LEARNING FROM INTERVAL TARGETS

Anonymous authors

Paper under double-blind review

Abstract

We consider regression problems where the exact real-valued targets are not directly available; instead, supervision is provided in the form of intervals around the targets—that is, only lower and upper bounds are known. Such a "learning from interval targets" setup arises in domains where labeling costs are high or there is inherent uncertainty in the target values. In these settings, traditional regression loss functions, which require exact target values, cannot be directly applied. To address this challenge, we propose two approaches: (i) modifying the regression loss function to be compatible with interval ground truths, and (ii) formulating a min-max problem where we minimize the typical regression loss with respect to the "worst-case" label within the interval. We provide theoretical guarantees for our methods, analyze their computational efficiency, and evaluate their practical performance on real-world datasets.

021 022

000

001 002 003

004

006 007

008 009

010

011

012

013

014

015

016

017

018

019

1 INTRODUCTION

Supervised learning has achieved significant empirical success, largely due to the availability of extensive labeled datasets. However, in many real-world tasks, obtaining target labels is challenging, which hampers the performance of these methods. This difficulty arises either from high labeling costs—for example, certain medical measurements are expensive—or from practical limitations, such as sensors that only record target values at discrete intervals (e.g., every hour), leaving intermediate values unobserved. Prior work has addressed this issue by incorporating additional information into the learning pipeline. For instance, some approaches encourage model outputs to be smooth over unlabeled data (Zhu, 2005; Chapelle et al.), while others enforce models to satisfy constraints derived from domain knowledge, such as physical laws (Willard et al., 2020; Swischuk et al., 2019).

033 In this work, we focus on regression tasks where only the lower and upper bounds of the target 034 values (intervals) are available. Our setting relates to both weak supervision and learning with side information. Learning with interval targets generalizes supervised learning, which corresponds to 035 the special case where the lower and upper bounds are equal. For many tasks, it is easier and 036 more practical for human labelers to provide interval targets instead of precise single values; thus, 037 these intervals can be viewed as a form of weak supervision. Additionally, in various settings, such intervals are readily available for unlabeled data, either from domain knowledge or inherent properties of the data, serving as side information. A prime example is bond pricing. Unlike stocks, 040 bonds are traded infrequently, so we may observe only a handful of trades over a given time span, 041 resulting in limited labeled data for bond prices. However, numerous bond quotes are available: the 042 bid and ask quotes represent the prices at which dealers are willing to buy and sell, respectively. 043 These bid and ask quotes can be treated as lower and upper bounds of the true bond prices when 044 actual trade prices are not observed.

A natural learning strategy for learning from interval targets is to learn a hypothesis whose outputs always lie within the provided intervals. Despite its simplicity, previous work (Cheng et al., 2023a) has shown that this method leads to a hypothesis that converges to the optimal one under two assumptions: (i) the true target function belongs to the hypothesis class, and (ii) the intervals have an ambiguity degree smaller than 1 (Section 2). However, these assumptions are unlikely to hold in practice. In particular, (ii) is often violated; for example, even in the simple case where the interval is a ball of radius ϵ around the target value y, the ambiguity degree equals 1. It is important to understand whether this approach can be effective under more relaxed assumptions. Our first contribution is a novel theoretical analysis of learning a hypothesis that lies within the target intervals. Our error bound, based on the Lipschitz constant of the hypothesis class, is applicable even when the ambiguity degree is 1 (Theorem 3.3) and can be extended to the agnostic setting where the target function may not belong to our hypothesis class (Theorem 3.8). Unlike prior bounds (Cheng et al., 2023a) that only consider the asymptotic case as the number of data points $n \to \infty$, our bound is non-asymptotic. The key insight is that, when the hypothesis class is smooth, the outputs for two close inputs cannot differ significantly. As a result, portions of the original intervals can be ruled out, leading to much smaller valid intervals (Figure 1b).

060 In our second contribution, we explore an alternative approach by learning a hypothesis that mini-061 mizes the loss with respect to the worst-case labels within the given intervals. Since we assume that 062 the true target values lie within these intervals, the worst-case loss serves as an upper bound on the 063 regression loss. We consider two variants of the second approach: i) we allow the worst-case labels 064 to be any points within the intervals, ii) we restrict the worst-case labels to be outputs of some hypothesis in our hypothesis class, thereby incorporating the smoothness property. We show that there 065 exists a distribution where the second variant can perform arbitrarily better than the first (Proposi-066 tion 4.4), indicating that constraining the worst-case labels to the hypothesis class is preferable in 067 the worst-case scenario. We demonstrate the effectiveness of both methods on real-world datasets. 068

069

071

0 1.1 RELATED WORK

Our problem is closely related to partial-label learning, where each training point is associated with 072 a set of candidate labels instead of a single target label (Cour et al., 2011; Ishida et al., 2017; Feng 073 et al., 2020a; Ishida et al., 2019; Yu et al., 2018). In classification with finite label sets, a popular 074 method is to minimize the average loss over the label set (Jin & Ghahramani, 2002), leading to 075 various extensions (Zhang et al., 2017; Wang et al., 2019; Xu et al., 2021; Wu et al., 2022; Gong 076 et al., 2022). Another important approach focuses on identifying the true label from the candidate 077 set (Lv et al., 2020; Zhang et al., 2016; Yu & Zhang, 2016). On the theoretical side, prior work has 078 studied learnability conditions (Liu & Dietterich, 2014; Cour et al., 2011) and developed statistically 079 consistent estimators (Lv et al., 2020; Feng et al., 2020b; Wen et al., 2021) based on the small 080 ambiguity degree assumption or specific label set generating distributions.

081 In regression, there has been less prior work. Cheng et al. (2023b) considers partial-label regression 082 with a finite label set and Cheng et al. (2023a) later extends it to the interval setting with infinitely 083 many labels. Both works provide statistically consistent estimators, but their theoretical results 084 heavily depend on the small ambiguity degree assumption. We note that the ambiguity degree, first 085 proposed for classification tasks in Cour et al. (2011), may not be suitable for regression tasks. In classification, a hypothesis is either correct or incorrect, and a small ambiguity degree ensures that, with enough observed label sets, we can recover the true label. However, in regression, we are 087 often satisfied with predictions that are sufficiently close to the target-for example, within an error 880 tolerance of ϵ —making the concept of ambiguity degree less applicable. 089

In our work, we study a projection loss, which is equivalent to the partial-label learning loss (PLL loss) in Lv et al. (2020), and can be seen as a generalized version of the limiting method in Cheng et al. (2023a). We provide a non-asymptotic error bound that does not rely on the ambiguity degree and extend our analysis to the agnostic setting. We provide additional related work in Appendix C.

094

096

1.2 PRELIMINARIES AND NOTATION

107 Let \mathcal{X} be the feature space and \mathcal{Y} be the label space. Let $f^*: \mathcal{X} \to \mathcal{Y}$ denote the target function. 108 We use uppercase letters (e.g., X) to represent random variables and lowercase letters (e.g., x) for 109 deterministic variables. We consider a regression problem where our goal is to learn a function 100 $f: \mathcal{X} \to \mathcal{Y}$ from a hypothesis class \mathcal{F} that approximates the target function f^* in the deterministic 101 label setting. Let \mathcal{D} be the distribution over $\mathcal{X} \times \mathcal{Y}$ where, for each $x \in \mathcal{X}$, the label y is determinis-102 tically given by $y = f^*(x)$ and let p be the pdf. **Our goal is to learn a function** f that minimizes 103 **the expected loss** $\operatorname{err}(f) := \mathbb{E}_{(X,Y)\sim\mathcal{D}} \left[\ell(f(X), Y)\right]$ for some loss function $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$.

Interval targets. However, we assume that we have access only to interval samples of the form $\{(x_i, l_i, u_i)\}_{i=1}^n$, where l_i and u_i are the lower and upper bounds of y_i , respectively. While we assume that the label is fixed to $f^*(x_i)$, we allow the intervals—that is, the bounds (l_i, u_i) —to be random and assume that each tuple (x_i, l_i, u_i) is sampled from some distribution \mathcal{D}_I . In this setting, we aim to explore learning strategies and determine what kinds of guarantees are possible.



Figure 1: (1a) An example of learning from intervals where the input is one dimension. The intervals are shown as gray boxes. A natural method is to learn a hypothesis that always lies within these intervals. Here, we illustrate two such hypotheses that are both valid but have different levels of smoothness. (1b) When the hypothesis is smooth (blue line), it lies within intervals much smaller than the original ones, depicted by the green region (Proposition 3.2). We can extend this result to hypotheses that approximately lie within the intervals (Theorem 3.3).

2 LEARNING FROM INTERVALS USING A PROJECTION LOSS

Since the target label y always lies within the interval [l, u], a natural strategy is to learn a hypothesis $f \in \mathcal{F}$ such that $f(x) \in [l, u]$ for all $x \in \mathcal{X}$ (Figure 1a). In previous work, Cheng et al. (2023a) analyzed the following strategy:

Learn f that minimizes the empirical risk of the 0-1 loss:
$$\sum_{i=1}^{n} \ell_{0-1}(f(x_i), l_i, u_i), \qquad (1)$$

where $\ell_{0-1}(f(x), l, u) := 1[f(x) < l] + 1[f(x) > u]$. Using ℓ_1 loss as the surrogate (equation (12)), they showed that f converges to f^* as $n \to \infty$ if two assumptions are satisfied, (i) Realizability, that is, $f^* \in \mathcal{F}$, (ii) Ambiguity degree is smaller than 1. Ambiguity degree is the maximum probability of a specific incorrect target y', belonging to the same interval [l, u] as the true target y:

Ambiguity degree
$$(\mathcal{D}, \mathcal{D}_I) := \sup_{(x,y)\sim\mathcal{D}, (x,l,u)\sim\mathcal{D}_I, y'\in\mathcal{Y}, y'\neq y} \Pr(y'\in[L,U]) < 1.$$
 (2)

139 These assumptions can be impractical and restrictive. First, our hypothesis class may not contain 140 f^* . Second, an ambiguity degree smaller than 1 implies that for any fixed x, if we keep sampling 141 the interval [l, u], the intersection of such intervals (in the limit) would only be the set of the true 142 target $\{y\}$; that is, we can recover the true y given an infinite number of intervals. However, this 143 assumption is unlikely to hold in practice because there is usually a gap between the upper and lower 144 bounds and the target y. For example, in the simple case where $[l, u] = [y - \epsilon, y + \epsilon]$ (a ball with 145 radius $\epsilon > 0$ around the true target y), the assumption fails since $y + \epsilon/2$ always lies within the interval at the same time with the true y. Finally, we note that the previous result is an asymptotic 146 bound that only applies when $n \to \infty$. In this work, we relax these assumptions and provide a 147 non-asymptotic generalization bound. 148

We begin by defining a suitable learning objective. Since the 0-1 loss above is not continuous, it is not suitable for gradient-based optimization techniques. To address this, we relax the loss by considering a projection

152 153

124 125

126 127

128

129

130 131 132

138

 $\pi_{\ell}(f(x),l,u) := \min_{\tilde{y} \in [l,u]} \ell(f(x),\tilde{y})$ (3)

for any surrogate loss function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$. We make the following assumption about ℓ ,

Assumption 1. The loss function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ can be written as $\ell(y, y') = \psi(|y - y'|)$ for some non-decreasing function ψ , and satisfies $\ell(y, y') = 0$ if and only if y = y'.

The following proposition shows that π_{ℓ} is a meaningful proxy for the 0-1 loss, and can be evaluated efficiently by only considering the boundaries of the interval.

Proposition 2.1. Suppose that $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ is a loss function that satisfies Assumption 1 then $\pi_{\ell}(f(x), l, u) = 0$ if and only if $f(x) \in [l, u]$, and we can write

$$\pi_{\ell}(f(x), l, u) = \mathbb{1}[f(x) < l]\ell(f(x), l) + \mathbb{1}[f(x) > u]\ell(f(x), u).$$
(4)

162 The proof is provided in Appendix D.1. In the rest of the paper, we refer to π_l as the **projection** 163 **loss.** Consequently, the informal goal given in (1) can be formalized as the following objective: 164

$$\min_{f} \sum_{i=1}^{n} \mathbb{1}[f(x_i) < l_i] \ell(f(x_i), l_i) + \mathbb{1}[f(x_i) > u_i] \ell(f(x_i), u_i).$$
(5)

3 THEORETICAL ANALYSIS OF THE PROJECTION APPROACH

170 Denote $\mathcal{F}_{\eta} := \{f \in \mathcal{F} \mid \mathbb{E}[\pi_{\ell}(f(X), L, U)] \leq \eta\}$ as a class of hypotheses for which the expected projection loss is smaller than η . When $\eta = 0$, we have $\hat{\mathcal{F}}_0 = \{f \in \mathcal{F} \mid \Pr(f(X) \in [L, U]) = 1\}$ which is a class of hypothesis that always lie within the interval, a property that we aim to achieve for our hypothesis f. However, since we only have access to a finite number of data points, we can only hope to learn $f \in \mathcal{F}_{\eta}$ for some small η . Specifically, by a standard uniform convergence argument (e.g. Mohri (2018)), we have that with probability at least $1 - \delta$ over the draws $(x_i, l_i, u_i) \sim \mathcal{D}_I$,

for all
$$f \in \mathcal{F}, \mathbb{E}[\pi_{\ell}(f(X), L, U)] \le \frac{1}{n} \sum_{i=1}^{n} \pi_{\ell}(f(x_i), l_i, u_i) + 2R_n(\Pi(\mathcal{F})) + M\sqrt{\frac{\ln(1/\delta)}{n}}.$$
 (6)

Here, $R_n(\Pi(\mathcal{F}))$ is the Rademacher complexity of the function class $\Pi(\mathcal{F}) := \{\pi_\ell(f(x), l, u) \mapsto$ $\mathbb{R} \mid f \in \mathcal{F}$ and we assume that the π_{ℓ} is uniformly bounded by M. Thus, given n, M, and the empirical loss on observed data (first term in R.H.S.), we have an **upper bound** of η which $f \in \mathcal{F}_{\eta}$.

