# AutoML Modeling Report

*Jimmey Jose*

## Binary Classifier with Clean/Balanced Data

| | |
|---|---|
| **Train/Test Split**<br><br>How much data was used for training? How much data was used for testing? | A total of 200 x-rays were used, of which 180 was used for training and 20 for testing. The exact breakup based on type of x-rays is given below:<br><br>**Break Up of Data Used**<br><br>|  | Training | Testing | **Total** |<br>|---|---|---|---|<br>| Normal chest x-rays | 90 | 10 | **100** |<br>| Pneumonia chest x-rays | 90 | 10 | **100** |<br>| **Total** | **180** | **20** | **200** | |
| **Confusion Matrix**<br><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? | Each cell in a confusion matrix captures the number of predictions that were either correct or wrong when testing an ML model.<br><br>The four different types of values in a confusion matrix from the perspective of this binary class model are described below. ('pneumonia' is being considered as the positive class for the purpose of the below descriptions.)<br>1. True Positives - Test cases in which the true label is 'pneumonia' and is correctly predicted to be the same.<br>2. False Positive - Test cases in which the true label is 'normal' and is predicted wrongly to be 'pneumonia'<br>3. True Negatives -Test cases in which the true label is 'normal' and is predicted to be the same<br>4. False Negatives - Test cases in which the true label is 'normal' and is predicted to 'pneumonia'<br><br>**Confusion matrix for this model**<br><br><br>• True positive rate for "pneumonia" class is 100 %<br>• False positive rate for "normal" class is 0 % |

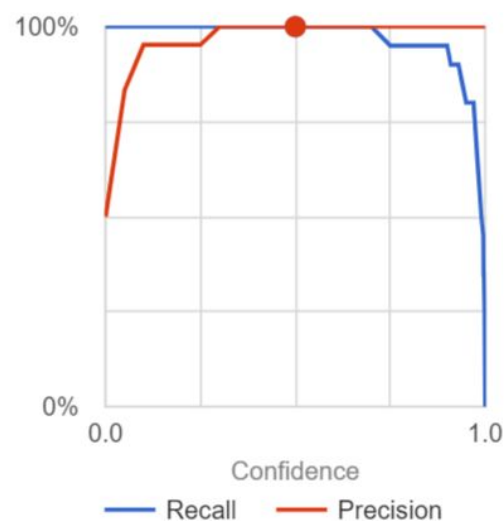| **Precision & Recall**<br><br>What does precision measure? What does recall measure? What precision and recall did the model achieve (report the values for a score threshold of 0.5)? | • Precision measures the ratio between the number of true positives and the number of predicted positives(True positives + False positives). This value gives a confidence level about how accurate true positive outcomes are.<br>• Recall measures the ratio between the number of true positives and the number of actual positives(True positives + False Negatives). This value gives a confidence level about how many of the true positives were identified.<br><br>The precision and recall were both 100 % at a threshold of 0.5 is this model.<br><br><table><tr><td>Total images</td><td>180</td></tr><tr><td>Test items</td><td>20</td></tr><tr><td>Precision ?</td><td>100%</td></tr><tr><td>Recall ?</td><td>100%</td></tr></table><br>Note: This Precision and recall are calculated using Micro-averaged method. |
| **Score Threshold**<br><br>When you increase the score 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 (which generally doesn't happen). In this context higher precision leads to lower recall as the threshold increases and a higher recall will lead to lower precision as threshold decreases.<br><br>With a low threshold, a model will classify everything even though the risk of getting it wrong is high; while with a high threshold a model might not classify anything because of the risk of getting it wrong.<br><br>In this model the precision and recall are 100 % when the threshold is in between .3 and .7; below .3, precision will drop and above .7, recall will drop.<br><br> |

