Exploring the Algorithm-Dependent Generalization of AUPRC Optimization with List Stability

Anonymous Author(s) Affiliation Address email

Abstract

Stochastic optimization of the Area Under the Precision-Recall Curve (AUPRC) is 1 a crucial problem for machine learning. Although various algorithms have been 2 extensively studied for AUPRC optimization, the generalization is only guaranteed 3 in the multi-query case. In this work, we present the first trial in the single-query 4 generalization of stochastic AUPRC optimization. For sharper generalization 5 bounds, we focus on algorithm-dependent generalization. There are both algorith-6 mic and theoretical obstacles to our destination. From an algorithmic perspective, 7 we notice that the majority of existing stochastic estimators are unbiased only when 8 the sampling strategy is unbiased, and is leave-one-out unstable due to the non-9 decomposability. To address these issues, we propose a sampling-rate-invariant 10 unbiased stochastic estimator with superior stability. On top of this, the AUPRC 11 optimization is formulated as a composition optimization problem, and a stochas-12 tic algorithm is proposed to solve this problem. From a theoretical perspective, 13 standard techniques of the algorithm-dependent generalization analysis cannot 14 be directly applied to such a listwise compositional optimization problem. To 15 fill this gap, we extend the model stability from instancewise losses to listwise 16 losses and bridge the corresponding generalization and stability. Additionally, we 17 construct state transition matrices to describe the recurrence of the stability, and 18 simplify calculations by matrix spectrum. Practically, experimental results on three 19 real-world datasets speak to the effectiveness and soundness of our framework. 20

21 **1 Introduction**

22 Area Under the Precision-Recall Curve (AUPRC) is a widely used metric in the machine learning community, especially in learning to rank, which effectively measures the trade-off between precision 23 and recall of a ranking model. Compared with threshold-specified metrics like accuracy and recall@k, 24 AUPRC reflects a more comprehensive performance by capturing all possible thresholds. In addition, 25 literature has shown that AUPRC is insensitive toward data distributions [20], making it adaptable 26 to largely skewed data. Benefiting from these appealing properties, AUPRC has become one of the 27 standard metrics in various applications, e.g., retrieval [54, 57, 22, 40], object detection [44, 48, 15], 28 medical diagnosis [49, 35], and recommendation systems [16, 71, 1, 63, 2]. 29

30 Over the past decades, the importance of AUPRC has prompted extensive researches on direct

AUPRC optimization. Early work focuses on full-batch optimization [44, 43, 26]. However, in the

³² era of deep learning, the rapidly growing scale of models and data makes these full-batch algorithms

infeasible. Therefore, in recent years, it has raised an increasing favor of the stochastic AUPRC
 optimization [9, 12, 31, 45, 53, 68]. See Appendix A for more on related work.

³⁵ Despite the promoting performance of these methods in various scenarios, the generalization of ³⁶ AUPRC optimization algorithms is still an open problem. Some studies [17, 62] provide provable

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

generalization for AUPRC optimization in information retrieval. In this scene, a dataset consists of 37 multiple queries, where each query corresponds to a set of positive and negative samples. However, 38 these results require sufficient queries to ensure small generalization errors, but leave the single-query 39 case alone, i.e., whether the generalization error tends to zero with the length of a single-query 40 increasing is still unclear. This limits the adaptation scope of these methods. To fill this gap, 41 in this paper we aim to design a stochastic optimization framework for AUPRC with a provable 42

algorithm-dependent generalization performance in the single-query case. 43

The target is challenging in three aspects: (a) Most AUPRC stochastic estimators are biased with 44 a biased sampling rate. Moreover, due to the non-decomposability, outputs of existing algorithms 45 might change a lot with slight changes in the training data, which is called *leave-one-out unstable* 46 in this paper. Such an unstability is harmful to the generalization. (b) The standard framework to 47 analyze the algorithm-dependent generalization requires the objective function to be expressed as a 48 sum of instancewise terms, while AUPRC involves a listwise loss. (c) The stochastic optimization of 49 AUPRC is a two-level compositional optimization problem, which is typically solved by alternate 50 updates. This brings more complicated stability calculations. 51

In search of a solution to (a), we propose a sampling-rate-invariant asymptotically unbiased stochastic 52 estimator based on a reformulation of AUPRC. Notably, to ensure the stability of the estimator, the 53 objective is formulated as a two-level compositional problem. To solve this problem, we propose 54 an algorithm that combines stochastic gradient descent (SGD), linear interpolation and exponential 55 moving average. Error analysis further supports the feasibility of our method, and inspires us to add a 56 semi-variance regularization term. 57

Facing challenge (b), we extend instancewise model stability to listwise model stability, and corre-58 spondingly put forward the generalization via stability of listwise problems. On top of this, we bridge 59 the generalization of AUPRC and the stability of the proposed optimization algorithm. 60

As for challenge (c), the key is to find an upper bound on the variation of model parameters with 61 slight jitter in the dataset. Since the variables to be optimized are typically updated alternately in 62 the compositional optimization problem, we propose state transition matrices of these variables, and 63 simplify the calculations of the stability with matrix spectrum. We also provide the convergence 64 analysis of the proposed method. 65

Last but not least, empirical studies on three real-world datasets further validate the effectiveness and 66 the soundness of the proposed framework. 67

In a nutshell, the main contributions of this paper are summarized as follows: 68

- Algorithmically, a stochastic learning algorithm is proposed for AUPRC optimization. The 69 core of the proposed algorithm is a stochastic estimator which is sampling-rate-invariant 70 asymptotically unbiased. 71
- 72 Theoretically, we present the first trial on the algorithm-dependent generalization of stochas-73 tic AUPRC optimization. To the best of our knowledge, it is also the first work to analyze the stability of stochastic compositional optimization problems. 74

• Technically, we extend the concept of the stability and generalization guarantee to listwise 75 non-convex losses. Then we simplify the stability analysis of compositional objective by matrix spectrum. These techniques might be instructive for other complicated metrics.

Problem Formulation 2 78

76

77

2.1 Preliminaries on AUPRC 79

Notations. Consider a set of N examples $S = \{(x_i, y_i)\}_{i=1}^N$ independently drawn from a sample space $\mathcal{D} = \mathcal{X} \times \mathcal{Y}$, where \mathcal{X} is the input space and $\mathcal{Y} = \{-1, 1\}$ is the label space. For sake of the 80 81 presentation, denote the set of positive examples of S as $S^+ = \{x_i^+\}_{i=1}^{N^+}$, and similarly the set of 82 negative examples is denoted as $S^- = \{x_i^-\}_{i=1}^{N^-}$, where $N^+ = |S^+|, N^- = |S^-|$. With a slight abuse of notation, we also denote $S = S^+ \cup S^-$ if there is no ambiguity. Generally, we assume 83 84 that the dataset is sufficiently large, such that $N^+/(N^+ + N^-) = \mathbb{P}(y = 1) := \pi$. Our target is to 85 learn a score function $h_{w}: \mathcal{X} \to \mathbb{R}$ with parameters $w \in \Omega \subseteq \mathbb{R}^{d}$, such that the scores of positive 86

- 87
- examples are higher than negative examples. Furthermore, when appling the score function to a dataset $S \in \mathcal{X}^N$, we denote $h_{w} : \mathcal{X}^N \mapsto \mathbb{R}^N$, where the k-th element of $h_{w}(S)$ has the top-k values of $\{h_{w}(x) | x \in S\}$. Denote the asymptotic upper bound on complexity as \mathcal{O} , and denote 88

