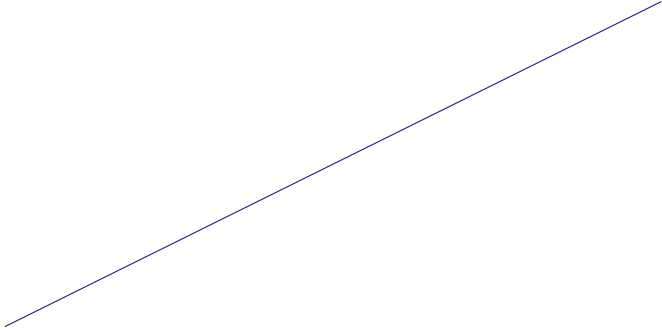




# ***Term Project***

# ***Google BigQuery***



201736055 한지안  
201833299 허서윤

- 1. Specs of BigQuery**
- 2. Key Features of BigQuery**
- 3. Pros and Cons of BigQuery**
- 4. Implementation of BigQuery**
- 5. References**

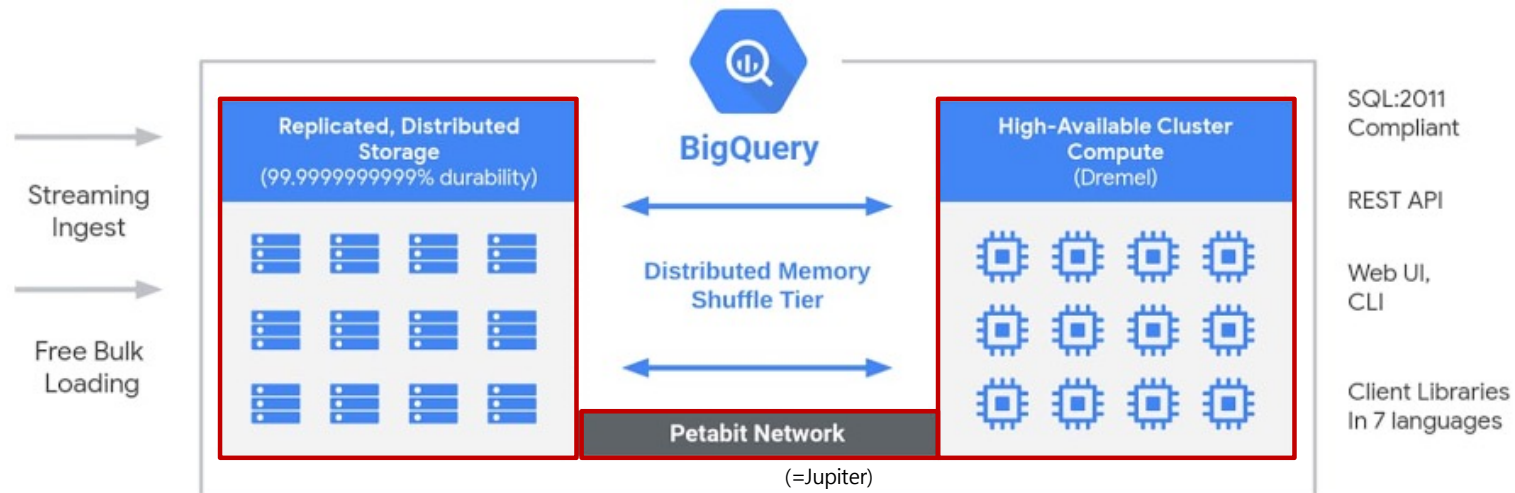
## *What is BigQuery?*

---



**Serverless** cloud service for users outside Google  
based on the **Dremel** project

## Architecture of BigQuery

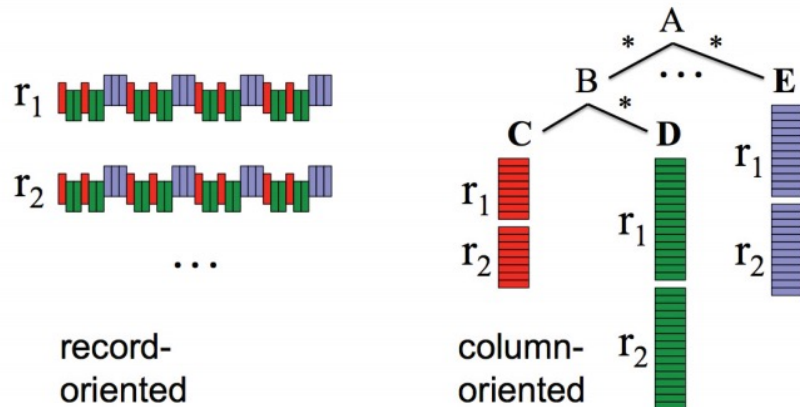


## Architecture of BigQuery



## Architecture of BigQuery

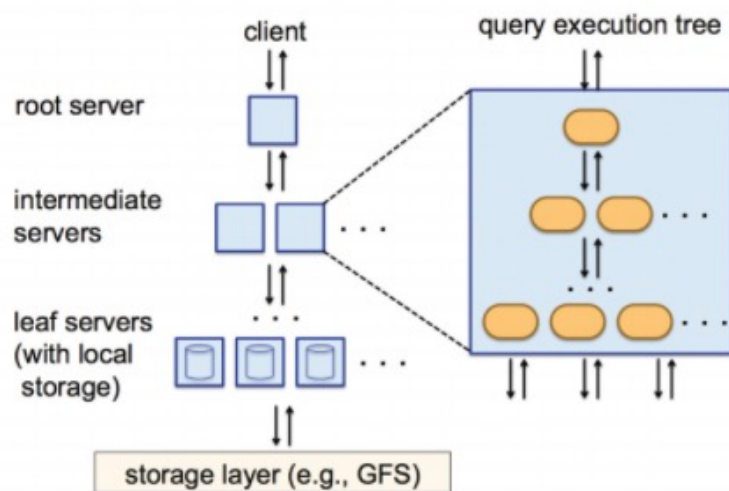
### ① Columnar Storage



- Save the table in an optimized column format
- Same data type data are gathered and stored
- Each table is compressed and encrypted on a disk
- Storage is durable, and each table replicates between data centers

## Architecture of BigQuery

### ② Tree Architecture Distribution



#### ❑ Root server

- Receives SQL query from the user
- Splits it into a small SQL query