183 **Plan of analysis.** In the rest of this section, we show general bounds for any $f \in \widetilde{\mathcal{F}}_{\eta}$. In Theorem 184 3.3 we show that for any x, f(x) belongs to a small interval that depends on M and η . This leads to 185 our main result: a generalization bound on the loss of f with respect to actual labels y, thus showing that regression can be done using interval targets (Section 3.3). 186

187 188

189

166 167

169

171

172

173

174

175

180

181

182

3.1 Effect of realizability and small ambiguity degree assumptions on $\widetilde{\mathcal{F}}_n$

We begin by examining the implications of the assumptions made in prior work (Section 2). The 190 realizability assumption can be restated as $f^* \in \mathcal{F}_0$ since the projection loss of f^* is always zero. 191 Second, the small ambiguity degree assumption implies that, for any x, the intersection of the inter-192 vals can only be the singleton set $\{y\}$. As a result, we have $\mathcal{F}_0 = \{f \in \mathcal{F} \mid \operatorname{err}(f) = 0\} \neq \emptyset$. 193

194 With these assumptions, we can show that minimizing the projection objective will converge to a hypothesis with zero error. The following informal argument summarizes the more elaborate 196 analysis of Cheng et al. (2023b). Here is the high-level idea: let f_n be the hypothesis that mini-197 mizes the empirical projection objective (5). Realizability implies that there exists $f^* \in \mathcal{F}$ with an expected loss of zero. Since f_n achieves the empirical risk no larger than that of f^* , it must 198 achieve an empirical risk of zero. From equation 6, we have $f_n \in \widetilde{\mathcal{F}}_{\eta_n}$ with high probability, where 199 200 201

 $\eta_n = 2R_n(\Pi(\mathcal{F})) + M\sqrt{\frac{\ln(1/\delta)}{n}}$. In general, $\eta_n = O(1/\sqrt{n})$, and as $n \to \infty$, we have $\eta_n \to 0$ which means that $\widetilde{\mathcal{F}}_{\eta_n} \to \widetilde{\mathcal{F}}_0$. Consequently, $\operatorname{err}(f_n) \to 0$ since any member of $\widetilde{\mathcal{F}}_0$ has zero error. 202

203 However, when the realizability and ambiguity degree assumptions do not hold, there may be $f \in \mathcal{F}_0$ 204 with $\operatorname{err}(f) > 0$. Additionally, with a finite amount of data, we can only learn a hypothesis $f \in \mathcal{F}_{\eta_n}$ 205 for some $\eta_n > 0$. In the next section, we will analyze $\widetilde{\mathcal{F}}_{\eta}$ without relying on the small ambiguity 206 degree assumption and in finite samples. Further, we also provide a formal argument that relates η 207 with n and provide an error bound depends on n in Appendix A. 208

209 3.2 PROPERTIES OF $\widetilde{\mathcal{F}}_{\eta}$ 210

211 Although our results extend to the probabilistic interval setting, where multiple intervals [l, u] are 212 drawn for each x, we focus on the deterministic interval setting for simplicity. In this case, each x is 213 associated with a fixed interval $[l_x, u_x]$. For a detailed discussion of the probabilistic interval setting, 214 please refer to Appendix E. Now, consider the following characterization of \mathcal{F}_0 . 215

Proposition 3.1. For any $f \in \mathcal{F}_0$, and ℓ that satisfies Assumption 1, we have $f(x) \in [l_x, u_x]$.

226

227

228 229

230 231

232

233

234 235 236

237

238 239 240

266

267 268



Figure 2: (2a) Based on the smoothness property, the difference between f(x) and f(x') cannot exceed m||x - x'||. As a result, the upper and lower bounds of f(x') imply the corresponding bounds for f(x). (2b) The lower bound gap of x' to x is defined as the difference between the lower bound of f(x) induced by x' and the largest lower bound $(\tilde{l}_{\mathcal{D}\to x}^{(m)})$; similarly for the upper bound gap. These gaps are crucial in bounding the size of $r_{\eta}(x)$ and $s_{\eta}(x)$ (how much we have to compensate when $f \in \tilde{\mathcal{F}}_{\eta}$) where larger gaps lead to larger values (Theorem 3.3).

We will show that the interval in which f(x) must lie can be shrunk (made smaller than $[l_x, u_x]$) if we assume that the class \mathcal{F} contains only *m*-Lipschitz functions. That is, for any $f \in \mathcal{F}$ and any $x, x' \in \mathcal{X}$, the condition $|f(x) - f(x')| \le m ||x - x'||$ holds (L_1 norm). We define for any x, x',

$$l_{x' \to x}^{(m)} := l_{x'} - m \|x - x'\|, u_{x' \to x}^{(m)} := u_{x'} + m \|x - x'\|.$$

Proposition 3.2. Let \mathcal{F} be a class of hypotheses that are *m*-Lipschitz and suppose that ℓ satisfies Assumption 1. Then for any $f \in \widetilde{\mathcal{F}}_0$ and for each x with p(x) > 0,

$$f(x) \in \bigcap [l_{x' \to x}^{(m)}, u_{x' \to x}^{(m)}] =: [l_{\mathcal{D} \to x}^{(m)}, u_{\mathcal{D} \to x}^{(m)}].$$
(7)

Proof. (Sketch) For $f \in \widetilde{\mathcal{F}}_0$, by Lipschitzness, for any $x, x' \in \mathcal{X}$, we have $|f(x) - f(x')| \leq m \|x - x'\|$ which implies $f(x') - m \|x - x'\| \leq f(x) \leq f(x') + m \|x - x'\|$. We can then replace f(x') with its lower and upper bound, $l_{x'}, u_{x'}$, respectively to achieve the result.

Interpretation. Note that $l_{x' \to x}^{(m)}$ and $u_{x' \to x}^{(m)}$ provide lower bound and upper bounds of f(x) induced by x', derived from the Lipschitz property of f (Figure 2a). We denote the intersection of all such intervals $[l_{x' \to x}^{(m)}, u_{x' \to x}^{(m)}]$ over all x' by $[l_{\mathcal{D} \to x}^{(m)}, u_{\mathcal{D} \to x}^{(m)}]$. First, we observe that $[l_{\mathcal{D} \to x}^{(m)}, u_{\mathcal{D} \to x}^{(m)}]$ is always smaller than $[l_x, u_x]$ because when we set x' = x, we have $[l_{x' \to x}^{(m)}, u_{x' \to x}^{(m)}] = [l_x, u_x]$. Second, as the hypothesis becomes more smooth, the interval $[l_{\mathcal{D} \to x}^{(m)}, u_{\mathcal{D} \to x}^{(m)}]$ gets smaller. To see this, f is more smooth when the Lipschitz constant m is smaller. This implies that $l_{x' \to x}^{(m)} = [l_{x'} - m ||x - x'||$ is larger and $u_{x' \to x}^{(m)} = u_{x'} + m ||x - x'||$ is smaller. As a result, we have a smaller interval $[l_{x' \to x}^{(m)}, u_{x' \to x}^{(m)}]$ for each x' and which implies a smaller $[l_{\mathcal{D} \to x}^{(m)}, u_{\mathcal{D} \to x}^{(m)}]$. This phenomenon can also be interpreted as implicitly "denoising" the original intervals by leveraging the smoothness of the hypothesis class.

255 Next, we extend Proposition 3.2 to $\widetilde{\mathcal{F}}_{\eta}$. The technical challenge is that for $f \in \widetilde{\mathcal{F}}_{\eta}$, f(x) may 256 lie outside the interval so we can't simply use $[l_{x'}, u_{x'}]$ as lower and upper bounds of f(x'). This 257 complicates the application of the Lipschitz property because f(x') can now be arbitrarily large or 258 small for any x', as long as the expected projection loss is smaller than η . The following result uses 259 a new notion of a bound gap of f(x) induced by x' which is the difference between the lower and 260 upper bounds induced by a given x' and the best lower and upper bounds from all x' (Figure 2b). 261 Formally, the lower bound gap for f(x) induced by x' is defined as $lg_{x'\to x}^{(m)} = l_{\mathcal{D}\to x} - l_{x'\to x}^{(m)}$ while 262 the upper bound gap is defined as $ug_{x' \to x}^{(m)} = u_{x' \to x}^{(m)} - u_{\mathcal{D} \to x}$. 263

Theorem 3.3. Let \mathcal{F} be a class of functions that are *m*-Lipschitz, and $\ell(y, y') = |y - y'|^p$ for any $p \ge 1$. For any $f \in \widetilde{\mathcal{F}}_\eta$ and for each x with p(x) > 0 we have,

 $f(x) \in [l_{\mathcal{D} \to x}^{(m)} - r_{\eta}(x), u_{\mathcal{D} \to x}^{(m)} + s_{\eta}(x)], \text{ where,}$ $\tag{8}$

$$r_{\eta}(x) = r \quad s.t. \quad \mathbb{E}_{X}[(r - lg_{X \to x}^{(m)})_{+}^{p}] = \eta,$$
(9)

$$s_{\eta}(x) = s \quad s.t. \quad \mathbb{E}_X[(s - ug_{X \to x}^{(m)})_+^p] = \eta.$$
 (10)

Proof. (Sketch) Our idea is based on the smoothness property of f. We can show that whenever f(x) is far from the interval of $[l_{D\to x}^{(m)}, u_{D\to x}^{(m)}]$, f(x') is also far away from $[l_{x'}, u_{x'}]$. However, this cannot happen frequently because the expected projection loss is smaller than η .

Interpretation. We compensate for $f \in \tilde{\mathcal{F}}_{\eta}$ by adding a buffer of size r and s to the interval derived in Proposition 3.2. If the average lower and upper bound gap is large, then we would have a larger compensation r, s. When $\eta = 0$, we have r = s = 0. In general, we can bound r, s in terms of $\eta^{1/p}$. **Proposition 3.4.** Under the conditions of Theorem 3.3, we can bound $r_{\eta}(x)$ and $s_{\eta}(x)$, as

$$r_{\eta}(x) \leq \inf_{\delta} \delta + (\eta / \Pr(lg_{X \to x}^{(m)} \leq \delta))^{1/p} \quad and \quad s_{\eta}(x) \leq \inf_{\delta} \delta + (\eta / \Pr(ug_{X \to x}^{(m)} \leq \delta))^{1/p}.$$
(11)

3.3 GENERALIZATION BOUNDS

274

275

276

277

278 279

281

282

284

285

286

291 292

293

294 295 296

297

298

299

300

301

302 303

304

305 306

307

321

In the previous section, we showed that the applicable interval of f(x) is smaller than the original interval $[l_x, u_x]$ when the hypothesis class \mathcal{F} is smooth. We can leverage this property to provide a generalization bound for any hypothesis $f \in \widetilde{\mathcal{F}}_{\eta}$. We denote the reduced interval from Theorem 3.3 as $I_{\eta}(x) := [l_{\mathcal{D}\to x}^{(m)} - r_{\eta}(x), u_{\mathcal{D}\to x}^{(m)} + s_{\eta}(x)]$. We will use the property that for intervals $I_1 = [l_1, u_1], I_2 = [l_2, u_2]$, and for any $y_1 \in I_1, y_2 \in I_2$, and any ℓ that satisfies Assumption 1,

$$\ell(y_1, y_2) \le \max(\ell(l_1, u_2), \ell(u_1, l_2)) =: d(\ell, I_1, I_2).$$
(12)

3.3.1 REALIZABLE SETTING

Theorem 3.5 (Error bound, Realizable setting). Let \mathcal{F} be a class of functions that are *m*-Lipschitz, assume that $f^* \in \widetilde{\mathcal{F}}_0$, then for any $f \in \widetilde{\mathcal{F}}_\eta$,

$$\operatorname{err}(f) \le \mathbb{E}[d(\ell, I_0(X), I_n(X))].$$
(13)

Notably, when we minimize the projection objective, f belongs to \mathcal{F}_{η} with a small η , where $\eta = O(1/\sqrt{n})$ (Section 3.1). While this bound is straightforward, we remark that it can be tight for certain hypothesis classes. For example, consider the case where \mathcal{F} consists of constant hypotheses and let $n \to \infty$. In this scenario, we have $r_{\eta}(x) \to r_0(x) = 0$ and $I_{\eta}(x) \to I_0(x)$. For each x, the error bound is given by

$$d(\ell, I_0(x), I_0(x)) = \ell(l_{\mathcal{D} \to x}^{(m)}, u_{\mathcal{D} \to x}^{(m)}) = \ell(\sup_{x'} l_{x'}, \inf_{x'} u_{x'}),$$
(14)

representing the loss between the boundaries of the intersected intervals. It is tight since the inequality holds when f^* and f each take values at the respective boundaries of the intersected interval.

3.3.2 AGNOSTIC SETTING

Now, we study the agnostic setting, where we do not assume the existence of such f^* in \mathcal{F} . Instead, we focus on comparing with $f_{\text{OPT}} = \arg \min_{f \in \mathcal{F}} \operatorname{err}(f)$, the hypothesis in \mathcal{F} with the smallest expected error. First, we show that, in contrast to the realizable setting, simply learning a hypothesis $f \in \mathcal{F}$ that always lies within the interval by minimizing the projection loss may not converge to f_{OPT} . This is because a smaller projection loss π does not imply a smaller standard loss ℓ .

Proposition 3.6. Let ℓ be an ℓ_p loss, for any hypothesis f_1, f_2 , there exists a distribution \mathcal{D}_I and \mathcal{D} such that $\mathbb{E}_{\mathcal{D}_I}[\pi_{\ell}(f_1(X), L, U)] < \mathbb{E}_{\mathcal{D}_I}[\pi_{\ell}(f_2(X), L, U)]$ but $\operatorname{err}(f_1) > \operatorname{err}(f_2)$.

While minimizing the projection loss, we might overlook a hypothesis that has a smaller standard loss but a higher projection loss. However, we remark that the projection loss is still useful since it is a lower bound of the standard loss.

Proposition 3.7. Let
$$\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$$
 be a loss function that satisfies Assumption 1, then for any f ,
 $\mathbb{E}[\pi_{\ell}(f(X), L, U)] \leq \operatorname{err}(f).$ (15)

Consequently, if we let $OPT = err(f_{OPT})$, we must have $f_{OPT} \in \mathcal{F}_{OPT}$ since the projection loss is upper bound by the standard loss. This means we can apply Theorem 3.3 for f_{OPT} and consequently achieve an error bound similar to what we obtained in the realizable setting. **Theorem 3.8** (Error bound, Agnostic setting). Let \mathcal{F} be a class of functions that are *m*-Lipschitz, and suppose ℓ satisfies Assumption 1 and the triangle inequality, then for any $f \in \widetilde{\mathcal{F}}_{\eta}$, we have

$$\operatorname{err}(f) \le \operatorname{OPT} + \mathbb{E}[d(\ell, I_{\eta}(X), I_{\operatorname{OPT}}(X))].$$
(16)

While it's not ideal to minimize the projection loss in the agnostic setting since we may not converge to f_{OPT} , our bound suggests that the expected error of f would not be much larger than that of f_{OPT} . This error bound becomes smaller when the intervals $I_{\eta}(x)$, $I_{\text{OPT}}(x)$ are small. Overall, our theoretical insight suggests that we can improve our error bound by (i) having a smoother hypothesis class (smaller m) (ii) increasing the number of data points n (which leads to smaller η), since both results in smaller intervals $I_{\eta}(x)$. However, if m is too small, \mathcal{F} may not contain a good hypothesis, causing OPT to be large.

