
Private Everlasting Prediction

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 A private learner is trained on a sample of labeled points and generates
2 a hypothesis that can be used for predicting the labels of newly sampled
3 points while protecting the privacy of the training set [Kasiviswanathan
4 et al., FOCS 2008]. Research uncovered that private learners may need to
5 exhibit significantly higher sample complexity than non-private learners
6 as is the case with, e.g., learning of one-dimensional threshold functions
7 [Bun et al., FOCS 2015, Alon et al., STOC 2019].

8 We explore prediction as an alternative to learning. Instead of putting
9 forward a hypothesis, a predictor answers a stream of classification queries.
10 Earlier work has considered a private prediction model with just a single
11 classification query [Dwork and Feldman, COLT 2018]. We observe that
12 when answering a stream of queries, a predictor must modify the hypothesis
13 it uses over time, and, furthermore, that it must use the queries for this
14 modification, hence introducing potential privacy risks with respect to the
15 queries themselves.

16 We introduce *private everlasting prediction* taking into account the privacy
17 of both the training set *and* the (adaptively chosen) queries made to the
18 predictor. We then present a generic construction of private everlasting
19 predictors in the PAC model. The sample complexity of the initial training
20 sample in our construction is quadratic (up to polylog factors) in the VC
21 dimension of the concept class. Our construction allows prediction for
22 all concept classes with finite VC dimension, and in particular threshold
23 functions with constant size initial training sample, even when considered
24 over infinite domains, whereas it is known that the sample complexity
25 of privately learning threshold functions must grow as a function of the
26 domain size and hence is impossible for infinite domains.

27 1 Introduction

28 A PAC learner is given labeled examples $S = \{(x_i, y_i)\}_{i \in [n]}$ drawn i.i.d. from an unknown
29 underlying probability distribution \mathcal{D} over a data domain X and outputs a hypothesis h
30 that can be used for predicting the label of fresh points x_{n+1}, x_{n+2}, \dots sampled from the same
31 underlying probability distribution \mathcal{D} [Valiant, 1984]. It is well known that when points
32 are labeled by a concept selected from a concept class $C = \{c : X \rightarrow \{0, 1\}\}$ then learning is
33 possible with sample complexity proportional to the VC dimension of the concept class.

34 Learning often happens in settings where the underlying training data is related to indi-
35 viduals and privacy-sensitive and where a learner is required, for legal, ethical, or other
36 reasons, to protect personal information from being leaked in the learned hypothesis h .
37 Private learning was introduced by Kasiviswanathan et al. [2011], as a theoretical model for

38 studying such tasks. A *private learner* is a PAC learner that preserves differential privacy
 39 with respect to its training set S . That is, the learner’s distribution on outcome hypotheses
 40 must not depend too strongly on any single example in S . Kasiviswanathan et al. showed
 41 via a generic construction that any finite concept class can be learned privately and with
 42 sample complexity $n = O(\log|C|)$. This value ($O(\log|C|)$) can be significantly higher than the
 43 VC dimension of the concept class C (see below).

44 It is now understood that the gap between the sample complexity of private and non-private
 45 learners is essential – an important example is private learning of threshold functions
 46 (defined over an ordered domain X as $C_{\text{thresh}} = \{c_t\}_{t \in X}$ where $c_t(x) = \mathbb{1}_{x \geq t}$), which requires
 47 sample complexity that is asymptotically higher than the (constant) VC dimension of C_{thresh} .
 48 In more detail, with *pure* differential privacy, the sample complexity of private learning is
 49 characterized by the representation dimension of the concept class [Beimel et al., 2013a].
 50 The representation dimension of C_{thresh} (hence, the sample complexity of private learning
 51 thresholds) is $\Theta(\log|X|)$ [Feldman and Xiao, 2015]. With *approximate* differential privacy,
 52 the sample complexity of learning threshold functions is $\Theta(\log^*|X|)$ [Beimel et al., 2013b,
 53 Bun et al., 2015, Alon et al., 2019, Kaplan et al., 2020, Cohen et al., 2022]. Hence, in
 54 both the pure and approximate differential privacy cases, the sample complexity grows
 55 with the cardinality of the domain $|X|$ and no private learner exists for threshold functions
 56 over infinite domains, such as the integers and the reals, whereas low sample complexity
 57 non-private learners exist for these tasks.

58 **Privacy preserving (black-box) prediction.** Dwork and Feldman [2018] proposed privacy-
 59 preserving prediction as an alternative for private learning. Noting that “[i]t is now known
 60 that for some basic learning problems [. . .] producing an accurate private model requires
 61 much more data than learning without privacy,” they considered a setting where “users
 62 may be allowed to query the prediction model on their inputs only through an appropriate
 63 interface”. That is, a setting where the learned hypothesis is not made public. Instead, it may
 64 be accessed in a “black-box” manner via a privacy-preserving query-answering prediction
 65 interface. The prediction interface is required to preserve the privacy of its training set S :
 66 **Definition 1.1** (private prediction interface [Dwork and Feldman, 2018] (rephrased)). A
 67 prediction interface M is (ϵ, δ) -differentially private if for every interactive query generating
 68 algorithm Q , the output of the interaction between Q and $M(S)$ is (ϵ, δ) -differentially private
 69 with respect to S .

