

---

# Unbiased Estimates for Multilabel Reductions of Extreme Classification with Missing Labels

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This paper considers the missing-labels problem in the extreme multilabel clas-  
2 sification (XMC) setting, i.e. a setting with a very large label space. The goal in  
3 XMC often is to maximize either precision or recall of the top-ranked predictions,  
4 which can be achieved by reducing the multilabel problem into a series of binary  
5 (One-vs-All) or multiclass (Pick-all-Labels) problems. Missing labels are a ubiqui-  
6 tous phenomenon in XMC tasks, yet the interaction between missing labels and  
7 multilabel reductions has hitherto only been investigated for the case of One-vs-All  
8 reduction. In this paper, we close this gap by providing unbiased estimates for  
9 general (non-decomposable) multilabel losses, which enables unbiased estimates  
10 of the Pick-all-Labels reduction, as well as the normalized reductions which are  
11 required for consistency with the recall metric. We show that these estimators  
12 suffer from increased variance and may lead to ill-posed optimization problems.  
13 To address this issue, we propose to use convex upper bounds which trade off an  
14 increase in bias against a strong decrease in variance.

## 15 1 Introduction

16 Extreme multilabel classification (XMC) is a machine learning setting in which the goal is to predict  
17 a small subset of positive (or relevant) labels for each data instance out of a very large (thousands to  
18 millions) set of possible labels. Such problems arise for example when annotating large encyclopedia  
19 [7, 28], in fine-grained image classification [9], and next-word prediction [25]. Further applications of  
20 XMC are recommendation systems, web-advertising and prediction of related searches [1, 29, 17, 6].

21 Typical datasets in these scenarios are very large, resulting in possibly billions of (data, label) pairs  
22 [4], making it impossible for human annotators to check each pair. Even annotating only a few  
23 samples fully in order to generate a clean test set can be prohibitively expensive. Therefore, both the  
24 available training- and test-data are likely to contain some errors. Fortunately, in many cases it is  
25 possible to constrain the structure of the labeling errors. Consider, for example, the case of tagging  
26 documents: Here, we can assume that each label with which the document has been tagged has been  
27 deemed relevant by the annotator, and thus is relatively surely a correct label. On the other hand, the  
28 annotator cannot possibly check hundreds of thousands of negative labels. This leads to the setting  
29 of missing labels investigated in this paper, in which only positive labels are affected by noise (they  
30 can go missing), whereas negative labels remain unchanged (no spurious labels). This model has  
31 been introduced to the XMC setting by Jain et al. [16], along with estimates for the *propensities*, the  
32 chance of a relevant label to be observed. Similar models are using in learning-to-rank[20, 27, 37]  
33 and recommendation systems[32, 14, 15]. For a formal definition of the setting we refer the reader to  
34 section 3, and for a more thorough discussion of prior works on missing labels and related settings to  
35 section 6.

36 A common strategy for learning XMC classifiers is to reduce the multilabel problem [34] into a series  
 37 of binary [8, 3, 40] or multiclass [18, 38, 31] problems, which then can be solved using existing  
 38 techniques. Such *loss reductions* can be shown to be consistent for the tasks of maximizing precision  
 39 at  $k$  or recall at  $k$ , but never both at the same time [24]. For one of these methods, One-vs-All,  
 40 adaptation to the missing labels setting has been shown to yield an improvement in propensity-scored  
 41 precision (an unbiased estimate of precision@k) metrics [30]. The reductions consistent for precision  
 42 lead to loss functions that can be decomposed into a sum of contributions from each label, which  
 43 means the results of Natarajan et al. [26] can be applied. In contrast, the reductions consistent for  
 44 recall contain a normalization term that is the inverse of the total number of true labels. This term is  
 45 also necessary for calculating the recall metric itself, demonstrating the need for unbiased estimates  
 46 for true, non-decomposable multilabel loss functions.

47 **Contributions** Our contributions are **1)** A mathematical model of the missing labels setting that  
 48 describes the observed labels as a product of an (unknown) mask variable with the true labels.  
 49 Crucially, this mask can be chosen to be *independent* of the labels (Theorem 1), enabling simple  
 50 proofs for our theorems. **2)** The unique unbiased estimate (Theorems 2, 3) for arbitrary multilabel  
 51 losses, and in particular for the loss functions arising from multilabel reductions. The unbiased  
 52 estimate of a lower-bounded loss need not be lower-bounded, and even for bounded losses the  
 53 unbiased estimate leads to an increase in variance. Therefore, we develop **3)** a convex upper-bound  
 54 (Theorem 4) for losses based on the normalized Pick-all-Labels reduction. In the missing-labels  
 55 setting, the generalization error is composed of two contributions: the error due to overfitting to  
 56 the specific, observed noise-pattern, and the error because only a finite sample has been observed.  
 57 We present empirical evidence **4)** that the former can be much stronger than the latter, and may be  
 58 reduced by switching to the upper bounds.

59 In the main paper, we provide shortened proofs that illustrate the key steps. Detailed step-by-step  
 60 proofs can be found in the appendix.

61 **Notation** Random variables will be denoted by capital letters  $X, Y, \dots$ , whereas calligraphic letters  
 62 denote sets and lower case letters their elements,  $x \in \mathcal{X}, \dots$ . Vectors will be denoted by bold font,  
 63  $\mathbf{y} \in \mathcal{Y}$ , if we plan to make use of the fact that they can be decomposed into components  $y_1, \dots, y_k$ ,  
 64 with  $\mathbf{y}_{-k}$  denoting the vector of all components except the  $k$ 'th. The letters  $f, g, h$  and  $\ell$  are reserved  
 65 for functions,  $i, j, k$  denote integers,  $[k]$  is the set  $\{1, \dots, k\}$ . We denote with  $\mathcal{X}$  the *data space*,  
 66  $\mathcal{Y} = \{0, 1\}^l$  the *label space* and  $\hat{\mathcal{Y}} = \mathbb{R}^l$  the *prediction space*. A dataset is defined through the three  
 67 random variables  $X \in \mathcal{X}$ ,  $\mathbf{Y} \in \mathcal{Y}$ , and  $\mathbf{Y}^* \in \mathcal{Y}$ , that represent the *data*, *observed label*, and *ground*  
 68 *truth label*. We mark quantities pertaining to the unobservable ground-truth with a superscript star  
 69 and call  $(X, \mathbf{Y}^*)$  the *clean data*.

## 70 2 Multilabel Reductions

71 In Menon et al. [24], five different reductions for turning the multilabel learning problem into a sum of  
 72 binary or multiclass problems are presented (cf. appendix). In the following, let  $\ell_{\text{BC}} : \{0, 1\} \times \mathbb{R} \rightarrow$   
 73  $\mathbb{R}$  be a binary loss and  $\ell_{\text{MC}} : [l] \times \mathbb{R}^l \rightarrow \mathbb{R}$  be a multiclass loss. Below, we present four of those  
 74 reductions, and rearrange their loss functions so that a common pattern emerges.

