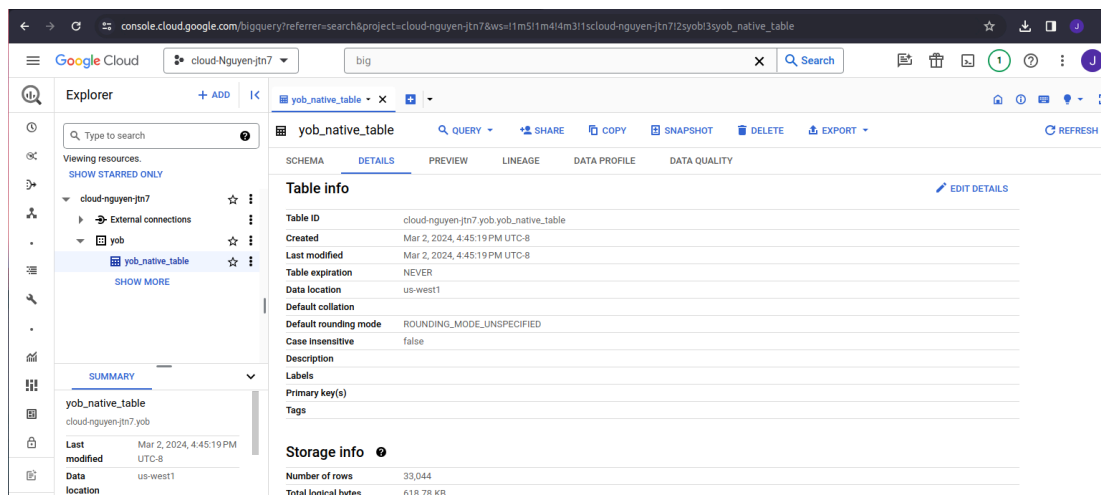


9.1g

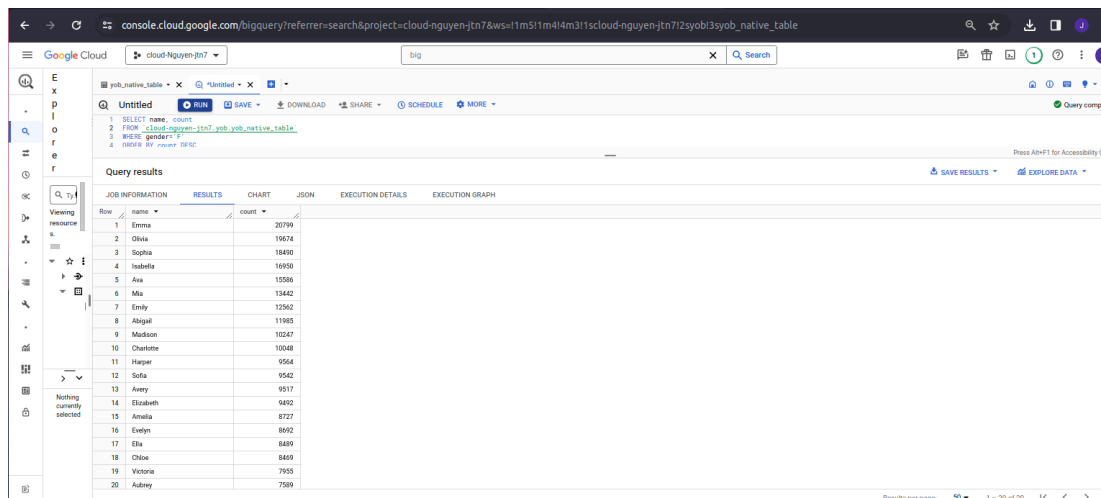
Then, click on the "Details" tab.

- Take a screenshot of the table's details that includes the number of rows in the table.



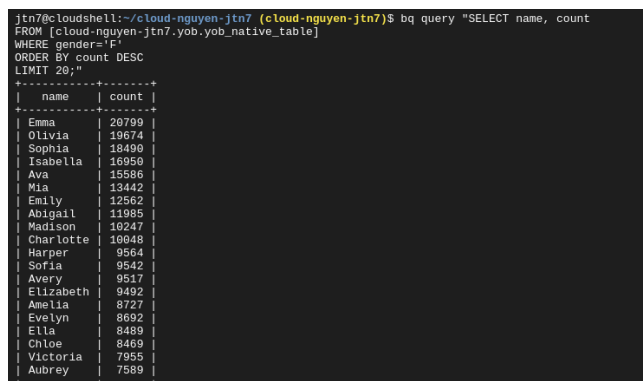
Run the query.

- Screenshot the query results and include it in your lab notebook



Via command-line bq command

- Screenshot your results and include it in your lab notebook



Via interactive bq session

- Screenshot your results and include it in your lab notebook

```
jtn7@cloudshell:~/cloud-nguyen-jtn7 [cloud-nguyen-jtn7]$ bq shell
Welcome to BigQuery! (Type help for more information.)
Cannot read termcap database;
using dumb terminal settings.
cloud-nguyen-jtn7> SELECT name, count FROM [cloud-nguyen-jtn7.yob.yob_native_table] WHERE gender='F' ORDER BY count DESC LIMIT 20;
+-----+-----+
| name | count |
+-----+-----+
| Emma | 20799 |
| Olivia | 19674 |
| Sophia | 18499 |
| Isabella | 16959 |
| Ava | 15586 |
| Mia | 13442 |
| Emily | 12562 |
| Abigail | 11985 |
| Madison | 10247 |
| Charlotte | 10048 |
| Harper | 9564 |
| Sofia | 9542 |
| Avery | 9517 |
| Elizabeth | 9492 |
| Amelia | 8727 |
| Evelyn | 8692 |
| Ella | 8489 |
| Chloe | 8469 |
| Victoria | 7955 |
| Aubrey | 7589 |
+-----+-----+
```

Finally, run a query on your name. How popular was it?

