



Using Image Provenance Network for Detection and Classification of Lung Diseases

Anoop Raju¹, Asha George², P. Mohammed Shameem³

¹PG Scholar, Department of Computer Science and Engineering, TKMIT anooppadippurayil@gmail.com

² Assistant Professor, Department of Computer Science and Engineering, TKMIT ashatec@gmail.com

³ Professor, Department of Computer Science and Engineering, TKMIT pms.tkmit@yahoo.in

ABSTRACT

Medical image analysis has been one of the most sought after areas by image processing researchers for long period due to its scope and relevance both in application level and in the effectiveness of the techniques used. Most of the techniques that proved efficient and effective in the twentieth century have turned to be obsolete now because of the enormous volume of data collected through various biomedical equipments. The intensive use of web enabled services for the Image provenance analysis is a recently emerged area after the widespread use of web enabled services for applications as well as development. Analysis of images based on the information over a time period can guide the investigations on disease identification and treatment in new dimensions. We propose an efficient mechanism for constructing a provenance network of images making use of prior information over various parameters influencing the content and connectivity of images. Experiments conducted over various data sets prove that the proposed method successfully generates guided provenance network with minimum false positives. This paper also proposes a tool for analyzing medical images for lung disease analysis. A provenance network is constructed for a sequence of images of infected cases at different time span after performing image filtering and interest point mapping using the proposed provenance network model. Experimental results show that the proposed technique can be used in effectively analyzing and diagnosing disease levels and the effect of medication. The technique has been compared with existing image analysis techniques specifically designed for medical analysis.

Key words: medical image analysis, provenance network, artificial neural network, lung disease detection, classification.

1. INTRODUCTION

Image processing has emerged to be an area in computer science and information technology that has undergone the most extensive research and that which has significant impact in easing the practices in many faces of day-to-day life.

Medical imaging places challenging problems to researchers as it demands the most accurate results in all cases. Medical Since then images are routinely acquired for medical diagnosis. Many techniques for the interpretation of medical images like registration, segmentation, classification, and reliable measurements have been developed, with sufficient focus on the evaluation of the algorithms. Medical image processing is now the most essential diagnosis tool integrated in a physicians' workflow. High-level image processing has become part of diagnosis, intervention planning and treatment. Imaging has turned out to be an inevitable component in health care, medical research and laboratories.

Medical image processing deals with the development of problem-specific approaches to the enhancement of raw medical image data for the purposes of selective visualization as well as further analysis. There are many topics in medical image processing: some emphasize general applicable theory and some focus on specific applications. Imaging of the lung is a mainstay of respiratory medicine. It provides local information about morphology and function of the lung parenchyma that is unchallenged by other noninvasive techniques [1]. X-rays of the chest are almost always done when doctors suspect a lung or heart disorder. Other imaging tests are done as needed to provide doctors with specific information to make a diagnosis.

Image provenance analysis is an area that is not well studied in the image processing literature. Image provenance analysis has been studied as a method for the detection of manipulated content in digital images [2]. Since provenance network generally depicts the relationships between a set of images, or the path of deriving a final image from a sequence of images, analysis of the same can be used in many applications other than mere identification of manipulated content. This thought has led to this work which uses a set of lung radiographs for the analysis and prediction of lung diseases. The determination of image provenance is a difficult task to solve. The complexity increases significantly when considering an end-to-end, fully-automatic provenance pipeline that performs at scale.

Chest radiology is the most common method used for diagnosis of lung diseases, the term lung disease refers to the abnormalities that affect the lung organ, diseases are such as asthma, COPD, lung cancer, pneumonia and many other breathing problems, in this paper develop a system that detects and classify the lung diseases as either pneumonia or lung cancer, this is accomplished by two stages they are feature extraction and classification, feature extraction is done through the use of Gabor filter, classification is through the use of neural network's like Feed Forward Neural Network(FFNN), Multi-Layer Perceptron Neural Network(MLPNN), Radial Basis Function(RBF) [3].

