

# Fraud Detection on Credit Card Transactions using BigQuery ML on GCP

## Dataset

Data used (also available as public dataset on bigquery) -

<https://github.com/jbrownlee/Datasets/blob/d20fcb6402ae34e653d4513b00f39257bb37ed7f/creditcard.csv.zip>

Dataset holds 28 feature which are obtained after PCA.

Feature time shows the time elapsed between first and that specific transaction.

Feature amount refers to the transaction amount.

Feature class refers to the transaction being fraudulent or not (1 for fraud and 0 otherwise)

## Ingesting the data-

Query - “`SELECT * FROM `bigquery-public-data.ml_datasets.ulb_fraud_detection` LIMIT 5`”

This query fetches the dataset stored in the GCP public repository and shows us the first 5 rows.

The screenshot displays the Google Cloud Platform BigQuery interface. The top navigation bar includes the Google Cloud Platform logo, a dropdown menu for 'vertex-ai-projects', a search bar, and a notification icon. The left sidebar contains a navigation menu with options: Analysis, SQL workspace, Data transfers, Scheduled queries, Migration, SQL translation, Administration, Monitoring, Capacity management, and BI Engine. The main content area shows a query editor with the following SQL query:

```
1 SELECT
2 *
3 FROM
4 `bigquery-public-data.ml_datasets.ulb_fraud_detection`
5 LIMIT
6 5
```

Below the query editor, the 'Query results' section is visible. It shows a status message: 'Query complete (2.1 sec elapsed, 67.4 MB processed)'. The results are displayed in a table with 11 columns: Row, Time, V1, V2, V3, V4, V5, V6, V7, V8, and V9. The table contains 5 rows of data.

Row	Time	V1	V2	V3	V4	V5	V6	V7	V8	V9
1	282.0	-0.356466189895633	0.725417515223392	1.97174888880699	0.8313427151031229	0.36968143455457897	-0.107776262637238	0.75160990401556	-0.120166	
2	380.0	-1.29983679037626	0.8818174638094001	1.4528418792573	-1.29369833187415	-0.025104981110272002	-1.170102626414	0.861610195375755	-0.1939	
3	403.0	1.23741280400293	0.512364829919811	0.687745559627426	1.69387239409344	-0.236322621630244	-0.650231739782802	0.118066209111006	-0.230545	
4	430.0	-1.8602576921529799	-0.629858920058775	0.9665704477622901	0.844632076311716	0.7599826624018221	-1.4811729034822199	-0.509681452204136	0.5407	
5	711.0	-0.431349344181742	1.0276943686964002	2.67081622589653	2.08478703855765	-0.27456734571551106	0.286856036075343	0.15210974930542498	0.200	

We can see the job description as well

## Query results

[SAVE RESULTS](#)[EXPLORE DATA](#) ▼[Job information](#)[Results](#)[JSON](#)[Execution details](#)

Query completed in 2.131 sec  
2:33 AM

```
1 SELECT
2 *
3 FROM
4 `bigquery-public-data.ml_datasets.ulb_fraud_detection`
5 LIMIT
6 5
```

Job ID	vertex-ai-projects:US.bqjob_54b62d14_17dd9a7dcd2
User	nikhilsanghi6@gmail.com
Location	United States (US)
Creation time	Dec 21, 2021, 2:33:17 AM
Start time	Dec 21, 2021, 2:33:17 AM
End time	Dec 21, 2021, 2:33:19 AM
Duration	2.1 sec
Bytes processed	67.36 MB
Bytes billed	68 MB
Job priority	INTERACTIVE
Destination table	<a href="#">Temporary table</a>
Use legacy SQL	false
Session ID	

## Data exploration

Let's check the data distribution-

Query = "SELECT COUNT(\*) FROM `bigquery-public-data.ml\_datasets.ulb\_fraud\_detection`  
where Class=0;"

This query fetches all the rows which are non-fraudulent and returns the total count.

