# AutoML Modeling Report

*Italo Farfán*

## Binary Classifier with Clean/Balanced Data

| | |
|---|---|
| **Train/Test Split**<br>How much data was used for training? How much data was used for testing? | 80% for train and 20 % testing (10% validation and 10% test) as show in the figure below.<br><br> |
| **Confusion Matrix**<br>What do each of the cells in the confusion matrix describe? What values did you observe (include a screenshot)? What is the true positive rate for the "pneumonia" class? What is the false positive rate for the "normal" class? | <br><br>The Confusion Matrix represents the outcome of model evaluation. We can see the corrects and incorrect predictions.<br>The True Positive Rate = 100%<br>The False Positive Rate for the "normal"= 10% |
| **Precision and Recall**<br>What does precision measure? | Precision measure the ratio of true positives to predicted positives. |

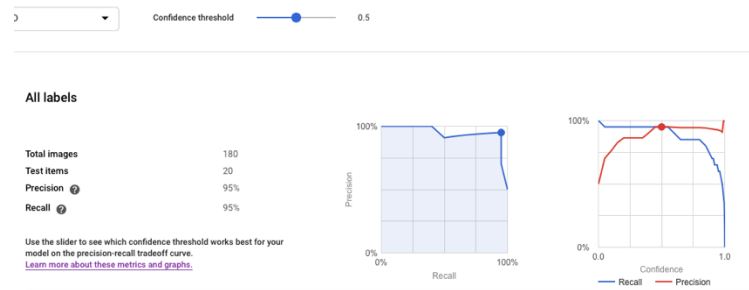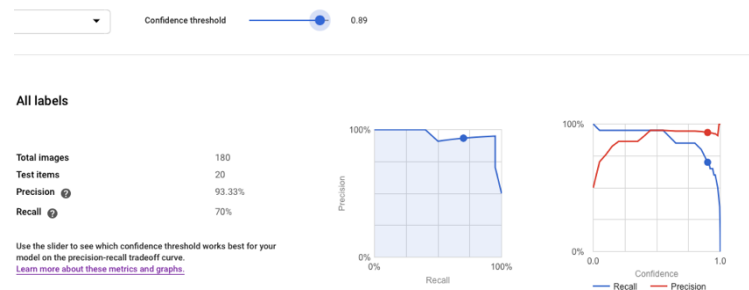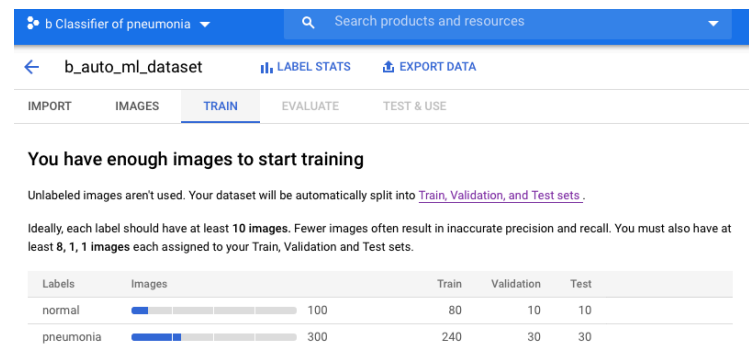| | |
|---|---|
| What does recall measure? What precision and recall did the model achieve (report the values for a score threshold of 0.5)? | Recall measure the ratio of true positives to actual positives.<br><br>Results of my model with Threshold of 0.5:<br>Precision = 95%<br>Recall = 95 %<br><br>For other hand I like so much the definition from Wikipedia: Recall in this context is also referred to as the true positive rate or sensitivity, and precision is also referred to as positive predictive value (PPV); other related measures used in classification include true negative rate and accuracy. True negative rate is also called specificity.<br><br><br><br>Note: A high precision model produces fewer false positives. A high recall model produces fewer false negatives |
| **Score Threshold**<br>When you increase the threshold what happens to precision? What happens to recall? Why? | Precision goes down and Recall goes down when I increase the threshold to 0.89 but if I continue increase the Threshold Precision goes up and Recall goes down as show in the figure below.<br><br> |

