

```
In [237]: 1 import xml.etree.ElementTree as ET
2 import pandas as pd
3 import numpy as np
4 import cv2
5 import matplotlib.pyplot as plt
6 #from tqdm.notebook import tqdm
7 import os
8 import re
9 import seaborn as sns
10 from tqdm import tqdm_notebook as tqdm
11 import warnings
12 warnings.filterwarnings('ignore')
13 from collections import defaultdict
14 #import tensorflow as tf
```

Loading Data

```
In [16]: 1 !gdown "https://drive.google.com/drive/u/0/my-drive"

/usr/local/lib/python2.7/dist-packages/gdown/parse_url.py:31: UserWarning: You
specified Google Drive Link but it is not the correct link to download the fil
e. Maybe you should try: https://drive.google.com/uc?id=None (https://drive.goo
gle.com/uc?id=None)
  .format(url='https://drive.google.com/uc?id={}'.format(file_id))
Downloading...
From: https://drive.google.com/drive/u/0/my-drive (https://drive.google.com/dri
ve/u/0/my-drive)
To: /content/my-drive
62.0kB [00:00, 1.99MB/s]
```

```
In [19]: 1 import shutil
2 shutil.unpack_archive("/content/NLMCXR_png.tgz", "/content/NLMCXR_png")
```

```
In [ ]: 1 shutil.unpack_archive("/content/ecgen-radiology.zip")
```

```
In [ ]: 1 columns = ["image_id", "caption", "comparison", "indication", "findings", "i
2 df = pd.DataFrame(columns = columns)
```

```
In [20]: 1 len(os.listdir("/content/NLMCXR_png/"))
```

Out[20]: 7471

Extracting data from XML files

```

In [ ]: 1 for file in tqdm(os.listdir("ecgen-radiology/")):
        2     if file.endswith(".xml"):
        3         k = "ecgen-radiology/"
        4         path = k + file
        5         mytree = ET.parse(path) # parsing xml report
        6         comparision = mytree.find(".*//AbstractText[@Label='COMPARISON']").text
        7         indication = mytree.find(".*//AbstractText[@Label='INDICATION']").text
        8         findings = mytree.find(".*//AbstractText[@Label='FINDINGS']").text #
        9         impression = mytree.find(".*//AbstractText[@Label='IMPRESSION']").text
        10
        11     mytree = ET.parse(path)
        12     for x in mytree.findall("parentImage"):
        13         image_id = x.attrib['id']+".png"
        14         filename = 'NLMCXR_png/' + image_id
        15         image = cv2.imread(filename) # reading image
        16
        17         height, width, channels = image.shape
        18         caption = '' if x.find('caption').text is None else x.find('caption').text
        19
        20         df = df.append(pd.Series([image_id, caption, comparision, indication,
        21                                 index = columns], index = columns), i

```

HBox(children=(FloatProgress(value=0.0, max=3956.0), HTML(value='')))

```

In [ ]: 1 df.shape

```

Out[64]: (7470, 8)

In []:

1 df

Out[65]:

	image_id	caption	comparison	indication	findings	impression
0	CXR1082_IM-0058-1001.png	Chest x-XXXX XXXX and lateral performed on XXX...	Chest x-XXXX XXXX and lateral from XXXX.	XXXX year old female with abdominal pain.	Stable cardiomegaly. Stable tortuosity of the ...	Stable cardiomegaly with clear lungs.
1	CXR473_IM-2101-1001.png	PA and lateral chest.	None	preop for XXXX	None	Heart size normal. Lungs clear.
2	CXR473_IM-2101-1002.png	PA and lateral chest.	None	preop for XXXX	None	Heart size normal. Lungs clear.
3	CXR1883_IM-0572-1001.png	Xray Chest PA and Lateral	None.	XXXX onset of right-sided weakness for one XXXX.	Frontal and lateral views of the chest show no...	No acute or active cardiac, pulmonary or pleur...
4	CXR1883_IM-0572-2001.png	Xray Chest PA and Lateral	None.	XXXX onset of right-sided weakness for one XXXX.	Frontal and lateral views of the chest show no...	No acute or active cardiac, pulmonary or pleur...
...
7465	CXR2421_IM-0965-2001.png	CHEST 2V FRONTAL/LATERAL	None	Dyspnea	The heart is enlarged. The left subclavian ICD...	Stable moderate to marked cardiomegaly.
7466	CXR3165_IM-1490-1001.png	Chest x-XXXX, 2 views, XXXX, XXXX XXXX PM	XXXX	Dyspnea	Normal and stable cardiomedialstinal contours. ...	Mildly improved XXXX XXXX opacities, which may...
7467	CXR3165_IM-1490-13013.png	Chest x-XXXX, 2 views, XXXX, XXXX XXXX PM	XXXX	Dyspnea	Normal and stable cardiomedialstinal contours. ...	Mildly improved XXXX XXXX opacities, which may...
7468	CXR1108_IM-0075-1001.png	Xray Chest PA and Lateral	None.	SHORTNESS OF BREATH;	The lungs are clear and hyperinflated. Heart s...	Hyperinflated lungs. No acute cardiopulmonary ...
7469	CXR1108_IM-0075-2001.png	Xray Chest PA and Lateral	None.	SHORTNESS OF BREATH;	The lungs are clear and hyperinflated. Heart s...	Hyperinflated lungs. No acute cardiopulmonary ...

