# 机器学习实验评分和分析教程-财务重点

## 8.任务6：ER：ROC

在"诊断"页面上，单击**ROC曲线**。将会出现类似于以下图像的图像：



回顾一下，ROC曲线显示以下内容：

- 它显示了灵敏度（真阳性率或TPR）和特异性（1-FPR或假阳性率）之间的权衡。敏感性的任何提高都将伴随特异性的降低。
- 曲线越接近ROC空间的左边界和上边界，模型越精确。
- 曲线越接近ROC空间的45度对角线，模型的精度就越差。
- 切点处切线的斜率给出了该测试值的似然比（LR）。您可以在上方的图表中查看。
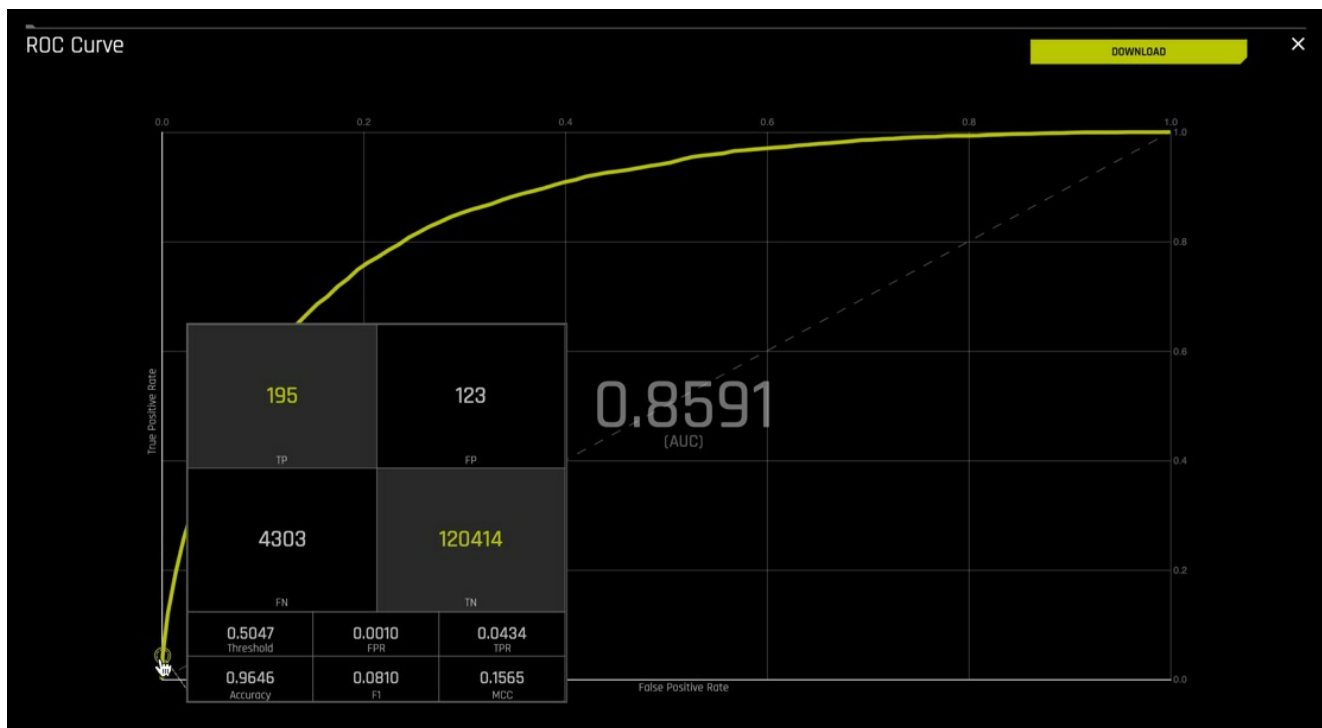- 曲线下的面积是模型准确性的度量。

回到Freddie Mac数据集，即使该模型使用对数损失进行评分以补偿误差，我们仍然可以查看ROC曲线结果，并查看其是否支持我们从混淆矩阵和分数部分的分析得出的结论。诊断页面。

