

Learning to quantify emphysema extent: What labels do we need?

Silas Nyboe Ørting, Jens Petersen, Laura H. Thomsen, Mathilde M. W. Wille and Marleen de Bruijne *Member, IEEE*,

Abstract—Accurate assessment of pulmonary emphysema is crucial to assess disease severity and subtype, to monitor disease progression and to predict lung cancer risk. However, visual assessment is time-consuming and subject to substantial inter-rater variability and standard densitometry approaches to quantify emphysema remain inferior to visual scoring. We explore if machine learning methods that learn from a large dataset of visually assessed CT scans can provide accurate estimates of emphysema extent. We further investigate if machine learning algorithms that learn from a scoring of emphysema extent can outperform algorithms that learn only from a scoring of emphysema presence. We compare four Multiple Instance Learning classifiers that are trained on emphysema presence labels, and five Learning with Label Proportions classifiers that are trained on emphysema extent labels. We evaluate performance on 600 low-dose CT scans from the Danish Lung Cancer Screening Trial and find that learning from emphysema presence labels, which are much easier to obtain, gives equally good performance to learning from emphysema extent labels. The best classifiers achieve intra-class correlation coefficients around 0.90 and average overall agreement with raters of 78% and 79% on six emphysema extent classes versus inter-rater agreement of 83%.

I. INTRODUCTION

EMPHYSEMA is a lung pathology characterized by destruction of lung tissue and enlargement of airspaces in the lung, causing shortness of breath. It is a main component of chronic obstructive pulmonary disease (COPD), a leading cause of mortality and morbidity world-wide [1]. Emphysema can be assessed on chest CT scans and its extent quantified by densitometry, where the amount of tissue affected by emphysema is estimated by measuring the percentage of lung volume with attenuation below a specific threshold. Although densitometry is simple and provides a single interpretable measurement of emphysema extent, it is also highly dependent on scanner hardware, reconstruction parameters [2] and software used for analysis [3].

This study was financially supported by the Danish Council for Independent Research (DFR) and the Netherlands Organization for Scientific Research (NWO). The sponsors had no involvement in the work.

S. Ørting and J. Petersen are with Department of Computer Science, University of Copenhagen, Copenhagen, Denmark

L. Thomsen is with Department of Internal Medicine, Hvidovre Hospital, Copenhagen Denmark

M. M. W. Wille is with Department of Diagnostic Imaging, Bispebjerg Hospital, Copenhagen, Denmark

M. de Bruijne is with Department of Computer Science, University of Copenhagen, Copenhagen, Denmark and Biomedical Imaging Group Rotterdam, Departments of Radiology and Medical Informatics, Erasmus MC - University Medical Center Rotterdam, The Netherlands

Manuscript received September 7th, 2018

An alternative to densitometry is visual assessment that can quantify extent and characterize emphysema subtype. The COPDGene CT Workshop Group [4] proposed a standard for visual assessment of COPD based on the characterization of emphysema appearance from the Fleischner society [5]. A slightly modified version of the standard was used for visual assessment in the Danish Lung Cancer Screening Trial (DLCST), where it was shown to be predictive of lung cancer [6]. A similar classification scheme defined in [7] was used in [8] where it was shown that visual presence and severity of emphysema is associated with increased mortality independent of densitometric measures of emphysema severity. The downside of visual assessment is that it is time-consuming and subject to inter-rater variability [4], [9].

Automated approaches based on the appearance of emphysema could provide fast and reproducible assessment of emphysema extent, location and sub-type, thus combining the superior disease characterization of visual assessment with the ease of densitometry. For instance [10] has shown that a shape-model of bullae-like structures can be used for emphysema detection. We have previously used machine learning algorithms based on texture features to predict regional emphysema presence [11] and emphysema extent [12]. Other learning based approaches have focused on discovery of emphysema patterns using supervised [13] and unsupervised [14], [15] learning, COPD detection and staging [16], [17] and emphysema detection in the more general context of interstitial lung disease classification [18], [19].

Multiple Instance Learning (MIL) has been used with success in a number of the prior works on emphysema and COPD detection [11], [16], [17] and for many related medical image analysis tasks as reviewed in [20]. MIL is a learning setting where the objects of interest are represented by a collection of samples. Each collection has a binary label and the goal is to learn which samples in a collection are “responsible” for the label. MIL has been very successful at detecting presence of abnormalities. However, visual assessment systems for lung disease, such as those developed for COPD [4], give estimates of affected lung tissue that is better captured by proportion labels. Label Proportions Learning (LLP) is the natural extension of MIL to cases where labels are proportions, but despite the success of MIL, LLP has seen almost no usage in medical imaging.

In this work we present the largest comparison yet of machine learning methods for assessing emphysema extent, extending our previous work on emphysema presence prediction [11], where a MIL method was used for regional

emphysema detection, and our work on extent prediction [12], where the LLP method Cluster Model Selection was used for regional emphysema extent prediction. We compare four MIL methods, of which three have not been used for emphysema detection before, and five LLP methods, of which four have not been used for emphysema detection or in medical imaging before. We investigate if learning from emphysema extent labels improves performance over learning from emphysema presence labels. Knowing what can be achieved by learning from labels of different quality and cost is paramount for cost-effective development and application of machine learning methods for clinical decision making.

II. MATERIALS AND METHODS

We view emphysema extent prediction as a bag learning problem. Bag learning is a machine learning setting where we are given a set of instances, a partition of the instances into bags and a labeling of the bags. The objective is to learn to predict both instance and bag labels for unseen data. In this work we view a region of the lung as a bag and patches sampled from the region as instances. The bag labels are regional emphysema extent scores, corresponding to estimated percentage of affected lung volume, and we wish to predict which patches contain emphysema, as well as the extent of emphysema in the region. Representing a scan as a set of patches provides a representation of local patterns in the lungs. By controlling the patch size we can focus on the scale at which patterns are expected to be distinct.

More formally, let \mathcal{X} be an instance space, \mathcal{Y} an instance label space, \mathcal{Z} a bag label space and $\mathbf{b} = (\mathbf{x} \subseteq \mathcal{X}, z \in \mathcal{Z})$ a labeled bag of instances. We use superscripts to refer to the label (\mathbf{b}^z), instances (\mathbf{b}^x) and instance labels (\mathbf{b}^y) associated with a bag \mathbf{b} . For a set of m bags $\mathbf{B} = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_m\}$, \mathbf{b}_i^x are the instances in the i 'th bag and \mathbf{b}_{ij}^x is the j 'th instance in the i 'th bag. We define the learning problem as

$$\arg \max_{\mathbf{Y}, h, \Theta} P(\mathbf{Y}, h, \Theta | \mathbf{B}), \quad (1)$$

where $\mathbf{Y} = \cup_{i=1}^m \mathbf{b}_i^y$ is a labeling of instances, $\Theta : \mathcal{Y} \mapsto \mathcal{Z}$ is a bag labeling function relating \mathbf{b}_i^y to \mathbf{b}_i^z and $h : \mathcal{X} \mapsto \mathcal{Y}$ is a hypothesis relating the instances \mathbf{b}_i^x to the corresponding instance labels \mathbf{b}_i^y , i.e. h is a method for predicting \mathbf{b}_i^y from \mathbf{b}_i^x .