7470 rows × 8 columns



```
In [ ]: 1 df[df['image_id']=='CXR1_1_IM-0001-3001.png']
```

```
Out[66]:
```

	image_id	caption	comparison	indication	findings	impression	height	width
6086	CXR1_1_IM-0001-3001.png	Xray Chest PA and Lateral	None.	Positive TB test	The cardiac silhouette and mediastinum size ar...	Normal chest x-XXXX.	624	512

```
In [ ]: 1 df.shape
```

```
Out[53]: (7470, 8)
```

```
In [ ]: 1 def absolute_path(x):
2         '''Makes the path absolute '''
3         x = 'NLMCXR_png/' + x
4         return x
5
6 df['Image_path'] = df['image_id'].apply(lambda x : absolute_path(x)) # makin
```

```
In [ ]: 1 df.head(5)
```

```
Out[69]:
```

	image_id	caption	comparison	indication	findings	impression	height	width
0	CXR1082_IM-0058-1001.png	Chest x-XXXX XXXX and lateral performed on XXX...	Chest x-XXXX XXXX and lateral from XXXX.	XXXX year old female with abdominal pain.	Stable cardiomegaly. Stable tortuosity of the ...	Stable cardiomegaly with clear lungs.	624	512
1	CXR473_IM-2101-1001.png	PA and lateral chest.	None	prep for XXXX	None	Heart size normal. Lungs clear.	510	512
2	CXR473_IM-2101-1002.png	PA and lateral chest.	None	prep for XXXX	None	Heart size normal. Lungs clear.	601	512
3	CXR1883_IM-0572-1001.png	Xray Chest PA and Lateral	None.	XXXX onset of right-sided weakness for one XXXX.	Frontal and lateral views of the chest show no...	No acute or active cardiac, pulmonary or pleur...	502	512
4	CXR1883_IM-0572-2001.png	Xray Chest PA and Lateral	None.	XXXX onset of right-sided weakness for one XXXX.	Frontal and lateral views of the chest show no...	No acute or active cardiac, pulmonary or pleur...	512	512

In []:

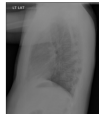
```

1 count = 1
2 fig = plt.figure(figsize=(15,35))
3
4 for filename in df['Image_path'].values[95:100]:
5     findings = list(df["findings"].loc[df["Image_path"] == filename].values)
6     img = cv2.imread(filename)
7     ax = fig.add_subplot(5, 2, count, xticks=[], yticks=[])
8     ax.imshow(img)
9     count += 1
10    ax = fig.add_subplot(5, 2, count)
11    plt.axis('off')
12    ax.plot()
13    ax.set_xlim(0,1)
14    ax.set_ylim(0, len(findings))
15    for i, f in enumerate(findings):
16        ax.text(0,i,f,fontsize=20)
17    count += 1
18 plt.show()

```



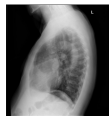
No focal lung consolidation. Heart size and pulmonary vascularity are within normal limits. No pneumothorax or pleural effusion. Osseous structures are grossly intact.



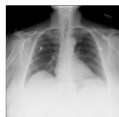
No focal lung consolidation. Heart size and pulmonary vascularity are within normal limits. No pneumothorax or pleural effusion. Osseous structures are grossly intact.



The heart and lungs have XXXX XXXX in the interval. Both lungs are clear and expanded. Heart and mediastinum normal. No change right anterior soft tissue surgical clips. Configuration of breast shadows on the PA view suggests prior right lumpectomy.



The heart and lungs have XXXX XXXX in the interval. Both lungs are clear and expanded. Heart and mediastinum normal. No change right anterior soft tissue surgical clips. Configuration of breast shadows on the PA view suggests prior right lumpectomy.



