

End Lab

01:23:28

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked.
[Learn more.](#)

[Open Google Console](#)

Username

student-02-3aad63b3217i



1 hour 30 minutes

Free



Password

rrQDl7TvzdxV



GCP Project ID

qwiklabs-gcp-01-0d0ccc1



Overview

[BigQuery](#) is Google's fully managed, NoOps, low cost analytics database. With BigQuery you can query terabytes and terabytes of data without having any infrastructure to manage or needing a database administrator. BigQuery uses SQL and can take advantage of the pay-as-you-go model. BigQuery allows you to focus on analyzing data to find meaningful insights.

The dataset you'll use is an [ecommerce dataset](#) that has millions of Google Analytics records for the [Google Merchandise Store](#) loaded into BigQuery. You have a copy of that dataset for this lab and will explore the available fields and row for insights.

In this lab you will query partitioned datasets and create your own dataset partitions to improve query performance and reduce cost.

[Overview](#)[Setup](#)[Create a new dataset](#)[Creating tables with date partitions](#)[View data processed with a partitioned table](#)[Creating an auto-expiring partitioned table](#)[Your turn: Create a Partitioned Table](#)[Confirm the oldest partition_age is at or below 60 days](#)[Congratulations!](#)[End your lab](#)

Setup

For each lab, you get a new GCP project and set of resources for a fixed time at no cost.

1. Make sure you signed into Qwiklabs using an [incognito window](#).
2. Note the lab's access time (for example, **02:00:00**) and make sure you can finish in that time block.

There is no pause feature. You can restart if needed, but you have to start at the beginning.

3. When ready, click **START LAB**.
4. Note your lab credentials. You will use them to sign in to Cloud Platform Console.

Caution: When you are in the console, do not deviate from the lab instructions. Doing so may cause your account to be blocked.
[Learn more.](#)

[Open Google Console](#)

Username

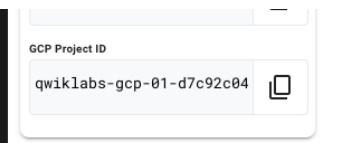
student-01-23efd9347325i



Password

gCXLv23N4fPN





5. Click **Open Google Console**.

6. Click **Use another account** and copy/paste credentials for **this** lab into the prompts.

If you use other credentials, you'll get errors or incur charges.

7. Accept the terms and skip the recovery resource page.

Do not click **End Lab** unless you are finished with the lab or want to restart it. This clears your work and removes the project.

Open BigQuery Console

1. In the Google Cloud Console, select **Navigation menu > BigQuery**.

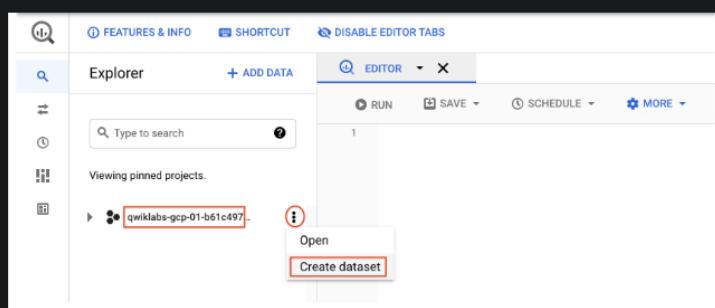
The **Welcome to BigQuery in the Cloud Console** message box opens. This message box provides a link to the quickstart guide and lists UI updates.

2. Click **Done**.

Create a new dataset

First, you will create a dataset to store your tables.

Create a new dataset within your project by clicking on the **View actions** icon next to your project ID in the Explorer section, and then selecting **CREATE DATASET**.



Set the *Dataset ID* to `ecommerce`. Leave the other options at their default values (Data Location, Default table Expiration). Click **CREATE DATASET**.

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 you will 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

In the **Query Editor**, add the below query. Before running, note the total amount of data it will process as indicated next to the query validator icon: "This query will process 1.74 GB when run".

```
#standardSQL
SELECT DISTINCT
    fullVisitorId,
    date,
    city,
    pageTitle
FROM `data-to-insights.ecommerce.all_sessions_raw`
WHERE date = '20170708'
LIMIT 5
```

Click **Run**.

The query returns 5 results.

Query webpage analytics for a sample of visitors in 2018

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

In the **Query Editor**, add the below query:

```
#standardSQL
SELECT DISTINCT
    fullVisitorId,
    date,
    city,
    pageTitle
FROM `data-to-insights.ecommerce.all_sessions_raw`
WHERE date = '20180708'
LIMIT 5
```

The **Query results** will tell you how much data this query will process.

Click **RUN**.

Notice that the query still processes 1.74 GB even though it returns 0 results. Why? The query engine needs to scan all records in the dataset to see if they satisfy the date matching condition in the WHERE clause. It must look at each record to compare the date against the condition of '20180708'.

