

CONSISTENT ALGORITHMS FOR MULTI-LABEL CLASSIFICATION WITH MACRO-AT- k METRICS

Erik Schultheis

Aalto University
Helsinki, Finland
erik.schultheis@aalto.fi

Wojciech Kołowski

Poznan University of Technology
Poznan, Poland
wkotlowski@cs.put.poznan.pl

Marek Wydmuch

Poznan University of Technology
Poznan, Poland
mwydmuch@cs.put.poznan.pl

Rohit Babbar

University of Bath / Aalto University
Bath, UK / Helsinki, Finland
rb2608@bath.ac.uk

Strom Borman

Yahoo Research
Champaign, USA
strom.borman@yahooinc.com

Krzysztof Dembczyński

Yahoo Research / Poznan University of Technology
New York, USA / Poznan, Poland
krzysztof.dembczynski@yahooinc.com

ABSTRACT

We consider the optimization of complex performance metrics in multi-label classification under the population utility framework. We mainly focus on metrics linearly decomposable into a sum of binary classification utilities applied separately to each label with an additional requirement of exactly k labels predicted for each instance. These “macro-at- k ” metrics possess desired properties for extreme classification problems with long tail labels. Unfortunately, the at- k constraint couples the otherwise independent binary classification tasks, leading to a much more challenging optimization problem than standard macro-averages. We provide a statistical framework to study this problem, prove the existence and the form of the optimal classifier, and propose a statistically consistent and practical learning algorithm based on the Frank-Wolfe method. Interestingly, our main results concern even more general metrics being non-linear functions of label-wise confusion matrices. Empirical results provide evidence for the competitive performance of the proposed approach.

1 INTRODUCTION

Various real-world applications of machine learning require performance measures of a complex structure, which, unlike misclassification error, do not decompose into an expectation over instance-wise quantities. Examples of such performance measures include the area under the ROC curve (AUC) (Drummond & Holte, 2005), geometric mean (Drummond & Holte, 2005; Wang & Yao, 2012; Menon et al., 2013; Cao et al., 2019), the F -measure (Lewis, 1995) or precision at the top (Kar et al., 2015). The theoretical analysis of such measures, as well as the design of consistent and efficient algorithms for them, is a non-trivial task.

In multi-label classification, one can consider a wide spectrum of measures that are usually divided into three categories based on the averaging scheme, namely instance-wise, micro, and macro averaging. Instance-wise measures are defined, as the name suggests, on the level of a single instance. Typical examples are Hamming loss, precision@ k , recall@ k , and the instance-wise F -measure. Micro-averages are defined on a confusion matrix that accumulates true positives, false positives, false negative, and true negatives from all the labels. Macro-averages require a binary metric to be applied to each label separately and then averaged over the labels. In general, any binary metric can be applied in any of the above averaging schemes. Not surprisingly, some of the metrics, for example misclassification error, lead to the same form of the final metric regardless of the scheme

used. One can also consider the wider class of measures that are defined as general aggregation functions of label-wise confusion matrices. This includes the measures described above, but also, e.g., the geometric mean of label-wise metrics or a specific variant of the F -measure (Opitz & Burst 2019) being a harmonic mean of macro-precision and macro-recall.

In this paper, we target the setting of prediction with a budget. Specifically, we require the predictions to be “budgeted-at- k ,” meaning that for each instance, exactly k labels need to be predicted. The budget of k requires the prediction algorithm to choose the labels “wisely”. It is also important in many real-world scenarios. For instance, in recommendation systems or extreme classification, there is a fixed number of slots (e.g., indicated by a user interface) required to be filled with related products/searches/ads (Cremonesi et al. 2010; Chang et al. 2021). Furthermore, having a fixed prediction budget is also interesting from a methodological perspective, as various metrics which lead to degenerate solutions without a budget, e.g., predict nothing (macro-precision) or everything (macro-recall), become meaningful when restricted to predict k labels per instance.

While all our theoretical results and algorithms apply to a general class of multi-label measures, we focus in this paper on macro-averaged metrics. If no additional requirements are imposed on the classifier, the linear nature of the macro-averaging means that a binary problem for each label can be solved independently, and existing techniques (Koyejo et al. 2015; Kotłowski & Dembczyński 2017) are sufficient. In turn, if we require predictions to be budgeted-at- k , the task becomes much more difficult, as this constraint tightly couples the different binary problems together. In general, they cannot be solved independently for each label, requiring instead more involved techniques to find the optimal classifier.

The macro-at- k metrics seem to be very attractive in the context of multi-label classification. Macro-averaging treats all the labels equally important. This prevents ignoring labels with a small number of positive examples (Schultheis et al. 2022), so-called tail labels, which are very common in applications of multi-label classification, particularly in the extreme setting when the number of all labels is very large (Jain et al. 2016; Babbar & Schölkopf 2019). Furthermore, it can be shown that one can remove tail labels from the training set with almost no drop of performance in terms of popular metrics, such as precision@ k and nDCG@ k , on extreme multi-label data sets (Wei & Li 2019; Schultheis et al. 2023). The macro-at- k metrics, on the other hand, are sensitive to the lack of tail labels in the training set.¹

We aim at delivering consistent algorithms for macro-at- k metrics, i.e., algorithms that converge in the limit of infinite training data to the optimal classifier for the metrics. Our main theoretical results are stated in a very general form, concerning the large class of aggregation functions of label-wise confusion matrices. Our starting point of the analysis are results obtained in the multi-class setting (Narasimhan et al. 2015; 2022), concerning consistent algorithms for complex performance measures with additional constraints. Nevertheless, they do not consider budgeted-at- k predictions, which do not apply to multi-class classification, while they play an important role in the multi-label setting. Furthermore, using arguments from functional analysis, we managed to significantly simplify the line of reasoning in the proofs. We first show that the problem can be transformed from optimizing over classifiers to optimizing over the set of feasible confusion matrices, and that the optimal classifier optimizes an unknown *linear* confusion-matrix metric. In the multi-label setting, interestingly, such a classifier corresponds to a prediction rule, which has the appealingly simple form: selecting the k highest-scoring labels based on an *affine transformation* of the marginal label probabilities. Combining this result with the optimization of confusion matrices, we state a Frank-Wolfe based algorithm that is consistent for finding the optimal classifier also for *nonlinear* metrics. Empirical studies provide evidence that the proposed approach can be applied in practical settings and obtains competitive performance in terms of the macro-at- k metrics.

2 RELATED WORK

The problem of optimizing complex performance metrics is well-known, with many articles published for a variety of metrics and different classification problems. It has been considered for binary (Ye et al. 2012; Koyejo et al. 2014; Busa-Fekete et al. 2015; Dembczynski et al. 2017), multi-class (Narasimhan et al. 2015; 2022), multi-label (Waegeman et al. 2014; Koyejo et al. 2015; Kotłowski & Dembczyński 2017), and multi-output (Wang et al. 2019) classification.

¹Results and description of such an experiment are given in Appendix I

Initially, the main focus was on designing algorithms, without a conscious emphasis on statistical consequences of choosing models and their asymptotic behavior. Notable examples of such contributions are the SVMperf algorithm (Joachims, 2005), approaches suited to different types of the F-measure (Dembczynski et al., 2011; Natarajan et al., 2016; Jasinska et al., 2016), or precision at the top (Kar et al., 2015). Wide use of such complex metrics has caused an increasing interest in investigating their theoretical properties, which can then serve as a guide to design practical algorithms.

The consistency of learning algorithms is a well-established problem. The seminal work of Bartlett et al. (2006) was studying this problem for binary classification under the misclassification error. Since then a wide spectrum of learning problems and performance metrics has been analyzed in terms of consistency. These results concern ranking (Duchi et al., 2010; Ravikumar et al., 2011; Calauzenes et al., 2012; Yang & Koyejo, 2020), multi-class (Zhang, 2004; Tewari & Bartlett, 2007) and multi-label classification (Koyejo et al., 2015; Kołowski & Dembczyński, 2017), classification with abstention (Yuan & Wegkamp, 2010; Ramaswamy et al., 2018), or constrained classification problems (Agarwal et al., 2018; Kearns et al., 2018; Narasimhan et al., 2022). Nevertheless, the problem of designing consistent algorithms for budgeted-at- k macro averages is relatively new.

