

Regret Analysis for Performance Metrics in Multi-Label Classification

The Case of Hamming and Subset Zero-One Loss

**Krzysztof Dembczyński¹, Willem Waegeman², Weiwei Cheng¹,
and Eyke Hüllermeier¹**

¹Knowledge Engineering and Bioinformatics (KEBI) Lab
Philipps-Universität Marburg

²Research Unit Knowledge-based Systems (KERMIT)
Universität Gent



- Given a vector $\mathbf{x} \in \mathcal{X}$ of features, the goal is to learn a function $\mathbf{h}(\mathbf{x})$ that predicts accurately a binary vector $\mathbf{y} = (y_1, \dots, y_m) \in \mathcal{Y}$ of labels.
- **Example:** Given a news report, the goal is to learn a machine that tags the news report with relevant categories.
- **The simple solution:** solve the problem for each of the label independently.

Since the prediction is made for all labels **simultaneously**, two interesting issues appear:

- A multitude of loss functions defined over multiple labels,
- Dependence/correlation between labels.

In recent years, a plenty of algorithms has been introduced ...

- A **large** number of **loss functions** is commonly applied as performance metrics, but a concrete **connection** between a **multi-label classifier** and a **loss function** is **rarely** established.

- A **large** number of **loss functions** is commonly applied as performance metrics, but a concrete **connection** between a **multi-label classifier** and a **loss function** is **rarely** established.
- This gives implicitly the misleading impression that the **same** method can be **optimal** for **different** loss functions.

- A **large** number of **loss functions** is commonly applied as performance metrics, but a concrete **connection** between a **multi-label classifier** and a **loss function** is **rarely** established.
- This gives implicitly the misleading impression that the **same** method can be **optimal** for **different** loss functions.
- It is assumed that **performance** can be improved by taking the **label dependence** into account, but this term is used in an **intuitive** manner, without any precise **formal** definition.

- The empirical results are given **on average** without investigation under which conditions a given algorithm benefits.

- The empirical results are given **on average** without investigation under which conditions a given algorithm benefits.
- The reasons for improvements are **not distinguished**.

Let us discuss multi-label loss functions . . .

- **Hamming loss** measures the fraction of labels whose relevance is incorrectly predicted:

$$L_H(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \frac{1}{m} \sum_{i=1}^m \llbracket y_i \neq h_i(\mathbf{x}) \rrbracket,$$

- while **subset 0/1 loss** measures whether the prediction totally agrees with the true labeling:

$$L_s(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \llbracket \mathbf{y} \neq \mathbf{h}(\mathbf{x}) \rrbracket.$$

Analysis contains:

- The form of risk minimizers,
 - Whether the risk minimizers coincide in some circumstances,
 - Bound analysis,
 - Regret analysis.
-
- The analysis is simplified by assuming an unconstrained hypothesis space.

- The risk minimizer defined by:

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

- is given for the Hamming loss by:

$$h_i^*(\mathbf{x}) = \arg \max_{b \in \{0,1\}} \mathbf{P}(y_i = b | \mathbf{x}), \quad i = 1, \dots, m,$$

- while for the subset 0/1 loss by:

$$\mathbf{h}^*(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}} \mathbf{P}(\mathbf{y} | \mathbf{x}).$$

- The minimizer of the Hamming loss is the **marginal mode**, while for the subset 0/1 loss the **joint mode**.

Proposition

The Hamming loss and subset 0/1 have **the same risk minimizer**,

$$\mathbf{h}_H^*(\mathbf{x}) = \mathbf{h}_s^*(\mathbf{x}),$$

if one of the following conditions holds:

- (1) Labels Y_1, \dots, Y_m are **conditionally m -independent**,

$$\mathbf{P}(\mathbf{Y}|\mathbf{x}) = \prod_{i=1}^m \mathbf{P}(Y_i|\mathbf{x}).$$

- (2) The probability of the **joint mode** satisfies,

$$\mathbf{P}(\mathbf{h}_s^*(\mathbf{x})|\mathbf{x}) \geq 0.5.$$

Proposition

For all distributions of \mathbf{Y} given \mathbf{x} , and for all models \mathbf{h} , the expectation of the **subset 0/1** loss can be **bounded** in terms of the expectation of the **Hamming** loss as follows:

$$\frac{1}{m} \mathbb{E}_{\mathbf{Y}}[L_s(\mathbf{Y}, \mathbf{h}(\mathbf{x}))] \leq \mathbb{E}_{\mathbf{Y}}[L_H(\mathbf{Y}, \mathbf{h}(\mathbf{x}))] \leq \mathbb{E}_{\mathbf{Y}}[L_s(\mathbf{Y}, \mathbf{h}(\mathbf{x}))]$$

