

Big Query Predictive Models Creation	1
Step 1. Creating Linear Regression model using BiqQuery without EDA – LR_Simple1	1
Step 2. Evaluating the LR Simple1 model using quality	2
Step 3. Creating and training the Linear Regression Model using BigQuery with EDA - LR_with_features	4
Step 4. Evaluating the LR_with_features model	4
Step 5 Creating a new model using a different algorithm - Boosted Tree Classifier – XGB.....	6
Step 6. Evaluating the model we created – XGB.	6
Step 7. Comparing the models (LR simple1, LR with features и XGB) using the ML.EVALUATE function.....	7
Step 8. Predicting the results using LR simple1, LR with features, XGB models, with the help of ML.PREDICT function.	8
Step 9. XGB model and ML.FEATURE_IMPORTANCE function to see the importance of each field for the created model.	10

Big Query Predictive Models Creation

The task is to create predictive models using different Data Science platforms and compare the results of prediction choosing the best model for this particular task. I used the dataset from Kaggle.com <https://www.kaggle.com/tejashvi14/employee-future-prediction>

The task is to predict if an employee will leave the company or will keep working in the same company in the next 2 years, the model will be created based on the dataset from Kaggle. The task is a Binary Classification task.

We have the following data for the model creation:

1. Education
2. JoiningYear
3. City
4. PaymentTier
5. Age
6. Gender
7. EverBenched
8. ExperienceInCurrentDomain
9. LeaveOrNot – target field

Step 1. Creating Linear Regression model using BiqQuery without EDA – LR_Simple1

```

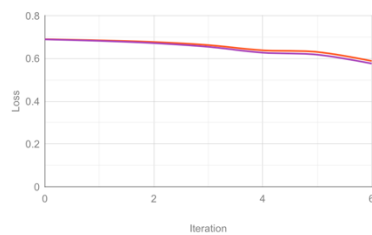
1 CREATE OR REPLACE MODEL `ML.Employee_LR_simple1`
2 OPTIONS(MODEL_TYPE = 'LOGISTIC_REG',
3         INPUT_LABEL_COLS = ['LeaveOrNot'],
4         MAX_ITERATIONS = 30,
5         LEARN_RATE_STRATEGY = 'LINE_SEARCH',
6         EARLY_STOP = TRUE,
7         MIN_REL_PROGRESS = 0.01,
8         DATA_SPLIT_METHOD = 'AUTO_SPLIT',
9         AUTO_CLASS_WEIGHTS = TRUE,
10        ENABLE_GLOBAL_EXPLAIN = TRUE,
11        CATEGORY_ENCODING_METHOD = 'DUMMY_ENCODING'
12 ) as
13
14 SELECT * FROM `my-project-dec9.ML.Train_Employee`

```

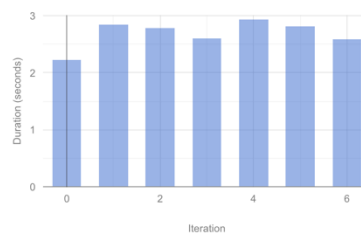
Step 2. Evaluating the LR Simple1 model using quality