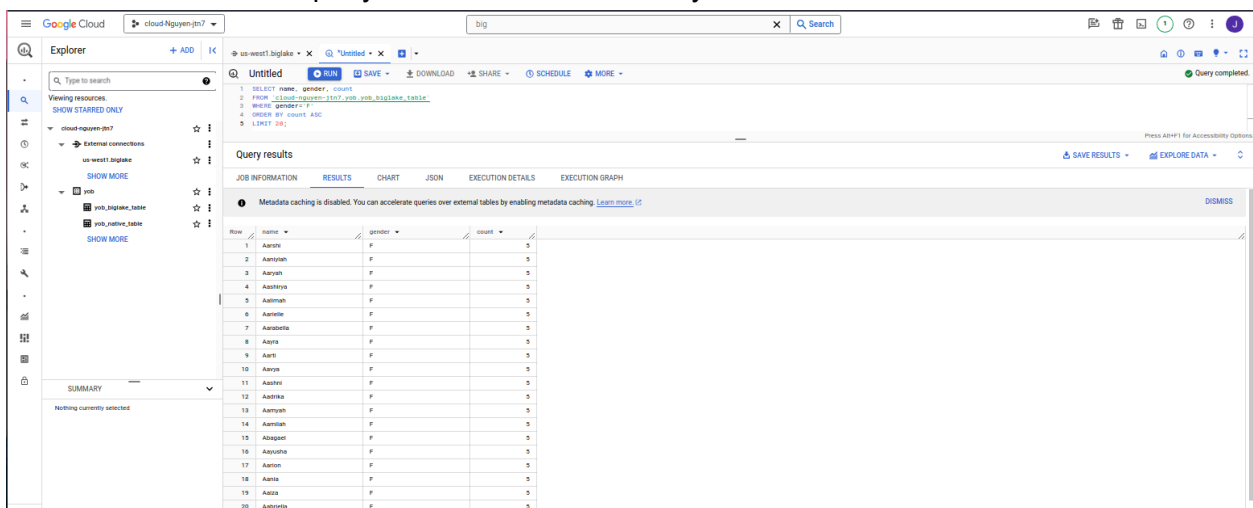
- Screenshot your results and include it in your lab notebook

```
cloud-nguyen-jtn7> SELECT name, gender, count FROM [cloud-nguyen-jtn7.yob.yob_native_table] WHERE name='Joseph';
+-----+-----+-----+
| name | gender | count |
+-----+-----+-----+
| Joseph | F | 6 |
| Joseph | M | 11995 |
+-----+-----+-----+
```

There were 6 women named Joseph and 11,995 men named Joseph in the dataset.

Run the query.

- Screenshot the query results and include it in your lab notebook



9.2g

Before running the query, answer the following question for your lab notebook:

- How much less data does this query process compared to the size of the table?

Selecting the data takes 21.94 GB while querying with the given query statement takes 3.05 GB. Thus the query process is 21.94 - 3.05 = 18.89 GB less than the size of the table.

Run the query and answer the following for your lab notebook:

- How many twins were born during this time range? **375362 twins were born**
- How much lighter on average are they compared to single babies? **~2.17116 lbs lighter**

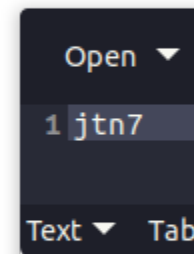
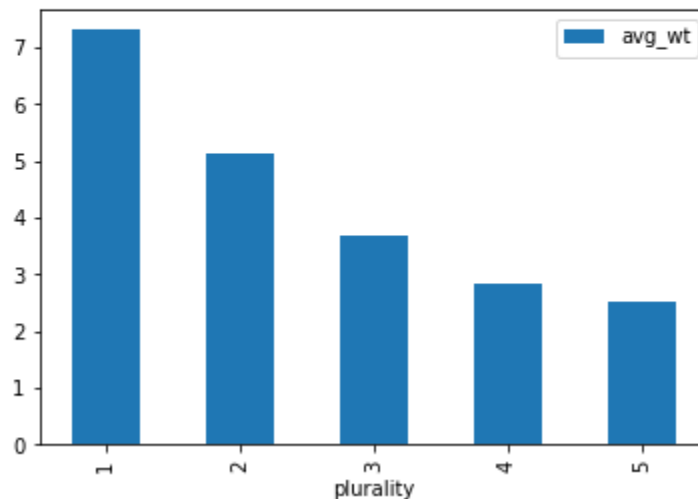
In examining the plots, which two features are the strongest predictors for a newborn baby's weight?

- Show the plots generated for the two most important features for your lab notebook

Plurality and Gestation time are the two most important features for predicting a newborn baby's weight.

