# AutoML Modeling Report

*Olivier POLACK*

## Binary Classifier with Clean/Balanced Data

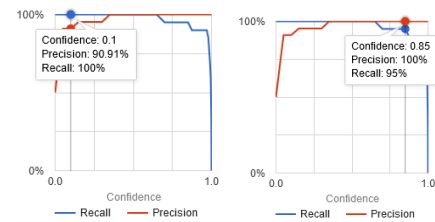| | |
|---|---|
| **Train/Test Split**<br>How much data was used for training? How much data was used for testing? | For this first model, 180 images were used training the model and the remaining 20 others were used for testing the model.<br><br>Total images ............ 180<br>Test items ............ 20<br>Precision ❓ ............ 100%<br>Recall ❓ ............ 100% |
| **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? | Figures in the confusion matrix describe the results of the tests performed on the model after it has been trained. The number of correctly predicted or classified items for each possible class is shown as well as the number of incorrect predictions for each class.<br><br>**Confusion matrix**<br>This table shows how often the model classified each label cor<br><br>In the first example, we have 100% of actual "normal" class images predicted as such, so 0% (or here "–" is displayed) of "normal" class images. For the "pneumonia" class, we also have 100% of actual "pneumonia" class detected as such, and no "pneumonia" class label classified as "normal". Assuming "Pneumonia" is being considered as the positive class, the true positive rate is 100%.<br>The false positive rate for the "normal" class is 0% as no "normal" cases are predicted as "pneumonia" cases. |

| **Precision and 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 is the ratio of true positives to the total of the true positives and false positives. It can be expressed in %. If the number of False Positives is 0, the model is as precise or accurate as it could be, thus Precision ratio equals 1 or 100%. The less accurate the model is, the more false positives are predicted, the more the Precision ratio decreases. On the opposite, instead of focusing on the number of False Positives the model predicted, Recall looks at the number of False Negatives. Recall is the ratio of True Positives to the actual Positives: True Positives plus False Negatives. Recall can be expressed in % too. So a model with a perfect Recall is a model that does not predict any False Negatives at all, regardless of the number of False Positives predictions. A model with a perfect Recall does not let any actually positive cases miss-classified as Negative. |
|---|---|
| **Score Threshold** When you increase the threshold what happens to precision? What happens to recall? Why? |  There will be a trade-off between Precision and Recall in almost all models unless the model is perfect. Higher precision leads to lower recall as the threshold increases and a higher recall will lead to lower precision as threshold decreases. |

# Binary Classifier with Clean/Unbalanced Data

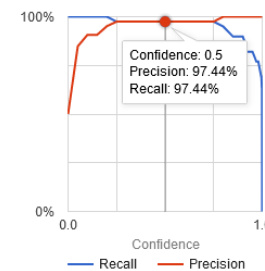| | |
|---|---|
| **Train/Test Split**<br>How much data was used for training? How much data was used for testing? | 400 images were imported, 100 normal images and 300 pneumonia images. 394 images were actually used, among which 355 were used for training the model and 39 were used for testing the model.<br>The reason why some images were not used is because these images are already considered as ground truth by Google. I have tried two times and I always had this warning for few images, please let me know in the review what I should do, I found in the FAQ an answer that seems to imply this is acceptable.<br><br>**All labels**<br><br>Total images     355<br>Test items     39 |
| **Confusion Matrix**<br>How has the confusion matrix been affected by the unbalanced data? Include a screenshot of the new confusion matrix. | The new confusion matrix shows among tested labels all Pneumonia labels are correctly predicted as such (there are no false positives). However, one Normal image is classified as Pneumonia, ie we have 1 false negative. This means that the much larger amount of Pneumonia labels among the dataset caused the model to be biased towards Pneumonia.<br><br>Predicted Label / PNEUMONIA / NORMAL<br>True Label<br>PNEUMONIA   29   -<br>NORMAL   1   9<br><br>Predicted Label / PNEUMONIA / NORMAL<br>True Label<br>PNEUMONIA   100%   -<br>NORMAL   10%   90% |
| **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)? | For a score threshold of 0.5, the model's precision and recall have dropped from 100% down to 97.44%.<br><br>Confidence: 0.5<br>Precision: 97.44%<br>Recall: 97.44%<br><br>— Recall   — Precision |
| **Unbalanced Classes**<br>From what you have observed, how do unbalanced classed affect a machine learning model? | Unbalanced classes decrease precision and recall rates as they tend to create a biased model predicting false positives in case more labels of the positive class have been thrown into the training dataset. |