Additionally, the LIMIT 5 does not reduce the total amount of data processed, which is a common misconception.

Why did the previous query return 0 records but still scan through 1.74GB of data?

- ✓ Before the query runs, the query engine does not know whether 2018 data exists to satisfy the WHERE clause condition and it needs to scan through all records in a non-partitioned table.
- ✗ The query engine has the metadata for each partition stored but still needs to scan all records even if the table is partitioned.
- ✗ The query was written incorrectly



Submit

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, we will now setup a date-partitioned table. This will allow us to completely ignore scanning records in certain partitions if they are irrelevant to our query.

Create a new partitioned table based on date

Click COMPOSE NEW QUERY and add the below query, then RUN:

```
#standardSQL
CREATE OR REPLACE TABLE ecommerce.partition_by_day
PARTITION BY date_formatted
OPTIONS(
    description="a table partitioned by date"
) AS
SELECT DISTINCT
PARSE_DATE("%Y%m%d", date) AS date_formatted,
fullvisitorId
FROM `data-to-insights.ecommerce.all_sessions_raw`
```

In this 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.

Click on the ecommerce dataset, then select the new partition_by_day table:

The screenshot shows the BigQuery interface. A project named "qwiklabs-gcp-d9cccd750c92585c3" is selected. Under the "ecommerce" dataset, the "partition_by_day" table is highlighted with a blue border, indicating it is the current selection.

Click on the Details tab.

Confirm that you see:

- Partitioned by: Day
- Partitioning on: date_formatted

The screenshot shows the BigQuery Table Editor for the "partition_by_day" table. The "DETAILS" tab is selected. Key information displayed includes:
Table ID: quickstart-project-01-b61e4976ae3b.ecommerce.partition_by_day
Table size: 13.81 MB
Number of rows: 478,323
Created: Jun 2, 2021, 5:04:03 PM UTC+5:30
Last modified: Jun 2, 2021, 5:04:03 PM UTC+5:30
Table expiration: NEVER
Data location: US
Description: a table partitioned by date
Table Type: Partitioned
Partitioned by: DAY
Partitioned on field: date_formatted
Partition expiration: Not required

Note: Partitions within partitioned tables on your Qwiklabs account will auto-expire after 60 days from the value in your date column. Your personal GCP account with billing-enabled will let you have partitioned tables that don't expire. For the purposes of this lab, the remaining queries will be ran against partitioned tables that have already been created.

View data processed with a partitioned table

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

```
#standardSQL  
SELECT *  
FROM `data-to-insights.ecommerce.partition_by_day`  
WHERE date_formatted = '2016-08-01'
```

This time ~25 KB or 0.025MB is processed, which is a fraction of what you queried.

Now run the below query, and note the total bytes to be processed:

```
#standardSQL  
SELECT *  
FROM `data-to-insights.ecommerce.partition_by_day`  
WHERE date_formatted = '2018-07-08'
```

You should see This query will process 0 B when run.

Why is there 0 bytes processed?

Why is there 0 bytes processed?

- The query engine knows which partitions already exist and knows that no partition exists for 2018-07-08 (the ecommerce dataset ranges from 2016-08-01 to 2017-08-01).
- The query engine has cached the results from a query we ran earlier and will return the same 10 records
- The query engine processes many fewer rows of data when you use partitions and caches each row for all users so 0 bytes are processed

Submit

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

In the left panel, click on **+ ADD DATA** and select **Explore public datasets**.

The screenshot shows the BigQuery interface. In the top navigation bar, there are tabs for 'FEATURES & INFO', 'SHORTCUT', and 'DISABLE EDITOR TABS'. Below the navigation bar, there's a toolbar with buttons for '+ ADD DATA', 'RUN', 'SAVE', 'SCHEDULE', and 'MORE'. A search bar says 'Type to search' and has a dropdown menu with 'Explore public datasets' and 'External data source'. On the left, there's a sidebar titled 'Viewing pinned projects.' showing a project named 'qwiklabs-gcp-01-b61c497...'. Underneath it, there's a tree view with 'ecommerce' expanded, showing 'partition_by_day' selected. There are also 'MORE RESULTS' links. The main pane shows a standardSQL query:

```
#standardSQL
SELECT *
FROM `data-to-insights.ecommerce.partition_by_day`
WHERE date_formatted = '2018-07-08'
```

Search for **GSOD NOAA**, and then select the dataset.

Click on **View dataset**.

Scroll through the tables in the **noaa_gsod** dataset (which are manually sharded and not partitioned).

