

---

# Adversarial Resilience in Sequential Prediction via Abstention

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We study the problem of sequential prediction in the stochastic setting with an  
2 adversary that is allowed to inject clean-label adversarial (or out-of-distribution)  
3 examples. Algorithms designed to handle purely stochastic data tend to fail in the  
4 presence of such adversarial examples, often leading to erroneous predictions. This  
5 is undesirable in many high-stakes applications such as medical recommendations,  
6 where abstaining from predictions on adversarial examples is preferable to mis-  
7 classification. On the other hand, assuming fully adversarial data leads to very  
8 pessimistic bounds that are often vacuous in practice.

9 To capture this motivation, we propose a new model of sequential prediction that  
10 sits between the purely stochastic and fully adversarial settings by allowing the  
11 learner to abstain from making a prediction at no cost on adversarial examples.  
12 Assuming access to the marginal distribution on the non-adversarial examples, we  
13 design a learner whose error scales with the VC dimension (mirroring the stochastic  
14 setting) of the hypothesis class, as opposed to the Littlestone dimension which  
15 characterizes the fully adversarial setting. Furthermore, we design a learner for VC  
16 dimension 1 classes, which works even in the absence of access to the marginal  
17 distribution. Our key technical contribution is a novel measure for quantifying  
18 uncertainty for learning VC classes, which may be of independent interest.

## 19 1 Introduction

20 Consider the problem of sequential prediction in the realizable setting, where labels are generated  
21 from an unknown  $f^*$  belonging to a hypothesis class  $\mathcal{F}$ . Sequential prediction is typically studied  
22 under two distributional assumptions on the input data: the stochastic setting where the data is  
23 assumed to be identically and independently distributed (i.i.d) according to some fixed (perhaps  
24 unknown) distribution, and the fully-adversarial setting where we make absolutely no assumptions  
25 on the data generation process. A simple empirical risk minimization (ERM) strategy works for the  
26 stochastic setting where the learner predicts according to the best hypothesis on the data seen so  
27 far. The number of mistakes of this strategy typically scales with the Vapnik-Chervonenkis (VC)  
28 dimension of the underlying hypothesis class  $\mathcal{F}$ . However, in the fully adversarial setting, this strategy  
29 can lead to infinite mistakes even for classes of VC dimension 1 even if the adversary is required to  
30 be consistent with labels from  $f^*$ . The Littlestone dimension, which characterizes the complexity of  
31 sequential prediction in fully-adversarial setting, can be very large and often unbounded compared to  
32 the VC dimension [Lit87]. This mismatch has led to the exploration of beyond worst-case analysis  
33 for sequential prediction [RST11, HRS20, RS13a, BCKP20].

34 In this work, we propose a new framework that sits in between the stochastic and fully-adversarial  
35 setting. In particular, we consider sequential prediction with an adversary that injects adversarial  
36 (or out-of-distribution) examples in a stream of i.i.d. examples, and a learner that is allowed to

37 abstain from predicting on adversarial examples. A natural motivation for our framework arises  
38 in medical diagnosis where the goal is to predict a patient’s illness based on symptoms. In cases  
39 where the symptoms are not among the commonly indicative ones for the specific disease, or the  
40 symptoms may suggest a disease that is outside the scope of the doctor’s knowledge, it is safer for  
41 the doctor to abstain from making a prediction rather than risk making an incorrect one. Similarly,  
42 for self-driving cars, in cases where the car encounters weather conditions outside of its training, or  
43 unknown information signs, it is better for the algorithm to hand over access to the driver instead of  
44 making a wrong decision which could end up being fatal.

45 In the proposed framework, the learner’s goal is to minimize erroneous predictions on examples  
46 that the learner chooses to predict on (i.i.d. or adversarial) while refraining from abstaining on too  
47 many i.i.d. examples. If the learner was required to predict on every example, then the adversary  
48 could produce a fully-adversarial sequence of examples which would force the learner to make many  
49 erroneous predictions. The abstention option allows us to circumvent this challenge and handle any  
50 number of adversarial injections without incurring error proportional to the number of injections.  
51 Now we can ask the natural question:

52 *Is there a statistical price for certainty in sequential prediction?*

53 In particular, can we recover stochastic-like guarantees in the presence of an adversary if we are  
54 allowed to abstain from predicting on adversarial examples? A priori, it is not clear where on the  
55 spectrum between the fully-adversarial and stochastic models, the complexity of this problem lies.  
56 The main challenges arise from the fact that adversary fully controls the injection levels and provides  
57 no feedback about which examples were adversarial, and the learner has to perform one-sample  
58 outlier detection algorithm, which is nearly impossible. Despite this, we show that it is possible to  
59 guarantee certainty in a statistically efficient manner.

## 60 **1.1 Main Contributions**

61 We summarize the main contributions of our work:

- 62 • We formalize a new model of beyond-worst case learning which captures online learning on  
63 stochastic input with a clean-label attack adversary. With the option of abstention, our model  
64 allows for any number of injections by the adversary without incurring error proportional to the  
65 number of injections. Our notion of error simultaneously guarantees few mistakes on classified  
66 data while ensuring low abstention rate on non-adversarial data. Our model also naturally connects  
67 to uncertainty quantification and testable learning.
- 68 • We propose an algorithm that achieves error  $O(d^2 \log T)$  for classes with VC dimension  $d$  for time  
69 horizon  $T$ , given access to the marginal distribution over the i.i.d. examples. This allows us to get  
70 (up to a factor of  $d$ ) the guarantees of the stochastic setting while allowing for adversarial injections.
- 71 • We further propose an algorithm that achieves  $O(\sqrt{T})$  error for the special (but important) case of  
72 VC dimension 1 classes without any access to the marginal distribution over the i.i.d. examples.

73 Our algorithms uses a novel measure of uncertainty for VC classes to identify regions of high  
74 uncertainty (where the learner abstains) or high information gain (where the learner predicts and learn  
75 from their mistakes). The measure uses structural properties of VC classes, in particular, shattered  
76 sets of varying sizes. In the known distribution setting, our measure is easy to compute, however for  
77 the unknown distribution setting, we show how to design a proxy using only the examples we have  
78 seen so far using a leave-one-out type strategy.

## 79 **1.2 Related Work**

80 **Beyond-worst case sequential prediction.** Due to pessimistic nature of bounds in adversarial online  
81 learning, there are several frameworks designed to address this issue. One approach is consider  
82 mild restrictions on the adversarial instances such as slight perturbation by noise. This has been  
83 formalized as the smoothed adversary model (see [RST11, HRS20, HRS22, HHSY22, BDGR22])  
84 and has been used to get statistical and computationally efficient algorithms. Another approach has  
85 been to make the future sequences more predictable given the past instances. Examples of such  
86 settings are predictable sequence [RS13b], online learning with hints [BCKP20], and notions of  
87 adaptive regret bounds [FRS20].

88 **Abstention-based learning.** Abstention has been considered in several other works in classification,  
 89 both in the online and offline settings. An early example of this is the Chow reject model [Cho70].  
 90 Various versions of this have been considered in offline learning (see e.g. [HW06, BW08, BZ20] and  
 91 references therein) and online learning (see e.g. [ZC16, CDG<sup>+</sup>19, NZ20] and references therein).  
 92 These results show that abstention can lead to algorithms with desired features, for example fast  
 93 rates without margins assumptions. Another line of work that is closely related to our setting is the  
 94 KWIK (*knows what it knows*) framework by [LLW08] which requires the learner to make predictions  
 95 only when it is absolutely confident, and abstain otherwise. This requirement was relaxed to allow  
 96 for mistakes in [SZB10]. The key difference from our work is that it assumes a fully-adversarial  
 97 stream thus, the error bounds can be as large as the domain size. Perhaps, the work that is closest to  
 98 our setting is the study of adversarial learning with clean-label injections by [GKKM20, KK21]. In  
 99 their *transductive* adversarial model, the learner is given labeled training examples and unlabeled  
 100 test examples on which it must predict, where the test examples may have been injected by an  
 101 adversary. They show how to abstain with few test misclassifications and few false abstentions on  
 102 non-adversarial examples. However, in many real-world scenarios, it is unrealistic to expect to have  
 103 the entire test set in advance, which further motivates the fully online setting that we consider.

104 **Adversarially robust learning.** Highly related to our setting is the problem of *inductive* learning  
 105 in the presence of adversarial corruptions. The literature on this is generally divided into two  
 106 scenarios: test-time attacks and training-time attacks. In the case of test-time attacks, the learning  
 107 algorithm is trained on an (uncorrupted) i.i.d. training set, but its *test* examples may be corrupted  
 108 by an adversary whose intention is to change the classification by corrupting the test example  
 109 [SZS<sup>+</sup>13, BCM<sup>+</sup>13, GSS15, FMS18, AKM19, MHS19, MHS20, MHS21, MHS22, MGDS20]. On  
 110 the other hand, in the case of training-time attacks, the training data the learning algorithm trains on  
 111 is corrupted by an adversary (subject to some constraints on what fraction it may corrupt and what  
 112 types of corruptions are possible), while the test examples are uncorrupted [Val85, KL93, BEK02,  
 113 BNL12, ABL17, SKL17, SHN<sup>+</sup>18, LF21, GKM21, HKM<sup>+</sup>22, BBHS22]. In particular, within  
 114 this literature, most relevant to the present work is the work on *clean-label* poisoning, where the  
 115 adversary’s corrupted examples are still labeled correctly by the target concept [SHN<sup>+</sup>18, BHQS21].  
 116 Comparing to the present work, it is interesting to note that the fact that our setting involves *sequential*  
 117 prediction (i.e., online), our problem may be viewed simultaneously as *both* training-time and test-  
 118 time corruption: that is, because on each round the point we are predicting (or abstaining) on may be  
 119 inserted by an adversary, this could be viewed as a test-time attack; on the other hand, since the prefix  
 120 of labeled examples we use to make this prediction may also contain adversarial examples, this can  
 121 also be viewed as a training-time attack. Thus, our setting requires reasoning accounting for issues  
 122 arising from both attack scenarios, representing a natural blending of the two types of scenarios.

123 **Adversarial Examples.** Our clean-label attack is very closely related and motivated by the setting of  
 124 adversarial examples [SZS<sup>+</sup>13, BCM<sup>+</sup>13, GSS14]. The goal in this setting is to learn a classifier that  
 125 predicts correctly on all adversarial examples, which is a very strong requirement. Empirical work in  
 126 this space has focused on designing methods to make training adversarially robust [MMS<sup>+</sup>17, WK18],  
 127 and also on detecting adversarial examples [PDDZ18, AHFD22]. Detecting adversarial examples  
 128 is a very challenging tasks and proposed solutions are often brittle [CW17]. In fact, our framework  
 129 does not explicitly require detection as long as we can predict correctly on these.