# Binary Classifier with Clean/Unbalanced Data

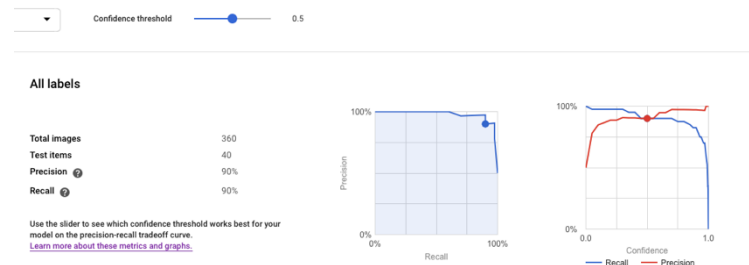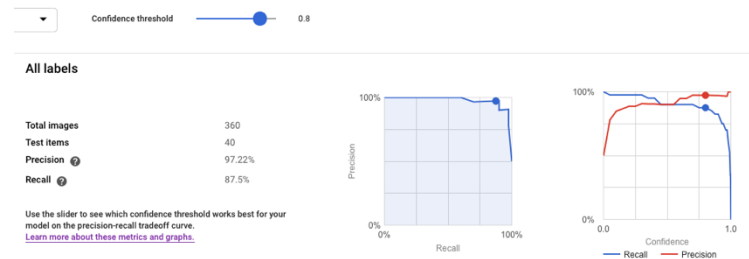| | |
|---|---|
| **Train/Test Split**<br>How much data was used for training? How much data was used for testing? | 80% for train and 20 % testing (10% validation and 10% test) as show in the figure below.<br><br> |
| **Confusion Matrix**<br>How has the confusion matrix been affected by the unbalanced data? Include a screenshot of the new confusion matrix. | I see two important changes: there is an increase of False Positives and there is an increase of True Positives.<br><br> |
| **Precision and Recall**<br>How have the model's precision and recall been affected by the unbalanced data (report the values for a score threshold of 0.5)? | Precision and Recall decreased 5% with Threshold of 0.5.<br>Precision = 90%<br>Recall = 90%<br><br> |
| **Unbalanced Classes**<br>From what you have observed, how do unbalanced classed affect a machine learning model? | The model improved the predictions of pneumonia cases This is easy to see in the matrix and it is obvious because now we have more samples of pneumonia cases but other hand unblanced classed increase the falses positives and |

this is very bad. If we were desing a model to predict cancer and help to doctor to take a decission to operate or not, have an increase of Falses Positives is so bad.

We may use a Threshold of 0.8 to help our model but the correct way is use a balanced classed.
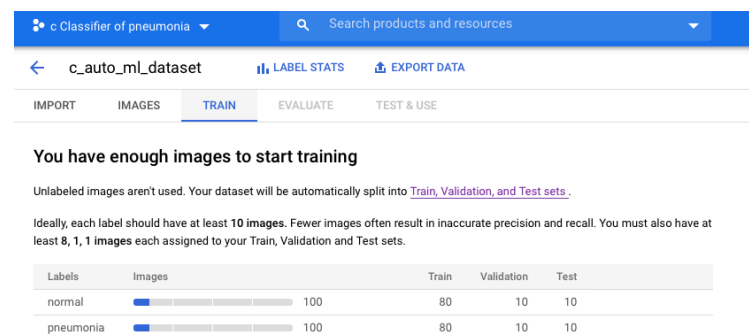


# Binary Classifier with Dirty/Balanced Data

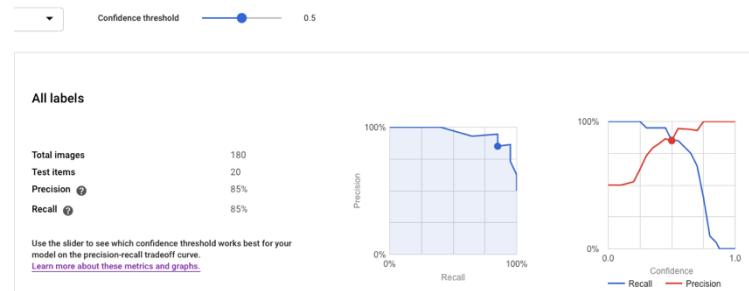| | |
|---|---|
| **Confusion Matrix**<br>How has the confusion matrix been affected by the dirty data? Include a screenshot of the new confusion matrix. | I see two important changes: there is an increase of False Positives and there is a decrease of True negative. This is very interesting because in my first model the False Negatives are 0%.<br><br><br><br>I used 80% for train and 20 % testing (10% validation and 10% test) as show in the figure below.<br><br> |