Chest x-rays are routinely taken from the back to front. Usually a view from the side is also taken. Chest x-rays provide a good outline of the heart and major blood vessels and usually can reveal a serious disorder in the lungs, the adjacent spaces, or the chest wall, including the ribs. For example, chest x-rays can show most pneumonias, lung tumors, chronic obstructive pulmonary disease, a collapsed lung (atelectasis) and air (pneumothorax) or fluid (pleural effusion) in the pleural space [4]. Although chest x-rays seldom give enough information to determine the exact cause of the abnormality, they can help a doctor determine whether and which other tests are needed to make a diagnosis.

Chest physicians begin their evaluation of a patient by the history and physical examinations, appropriate laboratory studies, and review of imaging studies such as standard erect posteroanterior and lateral chest radiographs, computed tomographic (CT) scans of the chest, or both. In doing so, they assess the differential diagnostic probabilities and possibilities based on the radiographic patterns. This approach serves to focus their efforts on history taking and the physical examination and facilitates their defining of proper diagnostic measures, which, in turn, guide therapeutic and management advice [5]. Parenchymal lung diseases are classified into nine imaging patterns on the basis of chest radiography and CT: focal pulmonary opacities, multifocal pulmonary opacities, segmental/lobar opacities, interstitial opacities, cavitory lesions, cystic lung disease, single small nodules, large masses, and multiple nodules. These radiographic patterns may be caused by infectious diseases, neoplastic diseases, or noninfectious/nonneoplastic disorders. The differential diagnoses for the nine imaging patterns are dissimilar but by no means mutually exclusive. Many diseases that usually cause focal opacities can produce multifocal opacities. Other disorders are nearly always present as multiple opacities, and the pathology only rarely localizes to one area. Opacities that conform perfectly to the segmental anatomy of the lung usually result from an abnormality of the bronchus or pulmonary artery leading to the opacity [6]. Interstitial, cavitory opacities, cysts, single small nodules,

large masses, and multiple nodules have distinct differential diagnoses and are discussed in a separate review. In this review, a simple approach is taken to divide causes of pulmonary abnormalities into three broad categories: infectious, neoplastic, and noninfectious/nonneoplastic.

Diagnosis of a focal or multifocal lung disorder starts with the abnormal chest radiograph or with abnormal findings on a chest CT scan. In many instances, the chest CT scan may be abnormal when the chest radiograph is normal or shows very indistinct abnormalities. When confronted with an abnormal chest radiograph, care should be taken to compare it with any previous chest radiographs. Each of the differential diagnoses can be further refined when the time course of an abnormality is considered, acute versus chronic. Generally, acute abnormalities are infectious and traumatic, whereas chronic abnormalities are neoplastic and diffuse parenchymal lung disease. In each category of radiographic pattern, the clinical features of the illness, the acuity/chronicity of the disease, the presence or absence of associated pleural or mediastinal abnormalities and ancillary laboratory tests also serve to narrow the differential diagnosis [7]. In some disorders, the combined radiographic, clinical, and laboratory presentation is virtually specific. In other disorders, cytologic, histopathologic, or microbiologic information is necessary to make a specific diagnosis.

In the diagnosis of lung disease, assessing the pattern, location, and regional distribution of involvement is the domain of radiology. Complementing classical chest radiography, the tomographic imaging techniques of computed tomography (CT, today in the form of multidetector CT, MDCT) and magnetic resonance imaging (MRI) can demonstrate and quantify, not just morphology, but increasingly also functional processes such as perfusion, respiration, and metabolism, region by region. MDCT and MRI are also often used in treatment monitoring or disease monitoring to identify progression, to allow appropriate treatment alterations or further interventions to be initiated [8].

Artificial neural networks (ANN) provide a powerful tool to help doctors to analyze, model and make sense of complex clinical data across a broad range of medical applications. Most applications of artificial neural networks to medicine are classification problems; that is, the task is on the basis of the measured features to assign the patient to one of a small set of classes [9]. ANN has regularly been used as a successful technique for classification and prediction problems. This work uses convolution neural networks for the analysis of Image provenance network in analysis lung disorders.

2. RELATED WORK

2.1 Provenance Network

Image provenance deals with identifying various types of transformations that have been accomplished on a set of selected images to obtain a new one. Moreira et.al. claim their research is the first work towards developing an end-to-end fully-automatic provenance pipeline that performs at scale [10]. Their work consists an image indexing scheme that utilizes a novel iterative filtering and distributed interest point selection. The work also proposes methods for provenance graph building that improve upon the methods of previous work in the field by Punto, and provides a novel clustering algorithm for further graph improvement [11]. The proposed approaches perform decently well in connecting the correct set of images. Although image content is the most reliable source of information connecting related images, other external information may be required to supplement the knowledge obtained from pixels. This external information can be obtained from file metadata, object detectors and compression factors, whenever available.

