

# Project Proposal



Clifton Zhuo

---

## Data Labeling Approach

<b>Project Overview and Goal</b>  What is the industry problem you are trying to solve? Why use ML in solving this task?	The problem is to detect the presence of pneumonia on X-ray scans. By using ML, healthcare workers can fasten their Pneumonia detection process, but also use it as an extra verification of the presence of Pneumonia.
<b>Choice of Data Labels</b>  What labels did you decide to add to your data? And why did you decide on these labels vs any other option?	<p>I went with the binary classification with additional options for uncertainty. There is the "yes" and "no" label, but I also included the uncertain option. Since this is a binary classification, even if annotators select randomly, they have a 50/50 chance to annotate correctly. To solve this uncertainty, I have put two additional labels. When the annotator chooses "yes" or "no", they must answer an additional question on how confident they are from 0-5.</p> <p>Optionally, there is also a blank open text box for annotators to provide their reasoning on their answer.</p> <p>The uncertain option also leaves room for unknown cases.</p>

## Test Questions & Quality Assurance

<h3>Number of Test Questions</h3> <p>Considering the size of this dataset, how many test questions did you develop to prepare for launching a data annotation job?</p>	<p>8 with an even answer distribution. (Yes and No)</p>												
<h3>Improving a Test Question</h3> <p>Given the following test question which almost 100% of annotators missed, statistics, what steps might you take to improve or redesign this question?</p>	<div><table><tr><th>ID</th><th>% CONTESTED</th><th>% MISSED</th><th>JUDGMENTS</th><th>LAST UPDATED</th><th>ENABLED</th></tr><tr><td>1881190030</td><td><div><div></div></div></td><td><div><div></div></div></td><td>2</td><td>2 days ago</td><td><input checked="" type="checkbox"/></td></tr></table></div> <p>Augment the instructions, include more examples or redesign the job.</p>	ID	% CONTESTED	% MISSED	JUDGMENTS	LAST UPDATED	ENABLED	1881190030	<div><div></div></div>	<div><div></div></div>	2	2 days ago	<input checked="" type="checkbox"/>
ID	% CONTESTED	% MISSED	JUDGMENTS	LAST UPDATED	ENABLED								
1881190030	<div><div></div></div>	<div><div></div></div>	2	2 days ago	<input checked="" type="checkbox"/>								
<h3>Contributor Satisfaction</h3> <p>Say you’ve run a test launch and gotten back results from your annotators; the instructions and test questions are rated below 3.5, what areas of your Instruction document would you try to improve (Examples, Test Questions, etc.)</p>	<div><h4>Contributor Satisfaction</h4><p>Number of participants: 20</p><p><b>3.2</b> / 5</p><p>Overall</p><div><div><b>3.3</b> / 5</div><div><b>2.9</b> / 5</div><div><b>2.8</b> / 5</div><div><b>3.7</b> / 5</div></div><div><div>Instructions Clear</div><div>Test Questions Fair</div><div>Ease Of Job</div><div>Pay</div></div></div> <p>All areas that are below 3.5 should be improved but starting with the lowest rated area. In this example, it would be ease of the job.</p>												

Limitations & Improvements

<b>Data Source</b>  Consider the size and source of your data; what biases are built into the data and how might the data be improved?	<p>I notice data contained 117 rows; this is an uneven number. This will ultimately lead to an unbalanced result.</p> <p>The data seems to specifically target Pneumonia. However, there may be similar unknown cases that are visually very similar to Pneumonia.</p>
<b>Designing for Longevity</b>  How might you improve your data labelling job, test questions, or product in the long term?	<p>As data increase, more likelihood increases on other diseases. I may have to include other potential diseases and discard the binary classification model. As a result, I will also need to provide more test questions and train on this new data point.</p>