## 130 2 Abstention Framework

131 In this section, we present the formal framework for sequential prediction with abstentions.

132 **Notation.** We will denote the domain with  $\mathcal{X}$  and the distribution over  $\mathcal{X}$  as  $\mathcal{D}$ . We let  $\Delta(\mathcal{X})$  denote  
 133 the set of all distributions over  $\mathcal{X}$ . We will work in the realizable setting where our label will be  
 134 according to some function in  $\mathcal{F}$  with VC dimension  $d$ . Given a class  $\mathcal{F}$  and a data set  $S = \{(x_i, y_i)\}$ ,  
 135 we will denote by  $\mathcal{F}|_S$ , the class  $\mathcal{F}|_S = \{f \in \mathcal{F} : \forall i \ f(x_i) = y_i\}$ . When the data set contains a  
 136 single point  $S = \{(x, y)\}$ , it will be convenient to denote  $\mathcal{F}|_S$  as  $\mathcal{F}^{x \rightarrow y}$ .

137 **Protocol.** At the start, the adversary (or nature) picks a distribution  $\mathcal{D}$  over the domain  $\mathcal{X}$  and the  
 138 labelling function  $f^* \in \mathcal{F}$ . We will be interested in both the setting where the learner knows the  
 139 distribution  $\mathcal{D}$  and the setting where the learner does not know the distribution  $\mathcal{D}$ . In the traditional  
 140 sequential prediction framework, the learner sees input  $x_t$  at time  $t$  and makes a prediction  $\hat{y}_t$  and  
 141 observes the true label  $y_t$ . The main departure of our setting from this is that an adversary also

142 decides before any round whether to inject an arbitrary element of  $\mathcal{X}$  (without seeing  $x_t$ ). We denote  
 143 by  $\hat{x}_t$  the instance after the adversarial injection ( $\hat{x}_t = x_t$  or  $\hat{x}_t \neq x_t$ ). The learner then observes  $\hat{x}_t$   
 144 and makes a prediction  $\hat{y}_t$  and observes the true label  $y_t$ , as in the traditional sequential prediction  
 145 framework. We present this formally as a protocol in 2.

---

**Protocol 1** Sequential Prediction with Adversarial Injections and Abstentions

---

Adversary (or nature) initially selects distribution  $\mathfrak{D} \in \Delta(\mathcal{X})$  and  $f^* \in \mathcal{F}$ . The learner does not have access to  $f^*$ . The learner may or may not have access to  $\mathfrak{D}$ .

**for**  $t = 1, \dots, T$  **do**

Adversary decides whether to inject an adversarial input in this the round ( $c_t = 1$ ) or not ( $c_t = 0$ ).

**if**  $c_t = 1$  **then** Adversary selects any  $\hat{x}_t \in \mathcal{X}$

**else** Nature selects  $\hat{x}_t \sim \mathfrak{D}$

Learner receives  $x_t$  and outputs  $\hat{y}_t \in \{0, 1, \perp\}$  where  $\perp$  implies that the learner abstains.

Learner receives clean label  $y_t = f^*(\hat{x}_t)$ .

---

146 It is important to note we are in the realizable setting, even after the adversarial injections since the  
 147 labels are always consistent with a hypothesis  $f^* \in \mathcal{F}$ . This model can naturally be extended to the  
 148 agnostic setting with adversarial labels.

149 **Objective.** The goal of the learner is to have low error rate on the rounds they decides to predict  
 150 (that is when  $\hat{y}_t \in \{0, 1\}$ ) while also ensuring that they do not abstain ( $\hat{y}_t = \perp$ ) on too many non-  
 151 adversarial rounds ( $c_t = 0$ ). More formally, the learner's objective is to minimize the following error  
 152 (or regret),

$$\text{Error} := \underbrace{\sum_{t=1}^T \mathbb{1}[\hat{y}_t = 1 - f^*(x_t)]}_{\text{MisclassificationError}} + \underbrace{\sum_{t=1}^T \mathbb{1}[c_t = 0 \wedge \hat{y}_t = \perp]}_{\text{AbstentionError}}.$$

153 We could formulate a relaxed cost-based version of this objective that allows us to trade-off these  
 154 errors, however, we will focus on the strong notion of error for this paper.

155 **Connections to Testable Learning.** A further interesting connection can be made by viewing our  
 156 model as an online version of testable learning framework of [RV22]. In order to see the analogy  
 157 more direct, we will focus on the setting where the learner knows the distribution  $\mathfrak{D}$ . In the setting, a  
 158 learning algorithm is seen as a tester-learner pair. The tester takes the data set as input and outputs  
 159 whether the data set passes the test. The algorithm then takes as input any data set that passes the test  
 160 and outputs a hypothesis. The soundness guarantee for the pair of algorithms is that the algorithm  
 161 run on any data set that passes the test must output a hypothesis that is good on the distribution. The  
 162 completeness requires that when the dataset is indeed from the "nice" distribution, then the tester  
 163 passes with high probability. We can see our framework in this light by noting that the decision  
 164 of whether to abstain or not serves as a test. Thus, in this light, completeness corresponds to the  
 165 abstention error being small when the data is non-adversarial i.e. is from the true distribution, while  
 166 the soundness corresponds to the misclassification error being small on points the algorithm decides  
 167 not to abstain. While the testable learning literature primarily focuses on the computational aspects  
 168 of learning algorithms, our focus is solely on the statistical aspects.

### 169 3 Warm-up: Disagreement-based Learners

170 As a first example to understand the framework, we consider the most natural learner for the problem.  
 171 Given the data  $S$  of the examples seen thus far, the learner predicts on examples  $\hat{x}$  whose labels it is cer-  
 172 tain of. That is, if there is a unique label for  $\hat{x}$  consistent with  $\mathcal{F}|_S$ , the learner predicts that label. Else,  
 173 it abstains from making a prediction. This region of uncertainty is known as the disagreement region.

174 **Example: thresholds in one dimension.** Consider learning a single-dimensional threshold in  $[0, 1]$   
 175 (that is, concepts  $x \mapsto \mathbb{1}[x \geq t]$  for any  $t \in [0, 1]$ ). While it is well known that ERM achieves  $\log(T)$   
 176 misclassification error for i.i.d. data sequences, in the case of an adversarially chosen sequence,  
 177 it is also well known that the adversary can select inputs in the disagreement region each time  
 178 (closer and closer to the decision boundary) and thereby force any non-abstaining learner to make  
 179 a linear number of mistakes (recall that the Littlestone dimension of thresholds is infinite) [Lit87].

180 Indeed, it is known that the function classes  $\mathcal{F}$  for which non-abstaining predictors can be forced  
 181 to have MisclassificationError =  $\Omega(T)$  are precisely those with embedded threshold problems of  
 182 unbounded size [Lit87, She78, Hod97, ALMM19]. Let us now consider the learner that abstains in  
 183 the disagreement region and predicts based on the consistent hypothesis outside of this region.

184 **Proposition 3.1.** *Disagreement-based learner for one dimensional thresholds has*  
 185 *MisclassificationError = 0 and AbstentionError  $\leq 2 \log(T)$ .*

186 To see this note that our learner only predicts when the input is not in the disagreement region  
 187 and thus it never predicts incorrectly (MisclassificationError = 0). As for the abstentions, a simple  
 188 exchangeability, argument shows that when there are  $n$  i.i.d. examples in the sequence, the probability  
 189 of the new non-adversarial example being in the disagreement region is  $1/n$ . Summing this over the  
 190 time horizon gives us the above proposition.

191 **Perfect Selective Classification and Active Learning.** The learner for the above thresholds problem  
 192 is a well-known strategy from the areas of perfect selective classification and active learning known  
 193 as *disagreement-based learning* [RS88, EYW10, CAL94, BBL09, DHM07, HY15]. In the perfect  
 194 selective classification setting [RS88, EYW10], the learner observes the examples sequentially, as in  
 195 our framework, and may predict or abstain on each, and must satisfy the requirement that whenever  
 196 it predicts its prediction is *always* correct. From this fact, it immediately follows that applying any  
 197 perfect selective classification strategy in our setting, we always have MisclassificationError = 0,  
 198 so that its performance is judged purely on its abstention rate on the iid examples. It was argued  
 199 by [EYW10] that the optimal abstention rate among perfect selective classification strategies is  
 200 obtained by the strategy that makes a prediction on the next example if and only if all classifiers in  
 201 the hypothesis class that are correct on all examples observed so far, *agree* on the example. Note that  
 202 this is precisely the same learner used above. This same strategy has also been studied in the related  
 203 setting of *stream-based active learning*<sup>1</sup> [CAL94, BBL09, Han09, Han14, DHM07, HY21]. The  
 204 abstention rate achievable by this strategy for general hypothesis classes is thoroughly understood  
 205 [Han07, Han11, Han09, Han12, Han14, Han16, EYW10, EYW12, WHE15, HY15]. In particular,  
 206 a complete characterization of the optimal distribution-free abstention rate of perfect selective  
 207 classification is given by the *star number* of the hypothesis class [HY15, Han16]. The star number  
 208  $\mathfrak{s}$  is the size of the largest number  $s$  such that there are examples  $\{x_1, \dots, x_s\}$  and hypotheses  
 209  $h_0, h_1, \dots, h_s$  such that  $h_i$  and  $h_0$  disagree exactly on  $x_i$ . For instance,  $\mathfrak{s} = 2$  for threshold classifiers  
 210 [HY15]. It was shown by [Han16] that the optimal distribution-free abstention rate for perfect  
 211 selective classification is sublinear if and only if the star number is finite (in which case it is always  
 212 at most  $\mathfrak{s} \log(T)$ ). One can show that  $\mathfrak{s}$  is always lower bounded by the VC dimension of the class.  
 213 Unfortunately, the star number is infinite for most hypothesis classes of interest, even including  
 214 simple VC classes such as *interval* classifiers [HY15].

215 **Beyond disagreement-based learning.** The learner that abstains whenever it sees an example that it  
 216 is not certain of may be too conservative. Furthermore, it does not exploit the possibility of learning  
 217 from mistakes. Let us consider another example to elucidate this failure. Consider the class of  $d$   
 218 intervals in one dimension where the positive examples form a union of intervals. This class has VC  
 219 dimension  $d$  but infinite star number. Suppose that  $d = 2$  but examples in the second interval are very  
 220 rarely selected by i.i.d. examples. Then the disagreement-based learner would suggest to abstain  
 221 on all examples to protect against the possibility that our new example is in the second interval.  
 222 However, consider the following simple strategy: if the new example lies between two positives (resp.  
 223 negatives), we predict positive (resp. negative), else we abstain.

224 **Proposition 3.2.** *The proposed strategy for the class of  $d$ -intervals in one dimension has*  
 225 *MisclassificationError  $\leq d$  and AbstentionError  $\leq 2d \log(T)$ .*

226 To see this, note that whenever we predict, either we are correct or we have identified the location of  
 227 a new interval, hence reducing the VC dimension of the class by 1. Since there are at most  $d$  intervals,  
 228 we will therefore make at most  $d$  errors when we predict, implying MisclassificationError  $\leq d$ . As  
 229 for abstaining on i.i.d. examples, the same argument for thresholds can be applied here by treating  
 230 the intervals as at most  $d$  thresholds.

