

### Main advantages of using BigQuery ML for model training, tuning and prediction:

- **Model training, evaluation and inference using SQL only**
- **Easy Access to data stored in BigQuery without the need to export it to local memory**
- **Centralized, ready to use model end-points, facilitating deployment and prototyping**
- **Distributed processing through BigQuery processing engine**

## CREATE MODEL

CREATE MODEL	CREATE MODEL
<b>Create ML model:</b> Generate a model to be stored in BigQuery	CREATE MODEL <model_name> CREATE MODEL IF NOT EXISTS <model_name> CREATE OR REPLACE MODEL <model_name>
<b>Define ML model type:</b> Type of model to be trained	OPTIONS( MODEL_TYPE = 'LINEAR_REG' )
<b>Define Hyperparameters (optional):</b> List of hyperparameters used for training	OPTIONS( <a href="#">list of hyperparameters</a> )
<b>Define training data:</b> Select the label and features from the data-source used to train the model	AS SELECT * FROM `my-project.my_dataset.my_table`

## DEFAULT DATA PREPROCESSING

Missing data	Null values are: <ul style="list-style-type: none"> <li>• Replaced by the mean in numerical columns</li> <li>• Mapped to an additional one-hot-encoded category in non numerical columns</li> </ul>
Standardization	Numerical columns are scaled using a standard scaler
One-hot encoding	Non-numerical non-array columns are one-hot-encoded

## MODEL LEARNING FLOW

Model Category	Model Type	Model Creation	Preprocessing	Feature and training infos	Evaluation	Inference	Explanation
Supervised learning	<b>Linear and logistic regression</b>	LINEAR_REG, LOGISTIC_REG					
	<b>Deep Neural Network</b>	DNN_CLASSIFIER, DNN_REGRESSOR					
	<b>Boosted Tree</b>	BOOSTED_TREE_CLASSIFIER					
Unsupervised learning	<b>Kmeans</b>	KMEANS	TRANSFORM	ML.TRAINING_INFO	ML.EVALUATE	ML.PREDICT, ML.DETECT_ANOMALIES	ML.EXPLAIN_PREDICT, ML.GLOBAL_EXPLAIN, ML.ADVANCED_WEIGHT
	<b>Matrix Factorization</b>	MATRIX_FACTORIZATION				ML.RECOMMEND	
	<b>Autoencoder</b>	AUTOENCODER		ML.FEATURE_INFO		ML.PREDICT, ML.DETECT_ANOMALIES, ML.RECONSTRUCTION_LOSS	
Time Series	<b>ARIMA +</b>	ARIMA_PLUS			ML.EVALUATE, ML.ARIMA_EVALUATE	ML.FORECAST, ML.DETECT_ANOMALIES	ML.EXPLAIN_FORECAST
Imported Models	<b>Tensorflow</b>	TENSORFLOW				ML.PREDICT	

## MANUAL DATA PREPROCESSING

	TRANSFORM()
<ul style="list-style-type: none"> <li>• Manual preprocessing can be added using the <b>TRANSFORM</b> clause</li> <li>• Features listed in the <b>TRANSFORM</b> clause are the only one used to train the model</li> <li>• Preprocessing parameters are saved within the model for prediction</li> </ul>	
<b>Bucketing:</b> Separate numerical values into buckets or groups	ML.BUCKETIZE(numerical_expression, array_split_points[, exclude_boundaries])
<b>MinMax Scaler:</b> Transform all the values into the range [0,1]	ML.MIN_MAX_SCALER(numerical_expression) OVER()
<b>Standard Scaler:</b> Transform all the values to center the mean to 0 and standard deviation to 1	ML.STANDARD_SCALER(numerical_expression) OVER()
Other transformations	<a href="#">Other preprocessing functions</a>

```
CREATE MODEL `mydataset.mymodel`
OPTIONS
( MODEL_TYPE='LINEAR_REG',
  LS_INIT_LEARN_RATE=.15,
  L1_REG=1,
  MAX_ITERATIONS=5 ) AS
SELECT
column1, column2, column3, label
FROM `mydataset.mytable`
WHERE
column4 < 10
```

```
SELECT
*
FROM ML.PREDICT(
MODEL `mydataset.mymodel`,
(
SELECT
label,
column1,
column2
FROM `mydataset.mytable`
))
```

## EVALUATION

	ML.EVALUATE()
<b>General evaluation:</b> Return a single row containing common evaluation metrics corresponding to the trained model	ML.EVALUATE(MODEL `mydataset.mymodel`)
<b>Classification models evaluation:</b> <ul style="list-style-type: none"> <li>• ROC curve</li> <li>• Confusion matrix</li> </ul>	ML.ROC_CURVE(MODEL `mydataset.mymodel`) ML.CONFUSION_MATRIX(MODEL `mydataset.mymodel`)
<b>Time series models evaluation:</b>	ML.ARIMA_EVALUATE(MODEL `mydataset.mymodel`)
Training iterations details:	ML.TRAINING_INFO(MODEL `mydataset.mymodel`)

## INFERENCE

	ML.PREDICT(), ML.FORECAST(), ML.RECOMMEND()
<b>Predict:</b> Predict outcomes using the trained ML model	SELECT * FROM ML.PREDICT(MODEL `my-project.my_dataset.my_model`, (SELECT * FROM input_data))
<b>Forecast (time series):</b> Forecasts time series values from a trained ARIMA_PLUS or ARIMA model	SELECT * FROM ML.FORECAST(MODEL `mydataset.my_model`, STRUCT(30 AS horizon, 0.8 AS confidence_level))
<b>Recommend:</b> Generates a predicted rating for every user-item row combination for a matrix factorization model	SELECT * FROM ML.RECOMMEND(MODEL `mydataset.mymodel`, (SELECT user,item FROM `mydataset.my-table`))

```
SELECT
*
FROM
ML.EVALUATE(MODEL `mydataset.mymodel`)
```

```
CREATE MODEL `mydataset.mymodel`
TRANSFORM(
column1,
ML.BUCKETIZE(column2, [1, 2, 3]) AS column2,
ML.MIN_MAX_SCALER(column3) AS column3
label
)
OPTIONS(...)
AS SELECT
column1, column2, column3, label
FROM `mydataset.mytable`
```