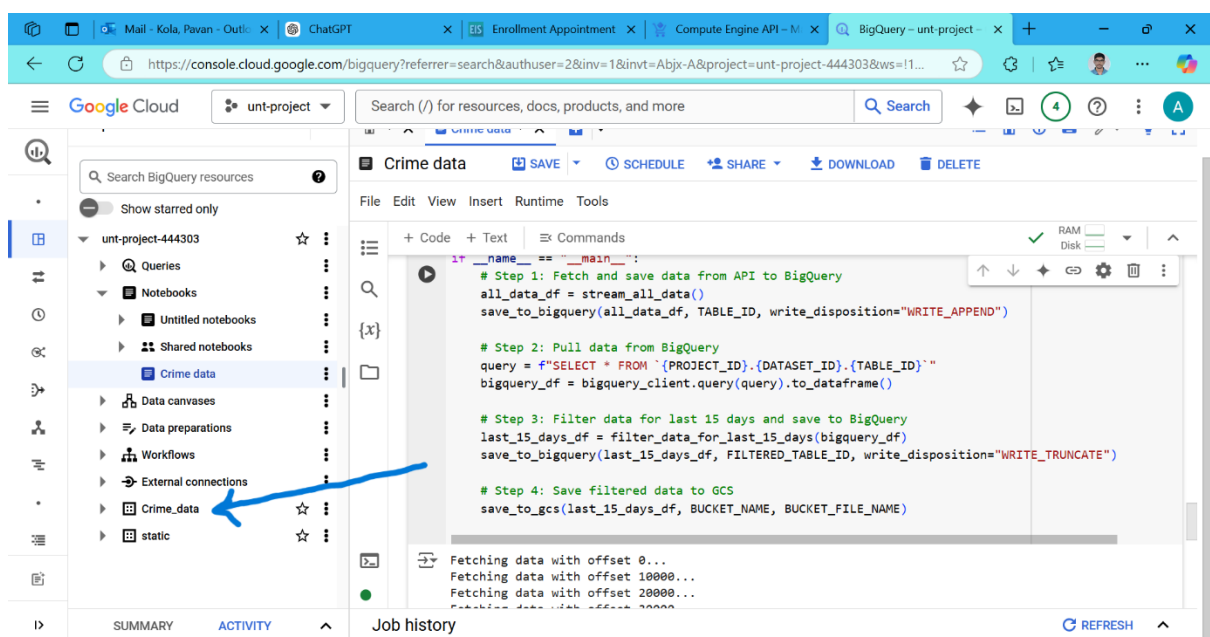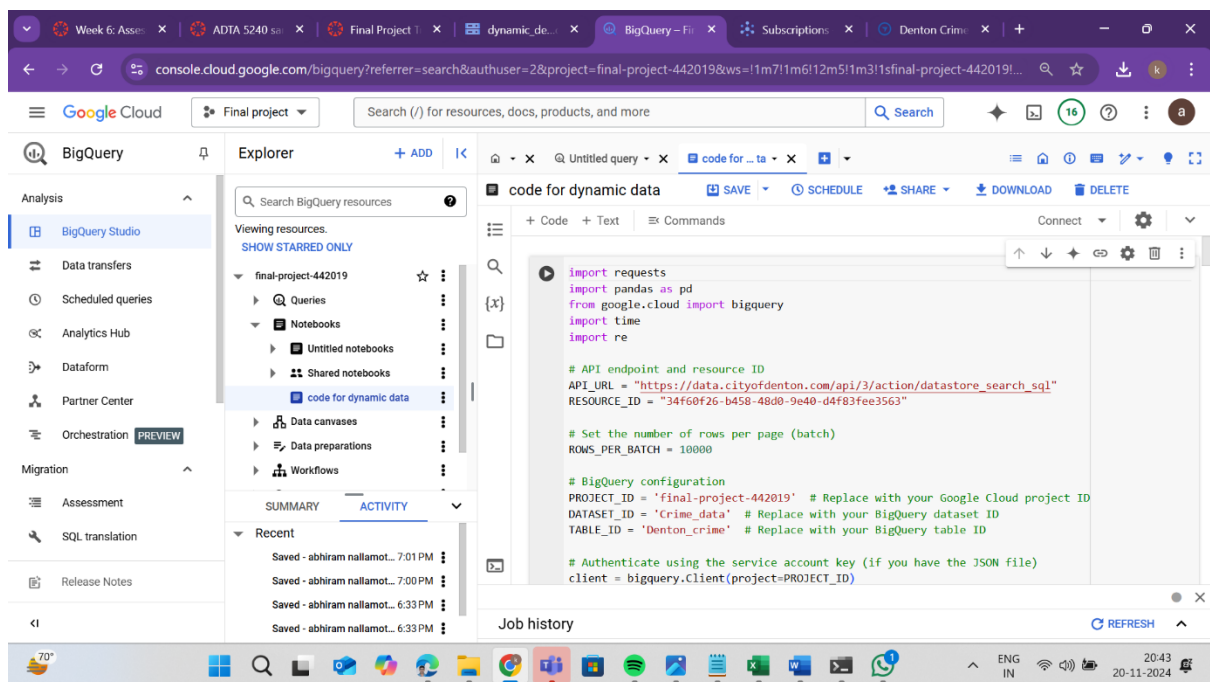
2. Send the GET Request:

   o After getting the SQL query ready, we need to send it to the API. We're doing this using Python's requests.get() function:

   o Here, API_URL is the link to the API, and params includes the SQL query (params = {'sql': SQL_QUERY}) we're sending.

   o This function sends the request to the API and waits for it to respond.

3. Check the Response:

o Once we get a response, we need to make sure everything went okay.

o A status code of 200 means it worked! If that's the case, we convert the response from JSON into a format that Python can work with easily.

o The records key is where we find all the rows of data we just pulled.

Outcome: By the end of these steps, we've got the data in a JSON format that contains all the crime records. Now it's ready for us to do something useful with it, like analysis or processing.

## DYNAMIC DATA PULLING IN BIGQUERY

**2. Pagination (Handling Large Data)**

**Type of Request**: GET Request with Pagination

**Explanation**:

When we're dealing with a lot of data, it's important not to try and grab everything at once. If we do, it could easily overwhelm the server or even cause our system to run out of memory. Instead, we use pagination—basically breaking the data into smaller chunks called pages. This way, we can pull the data piece by piece and make sure everything runs smoothly.

To do this, we use an OFFSET in our SQL query, which tells us where to start getting data from in the dataset.

Steps:

1. Set Initial Offset:

   o We start with an offset of 0, which means we're going to start pulling records from the beginning of the dataset.

2. Fetch Data in Batches:

   o To get the data in chunks, we use a loop (while True) to keep fetching 10,000 rows at a time.

3. Update Offset for Next Batch:

   o After pulling one batch, we update the offset by adding 10,000 (the number of rows in each batch) so that the next time around, we get the next set of rows.

4. Break on Last Page:

   o If the number of records we get is less than the batch size (i.e., fewer than 10,000 rows), that means we've reached the end, so we stop the loop.

5. Sleep to Simulate Streaming:

   o We add a time.sleep(2) call, which pauses the program for a couple of seconds between batches. This helps simulate a real-time streaming scenario where there's a small delay between each fetch.

Outcome: Using pagination like this, we can pull the data in manageable chunks without running into memory issues. This makes the whole process a lot more efficient and prevents overloading the system.

## 3. Cleaning Column Names

**Type of Data Transformation**: Data Normalization (Column Name Sanitization)

**Explanation**:

When we get data from the API, the column names might have spaces or special characters that don't work well in BigQuery. BigQuery has strict rules for column names, so we need to clean them up to make sure they're valid. To do that, we use a function called clean_column_names that replaces any invalid characters with an underscore (_).

Outcome: With the column names sanitized, the dataset becomes fully compatible with BigQuery and can be uploaded without issues. This step ensures that our data is clean, standardized, and ready for further processing.



## 4. Saving Data to BigQuery

**Type of Request**: BigQuery Load Job API (using load_table_from_dataframe)

**Explanation**:

Once we have the data cleaned and ready, the next step is to upload it to BigQuery. To do this, we use the BigQuery API, specifically the load_table_from_dataframe() function, which lets us take a pandas DataFrame and load it straight into a BigQuery table.

Outcome: Once everything is complete, the data is in BigQuery and ready to be queried or analyzed further. This makes it really easy to run all kinds of SQL queries on the data, creating reports or doing in-depth analysis.



## 5. Filtering Data by Date

**Type of Data Transformation**: Time-Based Filtering

**Explanation**:

Here, we're working on filtering the data to focus on records from a specific date. This means we're going to use the date_time column in our dataset and convert it to a proper datetime object. Once that's done, we can filter it down to match a particular date that we care about.