89

asymptotically equivalent as \asymp . 90

In this work, our main interest is to optimize a score function in the view of AUPRC: 91

$$AUPRC(\boldsymbol{w}; \mathcal{D}) = \int_0^1 \mathbb{P}(y = 1 | h_{\boldsymbol{w}}(\boldsymbol{x}) \ge c) \, d \, \mathbb{P}(h_{\boldsymbol{w}}(\boldsymbol{x}) \ge c | y = 1)$$

$$= \int_0^1 \frac{\pi TPR(c)}{\pi TPR(c) + (1 - \pi)FPR(c)} \, d \, \mathbb{P}(h_{\boldsymbol{w}}(\boldsymbol{x}) \ge c | y = 1),$$
(1)

where $(x, y) \sim \mathcal{D}$, c refers to a threshold, and $TPR(c) = \mathbb{P}(h_w(x) \geq c|y=1), FPR(c) = \mathcal{P}(h_w(x) \geq c|y=1)$ 92 $\mathbb{P}(h_w(x) \ge c|y=0)$. For a finite set S, AUPRC is typically approximated by replacing the 93 distribution function $\mathbb{P}(h_w(x) \ge c|y=1)$ with its empirical cumulative distribution function [8, 19]: 94 95

$$\widehat{\text{AUPRC}}(\boldsymbol{w}; \mathcal{S}) = \hat{\mathbb{E}}_{\boldsymbol{x}^+ \sim \mathcal{S}^+} \left[\frac{\pi \widehat{TPR}(h_{\boldsymbol{w}}(\boldsymbol{x}^+))}{\pi \widehat{TPR}(h_{\boldsymbol{w}}(\boldsymbol{x}^+)) + (1 - \pi) \widehat{FPR}(h_{\boldsymbol{w}}(\boldsymbol{x}^+))} \right],$$
(2)

where $\widehat{TPR}(c) = \hat{\mathbb{E}}_{\boldsymbol{x}\sim\mathcal{S}^+} \left[\ell_{0,1}(c-h_{\boldsymbol{w}}(\boldsymbol{x}))\right], \widehat{FPR}(c) = \hat{\mathbb{E}}_{\boldsymbol{x}\sim\mathcal{S}^-} \left[\ell_{0,1}(c-h_{\boldsymbol{w}}(\boldsymbol{x}))\right], \ell_{0,1}(x) = 1$ if 96 $x \leq 0$ or $\ell_{0,1}(x) = 0$ otherwise. It has been shown that AUPRC is an unbiased estimator when 97 $N^+/(N^+ + N^-) \to \pi$ and $N \to \infty$ [8]. With the above estimation, we have the following 98 optimization objective: 99

$$\min_{\boldsymbol{w}} \quad \widehat{\text{AUPRC}}^{\downarrow}(\boldsymbol{w}; \mathcal{S}) = 1 - \widehat{\text{AUPRC}}(\boldsymbol{w}; \mathcal{S}) = \hat{\mathbb{E}}_{\boldsymbol{x}^+ \sim \mathcal{S}^+} \left[\sigma \left(\frac{1 - \pi}{\pi} \cdot \frac{\widehat{FPR}(h_w(\boldsymbol{x}^+))}{\widehat{TPR}(h_w(\boldsymbol{x}^+))} \right) \right], \quad (3)$$

where $\sigma(x) = x/(1+x)$ is concave and monotonically increasing. To make it smooth, surrogate 100 losses ℓ_1, ℓ_2 are used to replace $\ell_{0,1}$ in \widehat{FPR} and \widehat{TPR} respectively, yielding the following surrogate 101 objective: 102

$$\min_{\boldsymbol{w}} f(\boldsymbol{w}; \mathcal{S}) = \hat{\mathbb{E}}_{\boldsymbol{x}^+ \sim \mathcal{S}^+} \left[\sigma \left(\frac{1 - \pi}{\pi} \cdot \frac{\widehat{FPR}(h_w(\boldsymbol{x}^+); \ell_1)}{\widehat{TPR}(h_w(\boldsymbol{x}^+); \ell_2)} \right) \right], \tag{4}$$

where $\widehat{TPR}(c; \ell_2) = \hat{\mathbb{E}}_{\boldsymbol{x} \sim S^+} \left[\ell_2(c - h_{\boldsymbol{w}}(\boldsymbol{x})) \right], \widehat{FPR}(c; \ell_1) = \hat{\mathbb{E}}_{\boldsymbol{x} \sim S^-} \left[\ell_1(c - h_{\boldsymbol{w}}(\boldsymbol{x})) \right]$. Specifically, when $N^+/(N^+ + N^-) = \pi$, it is equivalent to another commonly used formulation Average 103 104 Precision (AP) Loss: 105

$$\widehat{\operatorname{AP}}^{\downarrow}(\boldsymbol{w}; \mathcal{S}) = \mathop{\mathbb{E}}_{\boldsymbol{x}^{+} \sim \mathcal{S}^{+}} \left[\sigma \left(\frac{\sum_{\boldsymbol{x} \sim \mathcal{S}^{-}} \left[\ell_{1}(h_{\boldsymbol{w}}(\boldsymbol{x}^{+}) - h_{\boldsymbol{w}}(\boldsymbol{x})) \right]}{\sum_{\boldsymbol{x} \sim \mathcal{S}^{+}} \left[\ell_{2}(h_{\boldsymbol{w}}(\boldsymbol{x}^{+}) - h_{\boldsymbol{w}}(\boldsymbol{x})) \right]} \right) \right].$$
(5)

2.2 Stochastic Learning of AUPRC 106

Under the stochastic learning framework for instancewise losses, the empirical risk F(w; S) is 107 expressed as a sum of instancewise losses: $F(w; S) = \frac{1}{N} \sum_{x \sim S} \hat{f}(w; x)$, where $\hat{f}(w; x)$ is the 108 stochastic estimator of F(w; S). Different from instancewise losses, listwise losses like AUPRC 109 require a batch of samples to calculate the stochastic estimator. Specifically, at each step, a subset of 110 $\mathcal{S}: z = z^+ \cup z^-$ is randomly drawn, where z^+ consists of n^+ positive examples and z^- consists 111 of n^- negative examples. Then a stochastic estimator of the loss function, denoted as f(w; z), is 112 computed with z. Similar to the instancewise case, we consider a variant of the empirical/population 113 AUPRC risks as approximations, which is a sum of stochastic losses w.r.t. all posible z: 114

$$F(\boldsymbol{w}; \boldsymbol{\mathcal{S}}) = \frac{1}{M} \sum_{\boldsymbol{z}} \hat{f}(\boldsymbol{w}; \boldsymbol{z}), \quad F(\boldsymbol{w}) = \mathbb{E}_{\boldsymbol{\mathcal{S}} \sim \mathcal{D}}[F(\boldsymbol{w}; \boldsymbol{\mathcal{S}})], \quad (6)$$

where M is the number of all possible z. Unfortunately, due to the non-decomposability of the 115 empirical AUPRC risk f(w; S), it is tackle to determine the approximation errors between F(w; S)116 and f(w; S) in general. Nonetheless, in Sec. 3.3 we argue that by selecting proper $\hat{f}(w; z)$, F(w; S)117 can be asymptotically unbiased estimator of f(w; S), which naturally makes F(w) an asymptotically 118 unbiased estimator of 1 – AUPRC. In this case, \hat{f} is said to be an asymptotically unbiased 119 stochastic estimator. Moreover, if the unbiasedness holds under biased sampling rate, it is said 120 to be sampling-rate-invariant asymptotically unbiased. 121

122 **3** Asymptotically Unbiased Stochastic AUPRC Optimization

In this section, we will present our SGD-style stochastic optimization algorithm of AUPRC. In Sec. 3.1, we propose surrogate losses to make the objective function differentiable. In Sec. 3.2, we present details of the proposed stochastic estimator and the corresponding optimization algorithm. Analyses on approximation errors are provided in Sec. 3.3.

