

“Pneumothorax Detection: Feature Extraction, Classification, and Clustering of Chest X-ray Images.”

1. Abstract

Objective:

Pneumothorax is a **thoracic disease** that can lead to respiratory failure, cardiac arrest, or even death. **Chest X-ray (CXR) imaging** is the primary diagnostic method. This study aims to develop a **computerized diagnosis system** for pneumothorax detection using deep learning-based feature extraction and machine learning models.

Methods:

ResNet50 extracts embeddings from DICOM images, which are then used to train **Random Forest** and **XGBoost** classifiers. Additionally, **K-Means clustering** is applied to group similar images and identify patterns in pneumothorax cases.

Results:

The **Random Forest** model achieved **80.63% accuracy**, while **XGBoost** reached **78.75% accuracy**. Clustering results revealed distinct groupings of pneumothorax cases, highlighting similarities among images.

Conclusion:

Machine learning techniques, combined with deep learning-based feature extraction, show promising potential for **automated pneumothorax detection**, aiding in faster and more accurate medical diagnosis.

2. Dataset Analysis

The **Society for Imaging Informatics in Medicine**, in collaboration with **American College of Radiology (SIIM-ACR)**, collected the **CXR data for pneumothorax** and released it on Kaggle.

- The dataset contains: **DICOM images** and **run-length encoded files**. **DICOM (Digital Imaging and Communications in Medicine)** files are a **standard format** consisting of **header data** and an **image**, both of which are packed into a single file for storing **medical images**.
- The **header of the DICOM file** consists of a series of **tags** that provide information concerning the **patient's name, age, sex, demographics**, and various other parameters (**as shown in Figure 1**).
- This project aims to **process DICOM files to extract metadata, enrich the dataset, and perform preliminary analysis**.
- It focuses on **processing, filtering, and analyzing DICOM files for machine learning and medical analysis**.

```

Dataset.file_meta -----
(0002,0000) File Meta Information Group Length  UL: 200
(0002,0001) File Meta Information Version       OB: b'\x00\x01'
(0002,0002) Media Storage SOP Class UID         UI: Secondary Capture Image Storage
(0002,0003) Media Storage SOP Instance UID      UI: 1.2.276.0.7230010.3.1.4.8323329.300.1517875162.258081
(0002,0010) Transfer Syntax UID                UI: JPEG Baseline (Process 1)
(0002,0012) Implementation Class UID           UI: 1.2.276.0.7230010.3.0.3.6.0
(0002,0013) Implementation Version Name        SH: 'OFFIS_DCMTK_360'
-----
(0008,0005) Specific Character Set              CS: 'ISO_IR 100'
(0008,0016) SOP Class UID                      UI: Secondary Capture Image Storage
(0008,0018) SOP Instance UID                   UI: 1.2.276.0.7230010.3.1.4.8323329.300.1517875162.258081
(0008,0020) Study Date                        DA: '19010101'
(0008,0030) Study Time                       TM: '000000.00'
(0008,0050) Accession Number                  SH: ''
(0008,0060) Modality                          CS: 'CR'
(0008,0064) Conversion Type                   CS: 'WSD'
(0008,0090) Referring Physician's Name        PN: ''
(0008,103E) Series Description                 LO: 'view: AP'
(0010,0010) Patient's Name                    PN: '88c14312-3265-4a3f-b7bb-41818107d607'
(0010,0020) Patient ID                       LO: '88c14312-3265-4a3f-b7bb-41818107d607'
(0010,0030) Patient's Birth Date              DA: ''
(0010,0040) Patient's Sex                     CS: 'F'
(0010,1010) Patient's Age                     AS: '58'
(0018,0015) Body Part Examined                CS: 'CHEST'
...
(0028,0102) High Bit                          US: 7
(0028,0103) Pixel Representation              US: 0
(0028,2110) Lossy Image Compression           CS: '01'
(0028,2114) Lossy Image Compression Method    CS: 'ISO_10918_1'
(7FE0,0010) Pixel Data                        OB: Array of 154050 elements

```

Figure 1. Snapshot of metadata stored in a DICOM Image.