Two well known bag learning settings are multiple instance learning (MIL) and learning with label proportions (LLP). In the standard MIL setting bag labels are binary, instance labels are binary and bag labels are related to instance labels by the max rule, i.e. a bag is positive if at least one instance is positive

$$\mathbf{b}_i^z = \Theta_{\max}(\mathbf{b}_i^y) = \max_j \mathbf{b}_{ij}^y. \quad (2)$$

This MIL setting is powerful because it allow us to learn about instance labels when only little information about the relation between instance and bag labels is available. A potential issue with the max rule is that it focuses on the single most discriminative instance. This could lead to a situation with good bag-level detection but poor localization and extent prediction. Including information about the proportion of positive

instances could improve localization and extent prediction. In the standard LLP setting, bag labels are proportions, instance labels are binary and bag labels are related to instance labels by the mean rule, i.e. the bag label is the proportion of positive instances

$$\mathbf{b}_i^z = \Theta_{\text{mean}}(\mathbf{b}_i^y) = \frac{1}{|\mathbf{b}_i|} \sum_j \mathbf{b}_{ij}^y. \quad (3)$$

Although MIL methods require binary labels for training, i.e. $\Theta : \mathcal{Y} \mapsto \{0, 1\}$, we can use Θ_{mean} at test time to obtain proportion estimates of emphysema extent.

A. Methods

We compared four MIL methods (logistic, SVM, *mi*-logistic, *mi*-SVM) and five LLP methods (beta, Cluster Model Selection, α -SVM, α -logistic, Laplacian Mean Map). The methods can be grouped into three distinct strategies used to solve the bag learning problem: the simple strategy, the relabeling strategy and the mean strategy. Some methods have previously been successfully applied to emphysema and COPD prediction, logistic, SVM and *mi*-SVM in [17], [11] and Cluster Model Selection in [12]. The LLP methods, α -SVM [21] and Laplacian Mean Map [22], have been shown to perform well on a variety of datasets. The beta method [23] can be seen as an LLP version of logistic and the *mi*-logistic and α -logistic methods are logistic regression versions of their SVM counterparts.

a) Simple strategy: In the simple strategy the bag learning problem is solved by ignoring intra-bag dependencies. We assign each instance the label of the bag it came from, i.e. $\mathbf{b}_{ij}^y = \mathbf{b}_i^z$, and train a standard supervised method on the instance labels. Labels for unseen bags are predicted by predicting instance labels and using Θ_{mean} to derive a bag label. The learning problem now becomes

$$\arg \max_{\phi} P(h_{\phi} | \mathbf{Y}, \mathbf{X}), \quad (4)$$

where $\mathbf{X} = \cup_{i=1}^m \mathbf{b}_i^x$ is the set of instances and h is a model parameterized by ϕ . We consider two simple MIL models, logistic regression (log) and a support vector machine (svm); and one simple LLP model, beta regression [23] (beta). Beta regression is a generalized linear model where the outcome \mathbf{Y} follows a beta distribution allowing us to perform regression with proportion outcomes. Note that bag labels are only used for the initial instance labeling, so Θ plays no role in the simple strategy.

b) Relabeling strategy: In the relabeling strategy the bag learning problem is solved by splitting it into two sub problems that are solved separately, a standard learning problem (5) and an instance labeling problem (6),

$$\arg \max_{\phi} P(h_{\phi} | \mathbf{Y}, \mathbf{X}) \quad (5)$$

$$\arg \max_{\mathbf{Y}} P(\mathbf{Y} | h_{\phi}, \Theta, \mathbf{Z}), \quad (6)$$

where $\mathbf{Z} = \cup_{i=1}^m \{\mathbf{b}_i^z\}$ is the set of bag labels and $\Theta = \Theta_{\max}$ for MIL and $\Theta = \Theta_{\text{mean}}$ for LLP. The two sub problems are iterated until convergence, with the result of (5) being

used for (6) and the result of (6) being used for (5). We consider two relabeling MIL methods, *mi*-SVM [24] (*misvm*) and *mi*-logistic (*milog*); and three relabeling LLP methods, α -SVM [21] (*psvm*), α -logistic (*plog*) and Cluster Model Selection [25] (*cms*). The methods *milog* and *plog* have not previously been published, they are however very similar to their *svm* counterparts and we do not include the derivation here. Details can be found in Appendix B. The *cms* algorithm differs from the other relabeling methods in that it solves (5) by unsupervised clustering. We use a version of *cms* previously described in [12].

c) *Mean strategy*: In the mean strategy the bag learning problem is solved by replacing the direct dependence on instance labels with a dependence on a mean statistic $\boldsymbol{\mu}$ calculated over all instances

$$\arg \max_{\phi} P(h_{\phi} | \boldsymbol{\mu}, \mathbf{X}). \quad (7)$$

$\boldsymbol{\mu}$ is defined as

$$\boldsymbol{\mu} = \frac{1}{n} \sum_i \mathbf{Y}_i \mathbf{X}_i \quad (8)$$

where $\mathbf{Y}_i \in \{-1, 1\}$ and n is the number of instances. Knowing $\boldsymbol{\mu}$ allow us to minimize the expected risk of a large class of loss functions. However, since the instance labels \mathbf{Y} are still unknown $\boldsymbol{\mu}$ must be estimated. The basic idea for the mean strategy is to express $\boldsymbol{\mu}$ in terms of bag-wise averages and solve for these bag-wise averages

$$\boldsymbol{\mu} = \sum_{i=1}^m \frac{|\mathbf{b}_i|}{n} \boldsymbol{\mu}_i \quad (9)$$

$$\boldsymbol{\mu}_i = \mathbf{b}_i^z \boldsymbol{\mu}_i^+ - (1 - \mathbf{b}_i^z) \boldsymbol{\mu}_i^- \quad (10)$$

where $|\mathbf{b}_i|$ the number of instances in bag i and $\boldsymbol{\mu}_i, \boldsymbol{\mu}_i^+, \boldsymbol{\mu}_i^-$, are the unknown mean instance, mean positive instance and mean negative instance of bag i , respectively. Equation (10) yields an underdetermined system of equations. We consider a single mean LLP method, Laplacian Mean Map [22] (*lmm*), that solves the system of equations by regularizing with a bag similarity term. We refer to [22] for further details.

B. Measures

We measure agreement in the following way. Let n_k be the number of ratings for case k and $n_{c,k}$ the number of times label c is assigned to case k . Agreement on label c over all cases is defined as

$$\frac{\sum_k n_{c,k} (n_{c,k} - 1)}{\sum_k n_{c,k} (n_k - 1)}. \quad (11)$$

Overall agreement across labels is defined as

$$\frac{\sum_{c,k} n_{c,k} (n_{c,k} - 1)}{\sum_k n_k (n_k - 1)}. \quad (12)$$

When all cases have two ratings Equation 11 corresponds to the Jaccard similarity and Equation 12 corresponds to multi-class accuracy. For multiple raters these measures ensure that

partial agreement, e.g. two out of three, is counted appropriately. We measure prevalence of label c as the proportion of times a case is assigned label c out of all assignments.

$$\frac{\sum_k n_{c,k}}{\sum_{c,k} n_{c,k}} \quad (13)$$

C. Data

Examples of the appearance of emphysema in CT scans are provided in Appendix A.