70 Dwork and Feldman focused on the setting where the entire interaction between Q and
 71 $M(S)$ consists of issuing a single prediction query and answering it:

72 **Definition 1.2** (Single query prediction [Dwork and Feldman, 2018]). Let M be an algorithm
 73 that given a set of labeled examples S and an unlabeled point x produces a label y . M
 74 is an (ϵ, δ) -differentially private prediction algorithm if for every x , the output $M(S, x)$ is
 75 (ϵ, δ) -differentially private with respect to S .

76 W.r.t. answering a single prediction query, Dwork and Feldman showed that the sample
 77 complexity of such predictors is proportional to the VC dimension of the concept class.

78 1.1 Our contributions

79 In this work, we extend private prediction beyond a single query to answering any sequence
 80 – *unlimited in length* – of prediction queries. We refer to this as *private everlasting prediction*.
 81 Our goal is to present a generic private everlasting predictor with low training sample
 82 complexity $|S|$.

83 **Private prediction interfaces when applied to a large number of queries.** We begin by
 84 examining private everlasting prediction under the framework of Definition 1.1. We prove:
 85 **Theorem 1.3** (informal version of Theorem 3.3). *Let A be a private everlasting prediction*
 86 *interface for concept class C and assume A bases its predictions solely on the initial training set S ,*
 87 *then there exists a private learner for concept class C with sample complexity $|S|$.*

88 This means that everlasting predictors that base their prediction solely on the initial training
 89 set S are subject to the same complexity lowerbounds as private learners. Hence, to avoid

90 private learning lowerbounds, private everlasting predictors need to rely on more than
91 the initial training sample S as a source of information about the underlying probability
92 distribution and the labeling concept.

93 In this work, we choose to allow the everlasting predictor to rely on the queries made -
94 which are unlabeled points from the domain X , assuming the queries are drawn from the
95 same distribution the initial training S is sampled from. This requires changing the privacy
96 definition, as Definition 1.1 does not protect the queries made, yet the classification given to
97 a query can now depend on and hence reveal information provided in queries made earlier.

98 **A definition of private everlasting predictors.** Our definition of private everlasting
99 predictors is motivated by the observations above. Consider an algorithm \mathcal{A} that is first
100 fed with a training set S of labeled points and then executes for an unlimited number of
101 rounds, where in round i algorithm \mathcal{A} receives as input a query point x_i and produces
102 a label \hat{y}_i . We say that \mathcal{A} is an everlasting predictor if, when the (labeled) training set S
103 and the (unlabeled) query points are coming from the same underlying distribution, \mathcal{A}
104 answers each query points x_i with a good hypothesis h_i , and hence the label \hat{y}_i produced
105 by \mathcal{A} is correct with high probability. We say that \mathcal{A} is a *private* everlasting predictor if its
106 sequence of predictions $\hat{y}_1, \hat{y}_2, \hat{y}_3, \dots$ protects both the privacy of the training set S and the
107 query points x_1, x_2, x_3, \dots in face of any adversary that adaptively chooses the query points.

108 We emphasize that while private everlasting predictors need to exhibit average-case utility
109 – as good prediction is required only for the case where S and x_1, x_2, x_3, \dots are selected
110 i.i.d. from the same underlying distribution – our privacy requirement is worst-case, and
111 holds in face of an *adaptive* adversary that chooses each query point x_i after receiving the
112 prediction provided for (x_1, \dots, x_{i-1}) , and not necessarily in accordance with any probability
113 distribution.

114 **A generic construction of private everlasting predictors.** Our construction, called
115 `GenericBBL`, executes in rounds. The input to the first round is the initial labeled training
116 set S , where the number of samples in S is quadratic in the VC dimension of the concept
117 class. Each other round begins with a collection S_i of labeled examples and ends with newly
118 generated collection of labeled examples S_{i+1} . The set S is assumed to be consistent with
119 some concept $c \in C$ and our construction ensures that this is the case also for the sets S_i for
120 all i . We briefly describe the main computations performed in each round of `GenericBBL`.¹

- 121 • **Round initialization:** At the outset of a round, the labeled set S_i is partitioned into
122 sub-sets, each with number of samples which is proportional to the VC dimension
123 (so we have $\approx \frac{|S_i|}{\text{VC}(C)}$ sub-sets). Each of the sub-sets is used for training a classifier
124 non-privately, hence creating a collection of classifiers $F_i = \{f : X \rightarrow \{0, 1\}\}$ that are used
125 throughout the round.
- 126 • **Query answering:** Queries are issued to the predictor in an online manner. Each query
127 is first labeled by each of the classifiers in F_i . Then the predicted label is computed by
128 applying a privacy-preserving majority vote on these intermediate labels. (By standard
129 composition theorems for differential privacy, we could answer roughly $|F_i|^2 \approx \left(\frac{|S_i|}{\text{VC}(C)}\right)^2$
130 queries without exhausting our privacy budget.) To save on the privacy budget, the
131 majority vote is based on the `BetweenThresholds` mechanism of [Bun et al. \[2016\]](#)
132 (which in turn is based on the sparse vector technique). The algorithm fails when the
133 privacy budget is exhausted. However, when queries are sampled from the underlying
134 distribution then with a high enough probability the labels produced by the classifiers
135 in F_i would exhibit a clear majority.
- 136 • **Generating a labeled set for the following round:** The predictions provided in the
137 duration of a round are not guaranteed to be consistent with any concept in C and
138 hence cannot be used to set the following round. Instead, at the end of the round these
139 points are relabeled consistently with C using a technique developed by [Beimel et al.](#)

¹Important details, such as privacy amplification via sampling and management of the learning accuracy and error parameters are omitted from the description provided in this section.

