

Week 3 Data Warehouse

- Data warehouse.
- Partitioning and clustering.
- Best practices.
- Internals Big Query.
- BigQuery Machine learning.
- Deployment of BQ ML Models.
- Doing tasks using Airflow

| | OLTP | OLAP |
|---------------------|--------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------|
| Purpose | - Online Transaction Processing - Control and run essential business operations in <u>real time</u> | Online Analytical Processing <u>plan, solve problems</u> <u>support Decisions, Discover insights</u> |
| Data update | <u>short, fast updates initiated by user.</u> | <u>Data periodically refreshed with scheduled, long run batch job</u> |
| DB design | <u>Normalized DB for efficiency for fast response with user</u> | <u>Denormalized DB for analysis</u> |
| Size | <u>Generally small if historical data archived</u> | <u>Generally large due to data aggregating large DB</u> |
| Backup productivity | <u>Regular Backup happens increase productivity user</u> | <u>lost data can be loaded from OLTP</u> <u>increase productivity for business manager</u> |

Data warehouse (DW)

- DW is an OLAP solution meant for reporting and data analysis, unlike OL which flow ETL
DW Flows ETL
- DW receives data from different data sources, which then processed in staging areas before being ingested to actual warehouse (DB), DW may then feed data to separate Data Marts "smaller DB which end users can use"

Big Query (BQ)

- BQ is a DW solution offered by GCP.
- BQ is different than GCS
 - BQ is serverless: no server to manage or DB software to install
 - BQ is scalable and high availability: Google take care of underlying software and infrastructure.
- BQ has Built in features:
 - Machine learning
 - Geospatial analysis
 - Business intelligence
- BQ maximize flexibility by separating data to different engines.



Alternative

AWS Redshift, Azure Synapse

Pricing

GGP charges you in most of features like:

- processing
- storage $.02 \$/GB/month$
- ingestion
- extraction

Data processing:

- on demand [Default]: $5 \$/TB/month$ [first TB free]
- flat rate: $2000 \$/month$
cost is not with data, But with slots (CPU)
minimum 100 slots for $2k \$/M$.
- But if you run a lot of queries and 100 slot
not enough queries must wait till free slot.
on demand doesn't have this issue as you
deal with storage not slots.
- flat rate limited with queries due slots.
on demand limited with storage by process.
- when we run query on BQ on the top right
estimated values for [cost & time appear]

External tables:

BQ supports external data sources, you may query these sources directly from BQ, even data isn't in BQ.



Partitions

BQ tables can be partitioned into multiple smaller tables, Filter based on data, so we can query a specific sub-table based on data we are interested in.

- partition improves: - performance - reduce cost
- tables are partitioned based on
 - Time Stamp.
 - Data
 - Any numerical values.
 - partition may be [daily, monthly, yearly]
as a value for partition default daily.
- BQ limit number of partitions to 4000/table
 - Example run Query on partitioned and non partitioned
 - partitioned 106 MB of Data to process
 - Non partitioned 1600 MB // //
 - only one column cond used

Clustering

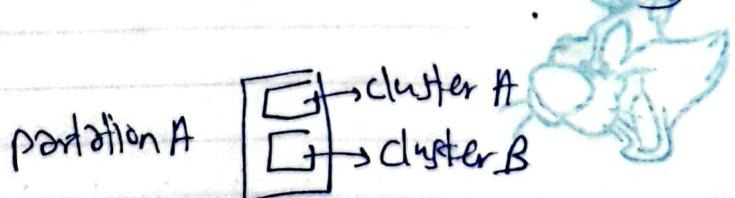
Rearranging table based on its values ex categorical we can cluster multiple column till (4).

