# AutoML Modeling Report

*<Ugur Uresin, December 16, 2023>*

## Binary Classifier with Clean/Balanced Data

| | |
|---|---|
| **Train/Test Split**<br>How much data was used for training? How much data was used for testing? | *<br>**How to split a dataset**<br><br>**TRAINING SET**<br>The subset of data used to train a machine learning model<br><br>**TEST SET**<br>The subset of data used to evaluate the performance of a trained machine learning model on unseen examples, simulating real-world data<br><br>**VALIDATION SET**<br>The intermediary subset of data used during the model development process to fine-tune hyperparameters<br><br>300 random images in total were selected from the Kaggle's dataset (150 normal & 150 pneumonia).<br>From this dataset, 120 images (60 normal & 60 pneumonia) were used for training, 30 images were used for testing (15 normal & 15 pneumonia) and 30 images were used for validation (15 normal & 15 pneumonia) as shown in the image below:<br><br>**Created** — Dec 15, 2023, 11:37:58 PM<br>**Total images** — 300<br>**Training images** — 240<br>**Validation images** — 30<br>**Test images** — 30<br><br>* Reference: https://medium.com/syntaxerrorpub/understanding-the-difference-between-training-test-and-validation-sets-in-machine-learning-c59feec6483b |
| **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? | The confusion matrix in the image is a tool used in machine learning to evaluate the performance of classification models. It shows the number of correct and incorrect predictions made by the model, compared to the actual labels. |

**Confusion matrix**

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).

| True label | Predicted label pneumonia | normal |
|---|---|---|
| pneumonia | 15 | 0 |
| normal | 2 | 13 |

Here's what each cell represents:

True Positive (TP): The model correctly predicted the positive class (in this case, "pneumonia"). Here, the TP for pneumonia is 15.

True Negative (TN): The model correctly predicted the negative class (in this case, "normal"). The TN for normal is 13.

False Positive (FP): The model incorrectly predicted the positive class. Here, there are 0 instances where the model predicted pneumonia when it was actually normal.

False Negative (FN): The model incorrectly predicted the negative class. Here, there are 2 instances where the model predicted normal when it was actually pneumonia.

Given these values, we can calculate the following:
True Positive Rate (TPR) for the "pneumonia" class, also known as sensitivity or recall, is TP / (TP + FN).

So it would be 15 / (15 + 2) = 15 / 17.

False Positive Rate (FPR) for the "normal" class is FP / (FP + TN), which would be 0 / (0 + 13) = 0,
since there are no false positives for the normal class.

**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** measures the accuracy of positive predictions. It is defined as TP / (TP + FP), the proportion of positive identifications that were actually correct.
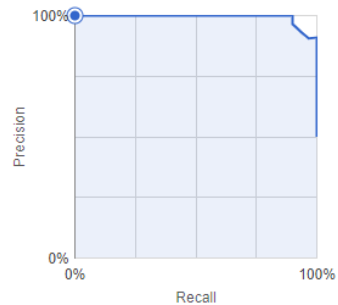
**Recall** measures the ability of a model to find all the relevant cases within a dataset. It is defined as TP / (TP + FN), the proportion of actual positives that were correctly identified.

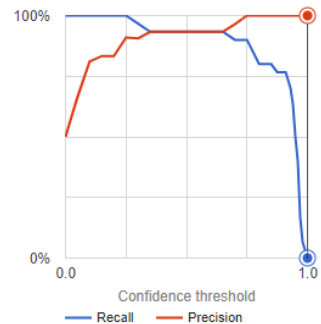| | |
|---|---|
| Average precision ❓ | 0.994 |
| Precision ❓ | 93.3% |
| Recall ❓ | 93.3% |
| Created | Dec 15, 2023, 11:37:58 PM |
| Total images | 300 |
| Training images | 240 |
| Validation images | 30 |
| Test images | 30 |

For this model, the precision and recall at the score threshold of 0.5 for the "pneumonia" class are:

**Precision**: It cannot be calculated directly from the confusion matrix because FP for the "normal" class is zero, and we do not have the number of false positives for the "pneumonia" class. If we assume there are no false positives for "pneumonia" as well, the precision would be 100%



**Recall**: As given below, the recall for "pneumonia" is 15 / 17 (93.3%)



## Score Threshold
When you increase the threshold what happens to precision? What happens to recall? Why?



When the threshold is set to 0.75, the precision of the model is at 100%.
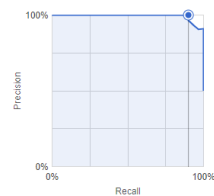The recall at this threshold is lower than the precision, though the exact value isn't visible in the graph. It appears to be less than 100% but higher than the recall at a threshold of 1.0.