2.1 Workflow Overview

The steps involved in this process are:

1. **Filtering DICOM files** To ensure only valid and available DICOM files are kept in the dataset.
2. **Extracting metadata** Filter rows to ensure only valid and available DICOM files are kept in the dataset.
3. **Creating image paths** Generate the full file paths for each DICOM file.
4. **Handling duplicates** Remove duplicate entries in the dataset based on the [ImageId](#).
5. **Calculating pneumothorax area** Use the provided Run-Length Encoding (RLE) to calculate the area affected by pneumothorax.
6. **Cleaning and preparing the dataset** Ensure the dataset is clean, with no missing or duplicate values, and ready for analysis.

The **annotation mask** was stored in the **run-length-encoded (RLE)** file with a **.csv extension**. The RLE file contained two columns, **image ID** and **encoded pixels**, for each figure. In Fig-2.

	ImageId	EncodedPixels
0	1.2.276.0.7230010.3.1.4.8323329.5597.1517875188.959090	-1
1	1.2.276.0.7230010.3.1.4.8323329.12515.1517875239.501137	
2	1.2.276.0.7230010.3.1.4.8323329.4904.1517875185.355709	175349 7 1013 12 1009 17 1005 19 1003 20 1002 22 1001 22 1000 23 1000 23 1000 23 999 24 999 23 1000 23 999 23 1000 22 1001 21 1001 22 1001 21 1002 21 1001 22 1001 21 1002 21 1001 22 1001 21 1002 ...
3	1.2.276.0.7230010.3.1.4.8323329.32579.1517875161.299312	407576 2 1021 7 1015 10 1013 12 1011 14 1008 17 1006 19 1005 20 1003 21 1003 22 1001 23 1001 24 999 25 999 25 999 26 997 27 997 27 996 28 996 28 996 29 994 30 994 30 994 30 993 31 993 32 992 32 9...
4	1.2.276.0.7230010.3.1.4.8323329.32579.1517875161.299312	252069 1 1021 3 1020 4 1018 5 1018 6 1016 7 1015 8 1014 9 1014 9 1013 10 1013 10 1012 11 1011 12 1010 12 1011 12 1010 13 1009 14 1008 15 1008 15 1007 16 1007 16 1006 17 1006 17 1005 17 1005 18 10...
...
11577	1.2.276.0.7230010.3.1.4.8323329.4461.1517875182.971843	592067 6 1016 10 1012 14 1007 18 1004 20 1003 22 1000 24 999 25 998 26 996 27 996 28 995 28 994 29 994 30 993 30 992 31 992 31 992 32 992 31 992 31 992 31 992 31 991 32 991 32 991 32 992 31 992 3...
11578	1.2.276.0.7230010.3.1.4.8323329.4461.1517875182.971843	610576 3 1001 38 981 53 966 63 956 73 947 87 932 97 922 107 912 117 902 127 893 135 888 146 877 152 873 156 868 161 864 165 861 167 859 167 859 168 858 170 857 169 857 170 856 169 857 168 859 167...
11579	1.2.276.0.7230010.3.1.4.8323329.32730.1517875162.25023	-1
11580	1.2.276.0.7230010.3.1.4.8323329.13252.1517875244.359912	-1
11581	1.2.276.0.7230010.3.1.4.8323329.12050.1517875237.113402	-1
11582 rows × 2 columns		

Figure 2. RLE file data for five images.

2.2 Resulting Dataset

- **Columns in the final dataset:**
 - **dicom_path:** Full path to the DICOM file
 - **PatientSex:** Gender of the patient (Male/Female)
 - **PatientAge:** Age of the patient
 - **PneumothArea:** Area affected by pneumothorax (calculated from RLE)
 - **Healthy:** Indicates if the patient is healthy or has pneumothorax
 - **EncodedPixels:** Extracted Features
- **Description:** This dataset is now enriched with important metadata and is ready for further analysis, such as classification or segmentation tasks.

