

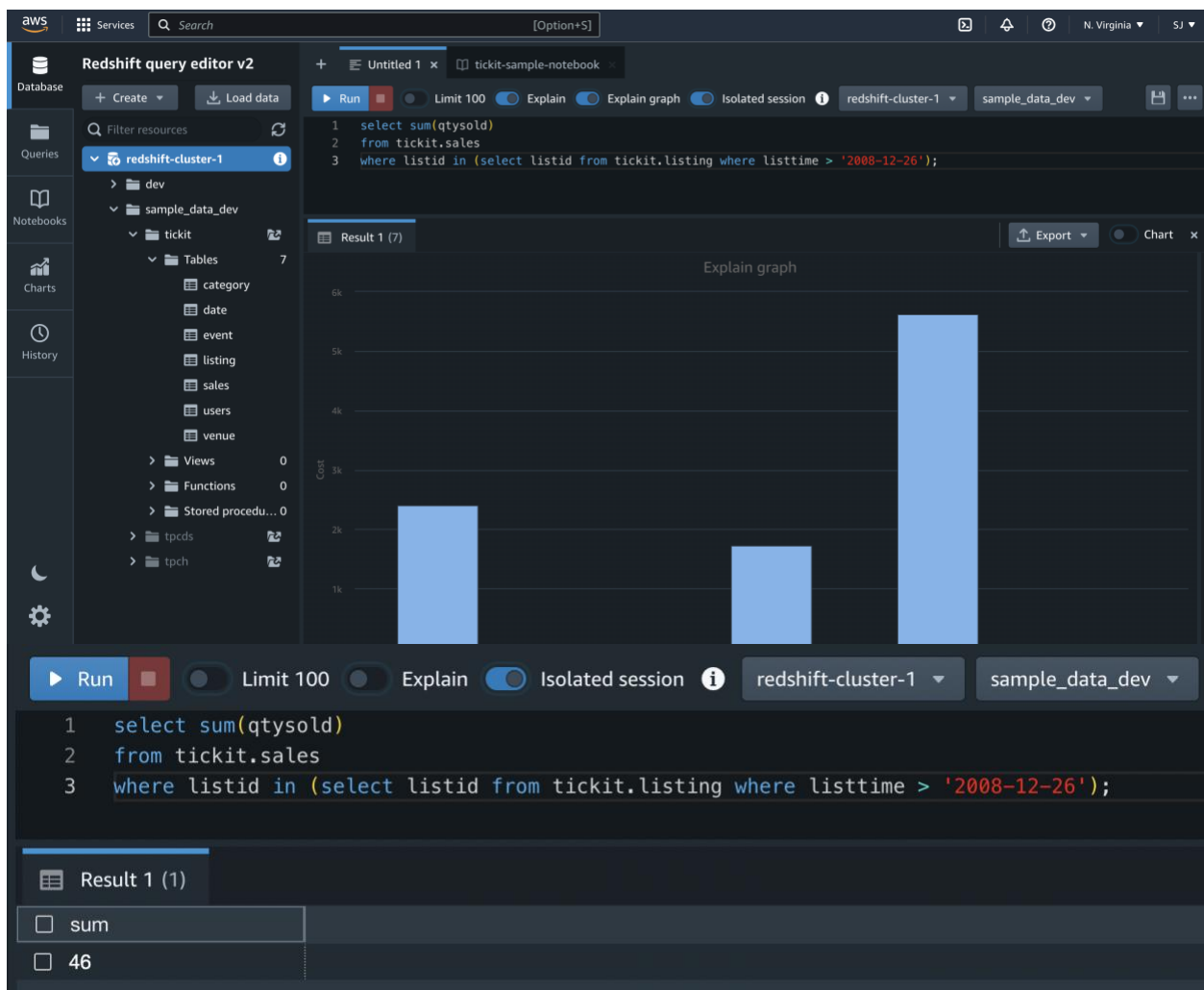
Cloud Data Warehousing

Part 1

Create an AWS Redshift Cluster with the tickitdb sample database.
Compare the costs of the following two queries

Query 1

```
select sum(qtysold)
from sales
where listid in (select listid from listing where listtime > '2008-12-26');
```



Step - 1

-> XN Seq Scan on listing (cost=0.00..2406.21 rows=1607 width=4)

Scans 1607 rows on "listing" table

Step - 2

-> XN Hash (cost=2406.21..2406.21 rows=1607 width=4)

A hash join and hash are used for inner joins, left outer joins, right outer joins. These operators are used when joining tables where the join columns aren't distribution keys and sort keys.

Step - 3

-> XN Seq Scan on sales (cost=0.00..1724.56 rows=172456 width=6)

Scans 172456 rows on "sales" table

Step - 4

Hash Cond: ("outer".listid = "inner".listid)

Hash condition: ("outer".listid = "inner".listid).