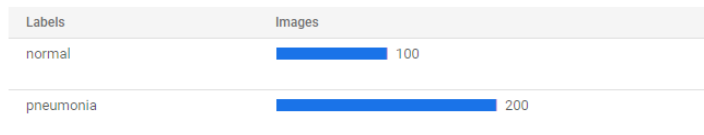| | The precision and recall are measures of a model's accuracy and are affected differently by the confidence threshold: |
|---|---|
| | Precision is the ratio of true positives to the total number of positive predictions made by the model. A higher confidence threshold generally leads to a higher precision because the model is more conservative in making positive predictions; it only makes a prediction when it is very sure. This means it makes fewer positive predictions, but those predictions are more likely to be correct. |
| | Recall is the ratio of true positives to the actual number of positive instances in the dataset. As the confidence threshold increases, the model requires higher confidence to make a positive prediction, which means it may miss some actual positives, leading to a lower recall. It's a measure of how many of the actual positives the model is able to capture. |
| | In summary, as you increase the confidence threshold: |
| | Precision may remain high or even increase because the model makes fewer false positive predictions. Recall tends to decrease because the model may miss more true positive cases, as it is more stringent in its prediction criteria. This is a trade-off often encountered in classification tasks: increasing precision typically comes at the cost of decreasing recall, and vice versa. The ideal balance depends on the specific application and the cost associated with false positives versus false negatives. |

# Binary Classifier with Clean/Unbalanced Data

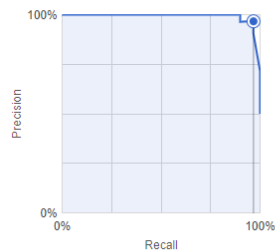| | |
|---|---|
| **Train/Test Split**<br>How much data was used for training? How much data was used for testing? | 300 images were used. 100 for normal, 200 for pneumonia as follows:<br><br>| Labels | Images |<br>|---|---|<br>| normal | 100 |<br>| pneumonia | 200 |<br><br>The split ratio is 80/10/10 for training/validation/test as follows:<br><br>| Dataset | x-ray-dataset |<br>|---|---|<br>| Annotation set | x-ray-dataset_icn |<br>| Data split | Randomly assigned (80/10/10) |<br><br>240 images for training, 30 images for validation, 30 images for testing<br><br>## All labels<br><br>| | |<br>|---|---|<br>| Average precision ❓ | 0.995 |<br>| Precision ❓ | 96.7% |<br>| Recall ❓ | 96.7% |<br>| Created | Dec 16, 2023, 2:27:47 AM |<br>| Total images | 300 |<br>| Training images | 240 |<br>| Validation images | 30 |<br>| Test images | 30 | |
| **Confusion Matrix**<br>How has the confusion matrix been affected by the unbalanced data? Include a screenshot of the new confusion matrix. | **Predicted label** / **True label**<br><br>| True label \ Predicted | pneumonia | normal |<br>|---|---|---|<br>| pneumonia | 19 | 1 |<br>| normal | 0 | 10 |<br><br>The imbalance can lead to a bias in the model where it may become better at identifying "pneumonia" simply because it has more examples of that class to learn from.<br><br>From the confusion matrix:<br>The model predicted "pneumonia" correctly 19 times and incorrectly only once when the true label was "normal".<br>It predicted "normal" correctly 10 times, and there are no instances where "normal" was predicted when the true label was "pneumonia".<br>Given the unbalanced dataset, the model might be overfitting to the "pneumonia" class, which could explain the high performance in predicting "pneumonia" correctly. The lack of false positives for "normal" is a good sign, indicating that when the model predicts "normal", it is highly likely to be correct. |

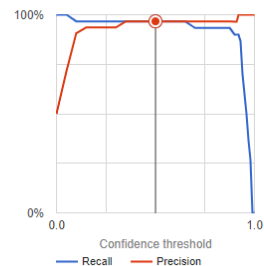| **Precision and Recall**<br>How have the model's precision and recall been affected by the unbalanced data (report the values for a score threshold of 0.5)? | Average precision ❓ | 0.995 |
| | Precision ❓ | 96.7% |
| | Recall ❓ | 96.7% |
| | Created | Dec 16, 2023, 2:27:47 AM |
| | Total images | 300 |
| | Training images | 240 |
| | Validation images | 30 |
| | Test images | 30 |

To evaluate your model, set the **confidence threshold** to see how precision and recall are affected. The best confidence threshold depends on your use case. Read some example scenarios 🔗 to learn how evaluation metrics can be used.

**Precision-recall curve** ❓          **Precision-recall by threshold** ❓



The precision and recall of the model are both 96.7%.
The average precision across all thresholds is 0.995.

| **Unbalanced Classes**<br>From what you have observed, how do unbalanced classed affect a machine learning model? | These high values for precision and recall suggest that the model is performing well on the test set. However, considering that the data is unbalanced, with twice as many "pneumonia" images as "normal", this could potentially lead to a bias in the model.<br>Moreover, if the model is overly confident in predicting the majority class, it could ignore the minority class, which would lead to a high number of false negatives for the minority class, thus affecting the recall. |

# Binary Classifier with Dirty/Balanced Data

<table>
<tr>
<td>

**Confusion Matrix**
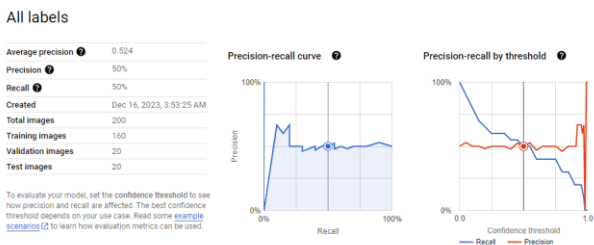How has the confusion matrix been affected by the dirty data? Include a screenshot of the new confusion matrix.

</td>
<td>

| Annotation set | x-ray-dataset-dirty_icn |
| --- | --- |
| Objective | Image classification (Single-label) |
| Items | 200 |
| Created | Dec 16, 2023, 1:09:27 AM |
| Last updated | Dec 16, 2023, 1:36:59 AM |

A dirty dataset was created by intentionally mislabeling 30% of the images for both labels (normal, pneumonia). A total of 200 images were used.



</td>
</tr>
<tr>
<td>

**Precision and 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 recall?

</td>
<td>



Both precision and recall are at 50%.
The average precision across all thresholds is 0.524.

For a score threshold of 0.5, if we assume the values shown are at this default threshold, the precision and recall are equal at 50%.

This typically indicates that the model is performing no better than random guessing for this threshold, which is a common occurrence when dealing with dirty data.

As for which binary classifier has the highest precision and recall:

The confusion matrix shows that both classes ("pneumonia" and "normal") have the same precision and recall.

Hence, both classifiers have equal precision and recall at 80% for correctly identifying true labels and 20% for mislabeling the other class.

However, if we consider the precision and recall at the threshold of 0.5 (assuming the provided values are for this threshold), both precision and recall are at 50%.

</td>
</tr>
</table>

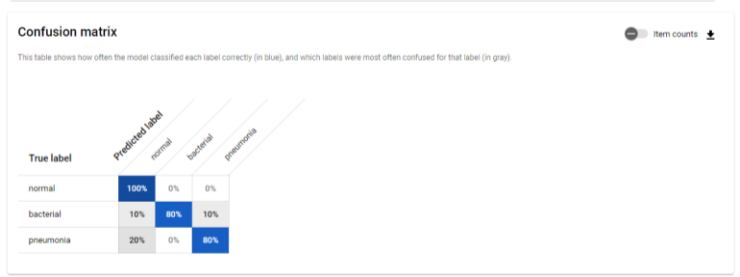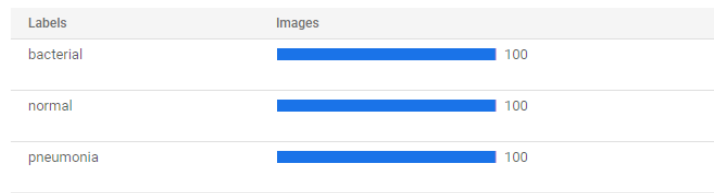| Dirty Data | Dirty data typically refers to data that is incorrect, misleading, or poorly formatted. This can include mislabeled examples, as in your case, where 30% of the images for both classes were intentionally mislabeled. Dirty data can affect machine learning models in several ways: |
|---|---|
| **Dirty Data** <br> From what you have observed, how does dirty data affect a machine learning model? | **Decreased Accuracy:** <br> Mislabeled data can lead to a model learning incorrect patterns, decreasing its overall accuracy. <br><br> **Overfitting:** <br> The model may overfit to noise and errors in the training data, leading to poor generalization to new, unseen data. <br><br> **Skewed Metrics:** <br> Performance metrics can become skewed, as the model may appear to perform well on the dirty data but fail on clean data. <br><br> **Increased Uncertainty:** <br> Dirty data increases the uncertainty of predictions and may require additional post-processing or manual review to ensure the quality of the model's outputs. <br><br> **Misguided Feature Importance:** <br> The model may assign importance to features that are related to the noise in the dirty data, rather than to the true underlying patterns that are predictive of the correct outcomes. <br><br> An Experience from an Industrial AI Use Case <br> I have had an experience with dirty data where we used incorrect labels while training the model for an error mode that was unlikely to occur. <br><br> This particular error mode did not occur in the business unit for a long time, leading us to believe that our model was performing well, as we did not observe any data drift. <br><br> However, when this low-probability but high-impact error mode began to occur in the business unit, we realized that the model failed to catch it for a considerable period. An error analysis revealed that the initial data was not correctly labeled for this specific error mode. <br><br> Since it was an industrial artificial intelligence application, it resulted in a loss of time for the business unit. <br><br> Therefore, it's very important to align with the business unit from the beginning to ensure that the labels are accurate. |

