# Creating Date-Partitioned Tables in BigQuery

**Task 1**. Create a new dataset
1. First, you will create a dataset to store your tables.
2. Click the three dots next to your Qwiklabs project ID and select **Create dataset**:
3. Name your dataset **ecommerce**.
4. Click **Create dataset**.


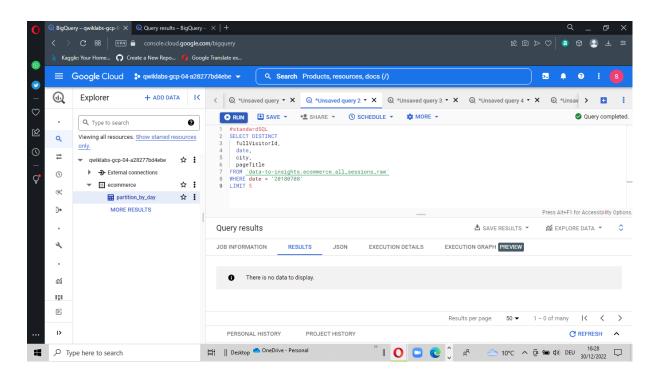**Task 2**. Creating tables with date partitions


A partitioned table is a table that is divided into segments, called partitions, that make it easier to manage and query your data. By dividing a large table into smaller partitions, you can improve query performance, and control costs by reducing the number of bytes read by a query.

Now create a new table and bind a date or timestamp column as a partition. Before we do that, let's explore the data in the non-partitioned table first.

**Query webpage analytics for a sample of visitors in 2017**
1. In the **Query Editor**, add the below query:
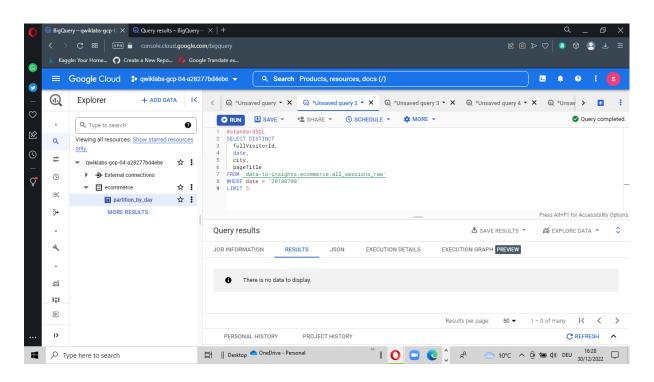
**Solution:**

**Query webpage analytics for a sample of visitors in 2018**

Let's modify the query to look at visitors for 2018 now.

1. Click **COMPOSE NEW QUERY** to clear the **Query Editor**, then add this new query. Note the WHERE date parameter is changed to 20180708:

**solution:**



*Common use-cases for date-partitioned tables*

Scanning through the entire dataset everytime to compare rows against a WHERE condition is wasteful. This is especially true if you only really care about records for a specific period of time like:
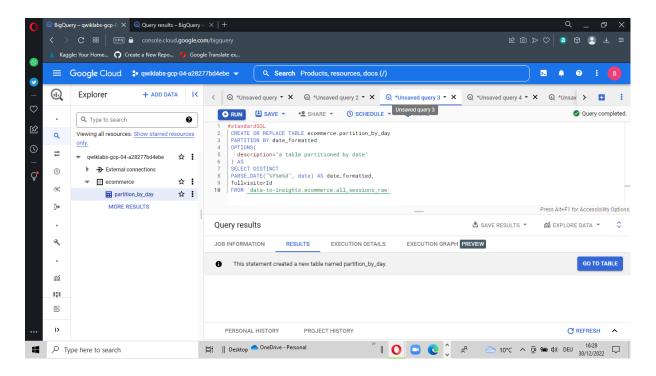
- All transactions for the last year
- All visitor interactions within the last 7 days
- All products sold in the last month

Instead of scanning the entire dataset and filtering on a date field like we did in the earlier queries, Now set up a date-partitioned table. This allows you to completely ignore scanning records in certain partitions if they are irrelevant to our query.

*Create a new partitioned table based on date*

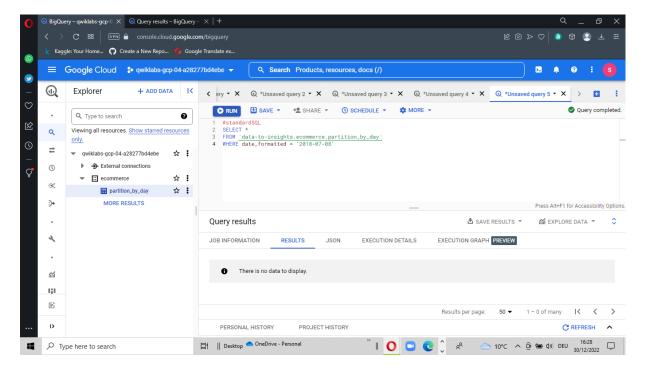1. Click **COMPOSE NEW QUERY** , add the query, then click **Run**:

**solution:**



In the above query, note the new option - PARTITION BY a field. The two options available to partition are DATE and TIMESTAMP. The PARSE_DATE function is used on the date field (stored as a string) to get it into the proper DATE type for partitioning.

