# AutoML Modeling Report

*Hussain Al-Balhareth*



Figure: AutoML five trained models

## Binary Classifier with Clean/Balanced Data

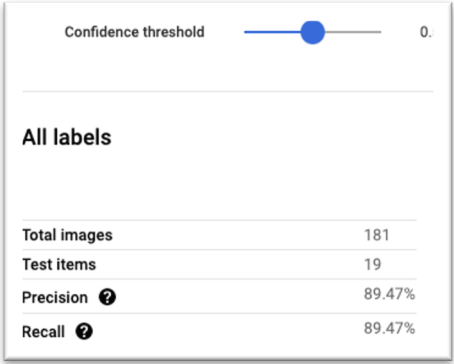| | |
|---|---|
| **Train/Test Split**<br>How much data was used for training? How much data was used for testing? |  |
| **Confusion Matrix**<br>What do each of the cells in the confusion matrix describe? What values did you observe (include a screenshot)? What is the true positive rate for the "pneumonia" class? What is the false positive rate for the "normal" class? | <br>The blue (diagonal) values refer to true predictions (positive or negative) while off-diagonal (grey) refers to false predictions (positive or negative). For example, in |

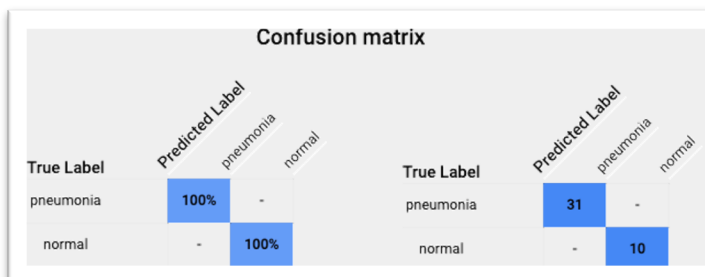| | the first row, a pneumonia (positive) case was predicted correctly (TP) 8 times (89% of the time) and was only once falsely predicted as normal (negative) or 11% of the time. With regard to the normal (negative) case, it was predicted correctly (TN) 8 times (90% of the time) and was only once falsely predicted (FP) as pneumonia 10% of the time. |
|---|---|
| **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 measures true predictions over total predications. While recall measures true predictions over total ground truth.<br><br> |
| **Score Threshold**<br>When you increase the threshold what happens to precision? What happens to recall? Why? | <br><br>Precision score is proportional to confidence and vice versa with recall case. For example, when the model predicts an image as normal with 60% confidence and 40% confidence, then in this case the outcome prediction would be normal and so on. |

# Binary Classifier with Clean/Unbalanced Data

| | |
|---|---|
| **Train/Test Split**<br>How much data was used for training? How much data was used for testing? | **Label Stats**<br><br>Labels — Images — Train — Validation — Test<br>normal — 100 — 80 — 10 — 10<br>pneumonia — 300 — 239 — 30 — 31 |
| **Confusion Matrix**<br>How has the confusion matrix been affected by the unbalanced data? Include a screenshot of the new confusion matrix. | **Confusion matrix**<br><br>True Label — Predicted Label (pneumonia, normal)<br>pneumonia — 100% — -<br>normal — - — 100%<br><br>True Label — Predicted Label (pneumonia, normal)<br>pneumonia — 31 — -<br>normal — - — 10<br><br>Yes, it improved to 100%. |
| **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)? | Confidence threshold — 0.5<br><br>**All labels**<br><br>Total images — 359<br>Test items — 41<br>Precision ❓ — 100%<br>Recall ❓ — 100%<br><br>It (unexpectedly) increased to 100%. |
| **Unbalanced Classes**<br>From what you have observed, how do unbalanced classed affect a machine learning model? | It weirdly gets higher scores probably due to higher exposure to training. Normally, such unbalance could create a bias or tendency in the model towards one prediction over another. |

