

Computer Aided Detection for Pneumoconiosis Screening on Digital Chest Radiographs

Horace Xu¹, Xiaodong Tao², Ramasubramanian Sundararajan³,
Weizhong Yan², Pavan Annangi³, Xiwen Sun⁴, Ling Mao⁴

¹ GE Global Research Shanghai, China

² GE Global Research Niskayuna, USA

³ GE Global Research Bangalore, India

⁴ Shanghai Pulmonary Hospital, China

Abstract: This paper presents a computer aided detection scheme on digital chest radiographs for pneumoconiosis screening. The scheme involves several medical image processing and analysis technologies, i.e. lung segmentation algorithm using the active shape model, image enhancement and features extraction from lung regions, feature down-selection by correlation analysis and clustering method, classification on lung regions by multi-scale classifiers support vector machine, and finally, chest based reporting out on the classification probabilities. Experiments are conducted on our digital chest x-ray database, and the result shows a good classification performance for screening application in clinic.

Keywords: Computer aided detection (CAD), active shape model (ASM), support vector machine (SVM), pneumoconiosis, digital radiography (DR).

1 Introduction

Pneumoconiosis is a lung disease caused by a long-term inhalation of operative dust, such as coal, asbestos and silica, and the local tissue reaction to the accumulated dust particles [1]. Chest radiography is the most practical tool according to the guidelines of China and International Labor Organizer (ILO) standard for pneumoconiosis diagnosis, which is based on the assessment of opacities on chest X-Ray (CXR) films and comparison with standard films. The opacities shape can be round or irregular, and the number of opacities in each lung zones can be described by a discrete score from 0/-, 0/0, 0/1, 1/0, 1/1, 1/2, 2/1, 2/2, 2/3, 3/2, 3/3, and 3/+, corresponding to twelve possible categories, which is named profusion level. Based on the perfusion level of the opacities on a patient's CXR, a stage between 0 and III is assigned to the patient to indicate the severity of the disease [2]. Figure 1 shows the X-ray images of normal chest and pneumoconiosis chests in different stage.

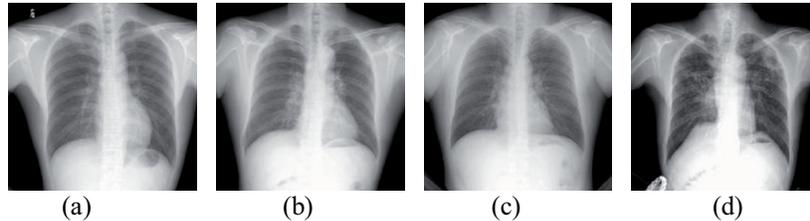


Fig. 1. Normal and pneumoconiosis chest X-ray images. (a) normal, (b) stage I, (c) stage II, (d) stage III.

In China, pneumoconiosis holds $\sim 90\%$ occupational diseases. The number of prevalence cases of pneumoconiosis is 638,234, up to the end of 2008, while there are over 610,000 suspicious cases, reported by Ministry of Health of the Peoples Republic of China [3]. However, the diagnosis of pneumoconiosis, which is based on the assessment of subtle opacities on chest X-rays films and comparison with standard films, is a challenge job in clinic. The diagnosis process is not only time consuming, subjective in staging, high experience depended, but also big variant in intra-observer and inter-observer. The radiologists' study in China shows that the intra-observer variance is $\sim 12.5\%$, and inter-observer variance is $\sim 22.5\%$ [1], which will cause different compensation to patients as well as the treatment.

In this paper, we present a computer aided detection (CAD) tool to detect the pneumoconiosis in digital Posterior-Anterior (PA) chest radiographs by image analysis automatically, which follows the process of lung segmentation, feature extraction in lung fields and sub-zones, classification and chest based report out. The CAD tool will help the radiologists in screening application as the second view, which will improve the work flow in mass screening, reduce the work load in x-ray images reading, and keep a consistent accuracy in pneumoconiosis diagnosis.

