

Motivation

- Insufficient theoretical analysis in multi-label classification (MLC) papers
- The notion of “label dependence” is often used in a purely *intuitive* manner, without a precise *formal* definition
- The results are given *on average* without investigation under which conditions a given algorithm benefits
- The reasons for improvements are *not carefully distinguished*
- It is implicitly assumed that *one* algorithm is going to be beneficial for *multiple* error measures

Main Question

The risk-minimizing model h^* is formally defined as:

$$h^*(\mathbf{x}) = \arg \min_h \mathbb{E}_{\mathbf{Y}|\mathbf{x}} L(\mathbf{Y}, h) = \arg \min_h \sum_{\mathbf{y}} \mathbf{P}(\mathbf{y} | \mathbf{x}) L(\mathbf{y}, h),$$

where $L(\mathbf{Y}, \mathbf{y})$ is a loss function defined on multi-label predictions.

Do we have to take into account the conditional dependence between labels in order to obtain a risk-minimizing model?

Loss Functions and Risk Minimizers

loss function	risk minimizer
Hamming loss: $L_H(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \sum_{i=1}^m [y_i \neq h_i(\mathbf{x})]$	$h_i^*(\mathbf{x}) = \arg \max_{b \in \{0,1\}} \mathbf{P}(y_i = b \mathbf{x})$
Rank loss: $L_r(\mathbf{y}, \mathbf{f}(\mathbf{x})) = \sum_{(i,j): y_i > y_j} ([f_i < f_j] + \frac{1}{2}[f_i = f_j])$	$f_i^*(\mathbf{x}) = \mathbf{P}(y_i = 1 \mathbf{x})$
Subset 0/1 loss: $L_s(\mathbf{y}, \mathbf{h}(\mathbf{x})) = [\mathbf{y} \neq \mathbf{h}(\mathbf{x})]$	$\mathbf{h}^*(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}} \mathbf{P}(\mathbf{y} \mathbf{x})$

Consequences and Conjectures

- The risk-minimizing prediction for the Hamming and the rank loss can be obtained from the marginal distributions $\mathbf{P}(Y_i | \mathbf{x})$, $i = 1, \dots, m$, alone
- It is not necessary to know the joint label distribution $\mathbf{P}(\mathbf{Y} | \mathbf{x})$ on \mathcal{Y} and take the conditional dependence into account
- As opposed to this, the modeling of conditional dependence has to be taken into account in order to minimize the subset zero-one loss
- In general, a specific learning and prediction strategy has to be tailored for a given performance measure.

Estimation of the Joint Distribution

- The conditional dependence can be fully described by the joint distribution
- The joint distribution enables to compute Bayes optimal predictions of any loss function

Probabilistic Classifier Chains (PCC)

Given a query instance \mathbf{x} , the (conditional) probability of each label combination $\mathbf{y} = (y_1, \dots, y_m) \in \mathcal{Y}$ can be computed using the *product rule of probability*:

$$\mathbf{P}(\mathbf{y} | \mathbf{x}) = \mathbf{P}(y_1 | \mathbf{x}) \cdot \prod_{i=2}^m \mathbf{P}(y_i | \mathbf{x}, y_1, \dots, y_{i-1}).$$

Algorithm:

- Learn m functions $g_i(\cdot)$ on an augmented input space $\mathcal{X} \times \{0, 1\}^{i-1}$, taking y_1, \dots, y_{i-1} as additional attributes:

$$g_i : \mathcal{X} \times \{0, 1\}^{i-1} \rightarrow [0, 1]$$

$$(\mathbf{x}, y_1, \dots, y_{i-1}) \mapsto \mathbf{P}(y_i = 1 | \mathbf{x}, y_1, \dots, y_{i-1})$$

- Assuming that the function $g_i(\cdot)$ can be interpreted as a *probabilistic* classifier whose prediction is the probability that $y_i = 1$, then:

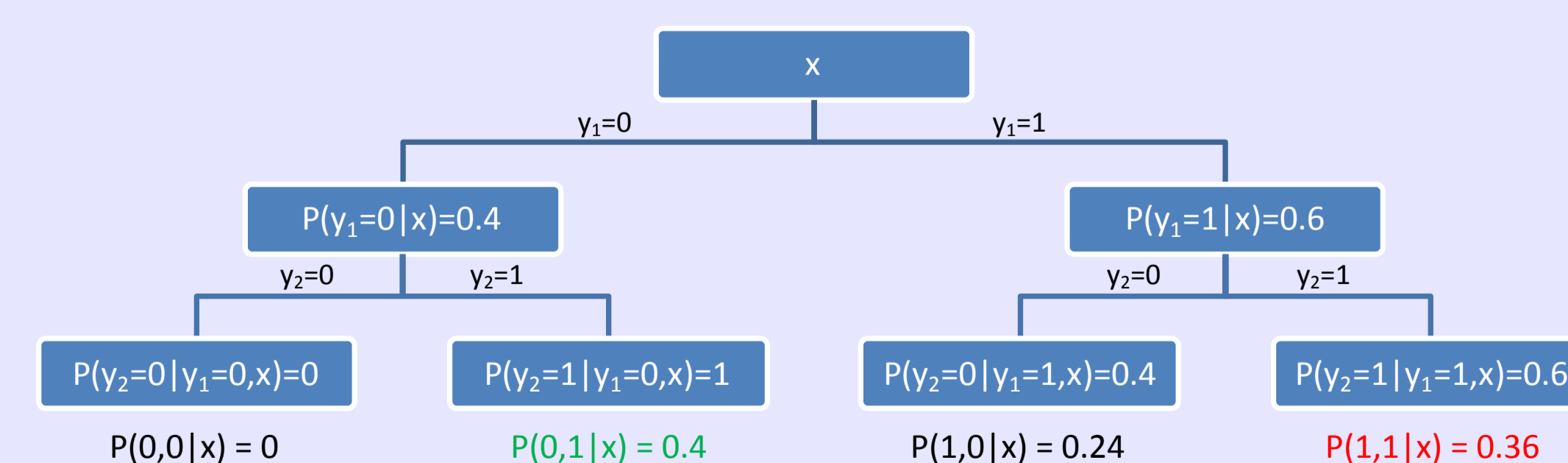
$$\hat{\mathbf{P}}(\mathbf{y} | \mathbf{x}) = g_1(\mathbf{x}) \cdot \prod_{i=2}^m g_i(\mathbf{x}, y_1, \dots, y_{i-1})$$

- Given $\mathbf{P}(\mathbf{y} | \mathbf{x})$ and a loss function $L(\cdot)$ to be minimized, an optimal prediction can then be derived in an explicit way:

$$h^*(\mathbf{x}) = \arg \min_h \mathbb{E}_{\mathbf{Y}|\mathbf{x}} L(\mathbf{Y}, h)$$

PCC vs. CC

- Original Classifier Chains (CC) (Read et al., 2009) follow a similar learning scheme
- The final prediction is computed by taking the predictions of consecutive models g_i , $i = 1, \dots, m$
- CC can be seen as a deterministic approximation of $\mathbf{P}(\mathbf{y} | \mathbf{x})$, in the sense of using $\{0, 1\}$ -valued probabilities
- CC estimates the joint mode in a greedy way



Considering the chaining classifiers as searching a path in a binary tree whose leaf nodes are associated with a labeling $\mathbf{y} \in \mathcal{Y}$, and with 0/1 branches for y_i on level i , CC follows a single path in this tree in a greedy manner.

Simulations

- Two artificial data sets: conditionally independent and conditionally dependent data
- 10 000 training examples, 3 labels, low-dimensional problems
- Three classifiers: Binary Relevance (BR), Classifier Chains (CC), Probabilistic Classifier Chains (PCC) (all used with logistic regression as base classifier), and Bayes optimal predictions (B-O)

Results on two artificial data sets: conditionally independent (left) and conditionally dependent (right).

classifier	Hamming loss	rank loss	subset 0/1 loss	classifier	Hamming loss	rank loss	subset 0/1 loss
BR	0.4178(1.5)	0.5527(1)	0.8108(2)	BR	0.3921(2)	0.5675(1)	0.7374(3)
CC	0.4189(3)	0.5934(3)	0.8124(3)	CC	0.4308(3)	0.6930(3)	0.6100(2)
PCC	0.4178(1.5)	0.5528(2)	0.8088(1)	PCC	0.3920(1)	0.5676(2)	0.6052(1)
B-O	0.4179	0.5532	0.8088	B-O	0.3920	0.5671	0.6057

Experimental Results on Benchmark Data

- Twelve benchmark data sets taken from <http://mlkd.csd.auth.gr/multilabel.html> and <http://www.cs.waikato.ac.nz/~jmr30/#datasets>
- Five classifiers: Binary Relevance (BR), Classifier Chains (CC), Probabilistic Classifier Chains (PCC), Ensembled Classifier Chains (ECC), Ensembled Probabilistic Classifier Chains (EPCC) (all used with logistic regression as base classifier)