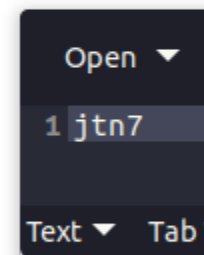
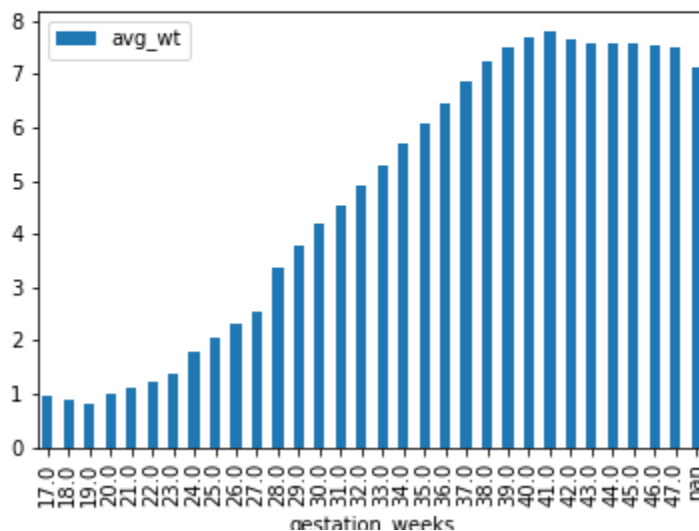
```
[6]: df = get_distinct_values('plurality')
df.plot(x='plurality', y='avg_wt', kind='bar')
```

```
[6]: <matplotlib.axes._subplots.AxesSubplot at 0x7f8d2d6c6f10>
```



```
[8]: df = get_distinct_values('gestation_weeks')
df.plot(x='gestation_weeks', y='avg_wt', kind='bar')
```

```
[8]: <matplotlib.axes._subplots.AxesSubplot at 0x7f8d2c997b50>
```



- What day saw the largest spike in trips to grocery and pharmacy stores?

On March 13, 2020, there was a 17% increase from baseline to grocery and pharmacy stores.

- On the day the stay-at-home order took effect (3/23/2020), what was the total impact on workplace trips?

There was a 49% decrease on workplace trips from the normal baseline.

- Which three airports were impacted the most in April 2020 (the month when lockdowns became widespread)?

Detroit Metropolitan Wayne County, McCarran International, and San Francisco International

- Run the query again using the month of August 2020. Which three airports were impacted the most?

McCarran International, Detroit Metropolitan Wayne County, and San Francisco International

- What table and columns identify the place name, the starting date, and the number of excess deaths from COVID-19?

Table: excess_death

Columns: placename, start_date, excess_deaths

- What table and columns identify the date, county, and deaths from COVID-19?

Table: us_counties

Columns: date, county, deaths

- What table and columns identify the date, state, and confirmed cases of COVID-19?

Table: us_states

Columns: date, state_name, confirmed_cases

- What table and columns identify a county code and the percentage of its residents that report they always wear masks?

Table: mask_use_by_county

Column: county_fips_code, always

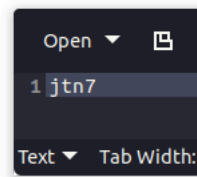
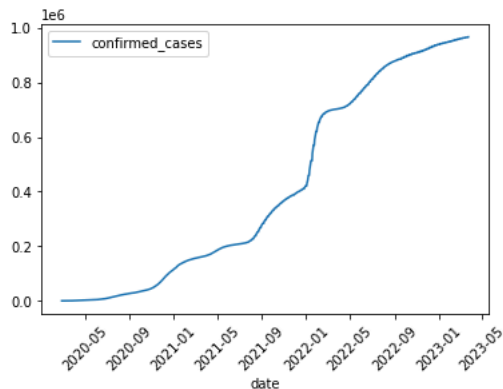
Within the notebook, perform the following query and plot the results using a line graph. The query pulls out the confirmed cases of COVID-19 across the state of Oregon by date.

- Show a screenshot of the plot and the code used to generate it for your lab notebook

```
[11]: query_string = """
SELECT date, confirmed_cases
FROM bigquery-public-data.covid19_nyt.us_states
WHERE state_name = 'Oregon'
ORDER BY date ASC
"""

from google.cloud import bigquery
df = bigquery.Client().query(query_string).to_dataframe()
df.plot(x='date', y='confirmed_cases', kind='line', rot=45)
```

[11]: <matplotlib.axes._subplots.AxesSubplot at 0x7f8d2c728f10>



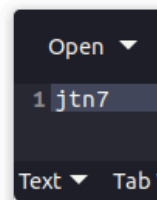
- From within your Jupyter notebook, run the query and write code that shows the first 10 states that reached 1000 deaths from COVID-19. Take a screenshot for your lab notebook.

```
[13]: query_string = """
SELECT state_name, MIN(date) as date_of_1000
FROM bigquery-public-data.covid19_nyt.us_states
WHERE deaths > 1000
GROUP BY state_name
ORDER BY date_of_1000 ASC
"""

from google.cloud import bigquery
df = bigquery.Client().query(query_string).to_dataframe()
df.head(10)
```

[13]:

| | state_name | date_of_1000 |
|---|---------------|--------------|
| 0 | New York | 2020-03-29 |
| 1 | New Jersey | 2020-04-06 |
| 2 | Michigan | 2020-04-09 |
| 3 | Louisiana | 2020-04-14 |
| 4 | Massachusetts | 2020-04-15 |
| 5 | Illinois | 2020-04-16 |
| 6 | California | 2020-04-17 |
| 7 | Connecticut | 2020-04-17 |
| 8 | Pennsylvania | 2020-04-17 |
| 9 | Florida | 2020-04-24 |