<sup>1</sup>In this setting, instead of observing a sequence of labelled examples, the learner only observes the examples without their target labels, and at each time may query to observe the target label. The disagreement-based strategy chooses to query precisely on the points for which the classifiers in the hypothesis class correct on the observed labels so far do not all agree on the label [CAL94]. The rate of querying for this strategy is precisely the same as the abstention rate in the perfect selection classification setting [Han11, Han14, EYW12].

231 **4 Higher-order Disagreement-based learner with Known marginals**

232 We will first focus on the setting when the marginal distribution  $\mathcal{D}$  is known. In this setting, the  
 233 algorithm that naturally suggests itself is to take a cover for the class under  $\mathcal{D}$  of accuracy  $\text{poly}(T^{-1})$   
 234 and use an adversarial algorithm for prediction. Since there are covers of size  $T^{O(d)}$  and Littlestone  
 235 dimension of any finite class is bounded by the logarithm of the size, this seems to indicate that this  
 236 algorithm will achieve our goal with misclassification error  $O(d \log T)$  and zero abstention error. But  
 237 unfortunately note that this algorithm competes only with the best classifier on the cover. The cover  
 238 is a good approximation only on the marginal distribution  $\mathcal{D}$  and not on the adversarial examples.  
 239 In fact, when we restrict to the hypothesis class being the cover, the data may no longer even be  
 240 realizable by the cover. Therefore, we need to use the access to the distribution in a completely  
 241 different way.

242 The inspiration for our approach comes from the work of [Han09, Han12] on active learning strategies  
 243 that go beyond disagreement-based learning by making use of higher-order notions of disagreement  
 244 based on *shattering*. We note, however, that while the work of [Han09, Han12] only yields asymptotic  
 245 and distribution-dependent guarantees (and necessarily so, in light of the minimax characterization of  
 246 [HY15] based on the star number), our analysis differs significantly in order to yield distribution-free  
 247 finite-sample bounds on the misclassification error and abstention rate.

248 As we saw earlier, just looking at the disagreement region does not lead to a good algorithm for  
 249 general VC classes (whenever the star number is large compared to the VC dimension). The main  
 250 algorithmic insight of this section is that certain higher-order versions of disagreement sets do indeed  
 251 lead to good notions of uncertainty for VC classes. Our measure uses the probability of shattering  $k$   
 252 examples (for different values of  $k$ ) freshly drawn from the underlying distribution, under the two  
 253 restrictions of the class corresponding to the two labels for the current test example, to make the  
 254 abstention decision. One can think of the probability of shattering as a proxy for the complexity of  
 255 the version space. This serves both as a method to quantify our uncertainty about whether the current  
 256 example is from the distribution or not, since we can understand the behavior of this quantity in the  
 257 i.i.d. case, and also as potential function which keeps track of the number of mistakes. In order to  
 258 formally state this, we will need some definitions.

259 **Definition 4.1** (Shattering and VC Dimension). Let  $\mathcal{X}$  be a domain and  $\mathcal{F}$  be a binary function class  
 260 on  $\mathcal{X}$  i.e.  $\mathcal{F} \subset \{0, 1\}^{\mathcal{X}}$ . A set  $\{x_1, \dots, x_k\} \subseteq \mathcal{X}$  is said to be shattered by  $\mathcal{F}$  if for all  $y \in \{0, 1\}^k$   
 261 there exists a function  $f \in \mathcal{F}$  such that  $f(x_i) = y_i$ . The VC dimension of  $\mathcal{F}$  is defined as the  
 262 maximum  $k$  such that there is a set of size  $k$  that is shattered.

263 **Definition 4.2** (Shattered  $k$ -tuples). Let  $k$  be a positive integer. The set of shattered  $k$ -tuples, denoted  
 264 by  $\mathcal{S}_k$ , for hypothesis class  $\mathcal{F}$  over a domain  $\mathcal{X}$  is defined as

$$\mathcal{S}_k(\mathcal{F}) = \{(x_1, \dots, x_k) : \{x_1, \dots, x_k\} \text{ is shattered by } \mathcal{F}\}.$$

265 Additionally, given a distribution  $\mathcal{D}$  on the domain, we will refer to as the  $k$  shattering probability of  
 266  $\mathcal{F}$  with respect to  $\mathcal{D}$ , denoted by  $\rho_k(\mathcal{F}, \mathcal{D})$ , as

$$\rho_k(\mathcal{F}, \mathcal{D}) = \mathcal{D}^{\otimes k}(\mathcal{S}_k(\mathcal{F})) = \Pr_{x_1, \dots, x_k \sim \mathcal{D}^{\otimes k}} [\{x_1, \dots, x_k\} \text{ is shattered by } \mathcal{F}].$$

267 Let us now describe the algorithm (see Algorithm 1). The algorithm maintains a state variable  $k$   
 268 which we will refer to as the level the algorithm is currently in. The level can be thought of as the  
 269 complexity of the current version space. At level  $k$ , we will work with shattered sets of size  $k$ . At  
 270 each round, the algorithm, upon receiving the example  $\hat{x}_t$ , computes the probabilities of shattering  
 271  $k$  examples (drawn i.i.d. from  $\mathcal{D}$ ) for each of the classes corresponding to sending  $\hat{x}_t$  to 0 and 1  
 272 respectively. The algorithm abstains if both these probabilities are large, else predicts according to  
 273 whichever one is larger. At the end of the round, after receiving the true label  $y_t$ , the algorithm checks  
 274 whether the probability of shattering  $k$  examples is below a threshold  $\alpha_k$ , in which case it changes  
 275 the level, that is, updates  $k$  to be  $k - 1$ .

276 Below, we state the main error bound of the algorithm. The theorem shows that both the misclassifica-  
 277 tion error and the abstention error are bounded in terms of the VC dimension.

278 **Theorem 4.1.** *Let  $\mathcal{F}$  be a hypothesis class with VC dimension  $d$ . Then, in the corruption model with*  
 279 *abstentions with time horizon  $T$ , Algorithm 1 with  $\alpha_k = T^{-k}$  gives the following guarantee*

$$\mathbb{E}[\text{MisclassificationError}] \leq d^2 \log T \quad \text{and} \quad \mathbb{E}[\text{AbstentionError}] \leq 6d.$$

---

**Algorithm 1:** Level-based learning for Prediction with Abstention

---

Set  $k = d$  and  $\mathcal{F}_1 = \mathcal{F}$   
**for**  $t = 1, \dots, T \wedge k > 1$  **do**  
  Receive  $\hat{x}_t$   
  **if**  $\min \left\{ \rho_k \left( \mathcal{F}_t^{\hat{x}_t \rightarrow 1} \right), \rho_k \left( \mathcal{F}_t^{\hat{x}_t \rightarrow 0} \right) \right\} \geq 0.6 \rho_k \left( \mathcal{F}_t \right)$  **then** predict  $\hat{y}_t = \perp$   
  **else** predict  $\hat{y}_t = \operatorname{argmax}_{j \in \{0,1\}} \left\{ \rho_k \left( \mathcal{F}_t^{\hat{x}_t \rightarrow j} \right) \right\}$   
  Upon receiving label  $y_t$ , update  $\mathcal{F}_{t+1} \leftarrow \mathcal{F}_t^{\hat{x}_t \rightarrow y_t}$   
  **if**  $\rho_k \left( \mathcal{F}_{t+1} \right) \leq \alpha_k$  **then** Set  $k = k - 1$   
**if**  $k = 1 \wedge \hat{x}_t \in \mathcal{S}_1$  **then**  $\hat{y}_t = \perp$   
**if**  $k = 1 \wedge \hat{x}_t \notin \mathcal{S}_1$  **then** Predict with the consistent label for  $\hat{x}_t$ 

---

280 We will now present the main technical idea that we will use to analyze the theorem. We defer the full  
281 proofs to Appendix A. Let  $k$  be a positive integer and  $x$  be any example in  $\mathcal{X}$ . The following lemma  
282 upper bounds the probability for a random example the shattering probability of the two restrictions  
283 of the class are both large compared to the original shattering probability of the class. That is to say  
284 for most examples, one of the two restrictions of the class will have a smaller shattering probability  
285 compared to the original class.

286 **Lemma 4.2.** *Let  $\mathcal{F}$  be a hypothesis class and  $\mathcal{D}$  be a distribution. Then for any  $k \in \mathbb{N}$  and any*  
287  *$\eta > 1/2$ , we have*

$$\Pr_{x \sim \mathcal{D}} \left[ \rho_k \left( \mathcal{F}^{x \rightarrow 1} \right) + \rho_k \left( \mathcal{F}^{x \rightarrow 0} \right) \geq 2\eta \rho_k \left( \mathcal{F} \right) \right] \leq \frac{1}{2\eta - 1} \cdot \frac{\rho_{k+1} \left( \mathcal{F} \right)}{\rho_k \left( \mathcal{F} \right)}. \quad (1)$$

288 With this in hand, we will first look at the abstention error. The intuition is that when the algorithm  
289 is at level  $k$ , an abstention occurs only if the condition from (1) is satisfied and Lemma 4.2 bounds  
290 the probability of this event. It remains to note that when the algorithm is at level  $k$  we both have an  
291 upper bound on  $\rho_{k+1}$  (since this is the condition to move down to level  $k$  from level  $k + 1$ ) and a  
292 lower bound on the  $\rho_k$  (since this is the condition to stay at level  $k$ ).

293 **Lemma 4.3** (Abstention error). *For any  $k \leq d$ , let  $[\ell_k, e_k]$  denote the interval of time when Algo-*  
294 *rithm 1 is at level  $k$ . Then, the expected number of non-adversarial rounds at level  $k$  on which the*  
295 *algorithm abstains satisfies*

$$\mathbb{E} \left[ \sum_{t=\ell_k}^{e_k} \mathbb{I} [c_t = 0 \wedge \hat{y}_t = \perp] \right] \leq 5T \cdot \frac{\alpha_{k+1}}{\alpha_k}. \quad (2)$$

296 Next, we bound the misclassification error. The main idea here is to note that every time a misclassi-  
297 cation occurs at level  $k$ , the  $k$ -th disagreement coefficient reduces by a constant factor. Since we have  
298 a lower bound on the disagreement coefficient at a fixed level, this leads to a logarithmic bound on  
299 the number of misclassifications at any given level.

