

Combining instance-based learning and logistic regression for multilabel classification

Weiwei Cheng · Eyke Hüllermeier

Received: 12 June 2009 / Revised: 12 June 2009 / Accepted: 16 June 2009 / Published online: 23 July 2009
Springer Science+Business Media, LLC 2009

Abstract Multilabel classification is an extension of conventional classification in which a single instance can be associated with multiple labels. Recent research has shown that, just like for conventional classification, instance-based learning algorithms relying on the nearest neighbor estimation principle can be used quite successfully in this context. However, since hitherto existing algorithms do not take correlations and interdependencies between labels into account, their potential has not yet been fully exploited. In this paper, we propose a new approach to multilabel classification, which is based on a framework that unifies instance-based learning and logistic regression, comprising both methods as special cases. This approach allows one to capture interdependencies between labels and, moreover, to combine model-based and similarity-based inference for multilabel classification. As will be shown by experimental studies, our approach is able to improve predictive accuracy in terms of several evaluation criteria for multilabel prediction.

Keywords Multilabel classification · Instance-based learning · Nearest neighbor classification · Logistic regression · Bayesian inference

1 Introduction

In conventional classification, each instance is assumed to belong to exactly one among a finite set of candidate classes. As opposed to this, the setting of multilabel classification allows an instance to belong to several classes simultaneously or, say, to attach more than one label to a single instance. Problems of this type are ubiquitous in everyday life: At IMDb, a movie can be categorized as *action*, *crime*, and *thriller*; a CNN news report can

Editors: Aleksander Kołcz, Dunja Mladenić, Wray Buntine, Marko Grobelnik, and John Shawe-Taylor.

W. Cheng · E. Hüllermeier (✉)
Department of Mathematics and Computer Science, University of Marburg, Marburg, Germany
e-mail: eyke@mathematik.uni-marburg.de

W. Cheng
e-mail: cheng@mathematik.uni-marburg.de

be tagged as *people* and *political* at the same time; in biology, a typical multilabel learning example is the gene functional prediction problem, where a gene can be associated with multiple functional classes, such as *metabolism*, *transcription*, and *protein synthesis*.

Multilabel classification has received increasing attention in machine learning in recent years, not only due to its practical relevance, but also as it is interesting from a theoretical point of view. In fact, even though it is possible to reduce the problem of multilabel classification to conventional classification in one way or the other and, hence, to apply existing methods for the latter to solve the former, straightforward solutions of this type are usually not optimal. In particular, since the presence or absence of the different class labels has to be predicted *simultaneously*, it is obviously important to exploit correlations and interdependencies between these labels. This is usually not accomplished by simple transformations to standard classification.

Even though quite a number of more sophisticated methods for multilabel classification has been proposed in the literature, the application of *instance-based learning* (IBL) has not been studied very deeply in this context so far. This is a bit surprising, given that IBL algorithms based on the nearest neighbor estimation principle have been applied quite successfully in classification and pattern recognition for a long time (Aha et al. 1991). A notable exception is the *multilabel k-nearest neighbor* (MLKNN) method that was recently proposed in Zhang and Zhou (2007), where it was shown to be competitive to state-of-the-art machine learning methods.

In this paper, we propose a novel approach to multilabel classification, which is based on a framework that unifies instance-based learning and logistic regression, comprising both methods as special cases. This approach overcomes some limitations of existing instance-based multilabel classification methods, including MLKNN. In particular, it allows one to capture interdependencies between the class labels in a proper way.

The rest of this paper is organized as follows: The problem of multilabel classification is introduced in a more formal way in Sect. 2, and related work is discussed in Sect. 3. Our novel method is then described in Sect. 4. Section 5 is devoted to experiments with several benchmark data sets. The paper ends with a summary and some concluding remarks in Sect. 6.

2 Multilabel classification

Let \mathbb{X} denote an instance space and let $\mathcal{L} = \{\lambda_1, \lambda_2, \dots, \lambda_m\}$ be a finite set of class labels. Moreover, suppose that each instance $\mathbf{x} \in \mathbb{X}$ can be associated with a subset of labels $L \in 2^{\mathcal{L}}$; this subset is often called the set of *relevant* labels, while the complement $\mathcal{L} \setminus L$ is considered as *irrelevant* for \mathbf{x} . Given training data in the form of a finite set T of observations in the form of tuples $(\mathbf{x}, L_{\mathbf{x}}) \in \mathbb{X} \times 2^{\mathcal{L}}$, typically assumed to be drawn independently from an (unknown) probability distribution on $\mathbb{X} \times 2^{\mathcal{L}}$, the goal in multilabel classification is to learn a classifier $h : \mathbb{X} \rightarrow 2^{\mathcal{L}}$ that generalizes well beyond these observations in the sense of minimizing the expected prediction loss with respect to a specific loss function; commonly used loss functions will be reviewed in Sect. 5.3.

Note that multilabel classification can be reduced to a conventional classification problem in a straightforward way, namely by considering each label subset $L \in 2^{\mathcal{L}}$ as a distinct (meta-)class. This approach is referred to as *label powerset* (LP) in the literature. An obvious drawback of this approach is the potentially large number of classes that one has to deal with in the newly generated problem; obviously, this number is $2^{|\mathcal{L}|}$ (or $2^{|\mathcal{L}|} - 1$ if the empty set is excluded as a prediction). This is the reason why LP typically works well if the original

label set \mathcal{L} is small but quickly deteriorates for larger label sets. Nevertheless, LP is often used as a benchmark, and we shall also include it in our experiments later on (cf. Sect. 5).

Another way of reducing multilabel to conventional classification is offered by the *binary relevance* approach. Here, a separate binary classifier h_i is trained for each label $\lambda_i \in \mathcal{L}$, reducing the supervision to information about the presence or absence of this label while ignoring the other ones. For a query instance \mathbf{x} , this classifier is supposed to predict whether λ_i is relevant for \mathbf{x} ($h_i(\mathbf{x}) = 1$) or not ($h_i(\mathbf{x}) = 0$). A multilabel prediction for \mathbf{x} is then given by $h(\mathbf{x}) = \{\lambda_i \in \mathcal{L} \mid h_i(\mathbf{x}) = 1\}$. Since binary relevance learning treats every label independently of all other labels, an obvious disadvantage of this approach is that it ignores correlations and interdependencies between labels.