# 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. | 300 images were used and there are equally distributed 3 labels as follows: |
|---|---|

300 images were used and there are equally distributed 3 labels as follows:

| Labels | Images | |
|---|---|---|
| bacterial | | 100 |
| normal | | 100 |
| pneumonia | | 100 |

**Confusion matrix**    Item counts ↓

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).

| True label | Predicted label | | |
|---|---|---|---|
| | normal | bacterial | pneumonia |
| normal | 100% | 0% | 0% |
| bacterial | 10% | 80% | 10% |
| pneumonia | 20% | 0% | 80% |

Normal: The model performs perfectly on this class, with 100% accuracy. No normal cases are confused with either bacterial or pneumonia.

Bacterial: The model correctly identifies 80% of bacterial cases, but confuses 10% of bacterial cases as normal and another 10% as pneumonia.

Pneumonia: The model correctly identifies 80% of pneumonia cases, but confuses 20% of pneumonia cases as normal.

Potential Remedies for Model Confusion:

Data Augmentation: You could increase the number of training examples for bacterial and pneumonia classes, especially focusing on cases that are typical for each class.

Feature Engineering: Identify and include more discriminative features that can help distinguish between bacterial and pneumonia cases.

Error Analysis: Perform an in-depth analysis of the misclassified cases to understand why the model is confusing them, and then take corrective measures based on the findings. Model Complexity: Increase the complexity of the model if it is too simple and unable to capture the nuances between the classes, or consider using a different model architecture.

Class Weighting: Adjust the weights of classes during training to penalize the model more for confusing bacterial and pneumonia cases.

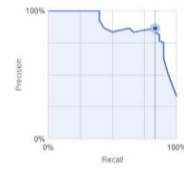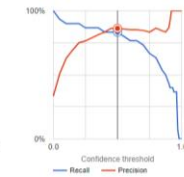| | |
|---|---|
| **Precision and Recall**<br>What are the model's precision and recall? How are these values calculated (report the values for a score threshold of 0.5)? | All labels<br><br>Average precision ❓ 0.877<br>Precision ❓ 86.2%<br>Recall ❓ 83.3%<br>Created Dec 16, 2023, 3:53:34 AM<br>Total images 300<br>Training images 240<br>Validation images 30<br>Test images 30<br><br>Precision-recall curve ❓    Precision-recall by threshold ❓<br><br>To evaluate your model, set the confidence threshold to see how precision and recall are affected. The best confidence threshold depends on your use case. Read some example scenarios ↗ to learn how evaluation metrics can be used.<br><br>**For Normal:**<br>True Positives (TP_normal): 100%<br>False Positives (FP_normal): 0% for bacterial + 0% for pneumonia<br>False Negatives (FN_normal): 0%<br><br>**For Bacterial:**<br>True Positives (TP_bacterial): 80%<br>False Positives (FP_bacterial): 10% for normal + 0% for pneumonia<br>False Negatives (FN_bacterial): 10% for normal + 10% for pneumonia<br><br>**For Pneumonia:**<br>True Positives (TP_pneumonia): 80%<br>False Positives (FP_pneumonia): 0% for normal + 20% for bacterial<br>False Negatives (FN_pneumonia): 20% for normal + 0% for bacterial |
| **F1 Score**<br>What is this model's F1 score? | F1=2×(PrecisionxRecall) / (Precision+Recall)<br>The F1 score for the model, given the precision of 86.2% and recall of 83.3%, is approximately 0.847 or 84.73% |