Dias et. al. [12] discusses the problem of identifying the image relationships within a set of near-duplicate images. They use the term Image Phylogeny Tree (IPT), due to its natural analogy with biological systems. The mechanism of building IPTs aims at finding the structure of transformations and their parameters if necessary, among a near-duplicate image set, and has immediate applications in security and law-enforcement, forensics, copyright enforcement, and news tracking services [13]. The work presents a method for calculating an asymmetric dissimilarity matrix from a set of near-duplicate images and formally introduces an efficient algorithm to build IPTs from such a matrix. The approach has been validated using both synthetic and real data, and show that using an appropriate dissimilarity function we can obtain good IPT reconstruction even when some pieces of information are missing. We also evaluate our solution when there are more than one near-duplicate sets in the pool of analysis and compare to other recent related approaches in the literature. One heuristic and one optimum branching algorithm are explored in [14] for reconstructing the evolutionary tree associated with a set of image documents.

Image provenance network is a graphical representation of relationships in between a set of collected images. The construction of such a graph includes various processes for identifying the associations among the images based on their matches or mismatches. The processes may include image phylogeny, object recognition and scene recognition, near duplicate detection or similar ones. There may be some common characteristics that make the images visually related. Image phylogeny solutions aim at finding kinship

relations between different versions of an image [15]. Similar to provenance analysis, image phylogeny limits its representation to a single-root tree with the original image as the root, even though there can be multiple original images contributing towards the creation of an image. The algorithm receives a query image and outputs the Image Phylogeny Tree (IPT). That method has also been extended to handle multiple (two) roots by taking spliced images into consideration [16]. Image description for provenance analysis typically avoids using computationally expensive methods such as deep learning because of scalability concerns [17, 18].

2.2. Medical Imaging

Pavithra & Pattar has presented a a system to detect lung cancer and pneumonia based on the classification of chest X-rays. First, a Gabor feature set is extracted and then classification is done using radial basis function, multi-layer perceptron and feed forward neural network. Research has been conducted on feature extraction of lung nodule in chest x-rays. In benign tumor, area and perimeter value are large as compare to malignant tumor [2]. But malignant tumors are more irregular. In GLCM properties, the average of contrast, correlation, energy and homogeneity values are calculated in different directions. Calculated properties gives a distinguish values for malignant or benign tumor, which helps in deciding the given image is of malignant tumor or benign tumor.

Lung nodules can be detected using multiscale wavelets and support vector machines [4]. This methodology uses three different types of kernels like linear, Radial Basis Function (RBF) and polynomial, among which the RBF kernel gives better class performance.

3. PROPOSED WORK

The proposed work develops in two stages: (A) Construction of provenance network, and (B) CNN based Lung disease detection and analysis. Before passing through these procedural stages, the chest X-ray is given for preprocessing operations to remove the irrelevant data and to enhance the region of interest. The intensity and contrast of the X-ray image is adjusted to make it fit for the further processing stages.

3.1. Construction of Image Provenance Network (IPN)

The provenance graph is a directed acyclic graph [19]. Each node in IPN corresponds to an image in the set of related images and the edges stand for the relationship. The first phase in provenance graph construction is creation of adjacency matrix.

Purpose of image pre-processing is to remove irrelevant data present on chest x-ray film. Recovery of useful information; strengthening region of interest and simplification of features on chest x-ray image is done with the help of pre-processing techniques [20]. It has two main steps such as image enhancement and image filtering. Chest x-ray image intensity and contrast is adjusted to make image more suitable for the further processing stages. Histogram equalization method is used to enhance image. It is nothing but the uniform distribution of image intensity pixel. Histogram equalization method for chest x-ray image where image intensity and contrast are adjusted. Image filtering is done to remove unwanted noise present in chest x-ray image. Lung segmentation means acquiring region of interest of lungs. To detect lung diseases it is required to identify lung boundaries. For lung boundary detection we have used intensity based method and discontinuity based method (Edge Detection). Intensity based Method: This method operates on individual pixels. It converts grayscale image into binary image [21]. Here threshold T minimum value is considered as criteria. If the image intensity is goes above T value then image is appears as white image otherwise it is black.