Some of the more sophisticated approaches learn a multilabel classifier h in an indirect way via a scoring function $f: \mathbb{X} \times \mathcal{L} \rightarrow \mathbb{R}$ that assigns a real number to each instance/label combination. The idea is that a score $f(\mathbf{x}, \lambda)$ is in direct correspondence with the probability that λ is relevant for \mathbf{x} . Given a scoring function of this type, multilabel prediction can be realized via thresholding:

$$h(\mathbf{x}) = \{\lambda \in \mathcal{L} \mid f(\mathbf{x}, \lambda) \geq t\},$$

where $t \in \mathbb{R}$ is a threshold. As a byproduct, a scoring function offers the possibility to produce a ranking of the class labels, simply by ordering them according to their score. Sometimes, this ranking is even more desirable as a prediction, and indeed, there are several evaluation metrics that compare a true label subset with a predicted ranking instead of a predicted label subset (cf. Sect. 5.3).

3 Related work

Multilabel classification has received a great deal of attention in machine learning in recent years, and a number of methods has been developed, often motivated by specific types of applications such as text categorization (Schapire and Singer 2000; Ueda and Saito 2003; Kazawa et al. 2005; Zhang and Zhou 2006), computer vision (Boutell et al. 2004), and bioinformatics (Clare and King 2001; Elisseeff and Weston 2002; Zhang and Zhou 2006). Besides, several well-established methods for conventional classification have been extended to the multilabel case, including support vector machines (Godbole and Sarawagi 2004; Elisseeff and Weston 2002; Boutell et al. 2004), neural networks (Zhang and Zhou 2006), and decision trees (Vens et al. 2008).

In this paper, we are especially interested in instance-based approaches to multilabel classification, i.e., methods based on the nearest neighbor estimation principle (Dasarathy 1991; Aha et al. 1991). This interest is largely motivated by the *multilabel k -nearest neighbor* (MLKNN) method that has recently been proposed in Zhang and Zhou (2007). In that paper, the authors show that MLKNN performs quite well in practice. In the concrete experiments presented, MLKNN even outperformed some state-of-the-art model-based approaches to multilabel classification, including RankSVM and AdaBoost.MH (Elisseeff and Weston 2002; Comite et al. 2003).

MLKNN is a binary relevance learner, i.e., it learns a single classifier h_i for each label $\lambda_i \in \mathcal{L}$. However, instead of using the standard k -nearest neighbor (KNN) classifier as a base learner, it implements the h_i by means of a combination of KNN and Bayesian inference: Given a query instance \mathbf{x} with unknown multilabel classification $L \subseteq \mathcal{L}$, it finds the k nearest neighbors of \mathbf{x} in the training data and counts the number of occurrences of λ_i among

these neighbors. Considering this number, y , as information in the form of a realization of a random variable Y , the posterior probability of $\lambda_i \in L$ is given by

$$\mathbf{P}(\lambda_i \in L|Y = y) = \frac{\mathbf{P}(Y = y|\lambda_i \in L) \cdot \mathbf{P}(\lambda_i \in L)}{\mathbf{P}(Y = y)}, \tag{1}$$

which leads to the decision rule

$$h_i(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{P}(Y = y|\lambda_i \in L)\mathbf{P}(\lambda_i \in L) \geq \mathbf{P}(Y = y|\lambda_i \notin L)\mathbf{P}(\lambda_i \notin L) \\ 0 & \text{otherwise.} \end{cases}$$

The prior probabilities $\mathbf{P}(\lambda_i \in L)$ and $\mathbf{P}(\lambda_i \notin L)$ as well as the conditional probabilities $\mathbf{P}(Y = y|\lambda_i \in L)$ and $\mathbf{P}(Y = y|\lambda_i \notin L)$ are estimated from the training data in terms of corresponding relative frequencies. As an aside, we note that these estimations come with a relatively high computational complexity, since they involve the consideration of all k -neighborhoods of all training instances.

4 Combining IBL and logistic regression

In this section, we introduce a machine learning method whose basic idea is to consider the information that derives from examples similar to a query instance as a feature of that instance, thereby blurring the distinction between instance-based and model-based learning to some extent. This idea is put into practice by means of a learning algorithm that realizes instance-based classification as logistic regression.

4.1 KNN classification

Suppose an instance \mathbf{x} to be described in terms of features ϕ_i , $i = 1, 2 \dots n$, where $\phi_i(\mathbf{x})$ denotes the value of the i -th feature for instance \mathbf{x} . The instance space \mathbb{X} is endowed with a distance measure: $\Delta(\mathbf{x}, \mathbf{x}')$ is the distance between instances \mathbf{x} and \mathbf{x}' . We shall first focus on the case of binary classification and hence define the set of class labels by $\mathcal{Y} = \{-1, +1\}$. A tuple $(\mathbf{x}, y) \in \mathbb{X} \times \mathcal{Y}$ is called a labeled instance or example. \mathcal{D} denotes a sample that consists of N labeled instances (\mathbf{x}_i, y_i) , $1 \leq i \leq N$. Finally, a new instance $\mathbf{x}_0 \in \mathbb{X}$ (a query) is given, whose label $y_0 \in \{-1, +1\}$ is to be estimated.

The nearest neighbor (NN) principle prescribes to estimate the label of the yet unclassified query \mathbf{x}_0 by the label of the nearest (least distant) sample instance. The KNN approach is a slight generalization, which takes the $k \geq 1$ nearest neighbors of \mathbf{x}_0 into account. That is, an estimation \hat{y}_0 of y_0 is derived from the set $\mathcal{N}_k(\mathbf{x}_0)$ of the k nearest neighbors of \mathbf{x}_0 , usually by means of a *majority vote*:

$$\hat{y}_0 = \arg \max_{y \in \mathcal{Y}} \#\{\mathbf{x}_i \in \mathcal{N}_k(\mathbf{x}_0) \mid y_i = y\}. \tag{2}$$

4.2 IBL as logistic regression

A key idea of our approach is to consider the labels of neighbored instances as “features” of the query \mathbf{x}_0 whose label is to be estimated. It is worth mentioning that similar ideas have recently been exploited in relational learning (Getoor and Taskar 2007) and collective classification (Lu and Getoor 2003; Ghamrawi and McCallum 2005).

