06b_customer_churn_console_execution.md

# Customer Churn using Serverless Spark through Google Cloud Console

**Goal** - Data Preparation and Model Training for Detecting Customer Churn.

Following are the lab modules:

1. Understanding Data
2. Solution Architecture
3. Parameter Requirements for the lab
4. Data Preparation
5. Model Training and Testing
6. Model Evaluation
7. Logging

## 1. Understanding Data

The dataset used for this project is customer churn data and customer test data..

The dataset contains the following features:

- Churn - Binary field which represents customers who left/were retained within the last month
- Services that each customer has signed up for – phone, multiple lines, internet, online security, online backup, device protection, tech support, and streaming TV and movies
- Customer account information – how long they've been a customer, contract, payment method, paperless billing, monthly charges, and total charges
- Demographic info about customers – gender, age range, and if they have partners and dependents

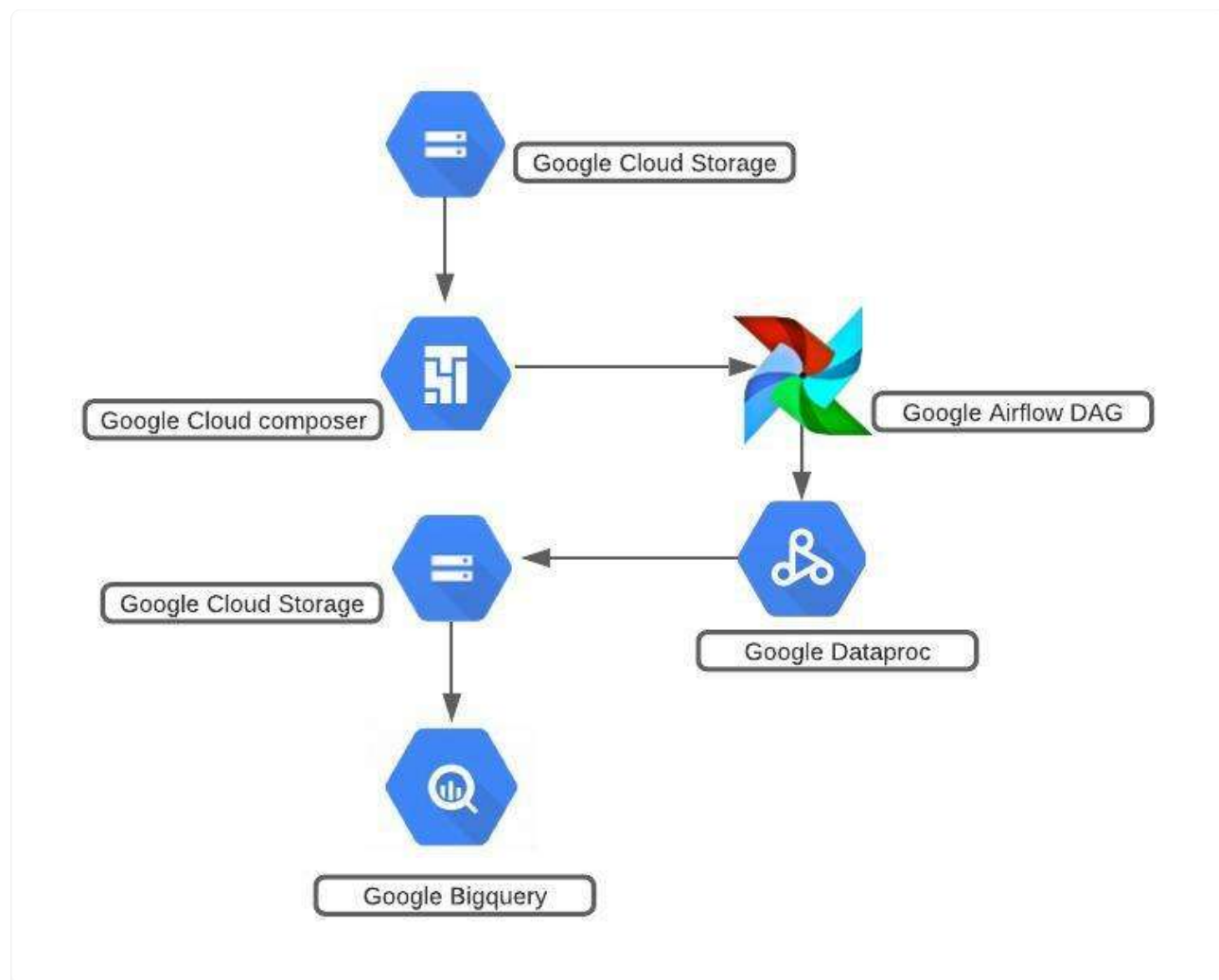**Note:** The following features refer to these same-host connections.