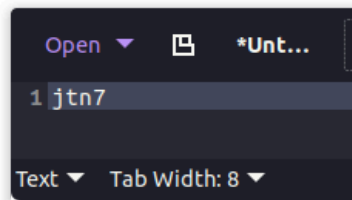
- Take a screenshot for your lab notebook of the Top 5 counties and the states they are located in.

```
[15]: query_string = """
SELECT DISTINCT mu.county_fips_code, mu.always, ct.county, ct.state_name
FROM bigquery-public-data.covid19_nyt.mask_use_by_county as mu
LEFT JOIN bigquery-public-data.covid19_nyt.us_counties as ct
ON mu.county_fips_code = ct.county_fips_code
ORDER BY mu.always DESC
"""

from google.cloud import bigquery
df = bigquery.Client().query(query_string).to_dataframe()
df.head(5)
```

```
[15]:
```

| | county_fips_code | always | county | state_name |
|---|------------------|--------|----------|------------|
| 0 | 06027 | 0.889 | Inyo | California |
| 1 | 36123 | 0.884 | Yates | New York |
| 2 | 48229 | 0.880 | Hudspeth | Texas |
| 3 | 06051 | 0.880 | Mono | California |
| 4 | 48141 | 0.877 | El Paso | Texas |



Deaths in Multnomah county

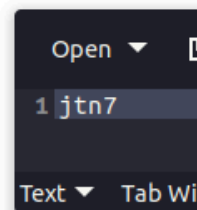
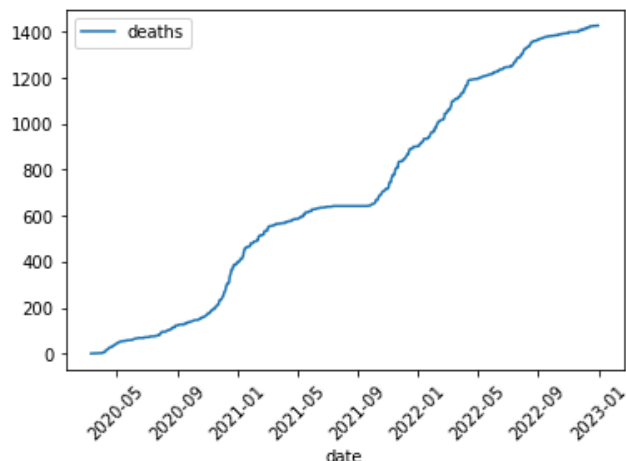
Construct a query string that obtains the number of deaths from COVID-19 that have occurred in Multnomah county for each day in the dataset, ensuring the data is returned in ascending order of date. Run the query and obtain the results.

- Plot the results and take a screenshot for your lab notebook.

```
[20]: query_string = """
SELECT date, deaths
FROM bigquery-public-data.covid19_nyt.us_counties as ct
WHERE ct.county = 'Multnomah'
ORDER BY date
"""

from google.cloud import bigquery
df = bigquery.Client().query(query_string).to_dataframe()
df.plot(x='date', y='deaths', kind='line', rot=45)
```

```
[20]: <matplotlib.axes._subplots.AxesSubplot at 0x7f8d2c568a50>
```



Deaths in Oregon

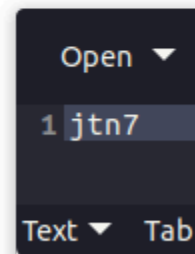
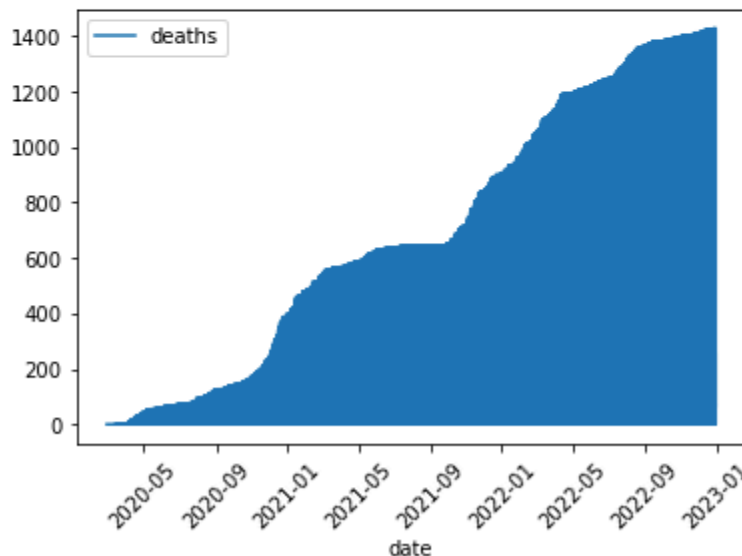
Construct a query string that obtains the number of deaths from COVID-19 that have occurred in Oregon for each day in the dataset, ensuring the data is returned in ascending order of date. Run the query and obtain the results.

- Plot the results and take a screenshot for your lab notebook.

```
[22]: query_string = """
      SELECT date, deaths
      FROM bigquery-public-data.covid19_nyt.us_counties as ct
      WHERE ct.state_name = 'Oregon'
      ORDER BY date
      """

      from google.cloud import bigquery
      df = bigquery.Client().query(query_string).to_dataframe()
      df.plot(x='date', y='deaths', kind='line', rot=45)
```

```
[22]: <matplotlib.axes._subplots.AxesSubplot at 0x7f8d2c6a8e90>
```



9.3g

- How long did the job take to execute?