300 **Lemma 4.4** (Misclassification error). *For any  $k \leq d$ , let  $[\ell_k, e_k]$  denote the interval of time when*  
301 *Algorithm 1 is at level  $k$ . For any threshold  $\alpha_k$  in Algorithm 1,*

$$\mathbb{E} \left[ \sum_{t=\ell_k}^{e_k} \mathbb{I} [\hat{y}_t = 1 - f^*(\hat{x}_t)] \right] \leq 2 \cdot \log \left( \frac{1}{\alpha_k} \right).$$

302 Putting together Lemma 4.4 and Lemma 4.3 along with a setting of  $\alpha_k = T^{-k}$ , gives us Theorem 4.1.

## 303 5 Structure-based Algorithm for VC Dimension 1 Classes

304 We move to the case of unknown distributions. In this setting, the example  $x$  are drawn from a  
305 distribution  $\mathcal{D}$  that is unknown to the learner. As we saw earlier, it is challenging to decide whether a

306 single point is out of distribution or not, even when the distribution is known. This current setting  
 307 is significantly more challenging since the learner needs to abstain on examples that are out of  
 308 distribution for a distribution that it doesn't know. A natural idea would be to use the historical data to  
 309 build a model for the distribution. The main difficulty is that, since we do not get feedback about what  
 310 examples are corrupted, our historical data has both in-distribution and out-of-distribution examples.  
 311 The only information we have about what examples are out of distribution is the prediction our  
 312 algorithm has made on them. Such issues are a major barrier in moving from the known distribution  
 313 case to the unknown distribution case.

314 A key quantity that we will use in our algorithm will be a version of the probability of the disagreement  
 315 region but computed on the historical data. The most natural version of this would be the leave-one-  
 316 out disagreement. That is, consider the set of examples which are in the disagreement region when  
 317 the class is restricted using the data set with the point under consideration removed. This estimate  
 318 would have been an unbiased estimator for the disagreement probability (referred to earlier as  $\rho_1$  in  
 319 Definition 4.2). As mentioned earlier, unfortunately, in the presence of adversarial injections, this  
 320 need not be a good estimate for the disagreement probability.

321 In order to remedy this, we consider a modified version of the leave-one-out disagreement estimate  
 322 which considers examples  $x$  in the disagreement region for the class  $\mathcal{F}|_{S_f \setminus (x,y)}$  where  $S_f$  is the  
 323 subset of the datapoints which disagrees with a fixed reference function  $f$ . It is important to note  
 324 that this function  $f$  is fixed independent of the true labelling function  $f^*$ . Though this seems a bit  
 325 artificial at first, we will see that this is a natural quantity to consider given the structure theorem for  
 326 VC dimension one classes which we will discuss subsequently.

327 **Definition 5.1.** Let  $\mathcal{F}$  be a class of functions and let  $f \in \mathcal{F}$  be a reference function. Let  $S =$   
 328  $\{(x_i, y_i)\}$  be a realizable data set. Define

$$\Gamma(S, \mathcal{F}, f) = \left\{ x : \exists y \quad (x, y) \in S \wedge x \in \mathcal{S}_1 \left( \mathcal{F}|_{S_f \setminus (x,y)} \right) \right\}$$

329 where  $S_f = \{(x, y) \in S : f(x) \neq y\}$ . Further, we will denote the size of this set by  $\gamma(S, \mathcal{F}) =$   
 330  $|\Gamma(S, \mathcal{F}, f)|$ . We will suppress the dependence on  $\mathcal{F}$  when it is clear from context.

331 Let us now present the main algorithm (see Algorithm 2). The main idea of the algorithm is similar to  
 332 Algorithm 1 in the known distribution case. We will make a prediction in the case when the difference  
 333 between  $\Gamma$  for the two classes corresponding to the two labels for the point  $x$  is large. The idea is that  
 334 when a prediction is made and a mistake happens, the size of  $\Gamma$  goes down similar to  $\rho_k$  in the known  
 335 distribution case. In the case when the difference is small, we will abstain.

---

**Algorithm 2:** Structure-based learning for Prediction with Abstention for Unknown Distribution

---

Let  $f \in \mathcal{F}$  be a reference function and  $\alpha$  be the abstention threshold.

Set  $\mathcal{F}_0 = \mathcal{F}$  and  $S_0 = \emptyset$

**for**  $t = 1, \dots, T$  **do**

  Receive  $\hat{x}_t$

336   Let  $a_0 = \left| \Gamma(S_{t-1}, \mathcal{F}^{\hat{x}_t \rightarrow 0}, f) \right|$  and  $a_1 = \left| \Gamma(S_{t-1}, \mathcal{F}^{\hat{x}_t \rightarrow 1}, f) \right|$

**if**  $\hat{x}_t \notin \mathcal{S}_1(\mathcal{F}_t)$  **then** Predict with the consistent label for  $\hat{x}_t$

**else if**  $\max\{a_0, a_1\} \geq \alpha$  **then**  $\hat{y}_t = \operatorname{argmax}_b a_b$

**else**  $\hat{y}_t = \perp$

  Upon receiving label  $y_t$ , update  $S_t \leftarrow S_{t-1} \cup \{(\hat{x}_t, y_t)\}$  and  $\mathcal{F}_t \leftarrow \mathcal{F}_{t-1}^{\hat{x}_t \rightarrow y_t}$

---

337 Though the algorithm is simple and can be made fairly general, our analysis is restricted to the case  
 338 of VC dimension one. The main reason for this restriction is that our analysis relies on a structure  
 339 theorem for VC dimension one classes which has no direct analogy for higher VC dimension. But  
 340 since the algorithm has a natural interpretation independent of this representation, we expect similar  
 341 algorithms to work for higher VC dimension classes as well.

342 **Theorem 5.1.** Let  $\mathcal{F}$  be a hypothesis class with VC dimension 1. Then, in the corruption model with  
 343 abstentions with time horizon  $T$ , Algorithm 2 with parameter  $\alpha = \sqrt{T}$  gives the following guarantee

$$\mathbb{E}[\text{MisclassificationError}] \leq 2\sqrt{T \log T} \quad \text{and} \quad \mathbb{E}[\text{AbstentionError}] \leq \sqrt{T \log T}.$$

344 We will now present the main technical ideas that we use to analyze the algorithm. We defer the full  
 345 proofs to Appendix B. We will keep track of the mistakes using the size of the disagreement region,

346 which we denote by  $\gamma_t = \gamma(S_t, \mathcal{F})$ . The main idea for the proof is to note that when we decide to  
347 predict the label with the larger value of  $\Gamma$  is bigger by an additive  $\alpha$ . Thus, when a mistake occurs the  
348 value of  $\gamma$  decreases by at least  $\alpha$ . Summing the errors over all time steps gives the following bound.

349 **Lemma 5.2.** *Algorithm 2 has  $\mathbb{E}[\text{MisclassificationError}] \leq 2T/\alpha$ .*

350 Next, we move on to the analysis of the number of abstentions. In fact, we will prove a stronger  
351 structural results that shows that the number of examples such that there is any set of adversarial  
352 injections that would make the algorithm abstain is small. We refer to examples on which algorithm  
353 can be made to abstain as attackable examples (formally defined in Definition B.2). The main idea  
354 is to prove that in any set of iid examples, there are only a few attackable ones. This is formally  
355 stated and proved as Lemma B.3. Using this claim, we can bound the number of abstentions using  
356 an exchangeability argument.

357 **Lemma 5.3.** *Algorithm 2 has  $\mathbb{E}[\text{AbstentionError}] \leq \alpha \log T$ .*

358 The proof of Theorem 5.1 follows from Lemma 5.2 and Lemma 5.3 by setting  $\alpha = \sqrt{T/\log T}$ .

## 359 6 Discussion and Future Directions

360 In this paper, we introduce a framework for beyond-worst case adversarial sequential predictions by  
361 using abstention. Our proposed framework falls at the intersection of several learning paradigms such  
362 as active learning, uncertainty quantification, and distribution testing. In this framework we show two  
363 main positive results, validating the beyond-worse case nature of our framework. However, our work  
364 has only scratched the surface of understanding learnability for VC classes in this new framework.  
365 Here we discuss several exciting future directions:

- 366 • Our techniques rely strongly on the realizability, however our framework can naturally be extended  
367 to settings with label noise. Immediate questions here would be to extend our results to simple  
368 noise models such as random classification noise and Massart noise or more ambitiously the  
369 agnostic noise model.
- 370 • Our framework can naturally be extended to more general forms of prediction such as multiclass  
371 classification, partial concept classes, and regression. It would be interesting to characterize  
372 learnability in these settings.
- 373 • In the known distribution case, it remains to find the optimal error bound. We conjecture that the  
374 correct dependence on VC dimension  $d$  should be linear.
- 375 • In the unknown distribution case, extending our result for VC dimension  $d > 1$  classes is wide  
376 open. Though, we have algorithms for certain classes of higher VC dimension, such as intervals  
377 and axis-aligned rectangles, they seem to heavily exploit the structure of the particular class. Thus,  
378 showing either an upper or lower bounds of the error dependence on both the class and the time  
379 horizon would be interesting.
- 380 • The unknown distribution case can be seen as form of distribution-free uncertainty quantification.  
381 It would be interesting to understand connections to other forms such as conformal prediction. On  
382 a technical level, our work exploits exchangeability of the i.i.d. sequence which is the foundation  
383 of conformal prediction, though the main challenge in our setting is the presence of adversarial  
384 inputs. It would be interesting to build on this connections and understand whether techniques  
385 can be ported over in either direction.
- 386 • Our focus for this paper has been entirely on the statistical benefits of abstention. Understanding  
387 the computational complexity in this setting is an exciting avenue for research. Concretely, for  
388 halfspaces in  $d$  dimensions, is there a polynomial time algorithm for learning with abstentions,  
389 even for well-behaved distributions such as Gaussians. On a related note, showing computational-  
390 statistical gaps in this specialized setting would be interesting, albeit disappointing.

391 **Broader Impact.** The aim of our work is design a theoretical framework for handling adversarial  
392 examples and guaranteeing certainty in sequential prediction. Our model assumes no cost for  
393 abstaining on adversarial inputs. However, in real-world scenarios, these costs can be high, consider  
394 the human cost of evaluating each of the abstained examples. Furthermore, as with any framework,  
395 application of our metric without taking into account the nuances of the real-world situation could lead  
396 to unintended consequences, for example, a high abstention rate on marginalized sub-populations.

