

Probably Approximately Precision and Recall Learning

Lee Cohen* Yishay Mansour† Shay Moran‡ Han Shao§

Abstract

Precision and *Recall* are foundational metrics in machine learning applications where both accurate predictions and comprehensive coverage are essential, such as in recommender systems and multi-label learning. In these tasks, balancing precision (the proportion of relevant items among those predicted) and recall (the proportion of relevant items successfully predicted) is crucial. A key challenge is that one-sided feedback—where only positive examples are observed during training—is inherent in many practical problems. For instance, in recommender systems like YouTube, training data only consists of songs or videos that a user has actively selected, while unselected items remain unseen. Despite this lack of negative feedback in training, it becomes crucial at test time; for example, in a recommender system, it’s essential to avoid recommending items the user would likely dislike.

We introduce a Probably Approximately Correct (PAC) learning framework where each hypothesis is represented by a graph, with edges indicating positive interactions, such as between users and items. This framework subsumes the classical binary and multi-class PAC learning models as well as multi-label learning with partial feedback, where only a single random correct label per example is observed, rather than all correct labels.

Our work uncovers a rich statistical and algorithmic landscape, with nuanced boundaries on what can and cannot be learned. Notably, classical methods like Empirical Risk Minimization fail in this setting, even for simple hypothesis classes with only two hypotheses. To address these challenges, we develop novel algorithms that learn exclusively from positive data, effectively minimizing both precision and recall losses. Specifically, in the realizable setting, we design algorithms that achieve optimal sample complexity guarantees. In the agnostic case, we show that it is impossible to achieve additive error guarantees (i.e., additive regret)—as is standard in PAC learning—and instead obtain meaningful multiplicative approximations.

*Stanford. Email: leecohencs@gmail.com. Authors are ordered alphabetically.

†Tel Aviv University and Google Research. Email: mansour.yishay@gmail.com.

‡Departments of Mathematics, Computer Science, and Data and Decision Sciences, Technion and Google Research.
Email: smoran@technion.ac.il.

§Harvard. Email: han@cmsa.fas.harvard.edu.

1 Introduction

Precision and *Recall* are fundamental metrics in a variety of machine learning applications where accurate prediction and comprehensive coverage are both critical. Recommender systems are one example of machine learning applications where precision and recall are essential. These systems, used in online platforms like streaming services and e-commerce websites, suggest items—such as movies, music, or products—that match user preferences. For instance, Netflix aims to recommend shows that a user will watch and enjoy, based on limited past interactions. Another example is multi-label learning, in which each input corresponds to multiple labels (for instance, an image may contain multiple objects), and the learning objective is to return all labels associated with each input.

A critical aspect of designing such systems is balancing two key metrics:

- Precision- The proportion of recommended items that the user likes. Low *precision loss* means most suggested items are liked by the users.
- Recall- The proportion of items the user would like that are successfully recommended. Low *recall loss* ensures the system does not miss out on suggesting items the user would have liked.

Precision and recall are often at odds. Increasing the number of recommendations can improve recall but may reduce precision. Beyond recommender systems, which will be our running example, these metrics play a vital role in multi-label learning, such as predicting all the objects present in an image, and in platforms like dating apps, where matching users effectively depends on balancing the relevance and variety of suggestions.

In an ideal scenario, we might consider a full information model where, for each user in the training set, the algorithm has access to all the items they like. This setup implicitly includes negative examples, as items not on a user’s list are considered unliked. At test time, the algorithm would then predict a list of recommendations for a given random user. Such a setting aligns well with the standard *Probably Approximately Correct* (PAC) framework [Val84], allowing us to apply standard PAC solutions (e.g., Empirical Risk Minimization).

However, this assumption of full information is often unrealistic in many applications. In real-world scenarios, we typically only observe a small fraction of the items a user liked in the training set, with no explicit information on items the user would not chosen. This setting is better characterized as a partial-information model. For instance, in Spotify, for each observed user in our training data, we might only know the songs they have listened to, without knowing about the vast majority they have not listened to. At test time, given a random user drawn from the same distribution as the training set, we present a full list of recommended items (and not just a single song).

Following standard practice in learning theory, we consider a hypothesis class \mathcal{H} . Each hypothesis in \mathcal{H} is modeled by a finite graph, g . For example, in the context of recommender systems, nodes represent users and items, and an edge between a user and an item indicates that the item is liked by the user (and should therefore be recommended). Our goal is to return a graph that minimizes both precision and recall losses. The recall and precision losses of graph g are measured with respect to the target graph g^{target} , which captures the true set of items liked by each user:

$$\begin{aligned} \ell^{\text{precision}}(g) &:= \mathbb{E}_{x \sim \mathcal{D}} \left[\frac{|N_g(x) \setminus N_{g^{\text{target}}}(x)|}{n_g(x)} \right], \\ \ell^{\text{recall}}(g) &:= \mathbb{E}_{x \sim \mathcal{D}} \left[\frac{|N_{g^{\text{target}}}(x) \setminus N_g(x)|}{n_{g^{\text{target}}}(x)} \right], \end{aligned}$$

where $N_{g'}(x)$ denotes the neighborhood of a node x in a graph g' and $n_{g'}(x) = |N_{g'}(x)|$ denotes the number of neighbors of x .

Note that binary and multiclass PAC learning can be viewed as special cases of our model, where the graphs are bipartite: one side contains nodes representing the inputs x , the other side contains nodes represent the labels y , and an edge (x, y) indicates that y is the label of x . We can also model multi-label PAC learning by allowing multiple edges (x, y) where each y represents one of the (potentially, multiple) labels of x .

We distinguish between two settings- *realizable* and *agnostic*. In the realizable setting, we assume g^{target} is in the class \mathcal{H} , and, aim to find an hypothesis with small precision and recall losses. In the agnostic setting, we do not assume that a perfect hypothesis is in the class. Instead, we aim to compete with the “best” hypothesis (graph) in the class, acknowledging that some error is unavoidable. In the context of the agnostic setting, defining the “best” hypothesis is subtle. One hypothesis might have high precision but low recall, while another has the opposite. Depending on the application’s needs, one may prefer higher precision over recall or vice versa.

This naturally leads us to the concept of *Pareto-loss* objective, which captures the trade-offs between precision and recall along the *Pareto frontier*.¹ Namely, the “best” graphs are on a Pareto frontier, which is the set of graphs where no other graphs are better in both precision and recall simultaneously (see, e.g., Figure 1). Here, we try, given desired precision and recall losses parameters (p, r) to return a graph whose precision and recall losses (p', r') satisfy $p' \lesssim p$ and $r' \lesssim r$. For example, one might aim to optimize precision while keeping recall below a specific threshold (e.g., at most 0.5).

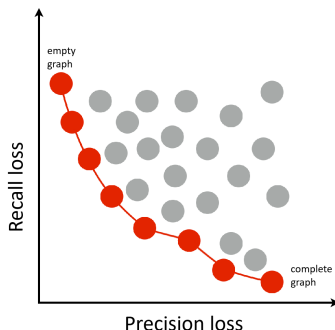


Figure 1: Example of graphs with varying precision and recall losses. Each point is a distinct graph, with red points on the Pareto frontier, showing optimal trade-offs between precision and recall losses. The empty graph always achieves zero precision loss (but has no guarantee on the recall loss), while the complete graph always achieves zero recall loss (but has no guarantee on the precision loss).

Our goal is to design algorithms whose sample complexity is polynomial in the log size of the class and the inverses of the accuracy and confidence parameters, and in some cases, on the maximum degree of the graph (which is arguably small in certain applications). We focus on finite hypothesis classes and aim for sample complexity bounds that depend logarithmically, rather than linearly, on the size of the class. This goal is motivated by standard sample complexity bounds in PAC learning and is particularly relevant when training data is costly to obtain, as is often the case with

¹The term “Pareto” originates from Vilfredo Pareto, an economist who observed that certain distributions followed a pattern where improvements in one dimension often involved trade-offs in another. The [Pareto frontier](#), inspired by this principle, represents optimal trade-offs between competing objectives. Here, our Pareto loss captures the balance between precision and recall, aiming to improve both while acknowledging the inherent trade-off.

human-provided feedback.² However, we aim to avoid dependencies on the number of vertices or edges in the graph.

Our learning problem is significantly more challenging than standard supervised learning tasks due to the absence of negative examples. Namely, we only observe positive examples, making it impossible to estimate precision loss directly. Without knowing what items users dislike, we cannot estimate how many irrelevant items a hypothesis might recommend. This limitation challenges approaches like standard supervised learning, which rely on both positive and negative examples.

In standard supervised learning, a classical solution known as *Empirical Risk Minimization* (ERM) involves finding a hypothesis that best fits the data by minimizing the average loss over all observed examples. However, in our case, the absence of negative examples means that ERM cannot be applied effectively, as there is no way to determine how well a hypothesis avoids irrelevant items. For example, the complete graph (where every possible recommendation is made) is consistent with every training set, since we have no negative examples to contradict it. However, such a hypothesis might have poor precision. Without negative examples in the training set, any hypothesis that covers all observed positive examples appears equally valid, even if it recommends many irrelevant items and incurs a high precision loss. The failure of the ERM principle is not unique to our learning problem; it also occurs in other learning problems, such as multi-class classification [DSBS15], density estimation [DL01; BKM19], and partial concept learning [AHHM22].

In fact, it is not just that ERM would fail; we provide an example in which two hypotheses have nearly the same recall loss but very different precision losses, making it impossible to determine which hypothesis has better precision loss based solely on the training data (see Example 1 for more details).

These challenges necessitate novel approaches to learning and evaluating hypotheses.

Contributions

- **Learning Model:** We propose a learning model that operates under partial information, where the algorithm observes only positive examples. For each user drawn from unknown distribution \mathcal{D} , it randomly observes one item they like. This model reflects real-world constraints and is more practical than assuming access to complete preference profiles as done in multi-label learning.
- **Realizable Setting:** We design algorithms that, given a sample of size

$$O\left(\frac{\log(|\mathcal{H}|/\delta)}{\varepsilon}\right),$$

achieve recall and precision losses of at most ε with probability at least $1 - \delta$. We propose two distinct approaches to achieve this goal: the first circumvents the challenge of estimating precision by minimizing an appropriate surrogate loss, while the second takes a more intuitive approach inspired by the maximum likelihood principle [SB14]. In essence, this second algorithm, when presented with two graphs consistent with the data, prioritizes the one with smaller degrees in a suitably quantified sense.

- **Agnostic Setting:** We demonstrate that achieving a vanishing additive error, as is standard in learning theory, is impossible in this setting by providing lower bounds on the sum of

²While it would be interesting to explore infinite hypothesis classes and characterize them via a combinatorial measure in the style of VC dimension, the findings in [LB24] suggest that such a dimension may not exist in this setting. We leave this as an intriguing open question, as addressing it falls beyond the scope of this work, which focuses on the finite case—a setting that is already challenging.

precision and recall losses, with multiplicative factors greater than 1. In the other direction, we show that constant multiplicative factor guarantees are indeed achievable by adapting our realizable setting algorithms to the agnostic case. Closing the gap between our upper and lower bounds on the best achievable multiplicative factor (5 vs. 1.05) remains an open question. For the Pareto-loss case, we establish both upper and lower bounds for the following question: Given that there exists a hypothesis in the class with precision and recall losses (p, r) , which guarantee pairs (p', r') are achievable? Finally, we pose open questions about determining the optimal factors achievable in the agnostic setting.

Related Work Precision and recall are natural and standard metrics used broadly in machine learning, spanning applications from binary classification [JL19], multi-class classification [GBV20], regression [TR09], and time series [TLZAG18] to information retrieval [AKV16] and generative models [SBLBG18]. Beyond precision-recall, another related metric—the area under the ROC curve (AUC)—has also been extensively studied in the history of binary classification [CM03; CM04; Ros04; AGHHR05], with a focus on generalization. Our work, however, studies a different problem of multi-label learning where the goal is to recommend a list of items to each user. Recommending a list of items has also been addressed in the context of cascading bandits [KSWA15]. However, while our objective is to identify the items that each user likes, their focus is on learning the top K items that are liked by most users. Another feature of our learning problem is that we can only learn from positive examples. PAC learning for binary classification from positive examples has been studied in the literature [Den98; DDGL99; LDG00; BD20].

Multi-label learning [McC99; SS00] has been an area of study in machine learning, with various, primarily experimental approaches (see, e.g., [EW01; PC11; KVJ12] and [ZZ14; BTDK22] for surveys). In multi-label learning, the training set consists of examples, each associated with multiple labels rather than just one. The goal is to train a model that can learn the relationships between the features of each example and all its labels. At test time, the learner predicts a list of labels for new examples, aiming to capture all the relevant labels that apply, rather than just a single one. Some works have examined multi-label learning from a theoretical standpoint, focusing in particular on the Bayes consistency of surrogate losses. Bayes-consistency in multi-label learning ensures that minimizing a surrogate loss also leads to minimizing the true target loss, which is crucial in multi-label settings where optimizing the actual loss is often computationally infeasible as it is non-convex, discrete losses in multi-label settings. Initiated by [GZ11] who first addressed Bayes-consistency for Hamming and ranking losses, showing binary relevance’s consistency with Hamming loss but highlighting ranking loss difficulties. Extensions include rank-based metrics like precision@ κ and recall@ κ [MRRK19], which are loss functions defined under the constraint that the number of labels predicted by the model is limited to κ . Recently, [MMZ24] established H -consistency bounds for multi-label learning, offering stronger guarantees than Bayes-consistency by providing non-asymptotic guarantees that apply to finite number of samples. Our model is inherently more challenging than traditional multi-label learning because our training set consists of examples, each associated with only a single correct label rather than all possible correct labels, with no negative feedback. Yet, at test time, the learner still needs to predict a list of relevant labels for new examples.