127 3.1 Differentiable Surrogate Losses

Since $\ell_{0,1}$ appears in both the numerator and denominator of Eq. (4), simply implementing ℓ_1, ℓ_2 with a single function [55, 9, 53] will bring difficulty to analyze the relationship between $\widehat{AUPRC}^{\downarrow}(w; S)$ and f(w; S). This motivates us to choose $\ell_1 \geq \ell_{0,1}, \ell_2 \leq \ell_{0,1}$, such that $\widehat{AUPRC}^{\downarrow}(w; S) \leq f(w; S)$, thus the original empirical risk could be optimized by minimizing its upper bound f(w; S). Concretely, ℓ_1 and ℓ_2 are defined as the one-side Huber loss and the one-side sigmoid loss:

$$\ell_1(x) = \begin{cases} -2x/\tau_1, & x < 0, \\ (1-x/\tau_1)^2, & 0 \le x < \tau_1, \\ 0, & x \ge \tau_1. \end{cases} \qquad \ell_2(x) = \begin{cases} \frac{\exp(-x/\tau_2) - 1}{\exp(-x/\tau_2) + 1}, & x < 0, \\ 0, & x \ge 0. \end{cases}$$
(7)

Here $\tau_1, \tau_2 > 0$ are hyperparameters. ℓ_1 is convex and decreasing, which ensures the gap between positive-negative pairs is effectively optimized. Additionally, compared with the square loss and the exponential loss, ℓ_1 is more robust to noises. ℓ_2 is Lipschitz continuous, and $\ell_2 \rightarrow \ell_{0,1}$ with $\tau_2 \rightarrow 0$.

136 3.2 Stochastic Estimator of AUPRC

The key to a stochastic learning framework is the design of the stochastic estimator (or the corresponding gradients), *i.e.*, $\hat{f}(w; z)$. Existing methods [9, 72, 12] implement it with $\widehat{AP}^{\downarrow}(w; z)$ (Eq. (5)), which might suffer from two problems:

(P1) Comparing Eq. (4) and Eq. (5), it can be seen that only when $n^+/(n^+ + n^-) \rightarrow \pi$, \widehat{AP}^+ is an asymptotically unbiased estimator. However, it is hardly satisfied since the sampling strategy is usually biased in practice.

(P2) Each term in the summation of $\widehat{AP}^{\downarrow}$ is related to all instances of a batch, leading to weak leave-one-out stability, *i.e.*, changing one instance might result in a relatively large fluctuation in the stochastic gradient, especially when changing a positive example.

To tackle the above problems, we first substitute $\widehat{FPR}(h_w(x^+); \ell_1)$ with $\hat{\mathbb{E}}_{x \sim z^-}[\ell_1(h_w(x^+) - h_w(x))]$, and then introduce an auxiliary vector $v \in \mathbb{R}^{N^+}$ to estimate \widehat{TPR} . Formally, we propose the following batch-based estimator:

$$\hat{f}(\boldsymbol{w};\boldsymbol{z}) = \hat{f}(\boldsymbol{w};\boldsymbol{z},\boldsymbol{v}) = \hat{\mathbb{E}}_{\boldsymbol{x}^{+}\sim\boldsymbol{z}^{+}} \left[\sigma \left(\frac{1-\pi}{\pi} \cdot \frac{\hat{\mathbb{E}}_{\boldsymbol{x}\sim\boldsymbol{z}^{-}} \left[\ell_{1}(h_{\boldsymbol{w}}(\boldsymbol{x}^{+}) - h_{\boldsymbol{w}}(\boldsymbol{x})) \right]}{\hat{\mathbb{E}}_{v\sim\boldsymbol{v}} \left[\ell_{2}(h_{\boldsymbol{w}}(\boldsymbol{x}^{+}) - v) \right]} \right) \right].$$
(8)

Such an estimator enjoys two advantages: in terms of **P1**, it is asymptotically unbiased regardless of the sampling rate (see Sec. 3.3 for detailed discussions); as for **P2**, we use v to substitute $h_w(S^+)$, such that each positive example in a mini-batch only appears in one term. Ideally, it can be considered as using all positive examples in the dataset to estimate \widehat{TPR} instead of that from a mini-batch. With the fact that $n^- \gg n^+$, this makes the corresponding algorithm more stable. Moreover, based on the model stability, generalization bounds are available (see Sec. 4).

155 3.3 Analyses on Approximation Errors

In this subsection, we analyze errors from two approximations in the above algorithm: 1) the gap between F(w; S) and the true AUPRC loss; 2) the gap between the interpolated scores $\phi(h_w(z^+))$ and the true scores $h_w(S^+)$. Proofs are provided in Appendix B.1.



Figure 1: Empirical analysis of estimation errors on simulation data.

- 159 Denote $\pi = N^+/(N^+ + N^-)$ and $\pi_0 = n^+/(n^+ + n^-)$. We would like to show that for all 160 $w \in \Omega, \mathbb{E}_{\boldsymbol{z}}[\hat{f}(\boldsymbol{w}; \boldsymbol{z})]$ is an unbiased estimator when $n \to \infty$, no matter how π_0 is chosen, while for 161 $\mathbb{E}_{\boldsymbol{z}}[\widehat{AP}^{\downarrow}(\boldsymbol{w}; \boldsymbol{z})]$, it holds only when $\pi_0 = \pi$. Since only one model \boldsymbol{w} is considered, we let $\boldsymbol{w}_t = \boldsymbol{w}$
- in the update rule of v (Eq. (10)), and we have the following proposition:
- **Proposition 1.** Consider updating v with Eq. (10) for T steps, then we have

$$\mathbb{E}[\boldsymbol{v}] = \mathbb{E}[\phi(h_{\boldsymbol{w}}(\boldsymbol{z}^{+}))] + (1-\beta)^{T} \left(\boldsymbol{v}_{1} - \mathbb{E}[\phi(h_{\boldsymbol{w}}(\boldsymbol{z}^{+}))]\right), \quad Var[\boldsymbol{v}] \leq Var[\phi(h_{\boldsymbol{w}}(\boldsymbol{z}^{+}))] \cdot \frac{\beta}{2-\beta}$$

- **Remark 1.** Two conclusions could be drawn from the above proposition: first, if the linear interpolation is asymptotically unbiased (see next subsection), by choosing a large T or setting $v_1 = \mathbb{E}[\phi(h_w(z^+))]$, we have $\mathbb{E}[v] \approx h_w(S^+)$; second, by choosing a smaller β , v is more likely to concentrate on $h_w(S^+)$.
- 168 **Proposition 2.** Assume the linear interpolation is asymptotically unbiased. Let $\kappa_1^2 = \hat{\mathbb{E}}_{c \sim h_w(\boldsymbol{z}^+)}[Var_{\boldsymbol{x} \sim \mathcal{S}^-}[\ell_1(c h_w(\boldsymbol{x}))]], \kappa_2^2 = \hat{\mathbb{E}}_{c \sim h_w(\boldsymbol{z}^+)}[Var_{v \sim \boldsymbol{v}}[\ell_2(c v))]].$ When $\kappa_1^2/n^- \rightarrow 0, \kappa_2^2/n^+ \rightarrow 0$, then there exists a positive scale H, such that

$$\hat{\mathbb{E}}_{\boldsymbol{z}\subseteq\mathcal{S}}[\hat{f}(\boldsymbol{w};\boldsymbol{z})] \xrightarrow{P} \widehat{AUPRC}^{\downarrow}(\boldsymbol{w};\mathcal{S}), \quad \hat{\mathbb{E}}_{\boldsymbol{z}\subseteq\mathcal{S}}\left[\widehat{AP}^{\downarrow}(\boldsymbol{w};\boldsymbol{z})\right] \xrightarrow{P} (1 + (\pi_0 - \pi)H) \cdot \widehat{AUPRC}^{\downarrow}(\boldsymbol{w};\mathcal{S}),$$

- where \xrightarrow{P} refers to convergence in probability, and $z \subseteq S$ refers to subsets described in Sec. 2.2.
- 172 **Remark 2.** The above proposition suggests that the proposed batch-based estimator is sampling-

rate-invariant asymptotically unbiased, while $\widehat{AP}^{\downarrow}$ tends to be larger when the sampling rate of the positive class is greater than the prior, and vice versa.

Simulation experiments are conducted as complementary to the theory. Following previous work
[8], the scores are drawn from three types of distributions, including binormal, bibeta and offset
uniform. The results of binormal distribution are visualized in Fig. 1, and detailed descriptions and
more results are available in Appendix B.2. These results are consistent with the above remark.