4 LEARNING FROM INTERVALS USING A MINMAX OBJECTIVE

In Section 2, our goal was to learn a function f that ideally lies within the given interval $(f \in \tilde{\mathcal{F}}_0)$, using an objective that penalizes values away from the given interval. In this section, we explore a different strategy: we aim to learn a function $f \in \mathcal{F}$ that minimizes the maximum loss with respect to the worst-case \tilde{y} within the interval. We demonstrate that this approach yields a closed-form solution that can be evaluated efficiently. First, we define the worst-case loss as

$$\rho_{\ell}(f(x), l, u) := \max_{\tilde{y} \in [l, u]} \ell(f(x), \tilde{y}).$$
(17)

Proposition 4.1. Let ℓ be a loss function that satisfies Assumption 1, then

$$\rho_{\ell}(f(x), l, u) = \mathbb{1}[f(x) \le \frac{l+u}{2}]\ell(f(x), u) + \mathbb{1}[f(x) > \frac{l+u}{2}]\ell(f(x), l).$$
(18)

The proof is provided in Appendix D.5. Since $y \in [l, u]$, this objective serves as an upper bound for the true loss: $\rho_{\ell}(f(x), l, u) \ge \ell(f(x), y)$. Consequently, if we have a hypothesis with a small expected value $\mathbb{E}[\rho_{\ell}(f(x), l, u)]$, then the error $\operatorname{err}(f)$ will also be small. Based on Proposition 4.1, we define the **Minmax** objective as

354 355 356

357

360

361

362

364 365

374

375

324

325

326 327 328

330

331

332

333

334

335 336

337 338

339

340

341

342

343 344

345 346

351

352

353

$$\min_{f} \sum_{i=1}^{n} \mathbb{1}[f(x_i) \le \frac{l_i + u_i}{2}] \ell(f(x_i), u_i) + \mathbb{1}[f(x_i) > \frac{l_i + u_i}{2}] \ell(f(x_i), l_i).$$

$$f = \frac{1}{i=1}$$
 $Z = Z$

In particular, when $\ell(y, y') = |y - y'|$, we can show that minimizing ρ is equivalent to performing supervised learning using the mid-point of each interval.

Corollary 4.2. Let $\ell(y, y') = |y - y'|$ then $\rho_{\ell}(f(x), l, u) = |f(x) - \frac{l+u}{2}| + \frac{u-l}{2}$ and the solution of equation 19 is equivalent to

$$f' = \arg\min_{f \in \mathcal{F}} \sum_{i=1}^{n} |f(x_i) - \frac{l_i + u_i}{2}|.$$
 (20)

(19)

366 The proof of this corollary is provided in Appendix D.6. This corollary establishes a connection 367 between the heuristic of using the midpoint as a target and our approach of minimizing the maximum 368 loss ρ . However, we note that ρ does not take the smoothness of the hypothesis class \mathcal{F} into account and may lead to the worst-case labels that are overly conservative and not reflective of the target 369 labels. While we aim to minimize the loss with respect to the worst-case labels, we also want them 370 to be realistic. Therefore, it would be beneficial to incorporate knowledge about certain properties 371 of the true labels. In particular, in the realizable setting, we know that $f^* \in \mathcal{F}_0$, so we may consider 372 the worst-case labels that can be generated by some $f \in \mathcal{F}_0$, 373

$$\min_{f \in \mathcal{F}} \max_{f' \in \widetilde{\mathcal{F}}_0} \mathbb{E}[\ell(f(X), f'(X)].$$
(21)

In the realizable setting, this method also provides an upper bound for $\operatorname{err}(f)$, but it is weaker than ρ because we are comparing against the worst-case $f' \in \widetilde{\mathcal{F}}_0$ rather than any possible $\tilde{y} \in [l, u]$.

Proposition 4.3. In the realizable setting where $f^* \in \widetilde{\mathcal{F}}_0$, for a bounded loss ℓ , for any $f \in \mathcal{F}$, 379

$$\operatorname{err}(f) \le \max_{f' \in \widetilde{\mathcal{F}}_0} \mathbb{E}[\ell(f(X), f'(X))] \le \mathbb{E}[\rho_\ell(f(X), L, U)].$$
(22)

The proof is provided in Appendix D.7. With this inequality, we can conclude that when a hypothesis has a small minmax objective, its expected loss would be small as well. Moreover, we demonstrate that restricting the worst-case labels to those that could be generated by some $f \in \tilde{\mathcal{F}}_0$ can lead to better performance than using all possible worst-case labels. This is due to worst-case labels being highly sensitive to the interval size.

Proposition 4.4. For any constant c > 0 and $\ell(y, y') = |y - y'|$, there exists a distribution \mathcal{D}_I and a hypothesis class \mathcal{F} and $f^* \in \mathcal{F}$ such that for $f_1 = \arg \min_{f \in \mathcal{F}} \max_{f' \in \widetilde{\mathcal{F}}_0} \mathbb{E}[\ell(f(X), f'(X)] \text{ and } f_2 = \arg \min_{f \in \mathcal{F}} \mathbb{E}[\rho_\ell(f(X), L, U)], \operatorname{err}(f_1) = 0$ while $\operatorname{err}(f_2) > c$.

The proof is in Appendix D.8. An empirical Minmax objective using labels from $\widetilde{\mathcal{F}}_0$ is given by

$$\min_{f \in \mathcal{F}} \max_{f' \in \widetilde{\mathcal{F}}_0} \sum_{i=1}^n \ell(f(x_i), f'(x_i)).$$
(23)

However, there is no closed-form solution for the inner maximization of objective in 23, making it less efficient to optimize than equation 19. To address this, we propose alternative approaches by approximately learning $f' \in \tilde{\mathcal{F}}_0$ to solve this objective.

1) Regularization. We keep track of two hypothesis $f, f' \in \mathcal{F}$ and introduce a regularization term based on the projection loss to ensure that f' is close $\tilde{\mathcal{F}}_0$. We call this method Minmax (reg),

397

398

399

380

387

388 389

390 391

392 393 394

$$\min_{f \in \mathcal{F}} \max_{f' \in \mathcal{F}} \sum_{i=1}^{n} \ell(f(x_i), f'(x_i)) - \lambda \sum_{i=1}^{n} \pi(f'(x_i), l_i, u_i).$$
(24)

Here the regularization term is always non-positive and depends only on f'. We can use a gradient descent ascent (Korpelevich, 1976; Chen & Rockafellar, 1997; Lin et al., 2020) algorithm that updates f and f' with one gradient step at a time to solve this objective.

2) Pseudo labels. We could replace a hypothesis class $\widetilde{\mathcal{F}}_0$ with a finite set of hypotheses $\{f_1, f_2, \ldots, f_k\}$ where $f_j \in \widetilde{\mathcal{F}}_\eta$ for some small η . We can get f_j by minimizing the empirical projection loss. We then relax our objective by learning f that minimizes the maximum loss with respect to f_j . We call this method **PL** (**Max**),

$$\min_{f \in \mathcal{F}} \max_{j \in \{1, \dots, k\}} \sum_{i=1}^{n} \ell(f(x_i), f_j(x_i)).$$
(25)

Since f_j are fixed, learning f becomes a minimization problem, which is more stable to solve compared to the original minmax problem. Alternatively, to further stabilize the learning objective, we can replace the max over f_j with mean. We refer to this variant as **PL** (Mean),

$$\min_{f \in \mathcal{F}} \sum_{j=1}^{k} \sum_{i=1}^{n} \ell(f(x_i), f_j(x_i)).$$
(26)

423

412 413 414

415

416

417 418

5 EXPERIMENTS

Following prior work (Cheng et al., 2023a), we conducted experiments on five public datasets from the UCI Machine Learning Repository: Abalone, Airfoil, Concrete, Housing, and Power Plant. Since these datasets are originally regression tasks with single target values, we transformed them into datasets with interval targets (described shortly). Dataset statistics are provided in Section F. For the experimental setup, we used the same configuration as (Cheng et al., 2023a): the model architecture is a MLP with hidden layers of sizes 10, 20, and 30. We trained the models using the Adam optimizer with a learning rate of 0.001 and a batch size of 512 for 1000 epochs.

431 Interval Data Generation Methodology. We propose a general approach for generating interval data for each target value y. This method depends on two factors: the interval size $q \in [0, \infty]$ and

32		Projection (equation 5)	Minmax (equation 19)	Minmax (reg) (equation 24)	PL (max) (equation 25)	PL (mean) (equation 26)
33	Abalone	$1.56_{0.01}$	$1.65_{0.02}$	$1.54_{0.01}$	$1.52_{0.01}$	$1.52_{0.01}$
34	Airfoil	$2.46_{0.08}$	$2.65_{0.07}$	$3.41_{0.04}$	$3.31_{0.04}$	$2.42_{0.07}$
35	Concrete	$5.75_{0.13}$	$7.34_{0.2}$	$6.23_{0.16}$	$5.86_{0.48}$	$5.43_{0.12}$
36	Housing	$5.17_{0.13}$	$6.88_{0.31}$	$5.42_{0.15}$	$5.07_{0.09}$	$5.05_{0.09}$
37	Power-plant	$3.4_{0.03}$	$3.47_{0.02}$	$3.48_{0.03}$	$3.33_{0.01}$	$3.33_{0.01}$
38	Average (rank)	2.8	4.4	4.2	2.2	1

Table 1: Test Mean Absolute Error (MAE) and the standard error (over 10 random seeds) for the uniform interval setting. PL (mean) is the best-performing method in this setting.

the interval location $p \in [0, 1]$. The interval is then defined as [l, u] = [y - pq, y + (1 - p)q]. When p = 0, the target value y is at the lower boundary of the interval whereas p = 1 places y at the upper boundary. In this work, we consider q and p to be generated from uniform distributions over specified ranges. The prior interval generation method in Cheng et al. (2023a) could be seen as a special case of our approach when $q \sim \text{Uniform}[0, q_{\max}]$ and $p \sim \text{Uniform}[0, 1]$.

5.1 RESULTS

439

440

441

447 448

449

450 Which method works best in the uniform setting? We begin by evaluating methods in the 451 uniform interval setting described in prior work (Cheng et al., 2023a), where the interval size 452 $q \sim \text{Uniform}[0, q_{\text{max}}]$ and the location of the interval $p \sim \text{Uniform}[0, 1]$. For each dataset, we set q_{max} to be approximately equal to the range of the target values, $y_{\text{max}} - y_{\text{min}}$. Specifically, we 453 set $q_{\text{max}} = 30$ (Abalone), 30 (Airfoil), 90 (Concrete), 120 (Housing), and 90 (Power Plant). Our 454 findings indicate that the PL (mean) method performs best in this uniform setting, with PL (max) 455 and the projection method ranking second and third, respectively (Table 1). Given the superior per-456 formance of PL (mean), we conducted an ablation study to better understand its effectiveness. We 457 explored the impact of varying the number of hypotheses k and compared it with an ensemble base-458 line that combines pseudo-labels *before* using them to train the model, for which we still find that 459 PL (mean) still performs better (Appendix I). 460

What about other interval settings? We conducted more detailed experiments to investigate which 461 factors impact the performance of each method. Specifically, we varied the interval size q and the 462 interval location p by 1) varying q_{max} , 2) varying q_{min} , 3) varying p with three settings designed 463 to position the true value y at: i) only one boundary of the interval, ii) both boundaries of the 464 interval, iii) the middle of the interval. Full details are provided in Appendix G. We found that: (1) 465 All methods are quite robust to changes in the interval size, except for the Minmax method, whose 466 performance decreases significantly as the interval size increases. This is consistent with our insights 467 from the proof of 4.4), (2) The location of the true value y can have a large impact on performance; 468 specifically, the Minmax method performs better when y is close to the middle of the interval. One 469 explanation is that Minmax is equivalent to supervised learning with the midpoint of the interval (Corollary 4.2). Conversely, the other methods perform better when y is close to *both* boundaries of 470 the interval but not when y is close to *only* one boundary. Finally, we conclude that if we only know 471 that the interval size is large, it is better to use the PL (pseudo-labeling). However, if we know the 472 true value y is close to the middle of the interval, then the Minmax method is more preferable. 473

474 475

476

5.2 CONNECTION TO OUR THEORETICAL ANALYSIS

To validate our theoretical findings in practice, we conducted experiments designed to test whether 477 our theory holds under empirical conditions. Recall that our main result (Theorem 3.3) states that 478 if a hypothesis f approximately lies within the intervals $(f \in \mathcal{F}_n)$ and is smooth, then f will lie 479 within intervals smaller than the original ones. To control the smoothness of our hypothesis, we 480 utilize a Lipschitz MLP, which is an MLP augmented with spectral normalization layers (Miyato 481 et al., 2018). The normalization ensures that the Lipschitz constant of the MLP is less than 1. We 482 then scale the output of the MLP by a constant factor m to ensure that the Lipschitz constant of the 483 hypothesis is less than m. 484

Reduced interval size Our first experiment aims to determine whether the intervals, within which our hypothesis $f \in \tilde{\mathcal{F}}_0$ lies, are smaller than the original intervals. Recall that the original in-



Figure 3: Test MAE of the projection method with Lipschitz MLP using different values of the Lipschitz constant. The vertical line is the Lipschitz constant approximated from the training set. (Top) The dashed horizontal lines are the test MAE of PL (Mean) and Projection approach with a standard MLP. (Bottom) Approximated interval size $I_{\eta}(x)$ for Lipschitz MLP with a different value of Lipschitz constant *m*. The dashed horizontal lines are the values from standard (non-Lipschitz) MLP. The figures for all datasets are in Appendix H.

505 tervals are given by [l, u], and our theorem suggests that they would reduce to $I_{\eta}(x) = [\tilde{l}_{\mathcal{D}\to x}^{(m)} -$ 506 $r_{\eta}(x), \tilde{u}_{\mathcal{D} \to x}^{(m)} + s_{\eta}(x)$]. However, it is not possible to calculate $I_{\eta}(x)$ directly since it requires access 507 to every $f \in \widetilde{\mathcal{F}}_0$. Instead, we approximate $I_n(x)$ using samples of hypotheses from $\widetilde{\mathcal{F}}_0$ by pro-508 ceeding as follows: 1) We train 10 models with the projection objective, each from different random 509 initializations (denoted by f_1, \ldots, f_{10}), 2) For each x, we approximate the reduced interval using the 510 minimum and maximum values of the outputs from these models, given by $[\min_i f_i(x), \max_i f_i(x)]$. 511 We set $m \in \{0.1, 0.1 \times 2^1, \dots, 0.1 \times 2^{13}\}$ and consider a uniform interval setting with $q_{\text{max}} = 90$. 512 As expected, when the hypothesis becomes smoother, we observe that the average interval size 513 decreases (Figure 3 (Bottom)). Moreover, we found that even when the Lipschitz constant is much 514 larger than the value estimated from the data (vertical line), the average reduced interval size remains 515 significantly smaller than the original interval (which is 45 since $q_{\text{max}} = 90$). We also observe that 516 the average interval sizes from the standard MLPs are smaller than the original values. 517

Test performance In addition to examining the average interval size, we also plot the test Mean 518 Absolute Error (MAE) of the Lipschitz MLP with the projection objective, compared with the test 519 MAE of the standard MLP (Figure 3 (Top)). We found that, with the right level of smoothness, 520 Lipschitz MLP can achieve better performance than the standard MLP. When the Lipschitz constant 521 is very small, the performance is poor for all datasets. However, performance improves as the 522 Lipschitz constant increases. We observe that the optimal Lipschitz constant is always larger than 523 the Lipschitz constant estimated from the training set (vertical line). For some datasets, performance 524 degrades when the Lipschitz constant becomes too large. This aligns with our insight from Theorem 3.8, which suggests that we can improve the error bound by ensuring that the hypothesis class is as 525 smooth as possible (smaller m so that $I_n(x)$ is small) while still containing a good hypothesis (i.e., 526 low OPT). Nevertheless, we do not need to know the Lipschitz constant of the dataset and can treat 527 it as a tunable hyperparameter in practice. 528

530 6 CONCLUSION

529

531 In this paper, we studied the problem of learning from interval targets where only the lower and upper 532 bounds of the target are known. We analyzed two approaches: i) learning a hypothesis that always 533 lies within the interval, and ii) minimizing with respect to the worst-case label (or pseudolabels). Our 534 results showed how smoothness can be beneficial by i) leading to a smaller interval, and ii) having a "regularized" worst-case label. On the experimental side, our proposed minmax/pseudolabel approach achieves good performance (where PL (mean) is the best-performing method) and validates 536 our theoretical insights in practice. For future work, it would be interesting to study this problem for 537 a more challenging setup with less assumption such as the true target may not always be inside the 538 given interval or the training data are not independently and identically distributed, for instance the 539 time-series setting.