Denote by p_0 the prior probability of $y_0 = +1$ and by π_0 the corresponding posterior probability. Moreover, let $\delta_i \stackrel{\text{df}}{=} \Delta(\mathbf{x}_0, \mathbf{x}_i)$ be the distance between \mathbf{x}_0 and \mathbf{x}_i . Taking the known label y_i as information about the unknown label y_0 , we can consider the posterior probability

$$\pi_0 \stackrel{\text{df}}{=} \mathbf{P}(y_0 = +1 | y_i).$$

More specifically, Bayes’ rule yields

$$\begin{aligned} \frac{\pi_0}{1 - \pi_0} &= \frac{\mathbf{P}(y_i | y_0 = +1)}{\mathbf{P}(y_i | y_0 = -1)} \cdot \frac{p_0}{1 - p_0} \\ &= \rho \cdot \frac{p_0}{1 - p_0}, \end{aligned}$$

where ρ is the likelihood ratio. Taking logarithms on both sides, we get

$$\log\left(\frac{\pi_0}{1 - \pi_0}\right) = \log(\rho) + \omega_0 \tag{3}$$

with $\omega_0 = \log(p_0) - \log(1 - p_0)$.

Model (3) still requires the specification of the likelihood ratio ρ . In order to obey the basic principle underlying IBL, the latter should be a function of the distance δ_i . In fact, ρ should become large for $\delta_i \rightarrow 0$ if $y_i = +1$ and small if $y_i = -1$: Observing a very close instance \mathbf{x}_i with label $y_i = +1$ ($y_i = -1$) makes $y_0 = +1$ more (un)likely in comparison to $y_i = -1$. Moreover, ρ should tend to 1 as $\delta_i \rightarrow \infty$: If \mathbf{x}_i is too far away, its label does not provide any evidence, neither in favor of $y_0 = +1$ nor in favor of $y_0 = -1$. A parameterized function satisfying these properties is

$$\rho = \rho(\delta) \stackrel{\text{df}}{=} \exp\left(y_i \cdot \frac{\alpha}{\delta}\right),$$

where $\alpha > 0$ is a constant. Note that the choice of a special functional form for ρ is quite comparable to the specification of the kernel function used in (non-parametric) kernel-based density estimation, as well as to the choice of the weight function in weighted NN estimation. $\rho(\delta)$ actually determines the probability that two instances whose distance is given by $\delta = \Delta(\mathbf{x}_0, \mathbf{x}_i)$ do have the same label.

Now, taking the complete sample neighborhood $\mathcal{N}(\mathbf{x}_0)$ of \mathbf{x}_0 into account and—as in the naive Bayes approach—making the simplifying assumption of conditional independence, we obtain

$$\begin{aligned} \log\left(\frac{\pi_0}{1 - \pi_0}\right) &= \omega_0 + \alpha \sum_{\mathbf{x}_i \in \mathcal{N}(\mathbf{x}_0)} \frac{y_i}{\delta_i} \\ &= \omega_0 + \alpha \cdot \omega_+(\mathbf{x}_0), \end{aligned} \tag{4}$$

where $\omega_+(\mathbf{x}_0)$ can be seen as a summary of the evidence in favor of label +1. As can be seen, the latter is simply given by the sum of neighbors with label +1, weighted by their distance, minus the weighted sum of neighbors with label -1.

As concerns the classification of the query \mathbf{x}_0 , the decision is determined by the sign of the right-hand side in (4). From this point of view, (4) does basically realize a weighted NN estimation, or, stated differently, it is a “model-based” version of instance-based learning.

Still, it differs from the simple NN scheme in that it includes a bias term ω_0 , which plays the same role as the prior probability in Bayesian inference.

From a statistical point of view, (4) is nothing else than a logistic regression equation. In other words, taking a “feature-based” view of instance-based learning and applying a Bayesian approach to inference comes down to realizing IBL as logistic regression.

By introducing a *similarity measure* κ , inversely related to the distance function Δ , (4) can be written in the form

$$\log\left(\frac{\pi_0}{1 - \pi_0}\right) = \omega_0 + \alpha \sum_{\mathbf{x}_i \in \mathcal{N}(\mathbf{x}_0)} \kappa(\mathbf{x}_0, \mathbf{x}_i) \cdot y_i. \tag{5}$$

Note that, as a special case, this approach can mimic the standard KNN classifier (2), namely by setting $\omega_0 = 0$ and defining κ in terms of the (data-dependent) “KNN kernel”

$$\kappa(\mathbf{x}_0, \mathbf{x}_i) = \begin{cases} 1 & \text{if } \mathbf{x}_i \in \mathcal{N}_k(\mathbf{x}_0) \\ 0 & \text{otherwise.} \end{cases} \tag{6}$$

4.3 Estimation and classification

The parameter α in (4) determines the weight of the evidence

$$\omega_+(\mathbf{x}_0) = \sum_{\mathbf{x}_i \in \mathcal{N}(\mathbf{x}_0)} \kappa(\mathbf{x}_0, \mathbf{x}_i) \cdot y_i \tag{7}$$

and, hence, its influence on the posterior probability estimation π_0 . In fact, α plays the role of a smoothing (regularization) parameter. The smaller α is chosen, the smoother an estimated probability function (obtained by applying (5) to all points $\mathbf{x}_0 \in \mathcal{X}$) will be. In the extreme case where $\alpha = 0$, one obtains a constant function (equal to ω_0).

An optimal specification of α can be accomplished by adapting this parameter to the data \mathcal{D} , using the method of maximum likelihood (ML). For each sample point \mathbf{x}_j denote by

$$\omega_+(\mathbf{x}_j) \stackrel{\text{df}}{=} \sum_{\mathbf{x}_i \neq \mathbf{x}_j \in \mathcal{N}(\mathbf{x}_j)} \kappa(\mathbf{x}_i, \mathbf{x}_j) \cdot y_i$$

the sample evidence in favor of $y_j = +1$. The log-likelihood function is then given by the mapping

$$\alpha \mapsto \sum_{j: y_j = +1} w_0 + \alpha \omega_+(\mathbf{x}_j) - \sum_{j=1}^N \log(1 + \exp(w_0 + \alpha \omega_+(\mathbf{x}_j))), \tag{8}$$

and the optimal parameter α^* is the maximizer of (8). The latter can be computed by means of standard methods from logistic regression. The posterior probability π_0 for the query is then given by

$$\pi_0 = \frac{\exp(\omega_0 + \alpha^* \omega_+(\mathbf{x}_0))}{1 + \exp(\omega_0 + \alpha^* \omega_+(\mathbf{x}_0))}.$$

To classify \mathbf{x}_0 , one applies the decision rule

$$\hat{y}_0 \stackrel{\text{df}}{=} \begin{cases} +1 & \text{if } \pi_0 \geq 1/2 \\ -1 & \text{if } \pi_0 < 1/2. \end{cases}$$

Subsequently, we shall refer to the method outlined above as IBLR (Instance-Based Learning by Logistic Regression).

