

AutoML Modeling Report



Svajune Klimauskaite

Binary Classifier with Clean/Balanced Data

Train/Test Split

How much data was used for training? How much data was used for testing?

The AutoML Vision uses the 80% of your content documents for training, 10% for validating, and 10% for testing.

- ✱ **Train** - Use the image to train the model.
- ✱ **Validation** - Use the image to validate the results that the model returns during training.
- ✱ **Test** - Use the image to verify the model's results after the model has been trained.

AutoML Vision automatically places images it in one of the three sets to ensure that there is enough training, validation, and testing content.

Here was the status for my dataset:

Training images	147
Validation images	25
Test images	28

Confusion Matrix

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?

Confusion matrix represents the percentage of times each label was predicted for each label in the training set during evaluation.

True label	Predicted label	
	Pneumonia	Normal
Pneumonia	76.9%	23.1%
Normal	12.5%	87.5%

- ✱ Within Pneumonia Images, the model labeled as Pneumonia (**true positive**) on the **76.9 %** of the cases, however on the **23.1** cases the model failed to label correctly (**false negative**).
- ✱ Within Healthy Images, the model labeled as Healthy (**false positive**) on the **87.5 %** of the cases, however on the **12.5** cases the model labeled as Pneumonia (**true negatives**).

Precision & Recall

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 are just different metrics for measuring the "success" or performance of a trained model.

- ✱ **precision** is defined as the number of true positives over all positives, and will be the higher when the amount of false positives is low.
- ✱ **recall** is defined as the number of true positives over true positives plus false negatives and will be higher when the number of false negatives is low.

Both take into account true positives and will be higher for high, positive accuracy, too.

Score threshold ?  0.50

Total images 200

Precision ? 81.0%

Recall ? 81.0%

Score Threshold

When you increase the score
threshold, what happens to
precision? What happens to recall?
Why?

The score threshold refers to the level of confidence the model must have to assign a Pneumonia case.

- ✱ If the score threshold is low, the model classifies more images, but runs the risk of misclassifying a few images in the process.
- ✱ If the score threshold is high, the model classifies fewer images, but it will have a lower risk of misclassifying images.

In this case, we want to ensure that Pneumonia case will not get labelled as Healthy, therefore we should be interested in lower confidence level and higher recall.