397 **References**

- 398 [ABL17] Pranjal Awasthi, Maria Florina Balcan, and Philip M Long. The power of localization  
399 for efficiently learning linear separators with noise. *Journal of the ACM (JACM)*,  
400 63(6):1–27, 2017.
- 401 [AHFD22] Ahmed Aldahdooh, Wassim Hamidouche, Sid Ahmed Fezza, and Olivier Déforges.  
402 Adversarial example detection for dnn models: A review and experimental comparison.  
403 *Artificial Intelligence Review*, 55(6):4403–4462, 2022.
- 404 [AKM19] Idan Attias, Aryeh Kontorovich, and Yishay Mansour. Improved generalization bounds  
405 for robust learning. In *Proceedings of the 30<sup>th</sup> International Conference on Algorithmic*  
406 *Learning Theory*, 2019.
- 407 [ALMM19] Noga Alon, Roi Livni, Maryanthe Malliaris, and Shay Moran. Private PAC learning  
408 implies finite littlestone dimension. In *Proceedings of the 51st Annual ACM Symposium*  
409 *on Theory of Computing (STOC)*, page 852–860, 2019.
- 410 [BBHS22] M.-F. Balcan, A. Blum, S. Hanneke, and D. Sharma. Robustly-reliable learners under  
411 poisoning attacks. In *Proceedings of the 35<sup>th</sup> Annual Conference on Learning Theory*,  
412 2022.
- 413 [BBL09] M.-F. Balcan, A. Beygelzimer, and J. Langford. Agnostic active learning. *Journal of*  
414 *Computer and System Sciences*, 75(1):78–89, 2009.
- 415 [BCKP20] Aditya Bhaskara, Ashok Cutkosky, Ravi Kumar, and Manish Purohit. Online learning  
416 with imperfect hints. In *Proceedings of the International Conference on Machine*  
417 *Learning (ICML)*, pages 822–831. PMLR, 2020.
- 418 [BCM<sup>+</sup>13] Battista Biggio, Iginio Corona, Davide Maiorca, Blaine Nelson, Nedim Šrndić, Pavel  
419 Laskov, Giorgio Giacinto, and Fabio Roli. Evasion attacks against machine learning at  
420 test time. In *Joint European conference on machine learning and knowledge discovery*  
421 *in databases*, pages 387–402. Springer, 2013.
- 422 [BD15] Shai Ben-David. 2 notes on classes with vapnik-chervonenkis dimension 1, 2015.
- 423 [BDGR22] Adam Block, Yuval Dagan, Noah Golowich, and Alexander Rakhlin. Smoothed online  
424 learning is as easy as statistical learning. In *Proceedings of the Conference on Learning*  
425 *Theory (COLT)*, pages 1716–1786. PMLR, 2022.
- 426 [BEK02] Nader H Bshouty, Nadav Eiron, and Eyal Kushilevitz. PAC learning with nasty noise.  
427 *Theoretical Computer Science*, 288(2):255–275, 2002.
- 428 [BHQ21] Avrim Blum, Steve Hanneke, Jian Qian, and Han Shao. Robust learning under clean-  
429 label attack. In *Conference on Learning Theory (COLT)*, 2021.
- 430 [BNL12] Battista Biggio, B Nelson, and P Laskov. Poisoning attacks against support vector  
431 machines. In *29th International Conference on Machine Learning*, pages 1807–1814.  
432 ArXiv e-prints, 2012.
- 433 [BW08] Peter L. Bartlett and Marten H. Wegkamp. Classification with a reject option using a  
434 hinge loss. *Journal of Machine Learning Research*, 9(59):1823–1840, 2008.
- 435 [BZ20] Olivier Bousquet and Nikita Zhivotovskiy. Fast classification rates without standard  
436 margin assumptions, 2020.
- 437 [CAL94] David Cohn, Les Atlas, and Richard Ladner. Improving generalization with active  
438 learning. *Machine Learning*, 15:201–221, 1994.
- 439 [CDG<sup>+</sup>19] Corinna Cortes, Giulia DeSalvo, Claudio Gentile, Mehryar Mohri, and Scott Yang.  
440 Online learning with abstention, 2019.
- 441 [Cho70] C. Chow. On optimum recognition error and reject tradeoff. *IEEE Transactions on*  
442 *Information Theory*, 16(1):41–46, 1970.

- 443 [CW17] Nicholas Carlini and David Wagner. Adversarial examples are not easily detected:  
444 Bypassing ten detection methods. In *Proceedings of the 10th ACM Workshop on*  
445 *Artificial Intelligence and Security*, pages 3–14. ACM, 2017.
- 446 [DHM07] S. Dasgupta, D. Hsu, and C. Monteleoni. A general agnostic active learning algorithm.  
447 In *Advances in Neural Information Processing Systems 20*, 2007.
- 448 [EYW10] Ran El-Yaniv and Yair Wiener. On the foundations of noise-free selective classification.  
449 *Journal of Machine Learning Research*, 11(5), 2010.
- 450 [EYW12] R. El-Yaniv and Y. Wiener. Active learning via perfect selective classification. *Journal*  
451 *of Machine Learning Research*, 13(2):255–279, 2012.
- 452 [FMS18] Uriel Feige, Yishay Mansour, and Robert E. Schapire. Robust inference for multiclass  
453 classification. In Firdaus Janoos, Mehryar Mohri, and Karthik Sridharan, editors,  
454 *Algorithmic Learning Theory, ALT 2018, 7-9 April 2018, Lanzarote, Canary Islands,*  
455 *Spain*, volume 83 of *Proceedings of Machine Learning Research*, pages 368–386. PMLR,  
456 2018.
- 457 [FRS20] Dylan J. Foster, Alexander Rakhlin, and Karthik Sridharan. Adaptive online learning,  
458 2020.
- 459 [GKKM20] Shafi Goldwasser, Adam Tauman Kalai, Yael Tauman Kalai, and Omar Montasser.  
460 Selective classification algorithms provably robust to arbitrary adversarial examples. In  
461 *Advances in Neural Information Processing Systems*, 2020.
- 462 [GKM21] Ji Gao, Amin Karbasi, and Mohammad Mahmoody. Learning and certification under  
463 instance-targeted poisoning. *The Conference on Uncertainty in Artificial Intelligence*  
464 *(UAI)*, 2021.
- 465 [GSS14] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing  
466 adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- 467 [GSS15] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing  
468 adversarial examples. In *Proceedings of the 3rd International Conference on Learning*  
469 *Representations, ICLR*, 2015.
- 470 [Han07] S. Hanneke. A bound on the label complexity of agnostic active learning. In *Proceedings*  
471 *of the 24<sup>th</sup> International Conference on Machine Learning*, 2007.
- 472 [Han09] S. Hanneke. *Theoretical Foundations of Active Learning*. PhD thesis, Machine Learning  
473 Department, School of Computer Science, Carnegie Mellon University, 2009.
- 474 [Han11] S. Hanneke. Rates of convergence in active learning. *The Annals of Statistics*, 39(1):333–  
475 361, 2011.
- 476 [Han12] S. Hanneke. Activized learning: Transforming passive to active with improved label  
477 complexity. *Journal of Machine Learning Research*, 13(5):1469–1587, 2012.
- 478 [Han14] S. Hanneke. Theory of disagreement-based active learning. *Foundations and Trends in*  
479 *Machine Learning*, 7(2–3):131–309, 2014.
- 480 [Han16] S. Hanneke. Refined error bounds for several learning algorithms. *Journal of Machine*  
481 *Learning Research*, 17(135):1–55, 2016.
- 482 [HHSY22] Nika Haghtalab, Yanjun Han, Abhishek Shetty, and Kunhe Yang. Oracle-efficient  
483 online learning for beyond worst-case adversaries. In *Advances in Neural Information*  
484 *Processing Systems (NeurIPS)*, 2022. to appear.
- 485 [HKM<sup>+</sup>22] S. Hanneke, A. Karbasi, M. Mahmoody, I. Mehalal, and S. Moran. On optimal learning  
486 under targeted data poisoning. In *Advances in Neural Information Processing Systems*  
487 *36*, 2022.
- 488 [Hod97] Wilfrid Hodges. *A Shorter Model Theory*. Cambridge University Press, 1997.

- 489 [HRS20] Nika Haghtalab, Tim Roughgarden, and Abhishek Shetty. Smoothed analysis of online  
490 and differentially private learning. *Advances in Neural Information Processing Systems*,  
491 33:9203–9215, 2020.
- 492 [HRS22] Nika Haghtalab, Tim Roughgarden, and Abhishek Shetty. Smoothed analysis with  
493 adaptive adversaries. In *Proceedings of the Annual Symposium on Foundations of*  
494 *Computer Science (FOCS)*, pages 942–953, 2022.
- 495 [HW06] Radu Herbei and Marten H. Wegkamp. Classification with reject option. *The Canadian*  
496 *Journal of Statistics / La Revue Canadienne de Statistique*, 34(4):709–721, 2006.
- 497 [HY15] Steve Hanneke and Liu Yang. Minimax analysis of active learning. *Journal of Machine*  
498 *Learning Research*, 16(109):3487–3602, 2015.
- 499 [HY21] S. Hanneke and L. Yang. Toward a general theory of online selective sampling: Trading  
500 off mistakes and queries. In *Proceedings of the 24<sup>th</sup> International Conference on*  
501 *Artificial Intelligence and Statistics*, 2021.
- 502 [KK21] Adam Tauman Kalai and Varun Kanade. Towards optimally abstaining from prediction  
503 with ood test examples, 2021.
- 504 [KL93] Michael Kearns and Ming Li. Learning in the presence of malicious errors. *SIAM*  
505 *Journal on Computing*, 22(4):807–837, 1993.
- 506 [LF21] Alexander Levine and Soheil Feizi. Deep partition aggregation: Provable defense against  
507 general poisoning attacks. *International Conference on Learning Representations*  
508 *(ICLR)*, 2021.
- 509 [Lit87] Nick Littlestone. Learning quickly when irrelevant attributes abound: A new linear-  
510 threshold algorithm. In *28th Annual Symposium on Foundations of Computer Science*  
511 *(sfcs 1987)*, pages 68–77, 1987.
- 512 [LLW08] Lihong Li, Michael L Littman, and Thomas J Walsh. Knows what it knows: a framework  
513 for self-aware learning. In *Proceedings of the 25th international conference on Machine*  
514 *learning*, pages 568–575, 2008.
- 515 [MGDS20] Omar Montasser, Surbhi Goel, Ilias Diakonikolas, and Nathan Srebro. Efficiently learn-  
516 ing adversarially robust halfspaces with noise. In *Proceedings of the 37th International*  
517 *Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume  
518 119 of *Proceedings of Machine Learning Research*, pages 7010–7021. PMLR, 2020.
- 519 [MHS19] Omar Montasser, Steve Hanneke, and Nathan Srebro. Vc classes are adversarially  
520 robustly learnable, but only improperly. In Alina Beygelzimer and Daniel Hsu, editors,  
521 *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99 of *Pro-*  
522 *ceedings of Machine Learning Research*, pages 2512–2530, Phoenix, USA, 25–28 Jun  
523 2019. PMLR.
- 524 [MHS20] Omar Montasser, Steve Hanneke, and Nati Srebro. Reducing adversarially robust learn-  
525 ing to non-robust PAC learning. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia  
526 Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Infor-*  
527 *mation Processing Systems 33: Annual Conference on Neural Information Processing*  
528 *Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- 529 [MHS21] Omar Montasser, Steve Hanneke, and Nathan Srebro. Adversarially robust learning  
530 with unknown perturbation sets. *CoRR*, abs/2102.02145, 2021.
- 531 [MHS22] O. Montasser, S. Hanneke, and N. Srebro. Adversarially robust learning: A generic  
532 minimax optimal learner and characterization. In *Advances in Neural Information*  
533 *Processing Systems 36*, 2022.
- 534 [MMS<sup>+</sup>17] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian  
535 Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint*  
536 *arXiv:1706.06083*, 2017.