It took about 43 seconds.

- Examine output.txt and show the estimate of π calculated.

Pi is roughly 3.141501591415016

After rescaling

- How long did the job take to execute? How much faster did it take?

It took 31 seconds so It was over 20% faster than previously.

- Examine output2.txt and show the estimate of π calculated.

Pi is roughly 3.1413466714134666

9.4

Reverse-engineer the code and make note of where the input is taken from and where the output goes to by default.

Answer the following questions for your lab notebook.

- Where is the input taken from by default?

'../javahelp/src/main/java/com/google/cloud/training/dataanalyst/javahelp/'

- Where does the output go by default?

'/tmp/output'

- Examine both the `getPackages()` function and the `splitPackageName()` function. What operation does the `'PackageUse()'` transform implement?

`PackageUse()` calls `getPackages()` as a generator function. `getPackages()` searches the keyword by line where `';` denotes the end of a line; it calls `splitPackageName()` which splits up the package names at each `'.'` by using python brackets for slicing the strings and returns an array of these substrings.

- Look up Beam's `CombinePerKey`. What operation does the `TotalUse` operation implement?

`CombinePerKey` is called with `sum`, meaning that it takes in an iterable of numbers and adds them together. In this case, we are summing up how many times the `PackageUse()` function is call, so `TotalUse` represents how many package uses there are.

The operations in the pipeline mimic a Map-Reduce pattern, demonstrating Beam's ability to support it.

Answer the following questions for your lab notebook.

- Which operations correspond to a "Map"?

`GetImport` and `PackageUse` call `beam.FlatMap()` which corresponds to Map.

- Which operation corresponds to a "Shuffle-Reduce"?

`TotalUse` calls `beam.CombinePerKey()` which shuffle and reduce the data.

- Which operation corresponds to a "Reduce"?

`Top_5` calls `beam.transforms.combiners.Top.Of()` which corresponds to a reduce.

Go to where the output file is written out and cat the file.

- Take a screenshot of its contents
- Explain what the data in this output file corresponds to based on your understanding of the program

```
jtn7@cloudshell:/tmp (cloud-nguyen-jtn7)$ cat output-00000-of-00001
[('org', 45), ('org.apache', 44), ('org.apache.beam', 44), ('org.apache.beam.sdk', 43), ('org.apache.beam.sdk.transforms', 16)]
jtn7@cloudshell:/tmp (cloud-nguyen-jtn7)$
```

The data in the file represents the name of the Java package and how many times that package was called.

When the code is invoked as a module, it executes the `run()` function. Within the function, a pipeline `p` is incrementally constructed using the Beam syntax (`lines`, `counts`, `output`) before the entire pipeline is executed via `p.run()`. Examine the code that implements the function and answer the following questions for your lab notebook:

- What are the names of the stages in the pipeline?

Split, PairWithOne, GroupAndSum

- Describe what each stage does.

Split: splits a line of text into individual words

PairWithOne: Each word is paired with a number

GroupAndSum: Matching words are grouped together and their number summed together. Thus we will get a list of unique words and their count.

After running the script, perform the following and show a screenshot of the results in your lab notebook:

- Use `wc` with an appropriate flag to determine the number of different words in King Lear.

```
(env) jtn7@cloudshell:~/training-data-analyst/courses/data_analysis/lab2/python (cloud-nguyen-jtn7)$ wc -w outputs-00000-of-00001
9568 outputs-00000-of-00001
(env) jtn7@cloudshell:~/training-data-analyst/courses/data_analysis/lab2/python (cloud-nguyen-jtn7)$
```

- Use `sort` with appropriate flags to perform a numeric sort on the key field containing the count for each word in descending order. Pipe the output into `head` to show the top 3 words in King Lear and the number of times they appear

```
sort: cannot read: 0: No such file or directory
(env) jtn7@cloudshell:~/training-data-analyst/courses/data_analysis/lab2/python (cloud-nguyen-jtn7)$ sort -t: -k 2nr,2 outputs-00000-of-00001 > sorted_output
(env) jtn7@cloudshell:~/training-data-analyst/courses/data_analysis/lab2/python (cloud-nguyen-jtn7)$ ls
env      grep.py  install_packages_OLD.sh  ls_popular.py  JavaProjectsThatNeedHelp.py  outputs-00000-of-00001  sorted_output
grep.py  head    install_packages.sh      JavaProjectsThatNeedHelp.py  OLD_grep.py
(env) jtn7@cloudshell:~/training-data-analyst/courses/data_analysis/lab2/python (cloud-nguyen-jtn7)$ head -n 3 sorted_output
the: 786
I: 622
and: 594
(env) jtn7@cloudshell:~/training-data-analyst/courses/data_analysis/lab2/python (cloud-nguyen-jtn7)$
```

Perform the following and show a screenshot of the results in your lab notebook:

- Use the previous method to show the top 3 words in King Lear, case-insensitive, and the number of times they appear.

```
(env) jtn7@cloudshell:~/training-data-analyst/courses/data_analysis/lab2/python (cloud-nguyen-jtn7)$ sort -t: -k 2nr,2 outputs-00000-of-00001 > sorted_output
(env) jtn7@cloudshell:~/training-data-analyst/courses/data_analysis/lab2/python (cloud-nguyen-jtn7)$ head -n 3 sorted_output
the: 908
and: 738
I: 622
(env) jtn7@cloudshell:~/training-data-analyst/courses/data_analysis/lab2/python (cloud-nguyen-jtn7)$
```

