
Optimizing the F-Measure in Multi-Label Classification: Plug-in Rule Approach versus Structured Loss Minimization

Krzysztof Dembczyński¹

Arkadiusz Jachnik¹

Wojciech Kotłowski¹

Willem Waegeman²

Eyke Hüllermeier³

¹Institute of Computing Science, Poznań University of Technology, Piotrowo 2, 60-965 Poznań, Poland

²NGDATA-Europe, Dok-Noord 7, 9000 Gent, Belgium

³Mathematics and Computer Science, Marburg University, Hans-Meerwein-Str., 35032 Marburg, Germany

KDEMBCZYNISKI@CS.PUT.POZNAN.PL

AJACHNIK@CS.PUT.POZNAN.PL

WKOTLOWSKI@CS.PUT.POZNAN.PL

WILLEMWAEGEMAN@GMAIL.COM

EYKE@MATHEMATIK.UNI-MARBURG.DE

Abstract

We compare the plug-in rule approach for optimizing the F_β -measure in multi-label classification with an approach based on structured loss minimization, such as the structured support vector machine (SSVM). Whereas the former derives an optimal prediction from a probabilistic model in a separate inference step, the latter seeks to optimize the F_β -measure directly during the training phase. We introduce a novel plug-in rule algorithm that estimates all parameters required for a Bayes-optimal prediction via a set of multinomial regression models, and we compare this algorithm with SSVMs in terms of computational complexity and statistical consistency. As a main theoretical result, we show that our plug-in rule algorithm is consistent, whereas the SSVM approaches are not. Finally, we present results of a large experimental study showing the benefits of the introduced algorithm.

1. Introduction

Motivated by applications such as tag recommendation in computer vision or gene function prediction in bioinformatics, the machine learning community has witnessed a rapid expansion of research on so-called multi-label classification (MLC) in recent years. MLC can be seen as specific type of structured output pre-

diction, and also shares commonalities with multi-task learning. The generalization from predicting a single binary label, like in conventional classification, to predicting a vector of such labels gives rise to a number of theoretical challenges; this includes the possibility to model statistical dependencies between different class labels as well as to define and minimize appropriate loss functions.

A large number of MLC loss functions has been proposed in the literature. In experimental studies, many of these losses are typically analyzed simultaneously. Yet, it is clear that a method performing optimally for one loss is likely to perform suboptimally for another loss. There are, however, a few general frameworks that can indeed be tailored for different loss functions in a generic way—an important example of such a framework is structured support vector machines (SSVMs).

For simple loss functions, analytic expressions of the Bayes (optimal) classifier can be derived. For example, it is known that the Hamming loss minimizer coincides with the marginal modes of the conditional distribution of the class labels given an instance, and methods such as binary relevance, stacking or M3L (Hariharan et al., 2012) perform particularly well in this case. Conversely, for the subset 0/1 loss, the risk minimizer is given by the joint mode of the conditional distribution, for which methods such as the label powerset classifier, conditional random fields and SSVMs without margin rescaling might be good choices.

For complex multi-label loss functions, the picture becomes more blurry, and the minimization of such losses requires more involved procedures. The F_β -measure is a specifically interesting example. Despite being

encountered in many application domains, algorithms suitable for optimizing this measure have been introduced only recently. They can be subdivided in two categories. Structured loss minimization approaches such as (Petterson & Caetano, 2010; 2011) intend to maximize the F_β -measure during the training phase in frameworks like SSVMs, whereas plug-in rule approaches (or decision-theoretic approaches) such as (Lewis, 1995; Chai, 2005; Jansche, 2007; Dembczyński et al., 2011; Quevedo et al., 2012; Ye et al., 2012) deduce F_β -measure maximizing predictions from a probabilistic model during an inference step. Let us notice that a similar distinction is considered by Ye et al. (2012); however, while our focus is on multi-label problems, their analysis is more relevant for binary classification.

As shown by Dembczyński et al. (2011), $m^2 + 1$ parameters of the conditional joint distribution (with m the number of labels) are needed to obtain the Bayes classifier for the latter type of methods. Departing from those results, we propose a novel algorithm that estimates these parameters directly via a reduction to a set of multinomial regression problems. Compared to the approach of Dembczyński et al. (2011) that constructs a model for the conditional joint distribution, our method avoids an approximation during the inference phase due to sampling from the probabilistic model. In addition, we analyze and compare the plug-in rule and structured loss minimization methods in terms of computational complexity and statistical consistency. Our main theoretical results show that our algorithm is consistent, while the SSVM approach followed in (Petterson & Caetano, 2010; 2011) is not. Finally, we present results of a large experimental study, in which we thoroughly compare both approaches.

2. Multi-Label Classification

We start with a more formal definition of the MLC problem. Let \mathcal{X} denote an instance space, and let $\mathcal{L} = \{\lambda_1, \lambda_2, \dots, \lambda_m\}$ be a finite set of class labels. We assume that an instance $\mathbf{x} \in \mathcal{X}$ is (non-deterministically) associated with a subset of labels $L \in 2^\mathcal{L}$; this subset is often called the set of relevant (positive) labels, while the complement $\mathcal{L} \setminus L$ is considered as irrelevant (negative) for \mathbf{x} . We identify a set L of relevant labels with a binary vector $\mathbf{y} = (y_1, y_2, \dots, y_m)$, in which $y_i = 1$ iff $\lambda_i \in L$. The set of possible labelings is denoted $\mathcal{Y} = \{0, 1\}^m$. We assume observations to be generated independently and randomly according to a probability distribution $P(\mathbf{X} = \mathbf{x}, \mathbf{Y} = \mathbf{y})$ (later denoted $P(\mathbf{x}, \mathbf{y})$) on $\mathcal{X} \times \mathcal{Y}$, i.e., an observation (\mathbf{x}, \mathbf{y}) is the realization of two random vectors,

$$\mathbf{X} = (X_1, X_2, \dots, X_q) \text{ and } \mathbf{Y} = (Y_1, Y_2, \dots, Y_m).$$

