

06a_customer_churn_gcloud_execution.md

Customer Churn using Serverless Spark through Google Cloud Shell

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

Following are the lab modules:

1. [Understanding Data](#)
2. [Solution Architecture](#)
3. [Declaring cloud shell Variables](#)
4. [Data Preparation](#)
5. [Model Training and Testing](#)
6. [Model Evaluation](#)
7. [Logging](#)

1. Understanding Data

The dataset used for this project are [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.

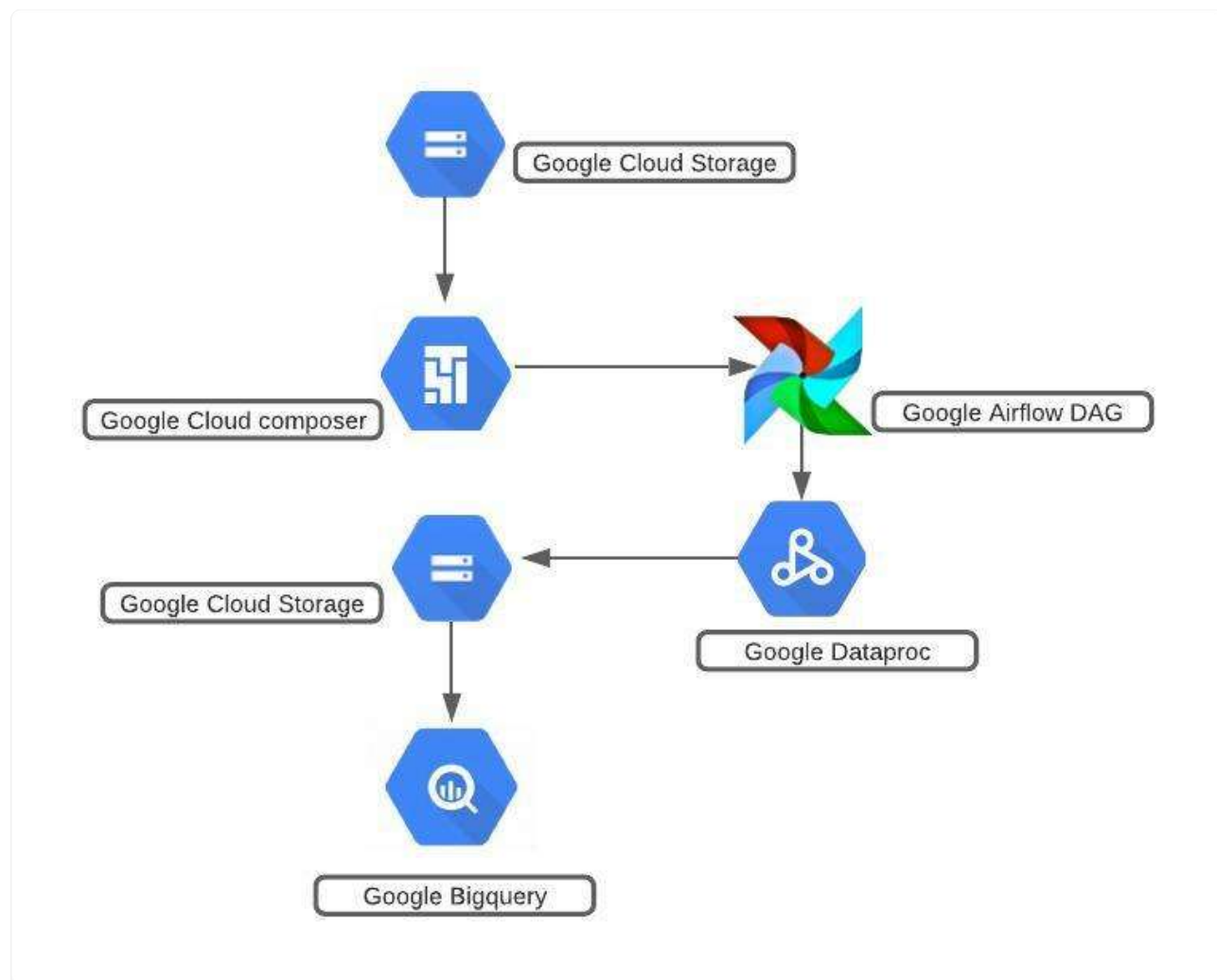
- `error_rate`

- rerror_rate
- same_srv_rate
- diff_srv_rate
- srv_count

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

- srv_error_rate
- srv_error_rate
- srv_diff_host_rate

2. Solution Architecture



Model Pipeline

The model pipeline involves the following steps:

- Data cleanup and preparation
- Building and training two Machine Learning Models (Logistic Regression and Random Forest Classifier) before saving them into cloud storage
- Using the model built in above step to evaluate test data