Once the GCM- and MI-based dissimilarity matrices are available, we rely on both for constructing the final provenance graph, by the means of a novel algorithm, named clustered provenance graph expansion. This algorithm, differently from Kruskal's algorithm, delivers directed provenance graphs rather than undirected minimum spanning trees [22, 23]. The main idea behind such solution is to group the available images in a way that only near duplicates of a common image are added to the same cluster.

Assuming that we are provided the query image Q, and a set of related images to the query, denoted by R, we build an $N \times N$ (here, $N = |R| + 1$), asymmetric matrix D, in which each indexed value $D[i; j]$ is the similarity (or dissimilarity) quotient between the images i and j. The full matrix is obtained by comparing $(n^2 - n) / 2$ pairs.

Before comparison, the images can be described using interest point detectors and descriptors. Once described, the k most relevant interest points of each image are then matched using brute-force pairwise comparison based on the L2 distance between the descriptors. The best matched correspondences are filtered to retain the geometrically consistent ones, as described in [24].

The regions in the images including these points are then extracted as shared Regions of Interest (ROIs) [25]. Finally, the similarity between the pixel distributions in these regions is computed based on a metric such as mutual information.

Based on the comparative values of the adjacency matrix, the final graph construction step chooses the most feasible set of directed edges (i.e., the set of edges that best represents the sequence of image operations). Each chosen directed edge denotes a parent-child relationship in the graph. any minimum (or maximum) cost spanning tree algorithm, such as Kruskal's Minimum Spanning Tree (MST) algorithm [26], can be used to build an undirected graph connecting relevant images and post-process this graph for directionality. Alternatively, a specialized algorithm for directly constructing a directed provenance graph, such as clustered provenance graph expansion, can be employed [27].

Table 1. Feature extraction Results for Chest Radiographs

SL. No	Area	Perimeter	Irregularity index	Diameter	Mean	SD	Entropy
1	7.80E+06	4	6.13E+06	3.15E+03	0.5621	0.4961	9.89E-01
2	8.99E+06	2	2.82E+07	3.38E+03	0.5632	0.496	0.9885
3	7.35E+06	2	2.31E+07	3.06E+03	0.5463	0.4978	0.9938
4	7.95E+06	4	6.24E+06	3.18E+03	0.5636	0.4959	0.9883
5	5.93E+06	7.4142	1.36E+06	2.75E+03	0.5237	0.4994	0.9984
6	5.43E+06	4	4.27E+06	2.63E+03	0.5225	0.9495	0.9985
7	4.85E+06	2	1.52E+07	2.48E+03	0.5679	0.4954	0.9866
8	5.26E+06	2	1.65E+07	2.59E+03	0.5334	0.4989	0.9968
9	4.86E+06	10.2426	5.82E+05	2.49E+03	0.5593	0.4965	0.9838
10	4.79E+06	5.36E+04	0.0209	2.47E+03	0.562	0.4961	0.9889

3.2. CNN for lung disease detection

For identification of lung disease as TB; lung cancer and pneumonia; here we used feed-forward neural network. The network proposed consists of seven input neuron of chest radiographs. Three hidden layer represents lung diseases as lung cancer; TB and pneumonia's output layer denotes lung disease classification. Activation of neuron is done by applying sigmoid function.

4. RESULT AND DISCUSSION

Experiments have been conducted on the chest radiographs of more than 700 patients. Features such as area of lung region, its perimeter and irregularity index, equivalent diameter, mean, standard deviation and image entropy are considered for further processing step of classification using neural network. Feature Extraction results for the chest radiographs are shown in Table 1. Classification using convolution neural networks has been implemented using Python as coding language. Experimental results for different sets of lung radiographs for the lung disease pneumonia is given in Table 2. The results can be further improved with the help of back propagation algorithm.