# Binary Classifier with Dirty/Balanced Data

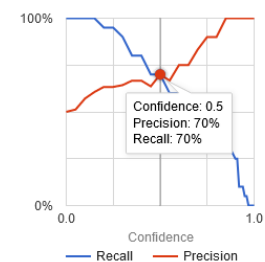| | |
|---|---|
| **Confusion Matrix**<br>How has the confusion matrix been affected by the dirty data? Include a screenshot of the new confusion matrix. | With the dirty data, the confusion matrix shows results are less precise and with poorer recall.<br><br> |
| **Precision and Recall**<br>How have 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 precision and recall are both down at 70% with the dirty data.<br><br><br><br>Results benchmark of the binary classifier models<br><br>Table:<br><br>{{TABLE}} |
| **Dirty Data**<br>From what you have observed, how does dirty data affect a machine learning model? | Dirty data obviously deteriorates the quality of the trained model, precision, recall and accuracy are affected. |

Results benchmark of the binary classifier models

| Model | Precision | Recall |
|---|---|---|
| 1: Clean/ Balanced Data | 100 % | 100 % |
| 2. Clean/Unbalanced Data | 97,4 % | 97,4 % |
| 3. Dirty/Balanced Data | 70 % | 70 % |

The model with balanced data has the highest precision and recall of 100% each at a threshold of 0.5 among the binary classifiers.

# 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. | The model is most likely the two Pneumonia classes labels and is most likely to get right the images of the Normal class. This is harder for the model to properly distinguish between different Pneumonia types as between Normal and Pneumonia, probably because differences between the two Pneumonia types are much subtle as the difference between Normal and Pneumonia.<br>New confusion matrix in % and items count<br><br><br><br><br><br>In order to remedy the model's confusion, I would try to increase the number of labels of each of the three classes. |
| **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)? | **Calculation of Precision and Recall for each classifier:**<br>**'Normal' Class**<br>True Positives = 10  / False Positives = 0<br>True Negatives = 8+9  = 17 / False Negatives = 2+1 = 3<br>True Positives + False Positives (TP+FP) = 10<br>True Positives + False Negatives (TP+FN) = 13<br>  Precision = TP/(TP+FP) = 10/10 = 100%<br>  Recall = TP/(TP+FN) = 10/13 = 76,9%<br><br>**'Viral Pneumonia' Class**<br>True Positives = 8  / False Positives = 1<br>True Negatives = 10+9  = 19 / False Negatives =  2<br>True Positives + False Positives (TP+FP) = 9<br>True Positives + False Negatives (TP+FN) = 10<br>  Precision = TP/(TP+FP) = 8/9 = 88,8%<br>  Recall = TP/(TP+FN) = 8/10 = 80% |

**'Bacteria Pneumonia' Class**
True Positives = 9  / False Positives = 2
True Negatives = 18 / False Negatives = 1
True Positives + False Positives (TP+FP) = 11
True Positives + False Negatives (TP+FN) = 10
  Precision = TP/(TP+FP) = 9/11 = 81,8%
  Recall = TP/(TP+FN) = 9/10 = 90%

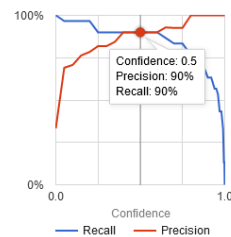**Precision and Recall calculated using micro-averaged method:**
TP = 10 + 8 + 9 = 27
FP = 0 + 1 + 2 = 3
  Precision = 27 / (27+3) = **90% Precision**

FN = 3 + 2 + 1
  Recall = 27 / (27 + 6) = **81,8% Recall**

Model's precision and recall for a score threshold of 0.5



| F1 Score | Formula of the F1 Score |
|---|---|
| What is this model's F1 score? | $F1\ Score = 2*(Precision*Recall)/(Precision+Recall)$ |

Formula of the F1 Score
 F1 Score = 2 *( Precision * Recall)/ (Precision + Recall)

The F1 score of this model can be calculated for each class:
'Normal' Class: 2*( 100% * 76,9%) / ( 100% + 76,9%)
  F1 Score = 86,9%

'Viral Pneumonia' Class:
     2*( 88,8% * 80%) / ( 88,8% + 80%)
  F1 Score = 84,1%

'Bacterial Pneumonia' Class:
     2*( 81,8% * 90%) / ( 91,8% + 90%)
  F1 Score  = 85,7%

**Then the F1 score can be calculated as the average score of the three classes:**

F1 Score = (86,9% + 84,1% + 85,2%) / 3
                    **F1 Score = 85,4%**