3. Declaring cloud shell variables

3.1 Set the PROJECT_ID in Cloud Shell

Open Cloud shell or navigate to shell.cloud.google.com

Run the below

```
gcloud config set project $PROJECT_ID
```

3.2 Verify the PROJECT_ID in Cloud Shell

Next, run the following command in cloud shell to ensure that the current project is set correctly:

```
gcloud config get-value project
```

3.3 Declare the variables

Based on the prereqs and checklist, declare the following variables in cloud shell by replacing with your values:

```
PROJECT_ID=$(gcloud config get-value project)    #current GCP project where we ar
REGION=                                           #GCP region where all our resour
SUBNET=                                           #subnet which has private google
BUCKET_CODE=                                     #GCP bucket where our code, data
BUCKET_PHS=                                      #bucket where our application lc
HISTORY_SERVER_NAME=                            #name of the history server whic
BQ_DATASET_NAME=                                #BigQuery dataset where all the
UMSA=serverless-spark                           #name of the user managed servic
SERVICE_ACCOUNT=$UMSA@$PROJECT_ID.iam.gserviceaccount.com
NAME=                                             #Your unique identifier
```

Note: For all the variables except 'NAME', please ensure to use the values provided by the admin team.

3.4 Update Cloud Shell SDK version

Run the below on cloud shell-

```
gcloud components update
```

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. Run PySpark Serverless Batch for Data Preparation

Run the below on cloud shell -

```
gcloud dataproc batches submit \  
  --project $PROJECT_ID \  
  --region $REGION \  
  pyspark --batch ${NAME}-batch-${RANDOM} \  
  gs://$BUCKET_CODE/customer_churn/00-scripts/customer_churn_data_prep.py \  
  --jars gs://spark-lib/bigquery/spark-bigquery-with-dependencies_2.12-0.22.2.jar \  
  --subnet $SUBNET \  
  --service-account $SERVICE_ACCOUNT \  
  --history-server-cluster projects/$PROJECT_ID/regions/$REGION/clusters/$HISTORY_SE  
  -- $PROJECT_ID $BQ_DATASET_NAME $BUCKET_CODE $NAME
```

4.2. Check the output table in BQ

Navigate to BigQuery Console, and check the **customer_churn_lab** dataset.

Once the data preparation batch is completed, a new table

'<your_name_here>_training_data' and '<your_name_here>_test_data' will be created.

To view the data in these tables -

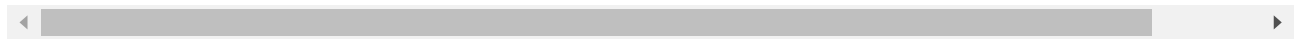
- Select the table from BigQuery Explorer by navigating 'project_id' > 'dataset' > 'table_name'
- Click on the **Preview** button to see the data in the table



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.

To query the table -

```
SELECT * FROM `<project_name>.<dataset_name>.<your_name_here>_training_data` LIMIT 1000
SELECT * FROM `<project_name>.<dataset_name>.<your_name_here>_test_data` LIMIT 1000
```



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

 A screenshot of the BigQuery interface. On the left, there is a sidebar with a search bar and a list of pinned projects. The main area shows a query editor with a query that has been executed. Below the editor, there is a 'Query results' section with a table of results. The table has columns for customerID, gender, SeniorCitizen, Partner, Dependents, tenure, tenure_group, PhoneService, MultipleLines, InternetService, and OnlineSecurity. The first two rows of data are visible.

Row	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	tenure_group	PhoneService	MultipleLines	InternetService	OnlineSecurity
1	0661-KQHNK	Female	0	true	true	6	Tenure_0-12	true	No	No	No
2	1269-FOYWN	Male	0	true	true	44	Tenure_24-48	true	No	No	No

5. Model Training and Testing

5.1. Run PySpark Serverless Batch for Model Training and Testing

The following script will train the model and save the model in the bucket.

Use the gcloud command below:

```
gcloud dataproc batches submit \
  --project $PROJECT_ID \
  --region $REGION \
  pyspark --batch ${NAME}-batch-${RANDOM} \
  gs://$BUCKET_CODE/customer_churn/00-scripts/customer_churn_model_building.py \
  --jars gs://spark-lib/bigquery/spark-bigquery-with-dependencies_2.12-0.22.2.jar \
  --subnet $SUBNET \
  --service-account $SERVICE_ACCOUNT \
  --history-server-cluster projects/$PROJECT_ID/regions/$REGION/clusters/$HISTORY_SE
  -- $PROJECT_ID $BQ_DATASET_NAME $BUCKET_CODE $NAME
```

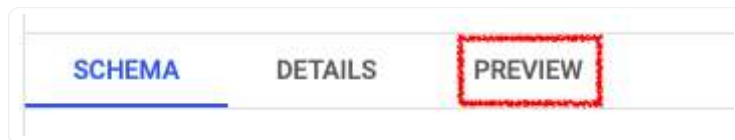
5.2. Query the model_test results BQ table

Navigate to BigQuery Console, and check the **customer_churn_lab** dataset.