540 REFERENCES

546

552

563

565

542	Dana Angluin and Philip Laird.	Learning from noisy examples.	Machine learning, 2:343–370,
543	1988.		

- Pranjal Awasthi, Maria Florina Balcan, and Philip M Long. The power of localization for efficiently learning linear separators with noise. *Journal of the ACM (JACM)*, 63(6):1–27, 2017.
- Jessa Bekker and Jesse Davis. Learning from positive and unlabeled data: A survey. *Machine Learning*, 109(4):719–760, 2020.
- 549 David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A
 550 Raffel. Mixmatch: A holistic approach to semi-supervised learning. Advances in neural information processing systems, 32, 2019.
- Mathis Brosowsky, Florian Keck, Olaf Dünkel, and Marius Zöllner. Sample-specific output constraints for neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 6812–6821, 2021.
- 556 Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. Semi-supervised learning.
- George HG Chen and R Tyrrell Rockafellar. Convergence rates in forward–backward splitting.
 SIAM Journal on Optimization, 7(2):421–444, 1997.
- Pengfei Chen, Ben Ben Liao, Guangyong Chen, and Shengyu Zhang. Understanding and utilizing
 deep neural networks trained with noisy labels. In *International conference on machine learning*,
 pp. 1062–1070. PMLR, 2019.
 - Xin Cheng, Yuzhou Cao, Ximing Li, Bo An, and Lei Feng. Weakly supervised regression with interval targets. In *International Conference on Machine Learning*, pp. 5428–5448. PMLR, 2023a.
- Xin Cheng, Deng-Bao Wang, Lei Feng, Min-Ling Zhang, and Bo An. Partial-label regression. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 37, pp. 7140–7147, 2023b.
- Timothee Cour, Ben Sapp, and Ben Taskar. Learning from partial labels. *The Journal of Machine Learning Research*, 12:1501–1536, 2011.
- Ilias Diakonikolas, Themis Gouleakis, and Christos Tzamos. Distribution-independent pac learning
 of halfspaces with massart noise. *Advances in Neural Information Processing Systems*, 32, 2019.
- Marthinus C Du Plessis, Gang Niu, and Masashi Sugiyama. Analysis of learning from positive and unlabeled data. *Advances in neural information processing systems*, 27, 2014.
- Charles Elkan and Keith Noto. Learning classifiers from only positive and unlabeled data. In
 Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 213–220, 2008.
- Gabriel Erion, Joseph D Janizek, Pascal Sturmfels, Scott M Lundberg, and Su-In Lee. Improving
 performance of deep learning models with axiomatic attribution priors and expected gradients.
 Nature machine intelligence, 3(7):620–631, 2021.
- Lei Feng, Takuo Kaneko, Bo Han, Gang Niu, Bo An, and Masashi Sugiyama. Learning with multiple complementary labels. In *International conference on machine learning*, pp. 3072–3081. PMLR, 2020a.
- Lei Feng, Jiaqi Lv, Bo Han, Miao Xu, Gang Niu, Xin Geng, Bo An, and Masashi Sugiyama. Provably consistent partial-label learning. *Advances in neural information processing systems*, 33: 10948–10960, 2020b.
- Daniel Fu, Mayee Chen, Frederic Sala, Sarah Hooper, Kayvon Fatahalian, and Christopher Ré. Fast and three-rious: Speeding up weak supervision with triplet methods. In *International conference on machine learning*, pp. 3280–3291. PMLR, 2020.
- 593 Xiuwen Gong, Dong Yuan, and Wei Bao. Partial label learning via label influence function. In *International Conference on Machine Learning*, pp. 7665–7678. PMLR, 2022.

594 Cho-Jui Hsieh, Nagarajan Natarajan, and Inderjit Dhillon. Pu learning for matrix completion. In 595 International conference on machine learning, pp. 2445–2453. PMLR, 2015. 596 Takashi Ishida, Gang Niu, Weihua Hu, and Masashi Sugiyama. Learning from complementary 597 labels. Advances in neural information processing systems, 30, 2017. 598 Takashi Ishida, Gang Niu, Aditya Menon, and Masashi Sugiyama. Complementary-label learning 600 for arbitrary losses and models. In International conference on machine learning, pp. 2971–2980. 601 PMLR, 2019. 602 Rong Jin and Zoubin Ghahramani. Learning with multiple labels. Advances in neural information 603 processing systems, 15, 2002. 604 605 Giannis Karamanolakis, Subhabrata Mukherjee, Guoqing Zheng, and Ahmed Hassan. Self-training 606 with weak supervision. In Proceedings of the 2021 Conference of the North American Chapter 607 of the Association for Computational Linguistics: Human Language Technologies, pp. 845–863, 608 2021. 609 Davood Karimi, Haoran Dou, Simon K Warfield, and Ali Gholipour. Deep learning with noisy 610 labels: Exploring techniques and remedies in medical image analysis. *Medical image analysis*, 611 65:101759, 2020. 612 613 George Em Karniadakis, Ioannis G Kevrekidis, Lu Lu, Paris Perdikaris, Sifan Wang, and Liu Yang. Physics-informed machine learning. Nature Reviews Physics, 3(6):422–440, 2021. 614 615 Karthik Kashinath, M Mustafa, Adrian Albert, JL Wu, C Jiang, Soheil Esmaeilzadeh, Kamyar Aziz-616 zadenesheli, R Wang, Ashesh Chattopadhyay, A Singh, et al. Physics-informed machine learning: 617 case studies for weather and climate modelling. Philosophical Transactions of the Royal Society 618 A, 379(2194):20200093, 2021. 619 Durk P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-supervised 620 learning with deep generative models. Advances in neural information processing systems, 27, 621 2014. 622 623 Ryuichi Kiryo, Gang Niu, Marthinus C Du Plessis, and Masashi Sugiyama. Positive-unlabeled 624 learning with non-negative risk estimator. Advances in neural information processing systems, 625 30, 2017. 626 Galina M Korpelevich. The extragradient method for finding saddle points and other problems. 627 Matecon, 12:747–756, 1976. 628 629 Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. arXiv preprint 630 arXiv:1610.02242, 2016. 631 Xiaoli Li and Bing Liu. Learning to classify texts using positive and unlabeled data. In IJCAI, 632 volume 3, pp. 587–592. Citeseer, 2003. 633 634 Yuncheng Li, Jianchao Yang, Yale Song, Liangliang Cao, Jiebo Luo, and Li-Jia Li. Learning from 635 noisy labels with distillation. In Proceedings of the IEEE international conference on computer vision, pp. 1910–1918, 2017. 636 637 Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave mini-638 max problems. In International Conference on Machine Learning, pp. 6083–6093. PMLR, 2020. 639 Liping Liu and Thomas Dietterich. Learnability of the superset label learning problem. In Interna-640 tional conference on machine learning, pp. 1629–1637. PMLR, 2014. 641 642 Philip M Long and Rocco A Servedio. Random classification noise defeats all convex potential 643 boosters. In Proceedings of the 25th international conference on Machine learning, pp. 608–615, 644 2008. 645 Jiaqi Lv, Miao Xu, Lei Feng, Gang Niu, Xin Geng, and Masashi Sugiyama. Progressive identifica-646 tion of true labels for partial-label learning. In *international conference on machine learning*, pp. 647

6500-6510. PMLR, 2020.

- Tengyu Ma. Lecture notes from machine learning theory, 2022. URL http://web.stanford.
 edu/class/stats214/.
- Ron Meir and Tong Zhang. Generalization error bounds for bayesian mixture algorithms. *Journal* of Machine Learning Research, 4(Oct):839–860, 2003.
- Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization
 for generative adversarial networks. In *International Conference on Learning Representations*, 2018.
- 657 Mehryar Mohri. Foundations of machine learning, 2018.

666

- ⁶⁵⁸ Nagarajan Natarajan, Inderjit S Dhillon, Pradeep K Ravikumar, and Ambuj Tewari. Learning with noisy labels. *Advances in neural information processing systems*, 26, 2013.
- Rattana Pukdee, Dylan Sam, J Zico Kolter, Maria-Florina Balcan, and Pradeep Ravikumar. Learning
 with explanation constraints. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, pp. 49883–49926, 2023a.
- Rattana Pukdee, Dylan Sam, Pradeep Kumar Ravikumar, and Nina Balcan. Label propagation with
 weak supervision. In *The Eleventh International Conference on Learning Representations*, 2023b.
- Alexander Ratner, Stephen H Bach, Henry Ehrenberg, Jason Fries, Sen Wu, and Christopher Ré.
 Snorkel: Rapid training data creation with weak supervision. In *Proceedings of the VLDB en- dowment. International conference on very large data bases*, volume 11, pp. 269. NIH Public
 Access, 2017.
- Alexander J Ratner, Christopher M De Sa, Sen Wu, Daniel Selsam, and Christopher Ré. Data
 programming: Creating large training sets, quickly. *Advances in neural information processing systems*, 29, 2016.
- Laura Rieger, Chandan Singh, William Murdoch, and Bin Yu. Interpretations are useful: penalizing explanations to align neural networks with prior knowledge. In *International conference on machine learning*, pp. 8116–8126. PMLR, 2020.
- Andrew Slavin Ross, Michael C Hughes, and Finale Doshi-Velez. Right for the right reasons:
 training differentiable models by constraining their explanations. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pp. 2662–2670, 2017.
- Salva Rühling Cachay, Benedikt Boecking, and Artur Dubrawski. End-to-end weak supervision.
 Advances in Neural Information Processing Systems, 34:1845–1857, 2021.
- Changho Shin, Winfred Li, Harit Vishwakarma, Nicholas Carl Roberts, and Frederic Sala. Universalizing weak supervision. In *The Tenth International Conference on Learning Representations*, 2022.
- Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel,
 Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fixmatch: Simplifying semi-supervised
 learning with consistency and confidence. *Advances in neural information processing systems*,
 33:596–608, 2020.
- Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Learning from noisy labels with deep neural networks: A survey. *IEEE transactions on neural networks and learning systems*, 34(11):8135–8153, 2022.
- Renee Swischuk, Laura Mainini, Benjamin Peherstorfer, and Karen Willcox. Projection-based
 model reduction: Formulations for physics-based machine learning. *Computers & Fluids*, 179:
 704–717, 2019.
- Ruth Urner and Shai Ben-David. Probabilistic lipschitzness a niceness assumption for deterministic labels. In *Learning Faster from Easy Data-Workshop@ NIPS*, volume 2, pp. 1, 2013.
- Jesper E Van Engelen and Holger H Hoos. A survey on semi-supervised learning. *Machine learning*, 109(2):373–440, 2020.

702 703 704	Qian-Wei Wang, Yu-Feng Li, and Zhi-Hua Zhou. Partial label learning with unlabeled data. In <i>Proceedings of the 28th International Joint Conference on Artificial Intelligence</i> , pp. 3755–3761, 2019.
705 706 707 708	Hongwei Wen, Jingyi Cui, Hanyuan Hang, Jiabin Liu, Yisen Wang, and Zhouchen Lin. Leveraged weighted loss for partial label learning. In <i>International conference on machine learning</i> , pp. 11091–11100. PMLR, 2021.
709 710 711	Jared Willard, Xiaowei Jia, Shaoming Xu, Michael Steinbach, and Vipin Kumar. Integrating physics-based modeling with machine learning: A survey. <i>arXiv preprint arXiv:2003.04919</i> , 1(1):1–34, 2020.
712 713 714 715	Dong-Dong Wu, Deng-Bao Wang, and Min-Ling Zhang. Revisiting consistency regularization for deep partial label learning. In <i>International conference on machine learning</i> , pp. 24212–24225. PMLR, 2022.
716 717 718	Jin-Long Wu, Heng Xiao, and Eric Paterson. Physics-informed machine learning approach for augmenting turbulence models: A comprehensive framework. <i>Physical Review Fluids</i> , 3(7): 074602, 2018.
719 720 721	Ning Xu, Congyu Qiao, Xin Geng, and Min-Ling Zhang. Instance-dependent partial label learning. Advances in Neural Information Processing Systems, 34:27119–27130, 2021.
722 723 724	Wanqian Yang, Lars Lorch, Moritz Graule, Himabindu Lakkaraju, and Finale Doshi-Velez. In- corporating interpretable output constraints in bayesian neural networks. <i>Advances in Neural</i> <i>Information Processing Systems</i> , 33:12721–12731, 2020.
725 726 727	Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning with graph embeddings. In <i>International conference on machine learning</i> , pp. 40–48. PMLR, 2016.
728 729 730	Fei Yu and Min-Ling Zhang. Maximum margin partial label learning. In Asian conference on machine learning, pp. 96–111. PMLR, 2016.
731 732	Xiyu Yu, Tongliang Liu, Mingming Gong, and Dacheng Tao. Learning with biased complementary labels. In <i>Proceedings of the European conference on computer vision (ECCV)</i> , pp. 68–83, 2018.
733 734 735	Xiaohua Zhai, Avital Oliver, Alexander Kolesnikov, and Lucas Beyer. S41: Self-supervised semi- supervised learning. In <i>Proceedings of the IEEE/CVF international conference on computer vi-</i> <i>sion</i> , pp. 1476–1485, 2019.
736 737 738	Jieyu Zhang, Yue Yu, Yinghao Li, Yujing Wang, Yaming Yang, Mao Yang, and Alexander Ratner. Wrench: A comprehensive benchmark for weak supervision. In <i>Thirty-fifth Conference on Neural</i> <i>Information Processing Systems Datasets and Benchmarks Track (Round 2).</i>
739 740 741	Jieyu Zhang, Cheng-Yu Hsieh, Yue Yu, Chao Zhang, and Alexander Ratner. A survey on program- matic weak supervision. <i>arXiv preprint arXiv:2202.05433</i> , 2022.
742 743 744	Min-Ling Zhang, Bin-Bin Zhou, and Xu-Ying Liu. Partial label learning via feature-aware disambiguation. In <i>Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining</i> , pp. 1335–1344, 2016.
745 746 747	Min-Ling Zhang, Fei Yu, and Cai-Zhi Tang. Disambiguation-free partial label learning. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 29(10):2155–2167, 2017.
748 749	Xiaojin Zhu and Andrew B Goldberg. Introduction to semi-supervised learning. Springer Nature, 2022.
751 752 753 754	Xiaojin Jerry Zhu. Semi-supervised learning literature survey. 2005.