2 Lung Segmentation and Sub-division

Segmenting the lungs from digital chest x-ray image is to get the lung fields for the following automatic analysis. Also, the automatic segmentation of lung fields from chest radiographs can be a useful tool for the morphology analysis based detection of abnormalities.

2.1 Active Shape Model Segmentation

Lung field segmentation on chest radiographs has attracted a number of researchers. Rule based schemes have been investigated by Ginneken and Romeny [4]. Ginneken [5] proposed a new hybrid segmentation scheme that combined the strengths of a rule-based approach and of a pixel classification approach. Vittitoe [6] developed a pixel classifier for the identification of lung regions using Markov random field modeling. Ginneken et al.[7] put forth a lung field algorithm based on ASM method with optimal features. Iglesias [8] investigated lung segmentation based on the detection of oriented edges and active contours models. As the efficiency and

accuracy are required in the segmentation process, ASM algorithm is employed to delineate the lung fields in this paper.

The ASM algorithm could be mainly divided into two steps, training and segmentation.

In the training step, a procedure is to be followed to build up the model, i.e. labeling lung contour landmarks on training images; aligning manual landmark shapes by scaling, rotating and translating; calculating the statistics of a set of aligned shapes; constructing the shape model.

The second step is to delineate the lung contour for an input image. The procedure includes initializing the process with the average shape obtained from training phase; searching the optimal point along the direction perpendicular to the contour for each landmark in the initial contour; finding the best fitted lung field shape and update the parameters and the shape model; repetition until the process converges.

Here, we extract a seed region of lung by histogram equalization and Ostu thresholding, and transform the average shape model with scaling and translating before searching the lung contour. The initial location of lung and the transfer of the average shape improve the accuracy of ASM in lung segmentation. Figure 2 shows the result of lung field segmentation with ASM algorithm.

Based on the segmentation result, we calculate couple of parameters as morphology analysis, i.e. the symmetry of left/right lung fields by comparing their areas, the costophrenic angle which illuminates pleural effusion, and the cardiothoracic ratio to detect cardiomegaly.

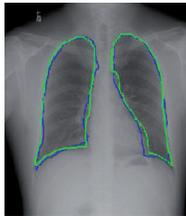


Fig. 2. The result of lung field segmentation using ASM algorithm. (Green contour is the ground truth.)

2.2 Subdividing the Lung Fields into Six Zones

According to the domain knowledge of pneumoconiosis diagnosis, the texture features are extracted from multiple regions to better capture the features of the diffuse abnormalities. On the basis of ILO guidance and radiologists' diagnostic measurement, six non-overlapping regions are subdivided in the lung fields, as shown in Figure 3. Each region will be analyzed with a separate classifier that is built by the features extracted solely from this region. The goal is able to capture the feature of abnormal variations in that particular part; in principle, this will reduce the influence of superimposed anatomical structures, which tend to make the abnormalities difficult to distinguish.

Given a lung field mask image, each lung field is subdivided into three zones - the upper, middle and lower zone. Right lung field subdivision is done by calculating the apex and diaphragm position and computing horizontal lines that divide the field into three parts of equal height. However, as the large variability of left lung shapes leads to difficulty in obtaining diaphragm position, the left lung field subdivision differs from the right lung field. Left lung field is subdivided by calculating the apex and bottom position and computing horizontal lines that divide the field into three parts of equal height. The outer point of the diaphragm is the bottom outer point of each lung field, and the inner point is found by tracing inward along the lung contour till there is a maximum change in curvature.

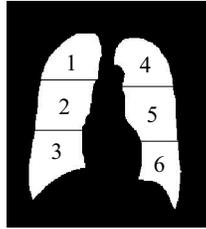


Fig. 3. The lung fields are subdivided into six regions.