- Values of categorical can be any thing string or even numerical.
- Tables with less than 1 GB do not show significant improvement, It may cost and increase time!



| <u>clustering</u> | <u>partitioning</u> |
|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| cost benefit unknown BQ Can't estimate cost for example we may have query with clusters from 4 columns each one could have $\frac{1}{x}$ clusters so time consuming to estimate. | cost known up front. Estimate cost before run "for ex we know that we have only <u>one column</u> and limit of <u>4000 partitions</u> so limited to estimate. |
| High granularity multiple criteria can used to sort | low granularity only single column, often time. |
| clustered fixed place | partitioned are flexible can be added. |
| Good for queries that uses filters or aggregation (ex: \sum) | Good when filter based on single column. |

- use of clustering if:
 - partitions > 4000 .
 - partition results small amount of data/partition.
- BQ has automatic re [clustering, partitioning] when we add, delete data.
- For partitioned data we may have clustering inside partition.



Best practices

1. Cost reduction

- Avoid SELECT * as reduce amount of processed data would reduce cost.
- price your queries.
- use clustered and partitioned.
- use streaming inserts with caution.
- Materialize query results. "Memorize it" Cache

2. Query performance

- Filter on partitioned columns.
- Denormalize data....?
 - Normalization is very important to remove any redundancy in the data.
 - Make relation between data.
 - store data in an atomic form [First, Mid, last]
- So why we do Denormalization ---?
 - Normalized data uses Joins to keep relationship between columns, but Joins takes a lot of resources and reduce performance of query.

e.g Normalized : Bringing Breakfast, lunch, dinner from supermarket, this supermarket store each element alone, so you start with breakfast and go through all element section and so on.

Denormalized : Bringing some thing, but here supermarket sections for meal section for lunch and so on, so you won't go to all sections.



→ missing some concepts in tricks

Date: _____ No. _____

For this example

- user would only take few minutes as every thing is organised by meal.
- we may have some objects repeated on different places like bread on [dinner, breakfast]

For supermarket

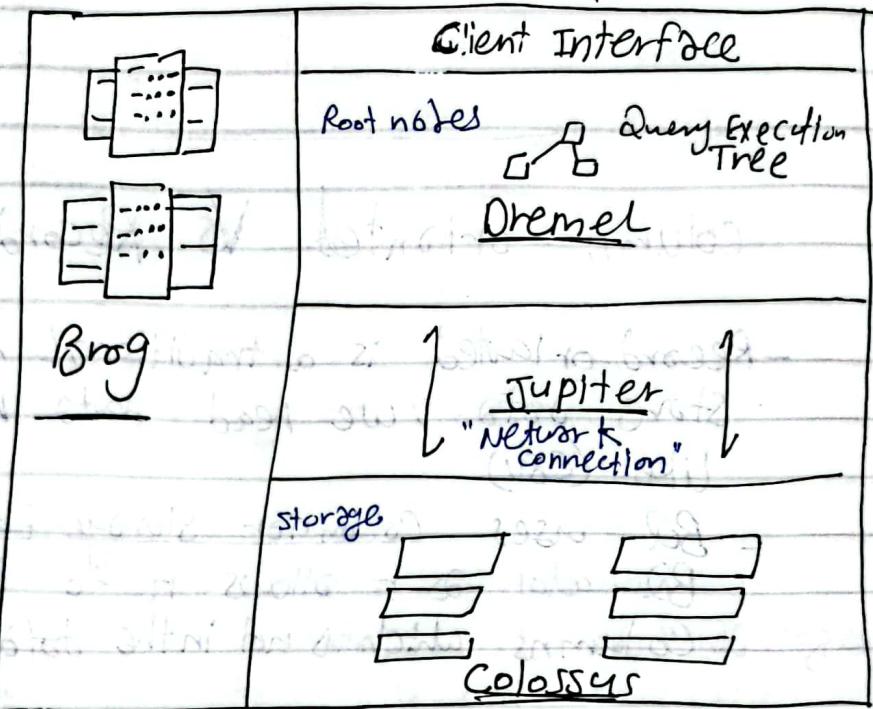
- that would reduce their speed and performance as we have redundancy bread repeated on different places.