Next we further study the interpolation error. For the sake of presentation, denote $p : [0, 1] \mapsto \mathbb{R}$ to be an increasing score function describing $h_{\boldsymbol{w}}(\mathcal{S}^+)$, where p(x) is the score in the bottom x-quantile of $h_{\boldsymbol{w}}(\mathcal{S}^+)$. Similarly, let \hat{p} to be the interpolation results of $\mathbb{E}_A[h_{\boldsymbol{w}}(\boldsymbol{z}^+)]$. Assume that $\mathbb{E}_A[h_{\boldsymbol{w}}(\boldsymbol{z}^+)]$ are located in the (i/n^+) -quantiles of p, where $i \in [n^+]$, such that $p(i/n^+) = \hat{p}(i/n^+)$ and all interpolation intervals are with length $1/n^+$. The following proposition provides an upper bound of the approximation error (see [60] for proof):

Proposition 3 (Linear Interpolation Error). Let p, \hat{p} be defined as above. Then we have

$$||p - \hat{p}||_{\infty} \le ||p''||_{\infty} / (8(n^+)^2).$$

Similar to the last subsection, simulation results are shown in Fig. 1(c), which shows the expected errors of linear interpolation are ignorable.

188 3.4 Optimization Algorithm

In the rest of this section, we focus on how to optimize F(w; S). The main challenge is to design update rules for v, such that it could efficiently and effectively approximate $h_w(S^+)$ without fullbatch scanning. To overcome the challenge, we propose an algorithm called **Stochastic Optimization** of AUPRC (SOPRC), which jointly updates model parameters w and the auxiliary vector v. A summary of the detailed process is shown as Alg. 1. At step t, a batch of data is sampled from the training set, and then compute the corresponding scores. Afterward, scores of positive examples are mapped into a N^+ -dimension vector with linear interpolation ϕ as shown in Alg. 2. v_{t+1} are updated with the interpolated scores in a moving average manner.

Practically, n^+ , n^- are finite, causing inevitable estimation errors in $f(w; z_{i_t}, v_{t+1})$. Notice that another factor influencing the stochastic estimation errors, *i.e.*, κ_1^2 and κ_2^2 . To reduce them, it is expected that the variance of positive (negative) scores are small, which motivates us to add a variance regularization term. However, it might force to reduce positive scores that higher than the mean value, which is contrary to our target. Therefore, we propose a **semi-variance regularization term** [4]:

$$\mathcal{L}_{var} = \frac{\lambda_1}{n^+} \sum_{\substack{\boldsymbol{x} \sim \boldsymbol{z}^+ \\ h_{\boldsymbol{w}}(\boldsymbol{x}) < \mu^+}} (h_{\boldsymbol{w}}(\boldsymbol{x}) - \mu^+)^2 + \frac{\lambda_2}{n^-} \sum_{\substack{\boldsymbol{x} \sim \boldsymbol{z}^- \\ h_{\boldsymbol{w}}(\boldsymbol{x}) > \mu^-}} (h_{\boldsymbol{w}}(\boldsymbol{x}) - \mu^-)^2, \tag{9}$$

where $\mu^+ = \frac{1}{n^+} \sum_{\boldsymbol{x} \sim \boldsymbol{z}^+} h_{\boldsymbol{w}}(\boldsymbol{x}), \ \mu^- = \frac{1}{n^-} \sum_{\boldsymbol{x} \sim \boldsymbol{z}^-} h_{\boldsymbol{w}}(\boldsymbol{x}), \ \lambda_1, \ \lambda_2$ are hyperparameters. Finally, we compute the gradients of $f(\boldsymbol{w}; \boldsymbol{z}_{i_t}, \boldsymbol{v}_{t+1}) + \mathcal{L}_{var}$, and update parameters \boldsymbol{w} with gradient descent.

Algorithm 1 SOPRC

Algorithm 2 Score Interpolation $\phi(\cdot)$

Input: Training dataset S, maximum iterations **Input:** A real value vector $\boldsymbol{u} \in \mathbb{R}^n$ where T, learning rate $\{\eta_t\}_{t=1}^T$ and $\{\beta_t\}_{t=1}^T$. $n < N^+$, range of target values [b, B]. **Output:** model parameters w_{T+1} . **Output:** Interpolated vector $\boldsymbol{m} = \phi(\boldsymbol{u})$. 1: Initialize model parameters w_1 and v_1 . 1: Sort *u* in descending order. 2: **for** t = 1 to T **do** 2: Initialize \boldsymbol{m} as $\boldsymbol{0}_{N^+}$, let $u_0 = max(2u_1 - u_2)$ Sample a subset z_{i_t} from S. $u_2, b), u_{n+1} = min(2u_n - 2u_{n-1}, B).$ 3: 3: **for** i = 1 to n **do** 4: **for** $j = \lceil \frac{N^+(i-1)}{n} \rceil$ to $\left\lfloor \frac{N^+ \cdot i}{n} \right\rfloor$ **do** Compute $h_{\boldsymbol{w}_t}(\boldsymbol{z}_{i_t}^+)$ and map the results 4: into $\phi(h_{\boldsymbol{w}_t}(\boldsymbol{z}_{i_t}^+))$ with Alg. 2. 5: Update v with 5: $m_i + = [(i - jn/N^+)u_{i-1}]$ $\boldsymbol{v}_{t+1} = (1 - \beta_t) \boldsymbol{v}_t$ $+(1+jn/N^+-i)u_i]/2$ (10) $+ \beta_t \phi(h_{\boldsymbol{w}_t}(\boldsymbol{z}_{i_t}^+)).$ end for 6: for $j = \left\lceil \frac{N^+ \cdot i}{n} \right\rceil$ to $\lfloor \frac{N^+ \cdot (i+1)}{n} \rfloor$ do Compute \mathcal{L}_{var} with Eq. (9). 7: 6: Update the model parameter: 7: 8: $m_i + = [(i+1-jn/N^+)u_{i-1}]$ $\boldsymbol{w}_{t+1} = \boldsymbol{w}_t - \eta_t \cdot \nabla \mathcal{L}_{var}$ $+(in/N^{+}-i)u_{i}]/2$ (11) $-\eta_t \cdot \nabla f(\boldsymbol{w}_t; \boldsymbol{z}_{i_t}, \boldsymbol{v}_{t+1}).$ end for 9: 10: end for 8: end for

204 Generalization of SOPRC via Stability

In this section, we turn to study the *excess generalization error* of the proposed algorithm. Formally, following standard settings [5], we consider the test error of the model A(S) trained on the training

set S. Our target is to seek an upper bound of the excess error $\mathbb{E}_{A,S}[F(A(S)) - F(\boldsymbol{w}^*)]$, where $\boldsymbol{w}^* \in \arg\min_{\boldsymbol{w} \in \Omega} \mathbb{E}_{A,S}[F(\boldsymbol{w}^*)]$. It can be decomposed as:

$$\mathbf{w}^{*} \in \arg\min_{\mathbf{w}\in\Omega} \mathbb{E}_{A,S}[F(\mathbf{w}^{*})]$$
. It can be decomposed as:
 $\mathbb{E}_{\mathcal{S},A}[F(A(\mathcal{S})) - F(\mathbf{w}^{*})] = \mathbb{E}_{\mathcal{S},A}[F(A(\mathcal{S})) - F(A(\mathcal{S});\mathcal{S})] + \mathbb{E}_{\mathcal{S},A}[F(A(\mathcal{S}))]$

$$\mathcal{S})) - F(\boldsymbol{w}^*)] = \underbrace{\mathbb{E}_{\mathcal{S},A}[F(A(\mathcal{S})) - F(A(\mathcal{S}); \mathcal{S})]}_{Estimation \ Error} + \underbrace{\mathbb{E}_{\mathcal{S},A}[F(A(\mathcal{S}); \mathcal{S}) - F(\boldsymbol{w}^*)]}_{Optimization \ Error}.$$

The estimation error sources from the gap of minimizing the empirical risk instead of the expected risk. In Sec. 4.1, we provide detailed discussion on the estimation error. The optimization error measures the gap between the minimum empirical risk and the results obtained by the optimization algorithm, which will be studied in Sec. 4.2. Detailed proofs of this section are available in Appendix C. Before the formal presentation, we show the main assumptions:

Assumption 1 (Bounded Scores & Gradient). $|\hat{f}(w; \cdot)| \le B$, $\|\nabla \hat{f}(w; \cdot)\|_2 \le G$ for all $w \in \Omega$.