- 537 [NZ20] Gergely Neu and Nikita Zhivotovskiy. Fast rates for online prediction with abstention,  
538 2020.
- 539 [PDDZ18] Tianyu Pang, Chao Du, Yinpeng Dong, and Jun Zhu. Towards robust detection of  
540 adversarial examples. In *Advances in Neural Information Processing Systems*, pages  
541 4579–4589, 2018.
- 542 [RS88] Ronald L. Rivest and Robert H. Sloan. Learning complicated concepts reliably and  
543 usefully (extended abstract). In Tom Mitchell and Reid Smith, editors, *Proceedings*  
544 *AAAI-88*, pages 635–640. AAAI, 1988.
- 545 [RS13a] Alexander Rakhlin and Karthik Sridharan. Online learning with predictable sequences.  
546 In *Proceedings of the Conference on Learning Theory (COLT)*, pages 993–1019. PMLR,  
547 2013.
- 548 [RS13b] Sasha Rakhlin and Karthik Sridharan. Optimization, learning, and games with pre-  
549 dictable sequences. *Advances in Neural Information Processing Systems*, 26, 2013.
- 550 [RST11] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. Online learning: Stochastic  
551 and constrained adversaries. *arXiv preprint arXiv:1104.5070*, 2011.
- 552 [RV22] Ronitt Rubinfeld and Arsen Vasilyan. Testing distributional assumptions of learning  
553 algorithms. *arXiv preprint arXiv:2204.07196*, 2022.
- 554 [She78] Saharon Shelah. *Classification Theory and the Number of Non-isomorphic Models*.  
555 North Holland, 1978.
- 556 [SHN<sup>+</sup>18] Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suci, Christoph Studer, Tudor  
557 Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks  
558 on neural networks. In *Advances in Neural Information Processing Systems*, pages  
559 6103–6113, 2018.
- 560 [SKL17] Jacob Steinhardt, Pang Wei Koh, and Percy Liang. Certified defenses for data poisoning  
561 attacks. *Advances in Neural Information Processing Systems*, 30, 2017.
- 562 [SZB10] Amin Sayedi, Morteza Zadimoghaddam, and Avrim Blum. Trading off mistakes and  
563 don’t-know predictions. In *Advances in Neural Information Processing Systems*, pages  
564 2092–2100, 2010.
- 565 [SZS<sup>+</sup>13] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian  
566 Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint*  
567 *arXiv:1312.6199*, 2013.
- 568 [Val85] Leslie G Valiant. Learning disjunctions of conjunctions. In *Proceedings of the 9th*  
569 *International Joint Conference on Artificial intelligence*, pages 560–566, 1985.
- 570 [WHE15] Y. Wiener, S. Hanneke, and R. El-Yaniv. A compression technique for analyzing  
571 disagreement-based active learning. *Journal of Machine Learning Research*, 16(4):713–  
572 745, 2015.
- 573 [WK18] Eric Wong and Zico Kolter. Provable defenses against adversarial examples via the  
574 convex outer adversarial polytope. In *International conference on machine learning*,  
575 pages 5286–5295. PMLR, 2018.
- 576 [ZC16] Chicheng Zhang and Kamalika Chaudhuri. The extended littlestone’s dimension for  
577 learning with mistakes and abstentions. In *Conference on Learning Theory*, pages  
578 1584–1616. PMLR, 2016.

579 **A Proofs from Section 4**

580 **A.1 Properties of Higher-Order Disagreement: Proof of Lemma 4.2**

581 We first begin by relating the probabilities of shattering  $k$  points for the classes gotten by restricting  
 582 to evaluating to 0 and 1 at  $x$  respectively, for an arbitrary point  $x$ . The proof of the following lemma  
 583 uses a simple inclusion exclusion argument.

584 **Lemma A.1.** *For all  $x \in \mathcal{X}$  and any hypothesis class  $\mathcal{F}$ , we have*

$$\mathfrak{D}^{\otimes k} \left( \mathcal{S}_k \left( \mathcal{F}^{x \rightarrow 1} \right) \right) + \mathfrak{D}^{\otimes k} \left( \mathcal{S}_k \left( \mathcal{F}^{x \rightarrow 0} \right) \right) \leq \mathfrak{D}^{\otimes k} \left( \mathcal{S}_k \left( \mathcal{F} \right) \right) + \mathfrak{D}^{\otimes k} \left( \mathcal{S}_k \left( \mathcal{F}^{x \rightarrow 1} \right) \cap \mathcal{S}_k \left( \mathcal{F}^{x \rightarrow 0} \right) \right).$$

585 *Equivalently,*

$$\rho_k \left( \mathcal{F}^{x \rightarrow 1} \right) + \rho_k \left( \mathcal{F}^{x \rightarrow 0} \right) \leq \rho_k \left( \mathcal{F} \right) + \Pr_{x_1, \dots, x_k \sim \mathfrak{D}} \left[ x_1, \dots, x_k \text{ is shattered by both } \mathcal{F}^{x \rightarrow 1} \text{ and } \mathcal{F}^{x \rightarrow 0} \right].$$

586 *Proof.* For  $x_1, \dots, x_n \in \mathcal{X}$ , consider the four indicator random variables given by

$$\begin{aligned} A_1 &= \mathbb{1} \left( x_1, \dots, x_n \text{ is shattered by } \mathcal{F}^{x \rightarrow 1} \right) \\ A_2 &= \mathbb{1} \left( x_1, \dots, x_n \text{ is shattered by } \mathcal{F}^{x \rightarrow 0} \right) \\ A_3 &= \mathbb{1} \left( x_1, \dots, x_n \text{ is shattered by } \mathcal{F}^{x \rightarrow 1} \text{ and } \mathcal{F}^{x \rightarrow 0} \right) \\ A_4 &= \mathbb{1} \left( x_1, \dots, x_n \text{ is shattered by } \mathcal{F} \right). \end{aligned}$$

587 Note that  $A_1 + A_2 \leq A_3 + A_4$ . Taking expectations gives the desired result.  $\square$

588 In order to prove the main lemma, we then take expectations with respect to the point  $x$  drawn  
 589 independently from  $\mathfrak{D}$ . The key observation is to relate the probability of  $k$  point being shattered by  
 590 both classes that evaluate to 0 and 1 at  $x$ , for a random point  $x$ , to the probability of shattering  $k + 1$   
 591 points.

592 *Proof of Lemma 4.2.* The proof follows by using Lemma A.1 to get (3) and Markov's inequality  
 593 to get (4). The final line follows by noting that since  $x_1, \dots, x_k$  and  $x$  are drawn from  $\mathfrak{D}$  and are  
 594 shattered, (4) computes the probability that  $k + 1$  points are shattered.

$$\begin{aligned} & \Pr_{x \sim \mathfrak{D}} \left[ \rho_k \left( \mathcal{F}^{x \rightarrow 0} \right) + \rho_k \left( \mathcal{F}^{x \rightarrow 1} \right) \geq 2\eta \rho_k \left( \mathcal{F} \right) \right] \\ & \leq \Pr_{x \sim \mathfrak{D}} \left[ \rho_k \left( \mathcal{F} \right) + \Pr_{x_1, \dots, x_k \sim \mathfrak{D}} \left[ x_1, \dots, x_k \text{ is shattered by both } \mathcal{F}^{x \rightarrow 1} \text{ and } \mathcal{F}^{x \rightarrow 0} \right] \geq 2\eta \rho_k \left( \mathcal{F} \right) \right] \quad (3) \\ & \leq \Pr_{x \sim \mathfrak{D}} \left[ \Pr_{x_1, \dots, x_k \sim \mathfrak{D}} \left[ x_1, \dots, x_k \text{ is shattered by both } \mathcal{F}^{x \rightarrow 1} \text{ and } \mathcal{F}^{x \rightarrow 0} \right] \geq (2\eta - 1) \rho_k \left( \mathcal{F} \right) \right] \\ & \leq \frac{\mathbb{E} \left[ \Pr_{x_1, \dots, x_k \sim \mathfrak{D}} \left[ x_1, \dots, x_k \text{ is shattered by both } \mathcal{F}^{x \rightarrow 1} \text{ and } \mathcal{F}^{x \rightarrow 0} \right] \right]}{(2\eta - 1) \rho_k \left( \mathcal{F} \right)} \quad (4) \\ & \leq \frac{1}{2\eta - 1} \cdot \frac{\rho_{k+1} \left( \mathcal{F} \right)}{\rho_k \left( \mathcal{F} \right)}. \end{aligned}$$

595  $\square$

596 **A.2 Proof of Lemma 4.4**

597 *Proof.* Note that from definition of Algorithm 1, we have that when  $\hat{y}_t \neq \perp$ ,  
 598  $\min \left\{ \rho_k \left( \mathcal{F}_t^{x_t \rightarrow 1} \right), \rho_k \left( \mathcal{F}_t^{x_t \rightarrow 0} \right) \right\} \leq 0.6 \rho_k \left( \mathcal{F}_t \right)$ . Thus, since we predict with the label corre-  
 599 sponding to  $\max \left\{ \rho_k \left( \mathcal{F}_t^{x_t \rightarrow 1} \right), \rho_k \left( \mathcal{F}_t^{x_t \rightarrow 0} \right) \right\}$ , if we make a mistake, we have

$$\rho_k \left( \mathcal{F}_{t+1} \right) \leq 0.6 \cdot \rho_k \left( \mathcal{F}_t \right).$$

600 Also, note once end of the phase corresponding to  $k$ ,  $e_k$ , is reached when  $\rho_k(\mathcal{F}_t) \leq \alpha_k$ . This leads  
 601 to the required bound.  $\square$

### 602 A.3 Proof of Lemma 4.3

603 *Proof.* Let  $t \in [\ell_k, m_k]$ . Recall that  $\mathcal{F}_t$  denotes the class consistent with the data seen till time  $t$ .  
 604 Recall that  $H_t$  denotes the history of the interaction till time  $t$ .

$$\mathbb{E} \left[ \mathbb{1} [c_t = 0 \wedge \hat{y}_t(\hat{x}_t) = \perp] \mid H_t \right] = \mathbb{E} \left[ \mathbb{1} [c_t = 0 \wedge \hat{y}_t(x_t) = \perp] \mid H_t \right] \quad (5)$$

$$\leq \mathbb{E} \left[ \mathbb{1} [\hat{y}_t(x_t) = \perp] \mid H_t \right]$$

$$\leq \Pr \left[ \min \left\{ \rho_k(\mathcal{F}_t^{x_t \rightarrow 1}), \rho_k(\mathcal{F}_t^{x_t \rightarrow 0}) \right\} \geq 0.6\rho_k(\mathcal{F}_t) \mid H_t \right] \quad (6)$$

$$\leq \Pr \left[ \rho_k(\mathcal{F}_t^{x_t \rightarrow 1}) + \rho_k(\mathcal{F}_t^{x_t \rightarrow 0}) \geq 1.2\rho_k(\mathcal{F}_t) \mid H_t \right] \quad (7)$$

$$\leq 5 \frac{\rho_{k+1}(\mathcal{F}_t)}{\rho_k(\mathcal{F}_t)} \quad (8)$$

$$\leq 5 \frac{\alpha_{k+1}}{\alpha_k}. \quad (9)$$

605 The equality in (5) follows from the fact that in the uncorrupted rounds,  $\hat{x}_t = x_t$ . The inequality in  
 606 (6) follows from the condition for abstention in Algorithm 1 and the fact that  $x_t \sim \mathcal{D}$ . (7) follows  
 607 from the fact that the min of two numbers is at most their average. The key step is (8) which follows  
 608 from Lemma 4.2. (9) follows from the fact that at level  $k$  in Algorithm 1,  $\rho_{k+1}(\mathcal{F}_t) \leq \alpha_k$  and  
 609  $\rho_k(\mathcal{F}_t) \geq \alpha_{k+1}$ . Summing this bound and noting that  $e_k - \ell_k$  is at most  $T$ , we get the required  
 610 bound.  $\square$