Step - 5

-> XN Hash IN Join DS_DIST_NONE (cost=2410.23..8022.29 rows=1441 width=2)

A hash join and hash are used for inner joins, left outer joins, right outer joins. These operators are used when joining tables where the join columns aren't distribution keys and sort keys.

Query 2

select sum(qtysold)

from sales

where saletime > '2008-12-26'

and listid in (select listid from listing where listtime > '2008-12-26');

Services

Search

[Option+S]

N. Virginia

SJ

Redshift query editor v2

Untitled 1 x

ticket-sample-notebook

Run

Limit 100

Explain

Explain graph

Isolated session

redshift-cluster-1

sample_data_dev

1 select sum(qtysold)

2 from tickit.sales

3 where saletime > '2008-12-26'


4 and listid in (select listid from tickit.listing where listtime > '2008-12-26');

Result 1 (8)

Export

Chart

Explain graph



Run

Limit 100

Explain

Isolated session

redshift-cluster-1

sample_data_dev

1 select sum(qtysold)

2 from tickit.sales

3 where saletime > '2008-12-26'

4 and listid in (select listid from tickit.listing where listtime > '2008-12-26');

Result 1 (1)

sum

46

Step - 1

-> XN Seq Scan on listing (cost=0.00..2406.21 rows=1607 width=4)

Scans 1607 rows on "listing" table

Step - 2

-> XN Hash (cost=2406.21..2406.21 rows=1607 width=4)

A hash join and hash are used for inner joins, left outer joins, right outer joins. These operators are used when joining tables where the join columns aren't distribution keys and sort keys.

Step - 3

-> XN Seq Scan on sales (cost=0.00..2155.70 rows=1794 width=6)

Scans 1794 rows on "sales" table

Step - 4

Hash Cond: ("outer".listid = "inner".listid)

Hash condition: ("outer".listid = "inner".listid).

Step - 5

-> XN Hash IN Join DS_DIST_NONE (cost=2410.23..4604.20 rows=15 width=2)

A hash join and hash are used for inner joins, left outer joins, right outer joins. These operators are used when joining tables where the join columns aren't distribution keys and sort keys.

1st Query is giving us the total quantity of items sold which were listed after 2008-12-26

2nd Query is giving us the total quantity of items which are listed and sold after 2008-12-26

We will get the same output with both the queries because items will be sold only after listing them first, so all the items listed after 2008-12-26 must be sold after 2008-12-26.

Total real cost of 1st query: 2406 + 1724 + 5612 = 9742

Total real cost of 2nd query: 2406 + 2155 + 2194 = 6755

So even if both the queries will give the same output 2nd query is more efficient because as we can see from the steps explained in the above screenshots for 2nd query, we are scanning only 1794 rows of the sales table as compared to scanning 172465 rows of sales table in case of 1st query, and these filtered rows from sales table are compared with the output of query on listing table instead of comparing all the rows from sales table as done in 1st query.

Part 2

Create a BigQuery project and dataset with listing.csv and sales.csv tables

Compare the costs of the following two queries

Query 1

select sum(qtysold)

from `msba-369223.BAN622.sales` as S1, `msba-369223.BAN622.listings` as L1

where S1.listid = L1.listid

and L1.listtime > '2008-12-26';

Google Cloud MSBA Search Products, resources, docs (/)

SANDBOX Set up billing to upgrade to the full BigQuery experience. [Learn more](#) DISMISS UPGRADE

Explorer + ADD DATA

Viewing all resources. [Show starred resources only.](#)

- msba-369223
 - External connections
 - BAN622
 - listings
 - sales

```

1 select sum(S1.qty sold)
2 from `msba-369223.BAN622.sales` as S1
3 where S1.saletime > "2008-12-25"
4 and S1.listid in (select L1.listid from `msba-369223.BAN622.listings` as L1
5 where L1.listtime > "2008-12-26")

```

Query results

JOB INFORMATION RESULTS JSON EXECUTION DETAILS EXECUTION GRAPH PREVIEW

Job ID	msba-369223:us-west2.bqjob_5ee054ae_1849794a7b9
User	swan11196@gmail.com
Location	us-west2
Creation time	Nov 20, 2022, 4:26:52 PM UTC-8
Start time	Nov 20, 2022, 4:26:53 PM UTC-8
End time	Nov 20, 2022, 4:26:53 PM UTC-8
Duration	0 sec
Bytes processed	6.88 MB
Bytes billed	20 MB
Job priority	INTERACTIVE
Use legacy SQL	false
Destination table	Temporary table

msba-369223.BAN622.li... table

192,500

S00: Input

msba-369223.BAN622.s... table

202

172,664

S01: Join+

1

S02: Output

Query results

JOB INFORMATION RESULTS JSON EXECUTION DETAILS EXECUTION GRAPH PREVIEW

Row	f0_
1	46

1st Query is giving us the total quantity of items sold which were listed after 2008-12-26

2nd Query is giving us the total quantity of items which are sold after 2008-12-25 and listed after 2008-12-26

[illegible]