1) Study population, CT scanning & visual assessment:

We used data collected in the Danish Lung Cancer Screening Trial (DLCST) [26]. The screening arm of the study enrolled 2052 participants for annual low dose CT screening. Scan parameters are reproduced below verbatim from [26].

All CT scans of the study were performed on a MDCT scanner (16 rows Philips Mx 8000, Philips Medical Systems, Eindhoven, The Netherlands). Scans were performed supine after full inspiration with caudocranial scan direction including the entire ribcage and upper abdomen with a low dose technique, 120kV and 40 mAs. Scans were performed with spiral data acquisition with the following acquisition parameters: Section collimation 16×0.75 mm, pitch 1.5, rotation time 0.5 second.

We used a 1mm reconstruction with pixel size of $0.78\text{mm} \times 0.78\text{mm}$.

We obtained visual assessment of emphysema from [9], where screening participants with at least two CT scans were selected for visual assessment ($n=1990$). The visual assessment used a slight modification of the assessment sheets from [4]. Baseline and final followup scan was assessed by two experts. Emphysema extent was assessed for the top, middle and lower regions of each lung. The regions were defined as above carina, between carina and lower pulmonary vein, and below lower pulmonary vein. Each region was assigned a score of 0%, 1-5%, 6-25%, 26-50%, 51-75% or 76-100% indicating the extent of emphysema in the region.

In general, prevalence was highest and rater agreement best in the upper regions. Prevalence and agreement for the upper right region are summarized in Table I. Prevalence for emphysema extent above 26% is low (≈ 36 of 1200 subjects). Agreement on the five categories indicating emphysema presence was around 50%. Using only two categories (0%, $\geq 1\%$) improves agreement to 82% on the emphysema category. Although the original six categories provide more information than presence/absence labels, they are noisier and likely harder to learn from.

2) *Patches*: We represented a lung region as a collection of 3D patches sampled from the region. Sampling was done by choosing patch center locations uniformly at random within the region. We used a fixed patch size of approximately 11mm^3 to match the size of the secondary lobule [5] and allowed overlapping patches. For each patch we extracted a set of multi-scale filter responses and used equalized histograms of the filter responses as the final representation of the patch.

Extent	All		Extent	Presence	
	Agreement	Prev		Agreement	Prev
0%	94 (93–95)	75.2	0%	94 (93–95)	75.2
1-5%	54 (47–60)	14.7			
6-25%	44 (34–53)	7.0			
26-50%	45 (26–61)	2.0	≥ 1%	81 (78–85)	24.8
51-75%	57 (26–80)	0.9			
≥ 76%	67 (00–99)	0.1			
Overall	83 (81–85)		91 (89–92)		

TABLE I: Agreement and mean prevalence in the upper right region of the training data. Numbers are percentages. First three columns are for all six categories, last three columns are for presence/absence. 95% confidence intervals for agreement estimated by bootstrapping are given in parenthesis.

The filters used were Gaussian blur, gradient magnitude, eigenvalues of the Hessian, Laplacian of Gaussian, Gaussian curvature and the Frobenius norm of the Hessian. All filters were calculated at scales 1mm, 2mm and 4mm. The filters and the patch sampling strategy have previously been used successfully for COPD texture analysis in [16].

III. EXPERIMENTS AND RESULTS

We created a set of 1800 bags by sampling patches from the upper right region of 1800 subjects, such that each bag corresponds to one unique subject. We chose the upper right region because it has the highest prevalence and agreement. Results in [11] indicate that although absolute performance decreases when training on regions with lower prevalence and agreement, this decrease is relatively smaller than the decrease in rater agreement and prevalence.

Each bag contained 100 patches from a single subject. The bags were split into three non-overlapping datasets of 400 training and 200 test bags. Each experiment was run on all three datasets. In each split, we used two-fold cross validation on the training bags for parameter tuning. The three separate sets of classifiers were finally trained on all 400 training bags and performance estimated on the corresponding 200 test bags.

All classifiers provide posterior instance label probabilities which were converted to binary predictions using a classifier specific instance threshold fitted on the training bags. Parameters are summarized in Appendix C.

To train and evaluate we derived point estimates of emphysema extent by converting visually assessed extent intervals to interval midpoints and taking the mean over both raters. As an example, for a region with ratings 6-25% and 1-5%, the ratings are converted to 15.5% and 3% and combined into 9.25%. The point estimates were used directly for training LLP classifiers and thresholded at zero to obtain binary labels for training MIL classifiers.

A. Extent prediction accuracy

The prediction performance of the nine classifiers is illustrated with correlation plots in Fig. 1. The numbers in the title of each plot are intra-class correlation coefficients (ICC, two-way model, agreement) for each replication. The average ICC coefficients over the three replications are shown

log	svm	milog	misvm	beta	plog	psvm	cms	lmm
0.88	0.86	0.90	0.89	0.89	0.69	0.71	0.78	0.91

TABLE II: Average ICC of of emphysema extent over the three replications. MIL on the left, LLP on the right.

	0	1-5	6-25	26-50	51-75	76-100	Overall
log	88	35	54	38	17	00	74
svm	89	37	49	27	12	00	74
milog	85	35	50	36	29	00	71
misvm	91	39	58	36	31	00	79
beta	91	35	54	24	47	00	78
plog	72	26	57	45	00	00	58
psvm	27	15	21	35	51	17	24
cms	62	22	49	28	37	00	49
lmm	81	31	49	30	44	17	66
Rater	95	49	53	47	32	00	83

TABLE III: Agreement percentages between classifiers and raters averaged over replications and raters. First four columns show MIL classifiers, next five columns show LLP classifiers, last column shows rater agreement.

in Table II. We see clear positive correlation between reference and predicted extent for all classifiers. It appears that plog and cms tend to underestimate extent, whereas psvm tends to overestimate for cases with low extent but seems to perform very well for larger extent. Most classifiers show the largest variation for 15% reference extent. For extent larger than 15% we see very few cases with 0% extent predicted. The ICC values across replications, also seen in Fig. 1, illustrate that the performance of some classifiers varies a lot, with a difference of 0.25 in the worst case (cms). The most stable ICC performance is seen for lmm, which also has the highest average ICC.