75 For *one-vs-all* (OVA) reduction, each label is considered independently, meaning that for each  
 76 instance  $l$  binary problems are to be solved. This leads to a loss function

$$\ell_{\text{OVA}}^*(\mathbf{y}^*, \hat{\mathbf{y}}) = \sum_{j=1}^l \ell_{\text{BC}}(y_j^*, \hat{y}_j) = \sum_{j=1}^l y_j^* (\ell_{\text{BC}}(1, \hat{y}_j) - \ell_{\text{BC}}(0, \hat{y}_j)) + \ell_{\text{BC}}(0, \hat{y}_j). \quad (1)$$

77 In contrast, *pick-all-labels* (PAL) considers all the positive labels for each instance and tries to  
 78 minimize their corresponding multiclass loss, leading to

$$\ell_{\text{PAL}}^*(\mathbf{y}^*, \hat{\mathbf{y}}) = \sum_{j: y_j^*=1} \ell_{\text{MC}}(j, \hat{\mathbf{y}}) = \sum_{j \in [l]} y_j^* \ell_{\text{MC}}(j, \hat{\mathbf{y}}). \quad (2)$$

79 Both approaches are consistent for precision at  $k$ . In order to make the reductions consistent for recall  
 80 instead of precision, the label value needs to be replaced with a normalized label

$$\tilde{y}_j^* := \frac{y_j^*}{\sum_{i=1}^l y_i^*} = \frac{y_j^*}{1 + \sum_{i \neq j}^l y_i^*}, \quad (3)$$

81 where the expression on the right has the advantage of being well defined even if there are no positives  
 82 for the sample. This leads to the OVA-N and PAL-N reductions. By moving label-independent parts  
 83 into functions  $f$  and  $g_j$ , the reductions get a common structure

$$\ell^*(\mathbf{y}^*, \hat{\mathbf{y}}) = f(\hat{\mathbf{y}}) + \sum_{j=1}^l z_j^* g_j(\hat{\mathbf{y}}), \quad (4)$$

84 where  $z_j = \tilde{y}_j^*$  for the normalized reductions and  $z_j^* = y_j^*$  otherwise. The functions  $f$  and  $g_j$  are the  
 85 same for the normalized and regular reduction (see appendix).

### 86 3 Unbiased Estimates with Missing Labels

87 We are interested in noisy labels where the noise is such that labels can only go missing. This is  
 88 described by the next two definitions, where the first gives a phenomenological characterization of  
 89 the setting, whereas the second defines the mathematical model used to describe it. For this setting  
 90 we then develop unbiased estimates for the preceding loss reductions, in the sense that for a given  
 91 loss  $\ell^*$  we are looking for a new loss function  $\ell$  such that  $\mathbb{E}[\ell(\mathbf{Y}, \hat{\mathbf{Y}})] = \mathbb{E}[\ell^*(\mathbf{Y}^*, \hat{\mathbf{Y}})]$ .

92 **Definition 1** (Propensity). The missing-labels setting we described informally in the introduction  
 93 leads to the following conditions on the  $l$  random variables

$$\mathbb{P}\{Y_j = 1 \mid Y_j^* = 1, \mathbf{Y}^*_{-j}, X\} =: p_j(X), \quad \mathbb{P}\{Y_j = 1 \mid Y_j^* = 0, \mathbf{Y}^*_{-j}, X\} = 0 \quad (5)$$

94 The value  $p_j(x) \in (0, 1]$  is called the *propensity* of the label  $j$  at point  $x$ .

95 Such propensity models have been used in extreme classification [30, 16, 39], learning-to-rank  
 96 [20, 27, 37], and recommendation systems [32, 14, 15].

97 The following proposition guarantees that a fixed-propensity unbiased estimator can be used to  
 98 construct a instance-dependent unbiased estimator

99 **Proposition 1.** Let  $f^*(X, Y^*)$  be some function such that for fixed propensity  $\mathbf{p}$ , an unbiased  
 100 estimate is given by  $f_{\mathbf{p}}$ , i.e.  $\mathbb{E}[f_{\mathbf{p}}(X, Y)] = \mathbb{E}[f^*(X, Y^*)]$ . For instance-dependent propensity  $\mathbf{p}(x)$ ,  
 101 an unbiased estimator of  $f^*$  is given by  $f_{\mathbf{p}(X)}$ .

102 *Proof.* Using the law of total expectation gives

$$\mathbb{E}[f^*(X, Y^*)] = \mathbb{E}[\mathbb{E}[f^*(X, Y^*) \mid X]] = \mathbb{E}[\mathbb{E}[f_{\mathbf{p}(X)}(X, Y^*) \mid X]] = \mathbb{E}[f_{\mathbf{p}(X)}(X, Y^*)]. \quad \square$$

103 Therefore, we will suppress the dependence of the propensity on the data point in the rest of the paper.

104 The relation between  $\mathbf{Y}^*$  and  $\mathbf{Y}$  can be modeled by a set of independent *mask* variables  $\mathbf{M}$ :

105 **Theorem 1** (Masking Model). Assuming  $\mathbf{Y}^*$  and  $\mathbf{Y}$  follow Definition 1, then then there exists a  
 106 random variable  $\mathbf{M} \in \{0, 1\}^l$  such that  $\mathbf{Y} = \mathbf{M} \odot \mathbf{Y}^*$  almost surely and  $M_j$  is independent of  
 107  $(\mathbf{Y}^*, X, \mathbf{M}_{-j})$  for all  $j \in [l]$ . It holds that  $\mathbb{E}[M_j] = p_j$ .

108 This can be seen as a multilabel generalization of the similar statement given in Teisseyre et al. [33].  
 109 The independent variables  $\mathbf{M}$  provide a convenient framework for proving the results that follow,  
 110 because the independence allows to factorize expectations containing  $\mathbf{M}$ .

111 **Proposition 2** (Unbiased Estimate for Decomposable Reductions). Assume the setting of Definition 1,  
 112 with the additional condition that the predictions  $\hat{\mathbf{Y}}$  are independent of the missing mask  $\mathbf{M}$ . Then  
 113 the unbiased estimate for the loss (4) with  $z = \mathbf{y}$ , denoted by  $\ell = \mathfrak{P}(\ell^*)$ , is given by

$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = f(\hat{\mathbf{y}}) + \sum_{j=1}^l \frac{y_j}{p_j} g_j(\hat{\mathbf{y}}). \quad (6)$$

114 The predictions have to be independent of the locations  $\mathbf{M}$  where the labels go missing. This is  
 115 fulfilled if the predictions  $\hat{\mathbf{Y}} = h(X, \mathbf{W})$  are the output of a classifier  $h$  whose weights  $\mathbf{W}$  are  
 116 independent of  $\mathbf{M}$ .<sup>1</sup>

117 For the normalized reductions, it would suffice to find an unbiased estimate of  $\tilde{Y}$  in order to apply  
 118 the same argument as above. However, we are not aware of a derivation for such an estimate that is  
 119 simpler than the fully generic case presented below.