140 [2021] in the context of private semi-supervised learning. Let S_{i+1} denote the query
141 points obtained during the i th round, after (re)labeling them. This is a collection of
142 size $|S_{i+1}| \approx \left(\frac{|S_i|}{\sqrt{\text{VC}(C)}}\right)^2$. Hence, provided that $|S_i| \gtrsim (\text{VC}(C))^2$ we get that $|S_{i+1}| > |S_i|$ which
143 allows us to continue to the next round with more data than we had in the previous
144 round.

145 **Theorem 1.4** (informal version of Theorem 5.1). *For every concept class C , Algorithm*
146 *GenericBBL is a private everlasting predictor requiring an initial set of labeled examples which is*
147 *(upto polylogarithmic factors) quadratic in the VC dimension of C .*

148 1.2 Related work

149 Beyond the work of Dwork and Feldman [2018] on private prediction mentioned above, our
150 work is related to private semi-supervised learning and joint differential privacy.

151 **Semi-supervised private learning.** As in the model of private semi-supervised learning of
152 Beimel et al. [2021], our predictors depend on both labeled and unlabeled sample. Beyond
153 the obvious difference between the models (outputting a hypothesis vs. providing black-box
154 prediction), a major difference between the settings is that in the work of Beimel et al. [2021]
155 all samples – labeled and unlabeled - are given at once at the outset of the learning process
156 whereas in the setting of everlasting predictors the unlabeled samples are supplied in an
157 online manner. Our construction of private everlasting predictors uses tools developed for
158 the semi-supervised setting, and in particular Algorithm LabelBoost of of Beimel et al.

159 **Joint differential privacy.** Kearns et al. [2015] introduced joint differential privacy (JDP)
160 as a relaxation of differential privacy applicable for mechanism design and games. For
161 every user u , JDP requires that the outputs jointly seen by all other users would preserve
162 differential privacy w.r.t. the input of u . Crucially, in JDP users select their inputs ahead of
163 the computation. In our settings, the inputs to a private everlasting predictor are prediction
164 queries which are chosen in an online manner, and hence a query can depend on previous
165 queries and their answers. Yet, similarly to JDP, the outputs provided to queries not
166 performed by a user u should jointly preserve differential privacy w.r.t. the query made by
167 u . Our privacy requirement hence extends JDP to an adaptive online setting.

168 **Additional works on private prediction.** Bassily et al. [2018] studied a variant of the
169 private prediction problem where the algorithm takes a labeled sample S and is then
170 required to answer m prediction queries (i.e., label a sequence of m unlabeled points
171 sampled from the same underlying distribution). They presented algorithms for this task
172 with sample complexity $|S| \gtrsim \sqrt{m}$. This should be contrasted with our model and results,
173 where the sample complexity is independent of m . The bounds presented by Dwork and
174 Feldman [2018] and Bassily et al. [2018] were improved by Dagan and Feldman [2020] and
175 by Nandi and Bassily [2020] who presented algorithms with improved dependency on the
176 accuracy parameter in the agnostic setting.

177 1.3 Discussion and open problems

178 We show how to transform any (non-private) learner for the class C (with sample complexity
179 proportional to the VC dimension of C) to a private everlasting predictor for \bar{C} . Our
180 construction is not polynomial time due to the use of Algorithm LabelBoost, and requires
181 an initial set S of labeled examples which is quadratic in the VC dimension. We leave open
182 the question whether $|S|$ can be reduced to be linear in the VC dimension and whether the
183 construction can be made polynomial time. A few remarks are in order:

- 184 1. Even though our generic construction is not computationally efficient, it does re-
185 sult in efficient learners for several interesting special cases. Specifically, algorithm
186 LabelBoost can be implemented efficiently whenever given an input sample S we
187 could efficiently enumerate all possible dichotomies from the target class C over the
188 points in S . In particular, this is the case for the class of 1-dim threshold functions
189 C_{thresh} , as well as additional classes with constant VC dimension. Another notable

190 example is the class C_{thresh}^{enc} which intuitively is an “encrypted” version of C_{thresh} . [Bun](#)
 191 [and Zhandry \[2016\]](#) showed that (under plausible cryptographic assumptions) the
 192 class C_{thresh}^{enc} cannot be learned privately and efficiently, while non-private learning is
 193 possible efficiently. Our construction can be implemented efficiently for this class. This
 194 provides an example where private everlasting prediction can be done efficiently, while
 195 (standard) private learning is possible but inefficient.