View as ☒ Graphs ☐ Table

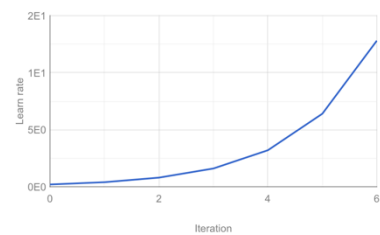
Loss

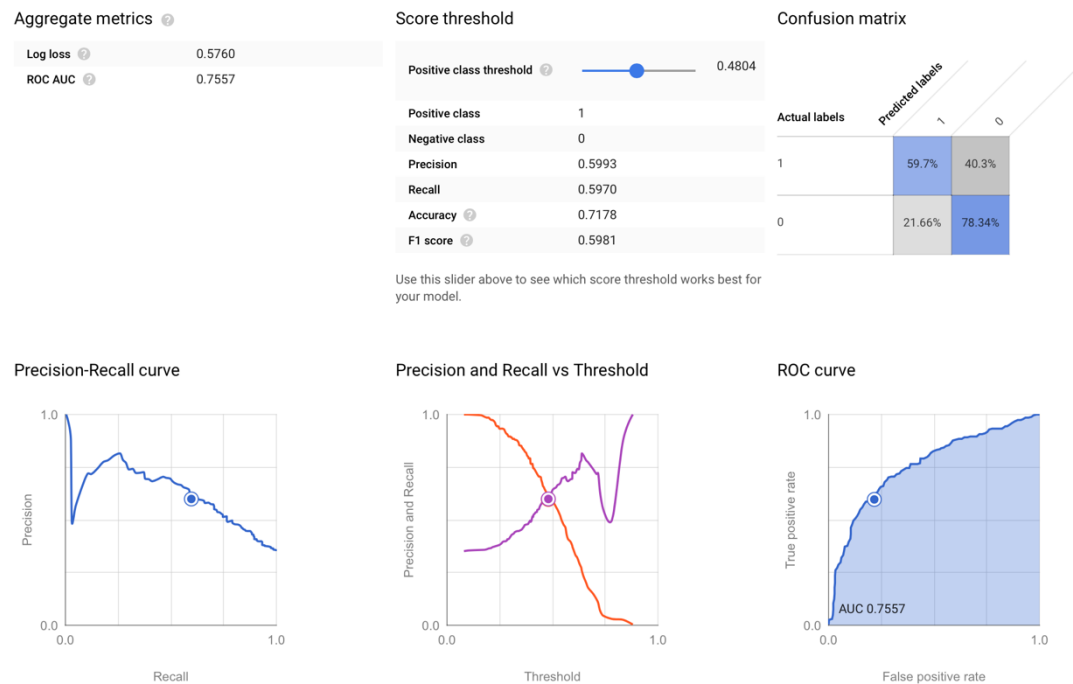


Duration (seconds)



Learn rate





Quality Metrics	Ideal Result	LR Simple1 Model's Result
Log Loss (Logarithmic Loss) Log-loss is indicative of how close the prediction probability is to the corresponding actual/true value (0 or 1 in case of binary classification). The more the predicted probability diverges from the actual value, the higher is the log-loss value.	The more the predicted probability diverges from the actual value, the higher is the log-loss value. Aiming for 0.	0.5760
ROC AUC (area under the curve of the receiver operating characteristic) AUC - ROC curve is a performance measurement for the classification problems at various threshold settings. ROC is a probability curve and AUC represents the degree or measure of separability. It tells how much the model is capable of distinguishing between classes. A receiver operating characteristic curve, or ROC curve, is a graphical plot that illustrates the diagnostic ability of a binary classifier system as its discrimination threshold is varied. The ROC curve is created by plotting the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings. The	The higher the AUC, the better the model is at predicting 0 classes as 0 and 1 classes as 1. Aiming to 1.	0.7557

<p>true-positive rate is also known as sensitivity, recall or probability of detection. The false-positive rate is also known as probability of false alarm and can be calculated as $(1 - \text{specificity})$. It can also be thought of as a plot of the power as a function of the Type I Error of the decision rule (when the performance is calculated from just a sample of the population, it can be thought of as estimators of these quantities). The ROC curve is thus the sensitivity or recall as a function of fall-out. In general, if the probability distributions for both detection and false alarm are known, the ROC curve can be generated by plotting the cumulative distribution function (area under the probability distribution from $-\infty$ to the discrimination threshold) of the detection probability in the y-axis versus the cumulative distribution function of the false-alarm probability on the x-axis.</p>		
<p><i>Model is not good enough yet, it needs to be fine-tuned.</i></p>		