2 Model

As is standard in learning theory, we assume a hypothesis class \mathcal{H} of graphs, our goal is to design algorithms whose sample complexity is polynomial in the log size of the class and the inverses of the accuracy and confidence parameters. More specifically, we are given a (possibly huge) set \mathcal{X} of

nodes and a hypothesis class \mathcal{H} of graphs on \mathcal{X} . We denote an unknown target graph g^{target} . The training set consists of a sequence $(x_i, v_i)_{i=1}^m$. The nodes x_1, \dots, x_m are drawn IID from unknown distribution \mathcal{D} . For each node x_i , a random neighbor v_i is drawn uniformly from its neighborhood $N_{g^{\text{target}}}(x_i)$ in the target graph. The algorithm then outputs a graph g^{output} , and the goal is to minimize the expected precision and recall losses:

$$\begin{aligned}\ell^{\text{precision}}(g^{\text{output}}) &:= \mathbb{E}_{x \sim \mathcal{D}} \left[\frac{|N_{g^{\text{output}}}(x) \setminus N_{g^{\text{target}}}(x)|}{n_{g^{\text{output}}}(x)} \right], \\ \ell^{\text{recall}}(g^{\text{output}}) &:= \mathbb{E}_{x \sim \mathcal{D}} \left[\frac{|N_{g^{\text{target}}}(x) \setminus N_{g^{\text{output}}}(x)|}{n_{g^{\text{target}}}(x)} \right],\end{aligned}$$

where for any graph g , $N_g(x)$ denotes the neighborhood of a node x in g and $n_g(x) = |N_g(x)|$ denotes the number of neighbors of x .

We focus on designing learning rules that can compete with the “best” graph in \mathcal{H} . Specifically, if there is a graph $g \in \mathcal{H}$ with precision and recall losses p and r , respectively, can we output a graph g^{output} whose precision and recall losses are comparable to p and r ? To answer this question we consider two natural metrics: scalar loss and Pareto loss.

Scalar-Loss Objective The scalar loss is defined as the average of precision and recall losses³

$$\ell^{\text{scalar}}(g) := \frac{\ell^{\text{precision}}(g) + \ell^{\text{recall}}(g)}{2}.$$

For any $\alpha > 0$, we say α -approximate optimal scalar loss is achievable if there exists a polynomial P such that, for any finite hypothesis class \mathcal{H} , there is an algorithm \mathcal{A} such that the following holds: For any $\varepsilon, \delta > 0$ and any distribution \mathcal{D} , if \mathcal{A} is given an IID training set of size $P(\log|\mathcal{H}|, 1/\varepsilon, 1/\delta)$, with probability at least $1 - \delta$, it outputs a graph with scalar loss satisfying

$$\ell^{\text{scalar}}(g^{\text{output}}) \leq \alpha \cdot \min_{g \in \mathcal{H}} \ell^{\text{scalar}}(g) + \varepsilon.$$

We emphasize that we aim to avoid dependencies on the number of vertices or edges in the graph, as allowing quadratic dependence on the number of nodes trivializes the problem. Our primary focus is on finite hypothesis classes, where we discuss dependencies on the cardinality of the hypothesis class, as is common in standard learning theory.

Pareto-Loss Objective Let $p, p', r, r' \in [0, 1]$. We write $(p, r) \implies (p', r')$ to denote the following statement: there exists a polynomial P such that, for any finite hypothesis class \mathcal{H} , there is an algorithm \mathcal{A} such that the following holds: If \mathcal{D} is a distribution for which there exists a graph in \mathcal{H} with precision and recall losses (p, r) , then for any $\varepsilon, \delta > 0$, if \mathcal{A} is given p, r and an IID training set of size $P(\log|\mathcal{H}|, 1/\varepsilon, 1/\delta)$, with probability at least $1 - \delta$, it outputs a graph with precision and recall losses at most $p' + \varepsilon$ and $r' + \varepsilon$.⁴

We are asking the following a question for each of our losses:

**What is the smallest α such that α -approximate optimal scalar loss is achievable?
Given $p, r \in [0, 1]$, which pairs (p', r') satisfy $(p, r) \implies (p', r')$?**

³The results can be generalized to any weighted sum of precision and recall losses via $w_1 \ell^{\text{precision}}(g) + w_2 \ell^{\text{recall}}(g) \leq 2 \max(w_1, w_2) \ell^{\text{scalar}}(g)$.

⁴Actually, our algorithms only requires the knowledge of r .

This graph learning problem is considerably more challenging than standard supervised learning. If the entire neighborhood were observed in the training set rather than a random neighbor $v_i \sim \text{Unif}(N_{g^{\text{target}}}(x_i))$, the task would reduce to standard supervised learning. However, observing only a random neighbor prevents an unbiased estimate of precision loss, complicating the problem. We demonstrate it in the following example.

Example 1. In Fig 2, in the target graph, user x_i likes a large number n of songs. Our hypothesis class contains two graphs, each with exactly one edge connecting to user x_i . In the red graph g_1 , the edge always connects to a true positive (i.e., the recommended song is one the user likes), while in the blue graph g_2 , it connects to a false positive (i.e., the recommended song is one the user dislikes). Both graphs have high recall loss, $\ell^{\text{recall}}(g_1) = \frac{n-1}{n}$ and $\ell^{\text{recall}}(g_2) = 1$. But graph g_1 has a precision loss of 0 as it always recommends a correct song, while graph g_2 has a precision loss of 1 as it always recommends a wrong song. Since n is large, it is unlikely that any song recommended by either graph will appear in the training set, making it impossible to distinguish between them, despite their drastically different precision losses.

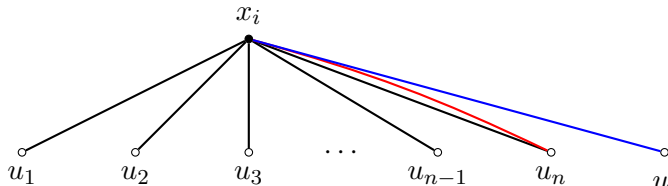


Figure 2: The target graph (black) has neighborhood $N_{g^{\text{target}}}(x_i) = \{u_1, \dots, u_n\}$ where n is huge. The graph g_1 (red) has only one neighbor $u_n \in N_{g^{\text{target}}}(x_i)$ while g_2 (blue) has only one neighbor $u' \notin N_{g^{\text{target}}}(x_i)$.

One might argue that the issue in the above example arises from the large degree of the target graph; however, we will later show that even when the target graph has a small degree, accurately estimating and optimizing precision remains impossible.

We emphasize that, unlike in other supervised learning settings, such as PAC learning, where minimizing empirical risk is often straightforward (e.g., by outputting any classifier that is consistent with the training set), here the learner only observes a single item v_i per user x_i , rather than the user’s entire neighborhood. As a result, minimizing empirical precision loss in this context is far from trivial. For instance, regardless of the target graph, a complete graph is always consistent with the training set but can still incur high precision loss.

3 Main Results

The Realizable Setting We begin by presenting our results in the realizable setting, where the target graph belongs to the hypothesis class. In this case, there is no distinction between optimizing the scalar-loss objective and the Pareto-loss objective. We propose two new algorithms that achieve both the scalar-loss and Pareto-loss objectives.

Theorem 1. In the realizable setting, there exist algorithms such that given an IID training set of size $m \geq O(\frac{\log(|\mathcal{H}|/\delta)}{\varepsilon})$, with probability at least $1 - \delta$, the output graph g^{output} satisfies

$$\ell^{\text{recall}}(g^{\text{output}}) \leq \varepsilon, \quad \ell^{\text{precision}}(g^{\text{output}}) \leq \varepsilon.$$

Since $\mathbb{1}(v_i \notin N_g(x_i))$ is an unbiased estimate of the recall loss $\ell^{\text{recall}}(g)$, any consistent graph (i.e., graph g with $\sum_{i=1}^m \mathbb{1}(v_i \notin N_g(x_i)) = 0$) will have low recall loss. But ERM does not work as the training set contains only positive examples, and a complete graph is consistent with any training set but can incur high precision loss. Hence, the main challenge lies in minimizing the precision loss. One of our proposed algorithms is based on the natural idea of maximum likelihood. At a high level, although multiple graphs may be consistent with the training set, for any observed edge (x_i, v_i) , if it is in graph g , the probability of observing this edge is $\frac{1}{n_g(x_i)}$ when g is the target graph. Consequently, we can rule out the complete graph, as its likelihood of generating any specific observed edge is low. The main challenge in the analysis, then, is how to connect the precision loss to the likelihood. The other algorithm is more directly aligned with the scalar-loss objective. While we cannot obtain an unbiased estimate of the precision loss, and hence the scalar loss, we introduce a surrogate loss that both upper- and lower-bounds the scalar loss within a constant multiplicative factor. Then we output a graph minimizing this surrogate loss.

The Agnostic Setting In the agnostic setting, we show in the next two theorems that it is impossible to achieve an additive error for both scalar-loss and Pareto-loss objectives as is standard in learning theory.

Theorem 2. *There exists a class $\mathcal{H} = \{g_1, g_2\}$ of two graphs, such that for any (possibly randomized improper) algorithm, there exists a target graph g^{target} with bounded degree and a data distribution \mathcal{D} with $\min_{g \in \mathcal{H}} (\ell^{\text{scalar}}(g)) > 0$ s.t. for any sample size $m > 0$, with probability 1 over the training set, the expected (over the randomness of the algorithm) loss of the output g^{output}*

$$\mathbb{E} [\ell^{\text{scalar}}(g^{\text{output}})] \geq 1.05 \cdot \min_{g \in \mathcal{H}} (\ell^{\text{scalar}}(g)).$$

Theorem 3. *There exists a class $\mathcal{H} = \{g_1, g_2\}$ of two graphs, such that for any (possibly randomized improper) algorithm given the knowledge of $(p, r) = (\frac{7}{16}, \frac{1}{4})$, there exists a target graph g^{target} with bounded degree and a data distribution \mathcal{D} for which there exists a graph $g^\dagger \in \mathcal{H}$ with $\ell^{\text{recall}}(g^\dagger) = \frac{1}{4}$ and $\ell^{\text{precision}}(g^\dagger) = \frac{7}{16}$ s.t. for any sample size $m > 0$, with probability 1 over the training set, the expected (over the randomness of the algorithm) precision and recall losses of the output g^{output} satisfy*

$$\mathbb{E} [\ell^{\text{recall}}(g^{\text{output}})] + \frac{12}{5} \mathbb{E} [\ell^{\text{precision}}(g^{\text{output}})] \geq \frac{7}{5}.$$

Remark 1. *Hence the output graph either suffers $\ell^{\text{recall}}(g^{\text{output}}) > \frac{1}{4} = \ell^{\text{recall}}(g^\dagger)$ or $\ell^{\text{precision}}(g^{\text{output}}) \geq \frac{23}{48} = \frac{23}{21} \ell^{\text{precision}}(g^\dagger)$. Thus $(\frac{7}{16}, \frac{1}{4}) \not\Rightarrow (\frac{7}{16} + 0.01, \frac{1}{4} + 0.01)$.*

Thus, in the scalar-loss case, we allow for a multiplicative factor α . In the Pareto-loss case, we ask a more general question: which pairs of guarantees (p', r') are achievable, given that there exists a hypothesis in the class with precision and recall losses (p, r) ? Since the recall loss is optimizable, for any given r , if there exists a hypothesis in the class with recall loss r , we can always achieve that recall loss. Therefore, we refine our question as follows: given any $p, r \in [0, 1]$, what is the minimum precision loss p' such that $(p, r) \Rightarrow (p', r)$?

Since the recall is estimable, when we have an algorithm achieving α -approximate scalar loss, we can first eliminate all graphs with recall loss higher than r and then run this algorithm over the remaining graphs. Then we can achieve $(p, r) \Rightarrow (\alpha(p + r), r)$.

Theorem 4. *There exist an algorithm such that given an IID training set of size $m \geq O(\frac{\log(|\mathcal{H}|/\delta)}{\varepsilon^2})$, with probability at least $1 - \delta$, the output graph g^{output} satisfies*

$$\ell^{\text{scalar}}(g^{\text{output}}) \leq 5 \cdot \min_{g \in \mathcal{H}} \ell^{\text{scalar}}(g) + \varepsilon.$$

This implies that for any $p, r \in [0, 1]$, $(p, r) \Rightarrow (5(p+r), r)$.

This result is achieved using the same surrogate loss idea in the realizable setting. However, in the agnostic setting, applying maximum likelihood directly no longer works. This is because it is possible that none of the graphs in the hypothesis class are consistent with the training set and thus all graphs have zero likelihood. Instead, we make some modifications to adapt the maximum likelihood idea work for Pareto-loss objective.

As we can see, there is a gap between the upper and lower bounds in the agnostic case, leaving an open question: What is the optimal multiplicative approximation factor α in the scalar case, and what is the optimal p' such that $(p, r) \Rightarrow (p', r')$?

The Semi-Realizable Setting The results in the agnostic setting fail to offer meaningful guarantees in certain natural scenarios, such as $p = 0$ and $r = \frac{1}{2}$ (i.e., when there exists a recommendation hypothesis that captures half of each user's liked items without recommending any items the user does not like). We are therefore interested in the following question: whether $(p = 0, r) \Rightarrow (p' = 0, r' = r)$?

We propose an algorithm with sample complexity depending on the target graph's degree and show that it is impossible to achieve zero precision loss with sample complexity independent of the target graph's degree.

As discussed previously, it's impossible for us to estimate the precision loss. However, we can still separate graphs with zero precision loss and non-zero precision loss if the target graph's degree is bounded. If the precision loss

$$\ell^{\text{precision}}(g, x) = \frac{|N_g(x) \setminus N_{g^{\text{target}}}(x)|}{n_g(x)} = 1 - \frac{|N_g(x) \cap N_{g^{\text{target}}}(x)|}{n_g(x)}$$

of graph g at user x is 0, we have

$$\frac{|N_g(x) \cap N_{g^{\text{target}}}(x)|}{n_g(x) \cdot n_{g^{\text{target}}}(x)} = \frac{1 - \ell^{\text{precision}}(g, x)}{n_{g^{\text{target}}}(x)} = \frac{1}{n_{g^{\text{target}}}(x)}.$$

If the precision loss $\ell^{\text{precision}}(g, x)$ is positive, we have

$$\frac{|N_g(x) \cap N_{g^{\text{target}}}(x)|}{n_g(x) \cdot n_{g^{\text{target}}}(x)} = \frac{1 - \ell^{\text{precision}}(g, x)}{n_{g^{\text{target}}}(x)} < \frac{1}{n_{g^{\text{target}}}(x)}.$$

Hence, we can use the quantity $\mathbb{E}_{x \sim \mathcal{D}} \left[\frac{|N_g(x) \cap N_{g^{\text{target}}}(x)|}{n_g(x) \cdot n_{g^{\text{target}}}(x)} \right]$ to separate graphs with zero and non-zero precision loss, and it is estimable. For each graph g , when the gap of this quantity between graph with zero precision loss and g is $\Delta_{g, \mathcal{D}} := \mathbb{E}_{x \sim \mathcal{D}} \left[\frac{1}{n_{g^{\text{target}}}(x)} \right] - \mathbb{E}_{x \sim \mathcal{D}} \left[\frac{|N_g(x) \cap N_{g^{\text{target}}}(x)|}{n_g(x) \cdot n_{g^{\text{target}}}(x)} \right] > 0$, then by obtaining $\frac{1}{\Delta_{g, \mathcal{D}}}$ samples of x , we can tell that g has nonzero precision loss. Let $\Delta_{\mathcal{D}}$ be the smallest gap of this quantity between graphs with zero and nonzero precision losses:

$$\Delta_{\mathcal{D}} = \min_{g \in \mathcal{H}: \ell^{\text{precision}}(g) > 0} \Delta_{g, \mathcal{D}}.$$

Then we can have sample complexity dependent on this gap.