- 196 2. It is now known that some learning tasks require the produced model to memorize
 197 parts of the training set in order to achieve good learning rates, which in particular
 198 disallows the learning algorithm from satisfying (standard) differential privacy [[Brown](#)
 199 [et al., 2021](#)]. Our notion of private everlasting prediction circumvents this issue, since
 200 the model is never publicly released and hence the fact that it must memorize parts
 201 of the sample is not of a direct privacy threat. In other words, our work puts forward
 202 a private learning model which, in principle, allows memorization. This could have
 203 additional applications in broader settings.
- 204 3. As we mentioned, in general, private everlasting predictors cannot base their predic-
 205 tions solely on the initial training set, and in this work we choose to rely on the *queries*
 206 presented to the algorithm (in addition to the training set). Our construction can be
 207 easily adapted to a setting where the content of the blackbox is updated based on *fresh*
 208 *unlabeled samples* (whose privacy would be preserved), instead of relying on the query
 209 points themselves. This might be beneficial to avoid poisoning attacks via the queries.

210 2 Preliminaries

211 2.1 Preliminaries from differential privacy

Definition 2.1 ((ϵ, δ) -indistinguishability). Let R_0, R_1 be two random variables over the
 same support. We say that R_0, R_1 are (ϵ, δ) -indistinguishable if for every event E defined
 over the support of R_0, R_1 ,

$$\Pr[R_0 \in E] \leq e^\epsilon \cdot \Pr[R_1 \in E] + \delta \quad \text{and} \quad \Pr[R_1 \in E] \leq e^\epsilon \cdot \Pr[R_0 \in E] + \delta.$$

212 **Definition 2.2.** Let X be a data domain. Two datasets $x, x' \in X^n$ are called *neighboring* if
 213 $|\{i : x_i \neq x'_i\}| = 1$.

214 **Definition 2.3** (differential privacy [[Dwork et al., 2006](#)]). A mechanism $M : X^n \rightarrow Y$ is
 215 (ϵ, δ) -differentially private if $M(x)$ and $M(x')$ are (ϵ, δ) -indistinguishable for all neighboring
 216 $x, x' \in X^n$.

217 In our analysis, we use the post-processing and composition properties of differential
 218 privacy, that we cite in their simplest form.

219 **Proposition 2.4** (post-processing). Let $M_1 : X^n \rightarrow Y$ be an (ϵ, δ) -differentially private algorithm
 220 and $M_2 : Y \rightarrow Z$ be any algorithm. Then the algorithm that on input $x \in X^n$ outputs $M_2(M_1(x))$
 221 is (ϵ, δ) -differentially private.

222 **Proposition 2.5** (composition). Let M_1 be a (ϵ_1, δ_1) -differentially private algorithm and let
 223 M_2 be (ϵ_2, δ_2) -differentially private algorithm. Then the algorithm that on input $x \in X^n$ outputs
 224 $(M_1(x), M_2(x))$ is $(\epsilon_1 + \epsilon_2, \delta_1 + \delta_2)$ -differentially private.

225 **Definition 2.6** (Exponential mechanism [[McSherry and Talwar, 2007](#)]). Let $q : X^n \times Y \rightarrow \mathbb{R}$
 226 be a score function defined over data domain X and output domain Y . Define $\Delta =$
 227 $\max(|q(x, r) - q(x', y)|)$ where the maximum is taken over all $y \in Y$ and neighbouring
 228 databases $x, x' \in X^n$. The exponential mechanism is the ϵ -differentially private mecha-
 229 nism which selects an output $y \in Y$ with probability proportional to $e^{\frac{\epsilon q(x, y)}{2\Delta}}$.

230 **Claim 2.7** (Privacy amplification by sub-sampling [[Kasiviswanathan et al., 2011](#)]). Let \mathcal{A}
 231 be an (ϵ', δ') -differentially private algorithm operating on a database of size n . Let $\epsilon \leq 1$ and
 232 let $t = \frac{n}{\epsilon} (3 + \exp(\epsilon'))$. Construct an algorithm \mathcal{B} operating the database $D = (z_i)_{i=1}^t$. Algorithm
 233 \mathcal{B} randomly selects a subset $J \subseteq \{1, 2, \dots, t\}$ of size n , and executes \mathcal{A} on $D_J = (z_i)_{i \in J}$. Then \mathcal{B} is
 234 $(\epsilon, \frac{4\epsilon'}{3 + \exp(\epsilon')} \delta')$ -differentially private.

235 **2.2 Preliminaries from PAC learning**

236 A concept class C over data domain X is a set of predicates $c : X \rightarrow \{0, 1\}$ (called concepts)
 237 which label points of the domain X by either 0 or 1. A learner \mathcal{A} for concept class C is
 238 given n examples sampled i.i.d. from an unknown probability distribution \mathcal{D} over the data
 239 domain X and labeled according to an unknown target concept $c \in C$. The learner should
 240 output a hypothesis $h : X \rightarrow [0, 1]$ that approximates c for the distribution \mathcal{D} . More formally,

241 **Definition 2.8** (generalization error). The *generalization error* of a hypothesis $h : X \rightarrow [0, 1]$
 242 with respect to concept c and distribution \mathcal{D} is defined as $\text{error}_{\mathcal{D}}(c, h) = \text{Exp}_{x \sim \mathcal{D}}[|h(x) - c(x)|]$.