120 **Theorem 2** (Unbiased Estimate for Non-Decomposable Loss). *For a generic multilabel loss function*  
 121  *$\ell^*$ , the unbiased estimate  $\ell = \mathfrak{P}(\ell^*)$  under the conditions of Theorem 2 is given by*

$$\ell(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{\mathbf{y}' \preceq \mathbf{y}} \prod_{j: y'_j=1} \left( \frac{y'_j(2-p_j) + p_j - 1}{p_j} \right) \ell^*(\mathbf{y}', \hat{\mathbf{y}}), \quad (7)$$

122 where  $\mathbf{y}' \preceq \mathbf{y}$  means  $\{0, 1\} \ni y'_j \leq y_j$ .

123 This means that for an instance with  $k$  positive labels, we need  $2^k$  evaluations of the original loss  
 124 function in order to calculate the unbiased estimate. This is only feasible because, despite having a  
 125 very large label space, typical extreme-classification datasets have only few positives per instance.

126 Unfortunately, the division by (products of) propensity values means that the unbiased estimates will  
 127 have much larger variance than the original loss function would have on clean data. As an illustrative  
 128 example, consider the binary case in the limit  $p \ll 1$ . We can show that in this case the variance  
 129 grows with  $p^{-1}$  compared to the evaluation on clean data.

130 **Proposition 3** (Increase in Variance). *Setting  $q^* := \mathbb{E}[Y^*]$  and  $\ell = \mathfrak{P}(\ell^*)$ , for small propensities*  
 131  *$p \ll 1$ , the variance increases with the inverse of the propensity,  $\mathbb{V}[\ell(Y, \hat{y})] \approx \frac{1}{p(1-q^*)} \mathbb{V}[\ell^*(Y^*, \hat{y})]$ .*

132 This means that in the binary case the variance increases linearly with inverse propensity. In the  
 133 multilabel case, this is amplified further due to the product of propensities.

134 The result above raises the question whether there might be other unbiased estimators with reduced  
 135 variance. For example, the conditional expectation  $\mathbb{E}[\ell^*(Y^*, X)|Y]$  also gives an unbiased estimate  
 136 with lower variance, but cannot be calculated without knowledge of the conditional probabilities  
 137  $\mathbb{P}\{Y | X\}$ . The following theorem states that  $\ell = \mathfrak{P}(\ell^*)$  is unique if we want the loss function to  
 138 work for all possible distributions of data. Thus we cannot reduce the variance.

139 **Theorem 3** (Uniqueness). *Let  $p_j \in (0, 1] \forall j \in [l]$ . For an arbitrary loss function  $\ell^*$ , let  $\ell$  and  $\ell'$  be*  
 140 *unbiased versions, in the sense that for all  $X, \mathbf{Y}, \mathbf{Y}^*$  that fulfill the masking model Theorem 1 with*  
 141 *propensity  $\mathbf{p}$ , it holds*

$$\mathbb{E}[\ell^*(\mathbf{Y}^*, X)] = \mathbb{E}[\ell(\mathbf{Y}, X)] = \mathbb{E}[\ell'(\mathbf{Y}, X)]. \quad (8)$$

142 Then,  $\ell' = \ell$ .

143 The unavoidable increase in variance indicates that there might be a bias-variance trade-off between  
 144 using the unbiased loss that may overfit more strongly on the observed noise, and using the original  
 145 loss function which gives wrong results even if  $n \rightarrow \infty$ . If one calculates a standard Rademacher  
 146 bound for generalization (see appendix), this error bound increases with a factor  $\frac{2-p}{p}$ .<sup>2</sup>

147 In a classical learning setup, the generalization error would be described by the difference between  
 148 the empirical risk and the true risk  $\hat{\mathbf{R}}_{\ell^*}^*[\hat{h}] - \mathbf{R}_{\ell^*}^*[\hat{h}]$ . However, in the case of missing labels, this  
 149 can be decomposed in two ways

$$\begin{aligned} \mathbf{R}_{\ell^*}^*[h] - \hat{\mathbf{R}}_{\ell}[h] &= \overbrace{\mathbf{R}_{\ell^*}^*[h] - \mathbf{R}_{\ell}[h]}^{=0} + \mathbf{R}_{\ell}[h] - \hat{\mathbf{R}}_{\ell}[h] & (9) \\ &= \underbrace{\mathbf{R}_{\ell^*}^*[h] - \hat{\mathbf{R}}_{\ell^*}^*[h]}_{\text{finite sample}} + \underbrace{\hat{\mathbf{R}}_{\ell^*}^*[h] - \hat{\mathbf{R}}_{\ell}[h]}_{\text{noise pattern}}, & (10) \end{aligned}$$

150 Whereas the first equation is just a restatement of the unbiasedness, the second contains some new  
 151 insight: The generalization error can be decomposed into the difference between the true risk  $\mathbf{R}_{\ell^*}^*[h]$

<sup>1</sup>In this sense, we will use the notation  $\ell(y, x)$  to evaluate a loss also on a data point.

<sup>2</sup>The bound in this paper corresponds to Natarajan et al. [26, Thm. 9], though that published result is wrong and missing the increase in the bound due to the increased range of the function.

152 and the empirical risk on clean training data  $\hat{R}_{\ell^*}^*[h]$ , and the difference between that and the estimated  
 153 empirical risk on observed data  $\hat{R}_\ell[h]$ . Because the classifier  $h$  depends (through  $Y = \mathbf{M} \odot Y^*$ ) on  
 154 the mask variables,  $\ell$  does not give an unbiased estimate (on training data) and thus the second term  
 155 is non-zero even in expectation. In fact, in the low-regularization regime this term may dominate the  
 156 entire error, as we will demonstrate in section 5.

## 157 4 Convex Upper-Bounds

158 The unbiased estimate allows us to calculate the loss even on data with missing labels, but can we  
 159 also use it for training? Ideally, the loss function should be lower-bounded, so the minimization is  
 160 well defined, it should be convex so the minimum is unique. Further, the variance of the unbiased  
 161 estimator should not be too large, so that a reasonable amount of training samples is sufficient.

162 If we assume  $\ell_{\text{BC}}$  and  $\ell_{\text{MC}}$  to be lower-bounded and convex, then only the PAL-reduction results in  
 163 an unbiased estimate that is guaranteed to have the same properties, as it is a positive combination  
 164 of  $\ell_{\text{MC}}$ . Due to the uniqueness result, it is not possible to find an unbiased estimate that is always  
 165 convex for the other reductions. Thus, in order to make them amenable for training, we propose to  
 166 switch from unbiased estimates to convex upper-bounds. Below we present solutions for the OvA  
 167 and normalized PAL-reduction. The normalized OVA-reduction remains an open problem.