3 Feature Extraction and Selection

Once a region of interest (ROI) – each of 6 zones and two lung fields – is segmented, the next task is to characterize it in terms of a set of features. Here, we are following two procedures to achieve the goal. First, we extract three types of features from each ROI, totally 247 features. Second, we select a set of features for the following classification by correlation analysis and clustering method.

3.1 Feature Extraction

1. Image enhancement by difference filtering

To enhance various size and various degree of opacities in pneumoconiosis chest, a multi-scale difference filter bank $L_n^\theta(x, y, d)$ is developed to improve the image contrast, where n , θ , and d denote the order of difference, the orientation in which the difference is computed and the difference scale, respectively.

2. Image intensity based feature

We extract a set of 6 features based on the histogram of intensity values – mean, standard deviation, skewness, kurtosis, energy and entropy.

Apart from calculating these on the original ROI, we also extract these features after applying difference filters on the image, where, two scales {1,2}, nine orientations $\{0^\circ, 30^\circ, 45^\circ, 60^\circ, 90^\circ, 120^\circ, 135^\circ, 150^\circ, 180^\circ\}$ and the first and second order are used.

3. Co-occurrence matrix based feature

We extract a set of 5 features based on the gray level co-occurrence matrix computed for the ROI, namely energy, entropy, local homogeneity, correlation and inertia. The co-occurrence matrix allows us to capture the level of similarity and dissimilarity among adjacent pixels in an ROI. Thus, an ROI with opacity will contain adjacent pixels with similarly high intensities, whereas a normal ROI will not contain such adjacent pixels. Computing these features for various orientations captures this information for various types of adjacency. In this paper, we use 4 orientations $\{0^\circ, 45^\circ, 90^\circ, 135^\circ\}$ to construct the co-occurrence matrix.

4. Frequency domain based feature

We also extract a set of 5 features based on a frequency domain transformation of the ROI. These features are the mean, root mean square (RMS), first moment of power spectrum (M1), horizontal filter and vertical filter.

In total, there are 247 features extracted as described above. Given the number of data points used for model building, this number is too large and needs to be pruned. We shall discuss the methods used for feature down-selection in the following subsection.

3.2 Feature Selection

We select a subset of features for classification by analyzing similarities (such as correlation) and eliminating redundancy. The objective is to retain the features that achieve maximum compression of the data with minimum information loss. However, the real measure of the selection method is simply whether the set of features provides a good enough description of the ROIs to allow for a good model to be built to classify ROI into twelve categories of profusion level, which is described in section 1. Therefore, we build classifiers with feature subsets of various sizes, and choose the smallest feature set where the classifier performance is acceptable.

Consider a data set with m rows (each representing an ROI) and n columns (each representing a feature that describes the ROI). We can now create an $n \times n$ matrix of values, where the value in cell (i, j) represents the similarity (absolute value of correlation) between feature i and feature j in the dataset.

If we consider each feature as a "point" in some space and the appropriate entry in this matrix as the similarity between pairs of points, we can cluster these points to arrive at groups of similar features. We perform clustering of features using a method known as hierarchical clustering [9]. First, features are grouped into a dendrogram, which is a tree where elements close to each other are clubbed together into branches. Depending upon the criterion provided by the user (number of clusters, distance between clusters etc), one can split the tree into branches and thus arrive at clusters. We can use this method to arrive at groups of similar features from the given dataset. Once these groups are given, we choose one feature from each group as a representative.

4 Classification

From machine learning point of view, pneumoconiosis detection is a binary classification problem, where classification is to determine a subject has or not has pneumoconiosis. This is the purpose of screening for workers under high risk of dust. On the other hand, pneumoconiosis staging, to determine the stage (or severity) of a patient's pneumoconiosis, is a multi-class classification problem. This can be regarded as the activity for quantitative diagnosis. This stage information coming out of pneumoconiosis staging is crucial because the determined stage is used as both legal evidence for workman's accident compensation and the means for selecting appropriate treatments.

4.1 Classification for Pneumoconiosis Screening

Normally, the task of building a classifier based on the available data can be viewed in terms of two major sub-tasks: finding the optimal level of complexity for the classifier, and building the classifier at the specified level of complexity.