Step 3. Creating and training the Linear Regression Model using BigQuery with EDA - LR_with_features

```

1 CREATE OR REPLACE MODEL `ML_LR_with_features`
2 TRANSFORM(
3 CASE WHEN Gender = 'male' THEN 1 ELSE 0 END AS Sex,
4 CASE WHEN Age BETWEEN 0 AND 21 THEN 1 ELSE 0 END AS Age_VeryYoung,
5 CASE WHEN Age BETWEEN 21 AND 25 THEN 1 ELSE 0 END AS Age_Young,
6 CASE WHEN Age BETWEEN 25 AND 28 THEN 1 ELSE 0 END AS Age_Young1,
7 CASE WHEN Age BETWEEN 28 AND 31 THEN 1 ELSE 0 END AS Age_Young2,
8 CASE WHEN Age BETWEEN 31 AND 36 THEN 1 ELSE 0 END AS Age_Middle,
9 CASE WHEN Age > 36 THEN 1 ELSE 0 END AS Age_Adult,
10 CASE WHEN ExperienceInCurrentDomain < 2 THEN 1 ELSE 0 END AS NoExperience,
11 CASE WHEN ExperienceInCurrentDomain BETWEEN 2 AND 4 THEN 1 ELSE 0 END AS LittleExperience,
12 CASE WHEN ExperienceInCurrentDomain > 4 THEN 1 ELSE 0 END AS PlentyExperience,
13 * EXCEPT( Gender, Age, ExperienceInCurrentDomain))
14 OPTIONS
15 (MODEL_TYPE = 'LOGISTIC_REG',
16 L1_REG = 0.2,
17 L2_REG = 0.6, |
18 LEARN_RATE_STRATEGY = 'CONSTANT',
19 LEARN_RATE = 0.5,
20 EARLY_STOP = TRUE,
21 AUTO_CLASS_WEIGHTS = TRUE,
22 ENABLE_GLOBAL_EXPLAIN = TRUE,
23 CATEGORY_ENCODING_METHOD = 'ONE_HOT_ENCODING'
24 ) AS
25 SELECT
26 LeaveOrNot AS label, Gender, Age, ExperienceInCurrentDomain
27 FROM
28 `my-project-dec9.ML.Train_Employee`

```

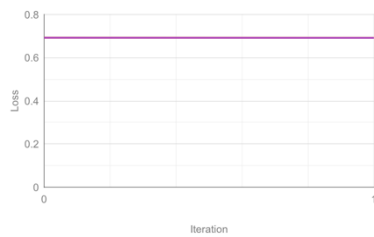
Run Save query Save view Schedule query More

This query will process 114.3 KB (ML) when run.

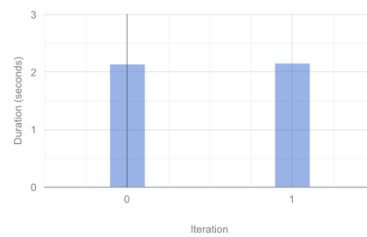
Step 4. Evaluating the LR_with_features model

View as ☒ Graphs ☐ Table

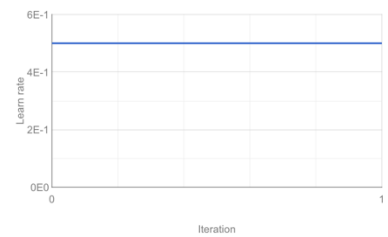
Loss



Duration (seconds)



Learn rate



Aggregate metrics

Log loss	0.6916
ROC AUC	0.4965

Score threshold

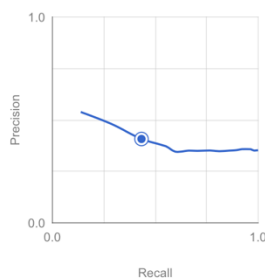
Positive class threshold	0.5079
Positive class	1
Negative class	0
Precision	0.4074
Recall	0.4331
Accuracy	0.5778
F1 score	0.4198

Use this slider above to see which score threshold works best for your model.

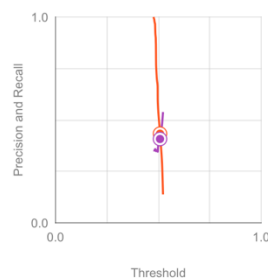
Confusion matrix

Actual labels	Predicted labels	
	1	0
1	43.31%	56.69%
0	34.33%	65.67%

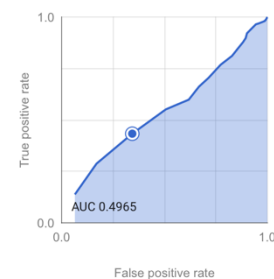
Precision-Recall curve



Precision and Recall vs Threshold



ROC curve



Quality Metrics	Ideal Result	LR with features Model's Result
Log Loss (Logarithmic Loss) Log-loss is indicative of how close the prediction probability is to the corresponding actual/true value (0 or 1 in case of binary classification). The more the predicted probability diverges from the actual value, the higher is the log-loss value.	The more the predicted probability diverges from the actual value, the higher is the log-loss value. Aiming for 0.	0.6916
ROC AUC (area under the curve of the receiver operating characteristic) AUC - ROC curve is a performance measurement for the classification problems at various threshold settings. ROC is a probability curve and AUC represents the degree or measure of separability. It	The higher the AUC, the better the model is at predicting 0 classes as 0 and 1 classes as 1. Aiming to 1.	0.4965

tells how much the model is capable of distinguishing between classes.		
<i>Metrics are worse, that means that the algorithm we selected is not optimal, we need to try different algorithms.</i>		