Theorem 5. *There exists an algorithm such that if there exists a graph $g' \in \mathcal{H}$ with $\ell^{\text{precision}}(g') = 0$ and $\ell^{\text{recall}}(g') = r$, then given an IID training set of size $O(\frac{\log(|\mathcal{H}|/\delta)}{\Delta_{\mathcal{D}}^2})$, with probability $1 - \delta$, it outputs a graph with $\ell^{\text{precision}}(g^{\text{output}}) = 0$ and $\ell^{\text{recall}}(g^{\text{output}}) = r$.*

When the target graph's degree $n_{g^{\text{target}}}(x)$ is bounded by C everywhere, we have

$$\mathbb{E} \left[\frac{|N_g(x) \cap N_{g^{\text{target}}}(x)|}{n_g(x) \cdot n_{g^{\text{target}}}(x)} \right] = \mathbb{E} \left[\frac{1 - \ell^{\text{precision}}(g, x)}{n_{g^{\text{target}}}(x)} \right] \leq \mathbb{E} \left[\frac{1}{n_{g^{\text{target}}}(x)} \right] - \mathbb{E} \left[\frac{\ell^{\text{precision}}(g, x)}{C} \right].$$

Therefore, we have $\Delta_{g, \mathcal{D}} \geq \frac{\ell^{\text{precision}}(g)}{C}$. This implies that when the target graph has bounded degree, we are able to find the graph with zero precision loss. However, when the target graph's degree becomes too large, we show that it is impossible to achieve precision-recall of $(0, r)$.

Theorem 6. *There exists a class $\mathcal{H} = \{g_1, g_2\}$ of two graphs, for any $m > 0$ and any (possibly randomized improper) algorithm \mathcal{A} , there exists a target graph g^{target} and a data distribution \mathcal{D} for which there exists a graph $g^\dagger \in \mathcal{H}$ with $\ell^{\text{precision}}(g^\dagger) = 0$ s.t. with probability $1 - \delta$ over the training set, the expected (over the randomness of the algorithm) precision and recall losses of the output g^{output} satisfy either $\mathbb{E}[\ell^{\text{recall}}(g^{\text{output}})] \geq \min_{g \in \mathcal{H}} \ell^{\text{recall}}(g) + \Omega(1)$ or $\mathbb{E}[\ell^{\text{precision}}(g^{\text{output}})] = \Omega(1)$.*

In the proof of the theorem, we construct a target graph with a degree much larger than the sample size m , as well as graphs g_1 and g_2 with degree 1. In this setup, regardless of whether g_1 has perfect precision (i.e., the neighborhood in g_1 is a subset of the true neighborhood) or very poor precision, we cannot distinguish between these two cases because we never observe a neighbor in g_1 being sampled. Therefore, the only way to achieve zero precision loss is to output the empty graph, which, however, results in a high recall loss.

4 Proof Overview

Given a sequence of IID users x_1, \dots, x_m , let $\widehat{\ell}^{\text{precision}}(g) = \frac{1}{m} \sum_{i=1}^m \frac{|N_g(x_i) \setminus N_{g^{\text{target}}}(x_i)|}{n_g(x_i)}$ and $\widehat{\ell}^{\text{recall}}(g) = \frac{1}{m} \sum_{i=1}^m \frac{|N_{g^{\text{target}}}(x_i) \cap N_g(x_i)|}{n_{g^{\text{target}}}(x_i)}$ denote the empirical precision and recall losses. It suffices to focus on empirical precision and recall losses minimization since by standard concentration bounds, minimizing these empirical losses leads to the minimization of the expected recall and precision losses.

Minimizing Precision Loss Through Maximum Likelihood The maximum likelihood method returns $g^{\text{output}} = \arg \max_{g \in \mathcal{H}} \prod_{i=1}^m \frac{\mathbb{1}(v_i \in N_g(x_i))}{n_g(x_i)}$, which is equivalent to returning the graph with the minimum sum of log degrees among all consistent graphs, i.e.,

$$g^{\text{output}} = \arg \min_{g: g \text{ is consistent}} \sum_{i=1}^m \log(n_g(x_i)).$$

Any consistent graph will have low empirical recall loss by applying standard concentration inequality and thus, $\widehat{\ell}^{\text{recall}}(g^{\text{output}})$ is small and we only to show that the empirical precision loss is small. The main technical challenge in the analysis is *how to connect precision loss with likelihood*.

Since g^{target} is contained in the hypothesis class in the realizable setting and it is consistent with the training data, due to our algorithm, we have

$$\sum_{i=1}^m \log(n_{g^{\text{output}}}(x_i)) \leq \sum_{i=1}^m \log(n_{g^{\text{target}}}(x_i)). \quad (1)$$

We first prove that for any graph g , its empirical precision loss can be bounded by a term of log degree and the empirical recall loss as follows:

$$\widehat{\ell}^{\text{precision}}(g) \leq \frac{2}{m} \sum_{i \in [m]: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \log \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} + 2\widehat{\ell}^{\text{recall}}(g). \quad (2)$$

By combining Eq (1) and (2), we have

$$\widehat{\ell}^{\text{precision}}(g^{\text{output}}) \leq -\frac{2}{m} \sum_{i: \frac{n_{g^{\text{output}}}(x_i)}{n_{g^{\text{target}}}(x_i)} < 1} \log \frac{n_{g^{\text{output}}}(x_i)}{n_{g^{\text{target}}}(x_i)} + 2\widehat{\ell}^{\text{recall}}(g^{\text{output}}).$$

However, with high probability, the first term $-\frac{2}{m} \sum_{i: \frac{n_{g^{\text{output}}}(x_i)}{n_{g^{\text{target}}}(x_i)} < 1} \log \frac{n_{g^{\text{output}}}(x_i)}{n_{g^{\text{target}}}(x_i)}$ is small. This is because, for any graph g , the probability of outputting g is

$$\mathbb{P}_{v_{1:m}}(g^{\text{output}} = g) \leq \mathbb{P}_{v_{1:m}}(g \text{ is consistent}) \leq \prod_{i: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} < 1} \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)}.$$

For any graph g with large $-\frac{2}{m} \sum_{i: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} < 1} \log \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)}$, the probability of outputting such a graph is low.

Modifying Maximum Likelihood in the Agnostic Setting In the agnostic setting, all graphs in the hypothesis class may have zero likelihood of being the true graph; thus, the standard maximum likelihood method doesn't work. However, we can make slight modifications to the maximum likelihood method to make it work for the Pareto-loss objective.

As mentioned earlier, the maximum likelihood method is equivalent to returning the graph with the minimum sum of log degrees among all consistent graphs. Hence, it can be decomposed into two steps: minimizing the recall loss and then regulating by minimizing the sum of log degrees. In the agnostic setting, given any r , we first find the set $\widehat{\mathcal{H}}$ of all graphs in the hypothesis class with recall loss at most $r + 2\varepsilon$, and then regulate by minimizing the sum of log truncated degrees by returning

$$g^{\text{output}} = \arg \min_{g' \in \widehat{\mathcal{H}}} \max_{g \in \widehat{\mathcal{H}}} \frac{1}{m} \sum_{i \in [m]} \log \frac{n_{g'}(x_i) \wedge 4n_g(x_i)}{n_g(x_i) \wedge 4n_{g'}(x_i)},$$

where $a \wedge b := \min(a, b)$. The truncation plays an important role here. Intuitively, let g^\dagger denote the graph with precision and recall losses p and r , respectively. If there exists an x_i such that $n_{g^\dagger}(x_i)$ is very large, minimizing the sum of log-untruncated degrees will never return g^\dagger . By applying truncation, we limit the effect of a single user with a very large degree.

Minimizing the Surrogate for the Scalar Loss Here we consider an alternative learning rule based on two simple principles for discarding sub-optimal hypotheses. We illustrate these principles with the following intuitive example: consider a music recommender system, and assume we are considering two candidate hypotheses, g' and g'' . Both hypotheses recommend classical music; however, g' recommends pieces by Bach 20% of the time and pieces by Mozart 10% of the time, while g'' never recommends any pieces by Mozart or Bach.

Now, suppose that in the training set, users frequently choose to listen to pieces by Mozart. This observation suggests that g'' should be discarded, as it never recommends Mozart. This leads to our first rule: if a hypothesis exhibits a high recall loss, it can be discarded. The second rule addresses precision loss, which is more challenging because it cannot be directly estimated from the data. To illustrate the second rule, imagine that in the training set, users tend to pick Bach pieces only 5% of the time. This suggests that g' is over-recommending Bach pieces, and therefore, g' might also be discarded based on its likely precision loss.

We formally capture this using a metric defined in the following. For any graph g , let U_i^g denote the uniform distribution over the neighbors $N_g(x_i)$ of x_i . Then, for any graph g , we define a vector $v_g : \mathcal{H} \times \mathcal{H} \rightarrow [0, 1]$ by

$$v_g(g', g'') = \frac{1}{m} \sum_{i=1}^m U_i^g(N_{g'}(x_i) \setminus N_{g''}(x_i)) = \frac{1}{m} \sum_{i=1}^m \mathbb{P}_{v \sim U_i^g} (v \in N_{g'}(x_i) \setminus N_{g''}(x_i)).$$

Intuitively, $v_g(g', g'')$ represents the fraction of items users like that are recommended by g' but not by g'' in the counterfactual scenario where g is the target graph. If g is indeed the target graph, then this quantity should be consistent with our training data, i.e.,

$$v_g(g', g'') \approx v_{\hat{g}}(g', g''),$$

where \hat{g} is the observed (empirical) graph; i.e., the graph in which every x_i is connected to the random number v_i which is observed in the training set. In the above example, $v_{g'}(g', g'')$ is 20% while $v_{\hat{g}}(g', g'')$ is 5% and thus g' is unlikely to be the target graph.

We define a metric $d_{\mathcal{H}}$ between two graphs g_1 and g_2 by

$$d_{\mathcal{H}}(g_1, g_2) = \|v_{g_1} - v_{g_2}\|_{\infty}.$$

Surprisingly, we show that $d_{\mathcal{H}}(g^{\text{target}}, g)$ is a surrogate for the scalar loss, providing both lower and upper bounds on the scalar loss with a constant multiplicative factor:

$$\frac{1}{2} d_{\mathcal{H}}(g^{\text{target}}, g) \leq \ell^{\text{scalar}}(g) \leq d_{\mathcal{H}}(g^{\text{target}}, g).$$

A standard union bound argument yields that with probability at least $1 - \delta$,

$$d_{\mathcal{H}}(\hat{g}, g^{\text{target}}) \leq O\left(\sqrt{\frac{\log |\mathcal{H}| + \log(1/\delta)}{m}}\right).$$

By triangle inequality, we have

$$d_{\mathcal{H}}(g^{\text{target}}, g) \leq d_{\mathcal{H}}(\hat{g}, g) + O\left(\sqrt{\frac{\log |\mathcal{H}| + \log(1/\delta)}{m}}\right).$$

Then, we return a graph $g^{\text{output}} \in \mathcal{H}$ such that

$$d_{\mathcal{H}}(\hat{g}, g^{\text{output}}) = \min_{g \in \mathcal{H}} d_{\mathcal{H}}(\hat{g}, g).$$

The Hardness of No Knowledge of the Target Graph’s Degree For any graph g , it’s precision loss at any user x is $\frac{|N_g(x) \setminus N_{g^{\text{target}}}(x)|}{n_g(x)}$ and we only get a random neighbor $v \sim N_{g^{\text{target}}}(x)$. If we are given the knowledge of the degree $n_{g^{\text{target}}}(x)$ of the target graph, then we can obtain an unbiased estimate of the precision loss, i.e., $1 - \frac{n_{g^{\text{target}}}(x)}{n_g(x)} \cdot \mathbf{1}(v \in N_g(x))$. But the difficulty lies in that we don’t know $n_{g^{\text{target}}}(x)$.

Consider the following example illustrated in Fig 3. For a given user x , there is a set of n nodes equally divided into two sets $N_1(x)$ and $N_2(x)$. Consider two graphs— g_1 with neighborhood $N_{g_1}(x) = N_1(x)$ and g_2 with neighborhood $N_{g_2}(x) = N_1(x) \cup N_2(x)$ being all n nodes.

In a world characterized by $\beta \in [\frac{1}{8}, \frac{2}{3}]$, the target graph is generated in the following random way: Randomly select $\frac{3}{4} \cdot \beta n$ nodes from $N_1(x)$ and $\frac{1}{4} \cdot \beta n$ nodes from $N_2(x)$. No matter what β is, w.p. $\frac{3}{4}$, v is sampled uniformly at random from $N_1(x)$ and w.p. $\frac{1}{4}$, v is sampled uniformly at random from $N_2(x)$. That is, every node in $N_1(x)$ has probability $\frac{3}{2n}$ of being sampled and every node in $N_2(x)$ has probability $\frac{1}{2n}$ of being sampled. Hence, if we have never seen the same user twice, we cannot distinguish between different β ’s.

For any g^{target} generated from the above process, the scalar loss of g_1 at x is

$$\begin{aligned} \ell^{\text{scalar}}(g_1, x) &= 1 - \left(\frac{|N_{g^{\text{target}}}(x) \cap N_{g_1}(x)|}{2|N_{g_1}(x)|} + \frac{|N_{g^{\text{target}}}(x) \cap N_{g_1}(x)|}{2|N_{g^{\text{target}}}(x)|} \right) = 1 - \left(\frac{3/4 \cdot \beta n}{n} + \frac{3/4 \cdot \beta n}{2\beta n} \right) \\ &= \frac{5}{8} - \frac{3}{4}\beta, \end{aligned}$$

and the scalar loss of g_2 at x is

$$\ell^{\text{scalar}}(g_2, x) = 1 - \left(\frac{|N_{g^{\text{target}}}(x) \cap N_{g_2}(x)|}{2|N_{g_2}(x)|} + \frac{|N_{g^{\text{target}}}(x) \cap N_{g_2}(x)|}{2|N_{g^{\text{target}}}(x)|} \right) = 1 - \left(\frac{\beta n}{2n} + \frac{\beta n}{2\beta n} \right) = \frac{1}{2} - \frac{1}{2}\beta.$$

When β is large, g_1 has a smaller loss; when β is small, g_2 has a smaller loss. With a huge number of users, we might never observe the same user twice and, therefore, cannot distinguish between different β values. Consequently, it’s impossible to determine which of the two graphs has a smaller loss. We show that, in this example, even if algorithms are allowed to be randomized and improper, it still impossible to compete with the best graph in the hypothesis class.