				ImageId	EncodedPixels	Healthy	Age	PneumothArea	Sex	HasFile		dicom_path
19	1.2.276.0.7230010.3.1.4.8323329.3604.1517875178.653360	0.3			218560 13 16 1 987 41	False	51	592.018999	M	True	dicom_files\1.2.276.0.7230010.3.1.4.8323329.3604.1517875178.653360.dcm	
					979 49 970 64 956 74 949							
					76 947 77 945 80 942 83							
					938 87 927 98 920 104							
					900 125 896 129 891 134							
20	1.2.276.0.7230010.3.1.4.8323329.3604.1517875178.653360	0.3			884 141 873 164 854 170	False	51	65.579943	M	True	dicom_files\1.2.276.0.7230010.3.1.4.8323329.3604.1517875178.653360.dcm	
					845 179 837 187 830 195							
					826 198 824 200 822 202							
					820 20...							
					184910 8 1013 13 1009 15							
52	1.2.276.0.7230010.3.1.4.8323329.335.1517875162.443624	0.3			1007 17 1006 17 1004 19	True	64	0.000000	F	True	dicom_files\1.2.276.0.7230010.3.1.4.8323329.335.1517875162.443624.dcm	
					1003 21 1002 21 1001 23							
					1000 23 999 24 999 23							
					999 23 1000 23 999 23							
					998 25 998 25 998 26 998							
54	1.2.276.0.7230010.3.1.4.8323329.368.1517875162.584967	0.3			25 998 24 999 24 999 24	False	48	154.165011	F	True	dicom_files\1.2.276.0.7230010.3.1.4.8323329.368.1517875162.584967.dcm	
					999 24 999 24 1000 23							
					1000 23...							
					205015 2 1021 3 1019 5							
					1017 6 1016 8 1015 8							
					1014 9 1014 9 1014 8							
					1015 8 1015 8 1015 9							
					1014 9 1014 9 1014 9							
					1013 10 1013 10 1013 10							
					1013 10 1013 10 1013 10							
					1013 10 1013 10 1013 10							
					1013 10 1013 10 1013 10							
					1013 10 1013 10 1013 10							
					1013 10 1013 10 1013 10							
					1013 10 1013 10 1013 10							

Figure 3. Snapshot of Preprocessed Dataset.

2.3 Challenges and Solutions

- **Challenge 1: Missing DICOM Files**
 - **Solution:** Rows with missing files were filtered out.
- **Challenge 2: Inconsistent Metadata**
 - **Solution:** Applied default values and filtered out rows with invalid data.
- **Challenge 3: Duplicate Entries**
 - **Solution:** Duplicate entries based on **ImageId** were removed.

2.4 Conclusion

Summary: The DICOM dataset has been cleaned, enriched with metadata, and is now ready for machine learning or medical analysis.

Next Steps:

- Train machine learning models on the dataset.
- Analyze the pneumothorax areas for clinical insights.
- Explore further medical imaging techniques using the data.

3. Feature Extraction Using ResNet50 for DICOM Images

To extract deep feature embeddings from DICOM (Digital Imaging and Communications in Medicine) images using a pre-trained ResNet50 model architecture consists of 50 layers. The architecture of ResNet50 is divided into four main parts: the **convolutional layers**, the **identity block**, the **convolutional block**, and the **fully connected layers**. The convolutional layers are responsible for extracting features from the input image, the identity block and convolutional block process and transform these features, and the fully connected layers make the final classification.

