# AutoML Modeling Report

*Georgi Kardzhaliyski*

## Binary Classifier with Clean/Balanced Data

| Train/Test Split<br>How much data was used for training? How much data was used for testing? | Total images: 180<br>Test items: 20<br><br>AutoML Vision automatically uses 80% of your images for training, 10% for validating, and 10% for testing. |
|---|---|
| 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>A confusions matrix shows how often the model classified each label correctly (in blue), and which labels were most often confused for that label (in gray). It is a table that is often used to describe the performance of a classification model on a set of test data for which the true values are known. The number of correct and incorrect predictions are summarized with count values and broken down by each class.<br><br>A confusion matrix has two rows and two columns, and shows how many normal samples were predicted as Normal or Pneumonia (the first column), and how many Pneumonia samples were predicted as Normal or Pneumonia (the second column) |

There are two labels in the confusion matrix: true / actual label and predicted label, with normal and pneumonia being the two categories.
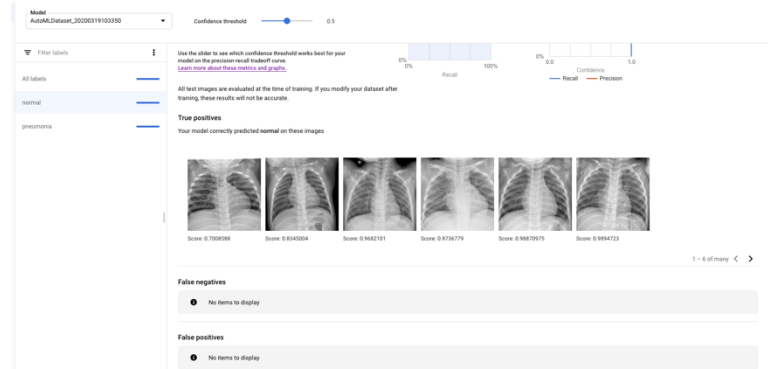
Every time the classifier makes a prediction, one of the cells in the table is incremented by one. By the end of the process, we can see exactly how our classifier performed.

https://towardsdatascience.com/multi-class-metrics-made-simple-part-i-precision-and-recall-9250280bddc2

We use the labels in the confusion matrix to understand the ways in which our classification model is confused when it makes predictions. In this case, when an image was incorrectly classified as normal when it should've been recognized as pneumonia. And the other way around.

The results from the confusion matrix show that there were no errors, all normal and pneumonia images were classified correctly. The errors are at 0%, as it can be seen from the screenshot of the confusion matrix. The classifier was perfect. All normal images were classified as normal, and all pneumonia images were classified as pneumonia.

The table gives us an insight into where we can improve our training to increase the model's accuracy.
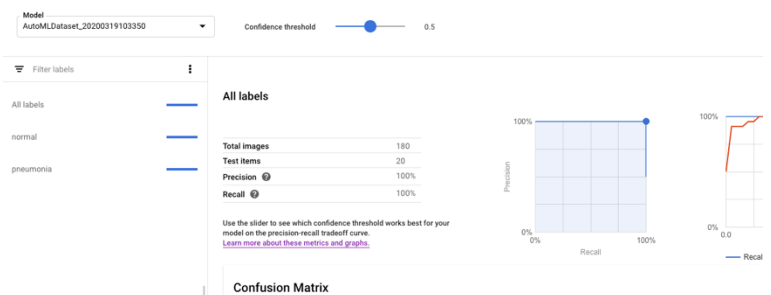
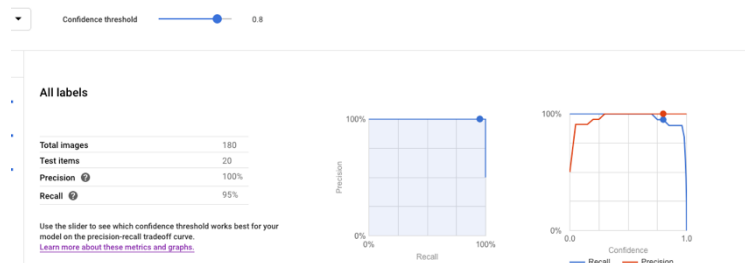| | |
|---|---|
| | <br><br>A high precision produces fewer false positives.<br>A high recall produces fewer false negatives, |
| **Precision and Recall**<br>What does precision measure? What does recall measure? What precision and recall did the model achieve (report the values for a score threshold of 0.5)? | Precision and recall help us understand how well our model is capturing information, and how much it's leaving out. Precision tells us, from all the test examples that were assigned a label, how many actually were supposed to be categorized with that label. Recall tells us, from all the test examples that should have had the label assigned, how many were actually assigned the label.<br><br> |