- use nested or repeated columns
- External resources reading from bucket affect performance badly.
- Reduce data before Joins.
- • Don't threat WITH clauses as prepared statement.
- Avoid oversharding tables ["Storing data in multiple tables"]
- Avoid Java Script user-defined functions.
- use approximate aggregation function.
- order statement should be last part.
- start with largest tabel and small ---.



BigQuery Architecture

Date _____ No. _____
 clients - web UI ...



Colossus: - Google global storage with speed 1TB/s using Jupiter.
 - Columnar storage, format & compression algorithms to store data.
 - Handle replication, replace and recovery and disrupted management.

Jupiter: - Network to connect Dremel and Colossus in house network to interconnect database.

Dremel: [Brain]
 - compute part of BQ - SQL query to tree leaves of tree slots branch mixers
 - slots reading data from storage
 - mixers perform aggregation.
 - Dremel split slots to queries as needed.



- Brog : - orchestration solution
- Brog is part of kubernetes.

column - orianted Vs Record - orianted

- Record oriented is a traditional method to store data , we read data row by row like (csv)
- BQ uses Columnar storage Format , this help BQ a lot as it allows it to discard other columns which is not in the data.
- Dremel modify queries to an execution tree parts of queries assigned to mixers then to even small parts to slots to access colossus.
- column storage allows fast retrieval from colossus which return to mixer to perform aggregation.



ML with BQ

Date: _____

No: _____

BQ-ML allows:

- Run and execute machine learning models.
- All of that using SQL only without of any experience with python.
- No need for experts just run SQL then:
 - Feature Engineering Auto.
 - Train of the model.
 - Export & Deploy of model.
- features like:
 - predict values - linear regression
 - - AutoML, DNN, wide regressor
 - Identify spam emails - logistic regression
 - - boosted Trees, DNN
 - Generate recommendation.
 - Reduce data dimensionality.
 - Find anomalies [Detect fraud, credit risk]
 - Group data to clusters [customer segmentation]
 - Time series forecasting [ARIMA]
- we will creat new table which contain certain columns casted.

CREATE OR REPLACE Model 'taxi-ridges-ny.tip-model'
options(

model-type = 'linear-regression'

input-label-cols = ['tip-amount']

DATA-SPLIT-METHOD = 'AUTO-SPLIT') AS

SELECT * FROM data set



- ③ Q show all details of model like:
- Graphs of train.
 - Evaluation metrics.
 - Model Graph.

• SELECT * FROM M.L. FEATURE-INFO
(MODEL 'taxi — .tip-model');

like describle

- M.L. FEATURE
- M.L. EVALUATE
- M.L. PREDICT
- M.L. Explain - predict
 - Return prediction with column n
 - preceded = 1.321 → ^{first} drop of
struct(3 as top-k-feature)

• Hyper parameter is available



BQ ML Deployment:

- we will Export model use Tensorflow serving which is Docker img and use HTTP to get prediction

• steps :

- 1- Authenticate GCP
- 2- Export model to Cloud Storage Bucket.
- 3- Download Exported model.
- 4- Run Docker image.
- 5- with image working, run prediction

Curl \ Features

```
-d '{ "instances": [ { "passenger_count": 1, "trip_distance": 12.2 } ] }' \ 
-X POST http://localhost:8501/v1/models/taxi-model:predict
```

returned value is prediction.

Integration BQ with Airflow

- we will use Airflow to automate process of download files and process it each certain time using Cron Job

- 1- GCS - 2 - GCS [Moving Files between folders]
- 2- GCS - 2 - bq-ext [GCS to BQ Extend]
3. bq - 2 - partitioned

- we will use Google-provided operators



we extend data in GCS for easier processing.

dated Nov 2018

it adds some pixels to the raw
with more info

: 768x256

400 at 1000x1000

original with 1000x1000
labelling added 1000x1000

1000x1000

noticing a lot of noise

but with the 20 different

class structures it makes no noise.
does it uses less pixel because of
big and tiny with noise

extended the original 768x256 to
1000x1000 to add some noise

noticing a lot of noise