# Binary Classifier with Dirty/Balanced Data

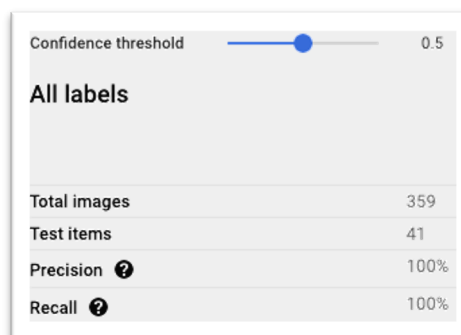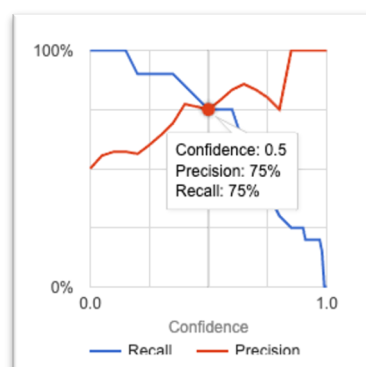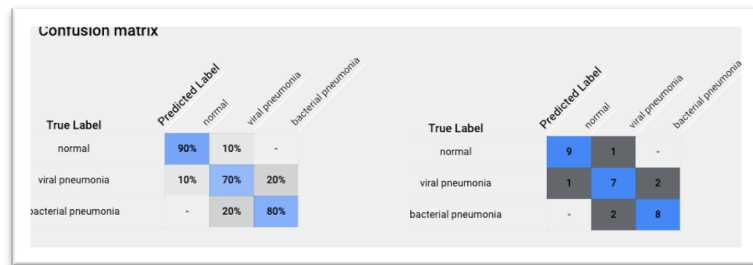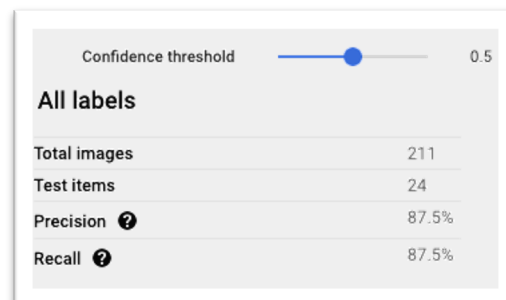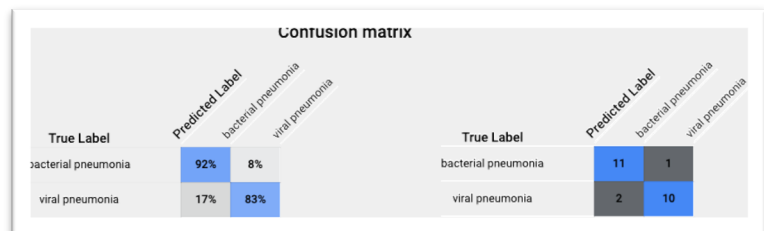| | |
|---|---|
| **Confusion Matrix**<br>How has the confusion matrix been affected by the dirty data? Include a screenshot of the new confusion matrix. | **Confusion matrix**<br><br>True Label / Predicted Label (pneumonia, normal)<br><br>pneumonia: 80%, 20%<br>normal: 30%, 70%<br><br>pneumonia: 8, 2<br>normal: 3, 7<br><br>It gets worse. |
| **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? | Confidence: 0.5<br>Precision: 75%<br>Recall: 75%<br><br>(Recall — Precision, Confidence 0.0 to 1.0)<br><br>It became very sensitive to confidence level with overall decrease in the performance by 15% to 75%. |
| **Dirty Data**<br>From what you have observed, how does dirty data affect a machine learning model? | It impacts the results negatively. |

# 3-Class Model

## Confusion Matrix

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.
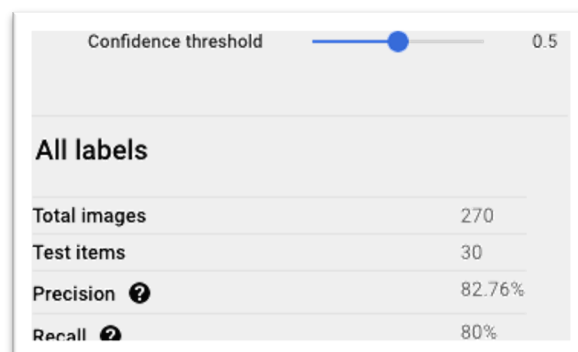


Viral pneumonia is most likely to confuse with least true prediction score (70%). While normal is most likely to get right with the highest prediction score (90%). To remedy that, I tried another model with a two-class pneumonia dataset (viral or bacterial) and got the following result:





| Confidence threshold | 0.5 |
| --- | --- |
| **All labels** | |
| Total images | 211 |
| Test items | 24 |
| Precision | 87.5% |
| Recall | 87.5% |

So, the overall precision and recall scores were both enhanced, e.g. increased by 7.5%.

## 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)?



| Confidence threshold | 0.5 |
| --- | --- |
| **All labels** | |
| Total images | 270 |
| Test items | 30 |
| Precision | 82.76% |
| Recall | 80% |

$$precision = \frac{\frac{9}{10} + \frac{7}{10} + \frac{8}{10}}{3} = 0.8$$

| | |
|---|---|
| | $$recall = \frac{\frac{9}{10} + \frac{7}{10} + \frac{8}{10}}{3} = 0.8$$ |
| **F1 Score**<br>What is this model's F1 score? | $$F1 = 2 \times \frac{82.76\% * 80\%}{82.76\% + 80\%} = 81.35\%$$ |