Optimizing non-decomposable metrics can be considered in two distinct frameworks (Dembczynski et al., 2017): population utility (PU) and expected test utility (ETU). The PU framework focuses on estimation, in the sense that a consistent PU classifier is one which correctly estimates the population optimal utility as the size of the training set increases. A consistent ETU classifier is one which optimizes the expected prediction error over a *given* test set. The latter might get better results, as the optimization is performed on the test set directly. Optimization of budgeted-at- k metrics in this framework has been recently considered in Schultheis et al. (2023). The former framework, which we focus on in this paper, has the advantage that prediction can be made for each test example separately, without knowing the entire test set in advance.

3 PROBLEM STATEMENT

Let $\mathbf{x} \in \mathcal{X}$ denote an input instance, and $\mathbf{y} \in \{0, 1\}^m$ the vector indicating the relevant labels, jointly distributed according to $(\mathbf{x}, \mathbf{y}) \sim \mathbb{P}$. Let $\mathbf{h}: \mathcal{X} \rightarrow [0, 1]^m$ be a *randomized multi-label classifier* which, given instance \mathbf{x} , predicts a possibly randomized class label vector $\hat{\mathbf{y}} \in \{0, 1\}^m$, such that $\mathbb{E}_{\hat{\mathbf{y}}|\mathbf{x}}[\hat{\mathbf{y}}] = \mathbf{h}(\mathbf{x})$. We assume that the predictions are *budgeted at k* , that is exactly k labels are always predicted as relevant, which means that $\|\hat{\mathbf{y}}\|_1 = \sum_{j=1}^m \hat{y}_j = k$ with probability 1. It turns out that this is *equivalent* to assuming $\|\mathbf{h}(\mathbf{x})\|_1 = \sum_{j=1}^m h_j(\mathbf{x}) = k$ for all $\mathbf{x} \in \mathcal{X}$. Indeed, $\|\mathbf{h}(\mathbf{x})\|_1 = k$ is *necessary*, because $k = \mathbb{E}_{\hat{\mathbf{y}}|\mathbf{x}}[\|\hat{\mathbf{y}}\|_1] = \|\mathbf{h}(\mathbf{x})\|_1$; but it also *suffices* as for any real-valued vector $\boldsymbol{\pi} \in [0, 1]^m$ with $\|\boldsymbol{\pi}\|_1 = k$, one can construct a distribution over binary vectors $\hat{\mathbf{y}} \in \{0, 1\}^m$ with $\|\hat{\mathbf{y}}\|_1 = k$ and marginals $\mathbb{E}_{\hat{\mathbf{y}}|\mathbf{x}}[\hat{\mathbf{y}}] = \boldsymbol{\pi}$; this can be accomplished using, e.g., *Madow's sampling scheme* (see Appendix A for the actual efficient algorithm). Thus, using notation $\Delta_m^k := \{\mathbf{v} \in [0, 1]^m : \|\mathbf{v}\|_1 = k\}$, the randomized classifiers budgeted at k are then all (measurable) functions of the form $\mathbf{h}: \mathcal{X} \rightarrow \Delta_m^k$. We denote the set of such functions as \mathcal{H} .

For any $\mathbf{x} \in \mathcal{X}$, let $\boldsymbol{\eta}(\mathbf{x}) := \mathbb{E}_{\mathbf{y}|\mathbf{x}}[\mathbf{y}]$ denote the vector of conditional label marginals. Given a randomized classifier $\mathbf{h} \in \mathcal{H}$, we define its *multi-label confusion tensor* $\mathbf{C}(\mathbf{h}) = (\mathbf{C}^1(h_1), \dots, \mathbf{C}^m(h_m))$ as a sequence of m binary classification confusion matrices associated with each label $j \in [m]$, that is $C_{uv}^j(h_j) = \mathbb{P}[y_j = u, \hat{y}_j = v]$ for $u, v \in \{0, 1\}$. Note that using the marginals and the definition of the randomized classifier,

$$C^j(h_j) = \begin{pmatrix} \mathbb{E}_{\mathbf{x}}[(1 - \eta_j(\mathbf{x}))(1 - h_j(\mathbf{x}))] & \mathbb{E}_{\mathbf{x}}[(1 - \eta_j(\mathbf{x}))h_j(\mathbf{x})] \\ \mathbb{E}_{\mathbf{x}}[\eta_j(\mathbf{x})(1 - h_j(\mathbf{x}))] & \mathbb{E}_{\mathbf{x}}[\eta_j(\mathbf{x})h_j(\mathbf{x})] \end{pmatrix}. \quad (1)$$

The set of all possible binary confusion matrices is written as $\mathcal{C} := \{\mathbf{C} \in [0, 1]^{2 \times 2} \mid \|\mathbf{C}\|_{1,1} = 1\}$, and is used to define the set of possible confusion tensors for predictions at k through $\mathcal{C}^k := \{\mathbf{C} \in [0, 1]^{m \times 2 \times 2} \mid \forall j \in [m] : C^j \in \mathcal{C}, \sum_{j=1}^m C_{01}^j + C_{11}^j = k\}$

In this work, we are interested in optimizing performance metrics that do not decompose over individual instances, but are general functions of the confusion tensor of the classifier \mathbf{h} . While in general, given two confusion matrices, we cannot say which one is better than the other without knowing the specific application, it is possible to impose a *partial* order that any reasonable performance metric should respect. To that end, define:

Table 1: Examples of binary confusion matrix measures, which can be used as building blocks for confusion tensor measures. For clarity, we denote $tn = C_{00}$, $fp = C_{01}$, $fn = C_{10}$, $tp = C_{11}$.

Metric	$\psi(\mathbf{C})$	Metric	$\psi(\mathbf{C})$
Accuracy	$tp + tn$	Recall	$\frac{tp}{tp+fn}$
Precision	$\frac{tp}{tp+fp}$	Balanced accuracy	$\frac{tp}{2(tp+fn)} + \frac{tn}{2(tn+fp)}$
F_β	$\frac{(1+\beta^2)tp}{(1+\beta^2)tp+\beta^2fn+fp}$	G-Mean	$\sqrt{\frac{tp \cdot tn}{(tp+fn)(tn+fp)}}$
Jaccard	$\frac{tp}{tp+fp+fn}$	AUC	$\frac{2 \cdot tp \cdot tn + tp \cdot fp + fn \cdot tn}{2(tp+fn)(fp+tn)}$

Definition 3.1 (Binary Confusion Matrix Measure). *Let $\mathcal{C} = \{\mathbf{C} \in [0, 1]^{2 \times 2} \mid \|\mathbf{C}\|_{1,1} = 1\}$ be the set of all possible binary confusion matrices, and $\mathbf{C}, \mathbf{C}' \in \mathcal{C}$. Then we say that \mathbf{C}' is at least as good as \mathbf{C} , $\mathbf{C}' \succeq \mathbf{C}$, if there exists constants ϵ_1, ϵ_2 such that*

$$\mathbf{C}' = \begin{pmatrix} C_{00} + \epsilon_1 & C_{01} - \epsilon_1 \\ C_{10} - \epsilon_2 & C_{11} + \epsilon_2 \end{pmatrix}, \quad (2)$$

i.e., if \mathbf{C}' can be generated from \mathbf{C} by turning some false positives to true negatives and false negatives to true positives. A function $\psi: \mathcal{C} \rightarrow [0, 1]$ is called a binary confusion matrix measure (Singh & Khim, 2022) if it respects that ordering, i.e., if for $\mathbf{C}' \succeq \mathbf{C}$ we have $\psi(\mathbf{C}') \geq \psi(\mathbf{C})$.