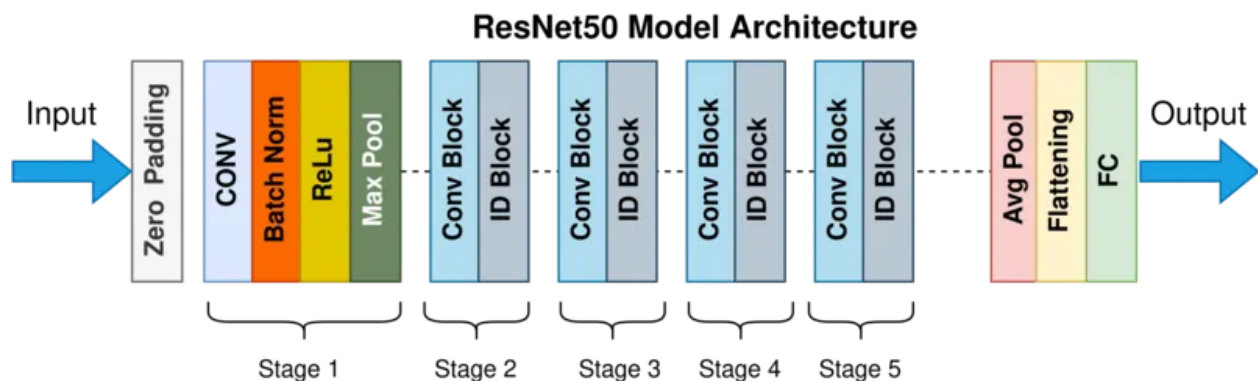


Figure 4. RestNet50 Architecture.

These feature embeddings can later be used for **classification tasks, clustering**, or further analysis.

- Used **ResNet50** for **deep feature extraction**.
- **Preprocessed grayscale DICOM images** to match the required RGB input format.
- **Extracted feature embeddings** for all the images.
- These features can be **used for classification, clustering, or further analysis** in machine learning tasks.

Steps:

Loading Pre-trained Model:

- The ResNet50 model is loaded without the top classification layer (**include_top=False**), using the **imagenet** weights.

Resizing and Preprocessing:

- Images are resized to 224x224 pixels (required input size for ResNet50).
- Images are converted from grayscale to RGB, and preprocessed for ResNet50 input.

Extracting Embeddings:

- Instead of using the final **fully connected (FC) layer**, we extract features from the **last convolutional layer (before the FC layer)**, which provides a **2048-dimensional feature vector** for each image.

4. Dataset with Embeddings

Embedding DataFrame:

- The embeddings are stored in a DataFrame with additional columns for patient features (age, sex, etc.).

Missing Embeddings:

- Rows without valid embeddings are removed from the DataFrame to ensure the dataset is clean.

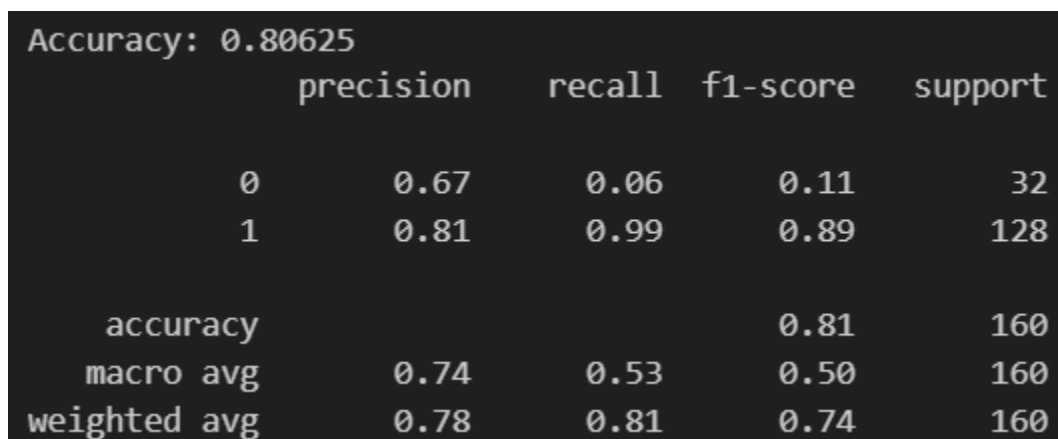
5. Machine Learning Model

Objective: Use the extracted embeddings to train a classification model to predict "Healthy" vs. "Unhealthy" status..

Steps:

1. **Data Splitting:**
 - Randomly splits the dataset into training (80%) and testing (20%) sets.
2. **Training the Model:**
 - A Random Forest Classifier is trained on the features (embeddings, age, sex) to predict the target variable **Healthy**.
3. **Model Evaluation:**
 - Accuracy and a classification report are generated to assess the model's performance.

5.1 Title: Training a Random Forest Classifier



Accuracy: 0.80625					
	precision	recall	f1-score	support	
0	0.67	0.06	0.11	32	
1	0.81	0.99	0.89	128	
accuracy			0.81	160	
macro avg	0.74	0.53	0.50	160	
weighted avg	0.78	0.81	0.74	160	

Figure 5. Snapshot of Random Forest Classification report.