- serror_rate

- rerror_rate
- same_srv_rate
- diff_srv_rate
- srv_count

**Note:** The following features refer to these same-service connections.

- srv_serror_rate
- srv_rerror_rate
- srv_diff_host_rate

# 2. Solution Architecture



**Model Pipeline**

The model pipeline involves the following steps:

- Data cleanup and preparation
- Building and training a Machine Learning Model (Random Forest Classifier) before saving it into a GCS bucket
- Using the model built in above step to evaluate test data

# 3. Parameters required for the lab

Keep the following details handy for configuring the serverless batch jobs:

```
PROJECT_ID=                                      #current GCP project where we ar
REGION=                                          #GCP region where all our resour
BQ_DATASET_NAME=                                 #BigQuery dataset where all the
BUCKET_CODE=                                      #GCP bucket where our code, data
HISTORY_SERVER_NAME=spark-phs                    #name of the history server which
VPC_NAME=                                        # Primary VPC containing the subn
SUBNET=                                          #subnet which has private google
UMSA=serverless-spark                            #name of the user managed servic
SERVICE_ACCOUNT=$UMSA@$PROJECT_ID.iam.gserviceaccount.com
NAME=                                            #Your unique identifier
```

**Note:** The values for all the above parameters will be provided by the admin team.

# 4. Data Preparation

Based on EDA, the data preparation script has been created. Among the 21 columns, relevant features have been selected and stored in BQ for the next step of model training.

## 4.1. Create a new batch

Navigate to Dataproc > Serverless > Batches and click on **+CREATE**

## 4.2. Provide the details for the batch

Next, fill in the following values in the batch creation window as shown in the images below:

- **Batch ID** - A unique identifier for your batch

- **Region** - The region name provided by the Admin team

- **Batch Type** - PySpark

- **Main Python File** - gs://<your_code_bucket_name>/customer_churn/00-scripts/customer_churn_data_prep.py

- **JAR Files** - gs://spark-lib/bigquery/spark-bigquery-with-dependencies_2.12-0.22.2.jar

- **Arguments** -
  Four Arguments needs to be provided.

  - <your_project_id>
  - <your_dataset_name>
  - <your_code_bucket_name>
  - <your_name> **Note:** Press RETURN after each argument

- **Service Account** - <UMSA_NAME>@<PROJECT_ID>.iam.gserviceaccount.com

- **Network Configuration** - select the network and subnetwork with Private Google Access Enabled Run PySpark Serverless Batch for Data Preparation

- **History Server Cluster** - <your_phs_cluster_name>

## Batch info

Batch ID *

batch-088e

Region *

us-central1 ▼

## Container

Batch type *

PySpark ▼

Main python file *

gs://ch.  .... 1/customer_churn/00-scripts/bq_code_files/code_files_bigquery_code

Can be a GCS file with the gs:// prefix, an HDFS file on the cluster with the hdfs:// prefix, or
a local file on the cluster with the file:// prefix

Additional python files

Custom container image

Specify a custom container image to add Java or Python dependencies not provided by the
default container image. You must host your custom container on Container Registry .

Jar files

gs://spark-lib/bigquery/spark-bigquery-with-dependencies_2.12-0.22.2.jar ⊗

## 4.3. Submit the Serverless batch

Once all the details are in, you can submit the batch. As the batch starts, you can see the execution details and logs on the console.