1.根据无人驾驶AI模型为实验生成的ROC曲线，确定AUC。回想一下，理想的分类模型的AUC为1。

2.对于曲线上的以下每个点，将鼠标悬停在下面的每个点上，即可确定正阳性率，假阳性率和阈值，如下图所示：

- 最佳精度

- 最佳F1
- 最佳我的客户中心



回想一下，对于二进制分类问题，准确性是做出的正确预测的数量与所做的所有预测的比率。概率将转换为预测的类别，以定义阈值。对于此模型，已确定在阈值.5098处发现最佳精度。

在此阈值下，模型预测：

- TP = 1 = 195个案例预计为违约和违约
- TN = 0 = 120,414个案例被预测为未违约且未违约
- FP = 1 = 123个案例被预测为违约且未违约
- FN = 0 = 4,303个案例，预计没有违约和被违约

3. From the AUC, Best MCC, F1, and Accuracy values from the ROC curve, how would you qualify your model, is it a good or bad model? Use the key points below to help you asses the ROC Curve.

Remember that for the **ROC** curve:

- The perfect classification model has an AUC of 1
- MCC is measured in the range between -1 and +1 where +1 is the perfect prediction, 0 no better than a random prediction and -1 all incorrect predictions.
- F1 is measured in the range of 0 to 1, where 0 means that there are no true positives, and 1 when there is neither false negatives nor false positives or perfect precision and recall.

- Accuracy is measured in the range of 0 to 1, where 1 is perfect accuracy or perfect classification, and 0 is poor accuracy or poor classification.

**Note:** If you are not sure what AUC, MCC, F1, and Accuracy are or how they are calculated review the concepts section of this tutorial.

## New Model with Same Parameters

In case you were curious and wanted to know if you could improve the accuracy of the model, this can be done by changing the scorer from Logloss to Accuracy.

1. To do this, click on the **Experiments** page.

2. Click on the experiment you did for task 1 and select **New Model With Same Params**



An image similar to the one below will appear. Note that this page has the same settings as the setting in Task 1. The only difference is that on the **Scorer** section **Logloss** was updated to **Accuracy**. Everything else should remain the same.

3. If you haven't done so, select **Accuracy** on the scorer section then select **Launch Experiment**

Similarly to the experiment in Task 1, wait for the experiment to run. After the experiment is done running, a similar page will appear. Note that on the summary located on the bottom right-side both the validation and test scores are no longer being scored by **Logloss** instead by **Accuracy**.



We are going to use this new experiment to run a new diagnostics test. You will need the name of the new experiment. In this case, the experiment name is **1.Freddie Mac Classification Tutorial**.

4. Go to the **Diagnostics** tab.

5. Once in the **Diagnostics** page, select **+Diagnose Model**

6. In the **Create new model diagnostics** :

1. Click on Diagnosed Experiment then select the experiment that you completed in Task in this case the experiment name is **1.Freddie Mac Classification Tutorial**
2. Click on Dataset then select the freddie_mac_500_test dataset
3. Initiate the diagnostics model by clicking on **Launch Diagnostics**



7. After the model diagnostics is done running a new diagnostic will appear

8. Click on the new diagnostics model. On the **Scores** section observe the accuracy value. Compare this Accuracy value to the Accuracy value from task 6.

9. Next, locate the new ROC curve and click on it. Hover over the **Best ACC** point on the curve. An image similar to the one below will appear:

How much improvement did we get from optimizing the accuracy via the scorer?

