# Project Proposal

*<Sung Jun Eun>*

## Data Labeling Approach

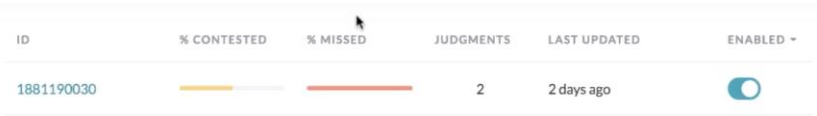| | |
|---|---|
| **Project Overview and Goal**<br><br>What is the industry problem you are trying to solve? Why use ML in solving this task? | It's kind of challenging task for doctors to identify pneumonia quick and accurately. With using ML technology we can train the model to distinguish normal and pneumonia patient easily. ML is used because ML technology can work as good as or better than human if it is trained with large dataset. |
| **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? | I chose label with 'normal', 'pneumonia', and 'uncertain'. 'Uncertain' label is for annotators who cannot distinguish between other two labels. If the annotator distinguished the image as 'pneumonia', I made them to choose their certainty between 1~5. This label is could be used in model to give different weight for different certainty. |

## Test Questions & Quality Assurance

| | |
|---|---|
| **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? | There is one test question per page and if we have 10 rows per page, it leads me to have about 12 test questions for our dataset |
| **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>I should add this question to 'Example' part and warn the annotators to not misunderstand the problem. |
| **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>I may add more examples because low score means that annotators didn't understand the problem well. Also, adding test questions could help them to acknowledge what makes them to make mistake in annotating data. |

# 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? | Size of the dataset Is quite small and all of the x-ray images are taken on the front of the body. If we get the x-ray image taken from the side part of body, It may help to increase the number of data and accuracy of dataset. |
| **Designing for Longevity**<br><br>How might you improve your data labeling job, test questions, or product in the long-term? | By watching the results from the annotators, I will update the 'example' part with the misannotated images by previous annotators. If we update the dataset with the x-ray image taken from the side part of body, I should update the instructions and test questions. |