Steps:

1. Convert date_time to Datetime:

   o The first thing we need to do is take the date_time column and convert it from a string into an actual datetime object that Python can work with.

2. Filter by Date:

- o Now that we have proper datetime objects, we can filter the data to keep only the rows where the date matches the specific date we're interested in.

Outcome: After filtering, we end up with a dataset that only contains records from the specified date. This makes it much more manageable and focuses our analysis on the relevant time period.



## 6. Saving Filtered Data to BigQuery

**Type of Request**: BigQuery Load Job API (using load_table_from_dataframe)

**Explanation**:

This step is pretty much like the previous upload to BigQuery, but instead of the entire dataset, we're focusing on just the filtered data this time. The idea is to take that smaller, more relevant subset and store it in a new table in BigQuery.

Steps:

1. Clean Column Names:

- o Before uploading, we go ahead and clean the column names again to make sure they're valid for BigQuery. This just means getting rid of any spaces or special characters that BigQuery doesn't like.

2. Load Filtered Data:

- o Once the column names are good to go, we use load_table_from_dataframe() to upload the filtered data to BigQuery, just like we did before.

Outcome: By the end of this step, the filtered dataset is stored in BigQuery. Now it's ready for us to run queries on or do some more analysis—focusing just on the specific slice of data we're interested in.



## 7. Uploading Filtered Data to Google Cloud Storage (GCS)

**Type of API Request**: GCS Blob Upload API

**Explanation**:

Once we have the filtered data, the next step is to save it as a CSV file and upload it to Google Cloud Storage (GCS). This step makes it easy to share the data with others or use it for further processing in different tools or applications. It's basically about having a convenient backup or making the data accessible outside of BigQuery.

Steps:

1. Save Filtered Data as CSV:

o The filtered data is saved locally as a CSV file, which is a simple, widely-used format that most tools can understand.

2. Upload to GCS:

o After saving it, the CSV is uploaded to a GCS bucket. This means we can easily access it later, share it, or use it in other applications for more analysis.

Outcome: Now the filtered data is safely stored in Google Cloud Storage, ready for sharing or further use. It gives us flexibility—whether it's for archiving purposes, sharing with team members, or using it in other workflows.

# WORKING WITH STATIC DATASET

## Dataset: Traffic closed cases in denton-septmber 2024

## STEP 1:



Here, we cleaned the datset in openrefine.

Step 2:

Now, we uploaded this datset in the bucket that we created in the google cloud platform.

## Step 3:

We created a cluster and then uploaded the bucket and then opened SSH window from VM instances.

**Queries in Hive (**Dynamic dataset)

**A total of 342 rows of data is pulled from the streaming data (Denton crime cases) from 25 November.**

1. Created a table

Query 1:

Crime count (count of different crimes in the city)



```
5 rows selected (13.619 seconds)
0: jdbc:hive2://localhost:10000> SELECT crime, COUNT(*) AS crime_count
. . . . . . . . . . . . . . .> FROM denton_crime
. . . . . . . . . . . . . . .> GROUP BY crime
. . . . . . . . . . . . . . .> ORDER BY crime_count DESC
. . . . . . . . . . . . . . .> LIMIT 10;
INFO  : Compiling command(queryId=hive_20241210173544_7798cb70-6507-406e-b964-fb3e76a3261d): SELECT crime, COUNT(*) AS crime_count
FROM denton_crime
GROUP BY crime
ORDER BY crime_count DESC
LIMIT 10
INFO  : Concurrency mode is disabled, not creating a lock manager
INFO  : Semantic Analysis Completed (retrial = false)
INFO  : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:crime, type:string, comment:null), FieldSchema(name:crime_count, type:bigint, comment:null)], properties:null)
INFO  : Completed compiling command(queryId=hive_20241210173544_7798cb70-6507-406e-b964-fb3e76a3261d); Time taken: 0.661 seconds
INFO  : Concurrency mode is disabled, not creating a lock manager
INFO  : Executing command(queryId=hive_20241210173544_7798cb70-6507-406e-b964-fb3e76a3261d): SELECT crime, COUNT(*) AS crime_count
FROM denton_crime
GROUP BY crime
ORDER BY crime_count DESC
LIMIT 10
INFO  : Query ID = hive_20241210173544_7798cb70-6507-406e-b964-fb3e76a3261d
INFO  : Total jobs = 1
INFO  : Launching Job 1 out of 1
INFO  : Starting task [Stage-1:MAPRED] in serial mode
INFO  : Subscribed to counters: [] for queryId: hive_20241210173544_7798cb70-6507-406e-b964-fb3e76a3261d
INFO  : Session is already open
INFO  : Dag name: SELECT crime, COUNT(*) AS crime_count
F...10 (Stage-1)
INFO  : Status: Running (Executing on YARN cluster with App id application_1733849209580_0002)

----------------------------------------------------------------------------------------------
      VERTICES      MODE        STATUS  TOTAL  COMPLETED  RUNNING  PENDING  FAILED  KILLED
----------------------------------------------------------------------------------------------
Map 1 .......... container     SUCCEEDED     1         1        0        0       0       0
Reducer 2 ...... container     SUCCEEDED     1         1        0        0       0       0
Reducer 3 ...... container     SUCCEEDED     1         1        0        0       0       0
----------------------------------------------------------------------------------------------
VERTICES: 03/03  [==========================>>] 100%  ELAPSED TIME: 5.27 s
----------------------------------------------------------------------------------------------
INFO  : Completed executing command(queryId=hive_20241210173544_7798cb70-6507-406e-b964-fb3e76a3261d); Time taken: 6.221 seconds
INFO  : OK
INFO  : Concurrency mode is disabled, not creating a lock manager
+---------------------+--------------+
|        crime        | crime_count  |
+---------------------+--------------+
| All Other Offenses  | 75           |
| Simple Assault      | 51           |
| Other Larceny       | 46           |
| Drive Under Infl    | 22           |
| Theft From Veh      | 22           |
| Drug/Narc Violation | 17           |
| Trespassing         | 15           |
| Burglary/B&E        | 13           |
| Intimidation        | 12           |
| Vandalism           | 12           |
+---------------------+--------------+
10 rows selected (6.936 seconds)
```

These are the top 10 crime incidents, in the Denton city: simple assaults is the highest with 51 total cases, Driver under infl cases are about 46.

Query 2:

Crime count by location name

There are a total of 116 crimes that took place at home/residence, followed by the incidents that took place at highway/road/alley/street/sidewalk : 99 total crimes, 43 crimes at parking spots.

Query 3:

At which address most of the crimes took place



Most of the cases took place at the south loop, followed by east mckinnet streer and west university drive.

West university drive has also have significant number of crime cases in the span of 12 days.

Query 4:

Crime count by date



The data is from 25 November 2024 when we look at this table, we can say that most of the cases took place on December 1$^{st}$ with a count of 35, then on December 7 there were 34 total cases.

## Queries in Spark ( Static dataset):



## Query 1:

Address where highest cases were recorded:



## Query 2:

City where highest cases were recorded:

```
                                    row(s)
spark-sql> SELECT nam_r_state AS state, COUNT(*) AS total_cases
        > FROM traffic_closed_cases_2024
        > GROUP BY nam_r_state
        > ORDER BY total_cases DESC
        > LIMIT 5;
DENTON                                          223
DALLAS                                          30
LEWISVILLE                                      23
LITTLE ELM                                      13
FORT WORTH                                      13
```

Query 3:

Date when highest cases were recorded:

```
Time taken: 2.11 seconds, Fetched 63 row(s)
spark-sql> SELECT cod_desc1 AS violation_description, COUNT(*) AS total_cases
        > FROM traffic_closed_cases_2024
        > GROUP BY cod_desc1
        > ORDER BY total_cases DESC
        > LIMIT 5;
09/26/2024 06:15AM      20
09/16/2024 05:34AM      20
09/04/2024 04:56AM      19
09/12/2024 04:46AM      14
09/10/2024 04:56AM      14
```

Query 4:

Reason for violation case

SSH-in-browser                                ⬆ UPLOAD FILE    ⬇ DOWNLOAD FILE   ▣  ▦  ⚙

```
Time taken: 0.93 seconds, Fetched 10 row(s)
spark-sql> SELECT plea, COUNT(*) AS total_cases
        > FROM traffic_closed_cases_2024
        > GROUP BY plea
        > ORDER BY total_cases DESC;
SPEEDING                                    91
DL NO DRIVER'S LICENSE                       86
FAIL TO MAINTAIN FINANCIAL RESPONSIBILIT     68
REGISTRATION EXPIRED REGISTRATION            38
DL DRIVING WHILE LICENSE INVALID             14
RAN RED LIGHT                                13
SPEED FAIL TO CONTROL SPEED                  13
DL EXPIRED OPERATORS LICENSE                 7
DISREGARD OFFICIAL TRAFFIC CONTROL DEVIC     7
CHANGED LANE WHEN UNSAFE                      6
SPEEDING IN A SCHOOL ZONE                     6
SPEEDING LESS THAN 10 % ABOVE POSTED SPE     5
RAN STOP SIGN                                5
WINDOW- UNAUTHORIZED GLASS TINT COATING       5
FOLLOWING TOO CLOSELY                         5
FAIL TO YIELD ROW TURNING LEFT               4
FAIL TO YIELD ROW ENTERING HIGHWAY FROM      4
DL FAIL TO DISPLAY DRIVER'S LICENSE          4
DROVE ONTO (FROM) CONTROLLED ACCESS HIGH     4
SAFETY SEAT- CHILD UNRESTRAINED UNDER 8      3
"LICENSE PLATE-WRONG      3
MADE U-TURN AT INTERSECTION                  3
DROVE WITHOUT LIGHTS WHEN REQUIRED           3
SAFETY BELT- DRIVER                          3
DISREGARD NO TURN ON RED LIGHT               3
SPEEDING INTERSTATE HWY 35                    3
TURN TURNED LEFT FROM WRONG LANE             3
REGISTRATION OPERATE UNREGISTERED MOTOR      3
DL DOMICILED IN TEXAS GREATER THAN 90 DA     3
DROVE WRONG WAY ON ONE-WAY ROADWAY           2
PASS DISREGARD NO PASSING ZONE               2
FAIL TO DRIVE IN A SINGLE LANE               2
SPEEDING IN 30 MILE HOUR ZONE                2
DL FAIL TO REPORT CHANGE OF ADDRESS OR N     2
```

Speeding violation cases top the list with 91, followed by no drivers license.


Query 5:

Top people with most violations

```
Time taken: 1.029 seconds, Fetched 102 row(s)
spark-sql> SELECT SPLIT(name, ' ')[0] AS first_name, COUNT(*) AS total_cases
        > FROM traffic_closed_cases_2024
        > GROUP BY SPLIT(name, ' ')[0]
        > ORDER BY total_cases DESC
        > LIMIT 5;
"ZEMCIK 9
"PEREZ   8
"GONZALEZ        8
"RODRIGUEZ       7
"MALDONADO       6
Time taken: 0.725 seconds, Fetched 5 row(s)
spark-sql> █
```

Zemick toped with 9, followed by Perez and Gonzalez with 8 each.

# Queries in BigQuery ( combined Queries of both datasets )

## City wise crime and total number of crime type count



## Most common crimes and their violations

```
28  FROM common_crimes
29  UNION ALL
30  SELECT
31      city,
32      violation AS type,
33      violation_count AS count,
34      "Traffic Violation" AS category
35  FROM common_violations
36  ORDER BY count DESC;
37
```



## Daily trends in Traffic cases:



## Highest number of violations for a citation number:

Denton_c...red ✕ | traffic-cases ✕ | *Untitled query ✕

Untitled query  ▶ RUN  🕐 SCHEDULE  OPEN IN ▾  ⚙ MORE ▾  💾 SAVE ▾  ⬇ DOWNLOAD  ➕ SHARE ▾          ✅ Query completed.

```
1  SELECT
2      cit_citation_no,
3      COUNT(*) AS violation_count
4  FROM Crime_data.`traffic-cases`
5  GROUP BY cit_citation_no
6  HAVING violation_count > 1
7  ORDER BY violation_count DESC
8  LIMIT 10;
9
```

Press Alt+F1 for Accessibility Options.

## Query results

SAVE RESULTS ▾   OPEN IN ▾

JOB INFORMATION    **RESULTS**    CHART    JSON    EXECUTION DETAILS    EXECUTION GRAPH

| Row | cit_citation_no ▾ | violation_count ▾ | |
|---|---|---|---|
| 1 | 10399984 | 3 | |
| 2 | 911451 | 3 | |
| 3 | 10364917 | 3 | |
| 4 | 10364927 | 3 | |
| 5 | 10444799 | 3 | |
| 6 | 10251290 | 3 | |
| 7 | 10446796 | 3 | |
| 8 | 80006206 | 3 | |

Results per page:  50 ▾   1 – 10 of 10   |<   <   >   >|