Definition 2.9 (PAC learning [Valiant, 1984]). Let C be a concept class over a domain X .
 Algorithm \mathcal{A} is an (α, β, n) -PAC learner for C if for all $c \in C$ and all distributions \mathcal{D} on X ,

$$\Pr[(x_1, \dots, x_n) \sim \mathcal{D}^n ; h \sim \mathcal{A}((x_1, c(x_1)), \dots, (x_n, c(x_n))) ; \text{error}_{\mathcal{D}}(c, h) \leq \alpha] \geq 1 - \beta,$$

243 where the probability is over the sampling of (x_1, \dots, x_n) from \mathcal{D} and the coin tosses of \mathcal{A} .
 244 The parameter n is the *sample complexity* of \mathcal{A} .

245 See Appendix A for additional preliminaries on PAC learning.

246 **2.3 Preliminaries from private learning**

247 **Definition 2.10** (private PAC learning [Kasiviswanathan et al., 2011]). Algorithm \mathcal{A} is
 248 a $(\alpha, \beta, \epsilon, \delta, n)$ -private PAC learner if (i) \mathcal{A} is an (α, β, n) -PAC learner and (ii) \mathcal{A} is (ϵ, δ)
 249 differentially private.

250 Kasiviswanathan et al. [2011] provided a generic private learner with $O(\text{VC}(C) \log(|X|))$
 251 labeled samples. Beimel et al. [2013a] introduced the representation dimension and showed
 252 that any concept class C can be privately learned with $\Theta(\text{RepDim}(C))$ samples.² For the
 253 sample complexity of (ϵ, δ) -differentially private learning of threshold functions over do-
 254 main X , Bun et al. [2015] give a lower bound of $\Omega(\log^* |X|)$. Recently, Cohen et al. [2022]
 255 give a (nearly) matching upper bound of $\tilde{O}(\log^* |X|)$.

256 **3 Towards private everlasting prediction**

257 In this work, we extend private prediction beyond a single query to answering any sequence
 258 – unlimited in length – of prediction queries. Our goal is to present a generic private
 259 everlasting predictor with low training sample complexity $|S|$.

260 **Definition 3.1** (everlasting prediction). Let \mathcal{A} be an algorithm with the following properties:

- 261 1. Algorithm \mathcal{A} receives as input n labeled examples $S = \{(x_i, y_i)\}_{i=1}^n \in (X \times \{0, 1\})^n$ and
 262 selects a hypothesis $h_0 : X \rightarrow \{0, 1\}$.
- 263 2. For round $r \in \mathbb{N}$, algorithm \mathcal{A} gets a query, which is an unlabeled element $x_{n+r} \in X$,
 264 outputs $h_{r-1}(x_{n+r})$ and selects a hypothesis $h_r : X \rightarrow \{0, 1\}$.

265 We say that \mathcal{A} is an (α, β, n) -*everlasting predictor* for a concept class C over a domain X if the
 266 following holds for every concept $c \in C$ and for every distribution \mathcal{D} over X . If x_1, x_2, \dots are
 267 sampled i.i.d. from \mathcal{D} , and the labels of the n initial samples S are correct, i.e., $y_i = c(x_i)$ for
 268 $i \in [n]$, then $\Pr[\exists r \geq 0$ s.t. $\text{error}_{\mathcal{D}}(c, h_r) > \alpha] \leq \beta$, where the probability is over the sampling
 269 of x_1, x_2, \dots from \mathcal{D} and the randomness of \mathcal{A} .

270 **Definition 3.2.** An algorithm \mathcal{A} is an $(\alpha, \beta, \epsilon, \delta, n)$ -*everlasting differentially private predic-*
 271 *tion interface* if (i) \mathcal{A} is a (ϵ, δ) -differentially private prediction interface M (as in Defini-
 272 tion 1.1), and (ii) \mathcal{A} is an (α, β, n) -*everlasting predictor*.

273 As a warmup, consider an $(\alpha, \beta, \epsilon, \delta, n)$ -*everlasting differentially private prediction interface*
 274 \mathcal{A} for concept class C over (finite) domain X (as in Definition 3.2 above). Assume that \mathcal{A} does
 275 not vary its hypotheses, i.e. (in the language of Definition 3.1) $h_r = h_0$ for all $r > 0$.³ Note

²We omit the dependency on $\epsilon, \delta, \alpha, \beta$ in this brief review.

³Formally, \mathcal{A} can be thought of as two mechanisms (M_0, M_1) where M_0 is (ϵ, δ) -differentially private. (i) On input a labeled training sample S mechanism M_0 computes a hypothesis h_0 . (ii) On a query $x \in X$ mechanism M_1 replies $h_0(x)$.

276 that a computationally unlimited adversarial querying algorithm can recover the hypothesis
 277 h_0 by issuing all queries $x \in X$. Hence, in using \mathcal{A} indefinitely we lose any potential benefits
 278 to sample complexity of restricting access to h_0 to being black-box and getting to the point
 279 where the lower-bounds on n from private learning apply. A consequence of this simple
 280 observation is that a private everlasting predictor cannot answer all prediction queries with
 281 a single hypothesis – it must modify its hypothesis over time as it processes new queries.

282 We now take this observation a step further, showing that a private everlasting predictor
 283 that answers prediction queries solely based on its training sample S is subject to the same
 284 sample complexity lowerbounds as private learners.