### 611 A.4 Proof of Theorem 4.1

612 *Proof.* First, let us look at the misclassification error. Note that Algorithm 1 will not misclassify  
 613 when  $k = 1$ . To see this, recall from Definition 4.2 that  $\hat{x}_t \in \mathcal{S}_1$  implies that there is a unique label  
 614 consistent with the history. For the remaining levels, we sum the errors from Lemma 4.4 and recall  
 615 that  $\alpha_k = T^{-k}$ , which gives us

$$\begin{aligned} \text{MisclassificationError} &= \sum_{t=1}^T \mathbb{1}[\hat{y}_t = 1 - f^*(x_t)] \\ &\leq 2 \sum_{k=2}^d \log \left( \frac{1}{\alpha_k} \right) \\ &\leq 2 \sum_{k=2}^d k \log T \\ &\leq d^2 \log T. \end{aligned}$$

616 For the abstention error, we again begin with the case of  $k = 1$ . Note that for  $t \geq \ell_1$ , we have

$$\Pr_{x \sim \mathcal{D}} [x \in \mathcal{S}_1(\mathcal{F}_t)] \leq \alpha_1.$$

617 Thus, we have a bound of  $T\alpha_1 \leq 1$  on the expected error in this case. For the remaining levels, we  
 618 sum the error from (2) over all  $k$ , which gives us the bound

$$\begin{aligned}
 \text{AbstentionError} &= \sum_{t=1}^T \mathbb{1}[c_t = 0 \wedge \hat{y}_t = \perp] \\
 &= \sum_{k=1}^d \mathbb{E} \left[ \sum_{t=\ell_k}^{e_k} \mathbb{1}[c_t = 0 \wedge \hat{y}_t = \perp] \right] \\
 &= 1 + 5T \sum_{k=2}^d \frac{\alpha_{k+1}}{\alpha_k} \\
 &= 1 + 5T \sum_{k=2}^d \frac{1}{T} \\
 &\leq 5d + 1 \\
 &\leq 6d.
 \end{aligned}$$

619 This gives us the desired bounds. □

## 620 B Proofs from Section 5

### 621 B.1 Structure for VC Dimension One Classes

622 **Definition B.1.** Consider a domain  $\mathcal{X}$  and a partial order  $\prec$  on  $\mathcal{X}$ . We say that a set  $I$  is an initial  
 623 segment of  $\prec$  if for all  $x \in I$  and  $y \in \mathcal{X}$  such that  $y \prec x$ , we have  $y \in I$ . We say that a partial order  
 624 is a tree ordering if every initial segment  $I$  is a linear order i.e. for all  $x, y \in I$ , either  $x \prec y$  or  $y \prec x$ .

625 **Theorem B.1** ([BD15]). *Let  $\mathcal{F}$  be a hypothesis class over the domain  $\mathcal{X}$ . Then, the following are*  
 626 *equivalent:*

627 a.  $\mathcal{F}$  has VC dimension 1.

628 b. There is a tree ordering  $\prec$  on  $\mathcal{X}$  and a hypothesis  $f \in \mathcal{F}$  such that every element of the set  

$$\mathcal{F}_f = \{h \oplus f : h \in \mathcal{F}\}$$

629 is an initial segment of  $\prec$ .

630 The result above was initially observed in [BD15].

### 631 B.2 Proof of Lemma 5.2

632 **Lemma B.2.** *For any  $t$ , we have that*

$$\gamma_{t+1} \leq \gamma_t - \alpha \cdot \mathbb{1}[\text{Misclassification at time } t] + 1. \quad (10)$$

633 *Proof.* First, note that in any round that a mistake was not made, we have that  $\gamma_{t+1} \leq \gamma_t + 1$ . This is  
 634 because at most one point is added to the data set in each round.

635 Note that  $|\Gamma(S_{t-1}, \mathcal{F}^{\hat{x}_t \rightarrow 1})| + |\Gamma(S_{t-1}, \mathcal{F}^{\hat{x}_t \rightarrow 0})| \leq |\Gamma(S_{t-1}, \mathcal{F})|$ . Further, we have  
 636  $|\Gamma(S_t, \mathcal{F})| \leq |\Gamma(S_{t-1}, \mathcal{F}^{\hat{x}_t \rightarrow y_t})| + 1$ . From the condition for predicting, we have that  
 637  $\max \left\{ |\Gamma(S_{t-1}, \mathcal{F}^{\hat{x}_t \rightarrow 1})|, |\Gamma(S_{t-1}, \mathcal{F}^{\hat{x}_t \rightarrow 0})| \right\} \geq \alpha$ . Since if we make a misclassification, we go to  
 638 the smaller value for  $\Gamma$ , we get the desired bound. □

639 *Proof.* Note that  $\text{MisclassificationError} = \sum_{t=1}^T \mathbb{1}[\text{Misclassification at time } t]$ . Rearranging (10)  
 640 and summing gives us

$$\text{MisclassificationError} \leq \frac{1}{\alpha} \sum_{t=1}^T (\gamma_t - \gamma_{t+1} + 1)$$

641 Note that for all  $i$ , we have that  $\gamma_i \leq T$  which gives us the desired bound.  $\square$

### 642 B.3 Proof of Lemma 5.3

643 **Definition B.2** (Attackable Point). Let  $\mathcal{F}$  be a hypothesis class and let  $f \in \mathcal{F}$  be a representative  
644 function. Let  $S$  be a realizable dataset. We say that a point  $x$  is attackable with respect to a data set  
645  $S$  if there exists a sequence of adversarial examples  $A_x$  such that algorithm abstains on example  
646  $x$  when the history is  $S \cup A_x \setminus \{x\}$ . In other words,  $x$  is attackable if there is a set of adversarial  
647 examples  $A_x$  such that  $x \notin \mathcal{S}_1 \left( \mathcal{F} |_{S \cup A_x \setminus \{x\}} \right)$  and

$$\max \left\{ \left| \Gamma \left( S \cup A_x \setminus \{x\}, \mathcal{F}^{x \rightarrow 0} \right) \right|, \left| \Gamma \left( S \cup A_x \setminus \{x\}, \mathcal{F}^{x \rightarrow 1} \right) \right| \right\} \leq \alpha.$$

648 The key lemma for our analysis is that the number of attackable examples is bounded.

649 **Lemma B.3.** *Let  $\mathcal{F}$  be a hypothesis and  $f \in \mathcal{F}$  be any reference function. Let  $S$  be any set of*  
650 *examples. Then, the number of attackable examples is at most  $\alpha$ .*

651 The proof uses a structure theorem for classes with VC dimension one which states that they can be  
652 represented as initial segments of a tree order.

653 *Proof.* Let  $f$  be the representative function from the characterization in Theorem B.1. Further, let  $\mathfrak{T}$   
654 be the tree corresponding to the tree order on  $\mathcal{X}$ . Since  $f$  is a fixed function that does not depend on  
655 the algorithm or the history of interaction with adversary, we can preprocess all points and labels  
656 to be xored with the labels of  $f$ . In other words, we transform the class to be such that  $f$  is the all  
657 zeros function. In this setting, the true hypothesis  $f^*$  corresponds to a path  $p$  and a threshold  $x^*$  on  
658  $\mathfrak{T}$  such that  $f^*(x) = 1$  if and only if  $x \in p$  and  $x \prec x^*$ . For this proof, it is important to consider  
659 only adversaries that do not get to choose the true hypothesis  $f^*$  adaptively. Thus the path is fixed  
660 throughout the history of the interaction.

661 Let  $S$  be the data set under consideration. First note that by definition, we only need to consider the  
662 points  $x$  in the disagreement region. The labels of all points that are not in the subtree of the deepest  
663 1 labeled point in  $S$  are fixed and thus cannot be in the disagreement region. Thus, we only need to  
664 consider the points in the subtree of the deepest 1 labeled point in  $S$ .

665 First note that points that have more than  $\alpha$  points in descendant subtree cannot be attacked. This is  
666 because if the point is labelled as a 0 then all its descendants are in the disagreement region. This  
667 remains true even for any points that are added to the data set. Thus,  $\left| \Gamma \left( S \cup A_x \setminus \{x\}, \mathcal{F}^{x \rightarrow 0} \right) \right| > \alpha$ .  
668 This is because in the definition of  $\gamma$  we remove all points labelled 0.

669 For any node  $u$ , denote by  $\text{pos}(v)$  the closest ancestor on the path  $p$  corresponding to the positive  
670 points. We claim that if  $u$  has fewer than  $\alpha$  points in its descendant subtree, then  $u$  is attackable only  
671 if

$$\left| \{v \in S : \text{pos}(u) \preceq \text{pos}(v) \wedge v \text{ is not a descendant of } u\} \right| \leq \alpha \quad (11)$$

672 First note that adding any 0 labelled points to  $S$  as  $A_u$  does not change the number of points in  $\Gamma$ .  
673 Further note that for  $v$  such that  $\text{pos}(u) \preceq \text{pos}(v)$ , adding a 1 labelled point must be on the path  
674 between  $\text{pos}(u)$  and  $\text{pos}(v)$ . But, this would remove  $u$  from the disagreement region. Thus, all point  
675  $v$  would be counted in  $\Gamma$  and if these are greater than  $\alpha$  then  $u$  is not attackable.

676 Consider the point  $w$  that is minimal amongst  $\text{pos}(u)$  for  $u$  satisfying (11) and  $u$  be a node such that  
677  $\text{pos}(u) = w$ . First note that all points that satisfy (11) are in the subtree of  $w$ . Second, note that this  
678 subtree has at most  $\alpha$  points. This is because if there were more than  $\alpha$  points,  $u$  would not satisfy  
679 (11) which is a contradiction.  $\square$

680 *Proof of Lemma 5.3.* Let  $i_T$  be the number of i.i.d. points in the data set at time  $T$ . Note that the only  
681 i.i.d. points that we abstain on are the attackable points. But, since the i.i.d. points are exchangeable  
682 if  $i$  i.i.d. points are seen so far, the probability of abstain is given by  $\frac{\alpha}{i}$ . Thus, the expected total

683 number of abstentions is at most

$$\sum_{i=1}^{i_T} \frac{\alpha}{i} \leq \alpha \log T$$

684 as required. □

## 685 C Discussion

686 In this section, we further discuss the proposed algorithms and propose potential ways to extend these  
687 beyond the current settings.