Insights:

- **Overall Accuracy:** 80.63% of test samples were correctly classified.
- **Class 0 (Unhealthy) Performance:**
 - Precision: **67%** → When the model predicts "Unhealthy," it's correct 67% of the time.
 - Recall: **6%** → The model only detects **6%** of actual unhealthy cases.
 - F1-score: **11%** → Poor performance for detecting unhealthy cases.
- **Class 1 (Healthy) Performance:**
 - Precision: **81%** → When the model predicts "Healthy," it's correct 81% of the time.
 - Recall: **99%** → The model identifies almost all healthy cases.

- F1-score: **89%** → Strong detection of healthy cases.
- **Imbalance Issue:**
 - The model is heavily biased towards classifying most patients as "Healthy" (Class 1). It performs poorly in detecting "Unhealthy" (Class 0).

5.2 Title: XGBoost Classifier Output

Accuracy: 0.7875					
	precision	recall	f1-score	support	
0	0.42	0.16	0.23	32	
1	0.82	0.95	0.88	128	
accuracy			0.79	160	
macro avg	0.62	0.55	0.55	160	
weighted avg	0.74	0.79	0.75	160	

Figure 6. Snapshot of XGBoost Classification report.

Insights:

Overall Accuracy: 78.75% of test samples were correctly classified.

Class 0 (Unhealthy) Performance:

- Precision: **42%** → When predicting "Unhealthy," it's correct 42% of the time.
- Recall: **16%** → Only 16% of actual unhealthy cases were detected.
- F1-score: **23%** → Still poor at detecting unhealthy cases, but slightly better than Random Forest.

Class 1 (Healthy) Performance:

- Precision: **82%** → When predicting "Healthy," it's correct 82% of the time.
- Recall: **95%** → The model correctly identifies 95% of healthy cases.
- F1-score: **88%** → Strong detection of healthy cases.

Imbalance Issue:

- The model still favors "Healthy" cases, but slightly improves in identifying "Unhealthy" compared to Random Forest.

Title: Comparison of Random Forest and XGBoost Models

Objective: Compare the performance of two machine learning models (Random Forest and XGBoost) in classifying healthy and unhealthy cases based on extracted embeddings.

Model Performance Comparison:		
Metric	Random Forest	XGBoost
Accuracy	80.63%	78.75%
Precision (Healthy)	0.81	0.82
Precision (Unhealthy)	0.67	0.42
Recall (Healthy)	0.99	0.95
Recall (Unhealthy)	0.06	0.16
F1-Score (Healthy)	0.89	0.88
F1-Score (Unhealthy)	0.11	0.23

Figure 7. Snapshot of Comparison of Random Forest and XGBoost Models

Observations

Random Forest:

- Higher overall accuracy (**80.63%**).
- Better performance in detecting unhealthy cases compared to XGBoost.
- The model is more balanced in terms of recall across both classes.

XGBoost:

- Slightly lower accuracy (78.75%).
- Higher precision for the "Healthy" class but poorer precision for the "Unhealthy" class.
- The model is biased towards detecting "Healthy" cases (higher recall), missing many "Unhealthy" cases.

3. Clustering with K-Means

Title : Dataset Creation

Steps

- Extract embeddings for all DICOM images.
- Store embeddings in a DataFrame.
- Add additional patient features (**Age, Sex**).
- Remove missing embeddings.

Process

1. Normalize extracted **feature embeddings**.
2. Apply **K-Means clustering** to group similar images.
3. Select **optimal number of clusters (K)**.
4. **Visualize cluster distributions.**