| | |
|---|---|
| **Score Threshold**<br>When you increase the threshold what happens to precision? What happens to recall? Why? | When I increased the confidence threshold, the prevision % stayed the same (100%), but the recall % went down from 100% to 95%.<br><br>The score threshold refers to the level of confidence the model must have to assign a category to a test item. If your score threshold is low, your model will classify more images, but runs the risk of misclassifying a few images in the process. If your score threshold is high, your model will classify fewer images, but it will have a lower risk of misclassifying images.<br><br> |

## Binary Classifier with Clean/Unbalanced Data

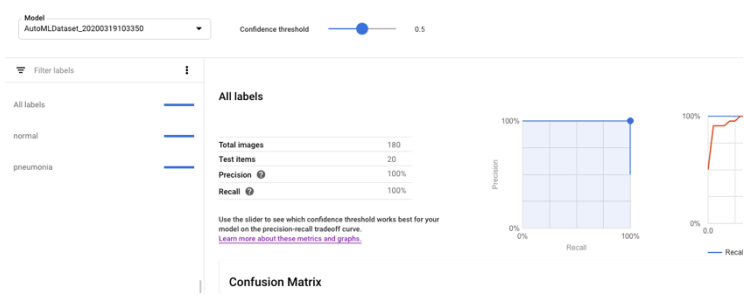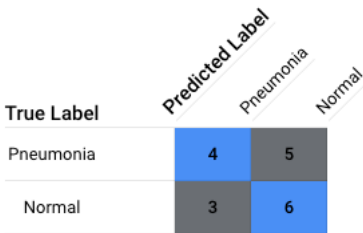| | |
|---|---|
| **Train/Test Split**<br>How much data was used for training? How much data was used for testing? |  |
| **Confusion Matrix**<br>How has the confusion matrix been affected by the unbalanced data? Include a screenshot of the new confusion matrix. | <br><br>There are more misclassified images when there is unbalanced data. Unbalanced data causes the model to skew towards the incorrect classification, because there |

| | is more noise in the data. |
|---|---|
| **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)? | <br><br>The precision and recalled are lower than when the data were clean and balanced. |
| **Unbalanced Classes**<br>From what you have observed, how do unbalanced classed affect a machine learning model? | Models will only learn about data that they are trained with. When we use unbalanced data, this causes the model to skew towards a particular outcome. |

# 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. | <br><br>More pneumonia images were misclassified as normal. More normal images were misclassified as pneumonia images. |
| **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? | <br><br>Precision and recall are both at 55.56%. The average precision is 0.759. It measures how well the model performs across all score thresholds. |

The precision % for the normal is 55.55, while the precision % for the pneumonia is 57.14%.

Of the binary classifiers (clean/balanced, dirty/balanced and clean/unbalanced), the precision and recall values at 0.5 score threshold for each of them were as follows:

Clean/balanced:
1. Precision: 100%
2. Recall: 100%

Dirty/balanced:
1. Precision: 55.56%
2. Recall: 55.56%

Clean/unbalanced:
1. Precision: 90%
2. Recall: 90%

The clean/balanced classifier had the highest score.

**Dity_Balanced_20200320065721**

Average precision
0.759
Precision* 55.56%
Recall* 55.56%
* Using a score threshold of 0.5

Model ID ICN8889573500801515520
Created Mar 20, 2020, 6:57:27 PM
Base model None
Data 182 images
Model type Cloud
Train cost 16 node hours
Deployment state Deployed

| **Dirty Data** | Mislabeled, dirty data, significantly impacts the model performance, the precision and recall metrics, confusion metrics. |
|---|---|
| From what you have observed, how does dirty data affect a machine learning model? | |