4.4 Including additional features

In the previous section, instance-based learning has been embedded into logistic regression, using the information coming from the neighbors of a query \mathbf{x}_0 as a “feature” of that query. In this section, we consider a possible generalization of this approach, namely the idea to extend the model (5) by taking further features of \mathbf{x}_0 into account:

$$\log\left(\frac{\pi_0}{1 - \pi_0}\right) = \alpha\omega_+(\mathbf{x}_0) + \sum_{\varphi_s \in \mathcal{F}} \beta_s \varphi_s(\mathbf{x}_0), \tag{9}$$

where $\mathcal{F} = \{\varphi_0, \varphi_1 \dots \varphi_r\}$ is a subset of the available features $\{\phi_0, \phi_1 \dots \phi_n\}$ and $\varphi_0 = \phi_0 \equiv 1$, which means that β_0 plays the role of ω_0 . Equation (9) is a common logistic regression model, except that $\omega_+(\mathbf{x}_0)$ is a “non-standard” feature.

The approach (9), that we shall call IBLR+, integrates instance-based and model-based (attribute-based) learning and, by estimating the regression coefficients in (9), achieves an optimal balance between both approaches. The extended model (9) can be interpreted as a logistic regression model of IBL, as outlined in Sect. 4.2, where the bias ω_0 is no longer constant:

$$\log\left(\frac{\pi_0}{1 - \pi_0}\right) = \omega_0(\mathbf{x}_0) + \alpha\omega_+(\mathbf{x}_0), \tag{10}$$

with $\omega_0(\mathbf{x}_0) \stackrel{\text{df}}{=} \sum \beta_s \varphi_s(\mathbf{x}_0)$ being an instance-specific bias determined by the model-based part of (9).

4.5 Extension to multilabel classification

So far, we only considered the case of binary classification. To extend the approach to multilabel classification with a label set $\mathcal{L} = \{\lambda_1, \lambda_2 \dots \lambda_m\}$, the idea is to train one classifier h_i for each label. For the i -th label λ_i , this classifier is derived from the model

$$\log\left(\frac{\pi_0^{(i)}}{1 - \pi_0^{(i)}}\right) = \omega_0^{(i)} + \sum_{j=1}^m \alpha_j^{(i)} \cdot \omega_{+j}^{(i)}(\mathbf{x}_0), \tag{11}$$

where $\pi_0^{(i)}$ denotes the (posterior) probability that λ_i is relevant for \mathbf{x}_0 , and

$$\omega_{+j}^{(i)}(\mathbf{x}_0) = \sum_{\mathbf{x} \in \mathcal{N}(\mathbf{x}_0)} \kappa(\mathbf{x}_0, \mathbf{x}) \cdot y_j(\mathbf{x}) \tag{12}$$

is a summary of the presence of the j -th label λ_j in the neighborhood of \mathbf{x}_0 ; here, $y_j(\mathbf{x}) = +1$ if λ_j is present (relevant) for the neighbor \mathbf{x} , and $y_j(\mathbf{x}) = -1$ in case it is absent (non-relevant).

Obviously, the approach (11) is able to take interdependencies between class labels into consideration. More specifically, the estimated coefficient $\alpha_j^{(i)}$ indicates to what extent the relevance of label λ_i is influenced by the relevance of λ_j . A value $\alpha_j^{(i)} \gg 0$ means that the presence of λ_j makes the relevance of λ_i more likely, i.e., there is a positive correlation. Correspondingly, a negative coefficient would indicate a negative correlation.

Note that the estimated probabilities $\pi_0^{(i)}$ can naturally be considered as scores for the labels λ_i . Therefore, a ranking of the labels is simply obtained by sorting them in decreasing order according to their probabilities. Moreover, a pure multilabel prediction for \mathbf{x}_0 is derived from this ranking via thresholding at $t = 0.5$.

Of course, it is also possible to combine the model (11) with the extension proposed in Sect. 4.4. This leads to a model

$$\log\left(\frac{\pi_0^{(i)}}{1 - \pi_0^{(i)}}\right) = \sum_{j=1}^m \alpha_j^{(i)} \cdot \omega_{+j}^{(i)}(\mathbf{x}_0) + \sum_{\varphi_s \in \mathcal{F}} \beta_s^{(i)} \varphi_s(\mathbf{x}_0). \quad (13)$$

We shall refer to the extensions (11) and (13) of IBLR to multilabel classification as IBLR-ML and IBLR-ML+, respectively.

5 Experimental results

This section is devoted to experimental studies that we conducted to get a concrete idea of the performance of our method. Before presenting the results of our experiments, we give some information about the learning algorithms and data sets included in the study, as well as the criteria used for evaluation.

5.1 Learning algorithms

For the reasons mentioned previously, our main interest is focused on MLKNN, which is arguably the state-of-the-art in instance-based multilabel ranking; we used its implementation in the MULAN package (Tsoumakas and Katakis 2007). MLKNN is parameterized by the size of the neighborhood, for which we adopted the value $k = 10$. This value is recommended in Zhang and Zhou (2007), where it was found to yield the best performance. For the sake of fairness, we use the same neighborhood size for our method, in conjunction with the KNN kernel (6). In both cases, the simple Euclidean metric (on the complete attribute space) was used as a distance function. For our method, we tried both variants, the pure instance-based version (11), and the extended model (13) with \mathcal{F} including all available features. Intuitively, one may expect the latter, IBLR-ML+, to be advantageous to the former, IBLR-ML, as it can use features in a more flexible way. Yet, one should note that, since we simply included all attributes in \mathcal{F} , each attribute will essentially be used twice in IBLR-ML+, thus producing a kind of redundancy. Besides, model induction will of course become more difficult, since a larger number of parameters needs to be estimated.

As an additional baseline we used binary relevance learning (BR) with three different base learners: logistic regression, C4.5 (the WEKA (Witten and Frank 2005) implementation J48 in its default setting), and KNN (again with $k = 10$). Finally, we also included label powerset (LP) with C4.5 as a base learner.