- Assumption 2 (L-Smooth Loss). $\|\nabla \hat{f}(\boldsymbol{w};\cdot) \nabla \hat{f}(\tilde{\boldsymbol{w}};\cdot)\|_2 \leq L \|\boldsymbol{w} \tilde{\boldsymbol{w}}\|_2$ for all $\boldsymbol{w}, \tilde{\boldsymbol{w}} \in \Omega$.
- 216 Assumption 3 (Lipschitz Continuous Functions). $|\ell_1(x) \ell_1(\tilde{x})| \le L_1|x \tilde{x}|, |\ell_2(x) \ell_2(\tilde{x})| \le L_1|x \ell_2(\tilde{x})| \le L_1|x$
- 217 $L_2|x-\tilde{x}|$ for all $x, \tilde{x} \in [-2B, 2B]$. $\|\phi(x) \phi(\tilde{x})\|_2 \le C_{\phi} \|x \tilde{x}\|_2$ for all $x, \tilde{x} \in \mathbb{R}^{N^+}$.

218 4.1 Generalization of AUPRC via Model Stability

The generalization of SGD-style algorithms for instancewise loss has been widely studied with stability measure [38, 21, 28]. However, these results could not be directly applied to listwise losses like AUPRC. The main reason is that the estimation of each stochastic gradient requires a list of examples, and the estimation is usually biased. Nonetheless, to bridge the optimization algorithm and the generalization of AUPRC, we propose a listwise variant of *on-average model stability* [38] as follows:

Definition 1 (Listwise On-average Model Stability). Let $S = \{(x_i, y_i)\}_{i=1}^N$ and $\widetilde{S} = \{(\widetilde{x}_i, y_i)\}_{i=1}^N$ be two sets of examples whose features are drawn independently from \mathcal{X} . For any $i = 1, \dots, N$, denote $S^{(i)} = \{(x_1, y_1), \dots, (x_{i-1}, y_{i-1}), (\widetilde{x}_i, y_i), (x_{i+1}, y_{i+1}), \dots, (x_n, y_n)\}$. A stochastic algorithm A is listwise on-average model (ϵ^+, ϵ^-) -stable if the following condition holds:

$$\mathbb{E}_{\mathcal{S},\widetilde{\mathcal{S}},A}\left[\frac{1}{N^{+}}\sum_{y_{i}=1}\left\|A(\mathcal{S})-A(\mathcal{S}^{(i)})\right\|_{2}\right] \leq \epsilon^{+}, \mathbb{E}_{\mathcal{S},\widetilde{\mathcal{S}},A}\left[\frac{1}{N^{-}}\sum_{y_{i}=-1}\left\|A(\mathcal{S})-A(\mathcal{S}^{(i)})\right\|_{2}\right] \leq \epsilon^{-}.$$

²²⁹ The following theorem shows that the estimation error is bounded by the above-defined stability:

Theorem 1 (Generalization via Model Stability). Let a stochastic algorithm A be listwise onaverage model (ϵ^+, ϵ^-) -stable and Asmp. 1 holds. Then we have

$$\mathbb{E}_{\mathcal{S},A}\left[F(A(\mathcal{S})) - F(A(\mathcal{S});\mathcal{S})\right] \le G(n^+\epsilon^+ + n^-\epsilon^-).$$
(12)

232 With the above theorem, now we only need to focus on the model stability of the proposed algorithm.

Notice that in Alg. 1, both w_t and v_t are updated at each step, thus we have to consider the stability of both simultaneously. The following lemma provides a recurrence for the stability w_t and v_t .

Lemma 1. Let
$$S, \tilde{S}, S^{(i)}$$
 be constructed as Def. 1 and Asmp. 1, 2, 3 hold. Let $\{w_t\}_t$ and $\{w_t^{(i)}\}_t$
be produced by Alg. 1 with S and $S^{(i)}$, respectively. Denote $L = \max\{L_w, L_v/n^+, C_{\phi}B, G/2, B'_{\ell}\},$
 $m^{(i)} = \begin{bmatrix} \|w_v - w^{(i)}\|_{\ell} & \|w_v - w^{(i)}\|_{\ell} \end{bmatrix}^{\top} m^+ = \begin{bmatrix} 1 \\ \sum w^+ - w^{(i)} \end{bmatrix}^{\top} m^- = \begin{bmatrix} \|w_v - w^{(i)}\|_{\ell} & \|w_v - w^{(i)}\|_{\ell} \end{bmatrix}^{\top}$

237
$$\boldsymbol{m}_{t}^{(i)} = \begin{bmatrix} \|\boldsymbol{w}_{t} - \boldsymbol{w}_{t}^{(i)}\|_{2} & \|\boldsymbol{v}_{t} - \boldsymbol{v}_{t}^{(i)}\|_{2} & 1 \end{bmatrix}$$
, $\boldsymbol{m}_{t}^{+} = \frac{1}{N^{\mp}} \sum_{y_{i}=1} \mathbb{E}_{\mathcal{S},A} \begin{bmatrix} \boldsymbol{m}_{t}^{(i)} \end{bmatrix}$, $\boldsymbol{m}_{t}^{-} = \frac{1}{N^{\mp}} \sum_{y_{i}=-1} \mathbb{E}_{\mathcal{S},A} \begin{bmatrix} \boldsymbol{m}_{t}^{(i)} \end{bmatrix}$. Then for all $t \in [T]$, by setting $\beta_{t} \leq 2C_{\phi}B/n^{+}$, we have

$$m_{t+1}^+ \le \frac{I_3 + R_t^+}{N^+} \cdot m_t^+, \quad m_{t+1}^- \le \frac{I_3 + R_t^-}{N^-} \cdot m_t^-,$$
 (13)

where I_3 is the 3×3 identity matrix and

$$R_t^+ = \begin{bmatrix} 2L\eta_t & \frac{L(1-\beta_t)\eta_t}{N^+} & \frac{L\eta_t}{N^+} \\ L\beta_t & 0 & \frac{1}{N^+} \\ 0 & 0 & 0 \end{bmatrix}, R_t^- = \begin{bmatrix} 2L\eta_t & \frac{L_v(1-\beta_t)\eta_t}{N^+} & \frac{L\eta_t \cdot n^+}{N^-} \\ L\beta_t & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$
(14)

Finally, we utilize the matrix spectrum of R_t^+ and R_t^- to show that the model stability w.r.t. Alg. 1 decreases as the number of training examples increases (see Appendix C.2 for details):

Theorem 2. Let $\lambda = LC_{\eta}(1 + \sqrt{1 - \beta^2 + \beta})$, and assumptions in Lem. 1 hold. By setting $\eta_t \leq \frac{C_{\eta}}{t}$, ₂₄₃ $\beta_t = \beta \approx 1/n^+$ and $T \leq N^+$, Alg. 1 is list on-average model stable with

$$\epsilon^{+} = \mathcal{O}\left(\frac{(Tn^{+})^{\frac{\lambda}{\lambda+1}}}{N^{+}}\right), \epsilon^{-} = \mathcal{O}\left(\frac{(Tn^{-})^{\frac{\lambda}{\lambda+1}}}{N^{-}}\right).$$
(15)

244 4.2 Convergence of AUPRC Stochastic Optimization

Following previous work [24, 34], we study the optimization error of the proposed algorithm under the *Polyak-Łojasiewicz (PL)* condition. It has been shown that the PL condition holds for several widely used models including some classes of neural networks [13, 41].