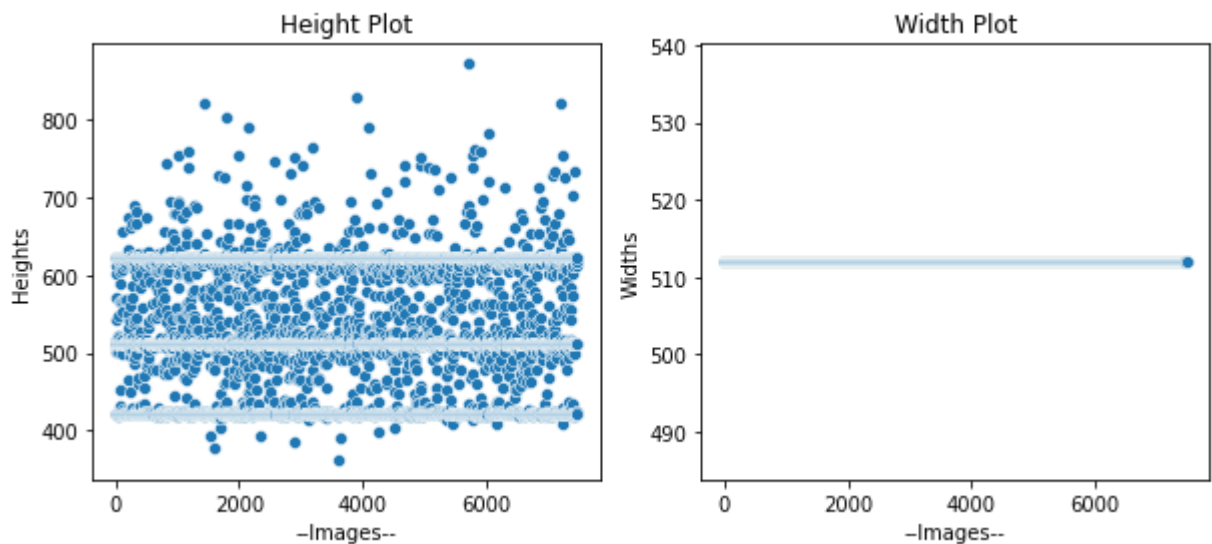
Negative for cardiac enlargement or vascular congestion. Minimal subsegmental atelectasis at the left base otherwise negative for focal confluent airspace disease. The visualized bony structures are intact. There are minimal degenerative disc changes of the mid/lower thoracic spine. No pneumothorax.

```

In [ ]: 1 plt.figure(figsize=(10,4))
        2 plt.subplot(121)
        3 plt.title('Height Plot')
        4 plt.ylabel('Heights')
        5 plt.xlabel('--Images--')
        6 sns.scatterplot(range(len(df.height.values)), df.height.values)
        7 plt.subplot(122)
        8 plt.title('Width Plot')
        9 plt.ylabel('Widths')
        10 plt.xlabel('--Images--')
        11 sns.scatterplot(range(len(df.width.values)), df.width.values)

```

Out[71]: <matplotlib.axes._subplots.AxesSubplot at 0x7ff023da0550>



```

In [ ]: 1 # number of missing values
        2 df.isnull().sum()

```

Out[72]:

image_id	0
caption	0
comparison	1157
indication	159
findings	997
impression	52
height	0
width	0
Image_path	0
dtype:	int64

```
In [ ]: 1 data = df[['image_id', 'findings', 'height', 'width', 'Image_path']]
```

```
In [ ]: 1 data.shape
```

```
Out[74]: (7470, 5)
```

```
In [ ]: 1 data.isnull().sum()
```

```
Out[75]: image_id      0
         findings    997
         height      0
         width       0
         Image_path   0
         dtype: int64
```

```
In [ ]: 1 data = data.dropna(axis=0) # drop all missing value rows
```

```
In [ ]: 1 data.shape
```

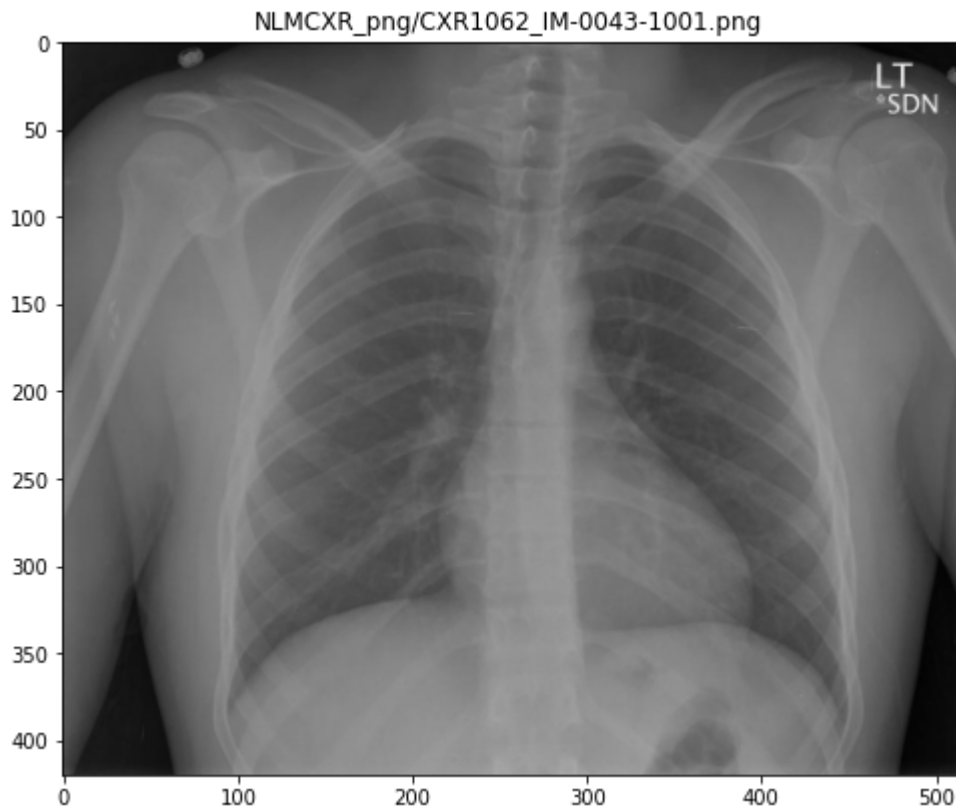
```
Out[77]: (6473, 5)
```

```
In [ ]: 1 data.isnull().sum()
```

```
Out[78]: image_id      0
         findings      0
         height      0
         width       0
         Image_path   0
         dtype: int64
```

```
In [ ]: 1 plt.figure(figsize=(8,7))
        2 img = cv2.imread(data['Image_path'].values[5])
        3 plt.imshow(img)
        4 plt.title(data['Image_path'].values[5])
        5
        6 data['findings'].values[5]
```