5.2 Data sets

Benchmark data for multilabel classification is not as abundant as for conventional classification, and indeed, experiments in this field are often restricted to a very few or even only a single data set. For our experimental study, we have collected a comparatively large number of seven data sets from different domains; an overview is given in Table 1.¹

¹All data sets are public available at <http://mlkd.csd.auth.gr/multilabel.html> and <http://lamda.nju.edu.cn/data.htm>.

Table 1 Statistics for the multilabel data sets used in the experiments. The symbol * indicates that the data set contains binary features; *cardinality* is the average number of labels per instance

Data set	Domain	#Instances	#Attributes	#Labels	Cardinality
<i>Emotions</i>	Music	593	72	6	1.87
<i>Image</i>	Vision	2000	135	5	1.24
<i>Genbase</i>	Biology	662	1186*	27	1.25
<i>Mediamill</i>	Multimedia	5000	120	101	4.27
<i>Reuters</i>	Text	7119	243	7	1.24
<i>Scene</i>	Vision	2407	294	6	1.07
<i>Yeast</i>	Biology	2417	103	14	4.24

The *emotions* data was created from a selection of songs from 233 musical albums (Trohidis et al. 2008). From each song, a sequence of 30 seconds after the initial 30 seconds was extracted. The resulting sound clips were stored and converted into wave files of 22050 Hz sampling rate, 16-bit per sample and mono. From each wave file, 72 features have been extracted, falling into two categories: rhythmic and timbre. Then, in the emotion labeling process, 6 main emotional clusters are retained corresponding to the Tellegen-Watson-Clark model of mood: amazed-surprised, happy-pleased, relaxing-clam, quiet-still, sad-lonely and angry-aggressive.

Image and *scene* are semantic scene classification data sets proposed, respectively, by Zhou and Zhang (2007) and Boutell et al. (2004), in which a picture can be categorized into one or more classes. In the scene data, for example, pictures can have the following classes: beach, sunset, foliage, field, mountain, and urban. Features of this data set correspond to spatial color moments in the LUV space. Color as well as spatial information have been shown to be fairly effective in distinguishing between certain types of outdoor scenes: bright and warm colors at the top of a picture may correspond to a sunset, while those at the bottom may correspond to a desert rock. Features of the image data set are generated by the SBN method (Maron and Ratan 1998) and essentially correspond to attributes in an RGB color space.

From the biological field, we have chosen the two data sets *yeast* and *genbase*. The yeast data set is about predicting the functional classes of genes in the Yeast *Saccharomyces cerevisiae*. Each gene is described by the concatenation of micro-array expression data and a phylogenetic profile, and is associated with a set of 14 functional classes. The data set contains 2417 genes in total, and each gene is represented by a 103-dimensional feature vector. In the *genbase* data, 27 important protein families are considered, including, for example, PDOC00064 (a class of oxydoreductases) and PDOC00154 (a class of isomerases). During the preprocessing, a training set was exported, consisting of 662 proteins that belong to one or more of these 27 classes.

From the text processing field, we have chosen a subset of the widely studied *Reuters-21578* collection (Sebastiani 2002). The seven most frequent categories are considered. After removing documents whose label sets or main texts are empty, 8866 documents are retained where only 3.37% of them are associated with more than one class label. After randomly removing documents with only one label, a text categorization data set containing 2,000 documents is obtained. Each document is represented as a bag of instances using the standard sliding window techniques, where each instance corresponds to a text segment enclosed in one sliding window of size 50 (overlapped with 25 words). “Function words” are removed from the vocabulary and the remaining words are stemmed. Instances in the bags

adopt the “bag-of-words” representation based on term frequency. Without loss of effectiveness, dimensionality reduction is performed by retaining the top 2% words with highest document frequency. Thereafter, each instance is represented as a 243-dimensional feature vector.

The *mediamill* data set is from the field of multimedia indexing and originates from the well-known TREC Video Retrieval Evaluation data (TRECVID 2005/2006) initiated by American National Institute of Standards and Technology (NIST), which contains 85 hours of international broadcast news data. The task in this data set is the automated detection of a lexicon of 101 semantic concepts in videos. Every instance of this data set has 120 numeric features including visual, textual, as well as fusion information. The trained classifier should be able to categorize an unseen instance to some of these 101 labels, e.g., face, car, male, soccer, and so on. More details about this data set can be found at Snoek et al. (2006).

5.3 Evaluation measures

To evaluate the performance of multilabel classification methods, a number of criteria and metrics have been proposed in the literature. For a classifier h , let $h(\mathbf{x}) \subseteq \mathcal{L}$ denote its multilabel prediction for an instance \mathbf{x} , and let L_x denote the true set of relevant labels. Moreover, in case a related scoring function f is also defined, let $f(\mathbf{x}, \lambda)$ denote the score assigned to label λ for instance \mathbf{x} . The most commonly used evaluation measures are defined as follows:

- *Hamming loss* computes the percentage of labels whose relevance is predicted incorrectly:

$$\text{HamLoss}(h) = \frac{1}{|\mathcal{L}|} |h(\mathbf{x}) \Delta L_x|, \tag{14}$$

where Δ is the symmetric difference between two sets.

- *One error* computes how many times the top-ranked label is not relevant:

$$\text{OneError}(f) = \begin{cases} 1 & \text{if } \arg \max_{\lambda \in \mathcal{L}} f(\mathbf{x}, \lambda) \notin L_x \\ 0 & \text{otherwise} \end{cases} \tag{15}$$

- *Coverage* determines how far one needs to go in the list of labels to cover all the relevant labels of an instance. This measure is loosely related to the precision at the level of perfect recall:

$$\text{Coverage}(f) = \max_{\lambda \in L_x} \text{rank}_f(\mathbf{x}, \lambda) - 1, \tag{16}$$

where $\text{rank}_f(\mathbf{x}, \lambda)$ denotes the position of label λ in the ordering induced by f .

- *Rank loss* computes the average fraction of label pairs that are not correctly ordered:

$$\text{RankLoss}(f) = \frac{\#\{(\lambda, \lambda') \mid f(\mathbf{x}, \lambda) \leq f(\mathbf{x}, \lambda'), (\lambda, \lambda') \in L_x \times \overline{L_x}\}}{|L_x| |\overline{L_x}|}, \tag{17}$$

where $\overline{L_x} = \mathcal{L} \setminus L_x$ is the set of irrelevant labels.