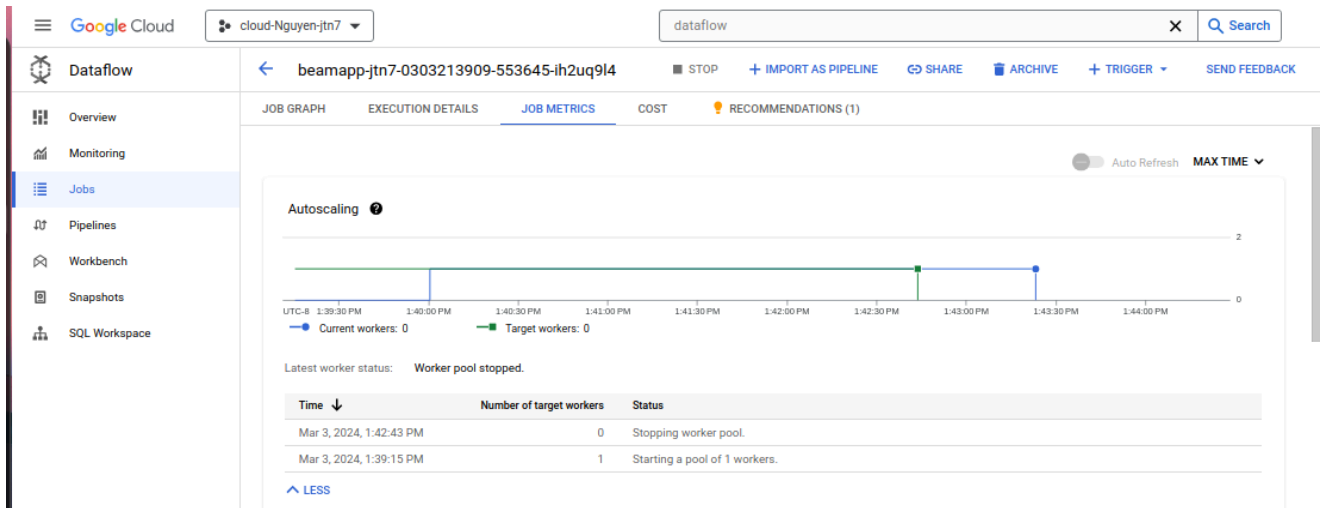
After executing the program, which takes around 5 minutes to complete, visit Dataflow in the web console and click on the Dataflow job that was executed. Examine both "Job Graph" and "Job Metrics". Include the following in your lab notebook:

- The part of the job graph that has taken the longest time to complete.

Writing the data took the longest time to complete.

| Step name | Status | Wall time | Stages | Input steps | Output steps |
|-------------|-----------|-----------|----------|------------------------------------|---------------------------------|
| Read | Succeeded | 1 second | F25, F30 | — | Split |
| Split | Succeeded | 0 seconds | F25 | Read.../ParDo(SplitBoundsSourceOf) | lowercase |
| lowercase | Succeeded | 0 seconds | F25 | Split | PairWithOne |
| PairWithOne | Succeeded | 0 seconds | F25 | lowercase | GroupAndSum/GroupByKey |
| GroupAndSum | Succeeded | 0 seconds | F25, F23 | PairWithOne | Format |
| Format | Succeeded | 0 seconds | F25 | GroupAndSum/Combine | Write.../WindowInto(WindowedOf) |
| Write | Succeeded | 3 seconds | F25, F26 | Format | — |

- The autoscaling graph showing when the worker was created and stopped.



- Examine the output directory in Cloud Storage. How many files has the final write stage in the pipeline created? **1 file, the output file which is stored in the results directory.**

cloud-nguyen-jtn7

Location: us (multiple regions in United States) Storage class: Standard Public access: Subject to object ACLs Protection: None

OBJECTS CONFIGURATION PERMISSIONS PROTECTION LIFECYCLE OBSERVABILITY INVENTORY REPORTS

Buckets > cloud-nguyen-jtn7 > results

UPLOAD FILES UPLOAD FOLDER CREATE FOLDER TRANSFER DATA MANAGE HOLDS EDIT RETENTION DOWNLOAD DELETE

Filter by name prefix only Filter objects and folders

| <input type="checkbox"/> | Name | Size | Type | Created | Storage class | Last modified | Public access | Version history |
|--------------------------|------------------------|---------|------------|-------------------------|---------------|-------------------------|---------------|-----------------|
| <input type="checkbox"/> | outputs-00000-of-00001 | 42.7 KB | text/plain | Mar 3, 2024, 1:42:42 PM | Standard | Mar 3, 2024, 1:42:42 PM | Not public | — |