Out[79]: 'XXXX XXXX and lateral chest examination was obtained. The heart silhouette is normal in size and contour. Aortic XXXX appear unremarkable. Lungs demonstrate no acute findings. There is no effusion or pneumothorax. No displaced rib fracture visualized.'




```
In [ ]: 1 data.Image_path
```

```
Out[81]: 0      NLMCXR_png/CXR1082_IM-0058-1001.png
        3      NLMCXR_png/CXR1883_IM-0572-1001.png
        4      NLMCXR_png/CXR1883_IM-0572-2001.png
        5      NLMCXR_png/CXR2431_IM-0973-1001.png
        6      NLMCXR_png/CXR2431_IM-0973-2001.png
        ...
        7465     NLMCXR_png/CXR2421_IM-0965-2001.png
        7466     NLMCXR_png/CXR3165_IM-1490-1001.png
        7467     NLMCXR_png/CXR3165_IM-1490-13013.png
        7468     NLMCXR_png/CXR1108_IM-0075-1001.png
        7469     NLMCXR_png/CXR1108_IM-0075-2001.png
        Name: Image_path, Length: 6473, dtype: object
```

```
In [ ]: 1 images = {}
        2 findings = {}
        3
        4 for img,fin in data[['Image_path','findings']].values:
        5     a = img.split('.')
        6     file_type = a[-1]
        7     a = a[0].split('-')
        8     a.pop(len(a)-1)
        9     a = ('-'.join(e for e in a))
       10     if a not in images.keys():
       11         images[a] = 1
       12         findings[a] = fin
       13     else:
       14         images[a] += 1
       15         findings[a] = fin
```

```
In [ ]: 1 images['NLMCXR_png/CXR1001_IM-0004'],findings['NLMCXR_png/CXR1001_IM-0004']
```

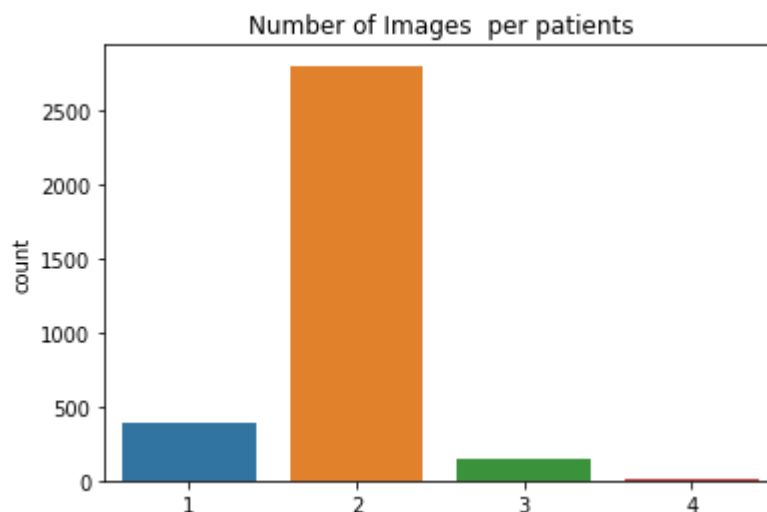
```
Out[83]: (2,
        'Interstitial markings are diffusely prominent throughout both lungs. Heart si
        ze is normal. Pulmonary XXXX normal.')
```

```
In [ ]: 1 print('Total Number of Unique_IDs :', len(images.keys()))
```

```
Total Number of Unique_IDs : 3350
```

```
In [ ]: 1 plt.title('Number of Images per patients')
        2 sns.countplot(list(images.values()))
```

Out[86]: <matplotlib.axes._subplots.AxesSubplot at 0x7ff02fbd7128>



```
In [ ]: 1 def train_test_split(data):
        2     persons = list(data.keys())
        3     persons_train = persons[:2500]
        4     persons_cv = persons[2500:3000]
        5     persons_test = persons[3000:3350]
        6     return persons_train, persons_cv, persons_test
        7
        8 images_train, images_cv, images_test = train_test_split(images)
```

```
In [ ]: 1 def combining_images(image_set):
2
3     image_per_person = defaultdict(list) # creating a list of dictionary to
4                                           #corresponding to a person_id
5     for pid in image_set:
6         for img in data['Image_path'].values:
7             if pid in img:
8                 image_per_person[pid].append(img)
9             else:
10                continue
11     return image_per_person
```

```
In [ ]: 1 img_per_person_train = combining_images(images_train)
2       2 img_per_person_cv = combining_images(images_cv)
3       3 img_per_person_test = combining_images(images_test)
```

```
In [ ]: 1 img_per_person_train['NLMCXR_png/CXR1001_IM-0004']
```

```
Out[95]: ['NLMCXR_png/CXR1001_IM-0004-1001.png', 'NLMCXR_png/CXR1001_IM-0004-1002.png']
```

```
In [ ]: 1 def load_image(file):
2       2 img = cv2.imread(file)
3       3 return img
```

```
In [ ]: 1 # just checking the ID which has 4 images
2       2 for k,v in images.items():
3       3     if v == 4:
4       4         print(k)
5       5         break
6
```