### 688 C.1 Generalizing Higher-order Disagreement to Unknown Distribution

689 In Section 4, we saw an algorithm that for any class  $\mathcal{F}$  with VC dimension  $d$ , achieves abstention  
690 error and misclassification error bounded only as a function of  $d$ . As mentioned earlier, an interesting  
691 open question is to extend this to the unknown distribution case. In this section, we will briefly  
692 discuss a natural algorithm extending the algorithm from the known distribution case. Recall that  
693 Algorithm 1 computed the probabilities of shattering  $k$  points using the knowledge of the distribution  
694 and made a prediction  $\hat{y}_t$  depending on the relative magnitudes of the probabilities corresponding  
695 to the two restricted classes. The main challenge in the unknown distribution case is that it is not  
696 immediately obvious how to compute these quantities.

697 One natural approach is to use the historical data as a proxy for the distribution. That is, given the  
698 data set  $S_t$  of size  $n$ , compute the leave- $k$ -out estimator for the probability as follows

$$\tilde{\rho}_k(S, \mathcal{F}) = \frac{1}{\binom{n}{k}} \sum_{T \subset S; |T|=k} \mathbb{1} [T \text{ is shattered by } \mathcal{F}|_{S \setminus T}].$$

699 There are a few things to observe about this estimator. First, in the case when the data is generated  
700 i.i.d., this estimator is unbiased. Further, though each of the summands is not independent, one  
701 can show concentration for estimators of this form using the theory of U-statistics. Additionally,  
702 recall that in Algorithm 1 required the thresholds  $\alpha_k$  to be set to  $T^{-O(k)}$  (we use  $T^{-k}$  but it is  
703 straightforward to extend this to  $T^{-ck}$  for  $c < 1$ ). This appears to be high precision but note that  
704 the “number of samples” one has for a data set of size  $n$  is  $n^{O(k)}$ . Thus, it is conceivable that such  
705 an estimator can give the necessary bounds. Unfortunately, the challenge with analyzing this in our  
706 setting is that we do not know which of the examples are corrupted. Thus, the adversary could inject  
707 examples to make our estimates arbitrarily inaccurate. Thus, as in the case of the VC dimension 1  
708 classes we saw in Section 5 and the case of the rectangles we will see subsequently, our analysis  
709 would need to not rely on the accuracy of the estimates but rather use these estimates to maintain  
710 progress or construct other versions of estimators that are unaffected by the adversarial injections.

### 711 C.2 Generalizing Structure-based Algorithm beyond VC Dimension 1 Classes

712 In Section 5 we saw an algorithm that for any class  $\mathcal{F}$  with VC dimension 1, achieves abstention  
713 and misclassification error bounded by  $O(\sqrt{T})$  without access to the underlying distribution. An  
714 interesting open question is to extend structure-based algorithms beyond the VC 1 classes. Here  
715 we will show that for the class of axis-aligned rectangles in dimension  $p$  (VC dimension is  $2p$ ),  
716 we can indeed design an algorithm that achieves abstention and misclassification error bounded  
717 by  $O(p\sqrt{T \log T})$ . This exhibits a class of VC dimension  $> 1$  for which we can attain the desired  
718 guarantees without access to the distribution.

719 Both our algorithms heavily utilize the structure of the underlying function class and analyze based  
720 on the notion of attackability. It would be interesting to characterize other structural hypothesis  
721 classes that enjoy similar guarantees. A natural extension to the axis-aligned rectangles would be any  
722 intersection closed hypothesis class.

723 **C.2.1 Structure-based Algorithm for Axis-aligned rectangles**

724 Recall that the class of axis-aligned rectangles consists of functions  $f$  parameterized by  
 725  $(a_1, b_1, \dots, a_p, b_p)$  such that

$$f_{(a_1, b_1, \dots, a_p, b_p)}(x) = \begin{cases} 1 & \text{if } \forall i \in [p], a_i \leq x_i \leq b_i \\ 0 & \text{otherwise.} \end{cases}$$

726 Now consider the following algorithm:

---

**Algorithm 3:** Structure-based learning for Axis-aligned Rectangles

---

Set  $a_1, \dots, a_p = -\infty$  and  $b_1, \dots, b_p = \infty$

**for**  $t = 1, \dots, T$  **do**

    Receive  $\hat{x}_t$

**if**  $\forall \tau < t, y_\tau = 0$  **then**  $\hat{y}_t = 0$

**else if**  $\hat{x}_t \notin \mathcal{S}_1(\mathcal{F}_{t-1})$  **then**  $\hat{y}_t = f(\hat{x}_t)$  for any  $f \in \mathcal{F}_{t-1}$

727 **else if**  $\exists s_1, \dots, s_\alpha < t$  and  $\exists i_1, \dots, i_\alpha \in [p]$  such that  $\hat{x}_{s_j, i_j} \in [\hat{x}_{t, i_j}, a_{i_j}] \cup (b_{i_j}, \hat{x}_{t, i_j}]$  **then**

$\hat{y}_t = 0$

**else**  $\hat{y}_t = \perp$

    Receiving label  $y_t$

    Update  $S_t \leftarrow S_{t-1} \cup \{(\hat{x}_t, y_t)\}$  and  $\mathcal{F}_t \leftarrow \mathcal{F}_{t-1}^{\hat{x}_t \rightarrow y_t}$

    Update  $a_1, \dots, a_p$  and  $b_1, \dots, b_p$  such that  $[a_1, b_1] \times \dots \times [a_p, b_p]$  is the smallest rectangle containing the points labelled positive so far. i.e.  $a_i = \min\{x_i : (x, 1) \in S\}$

---

728 If we have only seen 0 labels so far, the algorithm predicts 0. Otherwise, let  $[a_1, b_1] \times [a_2, b_2] \times \dots \times$   
 729  $[a_p, b_p]$  be the minimal rectangle enclosing the positive examples so far (i.e., the Closure hypothesis).  
 730 For the next point  $\hat{x}_t$ , if it isn't in the region of disagreement, we predict the agreed-upon label.  
 731 Otherwise, we check whether there exist at least  $\alpha$  examples  $\hat{x}_s, s < t$ , for each of which there exists  
 732 a coordinate  $i$  with  $x_{s,i} \in [x_{t,i}, a_i] \cup (b_i, x_{t,i}]$  (at most one of these sides is non-empty – or sometimes  
 733 both sides will be empty for some coordinates  $i$ ). If so, we predict 0. Otherwise, we abstain. (The  
 734 algorithm never predicts 1 in the region of disagreement, similar to the Closure algorithm). We  
 735 formally capture the guarantees of the algorithm below:

736 **Theorem C.1.** *Let  $\mathcal{F}$  be the class of axis aligned rectangles in  $\mathbb{R}^p$ . Then, Algorithm 3 with  $\alpha =$   
 737  $\sqrt{T}/\log T$  satisfies*

$$\begin{aligned} \text{MisclassificationError} &\leq p\sqrt{T \log T}, \\ \text{AbstentionError} &\leq 2p\sqrt{T \log T} + 2p \log T. \end{aligned}$$

738 *Proof.* For any example with target label 1, if the algorithm predicts 0, there are at least  $\alpha$  examples  
 739  $\hat{x}_s$  which each have some coordinate  $\hat{x}_{s,i}$  that was not in  $[a_i, b_i]$  before the update, but which will have  
 740  $\hat{x}_{s,i}$  in  $[a_i, b_i]$  after the update. For each example  $\hat{x}_s$ , this can happen at most  $p$  times (corresponding  
 741 to each coordinate) before it will never again be included in a future set of  $\alpha$  examples that convince  
 742 us to predict 0. So we make at most  $pT/\alpha$  misclassifications on adversarially injected examples.  
 743 Since the algorithm never predicts 1 unless the true label is 1, we never misclassify a negative example.  
 744 So it remains only to bound the number of abstentions on the i.i.d. examples.

745 For any  $n$ , suppose  $\hat{x}_t$  is the  $n$ -th i.i.d. example, and let  $\tilde{x}_1, \dots, \tilde{x}_n$  be these  $n$  i.i.d. examples (so  
 746  $\tilde{x}_n = x_t$ ). If  $x_t$  is "attackable" (same to Definition B.2 meaning that there is some set of examples  
 747 the adversary could add, knowing  $\hat{x}_t$ , to make us abstain) then it must be that either  $\hat{x}_t$  is a positive  
 748 example in the region of disagreement of the version space induced by the other  $n - 1$  points, or  
 749 else  $\hat{x}_t$  is a negative example such that there are  $< \alpha$  points  $\tilde{x}_s, s < n$ , for which there exists  $i$  with  
 750  $\tilde{x}_{s,i} \in [x_{t,i}, a_i] \cup (b_i, x_{t,i}]$ . In particular, in the latter case, it must be that each coordinate  $i$  has  $< \alpha$   
 751 examples  $\tilde{x}_s$  with  $\tilde{x}_{s,i} \in [x_{t,i}, a_i^*] \cup (b_i^*, x_{t,i}]$ , where the target concept is  $[a_1^*, b_1^*] \times \dots \times [a_p^*, b_p^*]$ .  
 752 This is because the current estimated rectangle will be inside the true rectangle.

753 We will use exchangeability to bound the probability that  $\tilde{x}_n$  is attackable by  $\frac{1}{n}$  times the number  
 754 of  $\tilde{x}_s, s \leq n$ , which would be attackable if they were swapped with  $\tilde{x}_n$ . Among  $\tilde{x}_1, \dots, \tilde{x}_n$  there  
 755 are at most  $2p$  positive examples in the region of disagreement of the version space induced by the  
 756 others (namely, the minimum spanning set of the positive examples). For each coordinate  $i$ , there are

757 at most  $2\alpha$  examples  $\tilde{x}_s, s \leq n$ , with  $< \alpha$  other examples  $\tilde{x}_{s',i}$  in  $[\tilde{x}_{s,i}, a_i^*) \cup (b_i^*, \tilde{x}_{s,i}]$  (namely, the  
758  $\leq \alpha$  examples with smallest  $\tilde{x}_{s,i}$  such that  $\tilde{x}_{s,i} > b_i^*$ , and the  $\leq \alpha$  examples with largest  $\tilde{x}_{s,i}$  such  
759 that  $\tilde{x}_{s,i} < a_i^*$ ). So there are at most  $2\alpha p$  negative examples  $\tilde{x}_s$  which would be attackable if they  
760 were swapped with  $\tilde{x}_n$ .

761 Altogether there are at most  $2p(\alpha + 1)$  examples  $\tilde{x}_s$  which would be attackable if they were swapped  
762 with  $\tilde{x}_n$ . Thus, the probability i.i.d. example  $\hat{x}_t$  is attackable is at most  $2p(\alpha + 1)/n$  where  $n$  is  
763 the number of i.i.d. points seen so far including  $\hat{x}_t$ . Summing, the expected number of abstentions  
764 on i.i.d. examples is at most  $2p(\alpha + 1) \log T$ . Now setting  $\alpha = \sqrt{T/\log T}$ , gives us the desired  
765 result.  $\square$