Assumption 4 (Polyak-Łojasiewicz Condition [34, 37]). Denote $w^* = \arg \min_{w \in \Omega} F(w)$. Assume F satisfy the expectation version of PL condition with parameter $\mu > 0$, i.e.,

$$\mathbb{E}_{\mathcal{S}}[F(\boldsymbol{w};\mathcal{S}) - F(\boldsymbol{w}^*)] \leq \frac{1}{\mu} \mathbb{E}_{\mathcal{S}}[\|\nabla F(\boldsymbol{w};\mathcal{S})\|_2^2].$$
(16)

	Stanford Online Products			iNaturalist			PKU VehicleID		
Methods	mAUPRC	R@1	R@10	mAUPRC	R@1	R@4	mAUPRC	R@1	R@5
Contrastive loss [27]	57.73	77.60	89.31	27.99	54.19	71.12	67.26	87.46	94.60
Triplet loss [32]	58.07	78.34	90.50	30.59	60.53	77.62	70.99	90.09	95.54
MS loss [69]	60.10	79.64	90.38	30.28	63.39	78.50	69.15	88.82	95.06
XBM [70]	61.29	80.66	91.08	27.46	59.12	75.18	71.24	92.78	95.83
SmoothAP [9]	61.65	81.13	92.02	33.92	66.13	80.93	72.28	91.31	96.05
DIR [57]	60.74	80.52	91.35	33.51	64.86	79.79	72.72	91.38	96.10
FastAP [12]	57.10	77.30	89.61	31.02	56.64	73.57	70.82	89.42	95.38
AUROC [25]	55.80	77.32	89.64	27.24	60.88	77.76	58.12	81.73	91.92
BlackBox [51]	59.74	79.48	90.74	29.28	56.88	74.10	70.92	90.14	95.52
Ours	62.75	81.91	92.50	36.16	68.22	82.86	74.92	92.56	96.43

Table 1: Quantitative results on SOP, iNaturalist, and VehicleID. All methods are trained with training sets. The best and the second best results are highlighted in **soft red** and **soft blue**, respectively.

The main difference to the existing convergence analysis on non-convex optimization is that the gradient estimation is biased. Nonetheless, we show that the bias terms from Alg. 1 tend to 0 with

²⁵² sufficient training data and training time (see Appendix C.3), leading to the following convergence:

Theorem 3. Let Asmp. 1, 3, 4 hold. By setting $\eta_t = \frac{2t+1}{\mu(t+1)^2}$ and $\beta_t = \beta \approx 1/n^+$, we have

$$\mathbb{E}_{A}[F(\boldsymbol{w}_{T+1}) - F(\boldsymbol{w}^{*})] = \mathcal{O}\left(n^{+}/T + 1/N^{+}\right).$$
(17)

Theorem 4. Let assumptions in Thm. 2 and 3 hold. By setting $T \simeq (N^+)^{\frac{\lambda+1}{2\lambda+1}} (n^+)^{-\frac{1}{2\lambda+1}}$, we have

$$\mathbb{E}_{\mathcal{S},A}[F(A(\mathcal{S})) - F(\boldsymbol{w}^*)] = \mathcal{O}\left((N^+)^{-\frac{\lambda+1}{2\lambda+1}} \cdot (n^+)^{\frac{3\lambda+1}{2\lambda+1}}\right) + \mathcal{O}\left((N^-)^{-\frac{\lambda+1}{2\lambda+1}} \cdot (n^-)^{\frac{3\lambda+1}{2\lambda+1}}\right).$$
(18)

Remark 3. Recall that $\lambda = LC_{\eta}(1 + \sqrt{1 - \beta^2 + \beta})$ and $C_{\eta} = 4/\mu$, when β is small, we have $\lambda \approx 4L/\mu$. Here L/μ is a condition number determined by the model and surrogate losses. Notice that $n^+ \ll N^+, n^- \ll N^-$, if $\lambda = 1$, the generalization bound is $\mathcal{O}\left((N^+)^{-2/3} \cdot (n^+)^{4/3} + (N^-)^{-2/3} \cdot (n^-)^{4/3}\right)$. As λ increases, it increases to $\mathcal{O}\left((N^+)^{-1/2} \cdot (n^+)^{3/2} + (N^-)^{-1/2} \cdot (n^-)^{3/2}\right)$.

260 5 Experiments

To validate the effectiveness of the proposed method, we conduct empirical studies on the image retrieval task, in which data distributions are largely skewed and AUPRC is commonly used as an evaluation metric. Detailed experimental settings are available in Appendix D.1.

264 5.1 Datasets

We evaluate the proposed method on three image retrieval benchmarks with various domains and scales, including **Stanford Online Products** (**SOP**)¹ [47], **PKU VehicleID**² [42] and **iNaturalist**³ [67]. We follow the official setting to split a test set from each dataset, and then further split the rest into a training set and a validation set by a ratio of 9: 1.

269 5.2 Main Results

We evaluate all methods with *mean AUPRC (mAUPRC)* and *Recall@k.* mAUPRC measures the mean value of the AUPRC over all queries, a.k.a. mean average precision (mAP). The performance comparisons on test sets are shown in Tab. 1. Consequently, we have the following observations: 1) In all datasets, the proposed method surpasses all competitors in the view of mAUPRC, especially in the large-scale long-tailed dataset iNaturalist. This validates the advantages of our method in boosting

¹https://github.com/rksltnl/Deep-Metric-Learning-CVPR16. Licensed MIT.

²https://www.pkuml.org/resources/pku-vehicleid.html. Data files © Original Authors.

³https://github.com/visipedia/inatcomp/tree/master/2018. Licensed MIT.



Figure 2: Qualitative results on iNaturalist. Left most: mean PR curves of different methods. Right two: convergence of different methods and batch sizes in terms of mAUPRC in the validation set.

Table 2: Ablation study over different components of our method on iNaturalist.

	III E (1.1		0.1		D O 1	DO1	D 0 1 (D C 22
No.	Unb. Est.	with \boldsymbol{v}_t	with \mathcal{L}_{var}	Opt.	mAUPRC	R@1	R@4	R@16	R@32
1	X	X	X	SGD	34.58	66.35	81.04	89.80	92.72
2	\checkmark	X	×	SGD	35.84	67.08	81.68	90.17	92.98
3	\checkmark	\checkmark	×	SGD	35.99	67.50	82.03	90.44	93.26
4	\checkmark	\checkmark	\checkmark	SGD	36.16	68.22	82.86	91.02	93.71
5	\checkmark	\checkmark	\checkmark	Adam	36.20	68.48	82.70	90.96	93.63

the AUPRC of models. 2) Compared to pairwise losses, the AUPRC/AP optimization methods enjoy 275 better performance generally. The main reason is that pairwise losses could only optimize models 276 indirectly by constraining relative scores between positive and negative example pairs, while ignoring 277 the overall ranking. 3) Although some pairwise methods like XBM have a satisfying performance on 278 Recall@1, their mAUPRC is relatively low. It is caused by the limitation of Recall@1, *i.e.*, it focuses 279 on the top-1 score while ignoring the ranking of other examples. What's more, this phenomenon 280 shows the inconsistency of Recall@k and AUPRC, revealing the necessity of studying AUPRC 281 optimization. More results are available in Appendix D.2. To qualitatively demonstrate the effect of 282 the proposed method, we also show the mean PR curves and convergence curves in Fig. 2. 283

284 5.3 Ablation Studies

We further investigate the effect of different components of the proposed method. Results are shown in Tab. 2, and more detailed statements and analyses are as follows.

Effect of Unbiased Estimator. To show the performance drop caused by the biased estimator, we replace the prior π in Eq. (8) with $n^+/(n^+ + n^-)$. Comparing line 1 and line 2, using the unbiased estimator increases the mAUPRC by 1.3%, which is consistent with our theoretical results in Sec. 3.3.