Similarly, in the multi-label case we cannot compare arbitrary confusion tensors, where one is better on some labels than on others,² but we can recognize if one is better on *all* labels:

Definition 3.2 (Confusion Tensor Measure). *For a given number of labels $m \in \mathbb{N}$, and two confusion tensors $\mathbf{C}, \mathbf{C}' \in \mathcal{C}^m$, we say that \mathbf{C}' is at least as good as \mathbf{C} , $\mathbf{C}' \succeq \mathbf{C}$, if for all labels $j \in [m]$ it holds that $\mathbf{C}'^{j'} \succeq \mathbf{C}^j$. A function $\Psi: \mathcal{C}^m \rightarrow [0, 1]$ is called a confusion tensor measure if it respects this ordering, i.e., if for $\mathbf{C}' \succeq \mathbf{C}$ we have $\Psi(\mathbf{C}') \geq \Psi(\mathbf{C})$.*

Of particular interest in this paper are functions which linearly decompose over the labels, that is *macro-averaged multi-label metrics* (Manning et al., 2008; Parambath et al., 2014; Koyejo et al., 2015; Kotłowski & Dembczyński, 2017) of the form:

$$\Psi(\mathbf{h}) = \Psi(\mathbf{C}(\mathbf{h})) = m^{-1} \sum_{j=1}^m \psi(\mathbf{C}^j(h_j)), \quad (3)$$

where ψ is some binary confusion matrix measure. If one takes a binary confusion matrix measure (e.g., any of those define in Table 1), then the resulting macro-average will be a valid confusion tensor measure. A more thorough discussion of these conditions can be found in Appendix H.

Macro-averaged metrics find numerous applications in multi-label classification, mainly due to their balanced emphasis across labels independent of their frequencies, and thus can potentially alleviate the “long-tail” issues in problems with many rare labels (Schultheis et al., 2022).

Denote the optimal value of the metric among all classifiers budgeted at k as:

$$\Psi^* := \sup_{\mathbf{h} \in \mathcal{H}} \Psi(\mathbf{h}), \quad (4)$$

and let $\mathbf{h}^* \in \operatorname{argmax}_{\mathbf{h}} \Psi(\mathbf{h})$ be an optimal (Bayes) classifier for which $\Psi(\mathbf{h}^*) = \Psi^*$, if it exists. For any classifier \mathbf{h} , define its Ψ -regret as $\Delta\Psi(\mathbf{h}) = \Psi^* - \Psi(\mathbf{h})$ to measure the suboptimality of \mathbf{h} with respect to Ψ : from the definition, $\Delta\Psi(\mathbf{h}) \geq 0$ for every classifier \mathbf{h} , and $\Delta\Psi(\mathbf{h}) = 0$ if and only if \mathbf{h} is optimal. If the Ψ -regret of a learning algorithm converges to zero with the sample size tending to infinity, it is called (statistically) consistent. We consider such algorithms in Section 5. Even though the objective (3) decomposes onto m binary problems, these are still coupled by the budget constraint, $\|\mathbf{h}(\mathbf{x})\|_1 = k$ for all $\mathbf{x} \in \mathcal{X}$, and cannot be optimized independently as we show later in the paper.

²This is specifically the trade-off we want to achieve for tail labels!

4 THE OPTIMAL CLASSIFIER

Finding the form of the optimal classifier for general macro-averaged performance metrics is difficult. For instance, when $\psi(\mathbf{C})$ is the F_β measure, the objective to be optimized is a sum of linear fractional functions, which is known to be NP-hard in general (Schaible & Shi, 2003). We are, however, able to fully determine the optimal classifier for the specific class of *linear utilities*, which are metrics depending linearly on the confusion tensor of the classifier. Furthermore, we also show that for a general class of macro-averaged metrics, under mild assumptions on the data distribution, the optimal classifier exists and turns out to also be the maximizer of some linear utility, whose coefficients, however, depend on its (unknown a priori) confusion tensor.

We start with a metric of the form $\Psi(\mathbf{C}) = \mathbf{G} \cdot \mathbf{C} = \sum_{j=1}^m \mathbf{G}^j \cdot \mathbf{C}^j$ for some vector of *gain matrices* (gain tensor) $\mathbf{G} = (\mathbf{G}^1, \dots, \mathbf{G}^m)$, possibly depending on the data distribution \mathbb{P} . We call such a utility *linear* as it linearly depends on the confusion matrices of the classifier. Note that we allow the gain matrix \mathbf{G} to be different for each label, making this more general than linear macro-averages. We need to consider this more general case, because it will appear as a subproblem when finding optimal predictions for non-linear macro-averages as presented below.

Linear metrics are decomposable over instances, and thus the optimal classifier has an appealingly simple form: It boils down to simply sorting the labels by an affine function of the marginals, and returning the top k elements.

Theorem 4.1. *The optimal classifier $\mathbf{h}^* := \operatorname{argmax}_{\mathbf{h} \in \mathcal{H}} \Psi(\mathbf{h})$ for $\Psi(\mathbf{h}) = \mathbf{G} \cdot \mathbf{C}(\mathbf{h})$ is given by*

$$\mathbf{h}^*(\mathbf{x}) = \operatorname{top}_k(\mathbf{a} \odot \boldsymbol{\eta}(\mathbf{x}) + \mathbf{b}), \quad (5)$$

where \odot denotes the coordinate-wise product of vectors, while the vectors \mathbf{a} and \mathbf{b} are given by:

$$a_j = G_{00}^j + G_{11}^j - G_{01}^j - G_{10}^j, \quad b_j = G_{01}^j - G_{00}^j, \quad (6)$$

and $\operatorname{top}_k(\mathbf{v})$ returns a k -hot vector extracting the top k largest entries of \mathbf{v} (ties broken arbitrarily).

Proof (sketch, full proof in Appendix B). After simple algebraic manipulations, the objective can be written as $\Psi(\mathbf{h}) = \mathbb{E} \left[\sum_{j=1}^m (a_j \eta_j(\mathbf{x}) + b_j) h_j(\mathbf{x}) \right] + R$, where a_j and b_j are as stated in the theorem, while R does not depend on the classifier. For each $\mathbf{x} \in \mathcal{X}$, the objective can thus be independently maximized by the choice of $\mathbf{h}(\mathbf{x}) \in \Delta_m^k$ which maximizes $\sum_{j=1}^m (a_j \eta_j(\mathbf{x}) + b_j) h_j(\mathbf{x})$. But this is achieved by sorting $a_j \eta_j(\mathbf{x}) + b_j$ in descending order, and setting $h_j(\mathbf{x}) = 1$ for the top k coordinates, and $h_j(\mathbf{x}) = 0$ for the remaining coordinates (with ties broken arbitrarily). \square

Examples of binary metrics for which their macro averages can be written in the linear form include:

- the accuracy $\psi(\mathbf{C}) = C_{00} + C_{11}$ (which leads to the *Hamming utility* after macro-averaging) with $a_j = 2, b_j = -1$, and thus for any $\mathbf{x} \in \mathcal{X}$, the optimal prediction $\mathbf{h}^*(\mathbf{x})$ returns k labels with the largest marginals $\eta_j(\mathbf{x})$;
- the same prediction rule is obtained for the *TP* metric $\psi(\mathbf{C}) = C_{00}$ (that leads to *precision@k*) with $a_j = 1, b_j = 0$;
- the recall $\psi(\mathbf{C}) = \mathbb{P}(y = 1)^{-1} C_{11}$ (macro-averaged to *recall@k*) has $a_j = \mathbb{P}(y_j = 1)^{-1}, b_j = 0$, so that the optimal classifiers returns top k labels sorted according to $\frac{\eta_j(\mathbf{x})}{\mathbb{P}(y_j=1)}$;
- the balanced accuracy $\psi(\mathbf{C}) = \frac{C_{11}}{2\mathbb{P}(y=1)} + \frac{C_{00}}{2\mathbb{P}(y=0)}$, gives $a_j = \frac{1}{2\mathbb{P}(y_j=1)} + \frac{1}{2\mathbb{P}(y_j=0)}, b_j = -\frac{1}{2\mathbb{P}(y_j=0)}$, with the optimal prediction sorting labels according to $\frac{\eta_j(\mathbf{x})}{\mathbb{P}(y_j=1)} - \frac{1-\eta_j(\mathbf{x})}{1-\mathbb{P}(y_j=1)}$.

