# Cloud Data Warehouse

The advantage of the cloud is infinite compute and infinite storage. Cloud-native data warehouse systems also allow for serverless workflows that can directly integrate Machine Learning on the data lake. They are also ideal for developing Business Intelligence solutions.

# GCP BigQuery

There is a lot to like about GCP BigQuery. It is serverless, it has integrated Machine Learning, and it is easy to use. This next section has a walkthrough of a k-means clustering tutorial.

The interface, when queried, intuitively gives back results. A key reason for this is the use of SQL and the direct integration with both Google Cloud and Google Data Studio
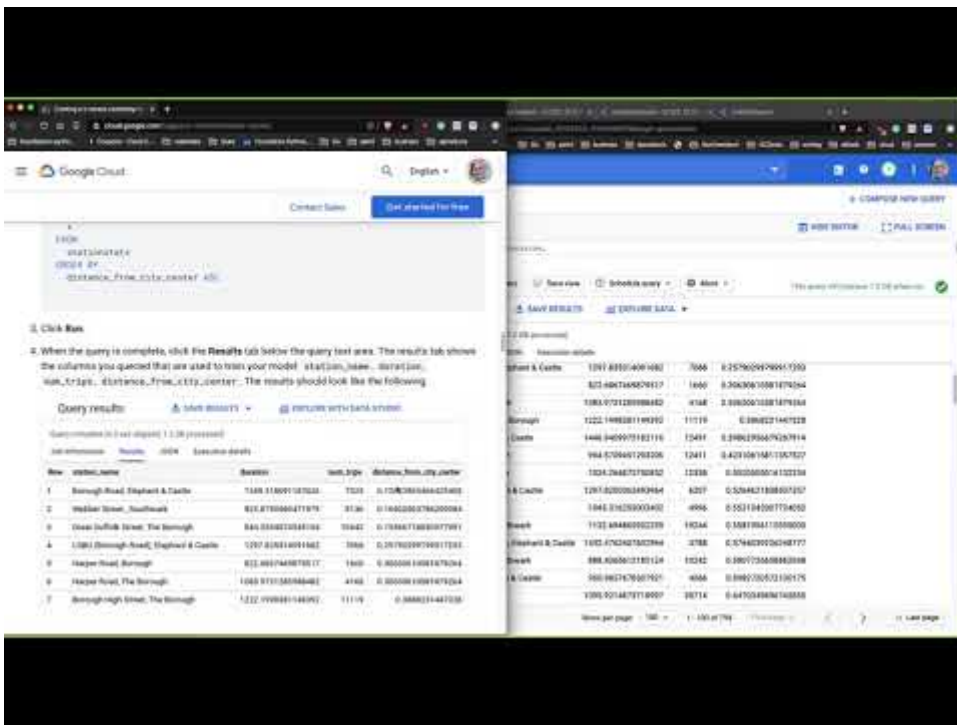


Learn to use Google BigQuery in the following screencast.

*Video Link: https://www.youtube.com/watch?v=eIec2DXqw3Q*

Even better, you can directly train Machine Learning models using a SQL statement. This workflow shows an emerging trend with Cloud Database services in that they let you both query the data and train the model. In this example, the `kmeans` section is where the magic happens.

```
CREATE OR REPLACE MODEL
  bqml_tutorial.london_station_clusters OPTIONS(model_type='kmeans',
    num_clusters=4) AS
WITH
  hs AS (
  SELECT
    h.start_station_name AS station_name,
    IF
      (EXTRACT(DAYOFWEEK
        FROM
          h.start_date) = 1
        OR EXTRACT(DAYOFWEEK
        FROM
          h.start_date) = 7,
        "weekend",
        "weekday") AS isweekday,
      h.duration,
      ST_DISTANCE(ST_GEOGPOINT(s.longitude,
          s.latitude),
        ST_GEOGPOINT(-0.1,
          51.5))/1000 AS distance_from_city_center
  FROM
    `bigquery-public-data.london_bicycles.cycle_hire` AS h
  JOIN
    `bigquery-public-data.london_bicycles.cycle_stations` AS s
  ON
    h.start_station_id = s.id
  WHERE
    h.start_date BETWEEN CAST('2015-01-01 00:00:00' AS TIMESTAMP)
    AND CAST('2016-01-01 00:00:00' AS TIMESTAMP) ),
```

```
  stationstats AS (
  SELECT
    station_name,
    isweekday,
    AVG(duration) AS duration,
    COUNT(duration) AS num_trips,
    MAX(distance_from_city_center) AS distance_from_city_center
  FROM
    hs
  GROUP BY
    station_name, isweekday)
SELECT
  * EXCEPT(station_name, isweekday)
FROM
  stationstats
```

Finally, when the k-means clustering model trains, the evaluation metrics appear as well in the console.



Often a meaningful final step is to take the result and then export it to their Business Intelligence (BI) tool, data studio.