Notably, the unbiased estimator is the main source of improvements in terms of mAUPRC.

Effect of v_t . To show the effect of introducing v_t to estimate $\phi(S^+)$, we directly use $\phi(z^+)$ instead in the first two lines. Comparing line 2 and line 3, using v_t could bring consistent improvements due to the better generalization ability.

Effect of \mathcal{L}_{var} . We show that shrinking variances could reduce the batch-based estimation errors.

²⁹⁵ Comparing line 3 and line 4, it can be seen that \mathcal{L}_{var} further boosts the proposed method.

Effect of Optimizer. Comparing line 4 and line 5, it can be seen that the choice of optimizer only has a slight influence.

298 6 Conclusion & Future Work

In this paper, we present a stochastic learning framework for AUPRC optimization. To begin with, we 299 propose a stochastic AUPRC optimization algorithm based on an asymptotically unbiased stochastic 300 estimator. By introducing an auxiliary vector to approximate the scores of positive examples, the 301 proposed algorithm is more stable. On top of this, we study algorithm-dependent generalization. First, 302 we propose list model stability to handle listwise losses like AUPRC, and bridge the generalization 303 and the stability. Afterward, we show that the proposed algorithm is stable, leading to an upper bound 304 of the generalization error. Experiments on three benchmarks validate the advantages of the proposed 305 framework. One limitation is the convergence rate is controlled by the scale of the dataset. In the 306 further, we will consider techniques like variance reduction to improve the convergence rate, and 307 jointly consider the corresponding algorithm-dependent generalization. 308

309 References

- [1] Shilong Bao, Qianqian Xu, Ke Ma, Zhiyong Yang, Xiaochun Cao, and Qingming Huang. Collaborative
 preference embedding against sparse labels. In *ACM International Conference on Multimedia*, pages
 2079–2087, 2019.
- Shilong Bao, Qianqian Xu, Zhiyong Yang, Xiaochun Cao, and Qingming Huang. Rethinking collaborative
 metric learning: Toward an efficient alternative without negative sampling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [3] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural
 results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [4] Shaun A Bond and Stephen E Satchell. Statistical properties of the sample semi-variance. *Applied Mathematical Finance*, 9(4):219–239, 2002.
- [5] Léon Bottou and Olivier Bousquet. The tradeoffs of large scale learning. Advances in Neural Information
 Processing Systems, 20, 2007.
- Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press, 2013.
- [7] Olivier Bousquet and André Elisseeff. Stability and generalization. *Journal of Machine Learning Research*,
 2:499–526, 2002.
- [8] Kendrick Boyd, Kevin H Eng, and C David Page. Area under the precision-recall curve: point estimates
 and confidence intervals. In *ECML PKDD*, pages 451–466. Springer, 2013.
- [9] Andrew Brown, Weidi Xie, Vicky Kalogeiton, and Andrew Zisserman. Smooth-ap: Smoothing the path
 towards large-scale image retrieval. In *European Conference on Computer Vision*, pages 677–694. Springer,
 2020.
- [10] Christopher Burges, Robert Ragno, and Quoc Le. Learning to rank with nonsmooth cost functions.
 Advances in Neural Information Processing Systems, 19:193–200, 2006.
- [11] Christopher JC Burges. From ranknet to lambdarank to lambdamart: An overview. *Learning*, 11(23-581):81, 2010.
- [12] Fatih Cakir, Kun He, Xide Xia, Brian Kulis, and Stan Sclaroff. Deep metric learning to rank. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1861–1870, 2019.
- [13] Zachary Charles and Dimitris Papailiopoulos. Stability and generalization of learning algorithms that
 converge to global optima. In *International Conference on Machine Learning*, pages 745–754. PMLR,
 2018.
- [14] Kean Chen, Jianguo Li, Weiyao Lin, John See, Ji Wang, Lingyu Duan, Zhibo Chen, Changwei He, and
 Junni Zou. Towards accurate one-stage object detection with ap-loss. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5119–5127, 2019.
- [15] Kean Chen, Weiyao Lin, John See, Ji Wang, Junni Zou, et al. Ap-loss for accurate one-stage object
 detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [16] Ting Chen, Yizhou Sun, Yue Shi, and Liangjie Hong. On sampling strategies for neural network-based
 collaborative filtering. In ACM SIGKDD International Conference on Knowledge Discovery and Data
 Mining, pages 767–776, 2017.
- [17] Wei Chen, Tie-Yan Liu, Yanyan Lan, Zhi-Ming Ma, and Hang Li. Ranking measures and loss functions in learning to rank. *Advances in Neural Information Processing Systems*, 22:315–323, 2009.
- [18] Yuansi Chen, Chi Jin, and Bin Yu. Stability and convergence trade-off of iterative optimization algorithms.
 arXiv preprint arXiv:1804.01619, 2018.
- [19] Stéphan Clémençon and Nicolas Vayatis. Nonparametric estimation of the precision-recall curve. In *International Conference on Machine Learning*, pages 185–192, 2009.
- [20] Jesse Davis and Mark Goadrich. The relationship between precision-recall and roc curves. In *International Conference on Machine Learning*, pages 233–240, 2006.

- [21] Andre Elisseeff, Theodoros Evgeniou, Massimiliano Pontil, and Leslie Pack Kaelbing. Stability of
 randomized learning algorithms. *Journal of Machine Learning Research*, 6(1), 2005.
- Martin Engilberge, Louis Chevallier, Patrick Pérez, and Matthieu Cord. Sodeep: a sorting deep net to learn
 ranking loss surrogates. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
 10792–10801, 2019.
- [23] Dylan J Foster, Spencer Greenberg, Satyen Kale, Haipeng Luo, Mehryar Mohri, and Karthik Sridharan.
 Hypothesis set stability and generalization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [24] Dylan J Foster, Ayush Sekhari, and Karthik Sridharan. Uniform convergence of gradients for non-convex
 learning and optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- Wei Gao and Zhi-Hua Zhou. On the consistency of auc pairwise optimization. In *International Conference* on *Machine Learning*, 2015.
- Mark Goadrich, Louis Oliphant, and Jude Shavlik. Gleaner: Creating ensembles of first-order clauses to
 improve recall-precision curves. *Machine Learning*, 64(1-3):231–261, 2006.
- Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping.
 In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, volume 2, pages 1735–1742. IEEE,
 2006.
- Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient
 descent. In *International Conference on Machine Learning*, pages 1225–1234, 2016.
- [29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition.
 In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [30] Kun He, Fatih Cakir, Sarah Adel Bargal, and Stan Sclaroff. Hashing as tie-aware learning to rank. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4023–4032, 2018.
- [31] Paul Henderson and Vittorio Ferrari. End-to-end training of object class detectors for mean average
 precision. In *Asian Conference on Computer Vision*, pages 198–213. Springer, 2016.
- [32] Elad Hoffer and Nir Ailon. Deep metric learning using triplet network. In *International workshop on similarity-based pattern recognition*, pages 84–92. Springer, 2015.
- [33] Qijia Jiang, Olaoluwa Adigun, Harikrishna Narasimhan, Mahdi Milani Fard, and Maya Gupta. Optimizing
 black-box metrics with adaptive surrogates. In *International Conference on Machine Learning*, pages
 4784–4793. PMLR, 2020.
- [34] Hamed Karimi, Julie Nutini, and Mark Schmidt. Linear convergence of gradient and proximal-gradient
 methods under the polyak-łojasiewicz condition. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 795–811. Springer, 2016.
- [35] Joon-myoung Kwon, Youngnam Lee, Yeha Lee, Seungwoo Lee, and Jinsik Park. An algorithm based
 on deep learning for predicting in-hospital cardiac arrest. *Journal of the American Heart Association*,
 7(13):e008678, 2018.
- [36] Yunwen Lei, Antoine Ledent, and Marius Kloft. Sharper generalization bounds for pairwise learning.
 Advances in Neural Information Processing Systems, 33:21236–21246, 2020.
- [37] Yunwen Lei, Mingrui Liu, and Yiming Ying. Generalization guarantee of sgd for pairwise learning.
 Advances in Neural Information Processing Systems, 34, 2021.
- [38] Yunwen Lei and Yiming Ying. Fine-grained analysis of stability and generalization for stochastic gradient
 descent. In *International Conference on Machine Learning*, pages 5809–5819, 2020.
- [39] Jian Li, Xuanyuan Luo, and Mingda Qiao. On generalization error bounds of noisy gradient methods for
 non-convex learning. In *International Conference on Learning Representations*, 2019.
- [40] Zhuo Li, Weiqing Min, Jiajun Song, Yaohui Zhu, Liping Kang, Xiaoming Wei, Xiaolin Wei, and Shuqiang
 Jiang. Rethinking the optimization of average precision: Only penalizing negative instances before positive
 ones is enough. *arXiv preprint arXiv:2102.04640*, 2021.
- [41] Chaoyue Liu, Libin Zhu, and Mikhail Belkin. Loss landscapes and optimization in over-parameterized
 non-linear systems and neural networks. *Applied and Computational Harmonic Analysis*, 2022.