We now switch to general case, in which the base binary metrics are not necessarily decomposable over instances, and optimizing their macro averages with budgeted predictors is a challenging task. We make the following mild assumptions on the data distribution and performance metric:

³We use $\mathbf{A} \cdot \mathbf{B} = \sum_{uv} A_{uv} B_{uv}$ to denote a dot product over matrices, and a concise notation $\mathbf{A} \cdot \mathbf{B} = \sum_j \mathbf{A}^j \cdot \mathbf{B}^j$ for ‘dot product’ over matrix sequences $\mathbf{A} = (\mathbf{A}^1, \dots, \mathbf{A}^m)$ and $\mathbf{B} = (\mathbf{B}^1, \dots, \mathbf{B}^m)$.

Assumption 4.2. The label conditional marginal vector $\boldsymbol{\eta}(\mathbf{x}) = \mathbb{E}_{\mathbf{y}|\mathbf{x}}[\mathbf{y}]$ is absolutely continuous with respect to the Lebesgue measure on $[0, 1]^m$ (i.e., has a density over $[0, 1]^m$).

A similar assumption was commonly used in the past works (Koyejo et al., 2014; Narasimhan et al., 2015; Dembczynski et al., 2017).

Assumption 4.3. The performance metric Ψ is differentiable and fulfills for all labels $j \in [m]$

$$\left. \frac{\partial}{\partial \epsilon} \Psi \left(\mathbf{C}^1, \dots, \mathbf{C}^j + \epsilon \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}, \dots, \mathbf{C}^m \right) \right|_{\epsilon=0} > 0. \quad (7)$$

Assumption 4.3 is essentially a ‘strictly monotonic and differentiable’ version of Definition 3.2 and is satisfied by all macro-averaged metrics given in Table 1.

Our first main result concerns the form of the optimal classifier for general confusion tensor measures, of which macro-averaged binary confusion matrix measures are special cases. To state the result, we define $\mathcal{C}_{\mathbb{P}} := \{\mathbf{C}(\mathbf{h}) : \mathbf{h} \in \mathcal{H}\}$, the set of confusion tensors achievable by randomized k -budgeted classifiers on distribution \mathbb{P} . Clearly, maximizing $\Psi(\mathbf{h})$ over $\mathbf{h} \in \mathcal{H}$ is equivalent to maximizing $\Psi(\mathbf{C})$ over $\mathbf{C} \in \mathcal{C}_{\mathbb{P}}$.

Theorem 4.4. Let the data distribution \mathbb{P} and metric Ψ satisfy Assumption 4.2 and Assumption 4.3 respectively. Then, there exists an optimal $\mathbf{C}^* \in \mathcal{C}_{\mathbb{P}}$, that is $\Psi(\mathbf{C}^*) = \Psi^*$. Moreover, any classifier \mathbf{h}^* maximizing the linear utility $\mathbf{G} \cdot \mathbf{C}(\mathbf{h})$ over $\mathbf{h} \in \mathcal{H}$ with $\mathbf{G} = (\mathbf{G}^1, \dots, \mathbf{G}^m)$ given by $\mathbf{G}^j = \nabla_{\mathbf{C}^j} \Psi(\mathbf{C}^*)$, also maximizes $\Psi(\mathbf{h})$ over $\mathbf{h} \in \mathcal{H}$.

Proof (sketch, full proof in Appendix C) We first prove that $\mathcal{C}_{\mathbb{P}}$ is a compact set by using certain properties of continuous linear operators in Hilbert space. Due to continuity of Ψ and the compactness of $\mathcal{C}_{\mathbb{P}}$, there exists a maximizer $\mathbf{C}^* = \operatorname{argmax}_{\mathbf{C} \in \mathcal{C}_{\mathbb{P}}} \Psi(\mathbf{C})$. By the first order optimality and convexity of $\mathcal{C}_{\mathbb{P}}$, $\nabla \Psi(\mathbf{C}^*) \cdot \mathbf{C}^* \geq \nabla \Psi(\mathbf{C}^*) \cdot \mathbf{C}$ for all $\mathbf{C} \in \mathcal{C}_{\mathbb{P}}$, so \mathbf{C}^* maximizes a linear utility $\mathbf{G} \cdot \mathbf{C}^*$ with gain matrices given by $\mathbf{G} = \nabla \Psi(\mathbf{C}^*)$. A careful analysis under Assumption 4.2 shows that \mathbf{C}^* uniquely maximizes $\mathbf{G} \cdot \mathbf{C}$ over $\mathbf{C} \in \mathcal{C}_{\mathbb{P}}$. \square

Theorem 4.4 reveals that Ψ -optimal classifier exists and can be found by maximizing a linear utility, that is, by predicting the top k labels sorted according to an affine function of the label marginals: $\mathbf{h}^*(\mathbf{x}) = \operatorname{top}_k(\mathbf{a}^* \odot \boldsymbol{\eta}(\mathbf{x}) + \mathbf{b}^*)$ for vectors \mathbf{a}^* and \mathbf{b}^* defined for gain matrices $\mathbf{G} = \nabla \Psi(\mathbf{C}^*)$ as in Theorem 4.1. Unfortunately, since \mathbf{C}^* is unknown in advance, the coefficients \mathbf{a}^* , \mathbf{b}^* are also unknown, and the optimal classifier is not directly available. However, knowing that \mathbf{h}^* optimizes a linear utility induced by the gradient of Ψ leads to a consistent algorithm described in the next section.

Although the optimal solution is expressed by affine functions of label marginals, in general, it cannot be obtained by solving the problem independently for each label, i.e., the values of a_j and b_j may depend on labels other than j . Let $\mathbf{h}^*(\mathbf{x})$ and $\mathbf{h}'^*(\mathbf{x})$ be optimal for distributions \mathbb{P} and \mathbb{P}' , respectively. Let \mathbb{P}' differ from \mathbb{P} only on a single label j . If we could solve the problem independently for each label, then $\mathbf{h}^*(\mathbf{x})$ and $\mathbf{h}'^*(\mathbf{x})$ would be the same up to label j , in the sense that the ordering relation between all other labels would not change. In Appendix E we show that this is not the case, presenting a simple counterexample showing that a different distribution on a single label changes the solution with respect to the other labels.

5 CONSISTENT ALGORITHMS

As any algorithm we propose has to operate on a finite sample, we need to introduce empirical counterparts for our quantities of interest. For example, we use $\hat{\boldsymbol{\eta}}(\mathbf{x})$ to denote the estimate of $\boldsymbol{\eta}(\mathbf{x})$ given by a label probability estimator trained on some set of training data $\mathcal{S} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$. We also define the empirical multi-label confusion tensor $\hat{\mathbf{C}}(\mathbf{h}, \mathcal{S})$ of a classifier \mathbf{h} with respect to some set \mathcal{S} of n instances. In this case, we have:

$$\hat{\mathbf{C}}_{uv}^j(h_j, \mathcal{S}) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}[y_{ij} = u, h_j(\mathbf{x}_i) = v]. \quad (8)$$

Following Narasimhan et al. (2015), we use the Frank-Wolfe algorithm (Frank & Wolfe, 1956) to perform an implicit optimization over feasible confusion tensors $\mathcal{C}_{\mathbb{P}}$, without having to explicitly

construct \mathcal{C}_P . This is possible, because Frank-Wolfe only requires us to be able to solve two sub-problems: First, given a classifier h , we need to calculate its empirical confusion tensor, which is straight-forward. Second, given a classifier and its corresponding confusion tensor, we need to solve a *linearized* version of the optimization problem, which is possible due to [Theorem 4.1](#).

Consequently, our algorithm, presented in [Algorithm 1](#), proceeds as follows: In the beginning, we split the available training data into two subsets. One for estimating label probabilities $\hat{\eta}$, and one for tuning the actual classifier. After determining $\hat{\eta}$, we initialize h to be the standard top-k classifier, which will get iteratively refined as follows. For the confusion tensor of the current classifier, we can determine a linear objective based on its gradient. Because we can linearly interpolate stochastic classifiers, which will lead to linearly interpolated confusion tensors, this gives us a descent direction over which we can optimize a step-size⁴ and the confusion tensor at this classifier. Based on this confusion tensor, we can do the next linearized optimization step, until we reach a fixed limit for the iteration count. We represent the randomized classifier as a set of deterministic classifiers h^i , and corresponding sampling weights α^i obtained across all iterations of the algorithm. The Frank-Wolfe algorithm scales to the larger problems as it only requires $\mathcal{O}(nm)$ time per iteration.