168 **Upper-Bound for OvA-Reduction** The OvA-reduction is based on a binary loss, which often  
 169 is a convex surrogate for the 0-1 loss. To get a convex loss in the missing-labels case, we thus  
 170 switch the order of operations [30, 5]: Instead of taking an unbiased estimate of a convex surrogate,  
 171 we form a convex surrogate of an unbiased estimate. Taking  $\theta$  to be a thresholding function (e.g.  
 172  $\theta(s) = \mathbb{1}\{s > 0\}$ ), the 0-1-loss can be written as

$$\ell_{0-1}^*(y, \hat{y}) = y\theta(\hat{y}) + (1 - y)(1 - \theta(\hat{y})) \quad (11)$$

173 with unbiased estimate

$$\ell_{0-1}(y, \hat{y}) = \left(\frac{2}{p_j} - 1\right) y\theta(\hat{y}) + (1 - y)(1 - \theta(\hat{y})) + y \left(\frac{p_j - 1}{p_j}\right). \quad (12)$$

174 As the last term does not depend on the predictions, it can be dropped for an optimization objective.  
 175 If  $\ell_{\text{BC}}(1, \hat{y})$  is a convex upper-bound on  $\theta(\hat{y})$  and  $\ell_{\text{BC}}(0, \hat{y})$  on  $(1 - \theta(\hat{y}))$ , so that overall  $\ell_{\text{BC}}$  is a  
 176 convex upper-bound on the 0-1 loss, then performing these substitutions gives a convex loss function  
 177 for the OvA-reduction:

$$\tilde{\ell}_{\text{OvA}}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{j=1}^l \left(\frac{2}{p_j} - 1\right) y_j \ell_{\text{BC}}(1, \hat{y}_j) + (1 - y_j) \ell_{\text{BC}}(0, \hat{y}_j) \quad (13)$$

178 **Upper-Bound for Normalized PAL-Reduction** We have formulated the normalized multilabel  
 179 reductions in terms of the variable  $\tilde{Y}^*$ . A naive attempt of correcting for the noisy labels by replacing  
 180  $Y^*$  with  $Y/p$  is not unbiased. However, the resulting estimator  $\tilde{Y}$  turns out to be an upper bound.  
 181 The two estimators are given by

$$\tilde{Y}_i^* = \frac{Y_i^*}{1 + \sum_{j \neq i} Y_j^*}, \quad \tilde{Y}_i := \frac{Y_i/p_i}{1 + \sum_{j \neq i} Y_j/p_j}. \quad (14)$$

182 **Theorem 4** (Normalized Label Upper-Bound). *Under the conditions of Theorem 2, replacing the*  
 183 *true label with the unbiased estimate of the observed label as shown in Equation 14 results in an*  
 184 *upper bound, whose error itself can be bounded by a data-dependent term*

$$\mathbb{E}[\tilde{Y}_i^*] + \sum_{j \neq i} \left(\frac{1 - p_j}{p_j}\right) \mathbb{E}\left[\frac{Y_i}{p_i} \cdot \frac{Y_j}{p_j}\right] \geq \mathbb{E}[\tilde{Y}_i] \geq \mathbb{E}[\tilde{Y}_i^*]. \quad (15)$$

185 *Proof.* For convenience denote  $S_i^* := \sum_{j \neq i} Y_j^*$  and  $S_i := \sum_{j \neq i} Y_j/p_j$ , and note that  $S_i$  is independ-  
 186 ent of  $M_i$ . By pulling out known factors and using the independence of  $M$  and  $\mathbf{Y}^*$  we can show  
 187 that

$$\mathbb{E}[S_i | \mathbf{Y}^*] = \sum_{j \neq i} \mathbb{E}[M_j Y_j^* / p_j | \mathbf{Y}^*] = \sum_{j \neq i} Y_j^* \mathbb{E}[M_j / p_j | \mathbf{Y}^*] = S_i^*. \quad (16)$$

188 Expanding terms and using independence of  $M_i$ , then applying the tower property and pulling out  
 189 the measurable factor results in

$$\mathbb{E}[\tilde{Y}_i] = \mathbb{E}\left[\frac{M_i Y_i^*/p_i}{1+S_i}\right] = \mathbb{E}\left[\frac{M_i}{p_i}\right] \mathbb{E}\left[\frac{Y_i^*}{1+S_i}\right] = \mathbb{E}\left[\mathbb{E}\left[\frac{Y_i^*}{1+S_i} \mid \mathbf{Y}^*\right]\right] = \mathbb{E}\left[Y_i^* \mathbb{E}\left[\frac{1}{1+S_i} \mid \mathbf{Y}^*\right]\right].$$

The function  $h : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}$  given by  $t \mapsto 1/(1+t)$  is convex, because its second derivative is  $2(1+t)^{-3}$ , which is larger than zero for non-negative  $t$ . Because  $S_i \geq 0$  almost surely, we can apply  
 190 Jensen's inequality to the inner expectation and use (16)

$$\mathbb{E}[\tilde{Y}_i] \geq \mathbb{E}\left[\frac{Y_i^*}{1+\mathbb{E}[S_i \mid \mathbf{Y}^*]}\right] = \mathbb{E}\left[\frac{Y_i^*}{1+S_i^*}\right] = \mathbb{E}[\tilde{Y}_i^*].$$

191 On the other hand, we can use the Taylor formula with intermediate point  $\zeta \in [S_i, S_i^*]$  to expand

$$\frac{1}{1+S_i} = \frac{1}{1+S_i^*} - \frac{S_i - S_i^*}{(1+S_i^*)^2} + \frac{(S_i - S_i^*)^2}{(1+\zeta)^3}. \quad (17)$$

192 Using  $\zeta \geq 0$  to bound the denominator, then multiplying with  $Y_i^*$  and taking the expectation gives

$$\mathbb{E}\left[\frac{Y_i^*}{1+S_i}\right] \leq \mathbb{E}\left[\frac{Y_i^*}{1+S_i^*}\right] + \mathbb{E}[Y_i^*(S_i - S_i^*)^2]. \quad (18)$$

193 The variance term can be calculated by substituting  $S_i$  and  $S_i^*$ , expanding the sum, and using the  
 194 independence of  $M$  to show that the mixed terms are zero:

$$\begin{aligned} \mathbb{E}[Y_i^*(S_i - S_i^*)^2] &= \mathbb{E}\left[Y_i^* \left(\sum_{j \neq i} Y_j^* \left(\frac{M_j}{p_j} - 1\right)\right)^2\right] \\ &= \sum_{j \neq i} \mathbb{E}\left[Y_i^*(Y_j^*)^2 \left(\frac{M_j}{p_j} - 1\right)^2\right] + \sum_{j \neq i} \sum_{k \notin \{i, j\}} \mathbb{E}[Y_i^* Y_j^* Y_k^*] \mathbb{E}\left[\frac{M_j}{p_j} - 1\right] \mathbb{E}\left[\frac{M_k}{p_k} - 1\right] \\ &= \sum_{j \neq i} \mathbb{E}[Y_i^* Y_j^*] \mathbb{E}\left[\frac{M_j}{p_j^2} - 2\frac{M_j}{p_j} + 1\right] = \sum_{j \neq i} \left(\frac{1-p_j}{p_j}\right) \mathbb{E}\left[\frac{Y_i}{p_i} \cdot \frac{Y_j}{p_j}\right]. \quad \square \quad (19) \end{aligned}$$

195 Note that the transformation of equation (3) was  
 196 crucial for this calculation, because it makes the  
 197 mask variables in the numerator and denomina-  
 198 tor independent.