## 4.4. Check the output table in BQ

Navigate to BigQuery Console, and check the **customer_churn_lab** dataset.
Once the data preparation batch is completed, new tables
'<your_name_here>_training_data' and '<your_name_here>_test_data' will be created.

To view the data in these tables -

- Select the table from BigQuery Explorer by navigating 'project_id' > 'dataset' > 'table_name'
- Click on the **Preview** button to see the data in the table



**Note:** If the **Preview** button is not visible, run the below queries to view the data. However, these queries will be charged for the full table scan.

```
SELECT * FROM `<project_name>.<dataset_name>.<your_name_here>_training_data` LIMIT
```



```
SELECT * FROM `<project_name>.<dataset_name>.<your_name_here>_test_data` LIMIT 100
```

<your_name_here>: _test_data table



**Note:** Edit all occurrences of <project_name> and <dataset_name> to match the values of the variables PROJECT_ID, and BQ_DATASET_NAME respectively

# 5. Model Training and Testing

Repeat the same steps as above to submit another batch for model training.

## 5.1. Create a new batch

Navigate to Dataproc > Serverless > Batches and click on **+CREATE**



## 5.2. Provide the details for the batch

Next, fill in the following values in the batch creation window as shown in the images below:

- **Batch ID** - A unique identifier for your batch
- **Region** - The region name provided by the Admin team
- **Batch Type** - PySpark
- **Main Python File** - gs://<your_code_bucket_name>/customer_churn/00-scripts/customer_churn_model_building.py
- **JAR Files** - gs://spark-lib/bigquery/spark-bigquery-with-dependencies_2.12-0.22.2.jar
- **Arguments** -
  Four Arguments needs to be provided.
    - <your_project_id>
    - <your_dataset_name>
    - <your_code_bucket_name>
    - <your_name> **Note:** Press RETURN after each argument
- **Service Account** - <UMSA_NAME>@<PROJECT_ID>.iam.gserviceaccount.com
- **Network Configuration -** select the network and subnetwork with Private Google Access Enabled Run PySpark Serverless Batch for Data Preparation
- **History Server Cluster -** <your_phs_cluster_name>

## 5.3. Query the model_test results BQ table

Navigate to BigQuery Console, and check the **customer_churn_lab** dataset.
Once the modeling batch is completed, a new table '<your_name_here>_predictions_data'
will be created.

To view the data in this table –

- Select the table from BigQuery Explorer by navigating 'project_id' > 'dataset' >
  'table_name'
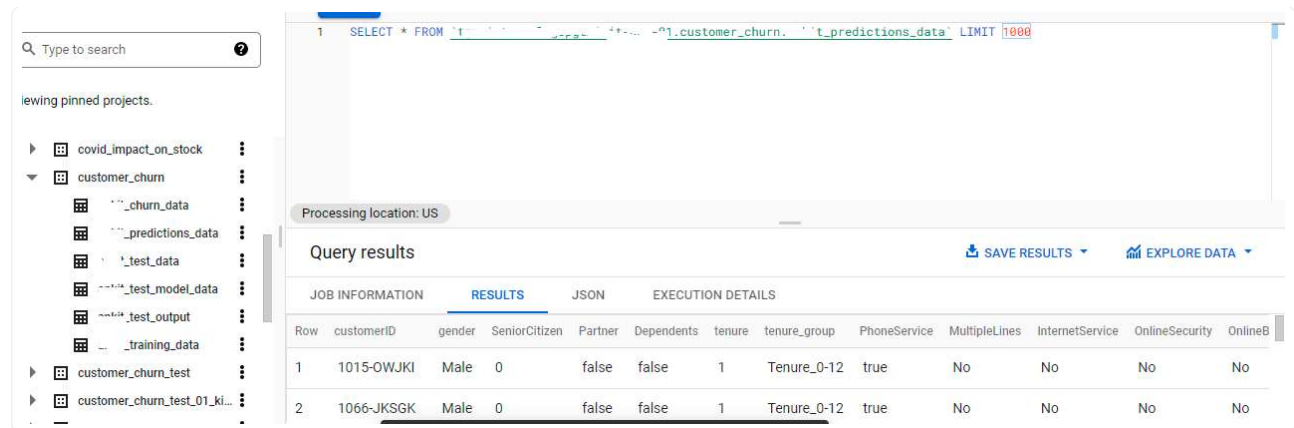- Click on the **Preview** button to see the data in the table