Algorithm 1 Multi-label Frank-Wolfe algorithm for complex performance measures

Require: Dataset $\mathcal{S} := \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$, number of iterations $t \in \mathbb{N}$, stopping condition $\epsilon \in \mathbb{R}$

- 1: Split dataset \mathcal{S} into \mathcal{S}_1 and \mathcal{S}_2
- 2: Learn label marginals model $\hat{\eta} : \mathcal{X} \rightarrow \mathbb{R}^m$ on \mathcal{S}_1
- 3: Initialize $h^0 : \mathcal{X} \rightarrow \mathcal{Y}_k$ ▷ Initial deterministic classifier
- 4: Initialize $\alpha^0 \leftarrow 1$ ▷ Initial probability of selecting the initial classifier h^0
- 5: $\mathbf{C}^0 \leftarrow \mathbf{C}(h^0, \mathcal{S}_2)$ ▷ Calculate the initial confusion tensor
- 6: **for** $i \in \{1, \dots, t\}$ **do** ▷ Perform t iterations
- 7: $\mathbf{G}^i \leftarrow \nabla_{\mathbf{C}} \Psi(\mathbf{C}^{i-1})$ ▷ Calculate tensor of gradients of \mathbf{C}^{i-1} in respect to Ψ (gain tensor)
- 8: $\mathbf{a}^i \leftarrow \mathbf{G}_{11}^i + \mathbf{G}_{00}^i - \mathbf{G}_{01}^i - \mathbf{G}_{10}^i$, $\mathbf{b}^i \leftarrow \mathbf{G}_{01}^i - \mathbf{G}_{00}^i$
- 9: $h^i(\mathbf{x}) \leftarrow \text{top}_k(\mathbf{a}^i \odot \hat{\eta}(\mathbf{x}) + \mathbf{b}^i)$ ▷ Construct the next classifier h^i
- 10: $\mathbf{C}' \leftarrow \mathbf{C}(h^i, \mathcal{S}_2)$ ▷ Calculate the confusion tensor of the next classifier h^i
- 11: $\alpha^i \leftarrow \arg\max_{\alpha \in [0,1]} \Psi((1-\alpha)\mathbf{C}^{i-1} + \alpha\mathbf{C}')$ ▷ Find the best combination of \mathbf{C}^{i-1} and \mathbf{C}' (step-size)
- 12: **if** $\alpha^i < \epsilon$ **then break** ▷ Stop if the step-size α^i is smaller then ϵ
- 13: $\mathbf{C}^i \leftarrow (1 - \alpha^i)\mathbf{C}^{i-1} + \alpha^i\mathbf{C}'$ ▷ Calculate a new confusion tensor based on the best α^i
- 14: **for** $j \in \{0, \dots, i-1\}$ **do** ▷ Update all the previous
- 15: $\alpha^j \leftarrow \alpha^j(1 - \alpha^i)$ ▷ probabilities of selecting corresponding h
- 16: **return** $(\{h^0, \dots, h^i\}, \{\alpha^0, \dots, \alpha^i\})$ ▷ Return randomized classifier

This algorithm can consistently optimize a confusion tensor measure if it fulfills certain conditions:

Theorem 5.1 (Consistency of Frank-Wolfe). *Assume the utility function $\Psi : [0, 1]^{m \times 2 \times 2} \rightarrow \mathbb{R}_{\geq 0}$ is concave over \mathcal{C}_P , L -Lipschitz, and β -smooth w.r.t. the 1-norm. Let $\mathcal{S} = (\mathcal{S}_1, \mathcal{S}_2)$ be a sample drawn i.i.d. from \mathbb{P} . Further, let $\hat{\eta}$ be a label probability estimator learned from \mathcal{S}_1 , and h_S^{FW} be the classifier obtained after κn iterations. Then, for any $\delta \in (0, 1]$, with probability of at least $1 - \delta$ over draws of \mathcal{S} ,*

$$\Delta \Psi(h_S^{\text{FW}}) \leq \mathcal{O}(\mathbb{E}_{\mathbf{x}}[\|\boldsymbol{\eta}(\mathbf{x}) - \hat{\boldsymbol{\eta}}(\mathbf{x})\|_1]) + \tilde{\mathcal{O}}\left(m^2 \sqrt{\frac{m \cdot \log m \cdot \log n - \log \delta}{n}}\right) + \frac{8\beta m}{\kappa n + 2}. \quad (9)$$

The proof of this theorem, given in [Appendix D](#), broadly follows [\(Narasimhan et al., 2015\)](#): First, we show that *linear* metrics can be estimated consistently with a regret growing with the L_1 error of the LPE. Then, we prove a uniform convergence result for estimating the multi-label confusion tensor. As a prerequisite, we derive the VC-dimension of the class of classifiers based on top-k scoring, i.e., those classifiers that minimize some linear confusion tensor metric as shown in [Theorem 4.1](#).

Lemma 5.2 (VC dimension for linear top-k classifiers). *For $\boldsymbol{\eta} : \mathcal{X} \rightarrow [0, 1]^m$, define*

$$\mathcal{H}_{\boldsymbol{\eta}}^j := \bigcup_{\mathbf{a}, \mathbf{b} \in \mathbb{R}^m} \{h : \mathcal{X} \rightarrow \{0, 1\} : h(\mathbf{x}) = \mathbb{1}[j \in \text{top}_k(\mathbf{a} \odot \boldsymbol{\eta} + \mathbf{b})]\}. \quad (10)$$

The VC-complexity of this class is $\text{VC}(\mathcal{H}_{\boldsymbol{\eta}}^j) \leq 6m \log(em)$.

⁴The classical version of FW uses a fixed step-size schedule of $\frac{2}{\epsilon+1}$ instead of an inner optimization, but we find the latter to give better results empirically. However, for the convergence result, fixed steps are assumed.

Table 2: Results of different inference strategies on measure calculated at $\{3, 5, 10\}$. Notation: P—precision, R—recall, F1—F1-measure. The green color indicates cells in which the strategy matches the metric. The best results are in **bold** and the second best are in *italic*. We additionally report basic statistics of the benchmarks: number of labels and instances in train and test sets, and average number of positive labels per instance, average number of positive instances per label.