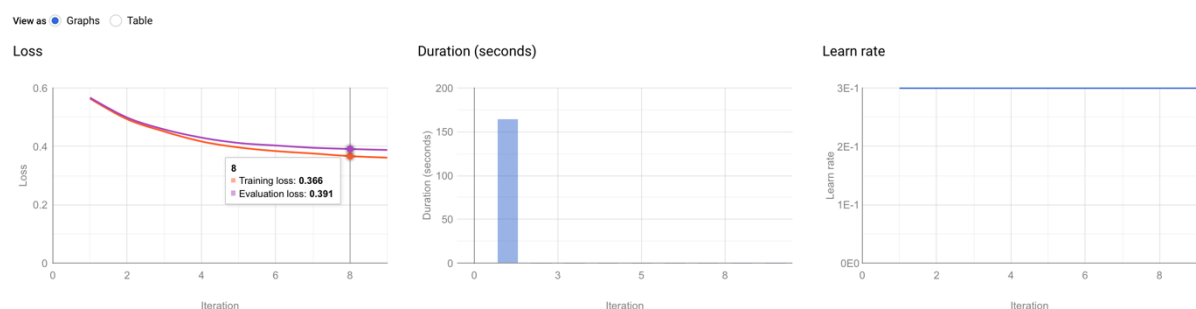
Step 5 Creating a new model using a different algorithm - Boosted Tree Classifier – XGB.

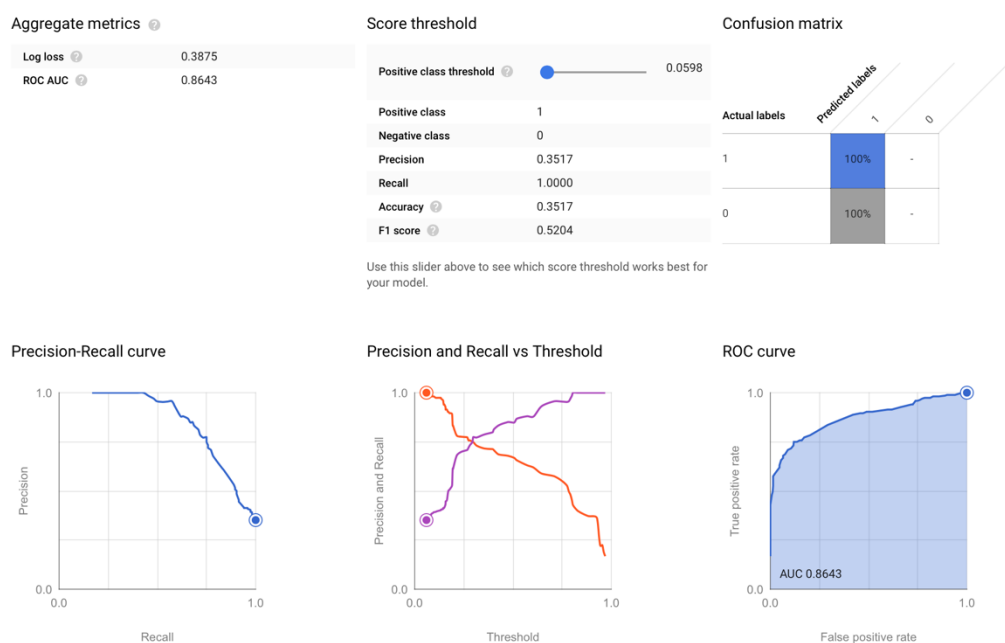
```

1 CREATE OR REPLACE MODEL `ML.Employee_XGB`
2 OPTIONS(MODEL_TYPE = 'BOOSTED_TREE_CLASSIFIER' ,
3         BOOSTER_TYPE = 'GBTREE',
4         TREE_METHOD = 'HIST',
5         MAX_TREE_DEPTH = 5,
6         L2_REG = 0.4,
7         EARLY_STOP = TRUE,
8         INPUT_LABEL_COLS = ['LeaveOrNot'],
9         MAX_ITERATIONS = 50,
10        ENABLE_GLOBAL_EXPLAIN = TRUE
11 ) as
12 SELECT *
13 FROM `my-project-dec9.ML.Train_Employee`

```

Step 6. Evaluating the model we created – XGB.