# Binary Classifier with Clean/Unbalanced Data

| | |
|---|---|
| **Train/Test Split**<br><br>How much data was used for training? How much data was used for testing? | A total of 400 x-rays were used, of which 360 was used for training and 40 for testing. The exact breakup based on type of x-rays is given below:<br><br>**Break Up of Data Used**<br><br><table><tr><td></td><td>Training</td><td>Testing</td><td>**Total**</td></tr><tr><td>Normal chest x-rays</td><td>90</td><td>10</td><td>**100**</td></tr><tr><td>Pneumonia chest x-rays</td><td>270</td><td>30</td><td>**300**</td></tr><tr><td>**Total**</td><td>**360**</td><td>**40**</td><td>**400**</td></tr></table> |
| **Confusion Matrix**<br><br>How has the confusion matrix been affected by the unbalanced data? Include a screenshot of the new confusion matrix. | The unbalanced data has resulted in higher accuracy for 'pneumonia' class which had the bigger data set. The confusion matrix shows the same through the following.<br><br>1. The true positive rate is higher by 7 % for 'pneumonia'. If you consider the size of the test data set, you will observe that for both classes,the model inaccurately predicted 1 record each which doesn't seem high; but if the test data were more equally distributed we may have seen a bigger true positive rate for 'pneumonia'<br>2. 'Pneumonia' class has a higher false positive of 10 % compared to 3 % for 'normal class which shows that there is a higher chance of being classified in the 'pneumonia' class<br><br> |
| **Precision & Recall**<br><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.) | The precision and recall at threshold of .5 are both 95%. The unbalanced data seems to give the same 95% precision and recall within a wide threshold between 0 .03 and 0.99 which appears to be very skewed and unreliable.<br><br><table><tr><td>Total images</td><td>360</td></tr><tr><td>Test items</td><td>40</td></tr><tr><td>Precision ❓</td><td>95%</td></tr><tr><td>Recall ❓</td><td>95%</td></tr></table><br>Note: The above Precision and recall is calculated using micro-averaged method. |

| Unbalanced Classes

From what you've observed, how do unbalanced classes affect a machine learning model? | 1. Unbalanced classes seem to skew predictions to the class with a larger data set which is 'pneumonia in this case
2. In this model, the precision and recall for 'pneumonia' is 96.67 each and for 'normal' is 90% each; this is a difference of only around 6% and for that reason, the combined precision and recall of 95% is accurate. There is generally a high chance to have a sizeable difference between the individual precision and recall because of unbalanced data which can make the model unreliable. |
|---|---|

# Binary Classifier with Dirty/Balanced Data

| Confusion Matrix

How has the confusion matrix been affected by the dirty data? Include a screenshot of the new confusion matrix. | The dirty data has brought down the True Positives and increased the False Positives substantially compared to the first model which had the same amount of data. This huge drop in reliability of data is because the dirty data resulted in the model training on the wrong parameters and noise.

 |
|---|---|
| Precision & Recall

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 reca*ll?* | The precision and recall has dropped from 100% each in the first model with the same amount of data to 65% each at a threshold of .5 for this model because of the dirty data.



- The first model with Clean/Balanced data has the highest precision and recall of 100% each at a threshold of .5 among the binary classifiers

**Summary of Precision and Recall in Binary Classifier Models**

| Binary Classifier | Precision | Recall |
|---|---|---|
| 1. Clean/ Balanced Data | 100 % | 100% |
| 2. Clean/Unbalanced Data | 95% | 95% |
| 3. Dirty/Balanced Data | 65% | 65% |

Note: The above Precision and recall is calculated using micro-averaged method. |

| | |
|---|---|
| **Dirty Data**<br><br>From what you've observed, how do dirty data affect a machine learning model? | 1. Dirty data confuses the machine learning model resulting in a substantial decrease in accuracy, precision and recall<br>2. The ML model may have fitted in more of the noise in the data making predictions unreliable<br>3. Equal amount of dirty data in each class does not result in correct output in the models. This model had an equal amount of bad data in both classes, but there is no relation between the precision and recall for each class |