- *Average precision* determines for each relevant label $\lambda \in L_x$ the percentage of relevant labels among all labels that are ranked above it, and averages these percentages over all relevant labels:

$$\text{AvePrec}(f) = \frac{1}{|L_x|} \sum_{\lambda \in L_x} \frac{\#\{\lambda' \mid \text{rank}_f(\mathbf{x}, \lambda') \leq \text{rank}_f(\mathbf{x}, \lambda), \lambda' \in L_x\}}{\text{rank}_f(\mathbf{x}, \lambda)}. \tag{18}$$

Notice that only Hamming loss evaluates mere multilabel predictions (i.e., the multilabel classifier h), while the others metrics evaluate the underlying ranking function f . Moreover, smaller values indicate better performance for all measures except average precision. Finally, except for coverage, all measures are normalized and assume values between 0 and 1.

5.4 Results and discussion

The results of a cross validation study (10-fold, 5 repeats) are summarized in Table 2. As can be seen, the baseline methods BR and LP are in general not competitive. Looking at the average ranks, IBLR-ML consistently outperforms all other methods, regardless of the evaluation metric, indicating that it is the strongest method overall. The ranking among the three instance-based methods is IBLR-ML > IBLR-ML+ > MLKNN for all measures except OneError, for which the latter two change the position.

To analyze the results more thoroughly, we followed the two-step statistical test procedure recommended in Demsar (2006), consisting of a Friedman test of the null hypothesis that all learners have equal performance and, in case this hypothesis is rejected, a Nemenyi test to compare learners in a pairwise way. Both tests are based on the average ranks as shown in the bottom line in Table 2. Even though the Friedman test suggests that there are significant differences between the methods, most of the pairwise comparisons remain statistically non-significant (at a significance level of 5%); see Fig. 1. This is not surprising, however, given that the number of data sets included in the experiments, despite being much higher than usual, is still quite limited from a statistical point of view. Nevertheless, the overall picture taken from the experiments is clearly in favor of IBLR-ML.

As to MLKNN, it is interesting to compare this method with the BR-version of KNN. In fact, since MLKNN is a binary relevance learner, too, the only difference between these

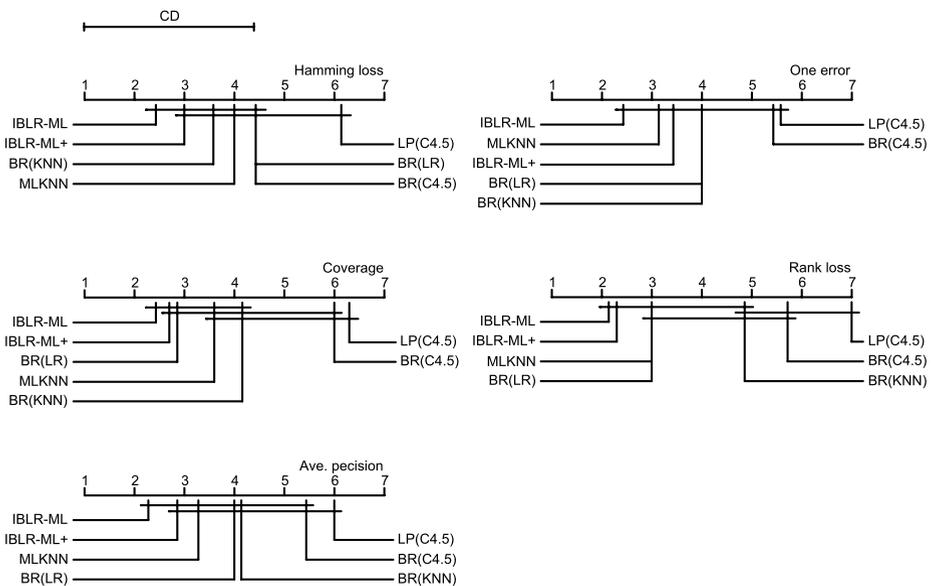


Fig. 1 Comparison of all classifiers against each other with the Nemenyi test. Groups of classifiers that are not significantly different (at $p = 0.05$) are connected

Table 2 Experimental results in terms of different evaluation measures. The number in brackets behind the performance value is the rank of the method on the corresponding data set (for each data set, the methods are ranked in decreasing order of performance). The average rank is the average of the ranks across all data sets