For screening application, we build up 8 classifiers for left/right lungs, and 6 zones by sub-division. The methodology used for building the models is described in [10], and is outlined here:

1. Search through the space of possible parameter values that describe classifier complexity.

2. For each possible set of parameter values, build and evaluate the performance of the classifier. We perform this evaluation using the k-fold cross validation method, whereby a dataset is divided into k parts, $k-1$ of which are used for training and the k^{th} is used for validation. This is done in k different ways, and the average performance across these validation cycles is reported as the performance of the classifier at that level of complexity.

3. Pick the parameter values where the best performance is achieved and build the final classifier using those values. The final classifier performance is evaluated using a method known as leave-one-out validation, which works the same way as k-fold validation but with k equal to the number of rows in the dataset. The leave-one-out performance is considered to be a good measure of the generalization ability of a classifier [11].

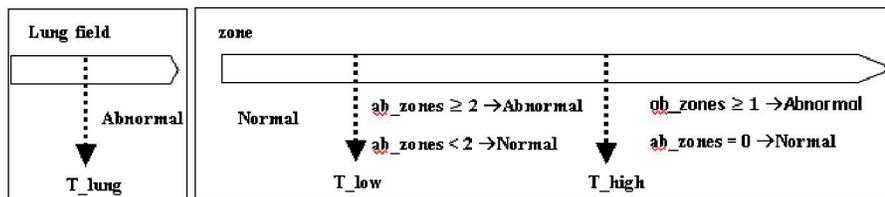


Fig. 4. The strategy of chest classification from two levels.

Once models for lung and zone level ROIs have been built, we evolve the chest based final report by leveraging the domain knowledge of pneumoconiosis screening,

for example, the classification probability of one lung ROI is larger than a threshold, the chest will be labeled as abnormal, but for zone ROIs, the final conclusion is generated from how many abnormal zones and how much the abnormal zone's probability. The strategy is shown in Figure 4. The thresholds are optimized by training in database.

For all the individual classifiers, we use Support Vector Machines (SVM), a machine technique pioneered by Vladimir Vapnik [11] and Nello Cristianini [12], and now extensively used in a number of applications. An SVM attempts to construct a hyperplane in a high dimensional space that has the largest distance (maximal margin) to the examples in the training set close to the hyperplane. Complexity is scaled by transforming the original features using an appropriate kernel function.

A key design choice to be made in building SVMs is the selection of the kernel to compute the similarity between two examples. We use the Radial Basis Function (RBF) kernel in our case, and control the complexity by varying the gamma(γ) parameter. The other parameter we control is the penalty for the error term C . These two parameters determine the complexity of the classifier being built. In order to choose the right level of complexity as defined by these two parameters, we measured the performance on unseen examples of an SVM model built at a particular complexity level using the k-fold validation technique.

4.2 Classification Study for Pneumoconiosis Staging

The pneumoconiosis detection problem has been actively studied. However, work on pneumoconiosis staging is very sparse, in spite of its importance.

Since the pneumoconiosis stages (0, I, II, III) are discrete categories ordered by disease severity, to be more exact, the pneumoconiosis staging problem is an ordinal classification problem. Ordinal classification is a special machine-learning task that falls between classification and regression [13]. Unlike nominal classification, where the target classes are un-ordered and classification performance is measured by classification accuracy, for ordinal classification both accuracy and the distance between the actual and predicted classes are relevant in evaluating classification performance. Hence, an important design strategy in developing ordinal classifiers is to leverage the ordering information to minimize the distance between actual and predicted classes.

In literature there are several approaches towards addressing ordinal classification problems. The simplest approach would be to solve ordinal problems by treating them as un-ordered nominal classification problems. This "passive" approach ignores the ordering information that could be used to improve the classification performance. More sophisticated approaches have been proposed in recent years, for example, [14-17].