A multilabel classifier $\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_m(\mathbf{x}))$ assigns a (predicted) label subset to each instance $\mathbf{x} \in \mathcal{X}$. The prediction accuracy of \mathbf{h} is measured in terms of its *risk*, that is, its expected loss

$$L(\mathbf{h}, P) = \mathbb{E} [\ell(\mathbf{Y}, \mathbf{h}(\mathbf{x}))] = \int \ell(\mathbf{y}, \mathbf{h}(\mathbf{x})) dP(\mathbf{x}, \mathbf{y}), \quad (1)$$

where $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a *loss function*. In addition, it will be convenient to use an expected loss conditioned on an instance $\mathbf{x} \in \mathcal{X}$:

$$L(\mathbf{h}, P | \mathbf{x}) = \mathbb{E} [\ell(\mathbf{Y}, \mathbf{h}(\mathbf{x})) | \mathbf{x}] = \sum_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}, \mathbf{h}(\mathbf{x})) P(\mathbf{y} | \mathbf{x}),$$

so that $L(\mathbf{h}, P) = \mathbb{E}[L(\mathbf{h}, P | \mathbf{X})]$.

The optimal classifier, commonly referred to as *Bayes classifier*, minimizes the risk conditioned on \mathbf{x} :

$$\mathbf{h}^*(\mathbf{x}) = \arg \min_{\mathbf{h} \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}} \ell(\mathbf{y}, \mathbf{h}) P(\mathbf{y} | \mathbf{x}). \quad (2)$$

We note that \mathbf{h}^* is in general not unique. However, the risk of \mathbf{h}^* , denoted $L^*(P)$, is unique, and is called the *Bayes risk*.

Label subsets can be compared in terms of a multitude of loss functions, many of which lead to intractable optimization problems. Therefore, the actual loss function, often referred to as the *task loss*, is usually replaced by a *surrogate loss* that is easier to cope with (e.g., a convex upper bound of the task loss). Alternatively, the original loss minimization problem can be decomposed into problems of simpler type, for which efficient algorithmic solutions are readily available.

3. The F_β -Measure

Given a prediction $\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), \dots, h_m(\mathbf{x})) \in \mathcal{Y}$ of a binary label vector $\mathbf{y} = (y_1, \dots, y_m)$, the F_β -measure is defined as follows:

$$F_\beta(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \frac{(1 + \beta^2) \sum_{i=1}^m y_i h_i(\mathbf{x})}{\beta^2 \sum_{i=1}^m y_i + \sum_{i=1}^m h_i(\mathbf{x})} \in [0, 1],$$

where $0/0 = 1$ by definition. This measure essentially corresponds to the weighted harmonic mean of precision and recall. Compared to measures like Hamming loss, the F_β -measure enforces a better balance between performance on relevant and irrelevant labels, and, therefore, it is more suitable for multi-label problems, in which irrelevant labels often prevail.

Since the F_β -measure is a score of adjustment between \mathbf{y} and \mathbf{h} , it is actually not a loss function but rather

a kind of utility measure. Therefore, we will either consider the F_β -based loss function (F_β -loss for short)

$$\ell_{F_\beta}(\mathbf{y}, \mathbf{h}(\mathbf{x})) = 1 - F_\beta(\mathbf{y}, \mathbf{h}(\mathbf{x})), \quad (3)$$

or speak of utility maximization instead of risk minimization.

4. Bayes Classifier for the F_β -Measure

Unfortunately, there is no closed-form of the Bayes classifier that maximizes the F_β -measure. However, as shown by Dembczyński et al. (2011), the Bayes classification can be computed efficiently, even in the general case without any assumption about the underlying probability distribution.

The Bayes classifier for the F_β -measure is given by (for the sake of clarity, we will suppress dependence on \mathbf{x} in the notation, whenever it is clear from the context):

$$\begin{aligned} \mathbf{h}^* &= \arg \max_{\mathbf{h} \in \mathcal{Y}} \mathbb{E}[F_\beta(\mathbf{Y}, \mathbf{h})] = \arg \max_{\mathbf{h} \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}} F_\beta(\mathbf{y}, \mathbf{h}) P(\mathbf{y}) \\ &= \arg \max_{\mathbf{h} \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}} \frac{(1 + \beta^2) \sum_{i=1}^m y_i h_i P(\mathbf{y})}{\beta^2 \sum_{i=1}^m y_i + \sum_{i=1}^m h_i}. \end{aligned} \quad (4)$$

Problem (4) can be solved via outer and inner maximization (Jansche, 2007). Namely, it can be transformed into an inner maximization

$$\mathbf{h}_k^* = \arg \max_{\mathbf{h} \in \mathcal{H}_k} \mathbb{E}[F_\beta(\mathbf{Y}, \mathbf{h})], \quad (5)$$

where $\mathcal{H}_k = \{\mathbf{h} \in \mathcal{Y} \mid \sum_{i=1}^m h_i = k\}$, followed by an outer maximization

$$\mathbf{h}^* = \arg \max_{\mathbf{h} \in \{\mathbf{h}_0^*, \dots, \mathbf{h}_m^*\}} \mathbb{E}[F_\beta(\mathbf{Y}, \mathbf{h})]. \quad (6)$$

The outer maximization (6) can be done by simply checking all $m + 1$ possibilities. The main effort is then required for solving the inner maximization (5).