Once the modelling 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

Query results

Processing location: US

SAVE RESULTS EXPLORE DATA

Row	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	tenure_group	PhoneService	MultipleLines	InternetService	OnlineSecurity	OnlineB
1	1015-OWJKI	Male	0	false	false	1	Tenure_0-12	true	No	No	No	No
2	1066-JKSGK	Male	0	false	false	1	Tenure_0-12	true	No	No	No	No

6. Model Evaluation

6.1. Run PySpark Serverless Batch for Model Evaluation

The following script will load the model and predict the new data.

Use the gcloud command below:

```
gcloud dataproc batches submit \
  --project $PROJECT_ID \
  --region $REGION \
  pyspark --batch ${NAME}-batch-${RANDOM} \
  gs://$BUCKET_CODE/customer_churn/00-scripts/customer_churn_model_testing.py \
  --jars gs://spark-lib/bigquery/spark-bigquery-with-dependencies_2.12-0.22.2.jar \
  --subnet $SUBNET \
  --service-account $SERVICE_ACCOUNT \
  --history-server-cluster projects/$PROJECT_ID/regions/$REGION/clusters/$HISTORY_SE
  -- $PROJECT_ID $BQ_DATASET_NAME $BUCKET_CODE $NAME
```

6.2. Query the model_test results BQ table

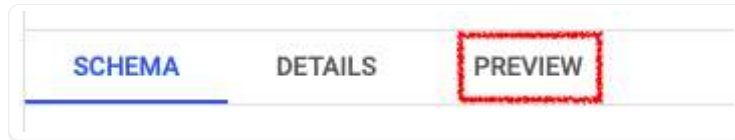
Navigate to BigQuery Console, and check the **customer_churn_lab** dataset.

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

Processing location: US

Query results

SAVE RESULTS EXPLORE DATA

Row	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	OnlineSecurity	OnlineBackup	DeviceProtection
1	5575-GNVDE	Male	0	false	false	34	true	No	DSL	Yes	No	Yes

7. Logging

7.1 Serverless Batch logs

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

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

Batch ID: batch-8459

Batch UUID: a1c152d3-4286-412f-ae68-36360f9c9a01

Resource type: Batch

Status: ✔ Succeeded

CLONE
DELETE
VIEW LOGS
REFRESH
VIEW SPARK HISTORY SERVER

Click here to view logs in Cloud Logging

MONITORING
DETAILS

i Metrics for a batch may lag behind the batch run by several minutes.

RESET ZOOM
 1 hour 6 hours 12 hours 1 day 2 days 4 days 7 days 14 days 30 days ✓ 11:43 - 11:48 ▼

Output LINE WRAP: OFF

History Server

Event log directory: gs://bo_bucket_phs/phs/*/spark-job-history

Last updated: 2022-03-16 14:24:27

Client local time zone: America/Los_Angeles

Search:

Version	App ID	App Name	Driver Host	Started	Completed	Duration	Spark User	Last Updated	Event Log
3.2.1	app-20220316212124-0000	[REDACTED]	10.0.0.8	2022-03-16 14:21:21	2022-03-16 14:22:01	40 s	spark	2022-03-16 14:22:01	Download

Showing 1 to 1 of 1 entries
[Show incomplete applications](#)

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:

Dataproc

obs on clusters

Clusters

Jobs

Workflows

Auto-scaling policies

Serverless

Batches

Utilities

Component exchange

Metastore

Workbench

Release notes

← batch-8459

CLONE

DELETE

VIEW LOGS

REFRESH

VIEW SPARK HISTORY SERVER

Batch ID

Batch UUID

Resource type

Status

MONITORING

DETAILS

Metrics for a batch may lag behind the batch run by several minutes.

RESET ZOOM


1 hour 6 hours 12 hours 1 day 2 days 4 days 7 days 14 days 30 days

✓ 11:43 - 11:48

Batch Spark Executors

Output

LINE WRAP: OFF

 **History Server**
3.1.2

Event log directory: gs://customer-churn-gcloud-execution-logs-1/phs/*/spark-job-history

Last updated: 2022-03-30 13:32:41

Client local time zone: Asia/Calcutta

Search:

Version	App ID	App Name	Driver Host	Started	Completed	Duration	Spark User	Last Updated	Event Log
3.2.1	app-20220330061551-0000	churn model	10.202.0.56	2022-03-30 11:45:49	2022-03-30 11:46:37	47 s	spark	2022-03-30 11:46:38	Download

Showing 1 to 1 of 1 entries
[Show incomplete applications](#)