199 In practice, most entries of the co-occurrence  
 200 matrix  $\mathbb{E}[Y_i \cdot Y_j]$  will be extremely small, caus-  
 201 ing only a minute contribution to the error bound. This can be illustrated by calculating, on two real  
 202 datasets, the upper-bound for the error of the proposed estimator, by approximating  $\mathbb{E}[Y_i \cdot Y_j]$  with  
 203 the label co-occurrence frequency. The propensities are estimated as in Jain et al. [16]. Looking at  
 204 the mean value, and the worst case for any label (Table 1), We can see that the error on average is  
 205 very small, indicating that the worst-case bound only applies to very few labels.

206 **Corollary 1** (PAL Upper-Bound). *Under the assumptions of Theorem 2, if the underlying multiclass*  
 207 *loss  $\ell_{MC}$  is a non-negative convex function, the expression*

$$\tilde{\ell}(\mathbf{y}, \hat{\mathbf{y}}) := \sum_{j=1}^l \frac{y_j/p_j}{1 + \sum_{j \neq i} y_j/p_j} \ell_{MC}(j, \hat{\mathbf{y}}) \quad (20)$$

208 *gives a nonnegative, convex upper-bound on the true normalized PAL loss in expectation.*

## 209 5 Experimental Results

210 In this section we present some empirical evidence that illustrates the influence of missing labels and  
 211 the unbiased estimates and upper bounds on overfitting and bias-variance trade-off. Additional results  
 212 and a more detailed description of the procedure can be found in the appendix.

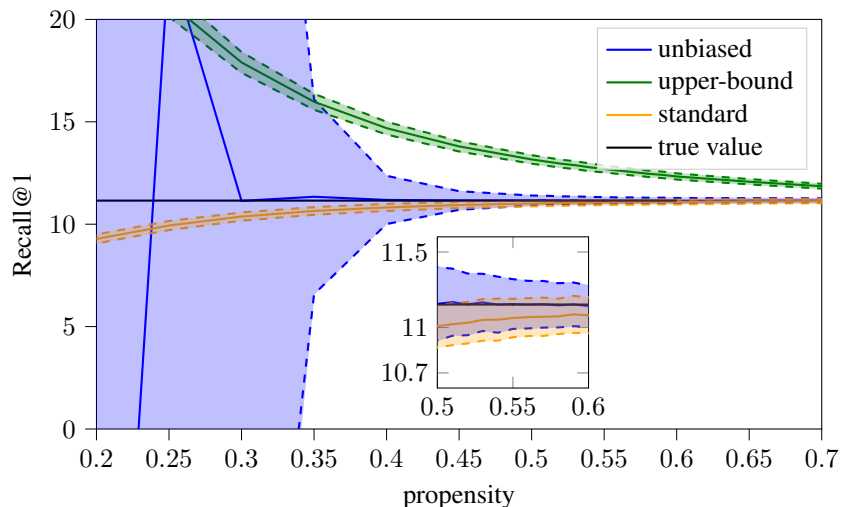


Figure 1: Unbiased estimate of per-example recall with artificial data as described in the main text. The shaded region corresponds to one standard deviation, estimated over 100 repetitions. The black line denotes the true recall.

213 **Prediction Setting** First, we want to demonstrate the variance problem in a simple prediction  
 214 setting, where the classifier is fixed and we want to determine its performance. Consider a setting  
 215 in which there are 100 different labels, which are independent and each has a probability of 10%.  
 216 We randomly draw 10 000 ground-truth label vectors, and generate observed labels by removing  
 217 according to a propensity  $p$  that is identical for all labels. The predictions are generated by randomly  
 218 choosing a label from the ground-truth. We calculate the average per-example recall using the  
 219 standard estimator, the unbiased estimator, and the upper bound, and plot the results in Figure 1.

220 As can be seen, for moderate propensities the unbiased estimator works well, but for propensities  
 221 below 0.45 the 10 000 samples are not sufficient to get an accurate estimate. In this setting, the  
 222 upper-bound results in a larger error than using the standard estimator.

223 **Training Setting** Ideally, we would benchmark our loss functions on a real XMC task. However,  
 224 for those we neither know the exact propensities, nor can we validate that the unbiased estimates and  
 225 upper bounds produce reasonable results, since the fully-labeled ground truth is unknown.

226 Instead of using fully artificial data, we chose to construct a dataset based on existing data: We took  
 227 AmazonCat-13k[22] and consider only the 100 most common labels, which are the ones with the  
 228 highest propensity according to Jain et al. [16]. We artificially remove labels according to inverse  
 229 propensity, which increases linearly based on the ordering of label frequencies, such that the most  
 230 common label has an inverse propensity of 2 and the 100th most common one of 20. This process  
 231 partially preserves the strong imbalances that are typical of extreme classification datasets.

232 On this data, we train a linear classifier with  $L_2$ -regularization using different basis loss functions  
 233 with **a)** the original (standard) loss on clean training data and **b)** noisy training data, as well as **c)** the  
 234 unbiased version and **d)** the upper-bound version on noisy data. For each training run, we evaluate  
 235 the loss on noisy and clean training and test data. For the evaluation on noisy data, the corresponding  
 236 unbiased estimators are used.