Let us consider the inner maximization problem for a given $k > 0$. We first introduce the following notation:

$$s_{\mathbf{y}} = \sum_{i=1}^m y_i, \quad \Delta_{ik}^u = \sum_{\mathbf{y}: y_i = u} \frac{P(\mathbf{y})}{\beta^2 s_{\mathbf{y}} + k}.$$

Then, we can write

$$\mathbf{h}_k^* = \arg \max_{\mathbf{h} \in \mathcal{H}_k} \sum_{\mathbf{y} \in \mathcal{Y}} \frac{(1 + \beta^2) \sum_{i=1}^m y_i h_i P(\mathbf{y})}{\beta^2 s_{\mathbf{y}} + k}. \quad (7)$$

By swapping sums in (7), we get

$$\begin{aligned} \mathbf{h}_k^* &= \arg \max_{\mathbf{h} \in \mathcal{H}_k} (1 + \beta^2) \sum_{i=1}^m h_i \sum_{\mathbf{y} \in \mathcal{Y}} \frac{y_i P(\mathbf{y})}{\beta^2 s_{\mathbf{y}} + k} \\ &= \arg \max_{\mathbf{h} \in \mathcal{H}_k} (1 + \beta^2) \sum_{i=1}^m h_i \Delta_{ik}^1. \end{aligned} \quad (8)$$

Since there are only k labels for which we can set $h_i = 1$, we obtain the optimal solution of the inner maximization by setting $h_i = 1$ for the top k values of Δ_{ik} . For the solution \mathbf{h}_k^* , we compute $\mathbb{E}[F_\beta(\mathbf{Y}, \mathbf{h}_k^*)]$, which is further used in the outer maximization. For the specific case of $\mathbf{h}_0 = \mathbf{0}$, $\mathbb{E}[F_\beta(\mathbf{Y}, \mathbf{h}_0)] = P(\mathbf{Y} = \mathbf{0})$.

Thus, we only need Δ_{ik}^1 for $1 \leq i, k \leq m$ and $P(\mathbf{0})$, that is, $m^2 + 1$ parameters to compute the optimal prediction. With these parameters, the solution can be obtained in $\mathcal{O}(m^2)$ time, i.e., the dominating part of the procedure is the inner maximization: For each k , a selection of the top k elements must be done, which can be accomplished in linear time. We will call this approach the *General F_β -Measure Maximizer* (GFM).

Under the assumption of independence of the labels Y_1, \dots, Y_m , the optimization problem (4) can be substantially simplified. Lewis (1995) and Jansche (2007) have shown independently that the optimal solution is either empty ($\mathbf{h}^* = \mathbf{0}$) or consists of those labels with the highest marginal probabilities $p_i = P(Y_i = 1)$. As a consequence, the form of the solution of the k -th inner maximization problem is known (the k labels with the highest marginal probabilities). The only missing element is to compute the value of the expected F_β -measure for a given k in order to select the best k .

Different approaches have been proposed to tackle this problem. Jansche (2007) introduced an algorithm that works in $\mathcal{O}(m^4)$ time. Chai (2005) and Quevedo et al. (2012) have independently derived $\mathcal{O}(m^3)$ algorithms based on dynamic programming. In a more recent follow-up paper, Ye et al. (2012) further improves the dynamic programming algorithm of Chai (2005) to an $\mathcal{O}(m^2)$ complexity for rational β by additional sharing of internal representations.

5. Plug-in Rule Classifier

One way to construct a classifier using training data $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ is to estimate all required parameters and then plug them into the form of the Bayes classifier. Such an approach is usually referred to as *plug-in rule classifier*. In the following, we show that all parameters required by the Bayes classifier for the F_β -measure can be efficiently obtained thanks to a specific reduction to a set of multinomial regression problems.

Let us start with the estimation of the Δ_{ik}^1 . Unfortunately, these quantities do not correspond to a proper probability distribution, i.e., Δ_{ik}^0 and Δ_{ik}^1 do not sum up to 1 or to any known constant; moreover, there is no other subset of these quantities either that may possess a property of that kind. Yet, there is a simple trick that we can use. Let us denote by \mathbf{P} and \mathbf{W} two

$m \times m$ matrices with elements

$$p_{is} = P(y_i = 1, s_y = s, | \mathbf{x}), \quad w_{ik} = (\beta^2 s + k)^{-1},$$

respectively. Then, the $m \times m$ matrix Δ with elements Δ_{ik} can be obtained by

$$\Delta = \mathbf{P} \mathbf{W}.$$

The matrix \mathbf{P} can be estimated by using a simple reduction, namely by solving m multi-class probability estimation problems (e.g., multinomial regression) with at most $m + 1$ classes. The scheme of the reduction for the i -th subproblem is the following:

$$(\mathbf{x}, \mathbf{y}) \rightarrow (\mathbf{x}, y = [\![y_i = 1]\!] \cdot s_y).$$

For a given \mathbf{x} , we estimate $P(y = [\![y_i = 1]\!] \cdot s_y | \mathbf{x})$, $y \in \{0, \dots, m\}$. We can model the probabilities $P(0 | \mathbf{x}), \dots, P(m | \mathbf{x})$ with a function $f(y, \mathbf{x})$ using the multinomial logistic transform:

$$P(y | \mathbf{x}) = \frac{\exp(f(y, \mathbf{x}))}{\sum_{j=0}^m \exp(f(j, \mathbf{x}))}$$

Then, the logistic loss has the form

$$\ell_{\log}(y, \mathbf{x}, f) = \log \left(\sum_{j=0}^m \exp(f(j, \mathbf{x})) \right) - f(y, \mathbf{x}).$$

The learning can be formulated as regularized minimization of the logistic loss:

$$f^*(y, \mathbf{x}) = \arg \min_f \frac{1}{n} \sum_{i=1}^n \ell_{\log}(y, \mathbf{x}, f) + \lambda J(f),$$

where λ controls the trade-off between the average loss over the training examples and the regularizer J that penalizes complex solutions. In a similar way, we can estimate $P(\mathbf{0} | \mathbf{x})$ by performing an additional reduction to the binary problem

$$(\mathbf{x}, \mathbf{y}) \rightarrow (\mathbf{x}, y = [\![\mathbf{y} = \mathbf{0}]\!]),$$

and solving it via logistic loss minimization.