In this paper, we propose using SVM for pneumoconiosis staging. In our model, the inputs are the selected features and outputs are the pneumoconiosis stages. To address the ordinal classification problem, we explore several design strategies, including categorical, ordinal, and regression.

5 Experiments and Results

In this section, we introduce the experiments by applying the algorithms described above on a database of digital chest X-Ray images comprising normal and pneumoconiosis chest x-ray images.

Normal chest images and pneumoconiosis chest images were collected for building up the database from clinic. The normal chest images are categorized by gender and age, while the pneumoconiosis chest images are categorized by stage, from 0+ to III.

A total of 427 images were collected, of which 252 were of normal subjects and 175 were of patients afflicted with pneumoconiosis.

A lung segmentation result is shown in Figure 2. Ground truth for segmentation validation was available on 400 images from this database. The performance of the algorithm was evaluated on the basis of DICE similarity coefficient (DSC) between the segmented lung mask and the ground-truth mask. DICE coefficient is computed as twice the area overlap (logical AND) of the two masks divided by the sum of the areas of the masks (logical OR). On the 400 cases tested, 89% of the images have > 85% DSC, and about 96% cases have DSC > 80%.

Of the 247 features extracted at each resolution, a smaller subset was selected for classification. At the lung level, the correlation threshold-based selection method was employed with a threshold of 80% to select 38 features. At the lung zone level, the hierarchical clustering method was used to select between 12 ~ 14 features for each zone. 2 models were built at the lung level (one for each lung) and 6 at the zone level (one for each zone). The feature list and the parameters of model are provided in Table 1 and 2 as example.

Table 1. Lung models

Lung	Left	Right
Log2 (C)	5	5
Log2 (Gamma)	-5	-7
Variables	'0_0_0_Mean'	
	'0_0_0_Skewness'	
	'0_0_0_Kurtosis'	
	'0_0_0_Entropy'	
	'2_45_1_Entropy'	
	

Table 2. Right lung zone models

Zone	Right Upper	Right Middle	Right Lower
Log2 (C)	15	1	5
Log2 (Gamma)	-9	-1	-3
Variables	2_0_2_Mean	0_0_0_Energy	1_180_2_Mean

	2_150_1_Energy	Deg90_Energy	2_150_1_Energy
	1_180_2_Mean	1_120_2_Skewness	0_0_0_Mean
	0_0_0_Mean	2_180_2_Entropy	Deg90_Correlation
	...		

The outputs of these models were combined to get the final result for the chest according to the strategy in Figure 4. And the thresholds we used in the experiment are below.

Threshold of lung level is 0.5;

Threshold of zone level is 0.64;

Low threshold of sum probability of chest is 0.43;

High threshold of sum probability of chest is 0.54.

We used the LIBSVM package to build and test the SVM classifiers that were used in the ensemble. This package uses a variant of the Sequential Minimal Optimization Algorithm [18, 19] to build the SVM. An RBF kernel was used in the SVM, and the γ and C parameters which determine classifier complexity were both varied in the range $[2^{-15}; 2^{15}]$ in order to arrive at the model with optimal complexity and good generalization. Table 3 gives a summary of the results in the database.

Table 3. Experimental results in the database.

	Sensitivity	Specificity	Accuracy
Test result	99.5%	92.7%	95.5%

6 Conclusions and Discussion

Pneumoconiosis detection, especially for coal workers, has attracted many researchers in the past decades [20-22]. The computer scheme we present here improved the lung segmentation accuracy in chest radiography, extracted the features created a feature vector by feature down-selection, built up multiple classifiers on lung fields and subdivision zones, and reported out based on patient's chest according to the clinic domain knowledge.

Our future job is to build up multi-class classifiers on each 6 zones of lung to get quantitative diagnosis for pneumoconiosis staging, which is one of the most challenge tasks in clinic. Moreover, we plan to train the parameters in larger training set or test the computer scheme in clinic data set, which is to optimize the feature vector, avoid overtraining issue, and make the algorithms more robust.