Inference strategy	Instance @3		Macro @3			Instance @5		Macro @5			Instance @10		Macro @10		
	P	R	P	R	F1	P	R	P	R	F1	P	R	P	R	F1
MEDIAMILL ($m = 101, n_{\text{train}} = 30993, n_{\text{test}} = 12914, \mathbb{E}[\ y\ _1] = 4.36, \mathbb{E}[y \times n_{\text{train}}] = 1338.8$)															
TOP-K	66.25	49.55	8.96	4.81	4.95	<i>51.96</i>	<i>62.04</i>	12.85	8.75	7.71	33.63	76.60	11.46	19.68	11.28
TOP-K+ w^{POW}	57.36	42.51	15.31	11.84	<i>10.54</i>	47.68	56.62	13.00	17.37	<i>12.64</i>	32.18	72.98	9.64	29.43	<i>13.07</i>
TOP-K+ w^{LOG}	39.72	27.32	14.43	10.10	9.41	35.40	39.96	11.38	15.33	10.95	28.45	63.36	9.86	26.25	12.26
TOP-K+ ℓ_{FOCAL}	65.87	49.60	10.08	4.87	4.94	52.08	62.16	11.99	8.93	7.90	<i>33.61</i>	<i>76.65</i>	10.76	20.08	11.37
TOP-K+ ℓ_{ASYM}	65.88	49.48	10.31	4.58	4.80	51.55	61.87	11.10	8.50	7.48	33.54	76.75	10.73	19.55	11.16
MACRO-P _{FW}	7.94	6.13	19.33	6.06	2.87	6.99	8.96	17.29	8.79	3.17	6.02	14.14	17.38	17.24	5.23
MACRO-R _{PRIOR}	6.37	3.67	8.81	19.82	5.31	7.38	7.25	8.91	26.50	6.71	8.31	17.42	10.53	39.24	8.85
MACRO-R _{FW}	6.37	3.67	8.81	19.82	5.31	7.38	7.25	8.91	26.50	6.71	8.31	17.42	10.53	39.24	8.85
MACRO-F1 _{FW}	45.20	33.05	<i>15.42</i>	11.17	12.21	43.57	51.60	<i>15.20</i>	15.05	13.82	28.12	64.23	<i>13.93</i>	23.32	14.81
FLICKR ($m = 195, n_{\text{train}} = 56359, n_{\text{test}} = 24154, \mathbb{E}[\ y\ _1] = 1.34, \mathbb{E}[y \times n_{\text{train}}] = 412.6$)															
TOP-K	23.94	56.96	23.04	38.41	26.56	16.99	66.01	17.12	47.03	23.49	10.16	77.35	10.72	59.37	17.24
TOP-K+ w^{POW}	22.35	53.44	17.96	44.26	24.21	16.10	62.80	13.76	52.39	20.68	9.77	74.54	9.08	63.98	15.08
TOP-K+ w^{LOG}	23.57	56.17	19.86	41.36	25.49	16.76	65.21	15.05	49.75	22.00	<i>10.06</i>	76.63	9.79	61.80	16.10
TOP-K+ ℓ_{FOCAL}	<i>23.64</i>	<i>56.27</i>	24.90	36.67	26.42	<i>16.89</i>	<i>65.62</i>	18.53	45.67	<i>24.16</i>	10.05	76.63	11.77	57.90	<i>18.14</i>
TOP-K+ ℓ_{ASYM}	23.37	55.65	23.09	37.00	26.12	16.74	65.04	17.39	45.61	23.60	10.06	76.63	10.91	58.36	17.48
MACRO-P _{FW}	4.65	11.49	39.34	6.63	8.06	5.66	22.75	41.74	9.70	10.57	2.83	22.26	37.59	10.68	8.50
MACRO-R _{PRIOR}	16.14	38.62	17.58	45.50	22.27	12.17	47.48	13.98	53.83	19.72	7.89	60.42	9.57	64.66	15.07
MACRO-R _{FW}	16.14	38.62	17.58	45.50	22.27	12.17	47.48	13.98	53.83	19.72	7.89	60.42	9.57	64.66	15.07
MACRO-F1 _{FW}	17.59	41.60	35.28	29.28	29.43	12.22	47.31	<i>34.13</i>	32.70	29.43	5.92	45.77	<i>34.55</i>	33.08	29.02
RCV1X ($m = 2456, n_{\text{train}} = 623847, n_{\text{test}} = 155962, \mathbb{E}[\ y\ _1] = 4.80, \mathbb{E}[y \times n_{\text{train}}] = 1218.6$)															
TOP-K	72.99	75.32	13.06	4.67	5.43	52.30	81.96	12.77	7.61	7.64	32.98	89.70	11.35	14.75	10.28
TOP-K+ w^{POW}	65.99	69.11	18.58	12.78	13.09	48.48	77.18	14.69	17.66	<i>13.64</i>	31.43	87.14	10.63	26.05	12.82
TOP-K+ w^{LOG}	70.70	73.37	19.97	8.10	9.80	51.18	80.49	16.03	11.75	11.29	32.66	<i>89.14</i>	11.96	19.01	12.06
TOP-K+ ℓ_{FOCAL}	<i>71.99</i>	<i>74.38</i>	14.06	4.83	5.76	<i>51.46</i>	<i>80.94</i>	12.49	7.65	7.75	32.38	88.75	10.59	14.42	10.06
TOP-K+ ℓ_{ASYM}	71.14	73.60	14.40	5.44	6.46	50.81	80.13	12.27	8.52	8.41	31.88	87.85	9.64	15.16	10.03
MACRO-P _{FW}	46.36	50.11	<i>21.11</i>	5.61	5.84	29.40	49.81	<i>21.69</i>	5.72	5.31	19.45	60.40	<i>21.66</i>	6.03	5.78
MACRO-R _{PRIOR}	44.26	46.10	14.60	<i>18.24</i>	12.04	34.77	56.28	13.13	<i>24.59</i>	12.77	24.08	70.51	10.66	<i>34.34</i>	12.39
MACRO-R _{FW}	43.28	44.99	14.56	18.41	11.95	34.15	55.24	13.15	24.89	12.73	23.78	69.71	10.76	34.66	12.44
MACRO-F1 _{FW}	58.20	61.22	21.45	10.37	<i>12.09</i>	44.42	71.86	21.96	12.25	13.68	27.26	78.88	22.10	14.86	15.12
AMAZONCAT ($m = 13330, n_{\text{train}} = 1186239, n_{\text{test}} = 306782, \mathbb{E}[\ y\ _1] = 5.04, \mathbb{E}[y \times n_{\text{train}}] = 448.6$)															
TOP-K	78.29	59.29	35.73	12.44	16.52	63.63	74.54	46.43	32.72	35.06	39.16	85.18	39.52	51.69	40.39
TOP-K+ w^{POW}	66.32	49.76	50.21	45.79	45.70	57.12	67.49	44.85	53.78	46.30	37.31	82.20	30.13	63.53	37.15
TOP-K+ w^{LOG}	72.56	<i>54.56</i>	50.30	32.06	36.94	<i>61.15</i>	<i>71.83</i>	48.93	42.87	<i>43.05</i>	<i>38.71</i>	<i>84.49</i>	36.84	56.71	<i>40.60</i>
MACRO-P _{FW}	47.00	35.57	<i>56.47</i>	23.74	29.62	41.04	50.74	55.85	27.45	30.23	30.66	69.67	55.27	29.09	34.51
MACRO-R _{PRIOR}	48.58	34.93	37.16	59.97	<i>42.02</i>	40.67	47.35	28.17	66.98	35.75	28.06	62.91	17.62	73.98	25.04
MACRO-R _{FW}	48.58	34.93	37.15	<i>59.97</i>	<i>42.02</i>	40.67	47.35	28.17	66.98	35.75	28.06	62.91	17.62	73.98	25.04
MACRO-F1 _{FW}	68.59	51.49	56.75	34.68	40.90	55.73	65.60	56.62	36.40	<i>41.92</i>	35.30	78.34	<i>54.67</i>	39.93	43.26

6 EXPERIMENTS

In this section, we empirically evaluate the proposed Frank-Wolfe algorithm on a variety of multi-label benchmark tasks that differ substantially in the number of labels and imbalance of the label distribution: MEDIAMILL (Snoek et al., 2006), FLICKR (Tang & Liu, 2009), RCV1X (Lewis et al., 2004), and AMAZONCAT (McAuley & Leskovec, 2013; Bhatia et al., 2016). For the first three datasets we use a multi-layer neural network for estimating $\hat{\eta}(x)$. For the last and largest dataset, we use a sparse linear label tree model, which is a common baseline in extreme multi-label classification (Jasinska-Kobus et al., 2020).⁵ In Appendix F we include all the details regarding the setup of probability estimators.

We evaluate the following classifiers optimizing the macro-at- k measures:

- MACRO-P_{FW}, MACRO-R_{FW}, MACRO-F1_{FW}: randomized classifiers found by the Frank-Wolfe algorithm (Algorithm 1) for optimizing macro precision, recall, and F_1 , respectively, based on $\hat{\eta}(x)$ coming from the model trained with binary cross-entropy loss (BCE).

⁵Code to reproduce the experiments: <https://github.com/mwydmuch/xCOLUMNS>

- **MACRO- R_{PRIOR}** : implements the optimal strategy for macro recall, which selects k labels with the highest $\hat{p}_j^{-1}\hat{\eta}_j$; \hat{p}_j s are estimates of label priors obtained on a training set and $\hat{\eta}(\mathbf{x})$ are given by the model trained with BCE loss.