	iblr-ml+	iblr-ml	mlknn	lp	br-lr	br-c4.5	br-knn
Hamming							
Emotions	0.213(3)	0.185(1)	0.263(6)	0.265(7)	0.214(4)	0.253(5)	0.191(2)
Genbase	0.002(2)	0.002(3)	0.005(7)	0.002(4)	0.002(5)	0.001(1)	0.004(6)
Image	0.182(1)	0.189(2)	0.195(4)	0.257(7)	0.202(5)	0.245(6)	0.193(3)
Mediamill	0.03(6)	0.028(3)	0.027(2)	0.039(7)	0.029(4)	0.032(5)	0.027(1)
Reuters	0.044(1)	0.084(6)	0.073(5)	0.067(4)	0.049(2)	0.058(3)	0.09(7)
Scene	0.126(4)	0.084(1)	0.087(2)	0.142(7)	0.14(6)	0.133(5)	0.093(3)
Yeast	0.199(4)	0.194(1)	0.194(2)	0.28(7)	0.206(5)	0.25(6)	0.196(3)
Average rank	3	2.43	4	6.14	4.43	4.43	3.57
One Error							
Emotions	0.278(3)	0.257(1)	0.393(5)	0.43(7)	0.278(4)	0.422(6)	0.265(2)
Genbase	0.014(5)	0.007(2)	0.009(3)	0.01(4)	0.015(6)	0.003(1)	0.017(7)
Image	0.328(1)	0.367(2)	0.382(4)	0.507(6)	0.37(3)	0.512(7)	0.386(5)
Mediamill	0.356(5)	0.185(3)	0.136(2)	0.367(6)	0.277(4)	0.381(7)	0.133(1)
Reuters	0.076(1)	0.22(6)	0.185(5)	0.162(4)	0.086(2)	0.145(3)	0.233(7)
Scene	0.349(4)	0.224(2)	0.223(1)	0.394(6)	0.364(5)	0.411(7)	0.26(3)
Yeast	0.249(5)	0.227(1)	0.228(2)	0.351(6)	0.241(4)	0.389(7)	0.234(3)
Average rank	3.43	2.43	3.14	5.57	4	5.43	4
Coverage							
Emotions	1.844(4)	1.689(1)	2.258(5)	2.576(6)	1.836(3)	2.608(7)	1.771(2)
Genbase	0.356(1)	0.422(4)	0.561(7)	0.529(6)	0.391(3)	0.372(2)	0.436(5)
Image	0.963(1)	1.056(3)	1.129(5)	1.589(6)	1.052(2)	1.615(7)	1.102(4)
Mediamill	16.681(4)	15.161(3)	12.757(1)	49.469(7)	14.323(2)	47.996(6)	21.344(5)
Reuters	0.411(1)	0.758(4)	0.676(3)	0.986(7)	0.44(2)	0.852(6)	0.82(5)
Scene	0.911(5)	0.466(1)	0.472(2)	1.145(6)	0.871(4)	1.288(7)	0.551(3)
Yeast	6.289(3)	6.203(1)	6.273(2)	9.204(6)	6.492(4)	9.353(7)	6.517(5)
Average rank	2.71	2.43	3.57	6.29	2.86	6	4.14
Rank Loss							
Emotions	0.168(2)	0.146(1)	0.258(5)	0.499(7)	0.168(3)	0.372(6)	0.183(4)
Genbase	0.002(1)	0.004(2)	0.006(4)	0.017(7)	0.005(3)	0.006(5)	0.01(6)
Image	0.175(1)	0.197(3)	0.214(4)	0.537(7)	0.196(2)	0.409(6)	0.252(5)
Mediamill	0.05(4)	0.043(3)	0.037(1)	0.451(7)	0.041(2)	0.187(6)	0.117(5)
Reuters	0.026(1)	0.083(4)	0.069(3)	0.18(7)	0.03(2)	0.092(5)	0.113(6)
Scene	0.15(4)	0.076(1)	0.077(2)	0.393(7)	0.157(5)	0.299(6)	0.109(3)
Yeast	0.168(3)	0.164(1)	0.167(2)	0.545(7)	0.176(4)	0.362(6)	0.204(5)
Average rank	2.29	2.14	3	7	3	5.71	4.86
Ave. Prec.							
Emotions	0.794(3)	0.816(1)	0.71(5)	0.683(6)	0.794(4)	0.683(7)	0.805(2)
Genbase	0.989(3)	0.99(2)	0.989(4)	0.986(6)	0.988(5)	0.993(1)	0.982(7)
Image	0.789(1)	0.763(2)	0.748(5)	0.653(6)	0.763(3)	0.649(7)	0.752(4)
Mediamill	0.694(5)	0.731(3)	0.751(1)	0.498(7)	0.722(4)	0.582(6)	0.739(2)
Reuters	0.951(1)	0.859(6)	0.881(4)	0.871(5)	0.944(2)	0.889(3)	0.848(7)
Scene	0.773(4)	0.867(1)	0.867(2)	0.734(6)	0.769(5)	0.715(7)	0.844(3)
Yeast	0.763(3)	0.769(1)	0.764(2)	0.621(6)	0.754(5)	0.619(7)	0.761(4)
Average rank	2.86	2.29	3.29	6	4	5.43	4.14

Table 3 Classification error on binary classification problems. The number in brackets behind the performance value is the rank of the method on the corresponding data set (for each data set, the methods are ranked in decreasing order of performance). The average rank is the average of the ranks across all data sets

Data set	IBLR-ML+	IBLR-ML	MLKNN	BR-KNN
breast-cancer	.280(4)	.252(1)	.259(2)	.262(3)
breast-w	.037(3.5)	.037(3.5)	.036(2)	.034(1)
colic	.195(3)	.176(1)	.350(4)	.182(2)
credit-a	.135(2)	.132(1)	.328(4)	.138(3)
credit-g	.229(1)	.265(3)	.306(4)	.261(2)
diabetes	.233(1)	.263(4)	.259(3)	.256(2)
heart-statlog	.170(1)	.193(2.5)	.363(4)	.193(2.5)
hepatitis	.175(1)	.192(2)	.204(4)	.199(3)
ionosphere	.117(2.5)	.117(2.5)	.108(1)	.171(4)
kr-vs-kp	.018(1)	.044(2.5)	.044(2.5)	.046(4)
labor	.210(3)	.130(1)	.270(4)	.150(2)
mushroom	.000(1.5)	.000(1.5)	.001(3.5)	.001(3.5)
sick	.030(1)	.039(2)	.061(4)	.040(3)
sonar	.250(2)	.245(1)	.327(4)	.284(3)
tic-tac-toe	.125(1)	.137(3)	.136(2)	.317(4)
vote	.044(1)	.060(2)	.074(3)	.076(4)
Average rank	1.84	2.09	3.19	2.88

two methods concerns the incorporation of global information in MLKNN, which is accomplished through the Bayesian updating (1) of local information about the relevance of labels. From Table 2, it can be seen that MLKNN is better than BR-KNN in terms of all ranking measures, but not in terms of the Hamming loss, for which it is even a bit worse. Thus, in terms of mere relevance prediction, MLKNN does not seem to offer special advantages. Our explanation for this finding is that the incorporation of global information is indeed not useful for a simple 0/1 prediction. In a sense, this is perhaps not very surprising, given that the use of global information is somehow in conflict with the basic principle of local estimation underlying nearest neighbor prediction. Exploiting such information does, however, offer a reasonable way to *break ties between class labels*, which in turn explains the positive effect on ranking performance. In fact, one should note that, when simply scoring labels by the number of occurrences among the k neighbors of a query, such ties are quite likely; in particular, all non-relevant labels that never occur will have a score of 0 and will hence be tied. Resorting to global information about their relevance is then clearly more reasonable than breaking ties at random.

To validate our conjecture that the incorporation of global information in MLKNN is actually not very useful for mere relevance prediction, we have conducted an additional experiments using 16 binary classification problems from the UCI repository. Using this type of data makes sense, since, for a binary relevance learner, minimizing Hamming loss is equivalent to minimizing 0/1 loss for m binary classification problems that are solved independently of each other. The results of a 5 times 10-fold cross validation, summarized in Table 3, are completely in agreement with our previous study. MLKNN does indeed show the worst performance and is even outperformed by the simple BR-KNN. Interestingly, IBLR-ML+ is now a bit better than IBLR-ML. A reasonable explanation for this finding is that, compared to the multilabel case, the relevance information that comes from the neighbors of a query in binary classification only concerns a single label and, therefore, is rather sparse. Correspondingly, information about additional features is reevaluated.