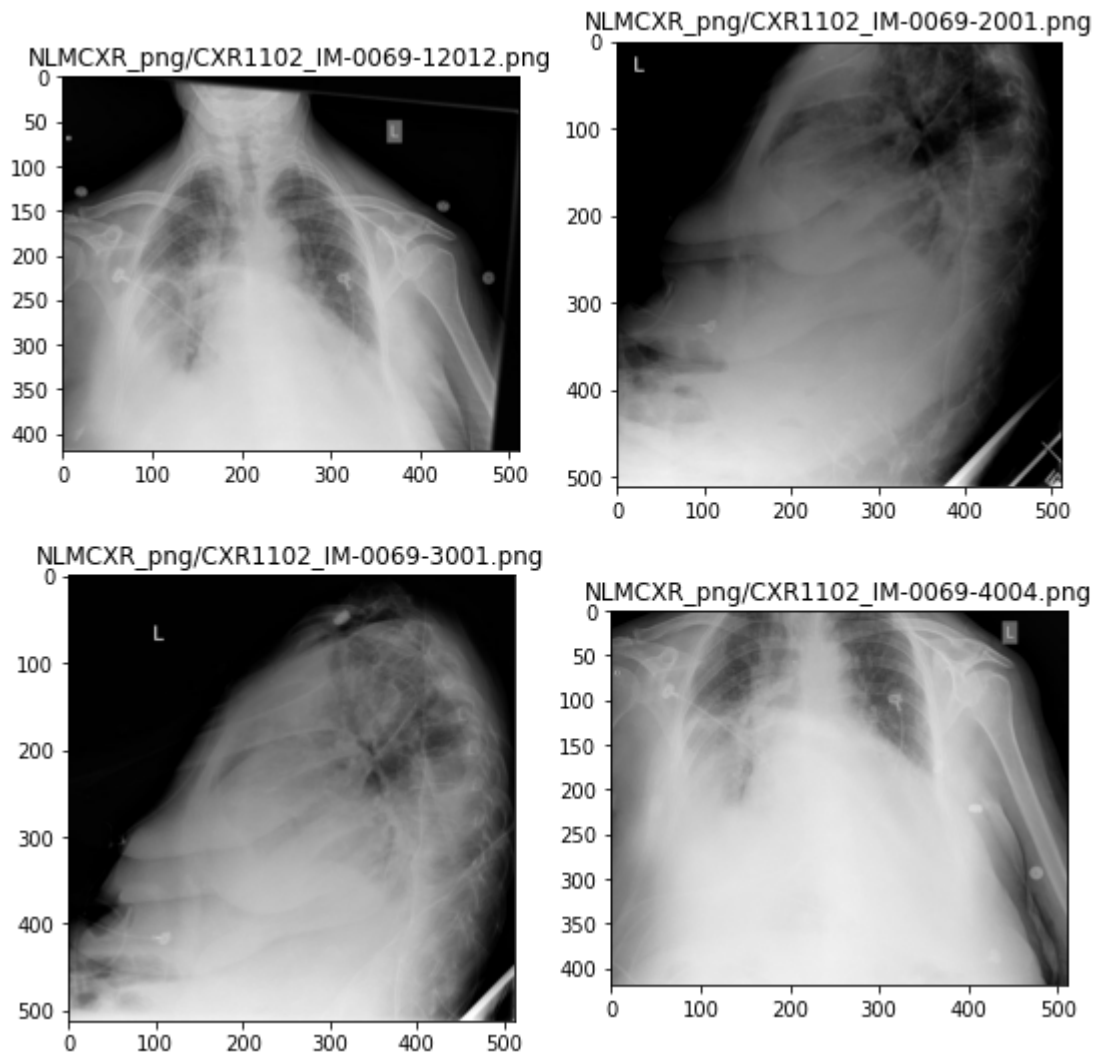
NLMCXR_png/CXR1102_IM-0069

```

In [ ]: 1 plt.figure(figsize=(9,9))
        2 plt.subplot(221)
        3 plt.imshow(load_image('NLMCXR_png/CXR1102_IM-0069-12012.png'))
        4 plt.title('NLMCXR_png/CXR1102_IM-0069-12012.png')
        5 plt.subplot(222)
        6 plt.imshow(load_image('NLMCXR_png/CXR1102_IM-0069-2001.png'))
        7 plt.title('NLMCXR_png/CXR1102_IM-0069-2001.png')
        8 plt.subplot(223)
        9 plt.imshow(load_image('NLMCXR_png/CXR1102_IM-0069-3001.png'))
       10 plt.title('NLMCXR_png/CXR1102_IM-0069-3001.png')
       11 plt.subplot(224)
       12 plt.imshow(load_image('NLMCXR_png/CXR1102_IM-0069-4004.png'))
       13 plt.title('NLMCXR_png/CXR1102_IM-0069-4004.png')

```

Out[98]: Text(0.5, 1.0, 'NLMCXR_png/CXR1102_IM-0069-4004.png')



2 side view and 2 front view images for the same ID

Sample chest scans of a person(4 images)

Now, we have multiple chest scans to produce a single report. Some person_ids have 1, some have 2 and the highest is 4. So we can take pairs of those images as input. If it has only one image, then it can be replicated.