285 Consider an $(\alpha, \beta < 1/8, \epsilon, \delta, n)$ -everlasting differentially private prediction interface \mathcal{A} for
 286 concept class C over (finite) domain X that upon receiving the training set $S \in (X \times \{0, 1\})^n$
 287 selects an infinite sequence of hypotheses $\{h_r\}_{r \geq 0}$ where $h_r : X \rightarrow \{0, 1\}$. Formally, we
 288 can think of \mathcal{A} as composed of three mechanisms $\mathcal{A} = (M_0, M_1, M_2)$ where M_0 is (ϵ, δ) -
 289 differentially private:

- 290 • On input a labeled training sample $S \in (X \times \{0, 1\})^n$ mechanism M_0 computes an
 291 initial state and an initial hypothesis $(\sigma_0, h_0) = M_0(S)$.
- 292 • On a query x_{n+r} mechanism M_1 produces an answer $M_1(x_{n+r}) = h_i(x_{n+r})$ and mech-
 293 anism M_2 updates the hypothesis-state pair $(h_{r+1}, \sigma_{r+1}) = M_2(\sigma_r)$.

294 Note that as M_0 and M_2 do not receive the sequence $\{x_{n+r}\}_{r \geq 0}$ as input, the sequence $\{h_r\}_{r \geq 0}$
 295 depends solely on S . Furthermore as M_1 and M_2 post-process the outcome of M_0 , i.e., the
 296 sequence of queries and predictions $\{(x_r, h_r(x_r))\}_{r \geq 0}$ preserves (ϵ, δ) -differential privacy with
 297 respect to the training set S . In Appendix B we prove:

298 **Theorem 3.3.** *\mathcal{A} can be transformed into a $(O(\alpha), O(\beta), \epsilon, \delta, O(n \log(1/\beta)))$ -private PAC learner
 299 for C .*

300 3.1 A definition of private everlasting prediction

301 Theorem 3.3 requires us to seek private predictors whose prediction relies on more infor-
 302 mation than what is provided by the initial labeled sample. Possibilities include requiring
 303 the input of additional labeled or unlabeled examples during the lifetime of the predictor,
 304 while protecting the privacy of these examples. In this work we choose to rely on the queries
 305 for updating the predictor’s internal state. This introduces a potential privacy risk for these
 306 queries as sensitive information about a query may be leaked in the predictions following it.
 307 Furthermore, we need take into account that a privacy attacker may choose their queries
 308 adversarially and adaptively.

309 **Definition 3.4** (private everlasting black-box prediction). An algorithm \mathcal{A} is an $(\alpha, \beta, \epsilon, \delta, n)$ -
 310 private everlasting black-box predictor for a concept class C if

- 311 1. **Prediction:** \mathcal{A} is an (α, β, n) -everlasting predictor for C (as in Definition 3.1).
- 312 2. **Privacy:** For every adversary \mathcal{B} and every $t \geq 1$, the random variables $\text{View}_{\mathcal{B}, t}^0$ and
 313 $\text{View}_{\mathcal{B}, t}^1$ (defined in Figure 1) are (ϵ, δ) -indistinguishable.

314 4 Tools from prior works

315 We briefly describe tools from prior works that we use in our construction. See Appendix C
 316 for a more detailed account.

317 **Algorithm LabelBoost [Beimel et al., 2021]:** Algorithm LabelBoost takes as input a
 318 partially labeled database $S \circ T \in (X \times \{0, 1, \perp\})^*$ (where the first portion of the database,
 319 S , contains labeled examples) and outputs a similar database where both S and T are
 320 (re)labeled. We use the following lemmata from Beimel et al. [2021]:

321 **Lemma 4.1** (privacy of Algorithm LabelBoost). *Let \mathcal{A} be an (ϵ, δ) -differentially private al-
 322 gorithm operating on labeled databases. Construct an algorithm \mathcal{B} that on input a partially*

Parameters: $b \in \{0, 1\}$, $t \in \mathbb{N}$.

Training Phase:

1. The adversary \mathcal{B} chooses two sets of n labeled elements $(x_1^0, y_1^0), \dots, (x_n^0, y_n^0)$ and $(x_1^1, y_1^1), \dots, (x_n^1, y_n^1)$, subject to the restriction $\left| \left\{ i \in [n] : (x_i^0, y_i^0) \neq (x_i^1, y_i^1) \right\} \right| \in \{0, 1\}$.
2. If $\exists i$ s.t. $(x_i^0, y_i^0) \neq (x_i^1, y_i^1)$ then set Flag = 1. Otherwise set Flag = 0.
3. Algorithm \mathcal{A} gets $(x_1^b, y_1^b), \dots, (x_n^b, y_n^b)$ and selects a hypothesis $h_0 : X \rightarrow \{0, 1\}$.
 \setminus * the adversary \mathcal{B} does not get to see the hypothesis h_0 * \setminus

Prediction phase:

4. For round $r = 1, 2, \dots, t$:
 - (a) If Flag = 1 then the adversary \mathcal{B} chooses two elements $x_{n+r}^0 = x_{n+r}^1 \in X$.
Otherwise, the adversary \mathcal{B} chooses two elements $x_{n+r}^0, x_{n+r}^1 \in X$.
 - (b) If $x_{n+r}^0 \neq x_{n+r}^1$ then Flag is set to 1.
 - (c) If $x_{n+r}^0 = x_{n+r}^1$ then the adversary \mathcal{B} gets $h_{r-1}(x_{n+r}^b)$.
 \setminus * the adversary \mathcal{B} does not get to see the label if $x_{n+r}^0 \neq x_{n+r}^1$ * \setminus
 - (d) Algorithm \mathcal{A} gets x_{n+r}^b and selects a hypothesis $h_r : X \rightarrow \{0, 1\}$.
 \setminus * the adversary \mathcal{B} does not get to see the hypothesis h_r * \setminus

Let $\text{View}_{\mathcal{B}, t}^b$ be \mathcal{B} 's entire view of the execution, i.e., the adversary's randomness and the sequence of predictions in Step 4c.

Figure 1: Definition of $\text{View}_{\mathcal{B}, t}^0$ and $\text{View}_{\mathcal{B}, t}^1$.

323 labeled database $S \circ T \in (X \times \{0, 1, \perp\})^*$ applies \mathcal{A} on the outcome of $\text{LabelBoost}(S \circ T)$. Then, \mathcal{B}
324 is $(\epsilon + 3, 4\epsilon\delta)$ -differentially private.

325 **Lemma 4.2** (Utility of Algorithm LabelBoost). Fix α and β , and let $S \circ T$ be s.t. S is labeled
326 by some target concept $c \in C$, and s.t. $|T| \leq \frac{\beta}{\epsilon} \text{VC}(C) \exp(\frac{\alpha|S|}{2\text{VC}(C)}) - |S|$. Consider the execution
327 of LabelBoost on $S \circ T$, and let h denote the hypothesis chosen by LabelBoost to relabel $S \circ T$.
328 With probability at least $(1 - \beta)$ we have that $\text{error}_S(h) \leq \alpha$.

329 **Algorithm BetweenThresholds** [Bun et al., 2016]: Algorithm BetweenThresholds takes
330 as input a database $S \in X^n$ and thresholds t_ℓ, t_u . It applies the sparse vector technique to
331 answer noisy threshold queries with L (below threshold) R (above threshold) and \top (halt).
332 We use the following lemmata by Bun et al. [2016] and observe that, using standard privacy
333 amplification theorems, Algorithm BetweenThresholds can be modified to allow for c times
334 of outputting \top before halting, with a (roughly) \sqrt{c} growth in its privacy parameter.

335 **Lemma 4.3** (Privacy for BetweenThresholds). Let $\epsilon, \delta \in (0, 1)$ and $n \in \mathbb{N}$. Then algorithm
336 BetweenThresholds is (ϵ, δ) -differentially private for any adaptively-chosen sequence of queries
337 as long as the gap between the thresholds t_ℓ, t_u satisfies $t_u - t_\ell \geq \frac{12}{\epsilon n} (\log(10/\epsilon) + \log(1/\delta) + 1)$.

338 **Lemma 4.4** (Accuracy of BetweenThresholds). Let $\alpha, \beta, \epsilon, t_\ell, t_u \in (0, 1)$ and $n, k \in \mathbb{N}$ satisfy
339 $n \geq \frac{8}{\alpha\epsilon} (\log(k+1) + \log(1/\beta))$. Then, for any input $x \in X^n$ and any adaptively-chosen sequence
340 of queries q_1, q_2, \dots, q_k , the answers $a_1, a_2, \dots, a_{\leq k}$ produced by BetweenThresholds on input x
341 satisfy the following with probability at least $1 - \beta$. For any $j \in [k]$ such that a_j is returned before
342 BetweenThresholds halts, (i) $a_j = L \implies q_j(x) \leq t_\ell + \alpha$, (ii) $a_j = R \implies q_j(x) \geq t_u - \alpha$, and (iii)
343 $a_j = \top \implies t_\ell - \alpha \leq q_j(x) \leq t_u + \alpha$.

344 **Observation 1.** Using standard composition theorems for differential privacy (see, e.g., Dwork
345 et al. [2010]), we can assume that algorithm BetweenThresholds takes another parameter c ,
346 and halts after c times of outputting \top . In this case, the algorithm satisfies $(\epsilon', 2c\delta)$ -differential
347 privacy, for $\epsilon' = \sqrt{2c \ln(\frac{1}{c\delta})} \epsilon + c\epsilon(e^\epsilon - 1)$.

348 **5 A Generic Construction**

349 Our generic construction Algorithm `GenericBBL` transforms a (non-private) learner for
 350 a concept class C into a private everlasting predictor for C . The proof of the following
 351 theorem follows from Theorem 5.2 and Claim 5.3 which are proved in Appendix E.

352 **Theorem 5.1.** *Given $\alpha, \beta, \delta < 1/16, \epsilon < 1$, Algorithm `GenericBBL` is a $(6\alpha, 4\beta, \epsilon, \delta, n)$ -private*
 353 *everlasting predictor, where n is set as in Algorithm `GenericBBL`.*

Algorithm `GenericBBL`

Initial input: A labeled database $S \in (X \times \{0, 1\})^n$ where $n = \frac{8\tau}{\alpha^3 \epsilon^2} \cdot \left(8VC(C) \log\left(\frac{26}{\alpha}\right) + 4 \log\left(\frac{4}{\beta}\right)\right)^2 \cdot \log\left(\frac{1}{\delta}\right) \cdot \log^2\left(\frac{64VC(C) \log\left(\frac{26}{\alpha}\right) + 32 \log\left(\frac{4}{\beta}\right)}{\epsilon \alpha^2 \beta \delta}\right) \cdot (3 + \exp(\epsilon + 4))$.

