

# Robi Datathon 2.0: Pre-Assessment

## Description

Welcome to Robi Datathon 2.0: Pre-assessment round.

In this round, you will be dealing with a classification problem for the given dataset.

Rules:

- Join the competition when it opens
- Add your team members to your team
- Get the dataset and submission templates
- Train your models on train set and predict on test set
- Submit the submission file in Kaggle competition page
- Upload the final submission file and code in the specific format

## Evaluation

The evaluation metric for this competition is

[AUC](<https://developers.google.com/machine-learning/crash-course/classification/roc-and-auc>).

AUC stands for **"Area under the ROC Curve."** That is, AUC measures the entire two-dimensional area underneath the entire ROC curve (think integral calculus) from (0,0) to (1,1). AUC ranges in value from 0 to 1. A model whose predictions are 100% wrong has an AUC of 0.0; one whose predictions are 100% correct has an AUC of 1.0. AUC will be calculated based on the predicted binary label(0 or 1) with respect to the test dataset.

**Kaggle Link :**

<https://www.kaggle.com/competitions/robi-datathon-2-pre-assessment/overview>

---

## ROC curve

An **ROC curve (receiver operating characteristic curve)** is a graph showing the performance of a classification model at all classification thresholds. This curve plots two parameters:

- True Positive Rate
- False Positive Rate

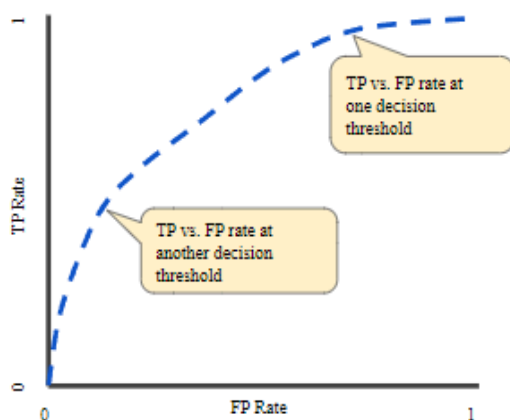
**True Positive Rate (TPR)** is a synonym for recall and is therefore defined as follows:

$$TPR = \frac{TP}{TP + FN}$$

**False Positive Rate (FPR)** is defined as follows:

$$FPR = \frac{FP}{FP + TN}$$

An ROC curve plots TPR vs. FPR at different classification thresholds. Lowering the classification threshold classifies more items as positive, thus increasing both False Positives and True Positives. The following figure shows a typical ROC curve.

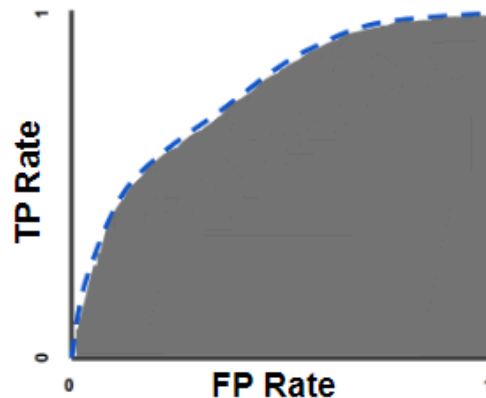


**Figure 4. TP vs. FP rate at different classification thresholds.**

To compute the points in an ROC curve, we could evaluate a logistic regression model many times with different classification thresholds, but this would be inefficient. Fortunately, there's an efficient, sorting-based algorithm that can provide this information for us, called AUC.

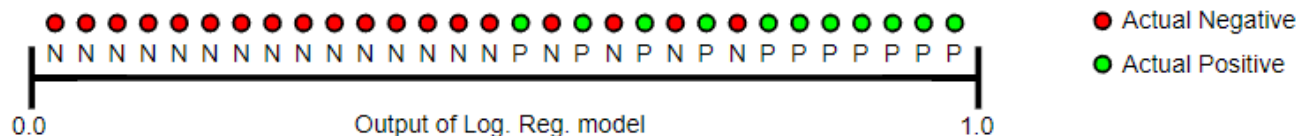
## AUC: Area Under the ROC Curve

**AUC** stands for "Area under the ROC Curve." That is, AUC measures the entire two-dimensional area underneath the entire ROC curve (think integral calculus) from (0,0) to (1,1).



**Figure 5. AUC (Area under the ROC Curve).**

AUC provides an aggregate measure of performance across all possible classification thresholds. One way of interpreting AUC is as the probability that the model ranks a random positive example more highly than a random negative example. For example, given the following examples, which are arranged from left to right in ascending order of logistic regression predictions:



**Figure 6. Predictions ranked in ascending order of logistic regression score.**

AUC represents the probability that a random positive (green) example is positioned to the right of a random negative (red) example.

AUC ranges in value from 0 to 1. A model whose predictions are 100% wrong has an AUC of 0.0; one whose predictions are 100% correct has an AUC of 1.0.

Source : <https://developers.google.com/machine-learning/crash-course/classification/roc-and-auc>