Data Preperation

```
In [ ]: 1 import itertools
2
3 def create_data(image_per_person):
4     # new dataset
5     person_id, image1, image2, report = [],[],[],[]
6     for pid, imgs in image_per_person.items(): #contains pid and the image
7
8         if len(imgs) == 1:
9             image1.append(imgs[0])
10            image2.append(imgs[0])
11            person_id.append(pid)
12            report.append(findings[pid])
13        else:
14            num = 0
15            a = itertools.combinations(imgs, 2)
16            for i in a:
17                image1.append(i[0])
18                image2.append(i[1])
19                person_id.append(pid + '_' + str(num))
20                report.append(findings[pid])
21                num += 1
22        data = pd.DataFrame()
23        data['Person_id'] = person_id
24        data['Image1'] = image1
25        data['Image2'] = image2
26        data['Report'] = report
27
28    return data
```

```
In [ ]: 1 train = create_data(img_per_person_train)
        2 test = create_data(img_per_person_test)
        3 cv = create_data(img_per_person_cv)
```

```
In [ ]: 1 train.head()
```

```
Out[101]:
```

	Person_id	Image1	Image2	Report
0	NLMCXR_png/CXR1082_IM-0058	NLMCXR_png/CXR1082_IM-0058-1001.png	NLMCXR_png/CXR1082_IM-0058-1001.png	Stable cardiomegaly. Stable tortuosity of the ...
1	NLMCXR_png/CXR1883_IM-0572_0	NLMCXR_png/CXR1883_IM-0572-1001.png	NLMCXR_png/CXR1883_IM-0572-2001.png	Frontal and lateral views of the chest show no...
2	NLMCXR_png/CXR2431_IM-0973_0	NLMCXR_png/CXR2431_IM-0973-1001.png	NLMCXR_png/CXR2431_IM-0973-2001.png	Lungs are clear. No focal airspace consolidati...
3	NLMCXR_png/CXR1062_IM-0043	NLMCXR_png/CXR1062_IM-0043-1001.png	NLMCXR_png/CXR1062_IM-0043-1001.png	XXXX XXXX and lateral chest examination was ob...
4	NLMCXR_png/CXR224_IM-0837_0	NLMCXR_png/CXR224_IM-0837-1001.png	NLMCXR_png/CXR224_IM-0837-2001.png	Clear lungs. Normal heart. No pneumothorax. No...

```
In [ ]: 1 train.to_csv('train.csv')
        2 test.to_csv('test.csv')
        3 cv.to_csv('cv.csv')
```

Preporcessing text data

```

In [ ]: 1 def lowercase(text):
        2     '''Converts to lowercase'''
        3     new_text = []
        4     for line in text:
        5         new_text.append(line.lower())
        6     return new_text
        7
        8 def decontractions(text):
        9     '''Performs decontractions in the doc'''
       10     new_text = []
       11     for phrase in text:
       12         phrase = re.sub(r"won't", "will not", phrase)
       13         phrase = re.sub(r"can't", "can not", phrase)
       14         phrase = re.sub(r"couldn't", "could not", phrase)
       15         phrase = re.sub(r"shouldn't", "should not", phrase)
       16         phrase = re.sub(r"wouldn't", "would not", phrase)
       17         # general
       18         phrase = re.sub(r"n't", " not", phrase)
       19         phrase = re.sub(r"\'re", " are", phrase)
       20         phrase = re.sub(r"\'s", " is", phrase)
       21         phrase = re.sub(r"\'d", " would", phrase)
       22         phrase = re.sub(r"\'ll", " will", phrase)
       23         phrase = re.sub(r"\'t", " not", phrase)
       24         phrase = re.sub(r"\'ve", " have", phrase)
       25         phrase = re.sub(r"\'m", " am", phrase)
       26         phrase = re.sub(r"\*+", "abuse", phrase)
       27         new_text.append(phrase)
       28
       29     return new_text
       30
       31 def rem_punctuations(text):
       32     '''Removes punctuations'''
       33     punctuations = '!'()-[]{};:'"\,<>/?@#$$%^&*~' # full stop is not removed
       34     new_text = []
       35     for line in text:
       36         for char in line:
       37             if char in punctuations:
       38                 line = line.replace(char, "")
       39         new_text.append(' '.join(e for e in line.split()))
       40     return new_text
       41
       42 def rem_numbers(text):
       43     '''Removes numbers and irrelevant text like xxxx'''
       44     new_text = []
       45     for line in text:
       46         temp = re.sub(r'x*', '', line)
       47         new_text.append(re.sub(r'\d', '', temp))
       48     return new_text
       49
       50 def words_filter(text):
       51     '''Removes words less than 2 characters except no and ct'''
       52     new_text = []
       53     for line in text:
       54         temp = line.split()
       55         temp2 = []
       56         for word in temp:

```

```
57         if len(word) <=2 and word != 'no' and word != 'ct':
58             continue
59         else:
60             temp2.append(word)
61             new_text.append(' '.join(e for e in temp2))
62     return new_text
63
64 def multiple_fullstops(text):
65     ''' Removes multiple full stops from the text'''
66     new_text = []
67     for line in text:
68         new_text.append(re.sub(r'\.\.+', '.', line))
69     return new_text
70
71 def fullstops(text):
72     new_text = []
73     for line in text:
74         new_text.append(re.sub('\.', ' .', line))
75     return new_text
76
77 def multiple_spaces(text):
78     new_text = []
79     for line in text:
80         new_text.append(' '.join(e for e in line.split()))
81     return new_text
82
83 def sepating_startg_words(text):
84     new_text = []
85     for line in text:
86         temp = []
87         words = line.split()
88         for i in words:
89             if i.startswith('.') == False:
90                 temp.append(i)
91             else:
92                 w = i.replace('.', '. ')
93                 temp.append(w)
94         new_text.append(' '.join(e for e in temp))
95     return new_text
96
97 def rem_apostrophes(text):
98     new_text = []
99     for line in text:
100         new_text.append(re.sub("'", '', line))
101     return new_text
```



```
In [ ]: 1 def text_preprocessing(text):
        2     '''Combines all the preprocess functions'''
        3     new_text = lowercase(text)
        4     new_text = decontractions(new_text)
        5     new_text = rem_punctuations(new_text)
        6     new_text = rem_numbers(new_text)
        7     new_text = words_filter(new_text)
        8     new_text = multiple_fullstops(new_text)
        9     new_text = fullstops(new_text)
       10     new_text = multiple_spaces(new_text)
       11     new_text = separting_startg_words(new_text)
       12     new_text = rem_apostrophes(new_text)
       13     return new_text
```

```
In [ ]: 1 train['Report'] = text_preprocessing(train['Report'])
        2 test['Report'] = text_preprocessing(test['Report'])
        3 cv['Report'] = text_preprocessing(cv['Report'])
```

In []: 1 train

Out[106]:

	Person_id	Image1	Image2	Report
0	NLMCXR_png/CXR1082_IM-0058	NLMCXR_png/CXR1082_IM-0058-1001.png	NLMCXR_png/CXR1082_IM-0058-1001.png	stable tortu
1	NLMCXR_png/CXR1883_IM-0572_0	NLMCXR_png/CXR1883_IM-0572-1001.png	NLMCXR_png/CXR1883_IM-0572-2001.png	frontal and la views the c show nor
2	NLMCXR_png/CXR2431_IM-0973_0	NLMCXR_png/CXR2431_IM-0973-1001.png	NLMCXR_png/CXR2431_IM-0973-2001.png	lungs are c no focal airs consol
3	NLMCXR_png/CXR1062_IM-0043	NLMCXR_png/CXR1062_IM-0043-1001.png	NLMCXR_png/CXR1062_IM-0043-1001.png	and lateral c eamination obtaine
4	NLMCXR_png/CXR224_IM-0837_0	NLMCXR_png/CXR224_IM-0837-1001.png	NLMCXR_png/CXR224_IM-0837-2001.png	clear lu normal hear pneumothc
...	
2759	NLMCXR_png/CXR2725_IM-1186_0	NLMCXR_png/CXR2725_IM-1186-1001.png	NLMCXR_png/CXR2725_IM-1186-2001.png	cardiac media: contour within
2760	NLMCXR_png/CXR728_IM-2287_0	NLMCXR_png/CXR728_IM-2287-1001.png	NLMCXR_png/CXR728_IM-2287-2001.png	the int diameter nc has devel
2761	NLMCXR_png/CXR943_IM-2439_0	NLMCXR_png/CXR943_IM-2439-3003.png	NLMCXR_png/CXR943_IM-2439-5005.png	heart cardiomedia: silhouette
2762	NLMCXR_png/CXR3379_IM-1627_0	NLMCXR_png/CXR3379_IM-1627-1001.png	NLMCXR_png/CXR3379_IM-1627-2001.png	no pneumoi large pl effusion . bo
2763	NLMCXR_png/CXR2565_IM-1068_0	NLMCXR_png/CXR2565_IM-1068-1002.png	NLMCXR_png/CXR2565_IM-1068-1003.png	the heart nc size . the i tortuou

2764 rows × 4 columns

```
In [ ]: 1 train.to_csv('processed_train.csv')
        2 test.to_csv('processed_test.csv')
        3 cv.to_csv('processed_cv.csv')
```