237 In this linear-classifier experiment, the noise-pattern overfitting is much stronger than the overfitting  
 238 due to finite sampling (10). Figure 2 shows this for the case of the BCE loss in OvA-reduction and  
 239 CCE loss in normalized PAL reduction. For the classifier trained on clean data (blue), the weights are  
 240 independent of the noise pattern and thus the dashed and dotted lines coincide in expectation. For  
 241 the case of OvA reduction using the BCE loss, the training loss gets reduced much further using the  
 242 unbiased loss function or the upper-bound loss function than using the standard loss. This decrease  
 243 more than compensates the increase in generalization gap, and as such the minimal test loss is better  
 244 with these two variants of the loss function. In contrast, in the non-decomposable case, even though

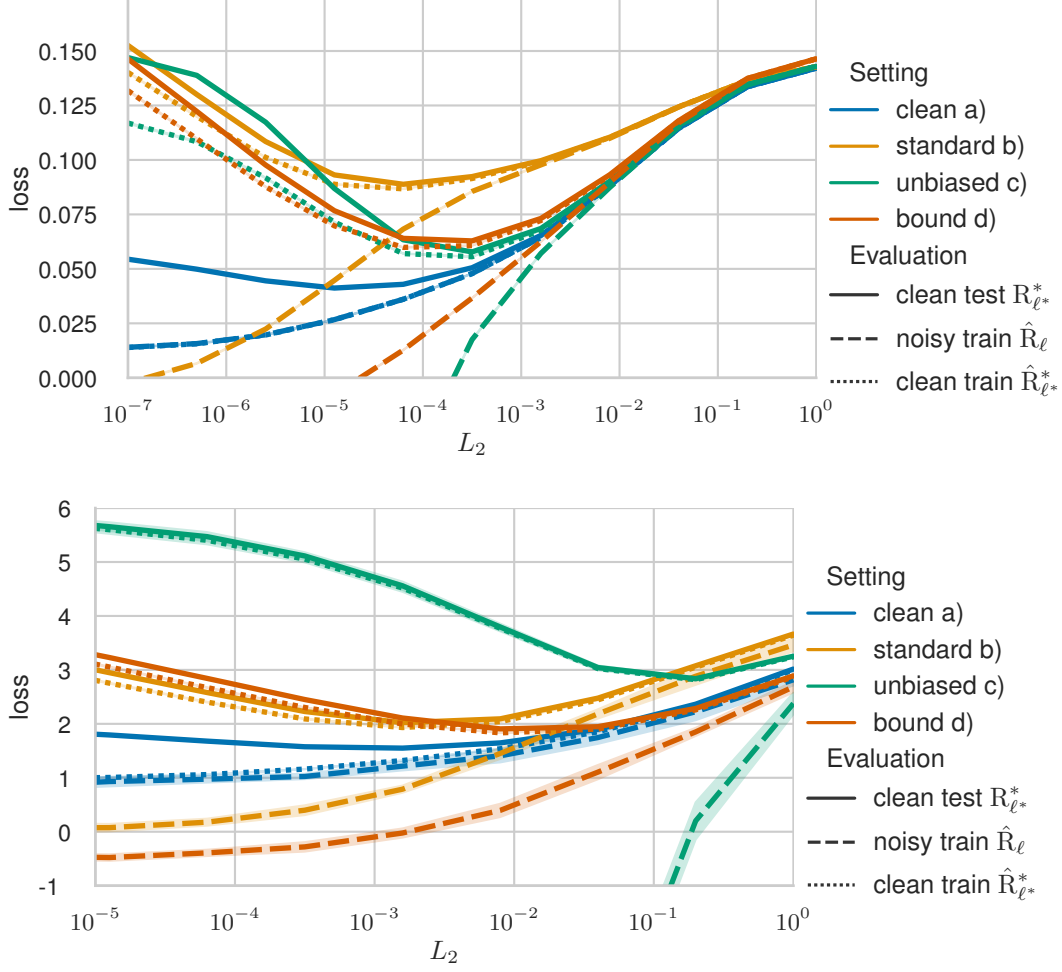


Figure 2: Binary cross-entropy (top) and normalized categorical cross-entropy (bottom) for different regularization strengths, evaluated on noisy training data, clean training data, and clean test data. The gaps between dashed and dotted lines correspond to overfitting to the noise pattern, the smaller gaps between dotted and solid lines show the generalization gaps due to the finite training sample. As the dashed lines are for noisy data, they are calculated using the unbiased estimate (6).

245 the observed training loss decreases drastically with the unbiased loss, the increase in overfitting  
 246 makes the test loss worse than using the biased standard loss function.

247 In this case, using the upper-bound (20) can mitigate the effect, though there is still significant  
 248 overfitting, as evidenced by the estimated training loss being less than zero. This is possible because  
 249 even though the loss we use for training is a non-negative upper bound on the expected unbiased loss,  
 250 the dashed curves show the value estimated for the loss using the unbiased estimator, which can be  
 251 negative due to overfitting. For the OVA case, the upper bound (13) also reduces overfitting, but does  
 252 not result in an overall better classifier on test data.

253 In terms of the bias-variance trade-off, the graphs show a clear trend: The optimal regularization  
 254 for training on noisy data is larger than on clean data. It is also larger when using the unbiased or  
 255 upper-bound loss as compared to standard loss. This is as expected from the variance analysis and  
 256 generalization bound presented in the theory.



## 257 6 Related Work

258 **Unbiased Estimates for Noisy Labels** Learning with missing labels is a specific instance of  
259 learning with class-conditional noise. For the case of binary labels, unbiased estimates of the loss  
260 function can be found in Natarajan et al. [26]. A more general approach is given in Van Rooyen and  
261 Williamson [35]. In their notation,  $f$  is a function and  $\mathbb{P}$  the probability distribution over clean data,  
262 that is transformed by the invertible operator  $T$  into a *corrupted* probability distribution. Let  $R$  be the  
263 inverse of  $T$ , and  $R^*$  its adjoint, then  $\langle \mathbb{P}, f \rangle = \langle R \circ T(\mathbb{P}), f \rangle = \langle T(\mathbb{P}), R^*(f) \rangle$ . This equation forms  
264 the basis for their “Theorem 5 (Corruption Corrected Loss)”, which states that a *corruption corrected*  
265 function  $l_R$  is given  $\forall a \in \mathcal{A}$  by  $l_R(\cdot, a) = R^*(l(\cdot, a))$ , where  $\mathcal{A}$  denotes the set of possible actions  
266 that will be evaluated by the loss functions. For a finite label space with  $n$  possible, the operator  
267  $R^*$  can be represented with an  $n \times n$  matrix. For the multilabel case here, applying this naively  
268 would require  $2^l$  evaluations of the original loss function. In contrast, the direct approach presented  
269 in section 3 is much more efficient.

270 **Alternatives** In some settings with noisy labels, it is possible to use a learning algorithm that is  
271 inherently noise tolerant [12, 36]. Certain performance objectives such as the balanced error or the  
272 AUC are noise robust even under the more general setting of mutually contaminated distributions as  
273 shown in Menon et al. [23]. A data re-calibration approach tries to identify from the training data  
274 which samples are corrupted, e.g. by looking at samples for which the network is very unsure, and  
275 adapt the training process correspondingly [13, 42, 19] It is also possible to first train a scorer on the  
276 noisy data naively, from which a classifier adapted to a given rate of missing labels can be constructed  
277 by choosing an appropriate threshold [23]. Similarly, the inference procedure of PLTs can be adapted  
278 to take into account a propensity model [39].

279 **Related Learning Settings** Learning with missing labels is highly related to learning from positive  
280 and unlabeled (PU) data [11]. An unbiased loss function for this setting is given in Du Plessis et al.  
281 [10]. The appearing difficulties, that non-negativity and convexity need not be preserved, are the same  
282 as in our setting [21]. A slightly different setting with missing labels is given by semi-supervised  
283 learning, where it is known for which labels are missing [41].