Title: Cluster Visualization

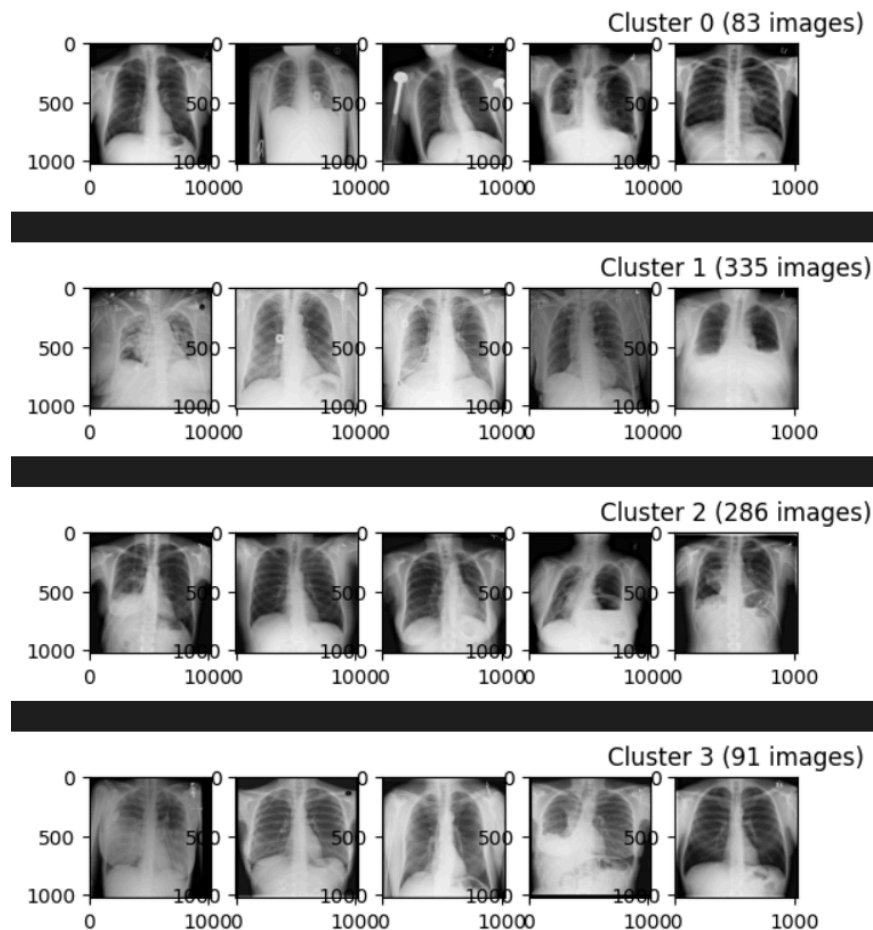
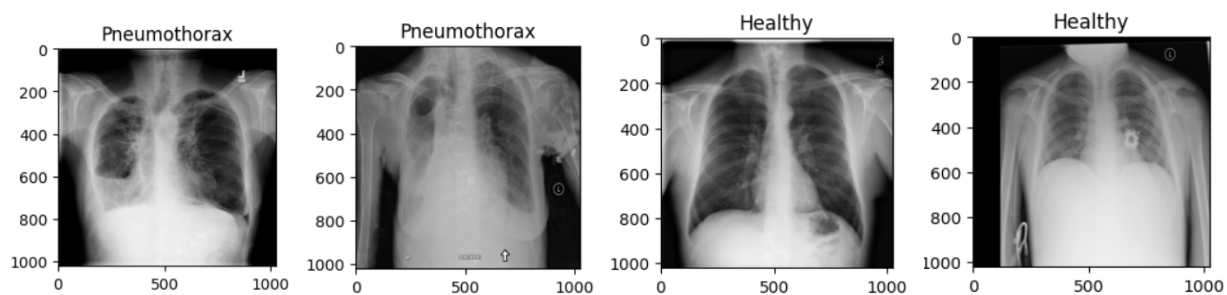
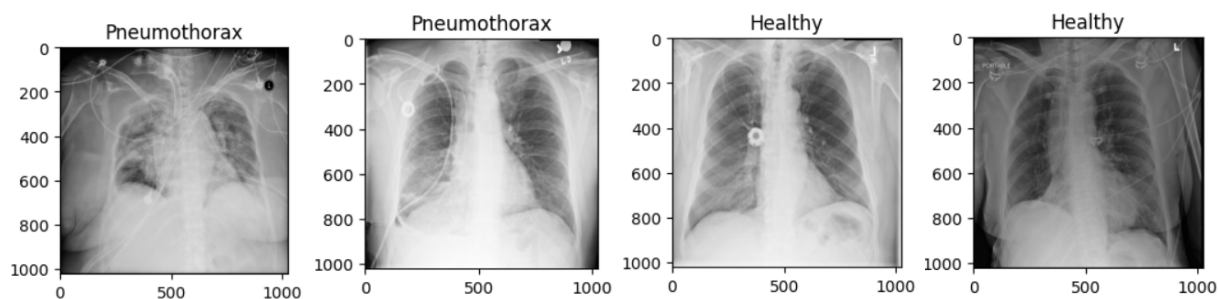