758 759

760

761

762

763

764

765

766

767

768

769

770 771

772

773 774

775

776 777

778

787

788 789

A GENERALIZATION BOUND BASED ON THE SAMPLE SIZE n

The generalization bounds of Theorem 3.5 and Theorem 3.8 are non-asymptotic. The error bound is applicable for any $f \in \tilde{\mathcal{F}}_{\eta}$ where the error depends on the value η . To improve the clarity of how generalization depends on the number of training sample n, we provide an explicit samplecomplexity generalization bound for hypothesis classes whose the Rademacher complexity decay as $O(1/\sqrt{n})$. This includes a class of linear models or a class of two-layer neural networks with a bounded weight (Ma, 2022). To simplify the Theorem, we will only present the statement and the proof for the case of L_1 loss. However, an extension for a general L_p loss is straightforward where we can replace the triangle inequality with the Minkowski's inequality.

Theorem A.1 (Generalization bound, Realizable Setting). *Take the conditions of Theorem 3.5 (realizability and m-Lipschitzness) and further assume*

- \mathcal{F} is hypothesis class whose the Rademacher complexity decays as $O(1/\sqrt{n})$ e.g. two-layer neural networks with bounded weights,
- the support of the distribution \mathcal{D}_I is a bounded set,
- the loss function is $\ell(y, y') = |y y'|$.

Then, with probability at least $1 - \delta$, for any f that minimize the empirical projection objective, for any $\tau > 0$,

$$\operatorname{err}(f) \leq \underbrace{\mathbb{E}_{X}[|u_{\mathcal{D}\to X}^{(m)} - l_{\mathcal{D}\to X}^{(m)}|]}_{(a)} + \underbrace{\tau + \left(\frac{D}{\sqrt{n}} + M\sqrt{\frac{\ln(1/\delta)}{n}}\right)\Gamma(\tau)}_{(b)},\tag{27}$$

784 785 where D, M are constants and $\Gamma(\tau) = \mathbb{E}_{\widetilde{X}} \left[1/\min(\Pr_X(lg_{X \to \widetilde{X}}^{(m)} \le \tau), \Pr_X(ug_{X \to \widetilde{X}}^{(m)} \le \tau)) \right]$ is a decreasing function of τ .

Interpretation: Our error bound is divided into two parts.

- Term (a) The first term represent an error term which depends on the smoothness property of our function class \mathcal{F} and the quality of the given intervals. The first error term is **irreducible** with the number of samples in the sense that it does not decrease as we have more training samples and this error term depends solely on the quality of the intervals and the smoothness of our hypothesis class. However, this term can be small. For example, in the case when the ambiguity degree is small, this error term would be zero, ensure a perfect recovery of the true labels.
- 796 Term (b) The second and the third term capture how well we can learn a hypothesis that belongs to 797 the intervals and these would decay as we have a larger sample size n. To see this, assume 798 that we have a fix value of τ , if one set $n \to \infty$ then the third term would converge to 799 zero. That is, (b) would converge to τ as $n \to \infty$. Since τ is arbitrary, we can set τ to be 800 small so that (b) would decay to zero as $n \to \infty$ and we are left with the first term (a). In 801 addition, the function $\Gamma(\tau)$ depends on the distribution of intervals \mathcal{D}_I . In particular, when \mathcal{D}_I has small lower/upper bound gaps, $\Gamma(\tau)$ would also be small which leads to a better 802 generalization bound for any fixed n. 803
- 804

805

806

807 *Proof.* Our result here can be derived from combining the results from Theorem 3.5, Proposition 808 3.4 and relate the Rademacher complexity of $\Pi(\mathcal{F})$ with \mathcal{F} . The proof here is divided into 2 steps; 809 i) Apply the Proposition 3.4 to the Theorem 3.5 to derive the bound in terms of η , ii) Bound η in terms of the sample size n.

Step 1: Derive the bound in term of η **.** Recall that from Theorem 3.5, we have

$$\operatorname{err}(f) \le \mathbb{E}[d(\ell, I_0(X), I_\eta(X))].$$
(28)

when $I_{\eta}(x) = [l_{\mathcal{D}\to x}^{(m)} - r_{\eta}(x), u_{\mathcal{D}\to x}^{(m)} + s_{\eta}(x)]$. Since we have an ℓ_1 loss, we have

$$d(\ell, I_0(x), I_\eta(x)) = |u_{\mathcal{D} \to x}^{(m)} - l_{\mathcal{D} \to x}^{(m)} + \max(r_\eta(x), s_\eta(x))|.$$
⁽²⁹⁾

Substitute this back in, we have an error bound

$$\operatorname{err}(f) \leq \mathbb{E}[|u_{\mathcal{D}\to X}^{(m)} - l_{\mathcal{D}\to X}^{(m)} + \max(r_{\eta}(X), s_{\eta}(X))|]$$
(30)

$$\leq \mathbb{E}[|u_{\mathcal{D}\to X}^{(m)} - l_{\mathcal{D}\to X}^{(m)}|] + \mathbb{E}[|\max(r_{\eta}(X), s_{\eta}(X))|] \quad \text{(triangle inequality)}. \tag{31}$$

Now, our goal is to bound the term $\mathbb{E}[|\max(r_n(X), s_n(X))|]$. From Proposition 3.4, we know that

$$r_{\eta}(x) \leq \inf_{\tau} \tau + (\eta / \Pr(lg_{X \to x}^{(m)} \leq \tau)) \quad \text{and} \quad s_{\eta}(x) \leq \inf_{\tau} \tau + (\eta / \Pr(ug_{X \to x}^{(m)} \leq \tau)).$$
(32)

We place δ with τ in the original statement because we will use δ as something else, later. This implies that

$$\max(r_{\eta}(x), s_{\eta}(x)) \leq \inf_{\tau} \tau + \left(\frac{\eta}{\min(\Pr(lg_{X \to x}^{(m)} \leq \tau), \Pr(ug_{X \to x}^{(m)} \leq \tau))}\right).$$
(33)

We define $\Lambda(\mathcal{D}, \tau) = \min(\Pr(lg_{X \to x}^{(m)} \leq \tau), \Pr(ug_{X \to x}^{(m)} \leq \tau))^{-1}$ so that $\max(r_{\eta}(x), s_{\eta}(x)) \leq \inf \tau + \eta \Lambda(\mathcal{D}, \tau).$ (34)

We can see that when $\Lambda(\mathcal{D},\tau) \ge 0$ and $\Lambda(\mathcal{D},\tau)$ is a decreasing function in τ . Substitue this back to the equation 31, for any $\tau > 0$, we would have

$$\operatorname{err}(f) \leq \mathbb{E}[|u_{\mathcal{D}\to X}^{(m)} - l_{\mathcal{D}\to X}^{(m)}|] + \mathbb{E}[|\tau + \eta\Lambda(\mathcal{D},\tau)|]$$
(35)

$$\leq \mathbb{E}[|u_{\mathcal{D}\to X}^{(m)} - l_{\mathcal{D}\to X}^{(m)}|] + \tau + \eta \mathbb{E}[\Lambda(\mathcal{D}, \tau)]$$
(36)

$$= \mathbb{E}[|u_{\mathcal{D}\to X}^{(m)} - l_{\mathcal{D}\to X}^{(m)}|] + \tau + \eta \Gamma(\mathcal{D}, \tau)$$
(37)

where we define $\Gamma(\mathcal{D},\tau) = \mathbb{E}[\Lambda(\mathcal{D},\tau)]$. We can see that every term in the equation above is independent of η , apart from the term η itself. This provide a more explicit error bound in term of η . Now, we will bound η in terms of the number of sample n.

Step 2: Bounding η in terms of the number of sample. Recall the result from equation 6, with probability at least $1 - \delta$ over the draws $(x_i, l_i, u_i) \sim \mathcal{D}_I$, for all $f \in \mathcal{F}$,

$$\mathbb{E}[\pi_{\ell}(f(X), L, U)] \le \frac{1}{n} \sum_{i=1}^{n} \pi_{\ell}(f(x_i), l_i, u_i) + 2R_n(\Pi(\mathcal{F})) + M\sqrt{\frac{\ln(1/\delta)}{n}}.$$
 (38)

Here, $R_n(\Pi(\mathcal{F}))$ is the Rademacher complexity of the function class $\Pi(\mathcal{F}) := \{\pi_\ell(f(x), l, u) \mapsto$ $\mathbb{R} \mid f \in \mathcal{F}$ and we assume that the π_{ℓ} is uniformly bounded by M. We recall that we learn \hat{f} by minimizing the empirical projection loss

$$\hat{f} = \arg\min_{f \in \mathcal{F}} \sum_{i=1}^{n} \pi_{\ell}(f(x_i), l_i, u_i).$$
(39)

Under the realizable setting, this objective would be zero since $f^* \in \mathcal{F}$ which implies that f^* has zero empirical projection $\sum_{i=1}^{n} \pi_{\ell}(f^*(x_i), l_i, u_i) = 0$ but \hat{f} also minimize the empirical projection loss so \hat{f} must also have a zero empirical projection loss. We write $\eta(f)$ to refer to the η value of f. Formally, defined as

$$\eta(f) = \mathbb{E}[\pi_{\ell}(f(X), L, U)].$$
(40)

(41)

Substituting \hat{f} to the bound above, we have

 $\eta(\hat{f}) \le 2R_n(\Pi(\mathcal{F})) + M\sqrt{\frac{\ln(1/\delta)}{n}}.$

The next step is to bound the Rademacher complexity $R_n(\Pi(\mathcal{F}))$ in terms of $R_n(\mathcal{F})$. We will do this by first showing that $\phi_i(f(x)) = \pi_\ell(f(x), l_i, u_i)$ is a Lipschitz continuous function and then reduce $R_n(\Pi(\mathcal{F}))$ to $R_n(\mathcal{F})$ with a variant of Talagrand's Lemma (Meir & Zhang, 2003). From our assumption that the support of \mathcal{D}_I is a bounded set, and our hypothesis class is a class of two-layer neural network with bounded weight, there exists a constant C for which, we have $|f(x)| \leq C$ almost surely. Here, we will show this property for L_p loss, recall that

$$\phi_i(f(x)) = \pi_\ell(f(x), l_i, u_i) \tag{42}$$

$$= (l_i - f(x))^p \mathbf{1}[f(x) < l_i] + (f(x) - u_i)^p \mathbf{1}[f(x) > u].$$
(43)

873 Differentiate with respect to f(x), we have

$$\nabla_{f(x)}\phi_i(f(x))| = p|(l_i - f(x))^{p-1}\mathbf{1}[f(x) < l_i] + (f(x) - u_i)^{p-1}\mathbf{1}[f(x) > u]|$$
(44)

$$\leq 2p(2C)^{p-1}.\tag{45}$$

Since this gradient is bounded for any f(x), we can conclude that $\phi_i(f(x))$ is *B*-Lipschitz for some constant *B*. Now, we unpack the definition of the Rademacher complexity,

$$R_n(\Pi(\mathcal{F})) = \mathbb{E}_{(x_i, l_i, u_i) \sim \mathcal{D}_I}[\mathbb{E}_{\sigma_i \sim \{-1, 1\}}[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \pi_\ell(f(x_i), l_i, u_i)\sigma_i]]$$
(46)

$$= \mathbb{E}_{(x_i, l_i, u_i) \sim \mathcal{D}_I} [\mathbb{E}_{\sigma_i \sim \{-1, 1\}} [\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \phi_i(f(x_i)) \sigma_i]].$$

$$(47)$$

We recall the following result from Meir & Zhang (2003) that when $\phi_1, \phi_2, \dots, \phi_n$ be functions where $\phi_i : \mathbb{R} \to \mathbb{R}$ are ϕ_i are L_i -Lipschitz, then

$$\mathbb{E}_{\sigma_i \sim \{-1,1\}} [\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \phi_i(f(x_i)) \sigma_i] \le \mathbb{E}_{\sigma_i \sim \{-1,1\}} [\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n L_i f(x_i) \sigma_i].$$
(48)

Applying this result with the fact that ϕ_i is *B*-Lipschitz for all i = 1, ..., n, we can conclude that

$$R_n(\Pi(\mathcal{F})) = \mathbb{E}_{(x_i, l_i, u_i) \sim \mathcal{D}_I}[\mathbb{E}_{\sigma_i \sim \{-1, 1\}}[\sup_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \phi_i(f(x_i))\sigma_i]]$$
(49)

$$\leq \mathbb{E}_{(x_i,l_i,u_i)\sim\mathcal{D}_I}[\mathbb{E}_{\sigma_i\sim\{-1,1\}}[\sup_{f\in\mathcal{F}}\frac{1}{n}\sum_{i=1}^n Bf(x_i)\sigma_i]]$$
(50)

$$=BR_n(\mathcal{F}).\tag{51}$$

We successfully reduce the Rademacher complexity of $\Pi(\mathcal{F})$ to \mathcal{F} . Since we assume that the Rademacher complexity of \mathcal{F} decays as $O(1/\sqrt{n})$, there exists a constant D such that

$$R_n(\Pi(\mathcal{F})) \le \frac{D}{\sqrt{n}} \tag{52}$$

and

$$\eta(\hat{f}) \le \frac{D}{\sqrt{n}} + M\sqrt{\frac{\ln(1/\delta)}{n}}$$
(53)

for some constant D, M. Substitute this back to the result from step 1 concludes our proof.