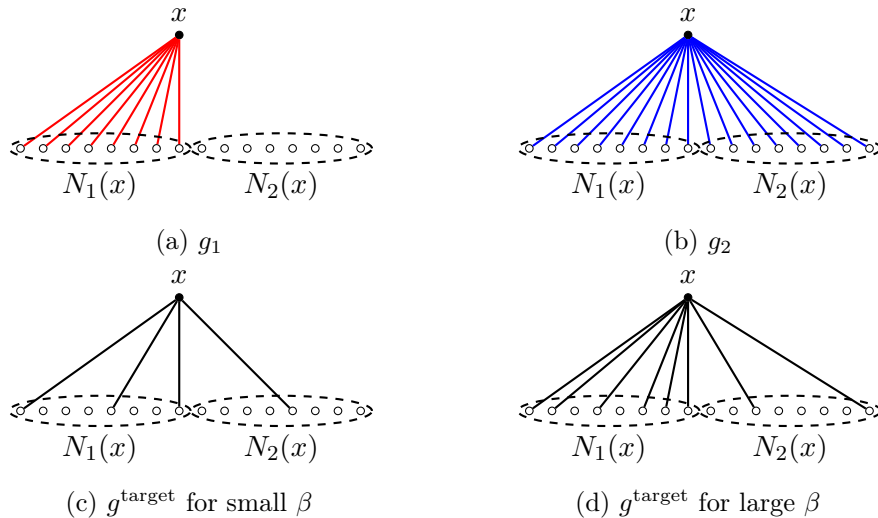


Figure 3: Illustration of g_1 , g_2 and randomly generated g^{target} .

5 Algorithms and Proofs

Notations Let $\ell^{\text{precision}}(g, x)$ and $\ell^{\text{recall}}(g, x)$ denote the precision loss and recall loss of graph g at node x :

$$\ell^{\text{precision}}(g, x) = \frac{|N_g(x) \setminus N_{g^{\text{target}}}(x)|}{n_g(x)},$$

$$\ell^{\text{recall}}(g, x) = \frac{|N_{g^{\text{target}}}(x) \setminus N_g(x)|}{n_{g^{\text{target}}}(x)}.$$

Let $\widehat{\ell}^{\text{precision}}(g) = \frac{1}{m} \sum \ell^{\text{precision}}(g, x_i)$ and $\widehat{\ell}^{\text{recall}}(g) = \frac{1}{m} \sum \ell^{\text{recall}}(g, x_i)$ denote the empirical precision and recall losses. Let $a \wedge b := \min(a, b)$.

5.1 Maximum Likelihood Method in the Realizable Case

In the realizable setting, the target graph g^{target} is in the hypothesis class. Given the IID training data $(x_1, v_1), \dots, (x_m, v_m)$, the maximum likelihood method returns the graph

$$g^{\text{output}} = \arg \max_{g \in \mathcal{H}} \prod_{i=1}^m \frac{\mathbb{1}(v_i \in N_g(x_i))}{n_g(x_i)}.$$

In other words, g^{output} is a graph in \mathcal{H} satisfying

- consistency: $\sum_{i=1}^m \mathbb{1}(v_i \notin N_{g^{\text{output}}}(x_i)) = 0$
- regulation: among all consistent graphs, $g^{\text{output}} = \arg \min_{g: g \text{ is consistent}} \sum_{i=1}^m \log(n_g(x_i))$.⁵

Theorem 1. *In the realizable setting, there exist algorithms such that given an IID training set of size $m \geq O(\frac{\log(|\mathcal{H}|/\delta)}{\varepsilon})$, with probability at least $1 - \delta$, the output graph g^{output} satisfies*

$$\ell^{\text{recall}}(g^{\text{output}}) \leq \varepsilon, \quad \ell^{\text{precision}}(g^{\text{output}}) \leq \varepsilon.$$

Proof of Theorem 1 It suffices to prove that for any fixed (x_1, \dots, x_m) , when $m \geq \frac{12 \log(4|\mathcal{H}|/\delta)}{\varepsilon}$, w.p. at least $1 - \delta/2$ over $v_{1:m}$, the empirical values of recall and precision are small, $\widehat{\ell}^{\text{recall}}(g^{\text{output}}) \leq \varepsilon/2$ and $\widehat{\ell}^{\text{precision}}(g^{\text{output}}) \leq \varepsilon/2$. Then we can show $\ell^{\text{recall}}(g^{\text{output}}) \leq \varepsilon$ and $\ell^{\text{precision}}(g^{\text{output}}) \leq \varepsilon$ by applying empirical Bernstein bounds to both precision and recall losses. Now we prove this statement. Bounding recall loss is easy as $\sum_{i=1}^m \mathbb{1}(v_i \notin N_g(x_i))$ is an unbiased estimate of recall loss. With probability $1 - \delta/4$ over the randomness of $v_i \sim \text{Unif}(N_{g^{\text{target}}}(x_i))$ for all $i \in [m]$, the empirical loss for recall is

$$\widehat{\ell}^{\text{recall}}(g) = \frac{1}{m} \sum_{i=1}^m \frac{|N_{g^{\text{target}}}(x_i) \setminus N_g(x_i)|}{n_{g^{\text{target}}}(x_i)} \leq \sum_{i=1}^m \mathbb{1}(v_i \notin N_g(x_i)) + \varepsilon/6 = \varepsilon/6, \quad (3)$$

for all consistent g with $\sum_{i=1}^m \mathbb{1}(v_i \notin N_g(x_i)) = 0$. Hence, w.p. $1 - \delta/4$, $\widehat{\ell}^{\text{recall}}(g^{\text{output}}) \leq \varepsilon/6$.

⁵the base of log in this work is 2.

Bounding precision loss is more challenging. Let $A_g = \{i \in [m] | n_{g^{\text{target}}}(x_i) \leq 2n_g(x_i)\}$. Then we can decompose the empirical precision loss as

$$\begin{aligned}
& \widehat{\ell}^{\text{precision}}(g) \\
&= \frac{1}{m} \sum_{i=1}^m \frac{|N_g(x_i) \setminus N_{g^{\text{target}}}(x_i)|}{n_g(x_i)} \\
&\leq \frac{1}{m} \sum_{i \in A_g} \frac{|N_g(x_i) \setminus N_{g^{\text{target}}}(x_i)|}{n_g(x_i)} + \frac{1}{m} \sum_{i=1}^m \mathbb{1}(i \notin A_g) \\
&\leq \frac{1}{m} \sum_{i \in A_g} \min\left(\frac{2|N_g(x_i) \setminus N_{g^{\text{target}}}(x_i)|}{n_{g^{\text{target}}}(x_i)}, 1\right) + \frac{2}{m} \sum_{i \notin A_g} \frac{|N_{g^{\text{target}}}(x_i) \setminus N_g(x_i)|}{n_{g^{\text{target}}}(x_i)} \\
&= \frac{2}{m} \sum_{i \in A_g} \min\left(\frac{n_g(x_i) + |N_{g^{\text{target}}}(x_i) \setminus N_g(x_i)| - n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)}, \frac{1}{2}\right) + \frac{2}{m} \sum_{i \notin A_g} \frac{|N_{g^{\text{target}}}(x_i) \setminus N_g(x_i)|}{n_{g^{\text{target}}}(x_i)} \\
&\leq \frac{2}{m} \sum_{i \in A_g} \min\left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} - 1, \frac{1}{2}\right) + \frac{2}{m} \sum_{i \in A_g} \frac{|N_{g^{\text{target}}}(x_i) \setminus N_g(x_i)|}{n_{g^{\text{target}}}(x_i)} + \frac{2}{m} \sum_{i \notin A_g} \frac{|N_{g^{\text{target}}}(x_i) \setminus N_g(x_i)|}{n_{g^{\text{target}}}(x_i)} \\
&\leq \frac{2}{m} \sum_{i \in A_g} \min\left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} - 1, \frac{1}{2}\right) + 2\widehat{\ell}^{\text{recall}}(g).
\end{aligned}$$

The second term is the empirical loss for recall, which is upper bounded by Eq (3). For the first term,

$$\begin{aligned}
& \frac{1}{m} \sum_{i \in A_g} \min\left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} - 1, \frac{1}{2}\right) \\
&= \frac{1}{m} \sum_{i: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq \frac{1}{2}} \min\left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} - 1, \frac{1}{2}\right) \\
&\leq \frac{1}{m} \sum_{i: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \min\left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} - 1, \frac{1}{2}\right) \\
&\leq \sum_{i: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \log \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)},
\end{aligned}$$

where the last inequality adopts the following fact: for all $z \geq 1$, $\min(z - 1, \frac{1}{2}) \leq \log z$. On the other hand, we have

$$\frac{1}{m} \sum_{i: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \log \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \leq \frac{1}{m} \sum_{i: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} < 1} \log \frac{n_{g^{\text{target}}}(x_i)}{n_g(x_i)}$$

when g satisfies $\sum_{i=1}^m \log(n_g(x_i)) \leq \sum_{i=1}^m \log(n_{g^{\text{target}}}(x_i))$. Hence, for any graph g with $\widehat{\ell}^{\text{precision}}(g) > \frac{\varepsilon}{2}$ and $\widehat{\ell}^{\text{recall}}(g) \leq \frac{\varepsilon}{6}$, we have

$$\sum_{i: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} < 1} \log \frac{n_{g^{\text{target}}}(x_i)}{n_g(x_i)} \geq \frac{m}{2} (\widehat{\ell}^{\text{precision}}(g) - 2\widehat{\ell}^{\text{recall}}(g)) > m \cdot \varepsilon / 12 \geq \log(4|\mathcal{H}|/\delta).$$

The probability of outputting such a graph g is

$$\mathbb{P}_{v_{1:m}}(g^{\text{output}} = g) \leq \mathbb{P}_{v_{1:m}}(g \text{ is consistent}) \leq \prod_{i: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} < 1} \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} < \frac{\delta}{4|\mathcal{H}|}.$$

Hence, with probability at least $1 - \delta/2$ over $v_{1:m}$, g^{output} will satisfy

$$\widehat{\ell}^{\text{recall}}(g^{\text{output}}) \leq \varepsilon/6, \widehat{\ell}^{\text{precision}}(g^{\text{output}}) \leq \varepsilon/2.$$

Then we are done with the proof. \square

5.2 Modified Maximum Likelihood Method in the Agnostic Case

In the agnostic setting, we shift our goal from finding a graph with nearly zero precision and recall losses to determining, given any $p, r \in [0, 1]$, the minimum precision loss p' such that $(p, r) \Rightarrow (p', r)$. In fact, we can show something stronger: we do not need to know p . Specifically, given any $r \in [0, 1]$, let $p = \min_{g: \ell^{\text{recall}}(g) \leq r} \ell^{\text{precision}}(g)$ be the optimal precision loss among all graphs with recall loss at most r . What is the smallest p' such that we achieve $\ell^{\text{recall}}(g^{\text{output}}) \leq r$ and $\ell^{\text{precision}}(g^{\text{output}}) \leq p'$?

Theorem 3. *There exists a class $\mathcal{H} = \{g_1, g_2\}$ of two graphs, such that for any (possibly randomized improper) algorithm given the knowledge of $(p, r) = (\frac{7}{16}, \frac{1}{4})$, there exists a target graph g^{target} with bounded degree and a data distribution \mathcal{D} for which there exists a graph $g^\dagger \in \mathcal{H}$ with $\ell^{\text{recall}}(g^\dagger) = \frac{1}{4}$ and $\ell^{\text{precision}}(g^\dagger) = \frac{7}{16}$ s.t. for any sample size $m > 0$, with probability 1 over the training set, the expected (over the randomness of the algorithm) precision and recall losses of the output g^{output} satisfy*

$$\mathbb{E}[\ell^{\text{recall}}(g^{\text{output}})] + \frac{12}{5}\mathbb{E}[\ell^{\text{precision}}(g^{\text{output}})] \geq \frac{7}{5}.$$

Recall that in the realizable case, the maximum likelihood method selects the consistent graph with the smallest empirical log degree, i.e., $g^{\text{output}} = \arg \min_{g: g \text{ is consistent}} \sum_{i=1}^m \log(n_g(x_i))$. Basically, the consistency guarantees small recall loss and the empirical log degree is used as a regulation term to bound precision loss. In the agnostic setting, the maximum likelihood method fails as there may be no consistent graph. We present a modified version of the maximum likelihood method, which still has the ‘‘consistency’’ component and adopts log degree as a regulation term. The algorithm operates as follows.

Algorithm:

1. **Finding a set of plausible graphs which make at most $m \cdot (r + \varepsilon)$ mistakes:** For any graph g , let $I_g = \{i | v_i \notin N_g(x_i)\}$ denote the indices of training points that the graph g is inconsistent with. Then let $\widehat{\mathcal{H}} = \{g \in \mathcal{H} | |I_g| \leq m \cdot (r + \varepsilon)\}$ be the set of graphs making at most $m \cdot (r + \varepsilon)$ mistakes.
2. **Returning the graph with low empirical log truncated degree:** Return

$$g^{\text{output}} = \arg \min_{g' \in \widehat{\mathcal{H}}} \max_{g \in \widehat{\mathcal{H}}} \frac{1}{m} \sum_{i \in [m]} \log \frac{n_{g'}(x_i) \wedge 4n_g(x_i)}{n_g(x_i) \wedge 4n_{g'}(x_i)}. \quad (4)$$

Theorem 7. Given any $r \in [0, 1]$ and $m \geq O(\frac{\log(|\mathcal{H}|) + \log(1/\delta)}{\varepsilon^2})$ IID training data, the modified maximum likelihood method can return a graph satisfying

$$\ell^{\text{recall}}(g^{\text{output}}) \leq r + \varepsilon, \quad \ell^{\text{precision}}(g^{\text{output}}) \leq 28r + 15p + \varepsilon,$$

where $p = \min_{g: \ell^{\text{recall}}(g) \leq r} \ell^{\text{precision}}(g)$.