As baselines, we use the following algorithms:

- **TOP-K**: selects k labels with the highest $\hat{\eta}_j$ coming from the model trained with BCE loss; the optimal strategy for instance-wise precision at k (Wydmuch et al., 2018).
- **TOP-K+ w^{POW}** , **TOP-K+ w^{LOG}** : similarly to TOP-K, selects k labels with the highest $w_j\hat{\eta}_j$, where w_j are calculated as a function of label priors corresponding to the power-law, $w_j^{\text{POW}} = \hat{p}_j^{-\beta}$, and log weights, $w_j^{\text{LOG}} = -\log \hat{p}_j$, with \hat{p} estimated on the training set. For power-law weights, we use $\beta = 0.5$. This kind of weighting aims to put more emphasis on less frequent labels.
- **TOP-K+ ℓ_{FOCAL}** , **TOP-K+ ℓ_{ASYM}** : multi-label focal loss and asymmetric loss (Lin et al., 2017; Ridnik et al., 2021) are variants of BCE loss, commonly used in multi-label classification to improve classification performance on harder, less frequent labels. Here, we train models using these losses and select k labels with the highest output scores.

For all baselines and **MACRO- R_{PRIOR}** , we always train the label probability estimator on the whole training set. For **MACRO- P_{FW}** , **MACRO- R_{FW}** , and **MACRO- $F1_{\text{FW}}$** , we tested different ratios (50/50 or 75/25) of splitting training data into sets used for training the label probability estimators and estimating confusion matrix C , as well as a variant where we use the whole training set for both steps. We also investigated two strategies for initialization of classifier h by either using equal weights (resulting in a TOP-K classifier) or random weights. Finally, we terminate the algorithm if we do not observe sufficient improvement in the objective. In practice, we found that Frank-Wolfe converges within 3–10 iterations. Because of the nature of the random classifier, we repeat the inference on the test set 10 times and report the mean results. In Table 2 we present the variant achieving the best results, and report all the results including standard deviations, running times, number of Frank-Wolfe iterations in Appendix G.

The randomized classifiers obtained via the Frank-Wolfe algorithm achieve, in most cases, the best results for the measure they aim to optimize, at the cost of loosing on some instance-wise measures. However, they sometimes fail to obtain the best results on the largest dataset, where the majority of labels have only a few (less than 10) positive instances in the training set, preventing them from obtaining accurate estimates of η and C . In this case, simple heuristics like **TOP-K+ w^{POW}** might work better. Popular Focal loss and Asymmetric loss preserve the performance on instance-wise metrics, but improvement on the macro measures is usually small. It is also worth noting that, as expected, **MACRO- R_{FW}** recovers the solution of **MACRO- R_{PRIOR}** in all cases.

7 CONCLUSIONS

In this paper, we have focused on developing a consistent algorithm for complex macro-at- k metrics in the framework of population utility (PU). Our main results have been obtained by following the line of research conducted in the context of multi-class classification with additional constraints. However, these previous works do not address the specific scenario of budgeted predictions at k , which commonly arises in multi-label classification problems. For the complex macro-at- k metrics, we have introduced a consistent Frank-Wolfe algorithm, which is capable of finding an optimal randomized classifier by transforming the problem of optimizing over classifiers to optimizing over the set of feasible confusion matrices and using the fact that the optimal classifier optimizes (unknown) linear confusion-matrix. Our empirical studies show that the introduced approach effectively optimizes macro-measures and it scales to even larger datasets with thousands of labels.

ACKNOWLEDGMENTS

A part of computational experiments for this paper had been performed in Poznan Supercomputing and Networking Center. We want to acknowledge the support of Academy of Finland via grants 347707 and 348215.

REFERENCES

- Alekh Agarwal, Alina Beygelzimer, Miroslav Dudik, John Langford, and Hanna Wallach. A reductions approach to fair classification. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 60–69. PMLR, 10–15 Jul 2018.
- Rohit Babbar and Bernhard Schölkopf. Dismec: Distributed sparse machines for extreme multi-label classification. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining, WSDM '17*, pp. 721—729, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450346757. doi:[10.1145/3018661.3018741](https://doi.org/10.1145/3018661.3018741).
- Rohit Babbar and Bernhard Schölkopf. Data scarcity, robustness and extreme multi-label classification. *Machine Learning*, 108(8):1329–1351, 2019.
- Peter L Bartlett, Michael I Jordan, and Jon D McAuliffe. Convexity, classification, and risk bounds. *Journal of the American Statistical Association*, 101(473):138–156, 2006.
- Eric Baum and David Haussler. What size net gives valid generalization? In D. Touretzky (ed.), *Advances in Neural Information Processing Systems*, volume 1. Morgan-Kaufmann, 1988.
- Kush. Bhatia, Kunal. Dahiya, Himanshu Jain, Anshul Mittal, Yashoteja Prabhu, and Manik Varma. The extreme classification repository: Multi-label datasets and code, 2016. URL <http://manikvarma.org/downloads/XC/XMLRepository.html>.
- Róbert Busa-Fekete, Balázs Szörényi, Krzysztof Dembczynski, and Eyke Hüllermeier. Online F-measure optimization. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pp. 595–603, 2015.
- Clément Calauzenes, Nicolas Usunier, and Patrick Gallinari. On the (non-) existence of convex, calibrated surrogate losses for ranking. *Advances in Neural Information Processing Systems*, 25, 2012.
- Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Wei-Cheng Chang, Daniel Jiang, Hsiang-Fu Yu, Choon Hui Teo, Jiong Zhang, Kai Zhong, Kedarnath Kolluri, Qie Hu, Nikhil Shandilya, Vyacheslav Ievgrafov, et al. Extreme multi-label learning for semantic matching in product search. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pp. 2643–2651, 2021.
- Paolo Cremonesi, Yehuda Koren, and Roberto Turrin. Performance of recommender algorithms on top-n recommendation tasks. In *Proceedings of the fourth ACM conference on Recommender systems*, pp. 39–46, 2010.
- Krzysztof Dembczynski, Willem Waegeman, Weiwei Cheng, and Eyke Hüllermeier. An exact algorithm for F-measure maximization. In *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain.*, pp. 1404–1412, 2011.
- Krzysztof Dembczynski, Wojciech Kotłowski, Oluwasanmi Koyejo, and Nagarajan Natarajan. Consistency analysis for binary classification revisited. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pp. 961–969. PMLR, 2017.
- Chris Drummond and Robert C. Holte. Severe class imbalance: Why better algorithms aren't the answer. In *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD)*, pp. 539–546. Springer, 2005.

- John C. Duchi, Lester W. Mackey, and Michael I. Jordan. On the consistency of ranking algorithms. In *Proceedings of the 27th International Conference on Machine Learning*, pp. 327–334, 2010.
- Stack Exchange. Continuous linear image of closed, bounded, and convex set of a hilbert space is compact. URL <https://math.stackexchange.com/q/908121>. Accessed: 2023-09-27.
- Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. LIBLINEAR: A library for large linear classification. *Journal of Machine Learning Research*, 9:1871–1874, 2008.
- Marguerite Frank and Philip Wolfe. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3(1-2):95–110, 1956. doi:[10.1002/nav.3800030109](https://doi.org/10.1002/nav.3800030109).
- Muhammad Hanif and K. R. W. Brewer. Sampling with unequal probabilities without replacement: a review. *International Statistical Review*, 48:317–335, 1980.
- Martin Jaggi. Revisiting frank-wolfe: Projection-free sparse convex optimization. In *International conference on machine learning*, pp. 427–435. PMLR, 2013.
- Himanshu Jain, Yashoteja Prabhu, and Manik Varma. Extreme multi-label loss functions for recommendation, tagging, ranking and other missing label applications. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16*, pp. 935–944. Association for Computing Machinery, 2016. ISBN 9781450342322.
- Kalina Jasinska, Krzysztof Dembczynski, Róbert Busa-Fekete, Karlson Pfannschmidt, Timo Klerx, and Eyke Hüllermeier. Extreme F-measure maximization using sparse probability estimates. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pp. 1435–1444. JMLR.org, 2016.
- Kalina Jasinska-Kobus, Marek Wydmuch, Krzysztof Dembczyński, Mikhail Kuznetsov, and Róbert Busa-Fekete. Probabilistic label trees for extreme multi-label classification. *CoRR*, abs/2009.11218, 2020.
- Thorsten Joachims. A support vector method for multivariate performance measures. In *Proceedings of the 22nd International Conference on Machine Learning (ICML 2005), August 7-11, 2005, Bonn, Germany, 2005*.
- Purushottam Kar, Harikrishna Narasimhan, and Prateek Jain. Surrogate functions for maximizing precision at the top. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pp. 189–198, 2015.
- Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. Preventing fairness gerrymandering: Auditing and learning for subgroup fairness. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 2564–2572. PMLR, 10–15 Jul 2018.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun (eds.), *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- Wojciech Kotłowski and Krzysztof Dembczyński. Surrogate regret bounds for generalized classification performance metrics. *Machine Learning*, 10:549–572, 2017.
- Oluwasanmi Koyejo, Nagarajan Natarajan, Pradeep Ravikumar, and Inderjit S. Dhillon. Consistent binary classification with generalized performance metrics. In *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pp. 2744–2752, 2014.
- Oluwasanmi Koyejo, Nagarajan Natarajan, Pradeep Ravikumar, and Inderjit S. Dhillon. Consistent multilabel classification. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pp. 3321–3329, 2015.

