

# Project Proposal



Jorge Luiz Figueira da Silva Junior

---

## Data Labeling Approach

### Project Overview and Goal

What is the industry problem you are trying to solve? Why use ML in solving this task?

Pneumonia is a respiratory disease with a high mortality rate, especially in children and the elderly. A quick diagnosis would facilitate early treatment, thus resulting in better clinical results.

In this project, data labeling job is designed to build a labeled dataset that distinguishes between healthy and pneumonia x-ray images.

Rapid radiological interpretation is not always available, especially in resource-poor settings. Machine Learning Techniques have the potential to contribute to health professionals, assisting in the diagnosis of diseases. They allow quick classification of complex cases for huge amounts of images, with an accuracy superior to manual detection.

### Choice of Data Labels

What labels did you decide to add to your data? And why did you decide on these labels vs any other option?

Given that the labeled data will be used to develop a product capable of identifying the presence or not of pneumonia on x-ray images, this is clearly a binary classification problem. Then we defined the labels with yes as present and not as absent. An unknown label was included to account for uncertainty in the annotation.

This method of data labeling was chosen because the objective of the problem is to find out whether pneumonia is indicated or not. On the other hand, a numerical scale approach (Likert Scale for example) could also be used, but taking into account non-experienced annotators, the first approach is simpler.

# Test Questions & Quality Assurance

<div><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></div>	<p>Considering using 100 rows for a standard launch. I believe that 20 test questions are enough to help annotators understand what they are doing and thus maintain their accuracy.</p>												
<div><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>	<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></td><td><div></div></td><td>2</td><td>2 days ago</td><td><input checked="" type="checkbox"/></td></tr></table></div> <p>These statistics reveal that we may have a problem with the instructions or the level of certainty of the annotators. Therefore, the steps to take are to detail the instructions or include more examples. Also in this specific case, we could take this question and bring it to the instructions as an example.</p>	ID	% CONTESTED	% MISSED	JUDGMENTS	LAST UPDATED	ENABLED	1881190030	<div></div>	<div></div>	2	2 days ago	<input checked="" type="checkbox"/>
ID	% CONTESTED	% MISSED	JUDGMENTS	LAST UPDATED	ENABLED								
1881190030	<div></div>	<div></div>	2	2 days ago	<input checked="" type="checkbox"/>								
<div><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>	<div><div><h4>Contributor Satisfaction</h4><p>Number of participants: 20</p><div><div>3.2 / 5</div><div>Overall</div></div><div><div>3.3 / 5</div><div>Instructions Clear</div></div><div><div>2.9 / 5</div><div>Test Questions Fair</div></div><div><div>2.8 / 5</div><div>Ease Of Job</div></div><div><div>3.7 / 5</div><div>Pay</div></div></div></div> <p>Ease of job was the most critical point. This may indicate that the instructions provided were not enough for the annotators. To mitigate this effect, the questions would be revised, and further commented examples would be provided.</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>As described by the data source, the images provided were taken with slightly different sizes and taken with slightly different exposure times. These effects on images can directly affect accuracy.</p> <p>In this scenario, different biases are being incorporated:</p> <p>Measurement bias: when incorrect measurements result in data distortion;</p> <p>Recall bias: Since some healthy lungs may appear cloudy, annotators may mislabel and data may be inconsistent.</p> <p>To improve the data, an image collection standardization process should be adopted, with future images being taken by the same devices and usage settings.</p>
<b>Designing for Longevity</b>  How might you improve your data labeling job, test questions, or product in the long-term?	<p>It is an iterative process. Tracking the satisfaction feedback, identifying which questions are frequently missed is an opportunity to improve the data labeling job, redesign questions, rules, tips, increasing the number of examples.</p> <p>Additionally, the inclusion of authorities in the field in the project (such as doctors and other specialists) could enhance the quality of the instructions.</p>