## 284 7 Summary and Discussion

285 This paper provides unbiased estimates for four cases of multilabel reductions given in Menon  
286 et al. [24]. Except for the PAL reduction, these estimators can be non-convex and even negatively  
287 unbounded. The unbiased estimates come with an increase in variance. This is unavoidable if  
288 unbiasedness is required, as the estimators can be shown to be unique. If sufficient training data is  
289 available, then the unbiased loss functions can be used, but for the normalized reductions we found  
290 that even 1.2 million instances in AmazonCat are not enough. Much fewer data points are needed  
291 in order to estimate the overall loss of a classifier. This is because for training, an accurate estimate  
292 for  $\mathbb{E}[\ell(Y^*, h(X)) | X]$  needs to be formed, whereas for evaluation this is averaged over the entire  
293 dataset,  $\mathbb{E}[\ell(Y^*, h(X))]$ . This indicates that the unbiased estimates can be useful for hyperparameter  
294 tuning and model selection.

295 For training, however, another approach is needed. A method that fixes the negative unboundedness  
296 and non-convexity and also reduces the variance is to switch to a convex upper-bound. We have  
297 shown that this can stabilize the training and improve the results.

298 Furthermore, the data in section 5 suggest training with missing labels requires more regularization,  
299 irrespective of whether training uses standard-, unbiased-, or convex upper-bound losses. Our findings  
300 agree with Arpit et al. [2] who found that typical regularizers prevent a deep network from memorizing  
301 *noisy* examples, while not hindering the learning of patterns from *clean* instances.

302 All in all, our results show that a) unbiasedness can be achieved for generic multilabel losses, and  
303 in particular the losses resulting from multilabel reduction, but also that b) these losses might not  
304 be suitable for optimization. We have presented one method that trades away unbiasedness for the  
305 ability to handle training with lower amounts of data. An exciting future research prospect would be  
306 to investigate families of loss functions that can continuously trade off bias and variance, and thus  
307 allow for optimal training with different amounts of available data.

## References

- [1] R. Agrawal, A. Gupta, Y. Prabhu, and M. Varma. Multi-label learning with millions of labels: Recommending advertiser bid phrases for web pages. In *Proceedings of the 22nd International Conference on World Wide Web*, WWW '13, page 13–24, New York, NY, USA, 2013. Association for Computing Machinery.
- [2] D. Arpit, S. Jastrzębski, N. Ballas, D. Krueger, E. Bengio, M. S. Kanwal, T. Maharaj, A. Fischer, A. Courville, Y. Bengio, and S. Lacoste-Julien. A Closer Look at Memorization in Deep Networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 233–242. PMLR, 06–11 Aug 2017.
- [3] R. Babbar and B. Schölkopf. Dismec: Distributed sparse machines for extreme multi-label classification. In *WSDM*, pages 721–729, 2017.
- [4] K. Bhatia, K. Dahiya, H. Jain, Y. Prabhu, and M. Varma. The extreme classification repository: Multi-label datasets and code. <http://manikvarma.org/downloads/XC/XMLRepository.html>, 2016.
- [5] Y.-T. Chou, G. Niu, H.-T. Lin, and M. Sugiyama. Unbiased Risk Estimators Can Mislead: A Case Study of Learning with Complementary Labels. In *International Conference on Machine Learning*, pages 1929–1938. PMLR, Nov. 2020.
- [6] K. Dahiya, A. Agarwal, D. Saini, K. Gururaj, J. Jiao, A. Singh, S. Agarwal, P. Kar, and M. Varma. Siamesexml: Siamese networks meet extreme classifiers with 100m labels. In *Proceedings of the International Conference on Machine Learning*, July 2021.
- [7] O. Dekel and O. Shamir. Multiclass-multilabel classification with more classes than examples. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 137–144, 2010.
- [8] K. Dembczynski, W. Cheng, and E. Hüllermeier. Bayes optimal multilabel classification via probabilistic classifier chains. In *ICML*, 2010.
- [9] J. Deng, A. C. Berg, K. Li, and L. Fei-Fei. What does classifying more than 10,000 image categories tell us? In *ECCV*, 2010.
- [10] M. Du Plessis, G. Niu, and M. Sugiyama. Convex formulation for learning from positive and unlabeled data. In *International conference on machine learning*, pages 1386–1394, 2015.
- [11] C. Elkan and K. Noto. Learning classifiers from only positive and unlabeled data. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '08, pages 213–220, New York, NY, USA, Aug. 2008. Association for Computing Machinery.
- [12] A. Ghosh, N. Manwani, and P. S. Sastry. Making risk minimization tolerant to label noise. *Neurocomputing*, 160:93–107, 2015. Publisher: Elsevier.
- [13] B. Han, Q. Yao, X. Yu, G. Niu, M. Xu, W. Hu, I. Tsang, and M. Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- [14] J. Huang, H. Oosterhuis, M. De Rijke, and H. Van Hoof. Keeping dataset biases out of the simulation: A debiased simulator for reinforcement learning based recommender systems. In *Fourteenth ACM conference on recommender systems*, pages 190–199, 2020.
- [15] J. Huang, H. Oosterhuis, and M. de Rijke. It is different when items are older: Debiasing recommendations when selection bias and user preferences are dynamic. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, pages 381–389, 2022.
- [16] H. Jain, Y. Prabhu, and M. Varma. Extreme multi-label loss functions for recommendation, tagging, ranking and other missing label applications. In *KDD*, August 2016.