- David D Lewis. Evaluating and optimizing autonomous text classification systems. In *Proceedings of the 18th Intl. ACM SIGIR Conf. on Research and Development in Information Retrieval*, pp. 246–254. ACM, 1995.
- David D Lewis, Yiming Yang, Tony Russell-Rose, and Fan Li. RCV1: A new benchmark collection for text categorization research. *Journal of Machine Learning Research*, 5(Apr):361–397, 2004.
- Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pp. 2980–2988, 2017.
- William G. Madow. On the theory of systematic sampling, II. *The Annals of Mathematical Statistics*, 20(3):333–354, 1949.
- Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, 2008.
- Julian McAuley and Jure Leskovec. Hidden factors and hidden topics: understanding rating dimensions with review text. In *Proceedings of the 7th ACM conference on Recommender systems*, pp. 165–172, 2013.
- Julian McAuley, Rahul Pandey, and Jure Leskovec. Inferring networks of substitutable and complementary products. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 785–794, 2015.
- Aditya K. Menon, Harikrishna Narasimhan, Shivani Agarwal, and Sanjay Chawla. On the statistical consistency of algorithms for binary classification under class imbalance. In *International Conference on Machine Learning (ICML)*, 2013.
- Samrat Mukhopadhyay, Sourav Sahoo, and Abhishek Sinha. k-experts - online policies and fundamental limits. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pp. 342–365. PMLR, 2022.
- Harikrishna Narasimhan, Harish Ramaswamy, Aadirupa Saha, and Shivani Agarwal. Consistent multiclass algorithms for complex performance measures. In Francis Bach and David Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pp. 2398–2407, Lille, France, 07 2015. PMLR.
- Harikrishna Narasimhan, Harish G. Ramaswamy, Shiv Kumar Tavker, Drona Khurana, Praneeth Netrapalli, and Shivani Agarwal. Consistent multiclass algorithms for complex metrics and constraints. 2022. URL <https://arxiv.org/abs/2210.09695>.
- Nagarajan Natarajan, Oluwasanmi Koyejo, Pradeep Ravikumar, and Inderjit Dhillon. Optimal classification with multivariate losses. In *Proceedings of The 33rd International Conference on Machine Learning (ICML)*, pp. 1530–1538, 2016.
- Nagarajan Natarajan, Inderjit S. Dhillon, Pradeep Ravikumar, and Ambuj Tewari. Cost-sensitive learning with noisy labels. *Journal of Machine Learning Research*, 18(155):1–33, 2018.
- Juri Opitz and Sebastian Burst. Macro F1 and macro F1. *CoRR*, abs/1911.03347, 2019.
- Shameem Puthiya Parambath, Nicolas Usunier, and Yves Grandvalet. Optimizing F-measures by cost-sensitive classification. In *Neural Information Processing Systems (NIPS)*, 2014.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp. 8024–8035. Curran Associates, Inc., 2019.

- Harish G. Ramaswamy, Ambuj Tewari, and Shivani Agarwal. Consistent algorithms for multiclass classification with an abstain option. *Electronic Journal of Statistics*, 12(1):530–554, 2018. doi:[10.1214/17-EJS1388](https://doi.org/10.1214/17-EJS1388).
- Pradeep Ravikumar, Ambuj Tewari, and Eunho Yang. On ndcg consistency of listwise ranking methods. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pp. 618–626. JMLR Workshop and Conference Proceedings, 2011.
- Tal Ridnik, Emanuel Ben-Baruch, Nadav Zamir, Asaf Noy, Itamar Friedman, Matan Protter, and Lihi Zelnik-Manor. Asymmetric loss for multi-label classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 82–91, 2021.
- Siegfried Schaible and Jianming Shi. Fractional programming: The sum-of-ratios case. *Optimization Methods and Software*, 18(2):219–229, 2003.
- Erik Schultheis, Marek Wydmuch, Rohit Babbar, and Krzysztof Dembczynski. On missing labels, long-tails and propensities in extreme multi-label classification. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 1547–1557, 2022.
- Erik Schultheis, Marek Wydmuch, Wojciech Kotlowski, Rohit Babbar, and Krzysztof Dembczynski. Generalized test utilities for long-tail performance in extreme multi-label classification. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 22269–22303. Curran Associates, Inc., 2023.
- Shashank Singh and Justin T Khim. Optimal binary classification beyond accuracy. In *Advances in Neural Information Processing Systems*, volume 35, pp. 18226–18240. Curran Associates, Inc., 2022.
- Cees G. M. Snoek, Marcel Worring, Jan C. van Gemert, Jan-Mark Geusebroek, and Arnold W. M. Smeulders. The challenge problem for automated detection of 101 semantic concepts in multimedia. In *ACM International Conference on Multimedia*, pp. 421–430. Association for Computing Machinery, 2006.
- Lei Tang and Huan Liu. Relational learning via latent social dimensions. In *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '09, pp. 817–826, New York, NY, USA, 2009. Association for Computing Machinery. ISBN 9781605584959. doi:[10.1145/1557019.1557109](https://doi.org/10.1145/1557019.1557109)
- Ambuj Tewari and Peter L Bartlett. On the consistency of multiclass classification methods. *Journal of Machine Learning Research*, 8(5), 2007.
- Willem Waegeman, Krzysztof Dembczynski, Arkadiusz Jachnik, Weiwei Cheng, and Eyke Hüllermeier. On the bayes-optimality of F-measure maximizers. *Journal of Machine Learning Research*, 15(1):3333–3388, 2014.
- Shuo Wang and Xin Yao. Multiclass imbalance problems: Analysis and potential solutions. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 42(4):1119–1130, 2012.
- Xiaoyan Wang, Ran Li, Bowei Yan, and Oluwasanmi Koyejo. Consistent classification with generalized metrics. 2019.
- Tong Wei and Yu-Feng Li. Does tail label help for large-scale multi-label learning? *IEEE transactions on neural networks and learning systems*, 31(7):2315–2324, 2019.
- Marek Wydmuch, Kalina Jasinska, Mikhail Kuznetsov, Róbert Busa-Fekete, and Krzysztof Dembczynski. A no-regret generalization of hierarchical softmax to extreme multi-label classification. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 31*, pp. 6355–6366. Curran Associates, Inc., 2018.
- Forest Yang and Sanmi Koyejo. On the consistency of top-k surrogate losses. In *International Conference on Machine Learning*, pp. 10727–10735. PMLR, 2020.

Nan Ye, Kian Ming Adam Chai, Wee Sun Lee, and Hai Leong Chieu. Optimizing F-measure: A tale of two approaches. In *Proceedings of the 29th International Conference on Machine Learning, ICML 2012, Edinburgh, Scotland, UK, June 26 - July 1, 2012*. icml.cc / Omnipress, 2012.

Ming Yuan and Marten Wegkamp. Classification methods with reject option based on convex risk minimization. *Journal of Machine Learning Research*, 11:111–130, mar 2010. ISSN 1532-4435.

Tong Zhang. Statistical analysis of some multi-category large margin classification methods. *Journal of Machine Learning Research*, 5(Oct):1225–1251, 2004.