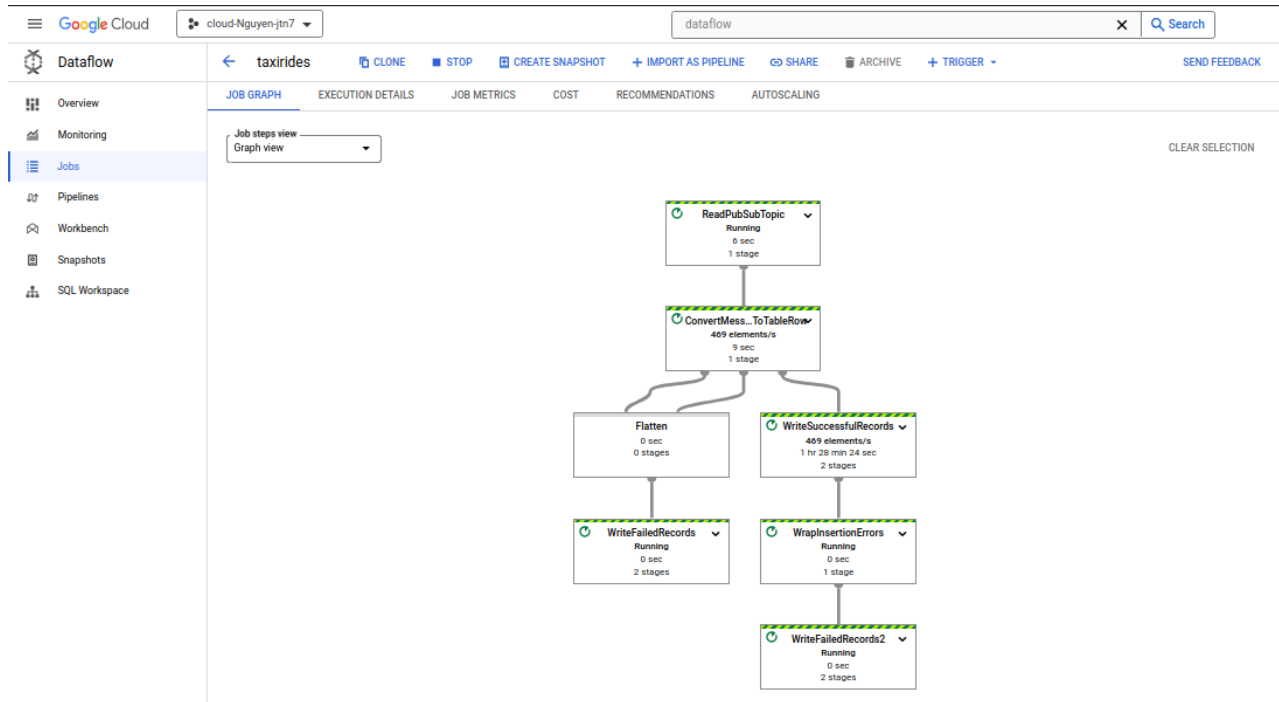
The data that is returned is a JSON object.

- Take a screenshot listing the different fields of this object.

```
(env) jtn7@cloudshell:~ (cloud-nguyen-jtn7) $ gcloud pubsub subscriptions pull taxisub --auto-ack
DATA: {"ride_id":"3b090b10-cc44-4ebd-aae7-41b5f993950a","point_id":1785,"latitude":40.60396,"longitude":-74.01998,"timestamp":"2024-03-03T17:08:42.45149-05:00","meter_reading":22.562826,"meter_increment":0.012648239,"ride_status":"enroute","passenger_count":3}
```

Then, visit the Cloud Dataflow console in the web UI and view the running job. Wait until the job is functioning and has processed 5-10 minutes worth of raw data from the stream..

- Take a screenshot of the pipeline that includes its stages and the number of elements per second being handled by individual stages.



Visit the BigQuery web console and navigate to the table we've created within our project at taxirides:realtime. Examine the table.

- Take a screenshot showing the number of passengers and the amount paid for the first ride

The screenshot shows the Google Cloud BigQuery console for a project named 'cloud-nguyen-jtn7'. The table 'realtime' in the dataset 'taxirides' is displayed. The table is partitioned and contains 18 rows of data. The columns are: ride_id, point_idx, latitude, longitude, timestamp, meter_reading, meter_increment, ride_status, and passenger_count.