The following is an excellent example of what a cluster visualization could look like in Google Big Query exported to Google Data Studio.

| cluster | MEDIAN_HOME_PRICE_COUNTY_... | VALUE_MILLIONS | ELO | WINNING_SEASON |
|---|---|---|---|---|
| 1. | 0 | | | | |
| 2. | 2 | | | | |
| 3. | 1 | | | | |

1 - 3 / 3  <  >

| | TEAM | TOTAL_ATTENDANCE_MILLIONS | COUNTY_POPULATION_MILLIONS |
|---|---|---|---|
| 3. | Golden State Warriors | | |
| 4. | Chicago Bulls | | |
| 5. | Boston Celtics | | |
| 6. | Los Angeles Clippers | | |
| 7. | Brooklyn Nets | | |
| 8. | Houston Rockets | | |
| 9. | Dallas Mavericks | | |
| 10. | Miami Heat | | |
| 11. | Cleveland Cavaliers | | |
| 12. | San Antonio Spurs | | |
| 13. | Toronto Raptors | | |
| 14. | Phoenix Suns | | |
| 15. | Sacramento Kings | | |
| 16. | Portland Trail Blazers | | |
| 17. | Oklahoma City Thunder | | |
| 18. | Washington Wizards | | |
| 19. | Orlando Magic | | |
| 20. | Utah Jazz | | |
| 21. | Detroit Pistons | | |

1 - 30 / 30  <  >

You can view the report using this direct URL.
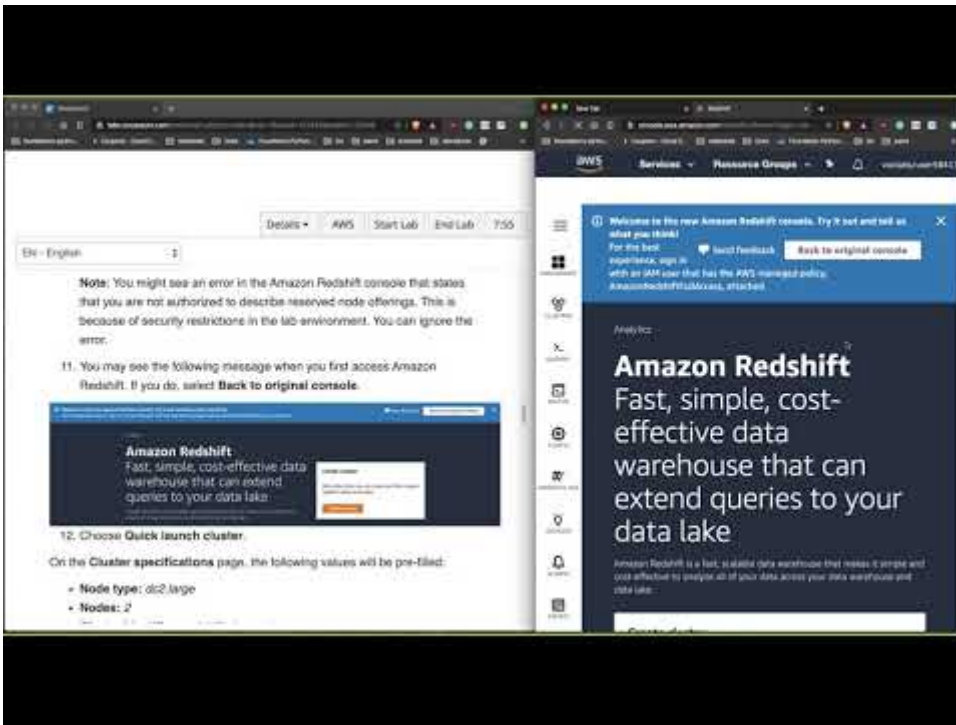
## Summary of GCP BigQuery

In a nutshell, GCP BigQuery is a useful tool for Data Science and Business Intelligence. Here are the key features.

- Serverless
- Large selection of Public Datasets
- Integrated Machine Learning
- Integration with Data Studio
- Intuitive
- SQL based

# AWS Redshift

AWS Redshift is a Cloud data warehouse designed by AWS. The key features of Redshift include the ability to query exabyte data in seconds through the columnar design. In practice, this means excellent performance regardless of the size of the data.

Learn to use AWS Redshift in the following screencast.

*Video Link: https://www.youtube.com/watch?v=vXSH24AJzrU*

## Key actions in a Redshift Workflow

In general, the key actions are as described in the Redshift getting started guide. These are the critical steps to setup a workflow.

- Cluster Setup

- IAM Role configuration (what can role do?)

- Setup Security Group (i.e. open port 5439)

- Setup Schema

```
create table users(
userid integer not null distkey sortkey,
username char(8),
```

- Copy data from S3

```
copy users from 's3://awssampledbuswest2/tickit/allusers_pipe.txt'
credentials 'aws_iam_role=<iam-role-arn>'
delimiter '|' region 'us-west-2';
```

- Query

```
SELECT firstname, lastname, total_quantity
FROM
(SELECT buyerid, sum(qtysold) total_quantity
```

```
FROM   sales
GROUP BY buyerid
ORDER BY total_quantity desc limit 10) Q, users
WHERE Q.buyerid = userid
ORDER BY Q.total_quantity desc;
```

## Summary of AWS Redshift

The high-level takeaway for AWS Redshift is the following.

- Mostly managed
- Deep Integration with AWS
- Columnar
- Competitor to Oracle and GCP Big Query
- Predictable performance on massive datasets

## Summary

This chapter covers storage, including object, block, filesystem, and Databases. A unique characteristic of Cloud Computing is the ability to use many tools at once to solve a problem. This advantageous trait is heavily at play with the topic of Cloud Storage and Cloud Databases.