# 3-Class Model

| | |
|---|---|
| **Confusion Matrix**<br><br>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. | This model is most likely to confuse the 'bacterial pneumonia' class and get 'normal' and 'viral pneumonia' classes correct as seen in the confusion matrix.<br><br><br><br>The following can be done to improve the model's confusion:<br>  1. Use more data with all 3 sets<br>  2. Improve quality of data through preprocessing, augmentation or any other technique<br>  3. Look for opportunities to adjust the threshold levels to improve results<br>  4. Look for opportunities to change or improve the modeling algorithm to improve results |
| **Precision & Recall**<br><br>What are the model's precision and recall? How are these values calculated? (Report the values for a score threshold of 0.5.) | This models precision and recall are 93.33% each at a threshold of 0.5<br><br><br><br>Note: The above Precision and recall is calculated using micro-averaged method.<br><br>There are 2 main ways to calculate Precision and Recall.<br>  1. Precision and Recall for each class<br>  2. Combined Precision and Recall using micro-averaged method |

| | |
|---|---|
| | **Precision and Recall for each classifier**<br><br>1. 'Normal' Class Precision and Recall calculation<br><br>True Positives (TP)                                       = 10<br>True Positives + False Positives (TP+FP) = 10+1+0 = 11<br>True Positives + True Negatives(TP+TN) = 10+0+0 =10<br><br>Precision =  TP / (TP+FP) = 10 / 11 =  90.91 %<br>Recall      = TP/ (TP+FN) = 10/10  = 100    %<br><br>2. 'Viral pneumonia' Class Precision and Recall calculation<br><br>True Positives (TP)                                       = 10<br>True Positives + False Positives (TP+FP) = 10+1+0   = 11<br>True Positives + True Negatives(TP+TN) = 10+0+0 = 10<br><br>Precision =  TP / (TP+FP) = 10 / 11 =  90.91 %<br>Recall      = TP/  (TP+FN) = 10/10   = 100    %<br><br>3. 'Bacterial pneumonia' Class Precision and Recall calculation<br><br>True Positives (TP)                                       =  8<br>True Positives + False Positives (TP+FP) = 8+0+0   =  8<br>True Positives + True Negatives(TP+TN) = 8+1+1  = 10<br><br>Precision =  TP / (TP+FP) =   8 / 8 =  100   %<br>Recall      = TP/  (TP+FN) =  8/10   =  80 %<br><br>**Precision and Recall using micro-averaged method**<br><br>True Positives (TP)                                       = 10 + 10 + 8 =28<br>True Positives + False Positives (TP+FP) =11+11 + 8    = 30<br>True Positives + True Negatives(TP+TN) =10+10+10  = 30<br><br>Precision =  TP / (TP+FP) = 28 / 30    = 93.33%<br>Recall      = TP/  (TP+FN)= 28 / 30    = 93.33% |
| **F1 Score**<br><br>What is this model's F1 score? | F1 scores can be calculated for each classifier and combined using multiple ways. Below are F1 scores by classifier and combined F1 scores in 2 popular different methods.<br><br>Formula:  F1 Score = 2 *( Precision * Recall)/ (Precision + Recall)<br><br>**F1 Score by Class**<br><br>'Normal' class F1 Score = 2 * (90.91 %* 100%)/  (90.91%+100%) = **95.24%**<br>'Viral pneumonia' class F1 Score = 2 * (90.91 %* 100%)/ (90.91%+100%) = **95.24%**<br>'Bacterial pneumonia' class F1 Score =  2 * (100 %* 80%)/ (100%+80) = **88.89%**<br><br>**Macr0-averaged F1 Score=** (95.24%+95.24%+88.89%)/ 3= **93.12%** |

| | **Micro-averaged F1 Scoro**=Micro-avg. precision/recall l=   **93.33%** |
|---|---|