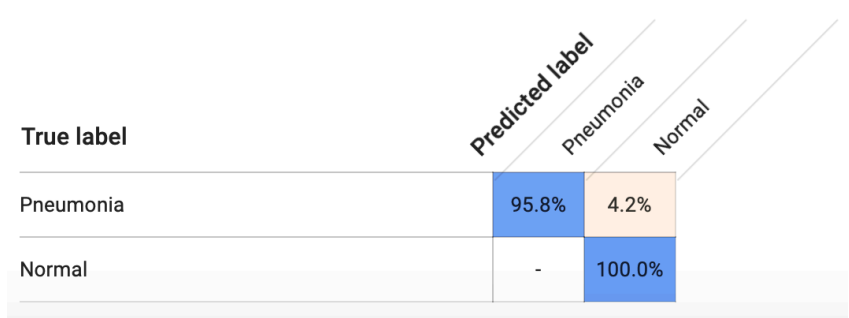
Score threshold ?  0.85

Total images 200

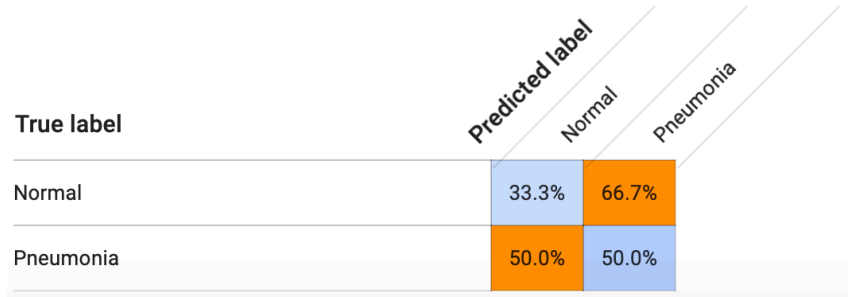







Precision ? 100.0%

Recall ? 61.9%

Binary Classifier with Clean/Unbalanced Data

<h3>Train/Test Split</h3> <p>How much data was used for training? How much data was used for testing?</p>	<table><tr><td>Training images</td><td>318</td></tr><tr><td>Validation images</td><td>42</td></tr><tr><td>Test images</td><td>39</td></tr></table>	Training images	318	Validation images	42	Test images	39						
Training images	318												
Validation images	42												
Test images	39												
<h3>Confusion Matrix</h3> <p>How has the confusion matrix been affected by the unbalanced data? Include a screenshot of the new confusion matrix.</p>	<p>Healthy cases had 100% success, while Pneumonia images had 4.2 % false negative.</p> <p>This way of unbalanced data shows very good results and it seems that it didn't bias Normal images to lean towards Pneumonia label.</p> <p>However, I assume the opposite training scenario, i.e. with 300 Normal images and 100 Pneumonia, would give completely opposite results.</p>  <table><tr><th>True label</th><th colspan="2">Predicted label</th></tr><tr><th></th><th>Pneumonia</th><th>Normal</th></tr><tr><th>Pneumonia</th><td>95.8%</td><td>4.2%</td></tr><tr><th>Normal</th><td>-</td><td>100.0%</td></tr></table>	True label	Predicted label			Pneumonia	Normal	Pneumonia	95.8%	4.2%	Normal	-	100.0%
True label	Predicted label												
	Pneumonia	Normal											
Pneumonia	95.8%	4.2%											
Normal	-	100.0%											
<h3>Precision & Recall</h3> <p>How have the model's precision and recall been affected by the unbalanced data? (Report the values for a score threshold of 0.5.)</p>	<p>The machine learns very well when we increased Pneumonia training data by 300%.</p> <table><tr><td>Score threshold ?</td><td><div><div></div></div>0.50</td></tr><tr><td>Total images</td><td>399</td></tr><tr><td>Precision ?</td><td>96.9%</td></tr><tr><td>Recall ?</td><td>96.9%</td></tr></table>	Score threshold ?	<div><div></div></div> 0.50	Total images	399	Precision ?	96.9%	Recall ?	96.9%				
Score threshold ?	<div><div></div></div> 0.50												
Total images	399												
Precision ?	96.9%												
Recall ?	96.9%												
<h3>Unbalanced Classes</h3> <p>From what you've observed, how do unbalanced classes affect a machine learning model?</p>	<p>The machine learns very well when we increased Pneumonia training data by 300%. Pneumonia images are much more complex than Normal ones, therefore the training, with all kind of Pneumonia variations, gave such a good output.</p>												

Binary Classifier with Dirty/Balanced Data

<h3>Confusion Matrix</h3> <p>How has the confusion matrix been affected by the dirty data? Include a screenshot of the new confusion matrix.</p>	<p>The confusion matrix says it all - machine didn't learn at all.</p>  <table><tr><th>True label</th><th>Predicted label</th><th>Normal</th><th>Pneumonia</th></tr><tr><th>Normal</th><td>33.3%</td><td>66.7%</td></tr><tr><th>Pneumonia</th><td>50.0%</td><td>50.0%</td></tr></table>	True label	Predicted label	Normal	Pneumonia	Normal	33.3%	66.7%	Pneumonia	50.0%	50.0%
True label	Predicted label	Normal	Pneumonia								
Normal	33.3%	66.7%									
Pneumonia	50.0%	50.0%									
<h3>Precision & Recall</h3> <p>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?</p>	<p>Score threshold  0.50</p> <table><tr><td>Total images</td><td>200</td></tr><tr><td>Precision </td><td>42.1%</td></tr><tr><td>Recall </td><td>42.1%</td></tr></table>	Total images	200	Precision 	42.1%	Recall 	42.1%				
Total images	200										
Precision 	42.1%										
Recall 	42.1%										
<h3>Dirty Data</h3> <p>From what you've observed, how do dirty data affect a machine learning model?</p>	<p>It shows the implication of having 30% dirty data. Machine didn't learn anything from this dataset.</p>										

3-Class Model

Confusion Matrix

Summarize the 3-class confusion matrix. What classes are the model most likely to confuse? What class(es) is the model most likely to get right? What might you do to try to remedy the model's "confusion"? Include a screenshot of the new confusion matrix.

True label	Predicted label		
	Bacteria	Normal	Viral
Bacteria	76.9%	15.4%	7.7%
Normal	5.9%	94.1%	-
Viral	14.3%	-	85.7%

- * Class Normal has the highest chances to get labelled correctly. It got confused only with the Pneumonia - Bacteria by 5.9%.
- * Class Viral has 85.7% chances to get it right and got confused with Pneumonia - Bacteria by 14.3%.
- * Class Bacteria is most likely to get confused with Normal or Viral.

Precision & Recall

What are the model's precision and recall? How are these values calculated? (Report the values for a score threshold of 0.5.)

Score threshold ?  0.50

Total images 300

Precision ? 86.5%

Recall ? 86.5%

- * First, the recall and precision are calculated for each class based on the predicted labels (absolute numbers are taken from the confusion matrix).
- * When we have the calculation for each class, then we calculate the average to get the precision/recall of the whole model.

F1 Score

What is this model's F1 score?

F1 measures how well the model performs across all score thresholds. In AutoML Vision, this metric is called Average Precision and is equal to **0.972**.

Analyzed
300 images
3 labels, 37 test images

Avg precision ?
0.972

Precision ?
86.486%

Recall ?
86.486%