- [42] Hongye Liu, Yonghong Tian, Yaowei Wang, Lu Pang, and Tiejun Huang. Deep relative distance learning:
 Tell the difference between similar vehicles. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2167–2175, 2016.
- [43] Donald Metzler and W Bruce Croft. A markov random field model for term dependencies. In *International* ACM SIGIR Conference on Research and Development in Information Retrieval, pages 472–479, 2005.
- [44] Pritish Mohapatra, CV Jawahar, and M Pawan Kumar. Efficient optimization for average precision svm.
 Advances in Neural Information Processing Systems, 27:2312–2320, 2014.
- [45] Pritish Mohapatra, Michal Rolinek, CV Jawahar, Vladimir Kolmogorov, and M Pawan Kumar. Efficient
 optimization for rank-based loss functions. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3693–3701, 2018.
- [46] Wenlong Mou, Liwei Wang, Xiyu Zhai, and Kai Zheng. Generalization bounds of sgld for non-convex
 learning: Two theoretical viewpoints. In *Conference on Learning Theory*, pages 605–638. PMLR, 2018.
- [47] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured
 feature embedding. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4004–
 4012, 2016.
- [48] Kemal Oksuz, Baris Can Cam, Emre Akbas, and Sinan Kalkan. A ranking-based, balanced loss function
 unifying classification and localisation in object detection. In *Advances in Neural Information Processing* Systems, 2020.
- [49] Brice Ozenne, Fabien Subtil, and Delphine Maucort-Boulch. The precision–recall curve overcame the
 optimism of the receiver operating characteristic curve in rare diseases. *Journal of clinical epidemiology*,
 68(8):855–859, 2015.
- [50] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,
 Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep
 learning library. Advances in Neural Information Processing Systems, 32:8026–8037, 2019.
- [51] Marin Vlastelica Pogančić, Anselm Paulus, Vit Musil, Georg Martius, and Michal Rolinek. Differentiation
 of blackbox combinatorial solvers. In *International Conference on Learning Representations*, 2019.
- [52] Tomaso Poggio and Christian R Shelton. On the mathematical foundations of learning. *American Mathematical Society*, 39(1):1–49, 2002.
- [53] Qi Qi, Youzhi Luo, Zhao Xu, Shuiwang Ji, and Tianbao Yang. Stochastic optimization of areas under
 precision-recall curves with provable convergence. *Advances in Neural Information Processing Systems*,
 34, 2021.
- [54] Tao Qin, Tie-Yan Liu, and Hang Li. A general approximation framework for direct optimization of
 information retrieval measures. *Information Retrieval*, 13(4):375–397, 2010.
- Tao Qin, Xu-Dong Zhang, Ming-Feng Tsai, De-Sheng Wang, Tie-Yan Liu, and Hang Li. Query-level loss
 functions for information retrieval. *Information Processing & Management*, 44(2):838–855, 2008.
- ⁴³⁹ [56] Vijay Raghavan, Peter Bollmann, and Gwang S Jung. A critical investigation of recall and precision as
 ⁴⁴⁰ measures of retrieval system performance. *ACM Transactions on Information Systems*, 7(3):205–229,
 ⁴⁴¹ 1989.
- [57] Jerome Revaud, Jon Almazán, Rafael S Rezende, and Cesar Roberto de Souza. Learning with average
 precision: Training image retrieval with a listwise loss. In *International Conference on Computer Vision*,
 pages 5107–5116, 2019.
- [58] William H Rogers and Terry J Wagner. A finite sample distribution-free performance bound for local discrimination rules. *The Annals of Statistics*, pages 506–514, 1978.
- ⁴⁴⁷ [59] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang,
 ⁴⁴⁸ Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge.
 ⁴⁴⁹ International Journal of Computer Vision, 115(3):211–252, 2015.
- 450 [60] Timothy Sauer. *Numerical analysis*. Addison-Wesley Publishing Company, 2011.
- [61] Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Learnability, stability and
 uniform convergence. *Journal of Machine Learning Research*, 11:2635–2670, 2010.

- [62] Yang Song, Alexander Schwing, Raquel Urtasun, et al. Training deep neural networks via direct loss
 minimization. In *International Conference on Machine Learning*, pages 2169–2177. PMLR, 2016.
- [63] Viet-Anh Tran, Romain Hennequin, Jimena Royo-Letelier, and Manuel Moussallam. Improving collabora tive metric learning with efficient negative sampling. In *International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1201–1204, 2019.
- [64] Evgeniya Ustinova and Victor Lempitsky. Learning deep embeddings with histogram loss. In D. Lee,
 M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- 461 [65] Leslie G Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- 462 [66] Aad W Van der Vaart. Asymptotic statistics, volume 3. Cambridge university press, 2000.
- [67] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro
 Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8769–8778, 2018.
- [68] Guanghui Wang, Ming Yang, Lijun Zhang, and Tianbao Yang. Momentum accelerates the convergence of
 stochastic aupre maximization. *arXiv preprint arXiv:2107.01173*, 2021.
- [69] Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R Scott. Multi-similarity loss with
 general pair weighting for deep metric learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5022–5030, 2019.
- [70] Xun Wang, Haozhi Zhang, Weilin Huang, and Matthew R Scott. Cross-batch memory for embedding
 learning. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6388–6397, 2020.
- [71] Zitai Wang, Qianqian Xu, Ke Ma, Yangbangyan Jiang, Xiaochun Cao, and Qingming Huang. Adversarial
 preference learning with pairwise comparisons. In *ACM International Conference on Multimedia*, pages
 656–664, 2019.
- [72] Fen Xia, Tie-Yan Liu, Jue Wang, Wensheng Zhang, and Hang Li. Listwise approach to learning to rank:
 theory and algorithm. In *International Conference on Machine Learning*, pages 1192–1199, 2008.
- [73] Yisong Yue, Thomas Finley, Filip Radlinski, and Thorsten Joachims. A support vector method for
 optimizing average precision. In *International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 271–278, 2007.

481 Checklist

483

484

485

486 487

488

489

491

492

494

495

496

497

- 482 1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
 - (b) Did you describe the limitations of your work? [Yes] Please see Sec. 6.
 - (c) Did you discuss any potential negative societal impacts of your work? [No] We haven't found any negative societal impact of our work.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 490 2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [Yes]
 - (b) Did you include complete proofs of all theoretical results? [Yes]
- 493 3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes]
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
- (c) Did you report error bars (e.g., with respect to the random seed after running experi ments multiple times)? [No] Because it would be too computationally expensive.

500 501	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
502	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
503	(a) If your work uses existing assets, did you cite the creators? [Yes]
504	(b) Did you mention the license of the assets? [Yes]
505	(c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
506 507	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [No] Assets we used are open source.
508 509	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
510	5. If you used crowdsourcing or conducted research with human subjects
511 512	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
513	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals if applicable? $[N/A]$
515 516	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]