The decomposition of the original problem into independent multinomial regression tasks has computational advantages. Moreover, since the number of distinct values of s_y is usually small, the number of classes in a single multinomial regression task is much smaller than $m + 1$; only in the worst case, we end up with a quadratic complexity in the number of labels m . Let us remark, however, that the quantities to be estimated across different tasks are not fully independent of each other (e.g., p_{im} is the same for all

i). Consequently, learning on a finite training set may lead to conflicting estimates of $P(\mathbf{0} | \mathbf{x})$ and the matrix \mathbf{P} , that is, estimates that are not in agreement with any valid distribution. To avoid such conflicts, one may include additional constraints in the learning problem or calibrate the estimates afterwards. Another alternative is to train a model that estimates the conditional joint distribution, for example, by using probabilistic classifier chains; then, estimates of the required parameters can be obtained by sampling from this distribution. This approach has been adopted by Dembczyński et al. (2011). One should note, however, that sampling usually dominates the complexity of the inference.

To summarize, the procedure for learning a probabilistic model has a time complexity that is at most quadratic in m . In the inference phase, we first obtain estimates of $P(\mathbf{0} | \mathbf{x})$ and \mathbf{P} for each test instance \mathbf{x} , again in at most quadratic time. The matrices \mathbf{P} and \mathbf{W} are multiplied to get Δ in at most cubic time. Finally, all parameters are plugged into the GFM procedure, which requires quadratic time in m . For a moderate number of labels (up to hundreds) and a small number of distinct values of s_y , this yields a feasible approach for MLC with F_β as performance measure. We refer to this method as *Exact- F_β -Plug-in classifier* (EFP).

Under the assumption of label independence, the plug-in rule approach simplifies a lot. Since we only need the marginal probabilities $p_i = P(Y_i = 1)$, we reduce the problem to m binary classification tasks. For label λ_i , the reduction takes the form

$$(\mathbf{x}, \mathbf{y}) \rightarrow (\mathbf{x}, y = y_i),$$

and we can learn a classifier in a similar way as above, namely through minimization of the logistic loss. Then, for each test instance \mathbf{x} , we obtain a vector of marginal probabilities p_i , to which we apply, for example, the method of Ye et al. (2012). This results in a procedure that is much faster than the exact one. The learning is linear in m , and the inference is quadratic for rational β or cubic in the general case. We refer to this method as *Label-independence- F_β -Plug-in classifier* (LFP). This method, however, may lead to suboptimal results if the assumption of label independence is violated. Theoretically, the difference between these two approaches can become arbitrarily large (Dembczyński et al., 2011).

6. Structured Loss Minimization

An alternative to the plug-in rule approach outlined above is to minimize the task loss directly. Since

this is intractable for the F_β -loss, one usually minimizes a convex upper bound. Here, we examine the general framework of structured hinge loss minimization (Tsochantaridis et al., 2005), mainly following the approach of Petterson & Caetano (2010; 2011).

The problem can be stated as learning a function $f(\mathbf{y}, \mathbf{x})$ such that a prediction $\mathbf{h}(\mathbf{x})$ is given by:

$$\mathbf{h}(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{y}, \mathbf{x}). \quad (9)$$

The training consists of minimizing the structural hinge loss, which can be defined by

$$\tilde{\ell}_h(\mathbf{y}, \mathbf{x}, f) = \max_{\mathbf{y}' \in \mathcal{Y}} \{\ell(\mathbf{y}, \mathbf{y}') + f(\mathbf{y}', \mathbf{x})\} - f(\mathbf{y}, \mathbf{x}), \quad (10)$$

where $\ell(\mathbf{y}, \mathbf{y}')$ is the corresponding task loss. Here, we consider the F_β -loss ℓ_{F_β} (3). The learning problem can then be stated as

$$f^*(\mathbf{y}, \mathbf{x}) = \arg \min_f \frac{1}{n} \sum_{i=1}^n \tilde{\ell}_h(\mathbf{y}_i, \mathbf{x}_i, f) + \lambda J(f),$$

where we trade-off the average hinge loss over the training examples and the regularization, similarly as for the logistic loss minimization in the plug-in rule approach. This is the so-called margin-rescaling estimator for SSVMs. The formulation leads to a quadratic program with exponentially many constraints. Therefore, one usually uses the cutting-plane algorithm (Kelley, 1960), which starts by solving the problem with no constraints and iteratively adds the most violated constraint for the current solution of the optimization problem. In each iteration, one thus needs to find

$$\mathbf{y}_i^* = \arg \max_{\mathbf{y}' \in \mathcal{Y}} \{\ell_{F_\beta}(\mathbf{y}_i, \mathbf{y}') + f(\mathbf{y}', \mathbf{x})\}. \quad (11)$$

Depending on the choice of $f(\mathbf{y}, \mathbf{x})$, one ends up with procedures of different complexity. Petterson & Caetano (2010) assume f to be decomposable over labels:

$$f(\mathbf{y}, \mathbf{x}; \mathbf{w}) = \sum_{i=1}^m f_i(y_i, \mathbf{x}). \quad (12)$$

This form of the function f leads to an effective formulation of the problem. It turns out that the constraint generation problem (11) can be solved in $\mathcal{O}(m^2)$ time. Moreover, the prediction problem (9) can be solved in linear time with respect to m :

$$h_i(\mathbf{x}) = \arg \max_{y \in \{0,1\}} f_i(y, \mathbf{x}), \quad i = 1, \dots, m.$$

We refer to this method as RML, as it was originally called in (Petterson & Caetano, 2010).¹

¹Here, we follow the interpretation of this method as given by Petterson & Caetano (2011); in the original paper, it was introduced in a different setting.

