LAB 01: Basic Feature Engineering in BQML

Learning Objectives

- · Setup up the environment
- · Create the project dataset
- Create the feature engineering training table
- · Create and evaluate the benchmark/baseline model
- · Extract temporal features
- · Perform a feature cross on temporal features
- · Evaluate model performance

Introduction

In this lab, we utilize feature engineering to improve the prediction of the fare amount for a taxi ride in New York City. We will use BigQuery ML to build a taxifare prediction model, using feature engineering to improve and create a final model.

In this Notebook we set up the environment, create the project dataset, create a feature engineering table, create and evaluate a benchmarkASZ model, extract temporal features, perform a feature cross on temporal features, and evaluate model performance throughout the process.

Each learning objective will correspond to a **#TODO** in this student lab notebook -- try to complete this notebook first and then review the <u>solution notebook (../solution/feateng-solution_bqml.ipynb)</u>. **NOTE TO SELF**: UPDATE HYPERLINK.

Set up environment variables and load necessary libraries

```
In [47]:
```

```
%%bash
export PROJECT=$(gcloud config list project --format "value(core.project)")
echo "Your current GCP Project Name is: "$PROJECT
```

Your current GCP Project Name is: cloud-training-demos

In [48]:

```
import os

PROJECT = "cloud-training-demos" # REPLACE WITH YOUR PROJECT NAME
REGION = "us-west1-b" # REPLACE WITH YOUR BUCKET REGION e.g. us-central1

# Do not change these
os.environ["PROJECT"] = PROJECT
os.environ["REGION"] = REGION
os.environ["BUCKET"] = PROJECT # DEFAULT BUCKET WILL BE PROJECT ID

if PROJECT == "your-gcp-project-here":
    print("Don't forget to update your PROJECT name! Currently:", PROJECT)
```

Check that the Google BigQuery library is installed and if not, install it.

```
In [ ]:
```

```
!pip freeze | grep google-cloud-bigquery==1.6.1 \mid\mid pip install google-cloud-bigquery==1.6.1
```

The source dataset

Our dataset is hosted in <u>BigQuery (https://cloud.google.com/bigquery/)</u>. The taxi fare data is a publically available dataset, meaning anyone with a GCP account has access. Click <u>here (https://console.cloud.google.com/bigquery?project=bigquery-public-data&p=nyc-tlc&d=yellow&t=trips&page=table)</u> to acess the dataset.

The Taxi Fare dataset is relatively large at 55 million training rows, but simple to understand, with only 6 features. The fare amount is the target, the continuous value we'll train a model to predict.

Create a BigQuery Dataset and Google Cloud Storage Bucket

A BigQuery dataset is a container for tables, views, and models built with BigQuery ML. Let's create one called **feat_eng** if we have not already done so in an earlier lab. We'll do the same for a GCS bucket for our project too.

In [75]:

```
%%bash
## Create a BiqQuery dataset for feat eng TEST if it doesn't exist
datasetexists=$(bq ls -d | grep -w feat eng)
if [ -n "$datasetexists" ]; then
    echo -e "BigQuery dataset already exists, let's not recreate it."
else
    echo "Creating BigQuery dataset titled: feat_eng"
    bq --location=US mk --dataset \
        --description 'Taxi Fare' \
        $PROJECT:feat eng
   echo "\nHere are your current datasets:"
   bq ls
fi
## Create GCS bucket if it doesn't exist already...
exists=$(gsutil ls -d | grep -w gs://${PROJECT}/)
if [ -n "$exists" ]; then
   echo -e "Bucket exists, let's not recreate it."
else
   echo "Creating a new GCS bucket."
   gsutil mb -1 ${REGION} gs://${PROJECT}
   echo "\nHere are your current buckets:"
   gsutil ls
fi
```

BigQuery dataset already exists, let's not recreate it. Bucket exists, let's not recreate it.

Create the training data table

Since there is already a publicly available dataset, we can simply create the training data table using this raw input data. Note the WHERE clause in the below query: This clause allows us to TRAIN a portion of the data (e.g. one hundred thousand rows versus one million rows), which keeps your query costs down. If you need a refresher on using MOD() for repeatable splits see this <u>post (https://www.oreilly.com/learning/repeatable-sampling-of-data-sets-in-bigquery-for-machine-learning)</u>.

• Note: The dataset in the create table code below is the one created previously, e.g. "feat_eng". The table name is "feateng training data".

Exercise: RUN the query to create the table.