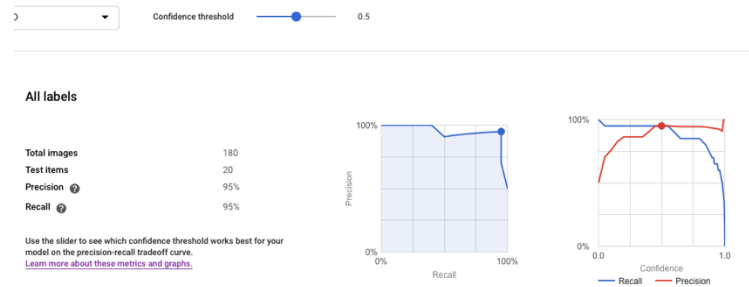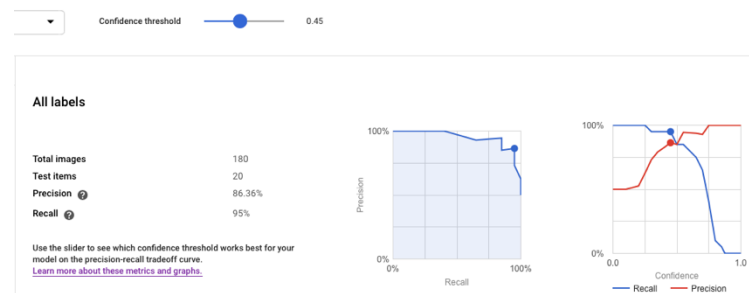| | |
|---|---|
| **Precision and Recall**<br>How have the model's precision and recall been affected by the dirty data (report the values for a score threshold of 0.5)? Of the binary classifiers, which has the highest precision? Which has the highest recall? | The results of my model with Threshold of 0.5:<br>Precision = 85%<br>Recall = 85 %<br><br>Binary Classifier with Dirty/Balanced Data<br><br>Of all binary classifiers the first model, The Binary Classifier with Clean/Balanced Data with 0.5 of Threshold, it has the highest precision and highest recall both 95%<br><br>Binary Classifier with Clean/Balanced Data |
| **Dirty Data**<br>From what you have observed, how does dirty data affect a machine learning model? | In general, the performance is worse and I can see and understand the importance of good labeled process. We may use a Threshold of 0.45 to increase the recall but the precision still is low. |

# 3-Class Model

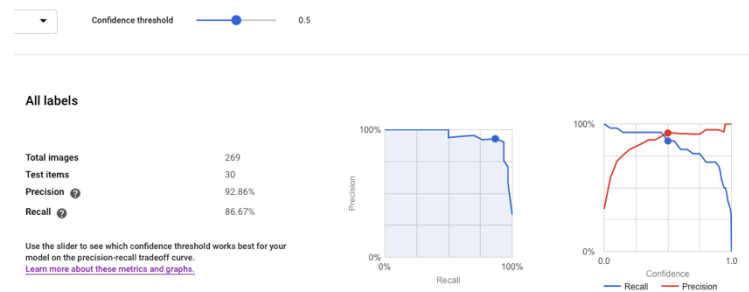| Confusion Matrix | |
|---|---|
| **Confusion Matrix**<br>Summarize the 3-class confusion matrix. Which classes is the model most likely to confuse? Which class(es) is the model most likely to get right? Why might you do to try to remedy the model's "confusion"? Include a screenshot of the new confusion matrix. | First it all I used 299 images. Something was wrong when I was uploading the 300 images and I had the same trouble many times, one of the images of viral pneumonia didn't upload correctly.<br><br><br><br>This is a good matrix because only two labels are bad predicted, but the first binary model still is better.<br><br>• The model most likely to confuse is the class "normal"<br>• The model most likely to get right is the class "viral pneumonia"<br><br><br><br>To try to remedy the model's "confusion" we may use a Threshold of 0.47. But I would like to meet with a specialist to analyze images of viral pneumonia and images of healthy children to try to find some characteristic that will help us to better label the images. |

| | | |
|---|---|---|
| Total images | 269 | |
| Test items | 30 | |
| Precision ⃝ | 90.32% | |
| Recall ⃝ | 93.33% | |

Use the slider to see which confidence threshold works best for your
model on the precision-recall tradeoff curve.
Learn more about these metrics and graphs.

---

**Precision and Recall**
What are the model's precision and recall? How are these values calculated (report the values for a score threshold of 0.5)?

The results of my model with Threshold of 0.5:
Precision = 92.86%
Recall = 86.67 %



**All labels**

| | | |
|---|---|---|
| Total images | 269 | |
| Test items | 30 | |
| Precision ⃝ | 92.86% | |
| Recall ⃝ | 86.67% | |

Use the slider to see which confidence threshold works best for your
model on the precision-recall tradeoff curve.
Learn more about these metrics and graphs.

These values are calculated automatically by Google Cloud Platform, but we can calculate how we are learned in class:

- **Precision:** Ratio of true positives to predicted positives.

$$Precision = \frac{TP}{(TP + FP)}$$

- **Recall:** Ratio of true positives to actual positives.

$$Recall = \frac{TP}{(TP + FN)}$$

---

**F1 Score**
What is this model's F1 score?

First, remember that the F1-score is a function of precision and recall. We already learned how to compute the per-class precision and recall.
Here is a summary of the precision and recall for our three classes:

| Class | Precission | Recall | F1-Score |
|---|---|---|---|
| bacterial pneumonia | 82.82% | 90% | 85.71% |
| viral pneumonia | 90.91% | 100% | 95.24% |
| normal | 100% | 80% | 88.89% |

The **macro-average F1-score, or F1** for short, is computed as a simple arithmetic mean of our per-class F1-scores.
So F1 = 89.95%