The screenshot shows the BigQuery interface with the 'noaa_gsod' dataset selected. The left sidebar shows the dataset is pinned. The main pane displays the 'SCHEMA' tab of the dataset. It lists 141 fields under the 'Table schema' section. The first few fields are: stn (STRING), when (STRING), year (STRING), mo (STRING), da (STRING), temp (FLOAT), count_temp (INTEGER), dewp (FLOAT), count_dewp (INTEGER), slp (FLOAT), and count_slp (INTEGER). Each field has a detailed description below it.

Next, copy and paste this below query to Query editor:

The screenshot shows the BigQuery Query Editor. The code area contains the following standardSQL query:

```
#standardSQL
SELECT
  DATE(CAST(year AS INT64), CAST(mo AS INT64), CAST(da AS
INT64)) AS date,
  (SELECT ANY_VALUE(name) FROM `bigquery-public-
data.noaa_gsod.stations` AS stations
  WHERE stations.usaf = stn) AS station_name, -- Stations may
have multiple names
  prcp
FROM `bigquery-public-data.noaa_gsod.gsod*` AS weather
WHERE prcp < 99.9 -- Filter unknown values
  AND prcp > 0 -- Filter stations/days with no
precipitation
  AND _TABLE_SUFFIX >= '2021'
ORDER BY date DESC -- Where has it rained/snowed recently
LIMIT 10
```

Note that the table wildcard * used in the FROM clause to limit the amount of tables referred to in the *TABLE_SUFFIX* filter.

Note that although a LIMIT 10 was added, this still does not reduce the total amount of data scanned (about 141.6 MB) since there are no partitions yet.

Click **Run**.

Confirm the date is properly formatted and the precipitation field is showing non-zero values.

Your turn: Create a Partitioned Table

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

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

Your query should look like this:

```
#standardSQL
CREATE OR REPLACE TABLE ecommerce.days_with_rain
PARTITION BY date
OPTIONS (
    partition_expiration_days=60,
    description="weather stations with precipitation, partitioned
by day"
) AS
SELECT
    DATE(CAST(year AS INT64), CAST(mo AS INT64), CAST(da AS
INT64)) AS date,
    (SELECT ANY_VALUE(name) FROM `bigquery-public-
data.noaa_gsod.stations` AS stations
     WHERE stations.usaf = stn) AS station_name, -- Stations may
have multiple names
    prcp
FROM `bigquery-public-data.noaa_gsod.gsod*` AS weather
WHERE prcp < 99.9 -- Filter unknown values
    AND prcp > 0 -- Filter
    AND _TABLE_SUFFIX >= '2021'
```

Confirm data partition expiration is working

To confirm you are only storing data from 60 days in the past up until today, run the DATE_DIFF query to get the age of your partitions, which are set to expire after 60 days.

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

Add this query and run it:

```
#standardSQL
# avg monthly precipitation
SELECT
    AVG(prcp) AS average,
    station_name,
    date,
    CURRENT_DATE() AS today,
    DATE_DIFF(CURRENT_DATE(), date, DAY) AS partition_age,
    EXTRACT(MONTH FROM date) AS month
FROM ecommerce.days_with_rain
WHERE station_name = 'WAKAYAMA' #Japan
GROUP BY station_name, date, today, month, partition_age
ORDER BY date DESC; # most recent days first
```

Confirm the oldest partition_age is at or below 60 days

Update the ORDER BY clause to show the oldest partitions first. The date you see there
Add this query and run it:

```
#standardSQL  
# avg monthly precipitation  
SELECT  
    AVG(prcp) AS average,  
    station_name,  
    date,  
    CURRENT_DATE() AS today,  
    DATE_DIFF(CURRENT_DATE(), date, DAY) AS partition_age,  
    EXTRACT(MONTH FROM date) AS month  
FROM ecommerce.days_with_rain  
WHERE station_name = 'WAKAYAMA' #Japan  
GROUP BY station_name, date, today, month, partition_age  
ORDER BY partition_age DESC
```



Note: Your results will vary if you re-run the query in the future, as the weather data, and your partitions, are continuously updated.

Congratulations!

You've successfully created and queried partitioned tables in BigQuery.

End your lab

When you have completed your lab, click **End Lab**. Qwiklabs removes the resources you've used and cleans the account for you.

You will be given an opportunity to rate the lab experience. Select the applicable number of stars, type a comment, and then click **Submit**.

The number of stars indicates the following:

- 1 star = Very dissatisfied
- 2 stars = Dissatisfied
- 3 stars = Neutral
- 4 stars = Satisfied
- 5 stars = Very satisfied

You can close the dialog box if you don't want to provide feedback.

For feedback, suggestions, or corrections, please use the **Support** tab.

Copyright 2021 Google LLC All rights reserved. Google and the Google logo are trademarks of Google LLC. All other company and product names may be trademarks of the respective companies with which they are associated.