Proof Sketch Let g^\dagger to be the graph with recall loss at most r and precision loss p . Let $\bar{r} = r + 2\varepsilon$ and $\bar{p} = p + 2\varepsilon$. Similar to the realizable setting, we build connection between precision and recall with empirical log degree (Lemma 2), i.e., for any constant $c \in (0, 1]$,

$$\widehat{\ell}^{\text{precision}}(g) \leq \frac{1+c}{m} \sum_{i \in [m]: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} + \frac{1+c}{c} \widehat{\ell}^{\text{recall}}(g).$$

It suffices to upper bound $\sum_{i \in [m]: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)}$. We first decompose it into

$$\begin{aligned} & \sum_{i \in [m]: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} \\ &= \underbrace{\sum_{i \in B} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)}}_{(a)} - \underbrace{\sum_{i \in B: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} < 1} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)}}_{(b)}, \end{aligned}$$

where $B = \{i \mid \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq \frac{1}{2}\}$. The term (b) is lower bounded by the recall loss in Lemma 4. Intuitively, if $\frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)}$ is small, the recall loss must be large while any graph in $\widehat{\mathcal{H}}$ has a small empirical recall loss.

For term (a), since $\widehat{\ell}^{\text{recall}}(g^\dagger) \leq \bar{r}$ and $\widehat{\ell}^{\text{precision}}(g^\dagger) \leq \bar{p}$, at most training points, $\frac{n_{g^\dagger}(x_i)}{n_{g^{\text{target}}}(x_i)}$ is in $[\frac{1}{2}, 2]$ (if it's too large at x_i , precision loss is large; if it's too small, recall loss is large). Also, for any graph in $\widehat{\mathcal{H}}$, the empirical recall loss is small and thus, at most training points, $\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq \frac{1}{2}$. At these training points satisfying $\frac{n_{g^\dagger}(x_i)}{n_{g^{\text{target}}}(x_i)}$ is in $[\frac{1}{2}, 2]$ and $\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq \frac{1}{2}$, we have

$$\begin{aligned} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} &\leq \log \frac{n_g(x_i) \wedge 4n_{g^\dagger}(x_i)}{n_{g^{\text{target}}}(x_i)} = \log \frac{n_g(x_i) \wedge 4n_{g^\dagger}(x_i)}{n_{g^\dagger}(x_i) \wedge 4n_g(x_i)} + \log \frac{n_{g^\dagger}(x_i) \wedge 4n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \\ &\leq \log \frac{n_g(x_i) \wedge 4n_{g^\dagger}(x_i)}{n_{g^\dagger}(x_i) \wedge 4n_g(x_i)} + \log \frac{n_{g^\dagger}(x_i)}{n_{g^{\text{target}}}(x_i)}. \end{aligned}$$

The first term is bounded due to our algorithm while the second term is bounded by the empirical precision (in Lemma 5). Intuitively, if $\frac{n_{g^\dagger}(x_i)}{n_{g^{\text{target}}}(x_i)}$ is large, the empirical precision of g^\dagger is large. \square

5.2.1 Proof of Theorem 3

Proof of Theorem 3 Suppose there are infinite users and \mathcal{D} be the uniform distribution. For each user x , there is an individual set of 12 nodes, denoted as $N(x) = \{v_{x,1}, v_{x,2}, \dots, v_{x,12}\}$. The graph g_1 's neighborhood is the first 8 nodes $N_{g_1}(x) = \{v_{x,1}, v_{x,2}, \dots, v_{x,8}\}$ and the graph g_2 's neighborhood is the last 8 nodes $N_{g_2}(x) = \{v_{x,5}, v_{x,6}, \dots, v_{x,12}\}$. There are two worlds in which the target graph g^{target} is generated differently:

- **World I:** W.p. $\frac{1}{2}$, $N_{g^{\text{target}}}(x) = N_{g_1}(x)$; w.p. $\frac{1}{2}$, $N_{g^{\text{target}}}(x) = \{u_1, u_2\}$ where u_1 is sampled uniformly from $\{v_{x,5}, v_{x,6}, \dots, v_{x,8}\}$ and u_2 is sampled uniformly from $\{v_{x,9}, v_{x,10}, \dots, v_{x,12}\}$.
- **World II:** W.p. $\frac{1}{2}$, $N_{g^{\text{target}}}(x) = N_{g_2}(x)$; w.p. $\frac{1}{2}$, $N_{g^{\text{target}}}(x) = \{u_1, u_2\}$ where u_1 is sampled uniformly from $\{v_{x,5}, v_{x,6}, \dots, v_{x,8}\}$ and u_2 is sampled uniformly from $\{v_{x,1}, v_{x,2}, \dots, v_{x,4}\}$.

These two worlds are symmetric. In world I,

$$\begin{aligned} \text{recall}(g_1) &= \frac{1}{2} \cdot 1 + \frac{1}{2} \cdot \frac{1}{2} = \frac{3}{4}, & \text{precision}(g_1) &= \frac{1}{2} \cdot 1 + \frac{1}{2} \cdot \frac{1}{8} = \frac{9}{16}, \\ \text{recall}(g_2) &= \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot 1 = \frac{3}{4}, & \text{precision}(g_2) &= \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{4} = \frac{3}{8}. \end{aligned}$$

Both g_1 and g_2 have the same recall and g_1 has a better precision in the world I. In world II, g_1 and g_2 switch their losses. Hence, in either world, we have

$$\min_{g \in \mathcal{H}} \ell^{\text{recall}}(g) = \frac{1}{4}, \quad \min_{g \in \mathcal{H}} \ell^{\text{precision}}(g) = \frac{7}{16}.$$

In both worlds, the distribution of v is identical, i.e., w.p. $\frac{1}{2}$, v is sampled from $\text{Unif}(N_{g_1}(x))$ and w.p. $\frac{1}{2}$, v is sampled from $\text{Unif}(N_{g_2}(x))$. Hence, if no x has been sampled twice, no algorithm can distinguish the two worlds. Since there are infinite users, we will not sample the same user twice almost surely. For any output graph g^{output} , we let

$$\begin{aligned} n_1 &= |N_{g^{\text{output}}}(x) \cap \{v_{x,1}, v_{x,2}, \dots, v_{x,4}\}| \in \{0, 1, \dots, 4\}, \\ n_2 &= |N_{g^{\text{output}}}(x) \cap \{v_{x,5}, v_{x,6}, \dots, v_{x,8}\}| \in \{0, 1, \dots, 4\}, \\ n_3 &= |N_{g^{\text{output}}}(x) \cap \{v_{x,9}, v_{x,10}, \dots, v_{x,12}\}| \in \{0, 1, \dots, 4\}. \end{aligned}$$

Then in world I, the expected recall and precision of g^{output} is

$$\begin{aligned} \mathbb{E}_{g^{\text{target}}} [\text{recall}(g^{\text{output}}, x)] &= \frac{1}{2} \left(\frac{n_1 + n_2}{8} + \frac{n_2 + n_3}{8} \right) = \frac{n_1 + 2n_2 + n_3}{16}, \\ \mathbb{E}_{g^{\text{target}}} [\text{precision}(g^{\text{output}}, x)] &= \frac{1}{2} \left(\frac{n_1 + n_2}{n_1 + n_2 + n_3} + \frac{n_2 + n_3}{4(n_1 + n_2 + n_3)} \right) = \frac{4n_1 + 5n_2 + n_3}{8(n_1 + n_2 + n_3)}. \end{aligned}$$

Then suppose we randomly choose one of the two world. By taking expectation over the world, x , and the randomness of the algorithm, we have

$$\begin{aligned} \mathbb{E} [(\text{recall}(g^{\text{output}}), \text{precision}(g^{\text{output}}))] &= \mathbb{E} \left[\left(\frac{n_1 + 2n_2 + n_3}{16}, \frac{5(n_1 + 2n_2 + n_3)}{16(n_1 + n_2 + n_3)} \right) \right] \\ &= \mathbb{E} \left[\left(\frac{n_1 + 2n_2 + n_3}{16}, \frac{5}{16} + \frac{5}{16} \cdot \frac{n_2}{n_1 + n_2 + n_3} \right) \right] \end{aligned}$$

Let's denote by $r(n_1, n_2, n_3) = \frac{n_1 + 2n_2 + n_3}{16}$ and $p(n_1, n_2, n_3) = \frac{5}{16} + \frac{5}{16} \cdot \frac{n_2}{n_1 + n_2 + n_3}$. Then for any $(n_1, n_2, n_3) \in \{0, \dots, 4\}^3$, we have

$$\begin{aligned}
& r(n_1, n_2, n_3) + \frac{12}{5}p(n_1, n_2, n_3) \\
&= \frac{n_1 + 2n_2 + n_3}{16} + \frac{3}{4} + \frac{3}{4} \cdot \frac{n_2}{n_1 + n_2 + n_3} \\
&\leq \frac{n_1 + 8 + n_3}{16} + \frac{3}{4} + \frac{3}{n_1 + 4 + n_3} && \text{(maximized at } n_2 = 4) \\
&\leq 2. && \text{(maximized at } n_1 + n_3 = 0 \text{ or } 8)
\end{aligned}$$

Hence, we have

$$\mathbb{E} [\text{recall}(g^{\text{output}})] + \frac{12}{5} \mathbb{E} [\text{precision}(g^{\text{output}})] \leq 2.$$

This is equivalent to

$$\mathbb{E} [\ell^{\text{recall}}(g^{\text{output}})] + \frac{12}{5} \mathbb{E} [\ell^{\text{precision}}(g^{\text{output}})] \geq \frac{7}{5}.$$

□

5.2.2 Proof of Theorem 7

Proof of Theorem 7 Let $\Delta = \sqrt{\frac{\log(|\mathcal{H}|/\delta)}{m}}$, $\bar{r} = r + 2\Delta$ and $\bar{p} = p + 2\Delta$. Let g^\dagger denote the graph with precision and recall losses (p, r) . We know that with probability at least $1 - \delta$, all graphs in $\widehat{\mathcal{H}}$ have empirical recall loss no greater than \bar{r} and g^\dagger also has $\widehat{\ell}^{\text{precision}}(g^\dagger) \leq \bar{p}$. Then the proof is divided into two parts:

(i) graph g^\dagger satisfies that

$$\max_{g \in \widehat{\mathcal{H}}} \frac{1}{m} \sum_{i \in [m]} \log \frac{n_{g^\dagger}(x_i) \wedge 4n_g(x_i)}{n_g(x_i) \wedge 4n_{g^\dagger}(x_i)} \leq 6\bar{r} + 4\bar{p} + \frac{2}{m}.$$

Therefore, we have $\max_{g \in \widehat{\mathcal{H}}} \frac{1}{m} \sum_{i \in [m]} \log \frac{n_{g^{\text{output}}}(x_i) \wedge 4n_g(x_i)}{n_g(x_i) \wedge 4n_{g^{\text{output}}}(x_i)} \leq 6\bar{r} + 4\bar{p} + \frac{2}{m}$.

(ii) any graph g' satisfying $\max_{g \in \widehat{\mathcal{H}}} \sum_{i \in [m]} \log \frac{n_{g'}(x_i) \wedge 4n_g(x_i)}{n_g(x_i) \wedge 4n_{g'}(x_i)} \leq 6m\bar{r} + 4m\bar{p} + 2$ has $\widehat{\ell}^{\text{precision}}(g') \leq 28r + 15p + o(1)$. Hence, g^{output} can achieve good precision.

We prove part (i) by Lemma 1 and part (ii) by combining Lemma 2 and 3. □

Lemma 1. For any x_1, \dots, x_m and graph g with $\widehat{\ell}^{\text{recall}}(g) \leq \bar{r}$, we have

$$\frac{1}{m} \sum_{i \in [m]} \log \frac{n_{g^\dagger}(x_i) \wedge 4n_g(x_i)}{n_g(x_i) \wedge 4n_{g^\dagger}(x_i)} \leq 6\bar{r} + 4\bar{p} + \frac{2}{m}.$$

Proof Let $E = \{i \mid \frac{n_g(x_i)}{n_{g^{\dagger}}(x_i)} \geq \frac{1}{2}, \frac{1}{2} \leq \frac{n_{g^{\dagger}}(x_i)}{n_{g^{\dagger}}(x_i)} \leq 2\}$. According to Lemma 6, we know $|\neg E| \leq 4m\bar{r} + 2m\bar{p}$.

$$\begin{aligned}
\sum_{i \in [m]} \log \frac{n_{g^{\dagger}}(x_i) \wedge 4n_g(x_i)}{n_g(x_i) \wedge 4n_{g^{\dagger}}(x_i)} &\leq \sum_{i \in E} \log \frac{n_{g^{\dagger}}(x_i) \wedge 4n_g(x_i)}{n_g(x_i) \wedge 4n_{g^{\dagger}}(x_i)} + |\neg E| \\
&\leq \sum_{i \in E} \log \frac{n_{g^{\dagger}}(x_i)}{n_g(x_i) \wedge n_{g^{\dagger}}(x_i)} + 4m\bar{r} + 2m\bar{p} \\
&= \sum_{i \in E} \log \frac{n_{g^{\dagger}}(x_i)}{n_{g^{\dagger}}(x_i)} - \sum_{i \in E} \log \frac{n_g(x_i) \wedge n_{g^{\dagger}}(x_i)}{n_{g^{\dagger}}(x_i)} + 4m\bar{r} + 2m\bar{p} \\
&\leq 6m\bar{r} + 4m\bar{p} + 2. \tag{Applying Lemmas 4 and 5}
\end{aligned}$$

□

Lemma 2. For any x_1, \dots, x_m , any graph g , and any positive constant $c \in (0, 1]$, the empirical loss for precision $\hat{\ell}^{\text{precision}}(g)$ satisfies

$$\hat{\ell}^{\text{precision}}(g) \leq \frac{1+c}{m} \sum_{i \in [m]: \frac{n_g(x_i)}{n_{g^{\dagger}}(x_i)} \geq 1} \log \frac{n_g(x_i) \wedge 2n_{g^{\dagger}}(x_i)}{n_{g^{\dagger}}(x_i)} + \frac{1+c}{c} \hat{\ell}^{\text{recall}}(g).$$