In the general setting where we have $\ell(y, y') = |y - y'|^p$, we would have the following error bound

$$\operatorname{err}(f) \le \left(\mathbb{E}_X[|u_{\mathcal{D}\to X}^{(m)} - l_{\mathcal{D}\to X}^{(m)}|^p]^{1/p} + \tau + \left(\frac{D}{\sqrt{n}} + M\sqrt{\frac{\ln(1/\delta)}{n}}\right)^{1/p} \Gamma(\tau)^{1/p} \right)^p$$
(54)

We can also write this generalization bound for the agnostic setting

Theorem A.2 (Generalization Bound, Agnostic Setting). Under the conditions of Theorem A.1 apart from realizability, with probability at least $1 - \delta$, for any f that minimize the empirical projection objective, for any $\tau > 0$,

$$\operatorname{err}(f) \leq \underbrace{\operatorname{OPT}}_{(a)} + \underbrace{\mathbb{E}_{X}[|u_{\mathcal{D} \to X}^{(m)} - l_{\mathcal{D} \to X}^{(m)}|]}_{(b)} + \underbrace{\tau + \left(\widehat{\operatorname{err}}(f) + \frac{D}{\sqrt{n}} + M\sqrt{\frac{\ln(1/\delta)}{n}} + \operatorname{OPT}\right)\Gamma(\tau),}_{(c)},$$
(55)

where D, M are constants and $\Gamma(\tau) = \mathbb{E}_{\widetilde{X}} \left[1/\min(\Pr_X(lg_{X \to \widetilde{X}}^{(m)} \leq \tau), \Pr_X(ug_{X \to \widetilde{X}}^{(m)} \leq \tau)) \right]$ is a decreasing function of τ , $\widehat{\operatorname{err}}(f)$ is an empirical projection error of f, and OPT is the expected error of the optimal hypothesis in \mathcal{F} .

Interpretation: In contrast to the realizability setting, our error bound for the agnostic setting is divided into three parts.

⁹³³ Term (a) The first term represent an error term of the optimal hypothesis in \mathcal{F} , given by OPT.

Term (b) The second term represent an error term which depends on the smoothness property of our function class \mathcal{F} and the quality of the given intervals similar to the realizability setting.

Term (c) The third and the fourth term capture how well we can learn a hypothesis that belongs to the intervals. The key difference between this agnostic setting and the realizability setting is that this term would not decay to zero anymore as $n \to \infty$. In particular, for a fixed τ , we can see that as $n \to \infty$, we would have $\widehat{\operatorname{err}}(f) \to \operatorname{OPT}$ since we are minimizing the empirical projection loss and as a result, this third part would converge to

$$\tau + 2 \operatorname{OPT} \Gamma(\tau). \tag{56}$$

Since this hold for any τ , the optimal τ would be the one such that $\tau = 2 \operatorname{OPT} \Gamma(\tau)$ and this value depends on the distribution \mathcal{D}_I .

Overall, when $n \to \infty$, the upper bound would converge to

$$OPT + \mathbb{E}_X[|u_{\mathcal{D}\to X}^{(m)} - l_{\mathcal{D}\to X}^{(m)}|] + \tau + 2 OPT \Gamma(\tau).$$
(57)

This can be small as long as the OPT is small, the expected lower/ upper bound gaps are small and when the noise in the given intervals are small. The proof of this theorem follows the same argument from the realizable setting.

B RELAXATION OF AMBIGUITY DEGREE FOR A REGRESSION SETTING

As noted in the related work section, the ambiguity degree is defined in the context of classification and it might not be suitable for regression tasks. This is due to the nature of the loss function, In classification, a hypothesis is either correct or incorrect, and a small ambiguity degree ensures that we can recover the true label. However, in regression, we are often satisfied with predictions that are sufficiently close to the target—for example, within an error tolerance of ϵ . This implies that we do not need to recover the exact true label, but a ball with a small radius around the true label might be sufficient.

In this section, we explore a relaxation of the original ambiguity degree to the regression setting. Motivated by the concept of a tolerable area around the true label y, we define an ambiguity radius

Definition B.1 (Ambiguity Radius). For distributions $\mathcal{D}, \mathcal{D}_I$ with a probability density function p, an ambiguity radius is defined as

$$\text{AmbiguityRadius}(\mathcal{D}, \mathcal{D}_I) := \min_{r \ge 0} r \quad \text{s.t.} \quad \Pr_{X, Y \sim \mathcal{D}}(\bigcap_{p(X, l, u) > 0} [l, u] \subseteq B(Y, r)) = 1$$
(58)

when $B(y,r) = \{y' \mid |y-y'| \le r\}$ is a ball of radius r around y.

The interpretation of this is that it is the smallest radius r for which we are guaranteed the intersection of all interval for a given x must lie within a radius of r from the true label y. As a direct consequence, we know that whenever the ambiguity degree is small the ambiguity radius must be zero since the intersection of all interval for a given x is just the true label $\{y\}$.

977 In fact, our analysis have captured the essence of this interval intersection for each x. We 978 recall that for any $f \in \tilde{\mathcal{F}}_0$ and for each x with p(x) > 0,

$$f(x) \in I_0(x) = [l_{\mathcal{D} \to x}^{(m)}, u_{\mathcal{D} \to x}^{(m)}] \subseteq B(y, r^*),$$
(59)

when r^* is the ambiguity radius. This follows directly from the definition of the ambiguity radius. As a result, we know that each interval $I_0(x)$ would have a size at most $2r^*$. The same technique as in the Section 3.3 would imply that the expected error of any $f \in \tilde{\mathcal{F}}_0$ would be at most $2r^*$ in the realizable setting (with L_1 loss).

Finally, we want to remark that our analysis not only is applicable to this extension of the ambiguity degree to the ambiguity radius, we further use the smooth property of \mathcal{F} and $I_0(x)$ might even be a proper subset of the ball $B(y, r^*)$, giving a result stronger than one based solely on the ambiguity radius.

990 991

976

979 980 981

C RELATED WORK

992 993

Weak supervision. Our setting is part of a sub-field of weak supervision where one learns from 994 noisy, limited, or imprecise sources of data rather than a large amount of labeled data. Learning 995 from noisy labels assumes that we only observe a noisy version of the true labels at the training time 996 where the noise follows different noise models (usually random noise) (Natarajan et al., 2013; Li 997 et al., 2017; Song et al., 2022; Angluin & Laird, 1988; Karimi et al., 2020; Awasthi et al., 2017; Chen 998 et al., 2019; Long & Servedio, 2008; Diakonikolas et al., 2019). Programmatic weak supervision, on 999 the other hand, assumes that we have access to multiple noisy weak labels (but deterministic noise) 1000 specified by domain experts, e.g. from logic rules or heuristics methods (Zhang et al., 2022; Zhang 1001 et al.; Ratner et al., 2016; 2017; Rühling Cachay et al., 2021; Shin et al., 2022; Karamanolakis et al., 1002 2021; Fu et al., 2020; Pukdee et al., 2023b). Positive-unlabeled learning is another type of weak supervision where the training set only contains positive examples and unlabeled examples (Kiryo 1003 et al., 2017; Du Plessis et al., 2014; Bekker & Davis, 2020; Elkan & Noto, 2008; Li & Liu, 2003; 1004 Hsieh et al., 2015). 1005

Learning with side information. In contrast to the weakly supervised setting, we have access to standard labeled data but also have access to some additional information. This could be unlabeled data which is studied in semi-supervised learning (Zhu, 2005; Chapelle et al.; Kingma et al., 2014; 1008 Van Engelen & Hoos, 2020; Berthelot et al., 2019; Zhu & Goldberg, 2022; Laine & Aila, 2016; 1009 Zhai et al., 2019; Sohn et al., 2020; Yang et al., 2016) or different constraints based on the domain 1010 knowledge such as physics rules (Willard et al., 2020; Swischuk et al., 2019; Karniadakis et al., 1011 2021; Wu et al., 2018; Kashinath et al., 2021) or explanations (Ross et al., 2017; Pukdee et al., 1012 2023a; Rieger et al., 2020; Erion et al., 2021) or output constraints (Yang et al., 2020; Brosowsky 1013 et al., 2021) which is similar to the interval targets. In some settings, interval targets are the best thing 1014 one could have (similar to the weak supervision setting) but in many cases such as in bond pricing, 1015 target intervals are readily available in the wild and could also be considered as a side information. 1016

1017

1018 D ADDITIONAL PROOFS

1019

1020 D.1 PROOF OF PROPOSITION 2.1

1022 *Proof.* First, we assume that $\pi_{\ell}(f(x), l, u) = 0$. This implies that there exists $\tilde{y} \in [l, u]$ such that $\ell(f(x), \tilde{y}) = 0$. From the assumption on ℓ that $\ell(y, y') = 0$ if and only if y = y', we must have $f(x) = \tilde{y} \in [l, u]$ as required. On the other hand, if $f(x) \in [l, u]$, it is clear that $\pi_{\ell}(f(x), l, u) = \ell(f(x), f(x)) = 0$ since $\ell(y, y') \ge 0$.

Now, assume that we can write $\ell(y, y') = \psi(|y - y'|)$ for some non-decreasing function ψ , we have

$$\pi_{\ell}(f(x), l, u) = \min_{\tilde{y} \in [l, u]} \psi(|f(x) - \tilde{y}|)$$
(60)

1029
1030
1031

$$y \in [l, u]$$

 $y \in [l, u]$
 $y \in [l, u]$
 $f(x) - \tilde{y}|)$
(61)

$$= \begin{cases} \psi(l - f(x)) & f(x) < l \\ \psi(0) & l \le f(x) \le u \end{cases}$$
(62)

1034
$$(\psi(f(x) - u) - f(x) > u)$$
1035
$$- 1[f(x) < l]\ell(f(x) - l) + 1[f(x) > u]\ell(f(x) - u)$$

$$= 1[f(x) < l]\ell(f(x), l) + 1[f(x) > u]\ell(f(x), u).$$
(63)

Here we rely on the assumption that ψ is non-decreasing so the minimum value of $\psi(x)$ happens when x is also at the minimum value.

1040 D.2 PROOF OF PROPOSITION 3.6

1042 Proof. Since $f_1 \neq f_2$, there exists x such that $f_1(x) \neq f_2(x)$. Without loss of generality, let 1043 $f_1(x) < f_2(x)$. Consider a simple one point distribution \mathcal{D} with only one data point (x, y) =1044 $(x, f_2(x) + \epsilon)$ with probability mass 1 and \mathcal{D}_I be another one point distribution with (x, l, u) =1045 $(x, f(x_1) - \epsilon, f(x_2) - \epsilon)$. We can see that $0 = \mathbb{E}_{\mathcal{D}_I}[\pi(f_1(X), L, U)] < \mathbb{E}_{\mathcal{D}_I}[\pi(f_2(X), L, U)] = \epsilon^p$ 1046 while $(f(x_2) - f(x_1) + \epsilon)^p = \operatorname{err}(f_1) > \operatorname{err}(f_2) = \epsilon^p$.

1048 D.3 PROOF OF PROPOSITION 3.7

Proof. From the Proposition 2.1,

$$\pi(f(x), l, u) = \mathbb{1}[f(x) < l]\ell(f(x), l) + \mathbb{1}[f(x) > u]\ell(f(x), u)$$
(64)

1053 Recall that $y \in [l, u]$, we consider 3 cases,

1.
$$f(x) < l, \pi(f(x), l, u) = \ell(f(x), l) = \psi(|l - f(x)|) \le \psi(|y - f(x)|) = \ell(f(x), y)$$

2. $f(x) > u, \pi(f(x), l, u) = \ell(f(x), u) = \psi(|f(x) - u|) \le \psi(|f(x) - y|) = \ell(f(x), y)$
3. $l \le f(x) \le u, \pi(f(x), l, u) = 0 \le \ell(f(x), y)$

1063 D.4 PROOF OF THEOREM 3.8

Proof. From the triangle inequality,

$$\ell(f(x), y) = \ell(f(x), f_{\text{OPT}}(x)) + \ell(f_{\text{OPT}}(x), y)$$
(65)

1068 We can take an expectation to have

$$\mathbb{E}[\ell(f(X), Y)] \le \mathbb{E}[\ell(f(X), f_{\mathsf{OPT}}(X)] + \mathsf{OPT}.$$
(66)

1071 Since $f_{\text{OPT}} \in \widetilde{\mathcal{F}}_{\text{OPT}}$ which from Theorem 3.3, we can bound

$$f_{\text{OPT}}(x) \in [l_{\mathcal{D}\to x}^{(m)} - r_{\text{OPT}}(x), l_{\mathcal{D}\to x}^{(m)} + s_{\text{OPT}}(x)].$$
 (67)

1075 Similarly, for any $f \in \widetilde{\mathcal{F}}_{\eta}$, we have

$$f(x) \in [l_{\mathcal{D} \to x}^{(m)} - r_{\eta}(x), u_{\mathcal{D} \to x}^{(m)} + s_{\eta}(x)]$$
(68)

1079 Finally, we can bound the error between any two intervals with the maximum loss between their boundaries. \Box

D.5 PROOF OF PROPOSITION 4.1

Proof. Since we can write $\ell(y, y') = \psi(|y - y'|)$ for some non-decreasing function ψ , we have

$$\rho_{\ell}(f(x), l, u) = \max_{\tilde{y} \in [l, u]} \psi(|f(x) - \tilde{y}|)$$
(69)

$$=\psi(\max_{\tilde{y}\in[l,u]}|f(x)-\tilde{y}|) \tag{70}$$

$$=\begin{cases} \psi(u - f(x)) & f(x) < \frac{l+u}{2} \\ \psi(f(x) - l) & f(x) \ge \frac{l+u}{2} \end{cases}$$
(71)

$$= 1[f(x) \le \frac{l+u}{2}]\ell(f(x), u) + 1[f(x) > \frac{l+u}{2}]\ell(f(x), l).$$
(72)

Here we rely on the assumption that ψ is non-decreasing so the maximum value of $\psi(x)$ happens when x is also at the maximum value.

D.6 PROOF OF COROLLARY 4.2

Proof. Since $\ell(y, y') = |y - y'|$, from Proposition 4.1, we have a closed form solution of ρ ,

$$\rho_{\ell}(f(x), l, u) = \mathbb{1}[f(x) \le \frac{l+u}{2}]\ell(f(x), u) + \mathbb{1}[f(x) > \frac{l+u}{2}]\ell(f(x), l)$$
(73)

$$1100 = 1[f(x) \le \frac{l+u}{2}](u-f(x)) + 1[f(x) > \frac{l+u}{2}](f(x)-l)$$
(74)

$$1102 \\ 1103 \\ 1104 \\ = 1[f(x) \le \frac{l+u}{2}](u - \frac{l+u}{2} + \frac{l+u}{2} - f(x)) + 1[f(x) > \frac{l+u}{2}](f(x) - \frac{l+u}{2} + \frac{l+u}{2} - l) \\ (75)$$

$$=\frac{u-l}{2}+1[f(x)\leq\frac{l+u}{2}](\frac{l+u}{2}-f(x))+1[f(x)>\frac{l+u}{2}](f(x)-\frac{l+u}{2})$$
(76)

$$\begin{aligned} & \begin{array}{c} 1107\\ 1108\\ 1109 \end{array} & = |f(x) - \frac{l+u}{2}| + \frac{u-l}{2}. \end{aligned} \tag{77} \end{aligned}$$

Since u_i, l_i are constants, $\frac{u_i - l_i}{2}$ would have no impact on the optimal solution of equation 19 and therefore, the optimal would also be the same as the one that minimizes $\sum_{i=1}^{n} |f(x_i) - \frac{l_i + u_i}{2}|$. \Box

D.7 PROOF OF PROPOSITION 4.3

Proof. From the realizability assumption, we know that $f^* \in \widetilde{\mathcal{F}}_0$, therefore,

$$\operatorname{err}(f) = \mathbb{E}[\ell(f(X), f^*(X))] \le \max_{f' \in \widetilde{\mathcal{F}}_0} \mathbb{E}[\ell(f(X), f'(X))].$$
(78)

On the other hand, Let $f'' \in \widetilde{\mathcal{F}}_0$, be a hypothesis that achieves the maximum value of $\mathbb{E}[\ell(f(X), f''(X))]$. Since $f'' \in \widetilde{\mathcal{F}}_0$ we know that

$$\mathbb{E}[\pi_{\ell}(f''(X), L, U)] = 0.$$
⁽⁷⁹⁾

Since the projection loss is always non-negative and is continuous, from Lemma E.1, we can con-clude that $\pi_{\ell}(f''(x), l, u) = 0$ for any x, l, u with positive density function p(x, l, u) > 0 which implies $f''(x) \in [l, u]$. Therefore, for any x with p(x) > 0,

$$\ell(f(x), f''(x)) \le \max_{\tilde{y} \in [l, u]} \ell(f(x), \tilde{y}) = \rho_{\ell}(f(x), l, u).$$
(80)

We can take an expectation over X, L, U and have the desired result.