In [76]:

```
%%bigquery
#Objective: Create the feature engineering training table
CREATE OR REPLACE TABLE feat eng.feateng training data
AS
SELECT
  (tolls_amount + fare_amount) AS fare_amount,
 passenger_count*1.0 AS passengers,
 pickup_datetime,
 pickup_longitude AS pickuplon,
 pickup latitude AS pickuplat,
 dropoff longitude AS dropofflon,
 dropoff latitude AS dropofflat
FROM `nyc-tlc.yellow.trips`
WHERE MOD(ABS(FARM_FINGERPRINT(CAST(pickup_datetime AS STRING))), 10000) = 1
 AND fare amount >= 2.5
 AND passenger_count > 0
 AND pickup longitude > -78
 AND pickup_longitude < -70
 AND dropoff longitude > -78
 AND dropoff longitude < -70
 AND pickup latitude > 37
 AND pickup_latitude < 45
 AND dropoff_latitude > 37
 AND dropoff latitude < 45
```

Out[76]:

Verify table creation

Verify that you created the dataset.

In [154]:

```
%%bigquery
-- LIMIT 0 is a free query; this allows us to check that the table exists.
SELECT * FROM feat_eng.feateng_training_data
LIMIT 0
```

Out[154]:

fare_amount passengers pickup_datetime pickuplon pickuplat dropofflon dropofflat

Benchmark Model: Create the benchmark/baseline Model

Next, you create a linear regression baseline model with no feature engineering. Recall that a model in BigQuery ML represents what an ML system has learned from the training data. A baseline model is a solution to a problem without applying any machine learning techniques.

When creating a BQML model, you must specify the model type (in our case linear regression) and the input label (fare amount). Note also that we are using the training data table as the data source.

Exercise: Create the SQL statement to create the model "Benchmark Model".

In []:

```
%%bigquery
# TODO: Your code goes here
# Objective: Create and evaluate the benchmark/baseline model
```

In [84]:

```
%%bigquery
# SOLUTION
CREATE OR REPLACE MODEL feat_eng.benchmark_model

OPTIONS
   (model_type='linear_reg',
        input_label_cols=['fare_amount'])
AS

SELECT
   fare_amount,
   passengers,
   pickup_datetime,
   pickuplon,
   pickuplat,
   dropofflon,
   dropofflon feat_eng.feateng_training_data
```

Out[84]:

REMINDER: The query takes several minutes to complete. After the first iteration is complete, your model (benchmark_model) appears in the navigation panel of the BigQuery web UI. Because the query uses a CREATE MODEL statement to create a model, you do not see query results.

You can observe the model as it's being trained by viewing the Model stats tab in the BigQuery web UI. As soon as the first iteration completes, the tab is updated. The stats continue to update as each iteration completes.

Once the training is done, visit the <u>BigQuery Cloud Console (https://console.cloud.google.com/bigquery)</u> and look at the model that has been trained. Then, come back to this notebook.

Evaluate the benchmark model

Note that BigQuery automatically split the data we gave it, and trained on only a part of the data and used the rest for evaluation. After creating your model, you evaluate the performance of the regressor using the ML.EVALUATE function. The ML.EVALUATE function evaluates the predicted values against the actual data.

NOTE: The results are also displayed in the <u>BigQuery Cloud Console</u> (https://console.cloud.google.com/bigguery) under the **Evaluation** tab.

Exercise: Review the learning and eval statistics for the benchmark_model.

In [115]:

```
%%bigquery
#Eval statistics on the held out data.
SELECT *, SQRT(loss) AS rmse FROM ML.TRAINING_INFO(MODEL feat_eng.benchmark_model)
```

Out[115]:

	training_run	iteration	loss	eval_loss	learning_rate	duration_ms	rmse
0	0	0	74.434806	68.880152	None	13164	8.627561

In [112]:

```
%%bigquery
SELECT * FROM ML.EVALUATE(MODEL feat_eng.benchmark_model)
```

Out[112]:

	mean_absolute_error	mean_squared_error	mean_squared_log_error	median_absolute_error	r2
0	5.213535	68.880152	0.258107	3.790249	0.
<	1			>	

NOTE: Because you performed a linear regression, the results include the following columns:

- · mean absolute error
- · mean squared error
- · mean squared log error
- median absolute error
- r2 score
- · explained variance

Resource for an explanation of the <u>Regression Metrics (https://towardsdatascience.com/metrics-to-evaluate-your-machine-learning-algorithm-f10ba6e38234)</u>.