RUN

SAVE

SCHEDULE

MORE

1

SELECT

2

COUNT(\*)

3

FROM

4

`bigquery-public-data.ml\_datasets.ulb\_fraud\_detection`

5

WHERE

6

Class=0;

Processing location: US

Query results

SAVE RESULTS

EXPLORE DATA

Query complete (0.4 sec elapsed, 2.2 MB processed)

Job information

Results

JSON

Execution details

Row	f0_
1	284315

Query = "SELECT COUNT(\*) FROM `bigquery-public-data.ml\_datasets.ulb\_fraud\_detection`  
where Class=1;"

This query fetches all the rows which are fraudulent and returns the total count.

<div> <div>RUN</div> <div>SAVE</div> <div>SCHEDULE</div> <div>MORE</div> </div>		
1	SELECT	
2	COUNT(*)	
3	FROM	
4	`bigquery-public-data.ml_datasets.ulb_fraud_detection`	
5	WHERE	
6	Class=1;	
Processing location: US		
<div> <div>Query results</div> <div>SAVE RESULTS</div> <div>EXPLORE DATA</div> </div>		
Query complete (0.3 sec elapsed, 2.2 MB processed)		
<div> <div>Job information</div> <div>Results</div> <div>JSON</div> <div>Execution details</div> </div>		
Row	f0_	
1	492	

## Observation -

There are 492 fraudulent Transactions out of 284315, which makes this dataset highly imbalanced.

Only 0.17% Fraudulent transactions.

## Model Building

Query = “

```
CREATE OR REPLACE MODEL fraud_detection.ulb_fraud_detection
TRANSFORM(
  * EXCEPT(Amount),
  SAFE.LOG(Amount) AS log_amount
)
OPTIONS(
  INPUT_LABEL_COLS=['class'],
  AUTO_CLASS_WEIGHTS = TRUE,
  DATA_SPLIT_METHOD='seq',
  DATA_SPLIT_COL='Time',
  MODEL_TYPE='logistic_reg'
) AS

SELECT
*
FROM `bigquery-public-data.ml_datasets.ulb_fraud_detection`
```

”

for more info - <https://cloud.google.com/bigquery-ml/docs/reference/standard-sql/bigqueryml-syntax-create>

This query creates a new model or replaces an already existing model, transforms the input variables except the amount feature column, and safely logs the amount column. (Safe makes sure that a null value is returned instead of any possible error).

### Options -

This needs the parameters we need to define for the model to be trained.

**Input\_label\_cols** refer to the prediction column which is “class” in our case.

**Auto\_class\_weights** refer to the weights assigned to different classes in our prediction column. It is set to True for imbalanced dataset such as ours.

**Model\_type** refer to the training algorithm used which is “Logistic Regression” in our case.

**Data\_split\_type** refer to the split between training and testing data.

```
1 CREATE OR REPLACE MODEL fraud_detection.ulb_fraud_detection
2 TRANSFORM(
3   * EXCEPT(Amount),
4   SAFE.LOG(Amount) AS log_amount
5 )
6 OPTIONS(
7   INPUT_LABEL_COLS=['class'],
8   AUTO_CLASS_WEIGHTS = TRUE,
9   DATA_SPLIT_METHOD='seq',
10  DATA_SPLIT_COL='Time',
11  MODEL_TYPE='logistic_reg'
12 ) AS
13
14 SELECT
15 *
16 FROM `bigquery-public-data.ml_datasets.ulb_fraud_detection`
```

Query results		
Job information   Results   Execution details		
Elapsed time	Slot time consumed ?	Stages ?
2 min	—	<div><div>✓</div> Preprocess 1 min 7.240 sec</div> <div><div>✓</div> Train</div>

RUN

SAVE

SCHEDULE

MORE

This query will process 67.4 MiB (ML) when

1CREATE OR REPLACE MODEL fraud\_detection.ulb\_fraud\_detection

2TRANSFORM(

3  \* EXCEPT(Amount),

4  SAFE\_LOG(Amount) AS log\_amount

Query results

Query complete (3 min 17 sec elapsed, 67.4 MB (ML) processed)