B. Replacing a rater

The ICC of predicted extent and average rater extent provides an overall measure of performance and a validation that the classifiers have learned what they are trained to do. We are also interested in how the classifiers compare against each rater on the original rater task, i.e assign one of six intervals of emphysema extent. We converted predicted extent into the six extent intervals and calculated agreement with each rater. Agreement was calculated as described in Section II-B and is reported in Table III as an average over raters and replications. The final column in Table III provides inter-rater agreement averaged over replications. We see that beta and misvm have the highest overall agreement (78 and 79), which is not far from the overall rater agreement of 83. The agreement pattern of misvm and beta also seem to match that of the raters to a large degree, with a large agreement on 0% extent cases. It is interesting that psvm has the worst overall performance yet seems to outperform the other classifiers and raters for 51-75% extent. However, we cannot rule out that this is just a random coincidence given the low prevalence of that class. Another interesting observation is that the best results relative to inter-rater agreement is seen for 6-25% and 51-75%, with four classifiers having better agreement scores than inter-rater agreement.

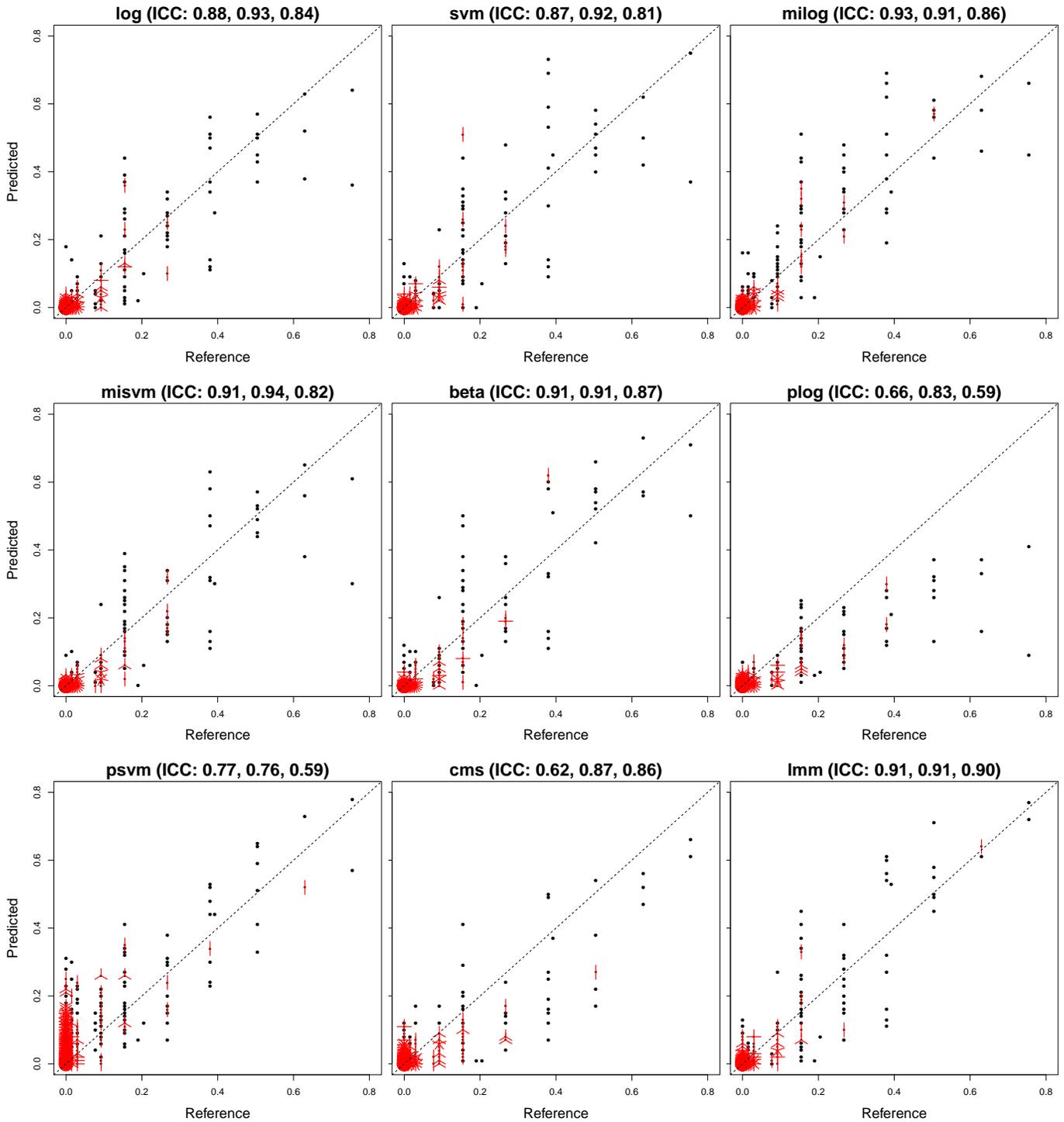


Fig. 1: Correlation between predicted and reference extent of emphysema. The x-axis is reference extent and the y-axis is predicted extent. The amount of red “petals” at a coordinate indicates the amount of coincident points. Plot titles show ICC coefficients for each replication.

C. Ranking classifiers

We use Friedman and Nemenyi test for comparing classifiers as suggested in [27]. We test the hypothesis H_0 : *All classifiers are equal* using Friedman test and significance level $\alpha = 0.05$. This test is based on the rank of the classifiers for each sample prediction. We use the absolute distance from

predicted extent to reference extent to assign ranks. In all three replications we get $p < 0.001$ for the Friedman test and reject the hypothesis that all classifiers have equal performance. We then test the pairwise hypothesis H_0 : *The classifiers are equal* for all pairs of classifiers using the Nemenyi test. The results of the Nemenyi tests are summarized in Fig. 2. Columns are

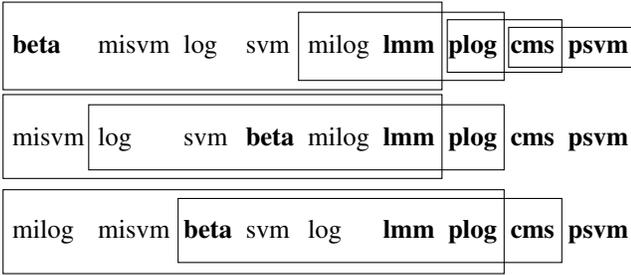


Fig. 2: Grouping of classifiers based on difference in extent prediction performance as decided by the Nemenyi test. Classifiers in the same box are not significantly different ($\alpha = 0.05$). Columns are sorted by mean rank over all test samples in descending order. **Bold** typeface indicates LLP methods.

sorted by average ranks and H_0 is rejected for classifiers that are not in the same box. We see that the LLP methods plog, cms and psvm are consistently ranked low, confirming the low ICC in Table II and the low overall agreement in Table III. Even though lmm is never significantly different from the best classifier, it is consistently ranked low. We also saw in Table III that lmm had low overall agreement with raters yet achieved the best average ICC. It is also interesting that misvm is consistently ranked in the top-2.