#### ❑ Intermediate server

- Query split and sent to leaf node (=slot)
- Gets through columnar storage process.

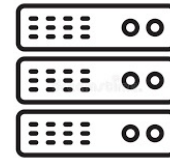
#### ❑ Leaf server

- Reads data from a storage layer
- Delivers results to the parent node.

## Key Features

### ① Serverless

Usually, there is no server, and computer resources are used only for analysis  
Only need to pay at this time.



### ② Cloud

No need to install/operate as a service.  
Large capacity support, fast performance support, low price

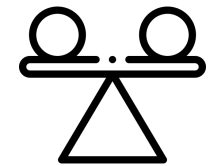




## Key Features

### ③ Stability

The risk of data loss is low because three copies are distributed and stored in different data centers



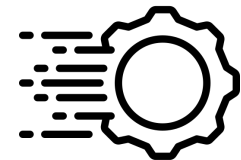
### ④ Batch & Streaming

It provides a batch that loads data at once, and a streaming function that allows you to input data in real time.



### ⑤ Dermal project

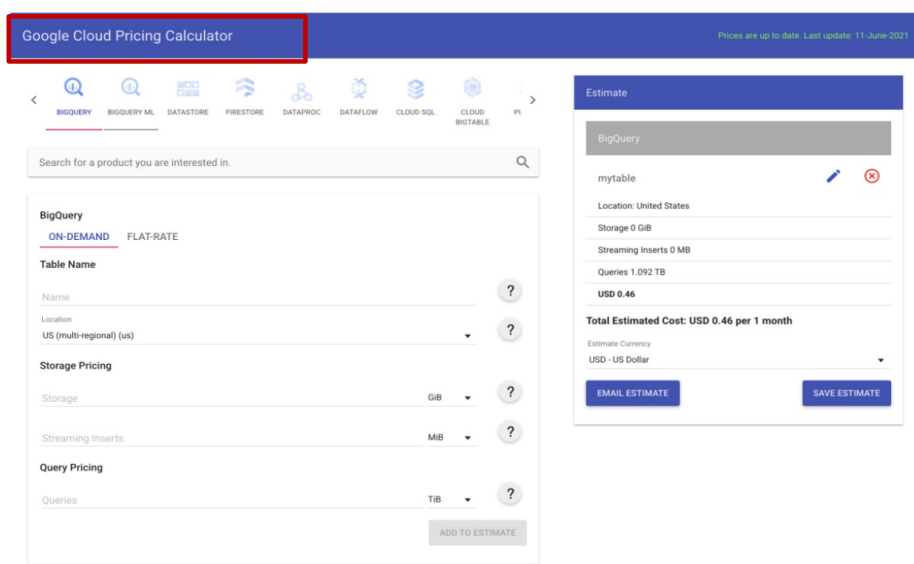
Structured data can be distributed and analyzed quickly