An alternative approach was introduced by Petterson & Caetano (2011), in which $f(\mathbf{y}, \mathbf{x})$ additionally models assortative (submodular) pairwise interactions between labels:

$$f(\mathbf{y}, \mathbf{x}) = \sum_{i=1}^m f_i(y_i, \mathbf{x}) + \sum_{1 \leq y_k < y_l \leq m} f_{k,l}(y_k, y_l), \quad (13)$$

where the $f_{k,l}$ are non-negative and take into account the normalized count of co-occurrences of labels λ_k and λ_l . In this model, the prediction (9) can be accomplished exactly and efficiently via graph-cuts. However, the worst-case bound for graph-cut algorithms is $\mathcal{O}(m^3)$. The constraint generation problem (11) becomes more involved, too, and requires the solution of an intractable optimization problem. The authors propose an approximate algorithm working in $\mathcal{O}(m^4)$, and they prove the non-trivial guarantee that each label which is part of the solution is also part of the optimal solution. We refer to this method as SML, as originally called in (Petterson & Caetano, 2011).

7. Statistical Consistency

The algorithms introduced in previous sections do not directly minimize the F_β -loss. Instead, they either minimize a surrogate loss or follow a reduction scheme. The performance of these algorithms can be evaluated empirically for finite samples, using synthetic or real benchmark data. However, it is also interesting to analyze their infinite sample performance, by verifying whether they converge to the Bayes classifier for the F_β -measure. This type of consistency analysis has already been performed for different learning frameworks in general (Bartlett et al., 2006; Tewari & Bartlett, 2007) and for multi-label classification in particular (Gao & Zhou, 2011), but not yet for the F_β -measure and the algorithms discussed in this paper. More formally, we will use the following definition of *statistical consistency*.

Definition 1. Given a surrogate loss $\tilde{\ell}(\mathbf{y}, f)$, a task loss $\ell(\mathbf{y}, \mathbf{h})$, and a prediction function $\mathbf{h} = \mathbf{h}(f)$, we say that $\tilde{\ell}$ is consistent with respect to ℓ , if for any distribution P and any sequence of classifiers f_1, f_2, \dots such that

$$\lim_{n \rightarrow \infty} \tilde{L}(f_n, P) = \min_f \tilde{L}(f, P),$$

we have

$$\lim_{n \rightarrow \infty} L(\mathbf{h}(f_n), P) = \min_{\mathbf{h}} L(\mathbf{h}, P),$$

where $\tilde{L}(f, P) = \mathbb{E}[\tilde{\ell}(\mathbf{Y}, f(\mathbf{X}))]$.

There is an equivalent condition for consistency which is sometimes more convenient to use:

Theorem 1 (Gao & Zhou (2011)). *For any distribution P , let $A(P) = \arg \min_{\mathbf{h}} L(\mathbf{h}, P)$ be the set of all Bayes classifiers for loss function ℓ . A surrogate loss $\tilde{\ell}$ is consistent with respect to ℓ if and only if for all P :*

$$\min_f \tilde{L}(f, P) < \inf_f \{ \tilde{L}(f, P) : \mathbf{h}(f) \notin A(P) \} \quad (14)$$

Although the above definition and theorem directly refer to the surrogate loss, we shall show later that they can also be applied to the plug-in rule approach. Starting with the analysis of the structured hinge loss, we find that both methods proposed in (Peterson & Caetano, 2010; 2011) are inconsistent, as stated formally in the theorem below. Since the feature vector \mathbf{x} does not play any role in the following, we omit the dependence on \mathbf{x} for the rest of the section.²

Theorem 2. *Consider the surrogate loss $\tilde{\ell}$ defined as in (10):*

$$\tilde{\ell}(\mathbf{y}, f) = \max_{\mathbf{y}' \in \mathcal{Y}} \{ \ell_{F_\beta}(\mathbf{y}, \mathbf{y}') + f(\mathbf{y}') \} - f(\mathbf{y}), \quad (15)$$

for any $\mathbf{y} \in \mathcal{Y}$, and any real valued function $f: \mathcal{Y} \rightarrow \mathbb{R}$. Let $\mathbf{h}(f) = \arg \max_{\mathbf{y}} f(\mathbf{y})$ be a label vector predicted by f . Then $\tilde{\ell}$ is inconsistent with respect to ℓ_{F_β} . The inconsistency remains even if additional structure given by (12) or (13) is imposed on the function f .