The new model predicted:

- Threshold = .5129
- TP = 1 = 194 cases predicted as defaulting and defaulted
- TN = 0 = 120,419 cases predicted as not defaulting and did not default
- FP = 1 = 118 cases predicted as defaulting and did not default
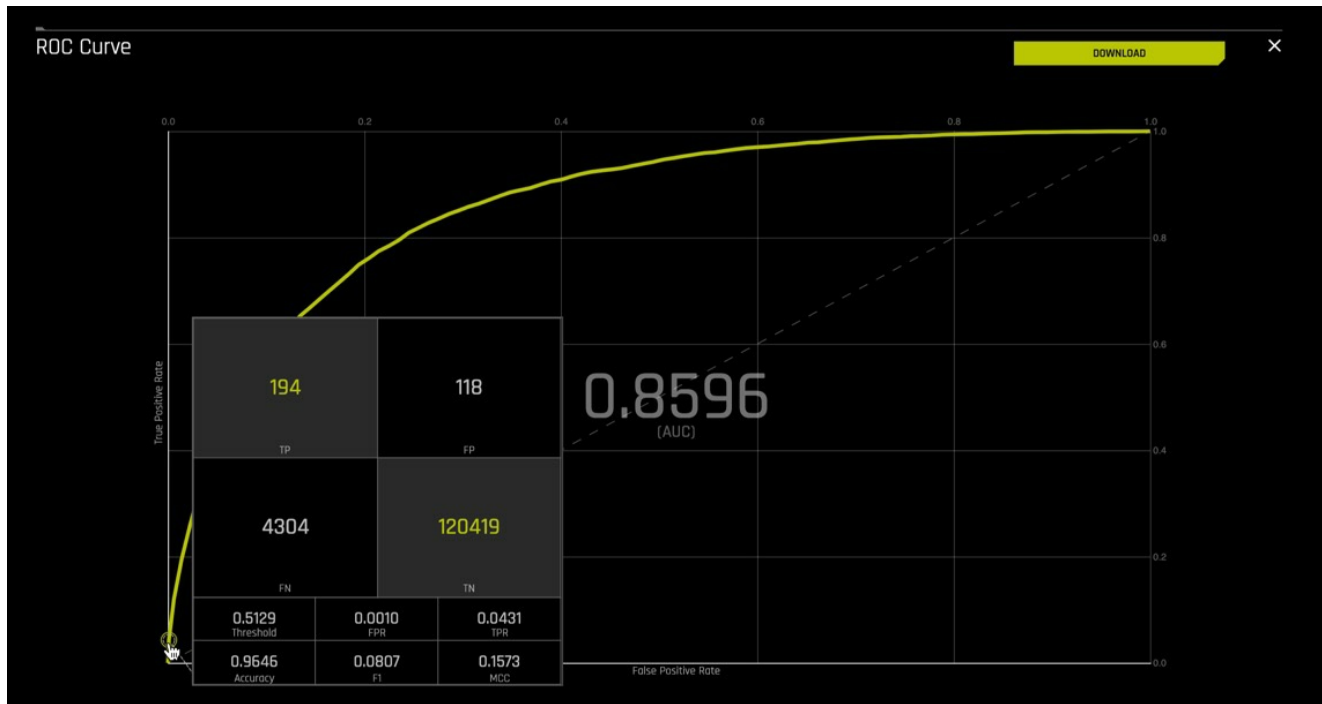- FN = 0 = 4,304 cases predicted not to default and defaulted

The first model predicted:

- Threshold = .5047
- TP = 1 = 195 cases predicted as defaulting and defaulted
- TN = 0 = 120,414 cases predicted as not defaulting and did not default
- FP = 1 = 123 cases predicted as defaulting and did not default
- FN = 0 = 4,303 cases predicted to not default and defaulted

The threshold for best accuracy changed from .5047 for the first diagnostics model to .5129 for the new model. This increase in threshold improved accuracy or the number of correct predictions made as a ratio of all predictions made. Note, however, that while the number of FP decreased the number of FN increased. We were able to reduce the number of cases that were predicted to falsy default, but in doing so, we increased the number of FN or cases that were predicted not to default and did.

The takeaway is that there is no win-win; sacrifices need to be made. In the case of accuracy, we increased the number of mortgage loans, especially for those who were denied a mortgage because they were predicted to default when, in reality, they did not. However, we also increased the number of cases that should not have been granted a loan and did. As a mortgage lender, would you prefer to reduce the number of False Positives or False Negatives?

10. Exit out of the ROC curve by clicking on the **x** located at the top-right corner of the plot, next to the **Download** option

## Deeper Dive and Resources