Proof of Lemma 2 Let $A_g = \{i \in [m] \mid n_{g^{\dagger}}(x_i) \leq (1+c)n_g(x_i)\}$. Then we have

$$\begin{aligned}
&\hat{\ell}^{\text{precision}}(g) \\
&= \frac{1}{m} \sum_{i=1}^m \frac{|N_g(x_i) \setminus N_{g^{\dagger}}(x_i)|}{n_g(x_i)} \\
&\leq \frac{1}{m} \sum_{i \in A_g} \frac{|N_g(x_i) \setminus N_{g^{\dagger}}(x_i)|}{n_g(x_i)} + \frac{1}{m} \sum_{i=1}^m \mathbb{1}(i \notin A_g) \\
&\leq \frac{1}{m} \sum_{i \in A_g} \min\left(\frac{(1+c)|N_g(x_i) \setminus N_{g^{\dagger}}(x_i)|}{n_{g^{\dagger}}(x_i)}, 1\right) + \frac{1+c}{c \cdot m} \sum_{i \notin A_g} \frac{|N_{g^{\dagger}}(x_i) \setminus N_g(x_i)|}{n_{g^{\dagger}}(x_i)} \\
&\leq \frac{1+c}{m} \sum_{i \in A_g} \min\left(\frac{n_g(x_i) + |N_{g^{\dagger}}(x_i) \setminus N_g(x_i)| - n_{g^{\dagger}}(x_i)}{n_{g^{\dagger}}(x_i)}, 1\right) + \frac{1+c}{c \cdot m} \sum_{i \notin A_g} \frac{|N_{g^{\dagger}}(x_i) \setminus N_g(x_i)|}{n_{g^{\dagger}}(x_i)} \\
&\leq \frac{1+c}{m} \sum_{i \in A_g} \min\left(\frac{n_g(x_i)}{n_{g^{\dagger}}(x_i)} - 1, 1\right) + \frac{1+c}{m} \sum_{i \in A_g} \frac{|N_{g^{\dagger}}(x_i) \setminus N_g(x_i)|}{n_{g^{\dagger}}(x_i)} + \frac{1+c}{c \cdot m} \sum_{i \notin A_g} \frac{|N_{g^{\dagger}}(x_i) \setminus N_g(x_i)|}{n_{g^{\dagger}}(x_i)} \\
&\leq \frac{1+c}{m} \sum_{i \in A_g} \min\left(\frac{n_g(x_i)}{n_{g^{\dagger}}(x_i)} - 1, 1\right) + \frac{1+c}{c} \hat{\ell}^{\text{recall}}(g).
\end{aligned}$$

Now we upper bound the first term.

$$\begin{aligned}
& \sum_{i \in A_g} \min \left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} - 1, 1 \right) \\
&= \sum_{i \in [m]: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq \frac{1}{1+c}} \min \left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} - 1, 1 \right) \\
&\leq \sum_{i \in [m]: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \min \left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} - 1, 1 \right) \\
&\leq \sum_{i \in [m]: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \min \left(\log \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)}, 1 \right) \\
&= \sum_{i \in [m]: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)},
\end{aligned}$$

where the last inequality adopts the fact: for all $z \geq 1$, $\min(z - 1, 1) \leq \log z$. \square

Lemma 3. For any g satisfying $\hat{\ell}^{\text{recall}}(g) \leq \bar{r}$ and $\frac{1}{m} \sum_{i \in [m]} \log \frac{n_g(x_i) \wedge 4n_{g^\dagger}(x_i)}{n_{g^\dagger}(x_i) \wedge 4n_g(x_i)} \leq 6\bar{r} + 4\bar{p} + \frac{2}{m}$, we have $\frac{1}{m} \sum_{i \in [m]: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} \leq 18\bar{r} + 12\bar{p} + \frac{4}{m}$.

Proof of Lemma 3 Let $B = \{i | \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq \frac{1}{2}\}$ and $C = \{i | \frac{1}{2} \leq \frac{n_{g^\dagger}(x_i)}{n_{g^{\text{target}}}(x_i)} \leq 2\}$. For any $i \in B$, the value of $\log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)}$ is in $[-1, 1]$. Then we have

$$\begin{aligned}
& \sum_{i \in [m]: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq 1} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} \\
&= \sum_{i \in B} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} - \sum_{i \in B: \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} < 1} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} \\
&\leq \sum_{i \in C \cap B} \log \frac{n_g(x_i) \wedge 2n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} + |\neg C| - \sum_{i \in B} \log \frac{n_g(x_i) \wedge n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} \\
&\leq \sum_{i \in C \cap B} \log \frac{n_g(x_i) \wedge 4n_{g^\dagger}(x_i)}{n_{g^{\text{target}}}(x_i)} + |\neg C| - \sum_{i \in B} \log \frac{n_g(x_i) \wedge n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} \\
&= \sum_{i \in C \cap B} \log \frac{n_g(x_i) \wedge 4n_{g^\dagger}(x_i)}{n_{g^\dagger}(x_i) \wedge 4n_g(x_i)} + \sum_{i \in C \cap B} \log \frac{n_{g^\dagger}(x_i) \wedge 4n_g(x_i)}{n_{g^{\text{target}}}(x_i)} + |\neg C| - \sum_{i \in B} \log \frac{n_g(x_i) \wedge n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} \\
&\leq \sum_{i=1}^m \log \frac{n_g(x_i) \wedge 4n_{g^\dagger}(x_i)}{n_{g^\dagger}(x_i) \wedge 4n_g(x_i)} + 2|\neg C| + 2|\neg B| + \sum_{i \in C \cap B} \log \frac{n_{g^\dagger}(x_i)}{n_{g^{\text{target}}}(x_i)} + |\neg C| - \sum_{i \in B} \log \frac{n_g(x_i) \wedge n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} \\
&\leq 18m\bar{r} + 12m\bar{p} + 4. \tag{Applying Lemmas 4, 5 and 6}
\end{aligned}$$

Note that in the second last inequality, we adopt the fact that $\frac{x \wedge 4y}{y \wedge 4x} \in [\frac{1}{4}, 4]$ for all $x, y > 0$. \square

Lemma 4. *For any graph g with $\widehat{\ell}^{\text{recall}}(g) \leq \bar{r}$ and any subset $S \subset [m]$, we have*

$$\sum_{i \in S \cap B} \log \frac{n_g(x_i) \wedge n_{g^{\text{target}}}}{n_{g^{\text{target}}}(x_i)} \geq -2m\bar{r} - 1.$$

where $B = \{i \mid \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \geq \frac{1}{2}\}$.

Proof of Lemma 4 Now let's focus on the rounds in $S \cap B$. We have

$$\sum_{i \in S \cap B} \left(1 - \frac{n_g(x_i) \wedge n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)}\right) \leq \sum_{i \in S \cap B} \ell^{\text{recall}}(g, x_i) \leq m\bar{r}.$$

Our problem becomes computing

$$\begin{aligned} & \min \sum_{i \in S \cap B} \log \frac{n_g(x_i) \wedge n_{g^{\text{target}}}}{n_{g^{\text{target}}}(x_i)} \\ \text{s.t. } & \sum_{i \in S \cap B} \frac{n_g(x_i) \wedge n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} \geq |S \cap B| - m\bar{r}. \end{aligned}$$

By applying Lemma 7, we know

$$\min \sum_{i \in S \cap B} \log \frac{n_g(x_i) \wedge n_{g^{\text{target}}}}{n_{g^{\text{target}}}(x_i)} \geq -2m\bar{r} - 1.$$

Thus, we have

$$\sum_{i \in S} \log \frac{n_g(x_i) \wedge n_{g^{\text{target}}}}{n_{g^{\text{target}}}(x_i)} \geq \min \sum_{i \in S \cap B} \log \frac{n_g(x_i) \wedge n_{g^{\text{target}}}}{n_{g^{\text{target}}}(x_i)} - |\neg B| \geq -4m\bar{r} - 1.$$

\square

Lemma 5. *For any graph g with $\widehat{\ell}^{\text{precision}}(g) \leq \bar{p}$ and any subset $S \subset [m]$, we have*

$$\sum_{i \in S \cap A} \log \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \leq 2m\bar{p} + 1.$$

where $A = \{i \mid \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \leq 2\}$.

Proof of Lemma 5 Since the empirical precision loss is bounded by \bar{p} , we have

$$\sum_{i \in S \cap A} \left(1 - \frac{n_{g^{\text{target}}}(x_i) \wedge n_g(x_i)}{n_g(x_i)}\right) \leq \sum_{i \in S \cap A} \ell^{\text{precision}}(g, x_i) \leq m\bar{p}.$$

By re-arranging terms, we have

$$\sum_{i \in S \cap A} \frac{n_{g^{\text{target}}}(x_i) \wedge n_g(x_i)}{n_g(x_i)} \geq |S \cap A| - m\bar{p}.$$

By applying Lemma 7, we have

$$\min \sum_{i \in S \cap A} \log \frac{n_{g^{\text{target}}}(x_i) \wedge n_g(x_i)}{n_g(x_i)} \geq -2m\bar{p} - 1.$$

Hence, we have

$$\sum_{i \in S \cap A} \log \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} \leq \sum_{i \in S \cap A} \log \frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i) \wedge n_g(x_i)} \leq 2m\bar{p} + 1.$$

□

Lemma 6. For any g with $\widehat{\ell}^{\text{recall}}(g) \leq \bar{r}$, we have

$$\sum_i \mathbf{1}\left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} < \frac{1}{2}\right) < 2m\bar{r}.$$

For any g with $\widehat{\ell}^{\text{precision}}(g) \leq \bar{p}$, we have

$$\sum_i \mathbf{1}\left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} > 2\right) < 2m\bar{p}.$$

Proof When $\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} < \frac{1}{2}$, we have

$$\ell^{\text{recall}}(g, x_i) \geq 1 - \frac{n_g(x_i) \wedge n_{g^{\text{target}}}(x_i)}{n_{g^{\text{target}}}(x_i)} > \frac{1}{2}.$$

Thus, we have $\sum_i \mathbf{1}\left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} < \frac{1}{2}\right) < 2m\bar{r}$. Similarly, when $\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} > 2$, we have

$$\ell^{\text{precision}}(g, x_i) \geq 1 - \frac{n_g(x_i) \wedge n_{g^{\text{target}}}(x_i)}{n_g(x_i)} > \frac{1}{2}.$$

Thus, we have $\sum_i \mathbf{1}\left(\frac{n_g(x_i)}{n_{g^{\text{target}}}(x_i)} > 2\right) < 2m\bar{p}$. □

Lemma 7. For any $k \in \mathbb{N}_+$, $c \geq 0$, let OPT denote the optimal value to the following constrained optimization problem:

$$\begin{aligned} & \min_{a_{1:k}} \sum_{i=1}^k \log a_i \\ & \text{s.t. } \sum_{i=1}^k a_i \geq k - c, \\ & \frac{1}{2} \leq a_i \leq 1, \forall i \in [k]. \end{aligned}$$

We have $OPT \geq -2c - 1$.

Proof of Lemma 7 We prove the lemma by showing that in the optimal solution, there will be at most one entry of $a_{1:k}$ not in $\{\frac{1}{2}, 1\}$. In this case, there are $\lfloor 2c \rfloor$ many $\frac{1}{2}$'s and one $c - \frac{\lfloor 2c \rfloor}{2}$. Then, we have

$$\sum_{i=1}^k \log a_i \geq -2c - 1.$$

Hence, it suffices to prove that in the optimal solution, there will be at most one entry of $a_{1:k}$ not in $\{\frac{1}{2}, 1\}$. Suppose that there are two entries $a_1 < a_2 \in (\frac{1}{2}, 1)$. For any $\Delta > 0$ s.t. $a_1 - \Delta, a_2 + \Delta \in [\frac{1}{2}, 1]$, we have

$$\log\left(\frac{a_2 + \Delta}{a_2}\right) < \log\left(\frac{a_1}{a_1 - \Delta}\right),$$

which is due to $\frac{x+\Delta}{x}$ is monotonically decreasing in x . By re-arranging terms, we have

$$\log(a_1 - \Delta) + \log(a_2 + \Delta) < \log(a_1) + \log(a_2).$$

Hence, we can always decrease the function value by changing a_1, a_2 to $a_1 - \Delta, a_2 + \Delta$. By setting $\Delta = (1 - a_2) \wedge (a_1 - \frac{1}{2})$, either a_1 is changed to $\frac{1}{2}$ or a_2 is changed to 1. We reduce the number of entries not being $\frac{1}{2}$ or 1. We are done with the proof. \square

5.3 Surrogate Loss Method in Both Realizable and Agnostic Cases

Again we focus empirical precision and recall losses minimization. Let $x_1, \dots, x_m \in \mathcal{X}$ denote a sequence of users. For each graph g , let U_i^g denote the uniform distribution over the neighbors $N_g(x_i)$. for any set $S \subset N_g(x_i)$. For any pair of graphs g', g'' , define the following:

$$\text{precision.loss}(g' | g'') = \text{recall.loss}(g'' | g') = \frac{1}{m} \sum_{i=1}^m U_i^{g'}(N_{g'}(x_i) \setminus N_{g''}(x_i)).$$

Here $\text{precision.loss}(g' | g'')$ is the precision loss of graph g' when the target graph is g'' and $\text{recall.loss}(g'' | g')$ is the recall loss of graph g'' when the target graph is g' . Thus, the goal is to output a graph g with small $\text{precision.loss}(g | g^{\text{target}})$ and $\text{recall.loss}(g | g^{\text{target}})$.

Our learning rule is based on two simple principles for discarding sub-optimal hypotheses. We illustrate these principles with the following intuitive example: consider a music recommendation system, and assume we are considering two candidate hypotheses, g_1 and g_2 . Both hypotheses recommend classical music; however, g_1 recommends pieces by Bach 20% of the time and pieces by Mozart 10% of the time, while g_2 never recommends any pieces by Mozart.

Now, suppose that in the training set, users frequently choose to listen to pieces by Mozart. This observation suggests that g_2 should be discarded, as it never recommends Mozart. This leads to our first rule: if a hypothesis exhibits a high recall loss, it can be discarded. The second rule addresses precision loss, which is more challenging because it cannot be directly estimated from the data. To illustrate the second rule, imagine that in the training set, users tend to pick Bach pieces only 5% of the time. This suggests that g_1 is over-recommending Bach pieces, and therefore, g_1 might also be discarded based on its likely precision loss.

We formally capture this using the following metric.

Definition 1. For a graph g define a vector $v_g : \mathcal{H} \times \mathcal{H} \rightarrow [0, 1]$ by

$$v_g(g', g'') = \frac{1}{m} \sum_{i=1}^m U_i^g(N_{g'}(x_i) \setminus N_{g''}(x_i)).$$

Define a metric $d_{\mathcal{H}}$ between graphs by $d_{\mathcal{H}}(h, k) = \|v_h - v_k\|_{\infty}$.