Proof. For the sake of clarity, we show the theorem for $\beta = 1$, but a small modification of the proof will suffice to cover any $\beta > 0$. Let us choose $m = 2$, and the joint label distribution $p_{00} = 0.5$, $p_{01} = 0.2$, $p_{10} = 0$, $p_{11} = 0.3$, where $p_{uv} = P(y_1 = u, y_2 = v)$. Given a prediction \mathbf{h} and a true label vector \mathbf{y} , the F_1 -loss (3) is summarized by the following table:

$\mathbf{h} = (0, 0)$	$\mathbf{h} = (0, 1)$	$\mathbf{h} = (1, 0)$	$\mathbf{h} = (1, 1)$
$\mathbf{y} = (0, 0)$	0	1	1
$\mathbf{y} = (0, 1)$	1	0	$1/3$
$\mathbf{y} = (1, 0)$	1	1	$1/3$
$\mathbf{y} = (1, 1)$	1	$1/3$	0

Thus, the Bayes classifier for the F_1 -loss is $(0, 0)$, whence $A(P) = \{(0, 0)\}$. Now, for any f , the expected surrogate loss (15), $\tilde{L}(f, P)$, is given by

$$\begin{aligned} & p_{00}(\max\{0 + f_{00}, 1 + f_{01}, 1 + f_{10}, 1 + f_{11}\} - f_{00}) \\ & + p_{01}(\max\{1 + f_{00}, 0 + f_{01}, 1 + f_{10}, \frac{1}{3} + f_{11}\} - f_{01}) \\ & + p_{11}(\max\{1 + f_{00}, \frac{1}{3} + f_{01}, \frac{1}{3} + f_{10}, 0 + f_{11}\} - f_{11}), \end{aligned}$$

where we denote $f_{uv} = f((u, v))$ and used the fact that $p_{10} = 0$. Now, we show that for any choice of f , the

²One may assume that everything is presented conditionally for a given \mathbf{x} .

surrogate loss is at least 1. Since we can bound max from below by any of its terms, we can write

$$\begin{aligned} \tilde{L}(f, P) & \geq p_{00}(1 + f_{01} - f_{00}) + p_{01}(1 + f_{00} - f_{01}) \\ & + p_{11}(1 + f_{00} - f_{11}) = 1 + 0.3(f_{01} - f_{11}) \end{aligned}$$

On the other hand,

$$\begin{aligned} \tilde{L}(f, P) & \geq p_{00}(1 + f_{11} - f_{00}) + p_{01}(1 + f_{00} - f_{01}) \\ & + p_{11}(1 + f_{00} - f_{11}) = 1 - 0.2(f_{01} - f_{11}). \end{aligned}$$

The two inequalities above imply that $\tilde{L}(f, P) \geq 1$ for any f . Now, if we choose f , such that $f_{00} = -0.1$, $f_{01} = 0.1$, $f_{10} = -0.1$, $f_{11} = 0.1$, we have $\tilde{L}(f, P) = 1$, which means that f is a Bayes classifier for $\tilde{\ell}$. However, $\mathbf{h}(f) \in \{(0, 1), (1, 1)\}$, and thus $\mathbf{h}(f) \notin A(P)$, which violates (14). Moreover, it is easy to see that f as constructed above satisfies (12) and (13) if we set $f_1(y_1) = 0$, $f_2(y_2) = 0.1(2y_2 - 1)$, and $f_{k,l} \equiv 0$, so that inconsistency even holds under these constraints. \square

Despite the fact that the plug-in rule approach defined in Section 5 may deliver conflicting estimates of \mathbf{P} and $P(\mathbf{0} | \mathbf{x})$ for finite training data, we can show that it is consistent if the sample size grows to infinity.

Theorem 3. *Consider the surrogate loss $\tilde{\ell}(\mathbf{y}, f)$ obtained for the reduction to independent multinomial regression tasks:*

$$\tilde{\ell}(\mathbf{y}, \mathbf{f}) = \ell_{\log}(\llbracket s_{\mathbf{y}} = 0 \rrbracket, f^0) + \sum_{i=1}^m \ell_{\log}(s_{\mathbf{y}} \llbracket y_i = 1 \rrbracket, f^i),$$

where $\mathbf{f} = (f^0, \dots, f^m)$ is a vector of real-valued functions $f^i: \{0, \dots, m\} \rightarrow \mathbb{R}$. Let $\mathbf{h}(\mathbf{f})$ be the prediction from the GFM procedure, where all probabilities are obtained from \mathbf{f} using the logistic transform as described in Section 5. Then $\tilde{\ell}$ is consistent with respect to ℓ_{F_β} .

Proof. Consider a standard multinomial logistic loss $\ell_{\log}(y, f)$ for multi-class classification, and let P_f be a probability estimate obtained from f through a logistic transform. It is easy to verify that $\mathbb{E}[\ell_{\log}(Y, f)] = H(P) + D(P \| P_f)$, where $H(P)$ is the entropy of the true distribution P , and $D(P \| P_f)$ is the relative entropy between P and P_f (Cover & Thomas, 1991). Thus, given any sequence f_1, f_2, \dots such that $\mathbb{E}[\ell_{\log}(Y, f_n)] \rightarrow \min_f \mathbb{E}[\ell_{\log}(Y, f)]$, we must have $D(P \| P_{f_n}) \rightarrow 0$, which, according to the Pinsker inequality, implies $P_{f_n} \rightarrow P$ in a total-variation sense.

We prove the theorem directly from Definition 1. Consider any sequence $\mathbf{f}_1, \mathbf{f}_2, \dots$ such that $\tilde{L}(\mathbf{f}_n, P) \rightarrow \min_f \tilde{L}(\mathbf{f}, P)$. Since $\tilde{L}(\mathbf{f}_n, P)$ is a sum of expected logistic losses, convergence of $\tilde{L}(\mathbf{f}_n, P)$ to its minimum implies convergence of each expected logistic

loss to its respective minimum. This in turn implies that all probability estimates $P_{f_n^i}$ converge to the true probabilities P in a total-variation sense. Since each term under $\arg \max$ in (6) is a continuous function of the probability estimates (and there are only a finite number of such terms), $\mathbf{h}(\mathbf{f}_n)$ will eventually be in $A(P)$ for sufficiently large n , i.e., the plug-in classifier will eventually agree with the Bayes classifier for the F_β -loss. This implies that $L_{F_\beta}(\mathbf{h}(\mathbf{f}_n), P) \rightarrow \min_{\mathbf{h}} L_{F_\beta}(\mathbf{h}, P)$. \square

8. Experimental Results

We complement our theoretical analysis by an empirical evaluation of the methods on finite data sets. More specifically, we compare two plug-in rule methods, namely the Exact- F_β -Plug-in classifier (EFP) and the Label-independence- F_β -Plug-in classifier (LFP), and two structured loss minimization methods, namely RML and SML.³ As an additional baseline, we include the so-called binary relevance (BR) approach that learns and predicts for each label independently. This algorithm essentially corresponds to LFP without the inference phase. We consider two performance measures, F_1 and Hamming loss, as well as running times of the training and inference procedures.