### Pros of BigQuery

① Cost-effective because it is possible to estimate query costs

② Models can be created & tested using SQL queries



The screenshot shows the Google Cloud Pricing Calculator interface. The top navigation bar includes links for BigQuery, BigQuery ML, Datastore, Firestore, DataProc, Dataflow, Cloud SQL, and Cloud Bigtable. The main section is titled "BigQuery" and has two tabs: "ON-DEMAND" (selected) and "FLAT-RATE". Under "Table Name", there are fields for "Name" and "Location" (set to "US (multi-regional) (us)"). Under "Storage Pricing", there are fields for "Storage" (set to "GiB") and "Streaming Inserts" (set to "MiB"). Under "Query Pricing", there is a field for "Queries" (set to "TiB"). A search bar at the top left says "Search for a product you are interested in." On the right, the "Estimate" section shows a summary for "mytable": Location: United States, Storage: 0 GiB, Streaming Inserts: 0 MB, Queries: 1.092 TiB, and a "Total Estimated Cost: USD 0.46 per 1 month". There are buttons for "EMAIL ESTIMATE" and "SAVE ESTIMATE".



### Cons of BigQuery

Specialized in analysis and OLAP  
Not suitable for OLTP

쿼리 편집기

```
1 insert into professional-data-engineering.bike_insight.sample (name) values ('hello');
```

실행 쿼리 저장 보기 저장 쿼리 예약 더보기

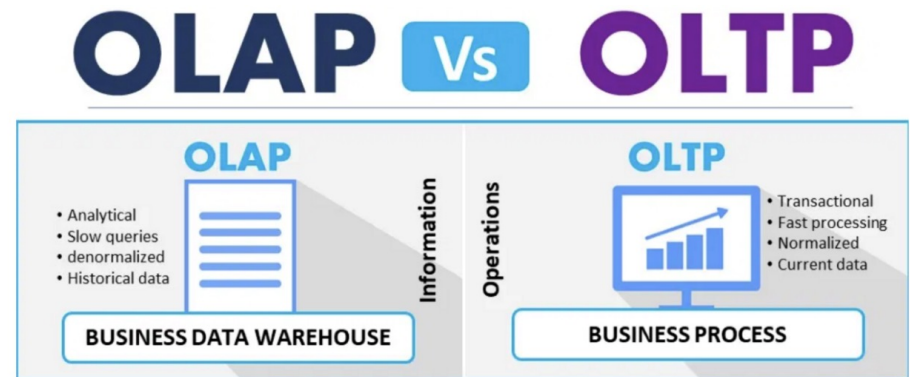
쿼리 결과

쿼리 완료(0.9초 경과, 0B 처리됨)

작업 정보 결과 JSON 실행 세부정보

쿼리 디버깅이나 최적화에 도움이 필요한 경우 문서를 확인하세요. [자세히 알아보기](#)