Table 2. Prediction Result Analysis

Set A (%)		Set B (%)		Set C (%)		Combined Set A, B, C (%)	
ACC%	AUC	ACC%	AUC	ACC%	AUC	ACC%	AUC
95.57	0.99	81.06	0.90	70.40	0.77	92.00	0.97

Chest radiographs database for 800 patients have experimented with the proposed method. The database has been taken from kaggle. Image preprocessing results are as follows which goes through process such as image grayscale conversion; image filtering to remove noise; image smoothing; binarization of image and edge detection respectively. Features such as area of lung region; its perimeter and irregularity index; equivalent diameter; mean; standard deviation and image entropy are considered for further processing step of classification using neural network. Feature Extraction results for the chest radiographs are shown in table. lung segmentation and feature extraction of chest x-ray. Feed forward neural network is used where network flows in only one direction. Using neural network in python results are computed. From the prediction obtained with neural network detection of disease is possible. TB; pneumonia and lung cancer are automatically detected using the proposed image processing method.

We trained a classifier model on feature vectors to discriminate between normal and abnormal chest X-Ray

images. A unique characteristic of a CNN is that it classifies by computing the hyper plane with the largest margin between two classes; i.e., the hyper plane with the largest distance to the nearest training data point of any class. Ideally, the feature vectors of abnormal CXRs will have a positive distance to the separating hyper plane, and feature vectors of normal CXRs will have a negative distance. The sign of the distance corresponds to the "side" of the hyper plane that the feature vector lies in hyperspace. The larger the distance the more confident we are of the class label. In our model we use this confidence measure to provide feedback to the operator. This enables confidence-based thresholding to find the optimal operating point of the classifier and the computation of the ROC curve. The optimal operating point of the classifier is the point on the ROC curve that provides the best ratio of sensitivity and specificity, given a cost function that describes the operational costs of missing an abnormal case in the field.

5. CONCLUSION

Image phylogeny has been studied for use in image forensics only. We have made an attempt to use image provenance analysis based on phylogeny for a positive purpose, to identify lung diseases from chest radiographs. Classification of radiographs of infected cases is done using convolution neural networks. Further investigations can be conducted to diagnose how intense the disease is being from details inspection of the provenance network. This approach has potential for further development because of this simplicity that will motivate to classify the types of lung diseases. The developed classification system is expected to provide valuable diagnosis for the physicians. The work can further be extended by including more feature extraction and selection methods for classifying more lung diseases. The limitation of this proposed method is that it is not robust when there are changes in the size and position of chest x-ray image. This limitation can be overcome by taking images over a similar conditions and digitalizing chest x-ray images.

REFERENCES

1. Putzu, Lorenzo, Giovanni Caocci, and Cecilia Di Ruberto. "Leucocyte classification for leukaemia detection using image processing techniques." *Artificial intelligence in medicine* 62.3 (2014): 179-191.
2. Berkowitz, Eugene A., Adam Bernheim, and Gerald W. Staton Jr. "Imaging Of Lung Disease, Part I: Focal And Diffuse Parenchymal Lung Diseases." *Scientific American* 2 (2018): 18.
3. Hussein, Sarfaraz, et al. "Tumornet: Lung nodule characterization using multi-view convolutional neural network with gaussian process." 2017 IEEE 14th International Symposium on Biomedical Imaging