8.1. Setting

All approaches included in the comparison use base functions that are linear in the feature space. We train BR, EFP and LFP by using regularized multinomial regression.⁴ We tune the regularization parameter λ for each base classifier independently by minimizing the logistic loss, which should provide better probability estimates. We use 5-fold cross-validation and choose λ from $\{10^{-4}, 10^{-3}, \dots, 10^3\}$.

Similarly, RML has a single parameter λ , and we tune it in 5-fold cross-validation using the same range of values. SML has an additional parameter c that determines the trade-off between the linear f_i and the label interaction terms $f_{j,k}$ in (13). To guarantee a fair comparison, we only tune λ using 5-fold cross-validation for each experiment, whereas the best value of the parameter c was selected in an earlier series of experiments. Moreover, we use 5% of the label pairs for all datasets—according to [Petterson & Caetano \(2011\)](#), the results with other settings are very similar. The

³For both RML and SML, we use the implementation offered by the authors ([Petterson & Caetano, 2011](#)) available at <http://users.cecs.anu.edu.au/~jpetterson/>.

⁴We use the implementation of multinomial regression by Mallet ([McCallum, 2002](#)) available at <http://mallet.cs.umass.edu/>.

maximal number of iterations in the cutting-plane algorithm is set to 1000 for both RML and SML.

It needs to be mentioned that the plug-in rule approaches are implemented in Java, while the RML and SML are C++ programs. Therefore, the evaluation times may not be fully comparable. We run the experiments on a Debian virtual machine with 8-core x64 processor and 5GB RAM.⁵ We use 6 benchmark datasets, which are publicly available from Mulan.⁶ Table 1 provides a summary of some basic statistics of these datasets.

Table 1. Datasets and their properties: the number of training (#TRAIN) and test (#TEST) examples, the number of labels (m) and features (d), and the average number of classes (q) in a single multinomial regression task of EFP.

DATASET	#TRAIN	#TEST	m	d	q
IMAGE	1200	800	5	135	4
SCENE	1211	1196	6	294	3.33
YEAST	1500	917	14	103	10.36
MEDICAL	333	645	45	1449	2.466
ENRON	1123	579	53	1001	6.58
MEDIAMILL	30993	12914	101	120	11.41

8.2. Results

The results are summarized in Table 2. As can be seen, all methods tailored for the F_1 -measure outperform the baseline on this measure, whereas BR achieves the best results for Hamming loss. This is coherent with the result of [Gao & Zhou \(2011\)](#), according to which this approach is consistent for this loss function. With respect to the F_1 -measure, the best method is EFP that wins on five out of six datasets. To some extent, this result can be explained by our theoretical results regarding the consistency of the methods.

In general, it seems that the plug-in rule methods are superior to the methods based on structured loss minimization. Among the latter, RML outperforms SML.⁷ This is in agreement with the original results of [Petterson & Caetano \(2011\)](#), suggesting that SML achieves better results only for datasets with a reduced number of features. One may conjecture that the label interaction terms are more helpful in properly modeling the feature space than in optimizing the F_β -measure.

⁵Remark that all methods can be easily parallelized. All subtasks of LFP and EFP can be trained independently. The implementation of RML and SML is also multi-threaded ([Petterson & Caetano, 2011](#)). We used, however, the sequential variants of the algorithms to simplify the comparison.

⁶<http://mulan.sourceforge.net/datasets.html>

⁷The results we obtained for RML are at least as good as the ones presented in ([Petterson & Caetano, 2011](#)), while for SML are slightly worse.

Table 2. Experimental results for Hamming loss (HL), F_1 , and running times (in seconds) of cross-validation (t_{cv}), training (t_{train}) (for the best set of parameters) and inference (t_{inf}). The best results are marked by a '*'.

	HL[%]	F [%]	t_{cv}	t_{train}	t_{inf}	HL[%]	F [%]	t_{cv}	t_{train}	t_{inf}	HL[%]	F [%]	t_{cv}	t_{train}	t_{inf}
	IMAGE					SCENE					YEAST				
BR	*19.90	43.63	*9	*0.392	0.087	10.51	55.73	*29	*0.733	0.241	*20.03	60.59	*26	*0.901	0.128
LFP	27.55	58.86	*9	*0.392	0.119	12.18	74.38	*29	*0.733	0.270	22.24	65.02	*26	*0.901	0.146
EFP	26.07	*59.77	24	0.606	0.183	12.22	*74.44	72	0.995	0.399	22.82	*65.47	101	2.004	0.367
RML	25.07	57.49	94	1.104	*0.051	*9.70	73.92	73	1.001	*0.118	22.82	64.78	206	5.194	*0.056
SML	28.82	56.99	156	7.116	0.052	15.65	68.50	52	1.129	0.123	24.52	63.96	319	4.385	0.070
	MEDICAL					ENRON					MEDIAMILL				
BR	*1.17	70.19	*9	*1	0.952	*4.54	55.49	*52	*4	1.016	*3.19	51.21	*3238	*118	13
LFP	1.18	*81.27	*9	*1	1.513	6.09	56.86	*52	*4	1.519	3.67	55.15	*3238	*118	20
EFP	1.23	80.39	16	1	1.883	5.34	*61.04	214	6	2.628	3.63	*55.16	24620	440	30
RML	1.20	80.63	1253	30	*0.144	6.35	57.69	3897	41	*0.143	4.12	49.35	—	1125	*7
SML	2.50	67.90	715	23	0.773	7.82	54.61	18780	62	0.887	4.18	50.02	—	10365	131