1) *Label stability*: We investigate label stability under changes in training data by predicting all test data with the trained models from each replication. For each classifier we got three sets of predictions of 60,000 instances and 600 bags. We converted bag predictions to the six extent intervals and measured agreement between replications for predicted bag and instance labels for each classifier. The stability results are summarized in Table IV. For bag labels, most classifiers have best agreement on 0% extent followed by 6-25%. Overall, beta and misvm are the most stable classifiers for both bag and instance labels, whereas milog is the most stable classifier on 6-25%, 26-50% and 51-75%. The missing scores for 51-75% and 76-100% are because there are no predictions of these classes in any of the replications. The inter-rater agreement on bag labels is included in the last row of Table IV. We see that misvm and beta always have equal or better agreement than the raters, and most methods have better agreement than raters on all non-zero extent scores.

IV. DISCUSSION & CONCLUSION

We have focused on comparing MIL methods, which have previously shown promising results for COPD and emphysema detection, with LLP methods that can learn directly from proportion labels. While end-to-end learning using CNNs have shown promising results for medical imaging tasks, and have just recently been used for emphysema quantification [28], we decided to use classic scale space features to focus on the aspects of learning from binary versus proportion labels, and to establish performance of classic feature engineering approaches.

Using the average rater as reference, the best classifiers achieve ICC coefficients around 0.9. Average overall agreement between the best classifiers and each rater on six em-

physema extent intervals is close to the inter-rater agreement (78-79% vs 83%). For some extent intervals the classifiers are better than the inter-rater agreement. These results show that that the presented approach to automatic emphysema extent prediction is viable and could be useful for routine assessment of emphysema extent.

The four best performing classifiers, beta, misvm, milog and lmm, have very similar ICC coefficients, with lmm being slightly more consistent across replications. However, beta and misvm show superior overall prediction of extent intervals with a much better discrimination of CT scans without visible emphysema compared to milog and lmm. Overall stability of beta and misvm is also superior to milog and lmm, although milog shows more stable predictions for the lower prevalence extent intervals 6-25%, 26-50% and 51-75%. Learning from scores indicating emphysema extent did not appear to be advantageous for extent prediction compared to learning based on emphysema presence alone. The MIL classifiers, misvm and milog, and the LLP classifiers, beta and lmm, show comparable performance.

One possible explanation for the lack of improved performance when training on extent labels, is that the extent labels are too noisy, as the relatively large disagreement between observers suggests. Obtaining more accurate and precise extent labels is costly and it is not clear if it is possible to improve the label quality significantly. In this work we have combined the emphysema estimates of two raters by simple averaging of point estimates. In [12] we showed that performance of the cms classifier improved when learning from labels incorporating rater uncertainty over learning from averaged point estimates. The approach used in [12] is not directly applicable to the other methods used here and we have used point estimates to keep the comparison fair. Recent work on classification of retinal images with a CNN-based method [29] show that modeling individual raters can improve performance over simple averaging of multiple raters. Although more than 30 raters were used in [29] it is possible that a more complex model of rater annotations could also improve performance when only two raters are used.

Another possible factor is that the model of proportion labels is too simple to exploit the additional information in the labels. The results in [28] indicate that learning from proportion labels can help more complex models based on CNNs to converge faster and to a better optima than learning from binary labels. A possible explanation for this is that explicitly modeling proportion labels has a regularizing effect on the feature learning part of CNNs, which would also explain why we do not see improved performance for LLP methods when features are fixed.

We considered three strategies for learning from bag labels, the simple strategy, the relabeling strategy and the mean strategy. For the MIL classifiers it appears that the relabeling strategy is best, whereas for the LLP classifiers it appears that the simple and mean strategies are best. One reason the simple strategy works better for LLP than for MIL could be that a proportion instance label as used in simple LLP is interpreted as a probability of emphysema in the patch, whereas the binary instance labels as used in simple MIL are

	Bag						Overall	Instance	
	0%	1-5%	6-25%	26-50%	51-75%	76-100%		E	NE
log	89	60	68	58	57	–	79	35	97
svm	89	57	60	51	58	–	78	44	97
milog	84	60	81	82	77	–	78	44	96
misvm	95	70	77	70	74	–	88	52	98
beta	96	71	72	58	62	50	89	56	98
plog	65	54	73	78	–	–	61	07	95
psvm	24	40	47	47	60	0	41	12	82
cms	65	60	59	45	12	–	61	05	93
lmm	82	60	68	58	67	17	73	43	97
Rater	95	49	53	47	32	0	83	–	–

TABLE IV: Label stability. Agreement percentages between predictions from each replication. Instance columns are for binary instance predictions (Emphysema / No Emphysema). A dash (–) indicates no predictions in that category.

interpreted as the probability the patch came from a CT scan with emphysema. In this sense, the proportion instance labels match the intended objective, predicting the proportion of patches with emphysema, much better than the binary instance labels.

A limitation of this study is that we have only trained and validated the classifiers on the upper right region of the lung. Due to the lower prevalence and agreement of visual scoring in the remaining five regions, we expect some decrease in extent prediction accuracy for these regions, similar to what was observed in [11] for regional emphysema detection. Investigating the performance over all regions should be considered in future work. However, the results in [11] show that a simple MIL classifier trained on subject-level presence/absence labels can provide the same performance as a classifier trained on region-level presence/absence labels. In light of the results here, this suggests that a MIL classifier, such as misvm, could provide accurate regional emphysema extent estimates even when trained only on subject-level presence of emphysema.

In conclusion, the best performing classifiers have close to human-level performance and are promising candidates for automatic quantification of emphysema extent. Furthermore, MIL classifiers having access to only emphysema presence labels perform just as well as LLP classifiers with access to emphysema extent labels. Reducing the labeling task from estimating emphysema extent to indicating presence, reduces the cost of training and makes it more feasible to implement in new settings.