- (ISBI 2017). IEEE, 2017.
4. L. R. Folio; "**Chest Imaging: An Algorithmic Approach to Learning**", New York: Springer;2012.
 5. Pavithra R and S.Y. Pattar, "**Detection and Classification of lung disease-pneumonia and lung cancer in chest radiology using artificial neural network**", International Journal of Scientific and Research Publications; Volume 5; Issue 10; October 2015.
 6. WHO Disease and injury country estimates; World Health Organization, 2009.
 7. A. Rocha, W. Scheirer, T. Boult, and S. Goldenstein, "**Vision of the unseen: Current trends and challenges in digital image and video forensics**," ACM Comput. Surv., vol. 43, no. 4, 2011, Art. no. 26.
 8. Z. Dias, A. Rocha, and S. Goldenstein, "**Image phylogeny by minimal spanning trees**," IEEE Trans. Inf. Forensics Security, vol. 7, no. 2, pp. 774–788, Apr. 2012.
 9. Z. Dias, A. Rocha, and S. Goldenstein, "**Exploring heuristic and optimum branching algorithms for image phylogeny**," J. Vis. Commun. Image Represent., vol. 24, no. 7, pp. 1124–1134, 2013.
 10. A. Pinto et al., "**Provenance filtering for multimedia phylogeny**," in Proc. IEEE Int. Conf. Image Process. (ICIP), Sep. 2017, pp. 1502–1506.
 11. S. Lowry et al., "**Visual place recognition: A survey**," IEEE Trans. Robot., vol. 32, no. 1, pp. 1–19, Feb. 2016.
 12. L. ˇCehovin, A. Leonardis, and M. Kristan, "**Visual object tracking performance measures revisited**," IEEE Trans. Image Process., vol. 25, no. 3, pp. 1261–1274, Mar. 2016.
 13. M. A. Oikawa, Z. Dias, A. de Rezende Rocha, and S. Goldenstein, "**Manifold learning and spectral clustering for image phylogeny forests**," IEEE Trans. Inf. Forensics Security, vol. 11, no. 1, pp. 5–18, Jan. 2016.
 14. L. Nanni and A. Lumini, "**Heterogeneous bag-of-features for object/scene recognition**," Appl. Soft. Comput. vol. 13, no. 4, pp. 2171–2178, 2013.
 15. A. Joly, O. Buisson, and C. Frélicot, "**Content-based copy retrieval using distortion-based probabilistic similarity search**," IEEE Trans. Multimedia, vol. 9, no. 2, pp. 293–306, Feb. 2007.
 16. H. Huang, W. Guo, and Y. Zhang, "**Detection of copy-move forgery in digital images using SIFT algorithm**," in Proc. IEEE Pacific-Asia Workshop Comput. Intell. Ind. Appl., 2008, pp. 272–276.
 17. E. Silva, T. Carvalho, A. Ferreira, and A. Rocha, "**Going deeper into copy-move forgery detection: Exploring image telltales via multi-scale analysis and voting processes**," J. Vis. Commun. Image Represent., vol. 29, pp. 16–32, May 2015.
 18. M. Milford et al., "**Condition-invariant, top-down visual place recognition**," in Proc. IEEE Int. Conf. Robot. Automat. (ICRA), May/Jun. 2014, pp. 5571–5577.
 19. F. Costa, A. Oliveira, P. Ferrara, Z. Dias, S. Goldenstein, and A. Rocha, "**New dissimilarity measures for image phylogeny reconstruction**," Pattern Anal. Appl., vol. 20, no. 4, pp. 1289–1305, 2017.
 20. T. Ge, K. He, Q. Ke, and J. Sun, "**Optimized product quantization for approximate nearest neighbor search**," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Jun. 2013, pp. 2946–2953.
 21. B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "**Places: A 10 million image database for scene recognition**," IEEE Trans. Pattern Anal. Mach. Intell., vol. 40, no. 6, pp. 1452–1464, Jun. 2018.
 22. F. de O. Costa, M. A. Oikawa, Z. Dias, S. Goldenstein, and A. R. de Rocha, "**Image phylogeny forests reconstruction**," IEEE Trans. Inf. Forensics Security, vol. 9, no. 10, pp. 1533–1546, Oct. 2014.
 23. A. Bharati et al., "**U-phylogeny: Undirected provenance graph construction in the wild**," in Proc. IEEE Int. Conf. Image Process. (ICIP), Sep. 2017, pp. 1517–1521.
 24. L. Kennedy and S.-F. Chang, "**Internet image archaeology: Automatically tracing the manipulation history of photographs on the Web**," in Proc. ACM Int. Conf. Multimedia, 2008, pp. 349–358.
 25. Z. Dias, A. Rocha, and S. Goldenstein, "**Video phylogeny: Recovering near-duplicate video relationships**," in Proc. IEEE Int. Workshop Inf. Forensics Secur. (WIFS), Nov./Dec. 2011, pp. 1–6.
 26. Z. Dias, S. Goldenstein, and A. Rocha, "**Large-scale image phylogeny: Tracing image ancestral relationships**," IEEE Multimedia, vol. 20, no. 3, pp. 58–70, Jul./Sep. 2013.
 27. Z. Dias, S. Goldenstein, and A. Rocha, "**Toward image phylogeny forests: Automatically recovering semantically similar image relationships**," Forensic Sci. Int., vol. 231, nos. 1–3, pp. 178–189, 2013.