D.8 PROOF OF PROPOSITION 4.4

Proof. Consider when $\mathcal{X} = \{0, 1\}$ and f^* such that $f^*(0) = f^*(1) = 0$. Consider a hypothesis class of constant functions $\mathcal{F} = \{f : \mathcal{X} \to \mathbb{R} \mid f(x) = d, \forall x \in \mathcal{X}\}$. We can see that $f^* \in \mathcal{F}$ \mathcal{F} . Assume that we have a uniform distribution over \mathcal{X} and we also have deterministic interval 1134 [l(x), u(x)]. Assume that for x = 0, we have an interval $[l(0), u(0)] = [-a, \epsilon]$ for some a > 0 and 1135 for x = 1, we have an interval $[l(1), u(1)] = [-\epsilon, 2\epsilon]$. Since \mathcal{F} is a class of constant hypothesis, for 1136 all x, we must have $f(x) \in [-a, \epsilon] \cap [-\epsilon, 2\epsilon] = [-\epsilon, \epsilon]$. This implies that

$$\widetilde{\mathcal{F}}_0 = \{ f \mid f(x) = c, \forall x \in \mathcal{X}, c \in [-\epsilon, \epsilon] \}.$$
(81)

1139 Therefore, 1140

$$f_1 = \arg\min_{f \in \mathcal{F}} \max_{f' \in \widetilde{\mathcal{F}}_0} \mathbb{E}[\ell(f(X), f'(X))]$$
(82)

$$= \arg\min_{f \in \mathcal{F}} \max_{f' \in \widetilde{\mathcal{F}}_0} \frac{1}{2} (|f(0) - f'(0)| + |f(1) - f'(1)|)$$
(83)

$$= \arg\min_{f \in \mathcal{F}} \max_{c \in [-\epsilon,\epsilon]} |f(0) - c|$$
(84)

(85)

1146 1147

1151 1152

1153 1154

1138

1148 By symmetry, we can see that the optimal $f_1(x) = 0$ which means that $err(f_1) = 0$. On the other 1149 hand, consider f_2 , from Corollary 4.2, f_2 is equivalent to the solution of supervised learning with 1150 the midpoint of each interval,

$$f_2 = \arg\min_{f \in \mathcal{F}} \mathbb{E}[\rho_\ell(f(X), L, U)]$$
(86)

$$= \arg\min_{f\in\mathcal{F}} \frac{1}{2} [|f(0) - \frac{-a+\epsilon}{2}| + |f(1) - \frac{-\epsilon+2\epsilon}{2}|].$$
(87)

By symmetry, the optimal f_2 should lie in the middle between these two points so that $f_2(x) = -a/2 + \epsilon$. We would have $\operatorname{err}(f_2) = |-a/2 + \epsilon|$ which can be arbitrarily large as $a \to \infty$.

1158

1159 E PROBABILISTIC INTERVAL SETTING

1160

In this section, we consider the probabilistic interval setting which is when, for each x, the corresponding interval is drawn from some distribution \mathcal{D}_I . We assume that \mathcal{D}_I is a nonatomic distribution i.e. it does not contain a point mass. We also use p to refer to the probability density function.

Assumption 2. A distribution P with a probability density function p(x) is a nonatomic distribution when for any x such that p(x) > 0 and for any $\epsilon > 0$, there exists a set $S_{x,\epsilon} \subseteq B(x,\epsilon)$ (a ball with radius ϵ) such that $\Pr(S_{x,\epsilon}) > 0$. We assume that the distribution D and D_I are nonatomic distributions.

Lemma E.1. Let P be a nonatomic distribution over \mathcal{X} with a probability density function p(x). For any continuous function $f : \mathcal{X} \to [0, \infty)$, if $\mathbb{E}_P[f(X)] = 0$ then f(x) = 0 for all x with p(x) > 0.

1172 1173 Proof. We will prove this by contradiction. Assume that there exists x with p(x) > 0 such that 1174 f(x) > 0. By the continuity of f, there exists $\delta_1 > 0$ such that for any $x' \in B(x, \delta_1)$ such 1175 that $|f(x) - f(x')| \le f(x)/2$ which implies that $f(x') \ge f(x)/2$. In addition, by the nonatomic 1176 assumption, there exists $S_{x,\delta_1} \subseteq B(x, \delta_1)$ such that $\Pr(S_{x,\delta_1}) > 0$. Therefore,

1177 1178

$$\mathbb{E}_{P}[f(X)] = \int_{w \in \mathcal{X}} f(w)p(w)dw$$
(88)

$$\geq \int_{w \in S_{x,\delta_1}} f(w) p(w) dw \tag{89}$$

1183
1184
$$\geq \int_{w \in S_{x,\delta_1}} \frac{f(x)p(w)}{2} dw$$
(90)

$$\frac{1185}{1186} = \frac{f(x)\Pr(S_{x,\delta_1})}{2} > 0.$$
(91)

This leads to a contradiction since $\mathbb{E}_P[f(X)] > 0$.

Similar to the deterministic interval setting, for any $f \in \widetilde{\mathcal{F}}_0$, f has to lie inside the interval as well. One difference would be that in the probabilistic interval setting, we can have multiple intervals for each x and since f has to lie inside all of them, f would also lie inside the intersection of all of them for which we denote as $[l_x, \tilde{u}_x]$ for each x.

Proposition E.2. For any $f \in \mathcal{F}_0$, and a loss function ℓ that satisfies Assumption 1, for any x with positive probability density p(x) > 0, we have

$$f(x) \in \bigcap_{p(x,l,u)>0} [l,u] := [\tilde{l}_x, \tilde{u}_x].$$

$$(92)$$

Proof. Let $f \in \mathcal{F}_0$ so we have $\mathbb{E}[\pi(f(X), L, U)] = 0$. From Lemma E.1, for any (x, l, u) such that p(x, l, u) > 0, we have $\pi(f(x), l, u) = 0$ which implies $f(x) \in [l, u]$ (From Proposition 2.1). There-fore, by taking an intersection over all possible intervals, we would have $f(x) \in \bigcap_{p(x,l,u)>0} [l,u] :=$ $[l_x, \tilde{u}_x].$

Proposition E.3. Let \mathcal{F} be a class of functions that are *m*-Lipschitz. For any x, x', denote $\tilde{l}_{x' \to x}^{(m)} =$ $\tilde{l}_{x'} - m \|x - x'\|, \ \tilde{u}_{x' \to x}^{(m)} = \tilde{u}_{x'} + m \|x - x'\|, \ \text{then for any } f \in \widetilde{\mathcal{F}}_0 \ \text{and for any } x \ \text{with positive probability density } p(x) > 0,$

$$f(x) \in \bigcap_{x'} [\tilde{l}_{x' \to x}^{(m)}, \tilde{u}_{x' \to x}^{(m)}] := [\tilde{l}_{\mathcal{D} \to x}^{(m)}, \tilde{u}_{\mathcal{D} \to x}^{(m)}]$$
(93)

Proof. Consider $f \in \mathcal{F}_0$, since f is m-Lipschitz, for any $x, x' \in \mathcal{X}$, we have $|f(x) - f(x')| \leq |f(x) - f(x)||$ m||x - x'|| which implies

$$f(x') - m \|x - x'\| \le f(x) \le f(x') + m \|x - x'\|$$
(94)

We illustrate this in Figure 2a. Then, from Proposition E.2, for $f \in \widetilde{\mathcal{F}}_0$, we have $\tilde{l}_{x'} \leq f(x') \leq \tilde{u}_{x'}$ which implies

1217
$$\tilde{l}_{x' \to x}^{(m)} = \tilde{l}_{x'} - m \|x - x'\| \le f(x') - m \|x - x'\|$$
(95)
1218 (m)

$$\tilde{u}_{x'\to x}^{(m)} = \tilde{u}_{x'} + m \|x - x'\| \ge f(x') - m \|x + x'\|.$$
(96)

Substitute back to equation equation 94 and take supremum over x', we have

$$\tilde{l}_{x' \to x}^{(m)} \le f(x) \le \tilde{u}_{x' \to x}^{(m)}$$
(97)

$$\sup_{x'} \tilde{l}_{x' \to x}^{(m)} \le f(x) \le \inf_{x'} \tilde{u}_{x' \to x}^{(m)}$$
(98)

$$\tilde{l}_{\mathcal{D}\to x}^{(m)} \le f(x) \le \tilde{u}_{\mathcal{D}\to x}^{(m)}.$$
(99)

Next, we present the probabilistic interval version of Theorem 3.3. Details of the proofs are the same, except that we use l, \tilde{u} instead of l, u.

Theorem E.4. Let \mathcal{F} be a class of functions that are *m*-Lipschitz. $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ is a loss function that satisfies Assumption 1. For any $f \in \widetilde{\mathcal{F}}_{\eta}$ and for any x with positive probability density p(x) > 0,

$$f(x) \in [\tilde{l}_{\mathcal{D} \to x}^{(m)} - r_{\eta}(x), \tilde{u}_{\mathcal{D} \to x}^{(m)} + s_{\eta}(x)]$$
(100)

where $\tilde{l}_{\mathcal{D}\to x}^{(m)}$, $\tilde{u}_{\mathcal{D}\to x}^{(m)}$ are defined as in Proposition E.3 and

1.
$$r_{\eta}(x) = r$$
 such that $\eta = \mathbb{E}[1[g(x, X, r) < L]\ell(g(x, X, r), L)]$ where $g(x, x', r) = \tilde{l}_{x'} - (r - (\tilde{l}_{\mathcal{D} \to x}^{(m)} - \tilde{l}_{x' \to x}^{(m)})).$

2.
$$s_{\eta}(x) = s$$
 such that $\eta = \mathbb{E}[1[h(x, X, s) > U]\ell(h(x, X, s), U)]$ where $h(x, x', s) = \tilde{u}_{x'} + (s - (\tilde{u}_{x' \to x}^{(m)} - \tilde{u}_{\mathcal{D} \to x}^{(m)})).$

1242 1243 1244 1245 Proof. Now, we will show that if $f \in \widetilde{\mathcal{F}}_{\eta}$ then we have $f(x) \in [\tilde{l}_{\mathcal{D}\to x}^{(m)} - r_{\eta}(x), \tilde{u}_{\mathcal{D}\to x}^{(m)} + s_{\eta}(x)]$ instead. First, we explore what would be a requirement to change the lower bound of f(x) from $\tilde{l}_{\mathcal{D}\to x}^{(m)}$ to $\tilde{l}_{\mathcal{D}\to x}^{(m)} - r$. Again, from Lipschitzness,

$$f(x') - m \|x - x'\| \le f(x)$$
(101)

1247 Taking a supremum here, we have

1246

1248 1249

1252 1253

1259

1260

1266

1268

1272

1277

1280

1281

1288 1289

$$\sup_{x'} f(x') - m \|x - x'\| \le f(x).$$
(102)

Here, we will use $\sup_{x'} f(x') - m ||x - x'||$ as a new lower bound for f(x). Assume that it is lower than $\tilde{l}_{D \to x}^{(m)}$, we can write

$$\sup_{x'} f(x') - m \|x - x'\| = \tilde{l}_{\mathcal{D} \to x}^{(m)} - r$$
(103)

for some r > 0, then it implies that for all $x' \in \mathcal{X}$, we must have

$$f(x') - m \|x - x'\| \le \tilde{l}_{\mathcal{D} \to x}^{(m)} - r$$
(104)

$$(f(x') - \tilde{l}_{x'} + (\tilde{l}_{x'} - m \|x - x'\|) \le \tilde{l}_{\mathcal{D} \to x}^{(m)} - r$$
(105)

$$f(x') \le \hat{l}_{x'} - \hat{l}_{x'\to x}^{(m)} + \hat{l}_{\mathcal{D}\to x}^{(m)} - r$$
(106)

$$f(x') \le \tilde{l}_{x'} - (r - (\tilde{l}_{\mathcal{D} \to x}^{(m)} - \tilde{l}_{x' \to x}^{(m)}))$$
(107)

1261 That is, if one can change the lower bound of f(x) from $\tilde{l}_{D\to x}^{(m)}$ to $\tilde{l}_{D\to x}^{(m)} - r$ then for all x', f(x') has 1262 to take value lower than $\tilde{l}_{x'}$ by at least $r - (\tilde{l}_{D\to x}^{(m)} - \tilde{l}_{x'\to x}^{(m)})$ whenever this term is positive. However, 1263 $f \in \widetilde{\mathcal{F}}_{\eta}$ so that f(x') can't be too far away from $\tilde{l}_{x'}$ since $\mathbb{E}[\pi_{\ell}(f(X), L, U)] \leq \eta$. From Proposition 1264 2.1, if one can write $\ell(y, y') = \psi(|y - y'|)$ for some non-decreasing function ψ then we have

$$\pi_{\ell}(f(x), l, u) = \mathbb{1}[f(x) < l]\ell(f(x), l) + \mathbb{1}[f(x) > u]\ell(f(x), u).$$
(108)

1267 Therefore,

$$\eta \ge \mathbb{E}[\pi_{\ell}(f(X), L, U)] \ge \mathbb{E}[\mathbb{1}[f(X) < L]\ell(f(X), L)].$$
(109)

Let $g(x, x', r) = \tilde{l}_{x'} - (r - (\tilde{l}_{\mathcal{D} \to x}^{(m)} - \tilde{l}_{x' \to x}^{(m)}))$ be the upper bound of f(x') for any x' as we derived in the equation equation 107. Since $1[a < L]\ell(a, L)]$ is a decreasing function over a, equation equation 109 implies

$$\eta \ge \mathbb{E}[\mathbf{1}[f(X) < L]\ell(f(X), L)] \ge \mathbb{E}[\mathbf{1}[g(x, X, r) < L]\ell(g(x, X, r), L)]$$
(110)