6 Summary and conclusions

We have presented a novel approach to instance-based learning, called IBLR, that can be used for classification in general and for multilabel classification in particular. Considering label information of neighbored examples as features of a query instance, the idea of IBLR is to reduce instance-based learning formally to logistic regression. An optimal balance between global and local inference, and in the extended version IBLR+ also between instance-based and model-based (attribute-oriented) learning, can then be achieved by the estimation of optimal regression coefficients.

For multilabel classification, this idea is especially appealing, as it allows one to take interdependencies between different labels into consideration. These dependencies are directly reflected by the sign and magnitude of related regression coefficients. This ability distinguishes IBLR from hitherto existing instance-based methods for multilabel classification, and is probably one of the main factors for its excellent performance. In fact, our extensive empirical study has clearly shown that IBLR improves upon existing methods, in particular the MLKNN method that can be considered as the state-of-the-art in instance-based multilabel classification.

Interestingly, our results also suggest that the basic idea underlying MLKNN, namely to combine instance-based learning and Bayesian inference, is beneficial for the ranking performance but not in terms of mere relevance prediction. Investigating the influence on specific performance measures in more detail, and elaborating on (instance-based) methods for minimizing specific loss functions, is an interesting topic of future work. Besides, for IBLR+, we plan to exploit the possibility to combine instance-based and model-based inference in a more sophisticated way, for example by selecting optimal feature subsets for both parts instead of simply using all features twice.

References

- Aha, D., Kibler, D., & Alber, M. (1991). Instance-based learning algorithms. *Machine Learning*, 6(1), 37–66.
- Boutell, M. R., Luo, J., Shen, X., & Brown, C. M. (2004). Learning multilabel scene classification. *Pattern Recognition*, 37(9), 1757–1771.
- Clare, A., & King, R. D. (2001). Knowledge discovery in multilabel phenotype data. In L. D. Raedt & A. Siebes (Eds.), *Lecture notes in computer science* (Vol. 2168, pp. 42–53). Berlin: Springer.
- Comite, F. D., Gilleron, R., & Tommasi, M. (2003). Learning multilabel alternating decision tree from texts and data. In P. Perner & A. Rosenfeld (Eds.), *Lecture notes in computer science* (Vol. 2734, pp. 35–49). Berlin: Springer.
- Dasarathy, B. V., editor (1991). *Nearest neighbor (NN) norms: NN pattern classification techniques*. Los Alamitos: IEEE Comput. Soc.
- Demsar, J. (2006). Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7, 1–30.
- Elisseeff, A., & Weston, J. (2002). A kernel method for multilabelled classification. In T. G. Dietterich, S. Becker, & Z. Ghahramani (Eds.), *Advances in neural information processing systems* (Vol. 14, pp. 681–687). Cambridge: MIT Press.
- Getoor, L., & Taskar, B., editors (2007). *Introduction to statistical relational learning*. Cambridge: MIT Press.
- Ghamrawi, N., & McCallum, A. (2005). Collective multilabel classification. In *Proc. CIKM-05*, Bremen, Germany.
- Godbole, S., & Sarawagi, S. (2004). Discriminative methods for multilabeled classification. In *LNCS: Vol. 3056. Advances in knowledge discovery and data mining* (pp. 20–33). Berlin: Springer.
- Kazawa, H., Izumitani, T., Taira, H., & Maeda, E. (2005). Maximal margin labeling for multi-topic text categorization. In L. K. Saul, Y. Weiss, & L. Bottou (Eds.), *Advances in neural inf. proc. syst.* (Vol. 17). Cambridge: MIT Press.
- Lu, Q., & Getoor, L. (2003). Link-based classification. In *Proc. ICML-03* (pp. 496–503) Washington.

- Maron, O., & Ratan, A. L. (1998). Multiple-instance learning for natural scene classification. In *Proc. ICML* (pp. 341–349), Madison, WI.
- Schapire, R. E., & Singer, Y. (2000). Boostexter: a boosting-based system for text categorization. *Machine Learning*, 39(2), 135–168.
- Sebastiani, F. (2002). Machine learning in automated text categorization. *ACM Computing Surveys*, 34(1), 1–47.
- Snoek, C. G. M., Worring, M., van Gemert, J. C., Geusebroek, J. M., & Smeulders, A. W. M. (2006). The challenge problem for automated detection of 101 semantic concepts in multimedia. In *Proc. ACM multimedia* (pp. 421–430), Santa Barbara, USA.
- Trohidis, K., Tsoumakas, G., Kalliris, G., & Vlahavas, I. (2008). Multilabel classification of music into emotions. In *Proc. int. conf. music information retrieval*.
- Tsoumakas, G., & Katakis, I. (2007). Multi-label classification: An overview. *International Journal of Data Warehousing and Mining*, 3(3), 1–13.
- Ueda, N., & Saito, K. (2003). Parametric mixture models for multilabel text. In S. Becker & S. Thrun (Eds.), *Advances in neural information processing* (Vol. 15, pp. 721–728). Cambridge: MIT Press.
- Vens, C., Struyf, J., Schietgat, L., Dzeroski, S., & Blockeel, H. (2008). Decision trees for hierarchical multi-label classification. *Machine Learning*, 73, 185–214.
- Witten, I., & Frank, E. (2005). *Data mining: practical machine learning tools and techniques* (2nd ed.). San Francisco: Morgan Kaufmann.
- Zhang, M.-L., & Zhou, Z.-H. (2006). Multi-label neural networks with applications to functional genomics and text categorization. In *IEEE transactions on knowledge and data engineering* (Vol. 18, pp. 1338–1351).
- Zhang, M.-L., & Zhou, Z.-H. (2007). ML-kNN: A lazy learning approach to multilabel learning. *Pattern Recognition*, 40(7), 2038–2048.
- Zhou, Z.-H., & Zhang, M.-L. (2007). Multi-instance multilabel learning with application to scene classification. In B. Schölkopf, J. Platt, & T. Hofmann (Eds.), *Advances in neural inf. proc. syst.* (Vol. 19, pp. 1609–1616). Cambridge: MIT Press.