- 356 [17] H. Jain, V. Balasubramanian, B. Chunduri, and M. Varma. Slice: Scalable linear extreme  
357 classifiers trained on 100 million labels for related searches. In *WSDM*, pages 528–536, 2019.
- 358 [18] Y. Jernite, A. Choromanska, and D. Sontag. Simultaneous learning of trees and representations  
359 for extreme classification and density estimation. In D. Precup and Y. W. Teh, editors, *Proceed-*  
360 *ings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of*  
361 *Machine Learning Research*, pages 1665–1674. PMLR, 06–11 Aug 2017.
- 362 [19] L. Jiang, Z. Zhou, T. Leung, L.-J. Li, and L. Fei-Fei. MentorNet: Learning data-driven  
363 curriculum for very deep neural networks on corrupted labels. In J. Dy and A. Krause, edi-  
364 tors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of  
365 *Proceedings of Machine Learning Research*, pages 2304–2313. PMLR, 10–15 Jul 2018.
- 366 [20] T. Joachims, A. Swaminathan, and T. Schnabel. Unbiased learning-to-rank with biased feedback.  
367 In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*,  
368 pages 781–789, 2017.
- 369 [21] R. Kiryo, G. Niu, M. C. du Plessis, and M. Sugiyama. Positive-Unlabeled Learning with  
370 Non-Negative Risk Estimator. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus,  
371 S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*  
372 *30*, pages 1675–1685. Curran Associates, Inc., 2017.
- 373 [22] J. McAuley and J. Leskovec. Hidden factors and hidden topics: understanding rating dimensions  
374 with review text. In *Proceedings of the 7th ACM conference on Recommender systems*, pages  
375 165–172, 2013.
- 376 [23] A. Menon, B. Van Rooyen, C. S. Ong, and B. Williamson. Learning from corrupted binary  
377 labels via class-probability estimation. In *International Conference on Machine Learning*, pages  
378 125–134, 2015.
- 379 [24] A. K. Menon, A. S. Rawat, S. Reddi, and S. Kumar. Multilabel reductions: what is my loss  
380 optimising? *Advances in Neural Information Processing Systems*, 32, 2019.
- 381 [25] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in  
382 vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- 383 [26] N. Natarajan, I. S. Dhillon, P. Ravikumar, and A. Tewari. Cost-sensitive learning with noisy  
384 labels. *The Journal of Machine Learning Research*, 18(1):5666–5698, 2017. Publisher: JMLR.  
385 org.
- 386 [27] H. Oosterhuis and M. de Rijke. Policy-aware unbiased learning to rank for top-k rankings. In  
387 *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in*  
388 *Information Retrieval*, pages 489–498, 2020.
- 389 [28] I. Partalas, A. Kosmopoulos, N. Baskiotis, T. Artieres, G. Paliouras, E. Gaussier, I. Androut-  
390 sopoulos, M.-R. Amini, and P. Galinari. Lshtc: A benchmark for large-scale text classification.  
391 *arXiv preprint arXiv:1503.08581*, 2015.
- 392 [29] Y. Prabhu and M. Varma. Fastxml: A fast, accurate and stable tree-classifier for extreme  
393 multi-label learning. In *KDD*, pages 263–272. ACM, 2014.
- 394 [30] M. Qaraei, E. Schultheis, P. Gupta, and R. Babbar. Convex Surrogates for Unbiased Loss  
395 Functions in Extreme Classification With Missing Labels. In *Proceedings of the Web Conference*  
396 *2021*, pages 3711–3720, Ljubljana Slovenia, Apr. 2021. ACM.
- 397 [31] S. J. Reddi, S. Kale, F. Yu, D. Holtmann-Rice, J. Chen, and S. Kumar. Stochastic negative  
398 mining for learning with large output spaces. In K. Chaudhuri and M. Sugiyama, editors,  
399 *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and*  
400 *Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 1940–1949. PMLR,  
401 16–18 Apr 2019.
- 402 [32] N. Sachdeva, C.-J. Wu, and J. McAuley. On sampling collaborative filtering datasets. In  
403 *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*,  
404 *WSDM ’22*, page 842–850, New York, NY, USA, 2022. Association for Computing Machinery.

- 405 [33] P. Teisseyre, J. Mielniczuk, and M. Łazęcka. Different strategies of fitting logistic regression  
406 for positive and unlabelled data. In V. V. Krzhizhanovskaya, G. Závodszy, M. H. Lees, J. J.  
407 Dongarra, P. M. A. Sloot, S. Brissos, and J. Teixeira, editors, *Computational Science – ICCS*  
408 *2020*, pages 3–17, Cham, 2020. Springer International Publishing. ISBN 978-3-030-50423-6.
- 409 [34] G. Tsoumakas and I. M. Katakis. Multi-label classification: An overview. *Int. J. Data Warehous.*  
410 *Min.*, 3:1–13, 2007.
- 411 [35] B. Van Rooyen and R. C. Williamson. A theory of learning with corrupted labels. *The Journal*  
412 *of Machine Learning Research*, 18(1):8501–8550, 2017. ISSN 1532-4435.
- 413 [36] B. Van Rooyen, A. Menon, and R. C. Williamson. Learning with symmetric label noise: The  
414 importance of being unhinged. In *Advances in Neural Information Processing Systems*, pages  
415 10–18, 2015.
- 416 [37] X. Wu, H. Chen, J. Zhao, L. He, D. Yin, and Y. Chang. Unbiased learning to rank in feeds  
417 recommendation. In *Proceedings of the 14th ACM International Conference on Web Search*  
418 *and Data Mining*, pages 490–498, 2021.
- 419 [38] M. Wydmuch, K. Jasinska, M. Kuznetsov, R. Busa-Fekete, and K. Dembczynski. A no-regret  
420 generalization of hierarchical softmax to extreme multi-label classification. In S. Bengio,  
421 H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in*  
422 *Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.
- 423 [39] M. Wydmuch, K. Jasinska-Kobus, R. Babbar, and K. Dembczynski. *Propensity-Scored Proba-*  
424 *bilistic Label Trees*, page 2252–2256. Association for Computing Machinery, New York, NY,  
425 USA, 2021.
- 426 [40] R. You, Z. Zhang, Z. Wang, S. Dai, H. Mamitsuka, and S. Zhu. Attentionxml: Label tree-based  
427 attention-aware deep model for high-performance extreme multi-label text classification. In  
428 H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors,  
429 *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 430 [41] H.-F. Yu, P. Jain, P. Kar, and I. Dhillon. Large-scale multi-label learning with missing labels. In  
431 *International conference on machine learning*, pages 593–601, 2014.
- 432 [42] S. Zheng, P. Wu, A. Goswami, M. Goswami, D. Metaxas, and C. Chen. Error-bounded correction  
433 of noisy labels. In *International Conference on Machine Learning*, pages 11447–11457. PMLR,  
434 2020.

## 435 Checklist

- 436 1. For all authors...
- 437 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's  
438 contributions and scope? [Yes]
- 439 (b) Did you describe the limitations of your work? [Yes]
- 440 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 441 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
442 them? [Yes]
- 443 2. If you are including theoretical results...
- 444 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 445 (b) Did you include complete proofs of all theoretical results? [Yes] Some proofs are in  
446 the appendix
- 447 3. If you ran experiments...
- 448 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
449 mental results (either in the supplemental material or as a URL)? [Yes]
- 450 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
451 were chosen)? [Yes]

- 452 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
453 ments multiple times)? [Yes]
- 454 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
455 of GPUs, internal cluster, or cloud provider)? [No]
- 456 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 457 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 458 (b) Did you mention the license of the assets? [No]
- 459 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 460
- 461 (d) Did you discuss whether and how consent was obtained from people whose data you're  
462 using/curating? [N/A]
- 463 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
464 information or offensive content? [N/A]
- 465 5. If you used crowdsourcing or conducted research with human subjects...
- 466 (a) Did you include the full text of instructions given to participants and screenshots, if  
467 applicable? [N/A]
- 468 (b) Did you describe any potential participant risks, with links to Institutional Review  
469 Board (IRB) approvals, if applicable? [N/A]
- 470 (c) Did you include the estimated hourly wage paid to participants and the total amount  
471 spent on participant compensation? [N/A]