V. REFERENCES

REFERENCES

- [1] “Global Strategy for the Diagnosis, Management and Prevention of COPD, Global Initiative for Chronic Obstructive Lung Disease (GOLD) 2017,” 2017. [Online]. Available: <http://www.goldcopd.org/>
- [2] J. Sieren, J. Newell, P. Judy, D. Lynch, K. Chan, J. Guo, and E. Hoffman, “Reference standard and statistical model for intersite and temporal comparisons of CT attenuation in a multicenter quantitative lung study,” *Medical physics*, vol. 39, no. 9, pp. 5757–5767, 2012.
- [3] M. O. Wielpütz, D. Bardarova, O. Weinheimer, H.-U. Kauczor, M. Eichinger, B. J. Jobst, R. Eberhardt, M. Koenigkam-Santos, M. Pud-erbach, and C. P. Heussel, “Variation of densitometry on computed tomography in COPD—influence of different software tools,” *PLoS one*, vol. 9, no. 11, p. e112898, 2014.
- [4] COPDGene CT Workshop Group, R. G. Barr, E. A. Berkowitz, F. Bigazzi, F. Bode, J. Bon, R. P. Bowler, C. Chiles, J. D. Crapo, G. J. Criner, J. L. Curtis, C. Dass, A. Dirksen, M. T. Dransfield, G. Edula, L. Eriksson, A. Friedlander, M. Galperin-Aizenberg, W. B. Gefter, D. S. Gierada, P. A. Grenier, J. Goldin, M. K. Han, N. A. Hanania, N. N. Hansel, F. L. Jacobson, H.-U. Kauczor, V. L. Kinnula, D. A. Lipson, D. A. Lynch, W. MacNee, B. J. Make, A. J. Mamary, H. Mann, N. Marchetti, M. Mascalchi, G. McLennan, J. R. Murphy, D. Naidich, H. Nath, J. D. N. Jr., M. Pistolesi, E. A. Regan, J. J. Reilly, R. Sandhaus, J. D. Schroeder, F. Sciruba, S. Shaker, A. Sharafkhaneh, E. K. Silverman, R. M. Steiner, C. Strange, N. Sverzellati, J. H. Tashjian, E. J. van Beek, L. Washington, G. R. Washko, G. Westney, S. A. Wood, and P. G. Woodruff, “A combined pulmonary-radiology workshop for visual evaluation of COPD: Study design, chest CT findings and concordance with quantitative evaluation,” *COPD: Journal of Chronic Obstructive Pulmonary Disease*, vol. 9, no. 2, pp. 151–159, 2012.
- [5] D. M. Hansell, A. A. Bankier, H. MacMahon, T. C. McLoud, N. L. Müller, and J. Remy, “Fleischner Society: Glossary of terms for thoracic imaging,” *Radiology*, vol. 246, no. 3, pp. 697–722, 2008.
- [6] M. M. W. Wille, L. H. Thomsen, J. Petersen, M. de Bruijne, A. Dirksen, J. H. Pedersen, and S. B. Shaker, “Visual assessment of early emphysema and interstitial abnormalities on CT is useful in lung cancer risk analysis,” *European Radiology*, pp. 1–8, 2015.
- [7] D. A. Lynch, J. H. Austin, J. C. Hogg, P. A. Grenier, H.-U. Kauczor, A. A. Bankier, R. G. Barr, T. V. Colby, J. R. Galvin, P. A. Gevenois *et al.*, “CT-definable subtypes of chronic obstructive pulmonary disease: a statement of the Fleischner Society,” *Radiology*, vol. 277, no. 1, pp. 192–205, 2015.
- [8] D. A. Lynch, C. M. Moore, C. Wilson, D. Nevrekar, T. Jennermann, S. M. Humphries, J. H. Austin, P. A. Grenier, H.-U. Kauczor, M. K. Han *et al.*, “CT-based visual classification of emphysema: Association with mortality in the COPDGen study,” *Radiology*, p. 172294, 2018.
- [9] M. M. W. Wille, L. H. Thomsen, A. Dirksen, J. Petersen, J. H. Pedersen, and S. B. Shaker, “Emphysema progression is visually detectable in low-dose CT in continuous but not in former smokers,” *European Radiology*, vol. 24, no. 11, pp. 2692–2699, Nov 2014.
- [10] R. Wiemker, M. Sevenster, H. MacMahon, F. Li, S. Dalal, A. Tahmasebi, and T. Klinder, “Automated assessment of imaging biomarkers for the PanCan lung cancer risk prediction model with validation on NLST data,” in *Medical Imaging 2017: Computer-Aided Diagnosis*, vol. 10134. International Society for Optics and Photonics, 2017, p. 1013421.
- [11] S. N. Ørting, J. Petersen, L. H. Thomsen, M. M. W. Wille, and M. de Bruijne, “Detecting emphysema using multiple instance learning,” in *2018 IEEE 15th International Symposium on Biomedical Imaging (ISBI 2018)*, 2018.
- [12] S. N. Ørting, J. Petersen, M. M. Wille, L. H. Thomsen, and M. de Bruijne, “Quantifying emphysema extent from weakly labeled CT scans of the lungs using label proportions learning,” *The Sixth International Workshop on Pulmonary Image Analysis*, pp. 31–42, 2016.
- [13] P. J. Castaldi, R. S. J. Estépar, C. S. Mendoza, C. P. Hersh, N. Laird, J. D. Crapo, D. A. Lynch, E. K. Silverman, and G. R. Washko, “Distinct quantitative computed tomography emphysema patterns are associated with physiology and function in smokers,” *American Journal of Respiratory and Critical Care Medicine*, vol. 188, no. 9, pp. 1083–1090, 2013.
- [14] P. Binder, N. K. Batmanghelich, R. S. J. Estépar, and P. Golland, “Unsupervised discovery of emphysema subtypes in a large clinical cohort,” in *International Workshop on Machine Learning in Medical Imaging*. Springer, 2016, pp. 180–187.
- [15] J. Yang, E. D. Angelini, P. P. Balte, E. A. Hoffman, J. H. Austin, B. M. Smith, J. Song, R. G. Barr, and A. F. Laine, “Unsupervised discovery of spatially-informed lung texture patterns for pulmonary emphysema: The MESA COPD study,” in *International Conference on Medical Image*

- Computing and Computer-Assisted Intervention*. Springer, 2017, pp. 116–124.
- [16] L. Sørensen, M. Nielsen, P. Lo, H. Ashraf, J. H. Pedersen, and M. de Bruijne, “Texture-based analysis of COPD: A data-driven approach,” *IEEE Transactions on Medical Imaging*, vol. 31, no. 1, pp. 70–78, Jan 2012.
- [17] V. Cheplygina, L. Sørensen, D. M. J. Tax, J. H. Pedersen, M. Loog, and M. de Bruijne, “Classification of COPD with multiple instance learning,” in *International Conference on Pattern Recognition*, 2014, pp. 1508–1513.
- [18] Q. Wang, Y. Zheng, G. Yang, W. Jin, X. Chen, and Y. Yin, “Multi-scale rotation-invariant convolutional neural networks for lung texture classification,” *IEEE Journal of Biomedical and Health Informatics*, vol. PP, no. 99, pp. 1–1, 2017.
- [19] M. Gao, U. Bagci, L. Lu, A. Wu, M. Buty, H.-C. Shin, H. Roth, G. Z. Papadakis, A. Depeursinge, R. M. Summers, Z. Xu, and D. J. Mollura, “Holistic classification of CT attenuation patterns for interstitial lung diseases via deep convolutional neural networks,” *Computer Methods in Biomechanics and Biomedical Engineering: Imaging & Visualization*, vol. 6, no. 1, pp. 1–6, 2018.
- [20] G. Quellec, G. Cazuguel, B. Cochener, and M. Lamard, “Multiple-instance learning for medical image and video analysis,” *IEEE reviews in biomedical engineering*, vol. 10, pp. 213–234, 2017.
- [21] F. X. Yu, D. Liu, S. Kumar, T. Jebara, and S. Chang, “ ∞ SVM for learning with label proportions,” *Proceedings of International Conference on Machine Learning*, 2013.
- [22] G. Patrini, R. Nock, T. Caetano, and P. Rivera, “(Almost) no label no cry,” in *Advances in Neural Information Processing Systems 27*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger, Eds. Curran Associates, Inc., 2014, pp. 190–198.
- [23] S. Ferrari and F. Cribari-Neto, “Beta regression for modelling rates and proportions,” *Journal of Applied Statistics*, vol. 31, no. 7, pp. 799–815, 2004.
- [24] S. Andrews, I. Tsochantaridis, and T. Hofmann, “Support vector machines for multiple-instance learning,” in *Advances in neural information processing systems*, 2003, pp. 577–584.
- [25] M. Stolpe and K. Morik, “Learning from label proportions by optimizing cluster model selection,” in *Machine Learning and Knowledge Discovery in Databases*, ser. Lecture Notes in Computer Science, D. Gunopulos, T. Hofmann, D. Malerba, and M. Vazirgiannis, Eds. Springer Berlin Heidelberg, 2011, vol. 6913, pp. 349–364.
- [26] J. H. Pedersen, H. Ashraf, A. Dirksen, K. Bach, H. Hansen, P. Toennesen, H. Thorsen, J. Brodersen, B. G. Skov, M. Døssing, J. Mortensen, K. Richter, P. Clementsen, and N. Seersholm, “The Danish randomized lung cancer CT screening trial—overall design and results of the prevalence round,” *Journal of Thoracic Oncology*, vol. 4, no. 5, 2009.
- [27] J. Demšar, “Statistical comparisons of classifiers over multiple data sets,” *Journal of Machine learning research*, vol. 7, no. Jan, pp. 1–30, 2006.
- [28] G. Bortsova, F. Dubost, S. Ørting, I. Katramados, L. Hogeweg, L. Thomsen, M. Wille, and M. de Bruijne, “Deep learning from label proportions for emphysema quantification,” *arXiv preprint arXiv:1807.08601*, 2018.
- [29] M. Y. Guan, V. Gulshan, A. M. Dai, and G. E. Hinton, “Who said what: Modeling individual labelers improves classification,” *Proceedings of AAAI Conference*, 2018.
- [30] H.-T. Lin, C.-J. Lin, and R. C. Weng, “A note on Platt’s probabilistic outputs for support vector machines,” *Machine learning*, vol. 68, no. 3, pp. 267–276, 2007.