Let \hat{g} be the observed (empirical) graph; i.e. the graph in which every x_i is connected to the random number v_i which is observed in the training set. A standard union bound argument yields:

Lemma 8. Let g^{target} denote the true graph (i.e. the data is generated from g^{target}). Then, with probability at least $1 - \delta$:

$$d_{\mathcal{H}}(\hat{g}, g^{\text{target}}) \leq O\left(\sqrt{\frac{\log|\mathcal{H}| + \log(1/\delta)}{m}}\right).$$

We now present our algorithm. We present two variants, one in the realizable setting (when $g^{\text{target}} \in \mathcal{H}$) and one in the general (agnostic) setting.

Algorithm (realizable case): Let ε denote the desired error. Output a graph $g^{\text{output}} \in \mathcal{H}$ such that

1. For all $g \in \mathcal{H}$, $v_{\hat{g}}(g, g^{\text{output}}) = 0$.
2. For all $g \in \mathcal{H}$, $v_{g^{\text{output}}}(g^{\text{output}}, g) \geq \varepsilon \implies v_{\hat{g}}(g^{\text{output}}, g) > 0$,

Notice that Item 1 corresponds to the first principle for discarding suboptimal graphs described earlier in this section, while Item 2 corresponds to the second principle.

Algorithm (agnostic case): output a graph $g^{\text{output}} \in \mathcal{H}$ such that

$$d_{\mathcal{H}}(\hat{g}, g^{\text{output}}) = \min_{g \in \mathcal{H}} d_{\mathcal{H}}(\hat{g}, g).$$

We prove that

Theorem 8. Let g^{target} denote the target graph. Then, for

$$m = O\left(\frac{\log|\mathcal{H}| + \log(1/\delta)}{\varepsilon^2}\right),$$

the agnostic-case algorithm outputs a graph g^{output} such that with probability at least $1 - \delta$,

$$\ell^{\text{scalar}}(g^{\text{output}}) \leq 5 \min_{g \in \mathcal{H}} \ell^{\text{scalar}}(g) + \varepsilon.$$

Remark 2. In the realizable setting, our algorithm achieves a quadratic improvement in sample complexity: learning with recall and precision losses at most ε can be achieved with $O\left(\frac{\log|\mathcal{H}| + \log(1/\delta)}{\varepsilon}\right)$ examples.

5.3.1 Proof of Theorem 2

Proof of Theorem 2 For simplicity, we adopt payoffs instead of losses here. The payoff of graph g at x is

$$u(g, x) = \frac{|N_{g^{\text{target}}}(x) \cap N_g(x)|}{2|N_g(x)|} + \frac{|N_{g^{\text{target}}}(x) \cap N_g(x)|}{2|N_{g^{\text{target}}}(x)|}.$$

If $N_g(x) = \emptyset$ and $N_{g^{\text{target}}}(x) \neq \emptyset$, $u(g, x) = \frac{1}{2}$; if both are empty set $u(g, x) = 1$. The expected payoff is $u(g) = \mathbb{E}_{x \sim \mathcal{D}} [u(g, x)] = 1 - \ell^{\text{scalar}}(g)$.

Construction of g^{target} and \mathcal{D} Let's start by focusing on one single point x and its neighborhood. There are n nodes $N_1(x) = [\frac{n}{2}]$ and $N_2(x) = \{\frac{n}{2} + 1, \dots, n\}$. Consider two graphs— g_1 with $N_{g_1}(x) = N_1(x)$ and g_2 with $N_{g_2}(x) = N_1(x) \cup N_2(x)$. So $N_{g_1}(x)$ contains half of the nodes in $N_{g_2}(x)$.

In a world characterized by $\beta \in [\frac{1}{8}, \frac{2}{3}]$, $N_{g^{\text{target}}}(x)$ is generated in the following random way: Randomly select $\frac{3}{4} \cdot \beta n$ nodes from $N_1(x)$ and $\frac{1}{4} \cdot \beta n$ nodes from $N_2(x)$. We denote this distribution by P_β . No matter what β is, w.p. $\frac{3}{4}$, v is sampled uniformly at random from $N_1(x)$ and w.p. $\frac{1}{4}$, v is sampled uniformly at random from $N_2(x)$. That is, every node in $N_1(x)$ has probability $\frac{3}{2n}$ of being sampled and every node in $N_2(x)$ has probability $\frac{1}{2n}$ of being sampled.

For any g^{target} generated from the above process, the payoff of g_1 at x is

$$u(g_1, x) = \frac{|N_{g^{\text{target}}}(x) \cap N_{g_1}(x)|}{2|N_{g_1}(x)|} + \frac{|N_{g^{\text{target}}}(x) \cap N_{g_1}(x)|}{2|N_{g^{\text{target}}}(x)|} = \frac{3/4 \cdot \beta n}{n} + \frac{3/4 \cdot \beta n}{2\beta n} = \frac{3}{4}\beta + \frac{3}{8},$$

and the payoff of g_2 at x is

$$u(g_2, x) = \frac{|N_{g^{\text{target}}}(x) \cap N_{g_2}(x)|}{2|N_{g_2}(x)|} + \frac{|N_{g^{\text{target}}}(x) \cap N_{g_2}(x)|}{2|N_{g^{\text{target}}}(x)|} = \frac{\beta n}{2n} + \frac{\beta n}{2\beta n} = \frac{1}{2}\beta + \frac{1}{2}.$$

We make infinite copies of $\{x, N_1(x), N_2(x)\}$. In each of the copy, $N_{g_1}(x) = N_1(x)$ and $N_{g_2}(x) = N_1(x) \cup N_2(x)$. For each x , we independently sample $N_{g^{\text{target}}}(x)$ from P_β . Let the data distribution over all of such copies of x . Then almost surely, there is no repentance in the training data, i.e., there does not exist $i \neq j$ such that $x_i = x_j$. And for any random sampled test point, w.p. 1, it has not been sampled in the training set.

Analysis For any unobserved $x \notin \{x_i | i \in [m]\}$, let $\alpha_1 = \frac{|N_{g^{\text{output}}}(x) \cap N_1(x)|}{n}$ and $\alpha_2 = \frac{|N_{g^{\text{output}}}(x) \cap N_2(x)|}{n}$. Note that α_1, α_2 are in $[0, \frac{1}{2}]$ and are possibly random variables if \mathcal{A} is randomized. Then the expected (over the randomness of g^{target}) payoff of g^{output} at x is

$$\begin{aligned} \mathbb{E}_{g^{\text{target}}} [u(g^{\text{output}}, x)] &= \mathbb{E}_{g^{\text{target}}} \left[\frac{|N_{g^{\text{target}}}(x) \cap N_{g^{\text{output}}}(x)|}{2|N_{g^{\text{output}}}(x)|} + \frac{|N_{g^{\text{target}}}(x) \cap N_{g^{\text{output}}}(x)|}{2|N_{g^{\text{target}}}(x)|} \right] \\ &= \frac{\alpha_1 n \cdot \frac{3}{2}\beta + \alpha_2 n \cdot \frac{1}{2}\beta}{2(\alpha_1 + \alpha_2)n} + \frac{\alpha_1 n \cdot \frac{3}{2}\beta + \alpha_2 n \cdot \frac{1}{2}\beta}{2\beta n} \\ &= \frac{\alpha_1 \beta}{2(\alpha_1 + \alpha_2)} + \frac{\beta}{4} + \frac{3}{4}\alpha_1 + \frac{1}{4}\alpha_2, \end{aligned} \tag{5}$$

which is monotonically increasing in α_1 . Hence $\mathbb{E}_{g^{\text{target}}} [u(g^{\text{output}}, x)]$ is maximized at $\alpha_1 = \frac{1}{2}$. Then

$$\mathbb{E}_{g^{\text{target}}} [u(g^{\text{output}}, x)] \leq \frac{\beta}{4} \cdot \left(\frac{1}{\frac{1}{2} + \alpha_2} + 1 \right) + \frac{3}{8} + \frac{1}{4}\alpha_2.$$

Note that β is not observable if we never sample the same x more than once (and thus the distribution of v conditional on β is identical for any β). Hence g^{output} is independent of β .

- If $\mathcal{P}_{x \sim \mathcal{D}}(\alpha_2(x) \leq \frac{1}{4}) \geq \frac{1}{2}$: when $\beta = \frac{1}{8}$, $\mathbb{E}_{g^{\text{target}}} [u(g^{\text{output}}, x)] \leq \frac{1}{4}(\frac{1}{4+8\alpha_2} + \alpha_2) + \frac{13}{32}$ is monotonically increasing in α_2 . Hence,

$$u(g_2) - \mathbb{E}_{g^{\text{target}}} [u(g^{\text{output}})] \geq \frac{1}{2}(\frac{9}{16} - \frac{49}{96}) = \frac{5}{192} = \frac{5}{84} \ell^{\text{scalar}}(g_2).$$

- If $\mathcal{P}_{x \sim \mathcal{D}}(\alpha_2(x) > \frac{1}{4}) \geq \frac{1}{2}$: when $\beta = \frac{2}{3}$, $\mathbb{E}_{g^{\text{target}}} [u(g^{\text{output}}, x)] = \frac{1}{3+6\alpha_2} + \frac{1}{4}\alpha_2 + \frac{13}{24}$ is maximized at $\alpha_2 = \frac{1}{2}$ for $\alpha \in [\frac{1}{4}, \frac{1}{2}]$. Hence,

$$u(g_1) - \mathbb{E}_{g^{\text{target}}} [u(g^{\text{output}})] \geq \frac{1}{2}(\frac{7}{8} - \frac{5}{6}) = \frac{1}{48} = \frac{1}{6} \ell^{\text{scalar}}(g_1).$$

Therefore, for any algorithm \mathcal{A} , for any $x_{1:m}, v_{1:m}$, $\mathbb{E}_{g^{\text{target}}} [\ell^{\text{scalar}}(g^{\text{output}})]$ is worse than $1.05 \cdot \min\{\ell^{\text{scalar}}(g_1), \ell^{\text{scalar}}(g_2)\}$ at either $\beta = \frac{1}{8}$ or $\beta = \frac{2}{3}$. So there exists a target graph such that $\ell^{\text{scalar}}(g^{\text{output}}) \geq 1.05 \cdot \min\{\ell^{\text{scalar}}(g_1), \ell^{\text{scalar}}(g_2)\}$. \square

5.3.2 Proof of Theorem 8

We use the following auxiliary metric between graphs:

Definition 2. For two graphs g', g'' define

$$\begin{aligned} d_{\text{p,r}}(g', g'') &= \text{precision.loss}(g'|g'') + \text{recall.loss}(g'|g'') \\ &= \text{precision.loss}(g''|g') + \text{recall.loss}(g''|g'). \end{aligned}$$

For any graph g , the scalar loss $\ell^{\text{scalar}}(g) = \frac{1}{2}d_{\text{p,r}}(g^{\text{output}}, g^{\text{target}})$. In the remainder of this section, we focus on proving Theorem 8. The basic idea is to show that $d_{\mathcal{H}}$ can be used as a surrogate for $d_{\text{p,r}}$. The following lemma plays a crucial role in our proof.

Lemma 9. For every pair of graphs h, k :

$$d_{\mathcal{H}}(h, k) \leq d_{\text{p,r}}(h, k).$$

If in addition $h, k \in \mathcal{H}$, we have:

$$d_{\text{p,r}}(h, k) \leq 2d_{\mathcal{H}}(h, k).$$

We first use Lemma 9 to prove Theorem 8, and later prove the Lemma.

Proof of Theorem 8 Assume $m = O(\frac{\log|\mathcal{H}| + \log(1/\delta)}{\varepsilon^2})$ is such that $d_{\mathcal{H}}(g^{\text{output}}, g^{\text{target}}) \leq \varepsilon/4$ with probability at least $1 - \delta$, and assume the latter event holds. Let $g \in \mathcal{H}$, by the triangle inequality:

$$d_{\text{p,r}}(g^{\text{output}}, g^{\text{target}}) \leq d_{\text{p,r}}(g^{\text{output}}, g) + d_{\text{p,r}}(g, g^{\text{target}}).$$

We upper bound the first term on the right-hand side as follows:

$$\begin{aligned} d_{\text{p,r}}(g^{\text{output}}, g) &\leq 2d_{\mathcal{H}}(g^{\text{output}}, g) && \text{(Lemma 9)} \\ &\leq 2d_{\mathcal{H}}(g^{\text{output}}, g^{\text{target}}) + 2d_{\mathcal{H}}(g^{\text{target}}, g) \\ &\leq 4d_{\mathcal{H}}(g^{\text{target}}, g) + \varepsilon && \text{(see below)} \\ &\leq 4d_{\text{p,r}}(g^{\text{target}}, g) + \varepsilon. && \text{(Lemma 9)} \end{aligned}$$

Altogether,

$$d_{\mathbf{p},\mathbf{r}}(g^{\text{output}}, g^{\text{target}}) \leq 5d_{\mathbf{p},\mathbf{r}}(g, g^{\text{target}}) + \varepsilon.$$

It remains to explain the second to last inequality above. It follows by two applications of the triangle inequality:

$$\begin{aligned} d_{\mathcal{H}}(g^{\text{output}}, g^{\text{target}}) &\leq d_{\mathcal{H}}(g^{\text{output}}, \widehat{g}) + \varepsilon/4 && (d_{\mathcal{H}}(g^{\text{target}}, \widehat{g}) \leq \varepsilon/4) \\ &\leq d_{\mathcal{H}}(g, \widehat{g}) + \varepsilon/4 && (g^{\text{output}} \in \arg \min_{g \in \mathcal{H}} d_{\mathcal{H}}(g, \widehat{g})) \\ &\leq d_{\mathcal{H}}(g, g^{\text{target}}) + \varepsilon/2. && (d_{\mathcal{H}}(g^{\text{target}}, \widehat{g}) \leq \varepsilon/4) \end{aligned}$$

□

Proof of Lemma 9 For the first inequality, note that both of the distributions U_i^h and U_i^k are uniform over their supports and hence $\text{TV}(U_i^h, U_i^k) = \max\{U_i^h(N_h(x_i) \setminus N_k(x_i)), U_i^k(N_k(x_i) \setminus N_h(x_i))\}$. Thus, for every $g', g'' \in \mathcal{H}$:

$$\begin{aligned} &|U_i^h(N_{g'}(x_i) \setminus N_{g''}(x_i)) - U_i^k(N_{g'}(x_i) \setminus N_{g''}(x_i))| \\ &\leq \text{TV}(U_i^h, U_i^k) \\ &\leq U_i^h(N_h(x_i) \setminus N_k(x_i)) + U_i^k(N_k(x_i) \setminus N_h(x_i)). \end{aligned}$$

Hence, by averaging the above inequalities over $i = 1, \dots, n$:

$$d_{\mathcal{H}}(h, k) \leq \frac{1}{m} \sum_{i=1}^m \text{TV}(U_i^h, U_i^k) \leq d_{\mathbf{p},\mathbf{r}}(h, k),$$

which yields the first inequality.