Job information

Results

Execution details

Elapsed time

3 min 17 sec

Slot time consumed

7 min 54.501 sec

Stages

Preprocess1 min 7.240 sec

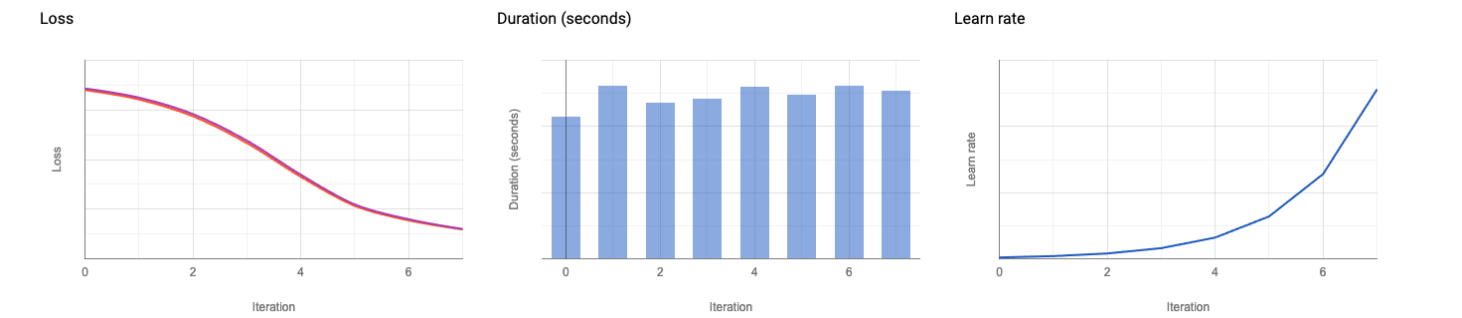
Train2 min 2.352 sec

Evaluate6.913 sec

Training iterations

Completed: 8

Planned: 20



The model training takes approximately 1 min 7.24 secs to preprocess the data and 2 mins 2.352 secs to train and finally 6.913 secs to evaluate the predictions.

We notice that the loss decreases with every iteration and the learning rate increases with each iteration

Schema of the data set -

ulb_fraud_detection			
DETAILS	TRAINING	EVALUATION	SCHEMA
Field name	Type	Mode	Description
V28	FLOAT64	NULLABLE	
V13	FLOAT64	NULLABLE	
V15	FLOAT64	NULLABLE	
V18	FLOAT64	NULLABLE	
V22	FLOAT64	NULLABLE	
V1	FLOAT64	NULLABLE	
V23	FLOAT64	NULLABLE	
V21	FLOAT64	NULLABLE	
Amount	FLOAT64	NULLABLE	
V9	FLOAT64	NULLABLE	
V17	FLOAT64	NULLABLE	
V19	FLOAT64	NULLABLE	
V27	FLOAT64	NULLABLE	
V20	FLOAT64	NULLABLE	
V8	FLOAT64	NULLABLE	
V24	FLOAT64	NULLABLE	
V5	FLOAT64	NULLABLE	
V10	FLOAT64	NULLABLE	
V6	FLOAT64	NULLABLE	
V16	FLOAT64	NULLABLE	
V26	FLOAT64	NULLABLE	
V25	FLOAT64	NULLABLE	
V14	FLOAT64	NULLABLE	
V2	FLOAT64	NULLABLE	
V11	FLOAT64	NULLABLE	
V12	FLOAT64	NULLABLE	
V3	FLOAT64	NULLABLE	
V7	FLOAT64	NULLABLE	
V4	FLOAT64	NULLABLE	

## Model Evaluation -

Observation -

At " 0 " Positive threshold we observe that the only 75 prediction are right.  
0 precision with 100 % recall, which is totally undesirable.