APPENDIX A
EMPHYSEMA

Fig. 3 shows slices from the upper right region of three CT scans. Background and airways have been masked. The left image is assessed as having no visible emphysema extent. The center image as having 6-25% and the right image as having 51-75% emphysema extent. For the center image, emphysema is predominately visible at the boundary of the lung, whereas it is distributed throughout the region in the right image.

APPENDIX B
METHODS

A. Notation

Let \mathcal{X} be an instance space, \mathcal{Y} an instance label space, \mathcal{Z} a bag label space and $\mathbf{b} = (\mathbf{x} \subseteq \mathcal{X}, z \in \mathcal{Z})$ a labeled bag of instances. We use superscripts to refer to the label (\mathbf{b}^z), instances (\mathbf{b}^x) and instance labels (\mathbf{b}^y) associated with a bag \mathbf{b} . For a set of m bags $\mathbf{B} = \{\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_m\}$, \mathbf{b}_i^x are the instances in the i 'th bag and $\mathbf{b}_{i,j}^x$ is the j 'th instance in the i 'th bag. For the set of all instances we use $\mathbf{X} = \cup_{i=1}^m \mathbf{b}_i^x$, for all instance labels we use $\mathbf{Y} = \cup_{i=1}^m \mathbf{b}_i^y$ and for all bag labels we use $\mathbf{Z} = \cup_{i=1}^m \{\mathbf{b}_i^z\}$.

B. *mi*-logistic

The bag learning problem for *mi*-logistic is a constrained optimization problem over model weights and unknown instance labels

$$\max_{\mathbf{w}, \mathbf{Y}} \prod_{i,j} p(\mathbf{b}_{i,j}^y | \mathbf{b}_{i,j}^x, \mathbf{w}) \quad (14)$$

$$\text{s.t. } \forall i : \Theta_{\max}(\mathbf{b}_i^y) = \mathbf{b}_i^z \in \{0, 1\}. \quad (15)$$

We use the heuristic for solving the *mi*-SVM problem from [24]. Initially, fix instance labels by setting them to bag labels, $\mathbf{b}_{i,j}^y = \mathbf{b}_i^z \forall i, j$. For fixed instance labels (14) reduces to standard logistic regression. Let $h(\cdot) = \sigma(\mathbf{w}^T \cdot)$ denote the fitted model. Instance labels are predicted as

$$\tilde{\mathbf{b}}_{i,j}^y = \mathbb{1}\{h(\mathbf{b}_{i,j}^x) > 0.5\} \quad (16)$$

and bag labels are predicted as

$$\tilde{\mathbf{b}}_i^z = \Theta_{\max}(\tilde{\mathbf{b}}_i^y). \quad (17)$$

Instance labels are then updated according to

$$\mathbf{b}_{i,j}^y = \begin{cases} 0 & \text{if } \mathbf{b}_i^z = 0 \\ 1 & \text{if } \mathbf{b}_i^z = 1, \tilde{\mathbf{b}}_i^z = 0, h(\mathbf{b}_{i,j}^x) > h(\mathbf{b}_{i,k}^x) \forall k \neq j \\ \tilde{\mathbf{b}}_{i,j}^y & \text{if otherwise} \end{cases} \quad (18)$$

The first clause ensures that instances from negative bags are always labeled negative. The second clause ensures that a positive bag predicted as negative will always have one positive instance by labeling the ‘‘most’’ positive instance as positive, and the third clause ensures all other instances in positive bags are relabeled to match the predicted class.

C. α -logistic

The α -logistic model can be derived by considering the joint probability over instances \mathbf{X} , bag labels \mathbf{Z} and instance labels \mathbf{Y}

$$P(\mathbf{Y}, \mathbf{X}, \mathbf{Z}) = P(\mathbf{Z} | \mathbf{Y}, \mathbf{X}) P(\mathbf{Y}, \mathbf{X}) \quad (19)$$

$$= P(\mathbf{Z} | \mathbf{Y}) P(\mathbf{Y}, \mathbf{X}) \quad \mathbf{Z} \perp \mathbf{X} | \mathbf{Y} \quad (20)$$

$$\propto P(\mathbf{Z} | \mathbf{Y}) P(\mathbf{Y} | \mathbf{X}) \quad P(\mathbf{X}) = \text{Constant} \quad (21)$$

We use a logistic model for instance labels and a binomial model for bag labels.