Figure 8. Snapshot of Cluster Visualization

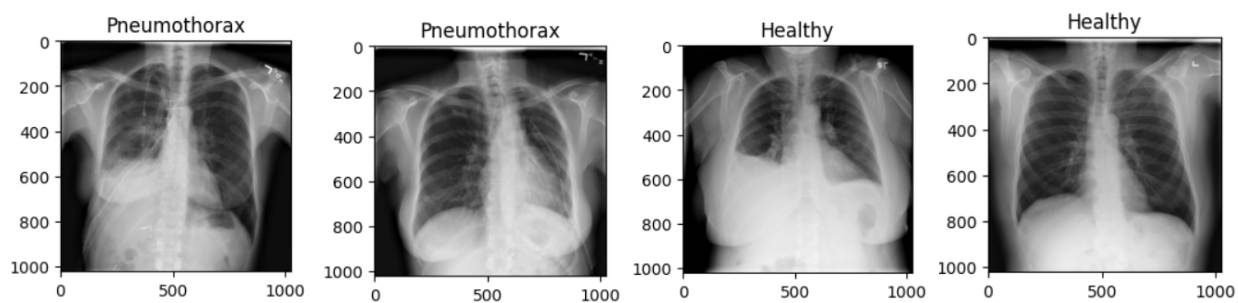
CLUSTER 0



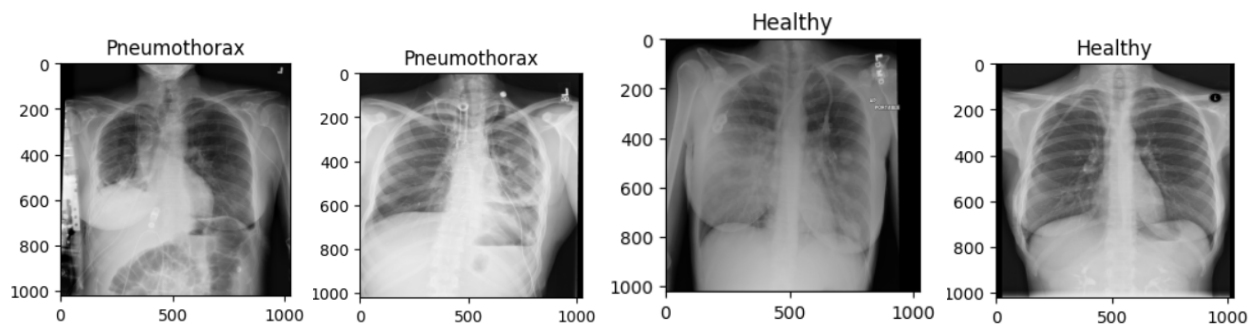
CLUSTER 1



CLUSTER 2



CLUSTER 3



CLUSTER 4

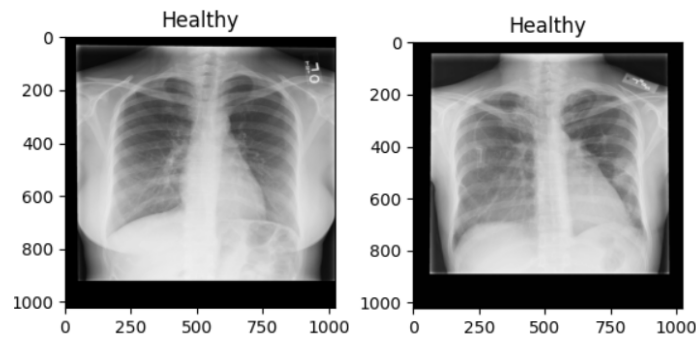


Figure 9. Snapshot of Five Clusters

Title: Cluster Distribution

- K-Means clustering on extracted image features and analyzes how the clusters are distributed among healthy vs. pneumothorax cases.
- **Cluster 1 and Cluster 2** contain the highest number of images, with a mix of **healthy and pneumothorax** cases.
- **Cluster 4** contains only **healthy cases** (no pneumothorax).
- **Cluster 0 and Cluster 3** contain a small number of pneumothorax cases.
- Some clusters have a **higher proportion of pneumothorax** cases (e.g., Clusters 1 & 2), while others are mostly **healthy**.
- This suggests that **K-Means may have grouped images based on visual similarities**, which could be useful for an **automated diagnostic system**.

Applies **K-Means clustering** to group images based on extracted features and then analyzes the **distribution of pneumothorax vs. healthy cases** in each cluster.

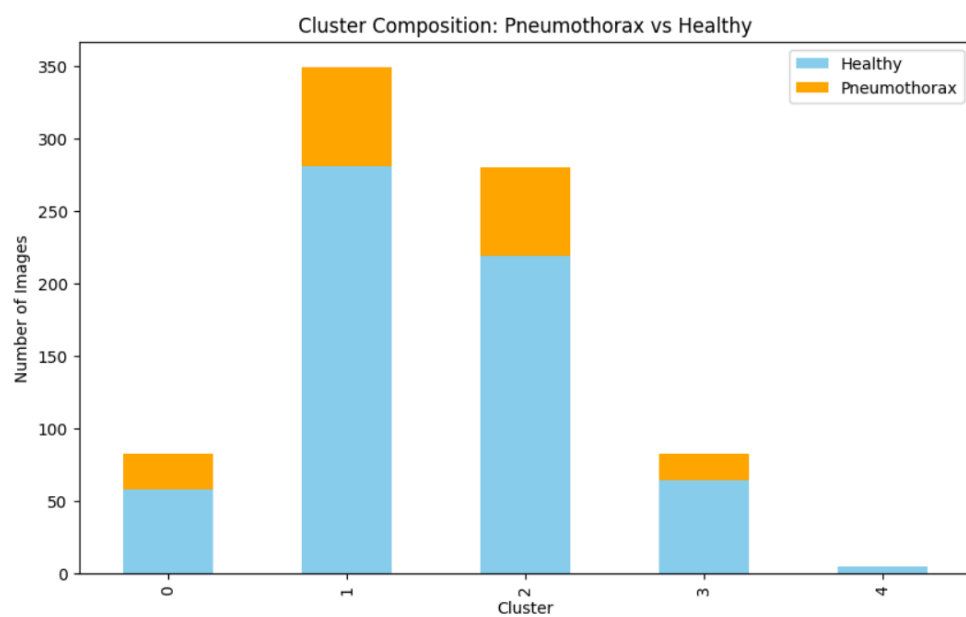


Figure 10. Snapshot of Bat Chart Cluster Distribution

Cluster ID	Healthy Images	Pneumothorax Images
0	58	24
1	281	68
2	219	61
3	64	18
4	5	0

Figure 11. Snapshot of Cluster Distribution

Conclusions and Future Scope

Deep learning algorithms have significantly enhanced the ability of machines to interpret medical images, revolutionizing AI-based disease diagnosis and prognosis. In this study, we utilized the **ResNet50** model for **feature extraction** and applied **Machine Learning** techniques for **automatic classification** of pneumothorax in chest X-ray images. Additionally, **K-Means clustering** was employed to group similar images based on extracted features.

Among the classification models, **Random Forest** demonstrated a better balance between **precision and recall**, while **XGBoost** showed a bias towards detecting the **Healthy** class and struggled with accurately identifying **Unhealthy** cases due to lower recall and precision. The clustering analysis successfully grouped images, allowing for further insights into their distribution across clusters.

Future Work:

- Explore alternative clustering techniques such as **DBSCAN** and **Agglomerative Clustering**.
- Utilize clustered data for further **classification tasks and analysis**.
- Compare feature extraction performance across different models, including **VGG16** and **InceptionV3**, to evaluate their effectiveness against **ResNet50**.