Our model have a right balance between precision and recall depending upon or business use case.

Now changing the positive thresholds->

## Aggregate Metrics ?

Threshold ?	0.5000
Precision ?	0.0518
Recall ?	0.8933
Accuracy ?	0.9784
F1 score ?	0.0980
Log loss ?	0.1197
ROC AUC ?	0.9853

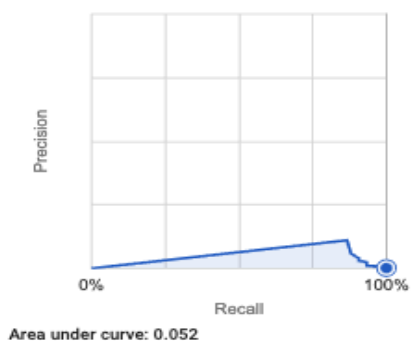
## Score threshold

Positive class threshold ?	<input type="text" value="0.0000"/>
Positive class	1
Negative class	0
Precision ?	0.0013
Recall ?	1.0000
Accuracy ?	0.0013
F1 score ?	0.0026

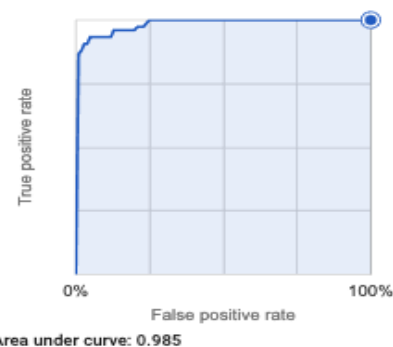
## Precision-recall by threshold ?



## Precision-recall curve ?



## ROC curve ?



## Confusion matrix

This table shows how often the model classified each label correctly (in blue), and which labels were most often confused for that label (in gray).

True label	Predicted label	
	1	0
1	75	0
0	56955	0

Changing the Positive threshold to 2 % doesn't change much of the metrics, there are few predictions which fall in the right category.



Aggregate Metrics

Threshold	0.5000
Precision	0.0518
Recall	0.8933
Accuracy	0.9784
F1 score	0.0980
Log loss	0.1197
ROC AUC	0.9853

Score threshold

Positive class threshold	0.0201
Positive class	1
Negative class	0
Precision	0.0015
Recall	1.0000
Accuracy	0.1314
F1 score	0.0030

Precision-recall by threshold

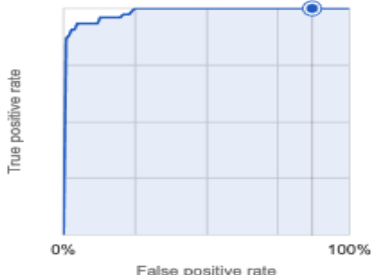


Precision-recall curve



Area under curve: 0.052

ROC curve



Area under curve: 0.985

Confusion matrix

This table shows how often the model classified each label correctly (in blue), and which labels were most often confused for that label (in gray).

True label	Predicted label	
	1	0
1	75	0
0	49538	7417

Changing the positive threshold to 20%, we notice a significant of the predictions falling in the categories.

## Aggregate Metrics ?

Threshold ?	0.5000
Precision ?	0.0518
Recall ?	0.8933
Accuracy ?	0.9784
F1 score ?	0.0980
Log loss ?	0.1197
ROC AUC ?	0.9853

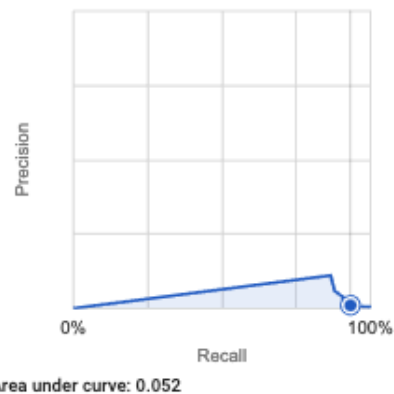
## Score threshold

Positive class threshold ?	<input type="text" value="0.2028"/>
Positive class	1
Negative class	0
Precision ?	0.0102
Recall ?	0.9333
Accuracy ?	0.8812
F1 score ?	0.0202