The previous results may suggest that one of the loss functions can be used as a proxy of the other:

- For some situations both risk minimizers coincide,
- One can provide mutual bounds for both loss functions,
- **However, in the worst case analysis, we will show that the regret is high . . .**

The **regret** of a **classifier** h with respect to a **loss function** L_z is defined as:

$$r_{L_z}(h) = \mathbb{E}_{\mathbf{X}\mathbf{Y}} L_z(\mathbf{Y}, h(\mathbf{X})) - \mathbb{E}_{\mathbf{X}\mathbf{Y}} L_z(\mathbf{Y}, h_z^*(\mathbf{X})),$$

where expectation is taken over the joint distribution $\mathbf{P}(\mathbf{X}, \mathbf{Y})$, and h_z^* is the Bayes-optimal classifier with respect to the loss function L_z .

Since both loss functions are decomposable with respect to individual instances, we analyze the expectation of \mathbf{Y} for a given x .

Proposition (Regret for subset 0/1 loss)

The following **upper bound** holds:

$$\mathbb{E}_{\mathbf{Y}} L_s(\mathbf{Y}, \mathbf{h}_H^*(\mathbf{x})) - \mathbb{E}_{\mathbf{Y}} L_s(\mathbf{Y}, \mathbf{h}_s^*(\mathbf{x})) < 0.5.$$

Moreover, this **bound is tight**, i.e.,

$$\sup_{\mathbf{P}} (\mathbb{E}_{\mathbf{Y}} L_s(\mathbf{Y}, \mathbf{h}_H^*(\mathbf{x})) - \mathbb{E}_{\mathbf{Y}} L_s(\mathbf{Y}, \mathbf{h}_s^*(\mathbf{x}))) = 0.5,$$

where the supremum is taken over all probability distributions on \mathcal{Y} .

Proposition (Regret for Hamming loss)

The following **upper bound** holds for $m > 3$:

$$\mathbb{E}_{\mathbf{Y}} L_H(\mathbf{Y}, \mathbf{h}_s^*(\mathbf{x})) - \mathbb{E}_{\mathbf{Y}} L_H(\mathbf{Y}, \mathbf{h}_H^*(\mathbf{x})) < \frac{m-2}{m+2}.$$

Moreover, this **bound is tight**, i.e.

$$\sup_{\mathbf{P}} (\mathbb{E}_{\mathbf{Y}} L_H(\mathbf{Y}, \mathbf{h}_s^*(\mathbf{x})) - \mathbb{E}_{\mathbf{Y}} L_H(\mathbf{Y}, \mathbf{h}_H^*(\mathbf{x}))) = \frac{m-2}{m+2},$$

where the supremum is taken over all probability distributions on \mathcal{Y} .

Summary:

- The risk minimizers of Hamming and subset 0/1 loss have a different structure: marginal mode vs. joint mode.
- Under specific conditions, these two types of loss minimizers are provably equivalent.
- These loss functions are mutually upper-bounded.
- Minimization of the subset 0/1 loss may cause a high regret for the Hamming loss and vice versa.

Let us empirically confirm theoretical results ...

- The simplest classifier in which a separate binary classifier $h_i(\cdot)$ is trained for each label λ_i :

$$\begin{aligned}h_i : \mathcal{X} &\rightarrow [0, 1] \\ \mathbf{x} &\mapsto y_i \in \{0, 1\}\end{aligned}$$

- It is often criticized for treating labels independently.
- However, it is still an **unbiased approach for the Hamming loss** (and other losses for which the marginal distribution is sufficient for obtaining the risk-minimizing model).

- The method reduces the problem to multi-class classification by considering each label subset $L \in \mathcal{L}$ as a distinct meta-class:

$$\begin{aligned}h : \mathcal{X} &\rightarrow [0, 1]^m \\x &\mapsto \mathbf{y} \in \{0, 1\}^m\end{aligned}$$

- It is often claimed to be a right approach to MLC, since it takes the *label dependence* into account.
- However, this approach is clearly **tailored for the subset 0/1 loss**.