$$P(\mathbf{Y} | \mathbf{X}) = \prod_{i,j} P(\mathbf{b}_{i,j}^y | \mathbf{b}_{i,j}^x, \mathbf{w}) \quad (22)$$

$$= \prod_{i,j} \sigma(\mathbf{w}^T \mathbf{b}_{i,j}^x)^{\mathbf{b}_{i,j}^y} (1 - \sigma(\mathbf{w}^T \mathbf{b}_{i,j}^x))^{1 - \mathbf{b}_{i,j}^y} \quad (23)$$

$$P(\mathbf{Z} | \mathbf{Y}) = \prod_i P(\mathbf{b}_i^z | \mathbf{b}_i^y) = \quad (24)$$

$$\prod_i \binom{|\mathbf{b}_i|}{|\mathbf{b}_i^z|} \Theta_{\text{mean}}(\mathbf{b}_i^y)^{|\mathbf{b}_i^z|} (1 - \Theta_{\text{mean}}(\mathbf{b}_i^y))^{|\mathbf{b}_i| - |\mathbf{b}_i^z|} \quad (25)$$

substituting into (21) gives us

$$P(\mathbf{Z} | \mathbf{Y}) P(\mathbf{Y} | \mathbf{X}) = \prod_i P(\mathbf{b}_i^z | \mathbf{b}_i^y) \prod_j P(\mathbf{b}_{i,j}^y | \mathbf{b}_{i,j}^x, \mathbf{w}) \quad (26)$$

We want to find the \mathbf{Y} and \mathbf{w} that maximize (26)

$$\arg \max_{\mathbf{Y}, \mathbf{w}} \prod_i P(\mathbf{b}_i^z | \mathbf{b}_i^y) \prod_j P(\mathbf{b}_{i,j}^y | \mathbf{b}_{i,j}^x, \mathbf{w}) \quad (27)$$

We do this by fixing \mathbf{Y} and \mathbf{w} iteratively. For fixed \mathbf{Y} we get standard logistic regression. For fixed \mathbf{w} we can optimize over each bag individually

$$\arg \max_{\mathbf{b}_i^y} P(\mathbf{b}_i^z | \mathbf{b}_i^y) \prod_{j=1} P(\mathbf{b}_{i,j}^y | \mathbf{b}_{i,j}^x, \mathbf{w}). \quad (28)$$

This can be done with the same greedy method used for α -SVM in [21].

APPENDIX C
PARAMETERS

All classifiers provide probability estimates of instance labels and a classifier-specific instance threshold was fitted on the training data by trying all thresholds in the range $[0, 0.01, 0.02, \dots, 0.99, 1]$. Fitted thresholds are reported in Table V. There is a large variation in fitted instance thresholds across classifiers, and for some classifiers there is a large variation across replications. Variation across replications is an indication that the classifier has learned substantially different decision rules for each replication. Variation between classifiers could just be a scaling issue, but is at least an indication that interpreting instance predictions as probability estimates is problematic.

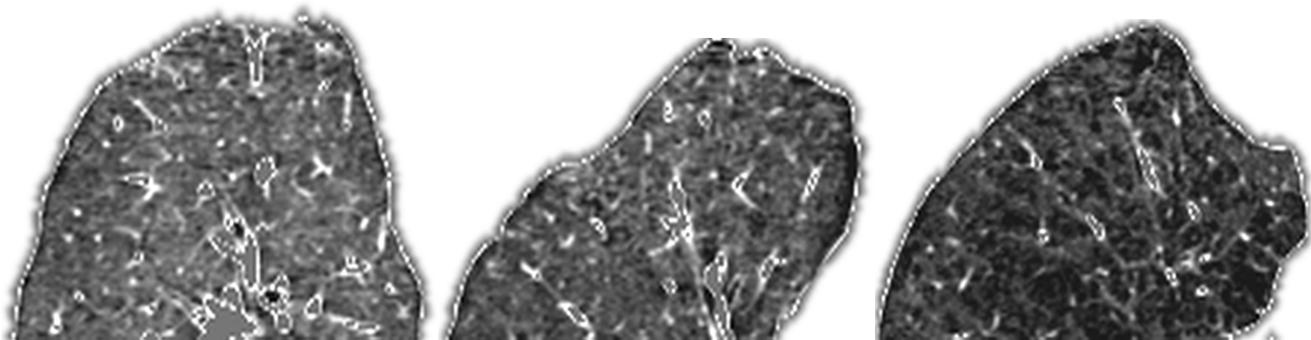


Fig. 3: Example slices. From left, visually assessed emphysema extent is 0%, 6-25% and 51-75%. Window level -780HU, window width 560HU.

Classifier	D1	D2	D3
log	0.78	0.79	0.60
beta	0.09	0.09	0.09
svm	0.77	0.76	0.70
misvm	0.85	0.94	0.78
milog	0.99	0.99	0.99
psvm	0.97	0.99	0.14
plog	0.01	0.01	0.01
cms	0.86	0.68	0.99
lmm	0.75	0.62	0.73

TABLE V: Fitted instance thresholds for each classifier and replication

A. beta

The implementation of beta regression requires uncorrelated features and we used the PCA algorithm to decorrelate features. We tried dimensionality reduction (only keep principal components with standard deviation ≥ 1). We tried two optimization methods, maximum likelihood estimation (ML) and bias correction (BC).

Fitted parameters

beta	D1	no dimensionality reduction, ML
	D2	no dimensionality reduction, BC
	D3	no dimensionality reduction, ML

B. svm, misvm, psvm

For all three classifiers we tried both linear and RBF kernels. In both cases we tried $C \in \{0.1, 1, 10, 100\}$. For psvm we tried $C_2 \in \{1, 10, 100, 1000\}$. For the RBF kernels we tried $\gamma \in \{0.1, 1\}$. We used Platt calibration [30] to obtain probability estimates from all three SVMs.

Fitted parameters

svm	D1	linear kernel, $C = 1$
	D2	linear kernel, $C = 0.1$
	D3	linear kernel, $C = 10$
misvm	D1	linear kernel, $C = 0.1$
	D2	linear kernel, $C = 0.1$
	D3	linear kernel, $C = 0.1$
psvm	D1	rbf kernel, $C = 1, C_2 = 1, \gamma = 0.1$
	D2	rbf kernel, $C = 0.1, C_2 = 1, \gamma = 1$
	D3	rbf kernel, $C = 1, C_2 = 1, \gamma = 0.1$

C. log, milog, plog

We tried dimensionality reduction using PCA for log. We did not use dimensionality reduction for milog and plog. We ran milog and plog until convergence of instance labels or for 20 iterations, whichever came first.

Fitted parameters

log	D1	no dimensionality reduction
	D2	no dimensionality reduction
	D3	dimensionality reduction

D. cms

We used the following fixed parameters, branching = 2, number of k -means iterations = 25, maximum iterations of CMA-ES = 1000, $\lambda = 13$. We tried number of clusters $k \in \{10, 20, 30, 40, 50, 60, 70, 80, 90, 100\}$.

Fitted parameters

cms	D1	$k = 70$
	D2	$k = 50$
	D3	$k = 100$

E. lmm

We tried

$$\lambda \in \{0, 1, 10, 100\}$$

$$\gamma \in \{0.00001, 0.0001, 0.001, 0.01, 0.1, 1\}$$

$$\sigma \in \{0.001, 0.01, 0.1, 0.125, 0.25, 0.5, 1.0\}$$

Fitted parameters

lmm	D1	$\lambda = 1, \gamma = 0.01, \sigma = 0.1$
	D2	$\lambda = 1, \gamma = 0.01, \sigma = 0.1$
	D3	$\lambda = 10, \gamma = 0.00001, \sigma = 0.25$