경과 시간	사용한 슬롯 시간
0.9초	3.611초



## Establishing and Using Classification Model in Census Data



### US Census Data

United States Census Bureau

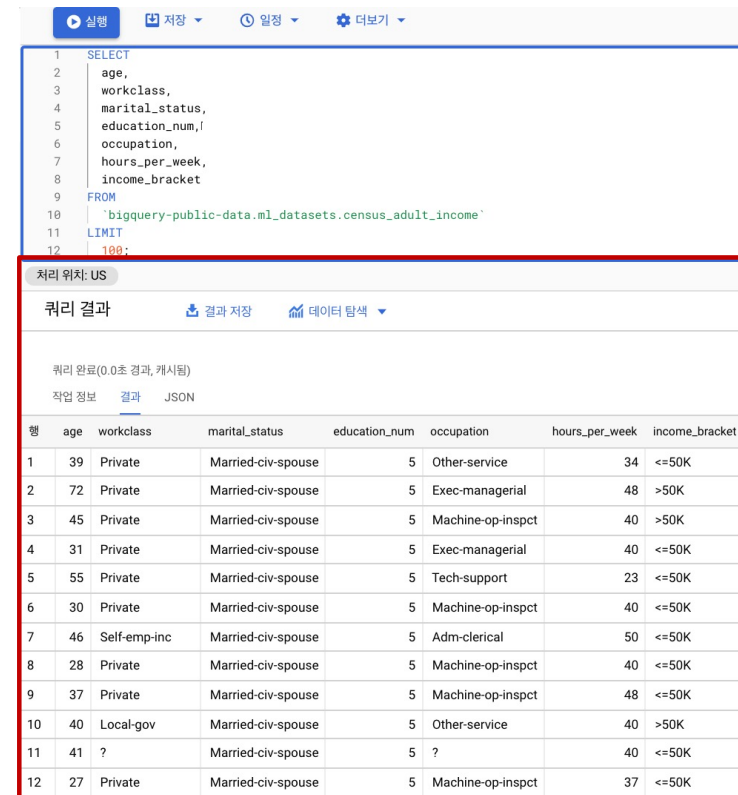
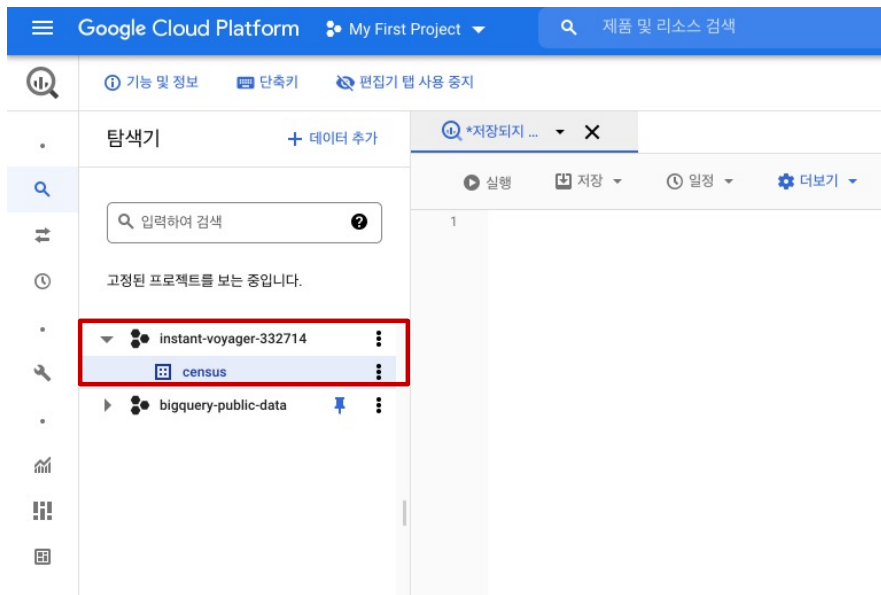
2000 and 2010 US Census data

Data Information
Age
Workclass
Marital_Status,
Education_num
Occupation
Hours_per_week
Income_bracket

## Implementation of BigQuery

1. Big query data sets are created to store models.

2. Returns 100 rows from a dataset



### 3. Create a view to compile training data

```
1 CREATE OR REPLACE VIEW
2   `census.input_view` AS
3 SELECT
4   age,
5   workclass,
6   marital_status,
7   education_num,
8   occupation,
9   hours_per_week,
10  income_bracket,
11  CASE
12    WHEN MOD(functional_weight, 10) < 8 THEN 'training'
13    WHEN MOD(functional_weight, 10) = 8 THEN 'evaluation'
14    WHEN MOD(functional_weight, 10) = 9 THEN 'prediction'
15  END AS dataframe
16 FROM
17   `bigquery-public-data.ml_datasets.census_adult_income`
```

#### 뷰 스키마

필터 속성 이름 또는 값 입력		
필드 이름	유형	모드
age	INTEGER	NULLABLE
workclass	STRING	NULLABLE
marital_status	STRING	NULLABLE
education_num	INTEGER	NULLABLE
occupation	STRING	NULLABLE
hours_per_week	INTEGER	NULLABLE
income_bracket	STRING	NULLABLE
dataframe	STRING	NULLABLE

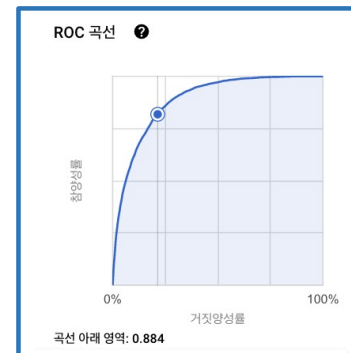
### 4. Create a logistic regression model

```
1 CREATE OR REPLACE MODEL
2   `census.census_model`
3 OPTIONS
4   ( model_type='LOGISTIC_REG',
5     auto_class_weights=TRUE,
6     input_label_cols=['income_bracket']
7   ) AS
8 SELECT
9   * EXCEPT(dataframe)
10 FROM
11   `census.input_view`
12 WHERE
13   dataframe = 'training'
```

## Implementation of BigQuery

### 5. Use the ML.EVALUATE function to evaluate the performance of the model

```
1 SELECT
2   *
3 FROM
4   ML.EVALUATE (MODEL `census.census_model`,
5   (
6     SELECT
7       *
8     FROM
9       `census.input_view`
10    WHERE
11      dataframe = 'evaluation'
12   )
13 )
```



혼동 행렬

이 테이블은 모델이 각 라벨을 올바르게 분류한 빈도(파란색) 및 가장 자주 혼동된 라벨(회색)을 보여줍니다.