Mean squared error (MSE) - Measures the difference between the values our model predicted using the test set and the actual values. You can also think of it as the distance between your regression (best fit) line and the predicted values.

Root mean squared error (RMSE) - The primary evaluation metric for this ML problem is the root mean-squared error. RMSE measures the difference between the predictions of a model, and the observed values. A large RMSE is equivalent to a large average error, so smaller values of RMSE are better. One nice property of RMSE is that the error is given in the units being measured, so you can tell very directly how incorrect the model might be on unseen data.

R2: An important metric in the evaluation results is the R2 score. The R2 score is a statistical measure that determines if the linear regression predictions approximate the actual data. Zero (0) indicates that the model explains none of the variability of the response data around the mean. One (1) indicates that the model explains all the variability of the response data around the mean.

Exercise: Write a SQL query to take the SQRT() of the mean squared error as your loss metric for evaluation for the benchmark_model.

In []:

```
%%bigquery
# TODO: Your code goes here
```

In [110]:

```
%%bigquery
#SOLUTION
SELECT SQRT(mean_squared_error) AS rmse FROM ML.EVALUATE(MODEL feat_eng.benchmark_model)
```

Out[110]:

rmse 0 8.299407

Model 1: EXTRACT dayofweek from the pickup_datetime feature.

- As you recall, dayofweek is an enum representing the 7 days of the week. This factory allows the enum to be obtained from the int value. The int value follows the ISO-8601 standard, from 1 (Monday) to 7 (Sunday).
- If you were to extract the dayofweek from pickup_datetime using BigQuery SQL, the dataype returned would be integer.

Exercise: EXTRACT dayofweek from the pickup_datetime feature.

Create a model titled "model_1" from the benchmark model and extract out the DayofWeek.

In []:

```
# TODO: Your code goes here
# Objective: Extract temporal features
```

In [103]:

```
%%bigquery
#SOLUTION
CREATE OR REPLACE MODEL feat eng.model 1
OPTIONS
  (model type='linear reg',
    input label cols=['fare amount'])
AS
SELECT
fare_amount,
  passengers,
  pickup datetime,
  EXTRACT(DAYOFWEEK FROM pickup_datetime) AS dayofweek,
  pickuplon,
  pickuplat,
  dropofflon,
  dropofflat FROM feat eng.feateng training data
```

Out[103]:

Once the training is done, visit the <u>BigQuery Cloud Console (https://console.cloud.google.com/bigquery)</u> and look at the model that has been trained. Then, come back to this notebook.

Exercise: Create two distinct SQL statements to see the TRAINING and EVALUATION metrics of model_1.

In []:

```
# TODO: Your code goes here
#Create the SQL statements to extract Model_1 TRAINING metrics.
```

In []:

```
# TODO: Your code goes here.
#Create the SQL statements to extract Model_1 EVALUATION metrics.
```

In [104]:

```
%%bigquery
#SOLUTION
SELECT *, SQRT(loss) AS rmse FROM ML.TRAINING_INFO(MODEL feat_eng.model_1)
```

Out[104]:

	training_run	iteration	loss	eval_loss	learning_rate	duration_ms	rmse
0	0	0	72.440724	88.953232	None	17251	8.511212

In [108]:

```
%%bigquery
#SOLUTION
SELECT * FROM ML.EVALUATE(MODEL feat_eng.model_1)
```

Out[108]:

	mean_absolute_error	mean_squared_error	mean_squared_log_error	median_absolute_error	r2
0	5.287076	88.953232	0.260932	3.713439	(
<	1			>	

Exercise: Write a SQL query to take the SQRT() of the mean squared error as your loss metric for evaluation for model_1.

In []:

```
# TODO: Your code goes here
#Create the SQL statement to EVALUATE Model_1 here.
```

In [117]:

```
%%bigquery
#SOLUTION
SELECT SQRT(mean_squared_error) AS rmse FROM ML.EVALUATE(MODEL feat_eng.model_1)
```

Out[117]:

rmse

9.431502

Model 2: EXTRACT hourofday from the pickup_datetime feature

As you recall, pickup_datetime is stored as a TIMESTAMP, where the Timestamp format is retrieved in the standard output format - year-month-day hour:minute:second (e.g. 2016-01-01 23:59:59). Hourofday returns the integer number representing the hour number of the given date.