# 3-Class Model

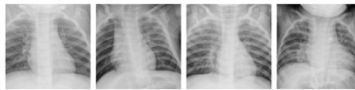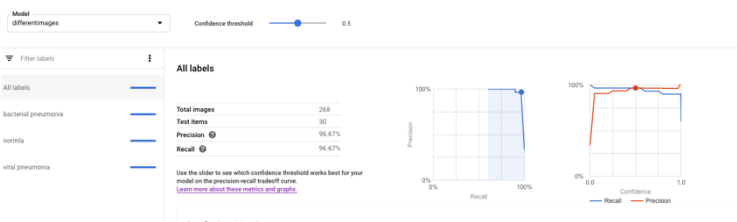| | |
|---|---|
| **Confusion Matrix**<br>Summarize the 3-class confusion matrix. Which classes is the model most likely to confuse? Which class(es) is the model most likely to get right? Why might you do to try to remedy the model's "confusion"? Include a screenshot of the new confusion matrix. | **Confusion Matrix** ☐ Show percentage ⬇<br><br>This table shows how often the model classified each label correctly (in blue), and which labels were most often confused for that label (in gray). Note that this table is limited to the 10 most confused labels. You can download the entire confusion matrix as a CSV file.<br><br>*Predicted Label*<br>*bacterial pneumonia / normla / viral pneumonia*<br><br>**True Label**<br>bacterial pneumonia — 9, -, 1<br>normla — -, 10, -<br>viral pneumonia — -, -, 10<br><br>The model is most likely to confuse the viral pneumonia. It classified the normal and viral pneumonia at 100%. In an ideal model, all the values on the diagonal will be high, and all the other values will be low. This shows that the desired categories are being identified correctly.<br>To remedy the confusion would be to collect more data. A larger dataset might expose a different and perhaps more balanced perspective on the classes. I can add more images of the under-represented class and delete images of the over-represented class. This could give a boost in the accuracy measures. I might also look at using a new platform with different algorithms to spot-check the results I'm getting. |
| **Precision and Recall**<br>What are the model's precision and recall? How are these values calculated (report the values for a score threshold of 0.5)? | Model: differentimages    Confidence threshold — 0.5<br><br>≡ Filter labels<br>All labels<br>All labels<br>bacterial pneumonia<br>normla<br>viral pneumonia<br><br>**All labels**<br>Total images 268<br>Test items 30<br>Precision ❓ 96.67%<br>Recall ❓ 96.67%<br>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.<br><br>The value of the precision in the multiclass classification is 96.67%; the value of the recall is 96.67% as well.<br><br>**All labels:**<br>Precision = 96.67%<br>Recall = 96.67% |

**Bacterial pneumonia:**
Precision = 100%
Recall = 90%

**Normal:**
Precision = 100%
Recall = 100%

**Viral pneumonia:**
Precision: 90.91%
Recall: 100%

Precision is calculated using this formula for a binary case problem, where:

TP = True Positive
FP = False Positive
FN = False Negative
TN = True Negative
P=Precision
R=Recall

$P = TP / (TP + FP)$

Recall is calculated using this formula:
$R= TP /(TP+FN)$

Similar to a binary case, we can calculate the precision and recall for each of the classes (normal, bacterial pneumonia, viral pneumonia) using the formulas above.

For a multiclass classification with 3 possible output labels, we would need to sum up the values by row and column, so we'd have for the sum of each column:

Total_Bacterial_Pneumonia=9
Total_Normal = 10
Total_Viral_Pneumonia = 11

For the sum of rows for the predicted label:
Total_Predicted_Bacterial_Pneumonia = 10
Total_Predicted_Normal = 10
Total_Predicted_Pneumonia = 10

In an imbalanced classification problem with more than two classes, precision is calculated as the sum of true positives across all classes divided by the sum of true positives and false positives across all classes.

- Precision = Sum c in C TruePositives_c / Sum c in C (TruePositives_c + FalsePositives_c)

In an imbalanced classification problem with more than two classes, recall is calculated as the sum of true positives across all classes divided by the sum of true positives and false negatives across all classes.

- Recall = Sum c in C TruePositives_c / Sum c in C (TruePositives_c + FalseNegatives_c)

| True Label | Predicted Label: bacterial pneumonia | normla | viral pneumonia |
|---|---|---|---|
| bacterial pneumonia | 9 | - | 1 |
| normla | - | 10 | - |
| viral pneumonia | - | - | 10 |

| **F1 Score**<br>What is this model's F1 score? | $F1 = 2 \cdot p \cdot r / p + r$<br>$F1 = 2*(0.97*0.97) / 0.97 + 0.97$<br>$F1 = 2*0.94 / 1.94$<br>$F1 = 1.88/1.94$<br>$F1 = 0.96$ |
|---|---|