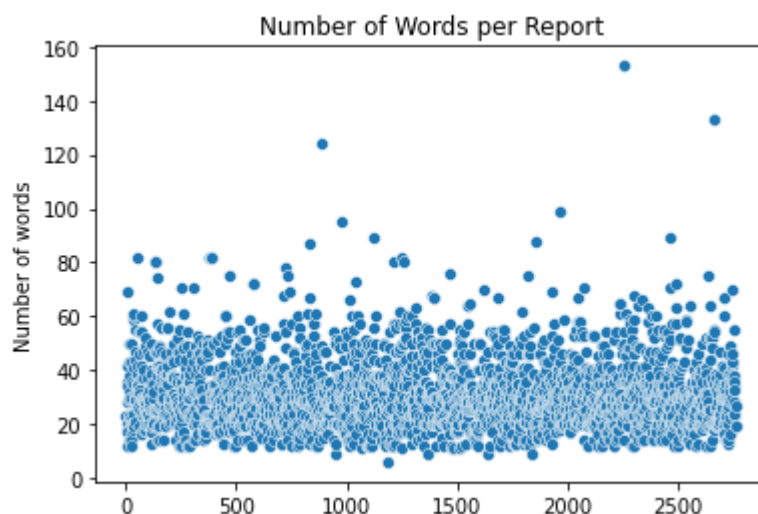
```
In [ ]: 1 l = [len(e.split()) for e in train['Report'].values] # Number of words in e
```

```
In [ ]: 1 print('maximum word in a report is :',max(1))
```

maximum word in a report is : 153

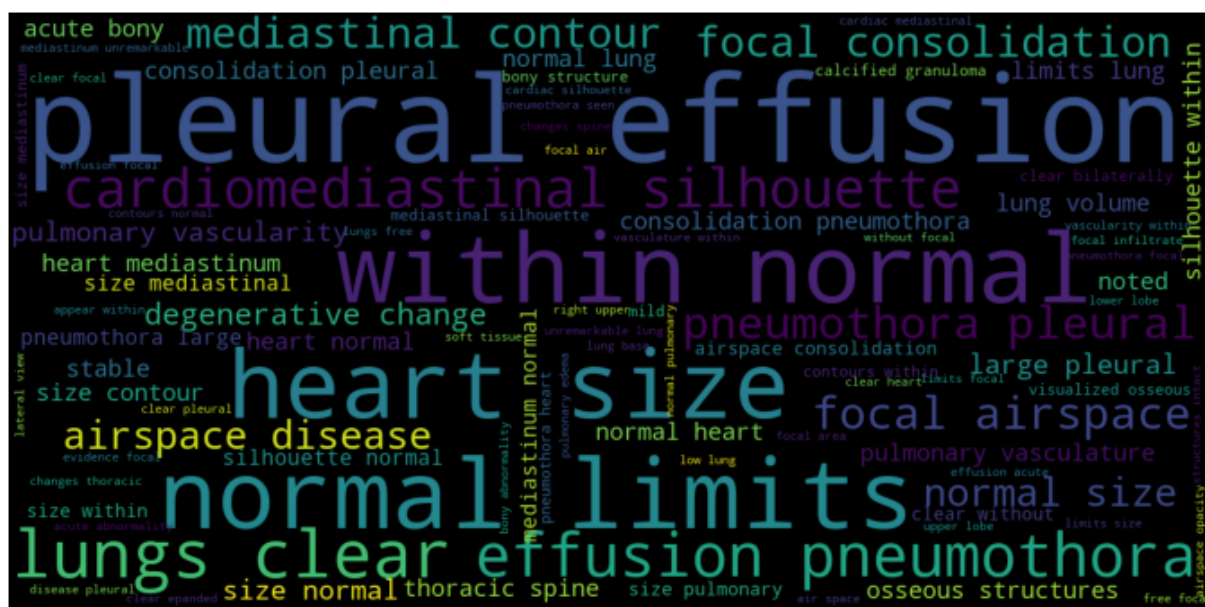
```
In [ ]: 1 plt.title('Number of Words per Report')
2 sns.scatterplot(range(train.shape[0]), 1)
3 plt.ylabel('Number of words')
```

Out[111]: Text(0, 0.5, 'Number of words')



Most of the reports contain word count below 100

```
In [ ]: 1 from wordcloud import WordCloud
2 def show_wordcloud(data, title = None):
3     wordcloud = WordCloud(background_color='black',max_words=800,max_font_si
4
5     fig = plt.figure(1, figsize=(12, 15))
6     plt.axis('off')
7     if title:
8         fig.suptitle(title, fontsize=20)
9         fig.subplots_adjust(top=2.3)
10
11     plt.imshow(wordcloud)
12     plt.show()
13
14 show_wordcloud(train['Report'])
```



```
In [ ]: 1 countword = train['Report'].str.split().apply(len).value_counts()
2 countword[:].plot(kind='bar',figsize=(20,5) , title = 'Words for each findin
```

Out[120]: <matplotlib.axes._subplots.AxesSubplot at 0x7ff022586358>



```
In [ ]: 1 def remodelling(x):
2     '''adds start and end tokens to a sentence '''
3     return 'startseq' + ' ' + x + ' ' + 'endseq'
```

```
In [ ]: 1 train['Report'] = train['Report'].apply(lambda x : remodelling(x))
        2 test['Report'] = test['Report'].apply(lambda x : remodelling(x))
        3 cv['Report'] = cv['Report'].apply(lambda x : remodelling(x))
```

Exporting final dataset to csv file

```
In [ ]: 1 # save the cleaned data(STRUCTURED DATA)
        2 train.to_csv('Final_Train_Data.csv', index=False)
        3 test.to_csv('Final_Test_Data.csv', index=False)
        4 cv.to_csv('Final_CV_Data.csv', index=False)
```