**Note:** If the **Preview** button is not visible, run the below queries to view the data. However, these queries will be charged for the full table scan.

```
SELECT * FROM `<project_name>.<dataset_name>.<your_name_here>_predictions_data` LI
```

**Note:** Edit all occurrences of <project_name> and <dataset_name> to match the values of the variables PROJECT_ID, and BQ_DATASET_NAME respectively



# 6. Model Evaluation

Repeat the same steps as above to submit another batch for model training.

## 6.1. Create a new batch

Navigate to Dataproc > Serverless > Batches and click on **+CREATE**



## 6.2. Provide the details for the batch

Next, fill in the following values in the batch creation window as shown in the images below:

- **Batch ID -** A unique identifier for your batch
- **Region -** The region name provided by the Admin team

- **Batch Type -** PySpark
- **Main Python File -** gs://<your_code_bucket_name>/customer_churn/00-scripts/customer_churn_model_testing.py
- **JAR Files -** gs://spark-lib/bigquery/spark-bigquery-with-dependencies_2.12-0.22.2.jar
- **Arguments -**
  Four Arguments needs to be provided.
  - <your_project_id>
  - <your_dataset_name>
  - <your_code_bucket_name>
  - <your_name> **Note:** Press RETURN after each argument
- **Service Account -** <UMSA_NAME>@<PROJECT_ID>.iam.gserviceaccount.com
- **Network Configuration -** select the network and subnetwork with Private Google Access Enabled Run PySpark Serverless Batch for Data Preparation
- **History Server Cluster -** <your_phs_cluster_name>

### Batch info

Batch ID *
batch-088e

Region *
us-central1 ▼

### Container

Batch type *
PySpark ▼

Main python file *
gs://ch·· ···· 1/customer_churn/00-scripts/bq_code_files/code_files_bigquery_code

Can be a GCS file with the gs:// prefix, an HDFS file on the cluster with the hdfs:// prefix, or a local file on the cluster with the file:// prefix

Additional python files

Custom container image

Specify a custom container image to add Java or Python dependencies not provided by the default container image. You must host your custom container on Container Registry.

Jar files
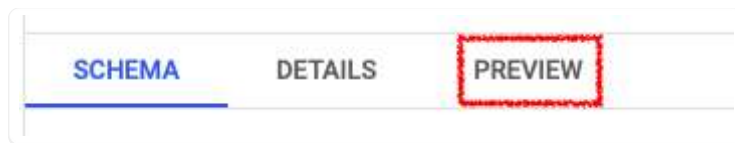gs://spark-lib/bigquery/spark-bigquery-with-dependencies_2.12-0.22.2.jar ⊗

## 6.3. Query the model_test results BQ table

Navigate to BigQuery Console, and check the **customer_churn_lab** dataset.
Once the model_testing batch is completed, a new table '<your_name_here>_test_output'
will be created.

To view the data in this table -

- Select the table from BigQuery Explorer by navigating 'project_id' > 'dataset' >
  'table_name'
- Click on the **Preview** button to see the data in the table



**Note:** If the **Preview** button is not visible, run the below queries to view the data. However,
these queries will be charged for the full table scan.

```
SELECT * FROM `<project_name>.<dataset_name>.<your_name_here>_test_output` LIMIT 1
```

**Note:** Edit all occurrences of <project_name> and <dataset_name> to match the values of
the variables PROJECT_ID, and BQ_DATASET_NAME respectively



# 7. Logging

## 7.1 Serverless Batch logs

Logs associated with the application can be found in the logging console under **Dataproc > Serverless > Batches > <batch_name>**.

You can also click on "View Logs" button on the Dataproc batches monitoring page to get to the logging page for the specific Spark job.



## 7.2 Persistent History Server logs

To view the Persistent History server logs, click the 'View History Server' button on the Dataproc batches monitoring page and the logs will be shown as below: