# Project Proposal

*Hussain Al-Balhareth*

## Data Labeling Approach

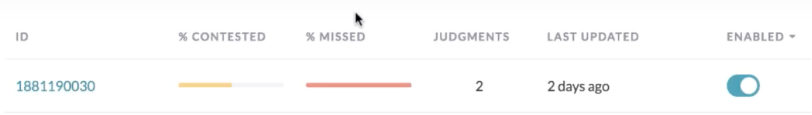| | |
|---|---|
| **Project Overview and Goal**<br><br>What is the industry problem you are trying to solve? Why use ML in solving this task? | - Build a product that helps doctors quickly identify cases of pneumonia in children<br>- Build a labeled dataset that distinguishes between healthy and pneumonia x-ray images that can be used by ML engineers later on down the line to build a classification product.<br>- Create a data labeling job using  Appen's platform. |
| **Choice of Data Labels**<br><br>What labels did you decide to add to your data? And why did you decide on these labels vs any other option? | - Label 0 for healthy and label 1 for pneumonia case.<br>- Such labels are numeric and helpful in binary classification using ML |

## Test Questions & Quality Assurance

### Test questions:

- Does this xray image indicate pneumonia case? (required)

### Quality assurance:

- How confident are you with your assessment? (required)

| | |
|---|---|
| **Number of Test Questions**<br><br>Considering the size of this dataset, how many test questions did you develop to prepare for launching a data annotation job? | 9 test questions out of 117 cases which is more than 5%. |
| **Improving a Test Question**<br><br>Given the following test question which almost 100% of annotators missed, statistics, what steps might you take to improve or redesign this question? | <br><br>We may augment the instructions or include more examples or such tricky cases. |
| **Contributor Satisfaction**<br><br>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.) | <br><br>**Contributor Satisfaction** ⓘ<br>Number of participants: 20<br><br>**3.2** / 5<br>Overall<br><br>**3.3** / 5 — Instructions Clear<br>**2.9** / 5 — Test Questions Fair<br>**2.8** / 5 — Ease Of Job<br>**3.7** / 5 — Pay<br><br>I will focus on all of them but on a priority basis starting with more clarifying examples andTest Questions, then Overview/Steps/Rules Tips. |

# Limitations & Improvements

| | |
|---|---|
| **Data Source**<br><br>Consider the size and source of your data; what biases are built into the data and how might the data be improved? | - Images are for chest x-ray with different sizes and exposure times.<br>- Classification job is a bit tricky as annotators are not specialist in such field like real doctors. So, the challenge is to make this task doable for non-experts as much as possible.<br>- Also, it's best that the data and the images be evenly distributed between Yes and No, High and Low confidence with variety and diversity. |
| **Designing for Longevity**<br><br>How might you improve your data labeling job, test questions, or product in the long-term? | By accounting for changes in data. I assume in this case the data does not change so a static model is adequate. |