The artificial data experiment:

- conditionally independent data,
- conditionally dependent data,
- non-linear data,
- low-dimensional problems with 2 or 3 labels,
- two classifiers: Binary Relevance (BR), Label Power-set (LP) with linear SVM as base learners.

Table: Results on two artificial data sets: conditionally independent (top) and conditionally dependent (down).

Conditional independence		
classifier	Hamming loss	subset 0/1 loss
BR	0.4208(\pm .0014)	0.8088(\pm .0020)
LP	0.4212(\pm .0011)	0.8101(\pm .0025)
B-O	0.4162	0.8016

Conditional dependence		
classifier	Hamming loss	subset 0/1 loss
BR	0.3900(\pm .0015)	0.7374(\pm .0021)
LP	0.4227(\pm .0019)	0.6102(\pm .0033)
B-O	0.3897	0.6029

B-O is the Bayes Optimal classifier.

Figure: Data set composed of two labels: the first label is obtained by a linear model, while the second label represents the XOR problem.

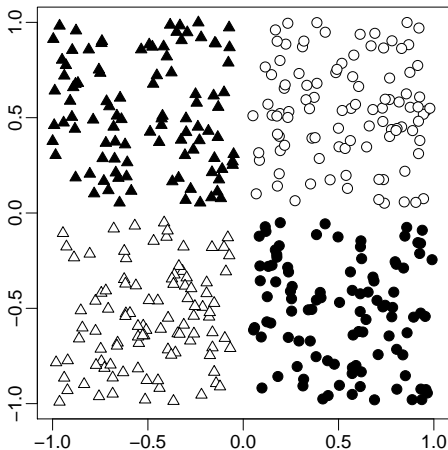


Table: Results of three classifiers on this data set.

classifier	Hamming loss	subset 0/1 loss
BR Linear SVM	0.2399(\pm .0097)	0.4751(\pm .0196)
LP Linear SVM	0.0143(\pm .0020)	0.0195(\pm .0011)
B-O	0	0

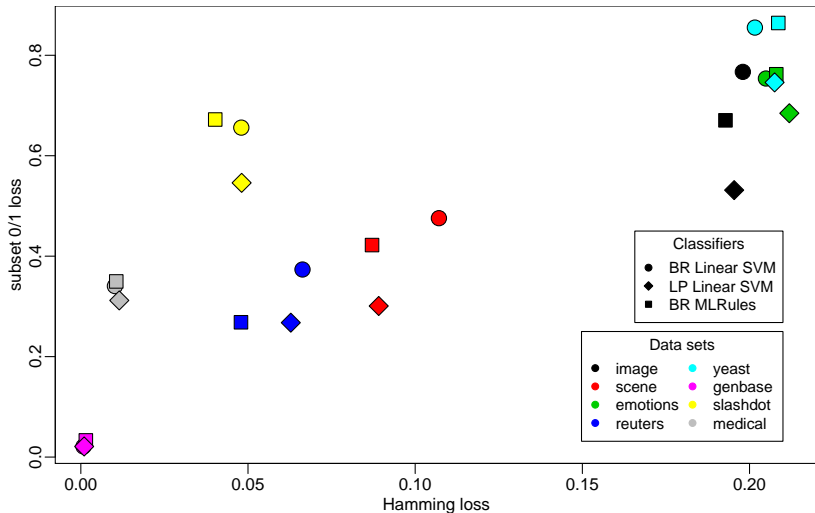
Table: Results of three classifiers on this data set.

classifier	Hamming loss	subset 0/1 loss
BR Linear SVM	0.2399(\pm .0097)	0.4751(\pm .0196)
LP Linear SVM	0.0143(\pm .0020)	0.0195(\pm .0011)
BR MLRules	0.0011(\pm .0002)	0.0020(\pm .0003)
B-O	0	0

Summary:

- LP takes the label dependence into account, but the conditional one: it is well-tailored for the subset 0/1 loss, but fails for the Hamming loss.
- LP may gain from the expansion of the feature or hypothesis space.
- One can easily tailor LP for solving the Hamming loss minimization problem, by marginalization of the joint probability distribution that is a by-product of this classifier.

Figure: Results of three classifiers on 8 benchmark data sets.



The experimental results on benchmark data confirm our claims.

The message to be taken home ...

- Surprisingly, new methods are often proposed without explicitly saying what loss they intend to minimize.
- A careful distinction between loss functions seems to be even more important for MLC than for standard classification.
- One cannot expect the same MLC method to be optimal for different types of losses.
- The reasons of improvements should be carefully distinguished.