**Task 3**. View data processed with a partitioned table

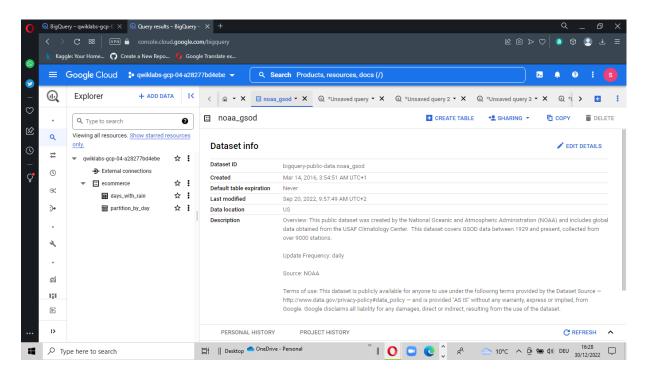1. Run the query, and note the total bytes to be processed:

**Task 4**. Creating an auto-expiring partitioned table

Auto-expiring partitioned tables are used to comply with data privacy statutes and can be used to avoid unnecessary storage (which you'll be charged for in a production environment). If you want to create a rolling window of data, add an expiration date so the partition disappears after you're finished using it.

**Explore the available NOAA weather data tables**
1. In the left menu, in Explorer, click on **Add Data** and select **Explore public datasets**.
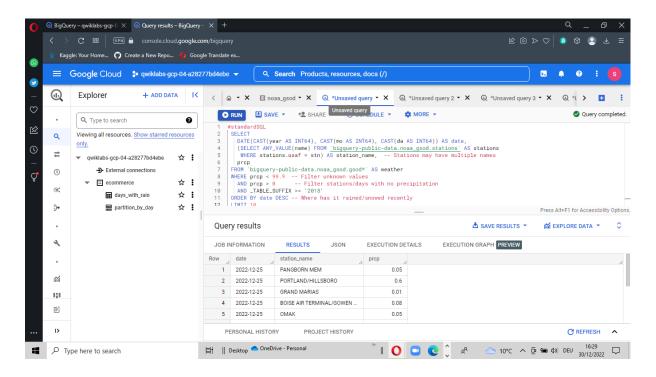
**solution:**



Your goal is to create a table that:

- Queries on weather data from 2018 onward
- Filters to only include days that have had some precipitation (rain, snow, etc.)
- Only stores each partition of data for 90 days from that partition's date (rolling window)
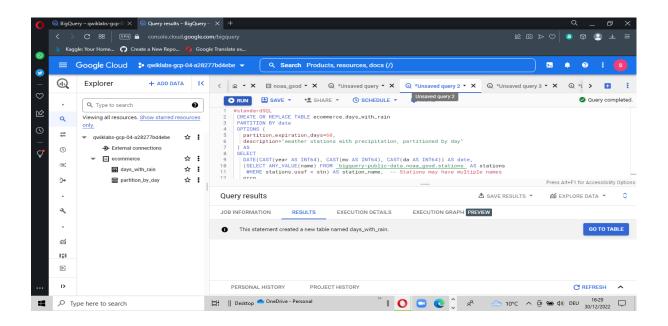
**solution:**



**NB:** The table wildcard * used in the FROM clause to limit the amount of tables referred to in the *TABLE_SUFFIX* filter.
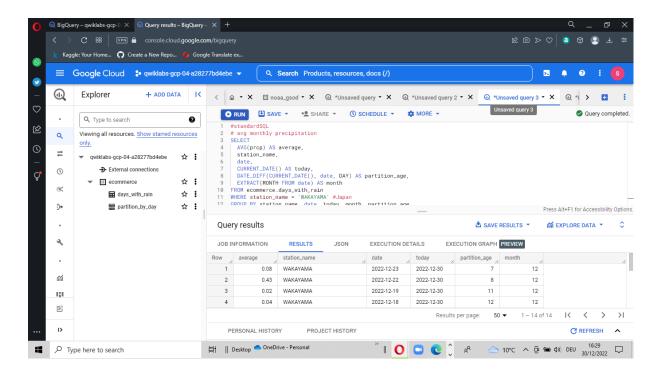
**Task 5**. create a partitioned table

- Modify the previous query to create a table with the below specifications:

  o Table name: ecommerce.days_with_rain
  o Use the date field as your PARTITION BY
  o For OPTIONS, specify partition_expiration_days = 60
  o Add the table description = "weather stations with precipitation, partitioned by day".

**solution:**



Below is a query which tracks the average rainfall for the NOAA weather station in Wakayama, Japan which has significant precipitation.

**solution:**

**Task 6**. Confirm the oldest partition_age is at or below 60 days

Update the ORDER BY clause to show the oldest partitions first

**Solution:**