	예측된 라벨	>50K	<=50K
True 라벨			
>50K		82%	18%
<=50K		21%	79%

쿼리 결과

쿼리 완료(0.9초 경과, 2.6MB 처리됨)

작업 정보 결과 JSON 실행 세부정보

행	precision	recall	accuracy	f1_score	log_loss	roc_auc
1	0.5740551583248212	0.7635869565217391	0.8168009919404836	0.6553935860058309	0.3944068834574654	0.8827432567432567

### 6. Predict the income class of all respondents

```
1 SELECT
2   *
3 FROM
4   ML.PREDICT (MODEL `census.census_model`,
5   (
6     SELECT
7       *
8     FROM
9       `census.input_view`
10    WHERE
11      dataframe = 'prediction'
12   )
13 )
```

작업 정보   결과   JSON   실행 세부정보

행	predicted_income_bracket	predicted_income_bracket_probs.label	predicted_income_bracket_probs.prob	age	workclass	marital_status	education_num	occupation	hours_per_week	income_bracket	dataframe
1	<=50K	>50K	0.05639289147442344	34	?	Married-civ-spouse	7	?	8	<=50K	prediction
		<=50K	0.9436071085255766								
2	<=50K	>50K	0.1311398392666556	21	?	Married-civ-spouse	9	?	30	<=50K	prediction
		<=50K	0.8688601607333444								
3	<=50K	>50K	0.06491298402568849	25	?	Married-civ-spouse	9	?	4	<=50K	prediction
		<=50K	0.9350870159743115								

### 7. More detailed analysis with Explainable AI method

```
1 #standardSQL
2 SELECT
3   *
4 FROM
5   ML.EXPLAIN_PREDICT(MODEL `census.census_model`,
6   (
7     SELECT
8       *
9     FROM
10      `census.input_view`
11    WHERE
12      dataframe = 'evaluation'),
13   STRUCT(3 as top_k_features))
```



## 6. Predict the income class of all respondents

```
1 SELECT
2   *
3 FROM
4   ML.PREDICT (MODEL `census.census_model`
5   (
6     SELECT
7       *
8     FROM
9       `census.input_view`
10    WHERE
11      dataframe = 'prediction'
12   )
13 )
```

## 7. More detailed analysis with Explainable AI method

실행

저장

일정

더보기

```
1 #standardSQL
2 SELECT
3 *
4 FROM
5 ML.EXPLAIN_PREDICT(MODEL `census.census_model`,
6 (
7   SELECT
8     *
9   FROM
10    `census.input_view`
11   WHERE
12     dataframe = 'evaluation'),
13   STRUCT(3 as top_k_features))
```

[illegible]

## ***Reference***

---

### Specs of BigQuery

<https://syujisu.tistory.com/190?category=907377>

<https://youtu.be/LhksTFvVriU>

### Architecture of BigQuery

<https://cloud.google.com/bigquery>

<https://velog.io/@jch9537/%ED%95%9C-%EC%A4%84-%EC%9A%A9%EC%96%B4%EB%B0%B0%EC%B9%98Batch%EB%9E%80>

### Pros of BigQuery

<https://xo.xello.com.au/blog/google-bigquery-5-benefits-cloud-data-warehouse>

<https://www.quora.com/What-are-the-pros-and-cons-of-using-Google-BigQuery-as-a-database>

<https://www.xplenty.com/blog/snowflake-vs-bigquery/>

### Cons of BigQuery

<https://www3.technologyevaluation.com/solutions/53566/google-bigquery>

<https://dzone.com/articles/introduction-to-google-bigquery>

### Big Query Implementation

<https://cloud.google.com/bigquery-ml/docs/logistic-regression-prediction>

### Data Source Site

<https://console.cloud.google.com/marketplace/product/united-states-census-bureau/us-census-data?project=instant-voyager-332714>



***Thank You***