We can also see that g(x, x', r) is a decreasing function of r which means $\mathbb{E}[1[g(x, X, r) < L]\ell(g(x, X, r), L)]$ is an increasing function of r. The largest possible value of r would then be the r such that the inequality holds,

$$\eta = \mathbb{E}[1[g(x, X, r) < L]\ell(g(x, X, r), L)].$$
(111)

which we denoted this as $r_{\eta}(x)$. Similarly, we can show that if the largest possible value of s such that we can change the upper bound of f(x) from $\tilde{u}_{\mathcal{D}\to x}^{(m)}$ to $\tilde{u}_{\mathcal{D}\to x}^{(m)} + s$ is given by

$$\eta = \mathbb{E}[\mathbf{1}[h(x, X, s) > U]\ell(h(x, X, s), U)]$$
(112)

where
$$h(x, x', s) = \tilde{u}_{x'} + (s - (\tilde{u}_{x' \to x}^{(m)} - \tilde{u}_{\mathcal{D} \to x}^{(m)})).$$

Theorem E.5. Under the conditions of Theorem E.4, if further assume that for each x, the lower and upper bound of y is given by deterministic function [l(x), u(x)] and ℓ is an ℓ_p loss $\ell(y, y') =$ $|y - y'|^p$ and denote the lower bound gap and upper bound gap of f(x) induced by x' as $lg_{x'\to x}^{(m)} =$ $\tilde{l}_{D\to x}^{(m)} - \tilde{l}_{x'\to x}^{(m)}$ and $ug_{x'\to x}^{(m)} = \tilde{u}_{D\to x}^{(m)}$ then we have

$$r_{\eta}(x) = r$$
 s.t. $\mathbb{E}[(r - lg_{X \to x}^{(m)})_{+}^{p}] = \eta$ (113)

$$s_{\eta}(x) = s \quad s.t. \quad \mathbb{E}[(s - ug_{X \to x}^{(m)})_{+}^{p}] = \eta$$
 (114)

where we denote $c_{+} = \max(0, c)$. Further, we can bound $r_{\eta}(x)$ and $s_{\eta}(x)$,

1292
1293
1294

$$r_{\eta}(x) \leq \inf_{\delta} \delta + \left(\frac{\eta}{\Pr(lg_{X \to x}^{(m)} \leq \delta)}\right)^{1/p}$$
(115)

1295
$$s_{\eta}(x) \leq \inf_{\delta} \delta + \left(\frac{\eta}{\Pr(ug_{X \to x}^{(m)} \leq \delta)}\right)^{1/p}.$$
 (116)

Proof. Since [l, u] is deterministic for each x, we have $\tilde{l}_x = l(x)$. By the property of squared loss,

$$\mathbb{E}[1[g(x, X, r) < L]\ell(g(x, X, r), L)] = \mathbb{E}[(L - g(x, X, r))_{+}^{p}]$$
(117)

$$= \mathbb{E}[(l(X) - g(x, X, r))_{+}^{p}]$$
(118)

$$= \mathbb{E}[(l(X) - (\tilde{l}_X - (r - (\tilde{l}_{\mathcal{D} \to x}^{(m)} - \tilde{l}_{X \to x}^{(m)}))))_+^p]$$
(119)

$$= \mathbb{E}[(r - lg_{X \to x}^{(m)})_{+}^{p}]$$
(120)

as required. We can use a similar argument for $s_{\eta}(x)$. Next, we can see that for any valid value of r, 1304

$$\eta \ge \mathbb{E}[(r - lg_{X \to x}^{(m)})_{+}^{p}] \ge \mathbb{E}[(r - \delta)_{+}^{p} \mathbb{1}[lg_{X \to x}^{(m)} \le \delta]] = (r - \delta)_{+}^{p} \Pr(lg_{X \to x}^{(m)} \le \delta).$$
(121)

1307 By rearranging, $r \le \delta + (\frac{\eta}{\Pr(lg_{X\to x}^{(m)} \le \delta)})^{1/p}$. Taking the infimum over δ , we have the desired inequality. Again, we can apply the same idea for $s_{\eta}(x)$.

¹³⁵⁰ F DATASET STATISTICS

We provide the statistics of the datasets including the number of data points, the number of features, the minimum and maximum values of the target value and the approximated Lipschitz constant in Table 2. The Lipschitz constant here is approximated by calculating the proportion $\frac{|y-y'|}{||x-x'||}$ for all pairs of data points then the value is given by the 95th percentiles of these proportions. We perform this procedure to avoid the outliers which have a size of around two orders of magnitude bigger than the 95th percentile value (Figure 4). This allows us to approximate the level of smoothness that does appear in the dataset rather than use the maximum Lipschitz constant. One could also think of this as a probabilistic Lipschitz value rather than the classical notion (Urner & Ben-David, 2013).

Dataset	# data points	# features	[y min, y max]	Lipschitz constant
Abalone	4177	10	[1,29]	3.23
Airfoil	1503	5	[103, 141]	7.75
Concrete	1030	8	[2,83]	13.8
Housing	414	6	[7, 118]	11.68
Power plant	9568	4	[420,496]	14.18

Table 2: Dataset statistics.



Figure 4: The value of $\frac{|y-y'|}{||x-x'||}$ by percentiles. We use the 95th percentile of this value as an approximated Lipschitz constant for each dataset.

1404 G IMPACTS OF THE INTERVAL SIZE AND INTERVAL LOCATION



G.1 IMPACT OF THE INTERVAL SIZE



Figure 5: Test MAE when varying the maximum interval size $q_{\text{max}} \in \{0, 30, 60, 90, 120\}$ while $q_{\text{min}} = 0$.



Figure 6: Test MAE when varying the minimum interval size $q_{\min} \in \{0, 15, 30, 45, 60, 75, 90\}$ while $q_{\max} = 90$.

1448 We want to investigate the impact of interval size on the performance of the proposed methods. 1449 Intuitively, a smaller interval would make the problem easier. In the extreme case when the interval 1450 size is zero, we recover the supervised learning setting. Here, we assume that the interval location p1451 is still drawn uniformly from [0, 1] and we consider two experiments. First, we vary the maximum 1452 interval size $q_{\text{max}} \in \{0, 30, 60, 90, 120\}$ while keeping the minimum interval size $q_{\text{min}} = 0$. As 1453 expected, a larger maximum interval size leads to the drop in test performance across the boards 1454 (Figure 5). Second, we vary the minimum inter val size $q_{\min} \in \{0, 15, 30, 45, 60, 75, 90\}$ while keeping $q_{\rm max}$ fixed at 90. We can see that the test performance also decreases for all methods as 1455 we increase the minimum interval size (Figure 6). Notably, the standard minmax approach is highly 1456 sensitive to the interval size where its performance degrades significantly much more than other 1457 approaches in both experiments. This is due to the nature of the approach that wants to minimize the loss with respect to the worst-case label, as we have a larger interval, these worst-case labels can be much stronger and may not represent the property of the true labels anymore. On the other hand, our other minmax approaches and the projection approach are more robust to the change in the minimum interval size and the error only went up slightly for both experiments.



G.2.1 When y is more likely to be on one side of the interval (vary p_{\min})



Figure 7: Test MAE when varying the minimum interval location $p_{\min} \in \{0, 0.2, 0.4, 0.6, 0.8, 1\}$. In this case, when $p_{\min} = 0$ we have the uniform interval setting while when $p_{\min} = 1$, y true always lie on the upper bound of the intervals.



Figure 8: Test MAE when varying the minimum interval location $p_{\min} \in \{0, 0.2, 0.4, 0.6, 0.8, 1\}$. In this case, when $p_{\min} = 0$ we have the uniform interval setting while when $p_{\min} = 1$, y true always lies on the upper bound of the intervals.(no minmax approach)

In the previous settings, we assume that the location of the interval p is drawn uniformly from U[0, 1], that is, when y true is equally likely to be located at anywhere on the intervals. Here, we explore what would happen when it is not the case. We assume that we fixed $q_{\min} = 0, q_{\max} = 90$ and consider three scenarios. First, we consider when y is more likely to be on one side of the

interval. Here, we consider when $p \sim U[p_{\min}, 1]$ where $p_{\min} \in \{0, 0.2, 0.4, 0.6, 0.8, 1\}$ (Figure 7). In this case, when $p_{\min} = 0$ we have the uniform interval setting while when $p_{\min} = 1$, y true always lies on the upper bound of the intervals. We can see that the test MAE of all approaches increases as p_{\min} is larger. Again, the minmax approach performs much worse than others. One explanation for this is that the minmax with respect to. the label would encourage the model to be close to the middle point of each interval (Corollary 4.2). However, the the y true is far away from the midpoint leads to his phenomenon. We also provide the test MAE with no minmax approach for better visualization (Figure 8)

G.2.2 WHEN *y* TRUE IS MORE LIKELY TO BE IN THE MIDDLE OF THE INTERVAL



Figure 9: Test MAE when varying the interval location, $p \sim U[0.5 - c, 0.5 + c]$ for $c \in \{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$. When c = 0, the true y is always in the middle of the interval and when c = 0.5, we recover the uniform interval setting.

Second, we consider when y true is more likely to be in the middle of the interval (p is close to 0.5). We capture this setting by considering $p \sim U[0.5 - c, 0.5 + c]$ for $c \in \{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$ (Figure 9). Intuitively, when c = 0, the true y is always in the middle of the interval and when c = 0.5, we recover the uniform interval setting. In contrast to the first setting, we can see that the minmax approach performs the best in this setting for a small value of c. Again, this is perhaps due to the nature of the minmax approach mentioned earlier which encourages the prediction to be close to the middle point of the interval, for which, in this case, close to the y true. Remarkably, minmax performs better until c = 0.2 which corresponds to $p \sim [0.3, 0.7]$ which is a reasonable location of y true in practice. However, when c is large we would recover the uniform interval setting and the minmax would go back to becoming the worst-performer. On the other hand, the performance of other approaches is better as c is larger, that is when y true is more spread out across the interval.



1566 G.2.3 When y is more likely to be on either side of the interval

Figure 10: Test MAE when varying the interval location, when p is drawn uniformly from $[0, 0.5 - c] \cup [0.5 + c, 1]$ when $c \in \{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$. Here, when c = 0 we have the uniform interval setting while when c = 0.5, y true is either on the upper or the lower bound of the intervals.



Figure 11: Test MAE when varying the interval location, when p is drawn uniformly from $[0, 0.5 - c] \cup [0.5 + c, 1]$ when $c \in \{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$. Here, when c = 0 we have the uniform interval setting while when c = 0.5, y true is either on the upper or the lower bound of the intervals.(no minmax approach)

Finally, we consider when y is more likely to be on either side of the interval where p is drawn uniformly from $[0, 0.5 - c] \cup [0.5 + c, 1]$ when $c \in \{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$. Here, when c = 0we have the uniform interval setting while when c = 0.5, y true is either on the upper or the lower bound of the intervals. We found that as c is larger where the y true is more likely to be near either of the boundaries, the minmax performance drop significantly (Figure 10). However, we found that the performance of other approaches increases (Figure 11). This is in contrast to the first setting where we see that when y is more likely to be near only one side of the boundary, the performance drops remarkably. 1620 Overall, from these experiments, we may conclude that for all approaches apart from the original 1621 minmax with respect to. labels, having y true that lies near both of the boundaries of the interval are 1623 beneficial to the test performance and lying on both sides is crucial.

1624 G.3 LARGE AMBIGUITY DEGREE SETTING

We consider a setting with large ambiguity degree where $q \sim \text{Uniform}[q_{\min}, 90]$ when $q_{\min} \in$ $\{30, 60, 90\}$ and $p \sim \text{Uniform}[0.5 - c, 0.5 + c]$ when $c \in \{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$. Here as c is smaller, y true would be located near the middle point of the interval while as c is larger, we would recover the uniform setting. These settings have a large ambiguity degree since when $q_{\min} > 0$, interval size can't be arbitrarily small and $[p_{\min}, p_{\max}] \subset [0, 1]$ implies that true y would not lie at the boundary of the constructed interval. As a result, the intersection of all possible intervals would no longer be just $\{y\}$ anymore which leads to the ambiguity degree of 1. We found that there is no single method that always performs well on every interval setting. The Minmax is the best performing method for all $c \le 0.3$ while when c > 0.3 the best-performing approaches are either PL (mean) or PL (max) (Figure 12).



Figure 12: The best performing approach for each c and q_{\min}



G.4 INTERVAL PADDING EXPERIMENT





¹⁷²⁸ H INTERVAL SIZE AND TEST PERFORMANCE OF LIPMLP

Figure 14: Approximated interval size $I_{\eta}(x)$ for Lipschitz MLP with a different value of Lipschitz constant *m*. The dashed horizontal lines are the values from standard MLP.



Figure 15: Test MAE of the projection method with Lipschitz MLP with different values of Lipschitz constant. The vertical line is the Lipschitz constant approximated from the training set. The dashed horizontal lines are the test MAE of PL (Mean) and Projection approach with a standard MLP.

¹⁷⁸² I ABLATION FOR PL (MEAN)

Since PL (mean) is the best-performing approach in the uniform interval setting, we also performed an ablation study to improve our understanding of this method. First, we explore the impact of the number of hypotheses k used to represent $\tilde{\mathcal{F}}_0$. We found that for every dataset, as k is larger, the test MAE becomes smaller. While we use k = 5 for all PL experiments, this ablation suggests that we can increase k to get better performance at the cost of more computation.



Figure 16: Test MAE for PL (mean) with different number of hypotheses k used to represent \mathcal{F}_0 . For almost every dataset, the test MAE decreases as k is larger.

Second, we also compare PL (mean) with a natural ensemble baseline where we combine pseudo labels by averaging them first and then train a model with respect to. the averaged labels. In particular, the objective for the ensemble baseline is given by

$$\min_{f} \sum_{i=1}^{n} \ell(f(x_i), \sum_{j=1}^{k} f_j(x_i)).$$
(122)

1817 We found that PL (mean) still performs better than this baseline on 2 out of 5 datasets while the 1818 other 3 datasets are similar.

	Abalone	Airfoil	Concrete	Housing	Power-plant
PL (mean) PL ensemble baseline	$\frac{1.52_{0.01}}{1.51_{0.01}}$	$\begin{array}{c} 2.42_{0.07} \\ 3.3_{0.04} \end{array}$	$\begin{array}{c} 5.43_{0.12} \\ 5.57_{0.19} \end{array}$	$5.05_{0.09}$ $5.06_{0.08}$	$\frac{3.33_{0.01}}{3.32_{0.01}}$

1824Table 3: Test Mean Absolute Error (MAE) and the standard error (over 10 random seeds) for PL1825(Mean) and a PL ensemble baseline