Currently the tool is expected to provide value in occupational disease detection, including mass screening and quantitative staging applications. However, it should be possible to extend it to interstitial diseases with diffuse character, such as lung tuberculosis.

References

1. Li, Dehong: Physicians textbook of occupational diseases pneumoconiosis. Publishing House of Peoples Daily, Beijing (2004)
2. ILO: Guidelines for the use of the ilo international classification of radiographs of pneumoconiosis - revised edition. International Labor Office, (2000)
3. Xinhua News Agency: Report on Ministry of Health of the Peoples Republic of China. Beijing (2009)
4. Ginneken, BV., Romeny BM, Viergever MA.: Computer-aided diagnosis in chest radiography: a survey. *IEEE Transactions on Medical Imaging* 20(12): 1228-1241 (2001)
5. Ginneken, BV.: Computer-aided diagnosis in chest radiography. *Medical Physics*, 2001; 28(6): 1144-1150.
6. Vittitoe, NF., Vargas-Voracek, R. and Floyd, CE: Identification of lung regions in chest radiographs using Markov random field modeling. *Medical Physics* 25(6): 976-985, 1998
7. Ginneken, BV., Alejandro F. Frangi, Joes J. Staal: Active shape model segmentation with optimal features. *IEEE Transactions on Medical Imaging* 21(8): 924-933, 2002
8. Iglesias, I., Souto, M., Alegria, A.M.: Lung segmentation on postero-anterior digital chest radiographs using active contours. *Lecture Notes in Computer Science* 3138: 538-546, 2004
9. Duda, Richard O., Hart, Peter E. and Stork, David G.: *Pattern Classification*. Wiley-Interscience (2000)
10. Hsu, Chih-Wei, Chang, Chih-Chung and Chih-Jen Lin: A practical guide to support vector classification. <http://www.csie.ntu.edu.tw/~cjlin/libsvm/index.html>
11. Vladimir Vapnik: *Statistical Learning Theory*. John Wiley (1998)
12. Nello Cristianini and John Shawe-Taylor: *An Introduction to Support Vector Machines*. Cambridge University Press (2000)
13. Kramer, S., Widmer, G., Pfahringer, B. and DeGroeve, M.: Prediction of ordinal classes using regression trees. *Fundamenta Informaticae* (2001)
14. McCullagh, P.: Regression model for ordinal data. *Journal of the Royal Statistical Society Series B* (42) 109-142 (1980)
15. Potharst, J.R. and Bioch, J.C.: Decision trees for ordinal classification. *Intelligent Data Analysis*, Vol. 4(2), pp. 97-112 (2000)
16. Chu & Keerthi: New approaches to support vector ordinal regression. *Proc. 22nd International Conference on Machine Learning*, Bonn, Germany (2005)
17. Mathieson, M. J.: Ordinal models for neural networks, in *Neural Networks in Financial Engineering*. World Scientific, Singapore (1995)
18. Platt, J. C.: Sequential minimal optimization: A fast algorithm for training support vector machines. *Tech. Rep. MSR-TR-98-14*, Microsoft Research, 1998.
19. Fan, R.-E., Chen, P.-H. and Lin, C.-J.: Working set selection using second order information for training svm. *Journal of Machine Learning Research*, vol. 6, pp. 1889~1918, 2005.
20. Kruger, R., Thompson, W. and Turner, A.: Computer diagnosis of pneumoconiosis. *IEEE Trans. Syst., Man, Cybern.*, vol. SMC-4, 1974.
21. Turner, A., Kruger, R. and Thompson, W.: Automated computer screening of chest radiographs for pneumoconiosis. *Investigation Radiology*, vol. 11, pp. 258--266, 1976.
22. Hiroshi Kondo and Takaharu Kouda: Computer-aided diagnosis for pneumoconiosis using neural network. *The 14th IEEE Symposium on Computer-Based Medical Systems (CMBS'01)*, 2001.