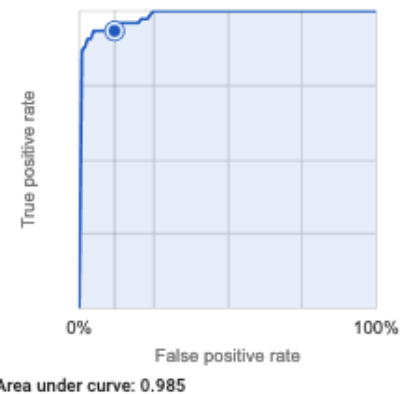
## Precision-recall by threshold ?



## Precision-recall curve ?



## ROC curve ?



## Confusion matrix

☒ Item counts 

This table shows how often the model classified each label correctly (in blue), and which labels were most often confused for that label (in gray).

True label	Predicted label	
	1	0
1	70	5
0	6770	50185

## Observation( Evaluation) -

Further increasing the Positive Threshold value to 0.6355 we obtain the maximum precision value of 0.11 with recall 0.866 and accuracy of 0.99 with a decent usable f1 score of 0.19.

There are 522 wrong predictions as fraudulent transactions and has 10 transactions which are predicted as non-fraudulent which were actually fraudulent which is more serious case, since we don't want any transaction that is fraudulent to go un-noticed, it will okay for us to flag transaction as fraudulent and then investigate (522) rather than flagging 10 cases as non-fraudulent for our banking use case.

This threshold can change the values of model predicted positive values depending upon the business use case, in some cases we need more precision such as our case and in some we need more recall.

We also can notice that the area under roc curve gives us the value of 0.985 which gives us good estimate of model performance.

ulb\_fraud\_detection

QUERY MODEL

DELETE MODEL

EXPORT MODEL

DETAILS

TRAINING

EVALUATION

SCHEMA

Aggregate Metrics

Threshold	0.5000
Precision	0.0518
Recall	0.8933
Accuracy	0.9784
F1 score	0.0980
Log loss	0.1197
ROC AUC	0.9853

Score threshold

Positive class threshold

0.6355

Positive class	1
Negative class	0
Precision	0.1107
Recall	0.8667
Accuracy	0.9907
F1 score	0.1964

Precision-recall by threshold

Precision-recall curve

ROC curve

Confusion matrix

Item counts

This table shows how often the model classified each label correctly (in blue), and which labels were most often confused for that label (in gray).

True label	Predicted label	
	1	0
1	65	10
0	522	56433

RUN

SAVE

SCHEDULE

MORE

This query will process 67.4 MiB when run.

```

1 SELECT
2   Amount,
3   predicted_class_probs,
4   Class
5 FROM
6   ML.PREDICT( MODEL fraud_detection.ulb_fraud_detection,
7   (
8     SELECT
9       *
10    FROM
11     `bigquery-public-data.ml_datasets.ulb_fraud_detection`
12   WHERE
13     Time = 80000.0) )

```

Processing location: US

Query results

SAVE RESULTS

EXPLORE DATA

Query complete (5.8 sec elapsed, 67.4 MB processed)

Job information   **Results**   JSON   Execution details

Row	Amount	predicted_class_probs.label	predicted_class_probs.prob	Class
1	105.62	1	0.01747670893573443	0
		0	0.9825232910642656	

Predicting class for a specific row with timestamp.

Class column indicates the actual class of the example and the predicted\_class\_probs\_label indicates the predicted label for that transaction example. This shows that the probability for the transaction to be fraudulent is 0.017 % while the probability for the transaction to be non-fraudulent is 0.98 %, which is correct in this case.

**RESULT -**

We created a Logistic Regression model using BigQuery ML to detect Fraudulent transactions among a highly imbalanced dataset.