

I. Milestone

The project and our goals

We wanted to find a topic where we can work with health-related image classification. Finally we chose the dataset of the National Institutes of Health (NIH) about Chest X-Ray images. Our goal is to make the model be able to predict whether the patient has any it of lung disease and if has, what type. In the dataset there are 14 types of diseases, so we have 15 classes (with No Finding):

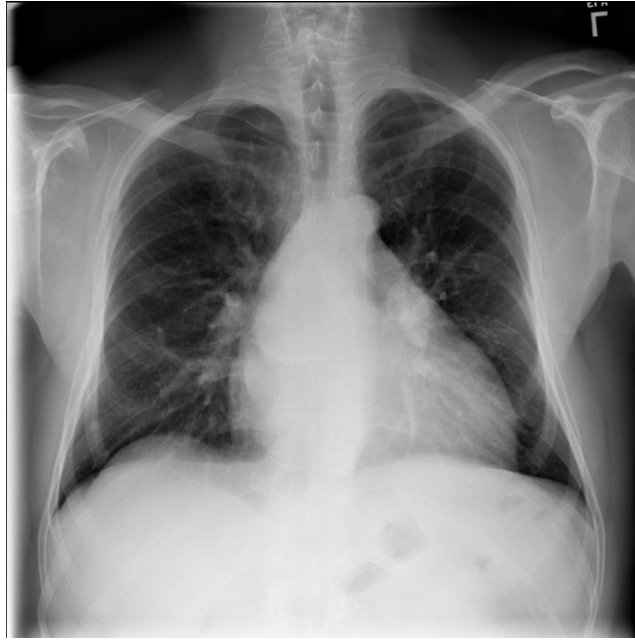
- Atelectasis
- Consolidation
- Infiltration
- Pneumothorax
- Edema
- Emphysema
- Fibrosis
- Effusion
- Pneumonia
- Pleural_thickening
- Cardiomegaly
- Nodule
- Mass
- Hernia
- No Finding

Dataset

We found the dataset on Kaggle (<https://www.kaggle.com/nih-chest-xrays/data>). The size of the dataset is 42 GB. There are 112,120 images with the size of 1024 x 1024. The dataset consists of scans from more than 30,000 patients. The labels were made by the radiologist of NIH Clinical Center. The labels are in a csv file (*Data_entry_2017.csv*) with some added features like age, number of follow up visits, AP vs PA scan, and the patient's gender.

The structure of the *Data_entry_2017.csv* file:

- Image Index: File name
- Finding Labels: Disease type (Class label)
- Follow-up #
- Patient ID
- Patient Age
- Patient Gender
- View Position: X-ray orientation
- OriginalImageWidth
- OriginalImageHeight
- OriginalImagePixelSpacing_x
- OriginalImagePixelSpacing_y



1. An image from the dataset

Existing implementations

We found an existing implementation (<https://github.com/gregwchase/nih-chest-xray>). We used their observations about the dataset during the preprocessing.

Preprocessing

CSV

preproCsv.py

We renamed the original csv file from `Data_entry_2017.csv` to `Data_Entry.csv`. We edited the column names by hand because they contained spaces. The owners of the linked gitHub repository mentioned that there are some corrupted data^s in the dataset: some ages were in days or months, so we deleted 27 rows. We deleted some columns (*Original_Image_Width*, *Original_Image_Height*, *Original_Image_Pixel_Spacing_X*, *Original_Image_Pixel_Spacing_Y*) because they don't have any effect on data.

In the original file there were appr. 709 different categories, because some images belonged to more than one class. We deleted these rows to reduce the number of categories, so 20,729 rows were deleted. There were some rows with unrealistic age (more than 120 years), so 16 rows became deleted.

No Finding	60387
Infiltration	9544
Atelectasis	4210
Effusion	3959
Nodule	2706
Pneumothorax	2198
Mass	2137
Consolidation	1314
Pleural_Thickening	1127
Cardiomegaly	1094
Emphysema	895
Fibrosis	727
Edema	633
Pneumonia	307
Hernia	110

2. The number of images in each categories

We made one-hot encoding on the 'Finding_Labels' column. Finally we shuffled the rows randomly, and created a new dataset, called *entry_data_edited.csv*.

Images

xray_preprocessor.py

This script's main responsibilities are image preprocessing and image selection. First, we had to load the preprocessed csv file (*entry_data_edited.csv*) into a dataframe. In this way, we could easily iterate through the whole dataset, row by row. Because the rows of the csv file had been already mixed, we could fill up the subsets (training, validation, test) one after another, without any kind of randomisation. This means, we could put the first "percent_of_training_data (50%) * len(dataset)" amount of data into the training subset, then the next "percent_of_validation_data (20%) * len(dataset)" into the validation subset, and the remaining "percent_of_test_data (30%) * len(dataset)" into the test set. We represented each subset as a pair of lists. One contains the input (*Image_Index*) and there is another one for the output (*Finding_Labels*). Our inputs were pictures, which had been stored as a list of numpy arrays.

Before we put our pictures to one of the subsets, we had to preprocess them. For image loading, we used openCV. We loaded them as grayscale images, than we resized them to 256px * 256px. The last step before the appendation was scaling. We used a minmax scaler for this purpose (featurerange = (0, 1)).

download_decompress.py

We downloaded our dataset manually from Kaggle, but you can use this script to get the same result. It uses the official Kaggle api (<https://github.com/Kaggle/kaggle-api>), so every download starts with a terminal command. You can specify the download location in the config file.

The dataset is compressed and it contains another 12 compressed files. The script extracts them separately, then merges all of these unzipped directories into a common one.

config.py

Contains different kind of paths and constants, which are usually used by multiple scripts/files. For example, you can change here the download and extract locations, image and subset sizes.