For the second inequality, assume $h, k \in \mathcal{H}$. Thus,

$$\begin{aligned} d_{\mathcal{H}}(h, k) &\geq \max\left\{\frac{1}{m} \sum_{i=1}^m U_i^h(N_h(x_i) \setminus N_k(x_i)), \frac{1}{m} \sum_{i=1}^m U_i^k(N_k(x_i) \setminus N_h(x_i))\right\} \\ &\geq \frac{1}{m} \sum_{i=1}^m \frac{U_i^h(N_h(x_i) \setminus N_k(x_i)) + U_i^k(N_k(x_i) \setminus N_h(x_i))}{2} \\ &= \frac{1}{2} d_{\mathbf{p},\mathbf{r}}(h, k). \end{aligned}$$

□

5.4 Algorithm and Proofs in the Semi-Realizable Case

In the semi-realizable case, there exists a hypothesis in the class with zero precision loss. The question is whether we achieve zero precision loss while allowing for the worst recall loss in the class.

Theorem 5. *There exists an algorithm such that if there exists a graph $g' \in \mathcal{H}$ with $\ell^{\text{precision}}(g') = 0$ and $\ell^{\text{recall}}(g') = r$, then given an IID training set of size $O(\frac{\log(|\mathcal{H}|/\delta)}{\Delta^2})$, with probability $1 - \delta$, it outputs a graph with $\ell^{\text{precision}}(g^{\text{output}}) = 0$ and $\ell^{\text{recall}}(g^{\text{output}}) = r$.*

The algorithm works as follows.

Algorithm: output

$$g^{\text{output}} = \arg \min_{g \in \mathcal{H}} \sum_{i=1}^m \frac{\mathbb{1}(v_i \in N_g(x_i))}{n_g(x_i)}.$$

If there are multiple solutions, we break ties by picking the graph with smallest empirical recall loss.

Proof For any graph g , $\mathbb{1}(v_i \in N_g(x_i))$ is an unbiased estimate of the recall $\frac{|N_g(x_i) \cap N_{g^{\text{target}}}(x_i)|}{n_{g^{\text{target}}}(x_i)}$. Thus, $\frac{\mathbb{1}(v_i \in N_g(x_i))}{n_g(x_i)}$ is an unbiased estimate of $\frac{|N_g(x_i) \cap N_{g^{\text{target}}}(x_i)|}{n_g(x_i) \cdot n_{g^{\text{target}}}(x_i)}$. Since g' has zero precision loss, $\frac{|N_{g'}(x_i) \cap N_{g^{\text{target}}}(x_i)|}{n_{g'}(x_i)} = 1$ almost everywhere. Thus, we have

$$\left| \frac{1}{m} \sum_{i=1}^m \frac{\mathbb{1}(v_i \in N_g(x_i))}{n_g(x_i)} - \mathbb{E} \left[\frac{1}{n_{g^{\text{target}}}(x)} \right] \right| \leq \sqrt{\frac{\log(|\mathcal{H}|) + \log(1/\delta)}{m}},$$

for all $g \in \mathcal{H}$. Then if $\Delta_{\mathcal{D}} > 0$, we need $\frac{1}{\Delta_{\mathcal{D}}^2}$ samples to separate g' from other graphs in the hypothesis class. \square

Theorem 6. *There exists a class $\mathcal{H} = \{g_1, g_2\}$ of two graphs, for any $m > 0$ and any (possibly randomized improper) algorithm \mathcal{A} , there exists a target graph g^{target} and a data distribution \mathcal{D} for which there exists a graph $g^\dagger \in \mathcal{H}$ with $\ell^{\text{precision}}(g^\dagger) = 0$ s.t. with probability $1 - \delta$ over the training set, the expected (over the randomness of the algorithm) precision and recall losses of the output g^{output} satisfy either $\mathbb{E}[\ell^{\text{recall}}(g^{\text{output}})] \geq \min_{g \in \mathcal{H}} \ell^{\text{recall}}(g) + \Omega(1)$ or $\mathbb{E}[\ell^{\text{precision}}(g^{\text{output}})] = \Omega(1)$.*

Proof Let's start by focusing on one single point x . Let $N = \{v_1, \dots, v_n\}$ for some $n \gg m$. Let $N_{g_1}(x) = \{v_1\}$ and $N_{g_2}(x) = \{v_2\}$. In world I, $N_{g^{\text{target}}}(x)$ is generated in the following way.

- w.p. $\frac{1}{2}$, $N_{g^{\text{target}}}(x) = N \setminus \{v_2\}$.
- w.p. $\frac{1}{2}$, $N_{g^{\text{target}}}(x) = \{v_1, v_2\}$.

We construct a symmetric world II by switching v_1 and v_2 , i.e.,

- w.p. $\frac{1}{2}$, $N_{g^{\text{target}}}(x) = N \setminus \{v_1\}$.
- w.p. $\frac{1}{2}$, $N_{g^{\text{target}}}(x) = \{v_1, v_2\}$.

We make infinite independent copies of (x, N) and let \mathcal{D} to be the uniform distribution over such x 's. Hence, in world I, $\ell^{\text{precision}}(g_1) = 0$ and $\ell^{\text{recall}}(g_1) = \frac{3}{4} - \frac{1}{2n}$; $\ell^{\text{precision}}(g_2) = \frac{1}{2}$ and $\ell^{\text{recall}}(g_2) = \frac{3}{4}$. When $n \rightarrow \infty$, we can't distinguish between two worlds. In order to achieve $\ell^{\text{precision}}(g^{\text{output}}) = 0$ in both worlds, we need to make $N_{g^{\text{output}}}(x) = \emptyset$ for almost every x . Then the recall loss would be 1. \square

6 Discussion

In this work, we study PAC learning guarantees for precision and recall. There are two natural open questions.

First, there is a gap between the upper and lower bounds. For the scalar-loss objective, we demonstrate that an $\alpha = 5$ approximate optimal scalar loss is achievable, while $\alpha = 1.05$ is not, leaving it unclear what the optimal α is. For the Pareto-loss objective, we establish an upper bound of $(p, r) \Rightarrow (5(p+r), r)$ and a lower bound of $(p, r) \not\Rightarrow (p+0.01, r+0.01)$, again suggesting a gap that we do not yet know how to close.

Second, it remains an open question whether there exists a combinatorial measure, similar to the VC dimension in standard PAC learning, that characterizes the learnability of precision and recall. Each graph implicitly defines a distribution at each node—specifically, a uniform distribution over its neighborhood. In Section 5.3, we also link the scalar loss to the total variation distance, thus reducing the scalar loss learning problem to a special case of distribution learning. However, as shown in [LB24], there is no such a dimension characterizing the sample complexity of learning certain distribution classes (in their case, a mixture of point mass and uniform distributions). This result suggests a potential limitation in identifying a combinatorial measure for our learning problem.

Acknowledgements

Lee Cohen is supported by the Simons Foundation Collaboration on the Theory of Algorithmic Fairness, the Sloan Foundation Grant 2020-13941, and the Simons Foundation investigators award 689988.

Yishay Mansour was supported by funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (grant agreement No. 882396), by the Israel Science Foundation, the Yandex Initiative for Machine Learning at Tel Aviv University and a grant from the Tel Aviv University Center for AI and Data Science (TAD).

Shay Moran is a Robert J. Shillman Fellow; he acknowledges support by ISF grant 1225/20, by BSF grant 2018385, by Israel PBC-VATAT, by the Technion Center for Machine Learning and Intelligent Systems (MLIS), and by the the European Union (ERC, GENERALIZATION, 101039692). Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the European Research Council Executive Agency. Neither the European Union nor the granting authority can be held responsible for them.

Han Shao was supported by Harvard CMSA.

References

- [AGHHR05] Shivani Agarwal, Thore Graepel, Ralf Herbrich, Sariel Har-Peled, and Dan Roth. Generalization bounds for the area under the ROC curve. *J. Mach. Learn. Res.*, 6:393–425, 2005 (cited on page 4).
- [AHHM22] Noga Alon, Steve Hanneke, Ron Holzman, and Shay Moran. A theory of pac learnability of partial concept classes. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*, pages 658–671. IEEE, 2022 (cited on page 3).
- [AKV16] Monika Arora, Uma Kanjilal, and Dinesh Varshney. Evaluation of information retrieval: precision and recall. *International Journal of Indian Culture and Business Management*, 12(2):224–236, 2016 (cited on page 4).

- [BD20] Jessa Bekker and Jesse Davis. Learning from positive and unlabeled data: a survey. *Machine Learning*, 109(4):719–760, 2020 (cited on page 4).
- [BTDK22] Jasmin Bogatinovski, Ljupčo Todorovski, Sašo Džeroski, and Dragi Kocev. Comprehensive comparative study of multi-label classification methods. *Expert Systems with Applications*, 2022 (cited on page 4).
- [BKM19] Olivier Bousquet, Daniel Kane, and Shay Moran. The optimal approximation factor in density estimation. In *Conference on Learning Theory*, pages 318–341. PMLR, 2019 (cited on page 3).
- [CM03] Corinna Cortes and Mehryar Mohri. AUC optimization vs. error rate minimization. In *Advances in Neural Information Processing Systems 16 [Neural Information Processing Systems, NIPS 2003, December 8-13, 2003, Vancouver and Whistler, British Columbia, Canada]*, pages 313–320. MIT Press, 2003 (cited on page 4).
- [CM04] Corinna Cortes and Mehryar Mohri. Confidence intervals for the area under the ROC curve. In *Advances in Neural Information Processing Systems 17 [Neural Information Processing Systems, NIPS 2004, December 13-18, 2004, Vancouver, British Columbia, Canada]*, pages 305–312, 2004 (cited on page 4).
- [DSBS15] Amit Daniely, Sivan Sabato, Shai Ben-David, and Shai Shalev-Shwartz. Multiclass learnability and the erm principle. *J. Mach. Learn. Res.*, 16(1):2377–2404, 2015 (cited on page 3).
- [DDGL99] Francesco De Comit e, Franois Denis, R emi Gilleron, and Fabien Letouzey. Positive and unlabeled examples help learning. In *Algorithmic Learning Theory: 10th International Conference, ALT’99 Tokyo, Japan, December 6–8, 1999 Proceedings 10*, pages 219–230. Springer, 1999 (cited on page 4).
- [Den98] Franois Denis. PAC learning from positive statistical queries. In *International conference on algorithmic learning theory*, pages 112–126. Springer, 1998 (cited on page 4).
- [DL01] Luc Devroye and G abor Lugosi. *Combinatorial methods in density estimation*. Springer Science & Business Media, 2001 (cited on page 3).
- [EW01] Andr e Elisseeff and Jason Weston. A kernel method for multi-labelled classification. In *Advances in Neural Information Processing Systems*, 2001 (cited on page 4).
- [GZ11] Wei Gao and Zhi-Hua Zhou. On the consistency of multi-label learning. In *Proceedings of the 24th Annual Conference on Learning Theory*, Proceedings of Machine Learning Research, 2011 (cited on page 4).
- [GBV20] Margherita Grandini, Enrico Bagli, and Giorgio Visani. Metrics for multi-class classification: an overview. *arXiv preprint arXiv:2008.05756*, 2020 (cited on page 4).
- [JL19] Brendan Juba and Hai S Le. Precision-recall versus accuracy and the role of large data sets. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33 of number 01, pages 4039–4048, 2019 (cited on page 4).
- [KVJ12] Ashish Kapoor, Raajay Viswanathan, and Prateek Jain. Multilabel classification using bayesian compressed sensing. In *Advances in Neural Information Processing Systems*, 2012 (cited on page 4).
- [KSWA15] Branislav Kveton, Csaba Szepesvari, Zheng Wen, and Azin Ashkan. Cascading bandits: learning to rank in the cascade model. In *International conference on machine learning*, pages 767–776. PMLR, 2015 (cited on page 4).

- [LB24] Tosca Lechner and Shai Ben-David. Inherent limitations of dimensions for characterizing learnability of distribution classes. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 3353–3374. PMLR, 2024 (cited on pages 3, 29).
- [LDG00] Fabien Letouzey, François Denis, and Rémi Gilleron. Learning from positive and unlabeled examples. In *International Conference on Algorithmic Learning Theory*, pages 71–85. Springer, 2000 (cited on page 4).
- [MMZ24] Anqi Mao, Mehryar Mohri, and Yutao Zhong. Multi-label learning with stronger consistency guarantees, 2024 (cited on page 4).
- [McC99] Andrew Kachites McCallum. Multi-label text classification with a mixture model trained by em. *AAAI’99 workshop on text learning*, 1999 (cited on page 4).
- [MRRK19] Aditya K Menon, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. Multilabel reductions: what is my loss optimising? In *Advances in Neural Information Processing Systems*, 2019 (cited on page 4).
- [PC11] James Petterson and Tibério Caetano. Submodular multi-label learning. In *Advances in Neural Information Processing Systems*, 2011 (cited on page 4).
- [Ros04] Saharon Rosset. Model selection via the AUC. In Carla E. Brodley, editor, *Machine Learning, Proceedings of the Twenty-first International Conference (ICML 2004), Banff, Alberta, Canada, July 4-8, 2004*, volume 69 of *ACM International Conference Proceeding Series*. ACM, 2004 (cited on page 4).
- [SBLBG18] Mehdi SM Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. Assessing generative models via precision and recall. *Advances in neural information processing systems*, 31, 2018 (cited on page 4).
- [SS00] Robert E. Schapire and Yoram Singer. Boostexter: a boosting-based system for text categorization. *Machine Learning*, 2000 (cited on page 4).
- [SB14] Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014 (cited on page 3).
- [TLZAG18] Nesime Tatbul, Tae Jun Lee, Stan Zdonik, Mejbah Alam, and Justin Gottschlich. Precision and recall for time series. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 1924–1934, 2018 (cited on page 4).
- [TR09] Luis Torgo and Rita Ribeiro. Precision and recall for regression. In *Discovery Science: 12th International Conference, DS 2009, Porto, Portugal, October 3-5, 2009 12*, pages 332–346. Springer, 2009 (cited on page 4).
- [Val84] L. G. Valiant. A theory of the learnable, 1984 (cited on page 1).
- [ZZ14] Min-Ling Zhang and Zhi-Hua Zhou. A review on multi-label learning algorithms. *IEEE Transactions on Knowledge and Data Engineering*, 2014 (cited on page 4).