Hourofday is best thought of as a discrete ordinal variable (and not a categorical feature), as the hours can be ranked (e.g. there is a natural ordering of the values). Hourofday has an added characteristic of being cyclic, since 12am follows 11pm and precedes 1am.

Exercise: EXTRACT hourofday from the pickup_datetime feature.

- · Create a model titled "model 2"
- EXTRACT the hourofday from the pickup datetime feature to improve our model's rmse.

In []:

```
# TODO: Your code goes here
# Objective: Extract temporal features
```

In [121]:

```
%%bigquery
#SOLUTION
CREATE OR REPLACE MODEL feat_eng.model_2
OPTIONS
  (model type='linear reg',
    input_label_cols=['fare_amount'])
AS
SELECT
  fare_amount,
  passengers,
  #pickup datetime,
  EXTRACT(DAYOFWEEK FROM pickup_datetime) AS dayofweek,
  EXTRACT(HOUR FROM pickup_datetime) AS hourofday,
  pickuplon,
  pickuplat,
  dropofflon,
  dropofflat
FROM `feat_eng.feateng_training_data`
```

Out[121]:

Exercise: Create two SQL statements to evaluate the model.

```
In [ ]:
```

```
# TODO: Your code goes here
```

In []:

```
# TODO: Your code goes here
```

In [122]:

```
%%bigquery
#SOLUTION
SELECT * FROM ML.EVALUATE(MODEL feat_eng.model_2)
```

Out[122]:

	mean_absolute_error	mean_squared_error	mean_squared_log_error	median_absolute_error	r2
0	5.256421	70.709518	0.262399	3.895509	0.
4				•	

In [123]:

```
%%bigquery
SELECT SQRT(mean squared error) AS rmse FROM ML.EVALUATE(MODEL feat eng.model 2)
```

Out[123]:

rmse

8.408895

Model 3: Feature cross dayofweek and hourofday using CONCAT

First, let's allow the model to learn traffic patterns by creating a new feature that combines the time of day and day of week (this is called a feature cross (https://developers.google.com/machine-learning/crashcourse/feature-crosses/video-lecture).

Note: BQML by default assumes that numbers are numeric features, and strings are categorical features. We need to convert both the dayofweek and hourofday features to strings because the model (Neural Network) will automatically treat any integer as a numerical value rather than a categorical value. Thus, if not cast as a string, the dayofweek feature will be interpreted as numeric values (e.g. 1,2,3,4,5,6,7) and hour ofday will also be interpreted as numeric values (e.g. the day begins at midnight, 00:00, and the last minute of the day begins at 23:59 and ends at 24:00). As such, there is no way to distinguish the "feature cross" of hourofday and dayofweek "numerically". Casting the dayofweek and hourofday as strings ensures that each element will be treated like a label and will get its own coefficient associated with it.

Exercise: Create the SQL statement to feature cross the dayofweek and hourofday using the CONCAT function. Name the model "model_3"

In []:

TODO: Your code goes here # Objective: Perform a feature cross on temporal features

In [1]:

```
%%bigquery
#SOLUTION
CREATE OR REPLACE MODEL feat_eng.model_3
OPTIONS
  (model type='linear reg',
    input_label_cols=['fare_amount'])
AS
SELECT
  fare_amount,
  passengers,
  #pickup datetime,
  #EXTRACT(DAYOFWEEK FROM pickup_datetime) AS dayofweek,
  #EXTRACT(HOUR FROM pickup_datetime) AS hourofday,
  CONCAT(CAST(EXTRACT(DAYOFWEEK FROM pickup datetime) AS STRING),
        CAST(EXTRACT(HOUR FROM pickup_datetime) AS STRING)) AS hourofday,
  pickuplon,
  pickuplat,
  dropofflon,
  dropofflat
FROM `feat_eng.feateng_training_data`
```

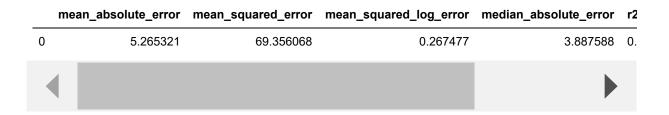
Out[1]:

Exercise: Create two SQL statements to evaluate the model.

In [2]:

```
%%bigquery
#SOLUTION
SELECT * FROM ML.EVALUATE(MODEL feat_eng.model_3)
```

Out[2]:



In [3]:

%%bigquery SELECT SQRT(mean_squared_error) AS rmse FROM ML.EVALUATE(MODEL feat_eng.model_3)

Out[3]:

rmse

8.328029