Table 2 also shows the runtimes for parameter tuning in cross-validation (t_{cv}), training for the best set of parameters, and inference. As we mentioned above, the comparison of the running times should be interpreted with caution, due to the use of different programming languages and differences in the implementations. For example, the inference times for BR and RML should basically be very similar, as in both cases we apply m linear models. Yet, the implementation of RML is much more efficient. Nevertheless, we are still able to derive several important conclusions.

RML is most efficient in inference, which is coherent with our analysis. Nevertheless, the inference times of the plug-in rule approaches are quite comparable to those of BR, despite their quadratic (for LFP) and cubic (for EFP) complexity. Admittedly, however, the datasets used in the experiments only contain a small to moderate number of labels (up to 100). For datasets with thousands of labels, the difference is likely to become substantially larger. The inference for SML is slower than for RML, but for all datasets except one, it is faster than the plug-in rule classifiers. For the MEDIAMILL dataset, this method takes the longest time.

The training of BR and LFP (these are exactly the same procedures) is the most effective. Training of EFP leads to m multinomial regression models. One should note, however, that the average number of classes (column q in Table 1) for all datasets is much smaller than the highest possible value $m + 1$. Therefore, the training of EFP is still quite effective and takes only a few times longer than the training of LFP. The cutting-plane algorithm and the constraint generation step slow down the training of RML and SML, and SML performs worst in this regard. Still, the cutting-plane algorithm converges very fast in several cases.

The plug-in rule approaches may tune the parameters internally for each subtask independently, without explicitly running the inference step. Therefore, the tuning times of LFP and EFP are also better than those of RML and SML. For the MEDIAMILL dataset, RML and SML were not able to perform the parameter tuning step in a reasonable amount of time. To get the results for this dataset, we trained these methods with different settings and selected the best result on the test set. We also reduced the maximum number of iterations to 200.

9. Conclusion

We discussed and analyzed two conceptually different approaches to F_β -measure maximization (F_β -loss minimization) in multi-label classification. The plug-in rule methods estimate all parameters that are needed to compute the prediction of the Bayes classifier, whereas methods based on structured loss minimization, such as structured SVMs, produce a classifier more directly through minimization of the loss on the training data.

Moreover, we introduced a novel plug-in rule algorithm that performs parameter estimation via a set of multinomial regression tasks. Theoretically, we have shown this algorithm to be consistent for the F_β -measure, whereas the SSVM approach is not consistent. This result is corroborated by our experimental studies, in which the plug-in rule approach performs particularly well.

Acknowledgments. The first three authors are supported by the Foundation of Polish Science under the Homing Plus programme, co-financed by the European Regional Development Fund. The last author is supported by German Research Foundation. We thank all anonymous reviewers for their valuable comments.

References

Bartlett, P., Jordan, M., and McAuliffe, J. Convexity, classification and risk bounds. *Journal of the American Statistical Association*, 101:138–156, 2006.

Chai, K. Expectation of F-measures: Tractable exact computation and some empirical observations of its properties. In *SIGIR*, 2005.

Cover, T. M. and Thomas, J. A. *Elements of Information Theory*. John Wiley, 1991.

Dembczyński, K., Waegeman, W., Cheng, W., and Hüllermeier, E. An exact algorithm for F-measure maximization. In *Advances in Neural Information Processing Systems*, volume 25, 2011.

Gao, W. and Zhou, Z. On the consistency of multi-label learning. In *COLT*, 2011.

Hariharan, B., Vishwanathan, S. V. N., and Varma, M. Efficient max-margin multi-label classification with applications to zero-shot learning. *Machine Learning Journal*, 88(1):127–155, 2012.

Jansche, M. A maximum expected utility framework for binary sequence labeling. In *ACL 2007*, pp. 736–743, 2007.

Kelley, J. E. The cutting-plane method for solving convex programs. *Journal of the Society for Industrial Applied Mathematics*, 8:704–712, 1960.

Lewis, D. Evaluating and optimizing autonomous text classification systems. In *SIGIR 1995*, pp. 246–254, 1995.

McCallum, A. Mallet: A machine learning for language toolkit. <http://mallet.cs.umass.edu>, 2002.

Petterson, J. and Caetano, T. S. Reverse multi-label learning. In *Advances in Neural Information Processing Systems 24*, pp. 1912–1920, 2010.

Petterson, J. and Caetano, T. S. Submodular multi-label learning. In *Advances in Neural Information Processing Systems 24*, pp. 1512–1520, 2011.

Quevedo, J., Luaces, O., and Bahamonde, A. Multilabel classifiers with a probabilistic thresholding strategy. *Pattern Recognition*, 45, 2012.

Tewari, A. and Bartlett, P. On the consistency of multi-class classification methods. *Journal of Machine Learning Research*, 8:1007–1025, May 2007.

Tsochantaridis, I., Joachims, T., Hofmann, T., and Altun, Y. Large margin methods for structured and interdependent output variables. *J. Mach. Learn. Res.*, 6: 1453–1484, 2005.

Ye, N., Chai, K., Lee, W., and Chieu, H. Optimizing F-measures: a tale of two approaches. In *ICML*, 2012.