Quality Metrics	LR simple1	LR with features	XGB
Log Loss (Logarithmic Loss) Log-loss is indicative of how close the prediction probability is to the corresponding actual/true value (0 or 1 in case of binary classification). The more the predicted probability diverges from the actual value, the higher is the log-loss value.	0.5760	0.6916	0.3875
ROC AUC (area under the curve of the receiver operating characteristic) AUC - ROC curve is a performance measurement for the classification problems at various threshold settings. ROC is a probability curve and AUC represents the degree or measure of separability. It tells how much the model is capable of distinguishing between classes.	0.7557	0.4965	0.8643
<i>XGB model is demonstrating better results comparing to the Linear Regression models.</i>			

Step 7. Comparing the models (LR simple1, LR with features и XGB) using the MLEVALUATE function.

```
1 SELECT
2 *
3 FROM
4 ML.EVALUATE (MODEL `my-project-dec9.ML.Employee_LR_simple1`)
```

Run Save query Save view Schedule query More

Query results SAVE RESULTS EXPLORE DATA

Query complete (0.2 sec elapsed, 0 B processed)

Job information Results JSON Execution details

Row	precision	recall	accuracy	f1_score	log_loss	roc_auc
1	0.6408163265306123	0.585820895522388	0.7388451443569554	0.6120857699805068	0.5760181762655407	0.7556803196803197

```
1 SELECT
2 *
3 FROM
4 ML.EVALUATE (MODEL `my-project-dec9.ML.LR_with_features`)
```

Run Save query Save view Schedule query More

Query results SAVE RESULTS EXPLORE DATA

Query complete (0.2 sec elapsed, 0 B processed)

Job information Results JSON Execution details

Row	precision	recall	accuracy	f1_score	log_loss	roc_auc
1	0.4074074074074074	0.4330708661417323	0.5777777777777777	0.4198473282442748	0.6916028862431451	0.4964885114885115

```
1 SELECT
2 *
3 FROM
4 ML.EVALUATE (MODEL `my-project-dec9.ML.Employee_XGB`)
```

Run Save query Save view Schedule query More

Query results SAVE RESULTS EXPLORE DATA

Query complete (0.4 sec elapsed, 0 B processed)

Job information Results JSON Execution details

Row	precision	recall	accuracy	f1_score	log_loss	roc_auc
1	0.8483412322274881	0.667910447761194	0.8412073490813649	0.7473903966597077	0.3875202999084784	0.8642887112887113

Step 8. Predicting the results using LR simple1, LR with features, XGB models, with the help of ML.PREDICT function.

This query will process 56.8 KB when run.

[EXPLORE DATA](#)

Job information **Results** JSON Execution details

Job information **Results** JSON Execution details

This query will process 56.6 KB when run.

[EXPLORE DATA](#)

[Job information](#) [Results](#) [JSON](#) [Execution details](#)

This query will process 81.8 KB when run.

[EXPLORE DATA](#)

Job information **Results** JSON Execution details

Step 9. XGB model and ML.FEATURE_IMPORTANCE function to see the importance of each field for the created model.

1 SELECT * FROM ML.FEATURE_IMPORTANCE (MODEL `my-project-dec9. ML.Employee_XGB`)

Run

Save query

Save view

Schedule query

More

This query will process 0 B when run.

Query results

SAVE RESULTS

EXPLORE DATA

Query complete (0.6 sec elapsed, 0 B processed)

Job information

Results

JSON

Execution details

Row	feature	importance_weight	importance_gain	importance_cover
1	Education	64	7.955924218750001	318.4442562499999
2	JoiningYear	42	31.35683533333333	631.0134000000002
3	City	50	10.45532284	250.85576800000007
4	PaymentTier	48	13.505746162499994	434.2616554166666
5	Age	28	2.9830879999999986	102.73426285714284
5	Gender	32	4.452454093750001	364.91915875
7	EverBenched	16	2.3583788750000005	89.7641
3	ExperienceInCurrentDomain	22	3.252678545454546	174.85523636363635

The most important features that have an impact of the employee’s decision to quit are the following features:

- **Education;**
- **Joining Year;**
- **City;**
- **Payment Tier;**
- **Age.**