| Row | ride_id | point_idx | latitude | longitude | timestamp | meter_reading | meter_increment | ride_status | passenger_count |
|-----|--------------------------------|-----------|---------------|---------------|---------------------------------|---------------|-----------------|-------------|-----------------|
| 1 | af31e3d-0789-485e-9088-da3... | 30 | 40.7365400... | -73.97878 | 2024-03-03 22:16:02.331600 U... | 1.39 | 0.04633333 | enroute | 1 |
| 2 | 06583be-f13d-4793-9d73-1fd... | 2052 | 40.67401 | -73.86390... | 2024-03-03 22:16:02.370450 U... | 47.043453 | 0.02292566 | enroute | 1 |
| 3 | ebd60345-6522-4c0a-90f5-480... | 366 | 40.71226 | -73.96808 | 2024-03-03 22:16:02.445890 U... | 10.511607 | 0.028720237 | enroute | 6 |
| 4 | ce8156c-9786-4d5d-90e9-27a... | 316 | 40.7749100... | -73.96317 | 2024-03-03 22:16:02.193140 U... | 9.931429 | 0.03142857 | enroute | 2 |
| 5 | 199f70e-0703-4cb7-ab04-4e3... | 280 | 40.7414200... | -73.953960... | 2024-03-03 22:16:02.282670 U... | 7.4489307 | 0.026603324 | enroute | 1 |
| 6 | f55149a3-dda5-43a4-ad5f-089... | 2039 | 40.8740300... | -73.905120... | 2024-03-03 22:16:02.203440 U... | 38.88864 | 0.019072408 | enroute | 1 |
| 7 | b9845d39-5996-453c-bd32-9a... | 157 | 40.71571 | -73.97945 | 2024-03-03 22:16:02.270930 U... | 3.0715137 | 0.019563781 | enroute | 1 |
| 8 | 48d01cb-2f12-4583-becb-370... | 1128 | 40.6645700... | -73.991750... | 2024-03-03 22:16:02.256690 U... | 24.161568 | 0.021419829 | enroute | 5 |
| 9 | 9551ade1-04a5-4b69-8c4d-da6... | 154 | 40.77454 | -73.957940... | 2024-03-03 22:16:02.327340 U... | 4.706383 | 0.03056093 | enroute | 1 |
| 10 | 9d8c92d0-ca34-4410-b08-835... | 936 | 40.7055600... | -73.9504 | 2024-03-03 22:16:02.407890 U... | 23.351733 | 0.024948454 | enroute | 1 |
| 11 | 31ee8635-7143-4619-b49f-007... | 19 | 40.72688 | -73.98556 | 2024-03-03 22:16:02.762330 U... | 0.687905 | 0.036173183 | enroute | 1 |
| 12 | 13270eb-2dff-4768-9ede-789... | 693 | 40.7452700... | -73.970580... | 2024-03-03 22:16:02.773750 U... | 12.803409 | 0.018475337 | enroute | 1 |
| 13 | 8eb479f1-cbf7-4550-eccd-635... | 630 | 40.74398 | -73.971910... | 2024-03-03 22:16:02.791300 U... | 12.955223 | 0.020563846 | enroute | 1 |
| 14 | 18c478b-f7d8-4832-bf65-7b0a... | 473 | 40.712 | -73.96724 | 2024-03-03 22:16:02.521880 U... | 11.376658 | 0.024052132 | enroute | 1 |
| 15 | 78369002-041d-4206-904c-24... | 1005 | 40.7255300... | -73.898130... | 2024-03-03 22:16:02.659240 U... | 25.808577 | 0.025680175 | enroute | 2 |
| 16 | 788b9b3d-a8f8-446f-ae2b-760... | 713 | 40.69154 | -73.81027 | 2024-03-03 22:16:02.691900 U... | 14.071288 | 0.019735327 | enroute | 1 |
| 17 | 6af3cd96-c5d6-4a38-a018-338... | 2955 | 40.84508 | -73.90815 | 2024-03-03 22:16:03.085660 U... | 64.689514 | 0.021891546 | enroute | 1 |
| 18 | 056df151-7116-4800-9471-def... | 278 | 40.80563 | -73.909810... | 2024-03-03 22:16:03.272030 U... | 8.754043 | 0.03148936 | enroute | 1 |

- Take a screenshot showing the estimated number of rows in the table.

The screenshot shows the Google Cloud BigQuery interface. On the left is the navigation menu with options like Analysis, Data transfers, Scheduled queries, Analytics Hub, Dataform, Partner Center, Migration, Assessment, SQL translation, Administration, Monitoring, Capacity management, BI Engine, and Policy tags. The main area displays the 'realtime' table details. The 'DETAILS' tab is active, showing partitioning information: Partitioned by DAY, Partitioned on field timestamp, Partition expiration Partitions do not expire, and Partition filter Not required. The 'Storage info' section shows: Number of rows 6,860, Number of partitions 0, Total logical bytes 690.01 KB, Active logical bytes 690.01 KB, Long term logical bytes 0 B, Total physical bytes 0 B, Active physical bytes 0 B, Long term physical bytes 0 B, and Time travel physical bytes 0 B. The 'Streaming buffer statistics' section shows: Estimated size 59.62 MB, Estimated rows 285,659, and Earliest entry time Mar 3, 2024, 2:15:39 PM UTC-8.

- Take a screenshot showing the per-minute number of rides, passengers, and revenue for the data collected

The screenshot shows the Google Cloud BigQuery interface with a SQL query executed. The query is:

```
1 SELECT
2   FORMAT_TIMESTAMP('%t', timestamp, 'America/Los_Angeles') AS minute,
3   COUNT(DISTINCT ride_id) AS total_rides,
4   SUM(passenger_count) AS total_passengers,
5   SUM(revenue_reading) AS total_revenue
6 FROM
7   taxirides realtime
8 WHERE
9   ride_status = 'dropoff'
10 GROUP BY
11   minute
12 ORDER BY
13   minute ASC
```

 The 'Query results' section is active, displaying a table with columns: minute, total_rides, total_passengers, and total_revenue. The results show data for minutes from 14:13 to 14:23. The table has 11 rows of data.

| minute | total_rides | total_passengers | total_revenue |
|--------|-------------|------------------|-----------------|
| 14:13 | 17 | 31 | 252.7700000000 |
| 14:14 | 51 | 81 | 875.4200000000 |
| 14:15 | 44 | 126 | 1039.7199979999 |
| 14:16 | 62 | 117 | 1045.8799969999 |
| 14:17 | 52 | 85 | 824.9999937000 |
| 14:18 | 78 | 143 | 1163.9300021 |
| 14:19 | 52 | 97 | 757.9000016999 |
| 14:20 | 55 | 91 | 820.9799986999 |
| 14:21 | 44 | 89 | 715.3800051000 |
| 14:22 | 43 | 87 | 675.2199960999 |
| 14:23 | 10 | 15 | 210.120002 |

- Take a screenshot showing the plot for your data for your lab notebook

