

# Project Proposal



Mohamed Amr

---

## 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?	I am trying to build a ML system for the healthcare industry to aid in quickly identifying healthy patients and surfacing potential cases 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?	Basically, I added two main labels that annotators could make their decision on. The first label is for <b>pneumonia</b> images, and the second for <b>normal</b> images. Moreover, I added a third label ( <b>Not Sure</b> ) to guide annotators. I ask annotators to describe how likely they think a case of pneumonia is in a given image, by rating it from 1 to 5.

# Test Questions & Quality Assurance

<div><b>Number of Test Questions</b></div> <div>Considering the size of this dataset, how many test questions did you develop to prepare for launching a data annotation job?</div>	<div>I will consider the ground truth data. So, I will add sixteen test questions.</div>												
<div><b>Improving a Test Question</b></div> <div>Given the following test question which almost 100% of annotators missed, statistics, what steps might you take to improve or redesign this question?</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></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> <div><div>I may need to augment the instruction or include more examples or redesign the job.</div></div>	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"/>								
<div><b>Contributor Satisfaction</b></div> <div>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.)</div>	<div><div><div><div><b>Contributor Satisfaction</b></div><div>Number of participants: 20</div><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><div><div>I will add more reasons on the test questions, add more tricky cases in the examples, and I will add more examples for each possible data annotation. Also, I will improve the steps and instructions to make the annotation job more clear.</div></div></div>												

## 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>The original dataset contains subfolders for each image category (Pneumonia/Normal). There are 5,863 X-Ray images and 2 categories (Pneumonia/Normal).</p> <p>It's biased for children X-Ray. So, we can improve it by collecting more X-Ray images for different age groups. Moreover, we may consider gender effect, and we can add more labels to the images with (Male, Female, Other .. etc. )</p>
<b>Designing for Longevity</b>  How might you improve your data labeling job, test questions, or product in the long-term?	<p>I can improve data labeling job by designing for failures and longevity. As I mentioned above, we can improve it by collecting more X-Ray images which is labeled for different age groups and genders.</p>