1. Let $\tau > 1.1 * 10^{10}$. Set $\alpha_1 = \alpha/2, \beta_1 = \beta/2$. Define $\lambda_i = \frac{8VC(C) \log\left(\frac{13}{\alpha_i}\right) + 4 \log\left(\frac{2}{\beta_i}\right)}{\alpha_i}$.
 /* by Theorem A.2 λ_i samples suffice for PAC learning C with parameters α_i, β_i */
 2. Let $S_1 \subseteq S$ be a random subset of size $n \cdot \frac{\epsilon}{3 + \exp(\epsilon + 4)} = \frac{\tau \cdot \lambda_1^2 \cdot \log\left(\frac{1}{\delta}\right) \cdot \log^2\left(\frac{\lambda_1}{\epsilon \alpha_1 \beta_1 \delta}\right)}{\alpha_1 \epsilon}$.
 3. Repeat for $i = 1, 2, 3, \dots$
 - (a) Divide S_i into $T_i = \frac{\tau \cdot \lambda_i \cdot \log\left(\frac{1}{\delta}\right) \cdot \log^2\left(\frac{\lambda_i}{\epsilon \alpha_i \beta_i \delta}\right)}{\alpha_i \epsilon}$ disjoint databases $S_{i,1}, \dots, S_{i,T_i}$ of size λ_i .
 - (b) For $t \in [T_i]$ let $f_t \in C$ be a hypothesis minimizing error $_{S_{i,t}}(\cdot)$. Define $F_i = (f_1, \dots, f_{T_i})$.
 - (c) Set $R_i = \frac{25600|S_i|}{\epsilon}$. Set $t_u = 1/2 + \alpha_i, t_\ell = 1/2 - \alpha_i$. Set the privacy parameters $\epsilon'_i = \frac{1}{3\sqrt{c_i \ln\left(\frac{2}{\delta}\right)}}$ and $\delta'_i = \frac{\delta}{2c_i}$, where $c_i = 64\alpha_i R_i$. Instantiate algorithm `BetweenThresholds` on the database of hypotheses F_i allowing for $c_i = 64\alpha_i R_i$ rounds of Υ while satisfying $(1, \delta)$ -differential privacy (as in Observation 2).
 - (d) For $\ell = 1$ to R_i :
 - i. Receive as input a prediction query $x_{i,\ell} \in X$.
 - ii. Give `BetweenThresholds` the query $q_{x_{i,\ell}}$ where $q_{x_{i,\ell}}(F_i) = \sum_{t \in [T_i]} f_t(x_{i,\ell})$, and obtain an outcome $y_{i,\ell} \in \{L, \Upsilon, R\}$.
 - iii. Respond with the label 0 if $y_{i,\ell} = L$ and 1 if $y_{i,\ell} \in \{R, \Upsilon\}$.
 - iv. If `BetweenThresholds` halts, then halt and fail (recall that `BetweenThresholds` only halts if c_i copies of Υ were encountered during the current iteration).
 - (e) Denote $D_i = (x_{i,1}, \dots, x_{i,R_i})$.
 - (f) Let $\hat{S}_i \subseteq S_i$ and $\hat{D}_i \subseteq D_i$ be random subsets of size $\frac{\epsilon |S_i|}{3 + \exp(\epsilon + 4)}$ and $\frac{\epsilon |D_i|}{3 + \exp(\epsilon + 4)}$ respectively, and let $\hat{S}'_i \circ \hat{D}'_i \leftarrow \text{LabelBoost}(\hat{S}_i \circ \hat{D}_i)$. Let $S_{i+1} \subseteq \hat{D}'_i$ be a random subset of size $\lambda_{i+1} T_{i+1}$.
 - (g) Set $\alpha_{i+1} \leftarrow \alpha_i/2$ and $\beta_{i+1} \leftarrow \beta_i/2$.
-

354 **Theorem 5.2** (accuracy of algorithm `GenericBBL`). *Given $\alpha, \beta, \delta < 1/16, \epsilon < 1$, for any con-*
 355 *cept c and any round r , algorithm `GenericBBL` can predict the label of x_r as $h_r(x_r)$, such that*
 356 *$\Pr[\text{error}_{\mathcal{D}}(c(x_r) \neq h_r(x_r)) \leq 6\alpha] \geq 1 - 4\beta$.*

357 **Claim 5.3.** *`GenericBBL` is (ϵ, δ) -differentially private.*

358 **Remark 5.4.** *For simplicity, we analyzed `GenericBBL` in the realizable setting, i.e., under the*
 359 *assumption that the training set S is consistent with the target class C . Our construction carries*
 360 *over to the agnostic setting via standard arguments (ignoring computational efficiency). We*
 361 *refer the reader to [Beimel et al., 2021] and [Alon et al., 2020] for generic agnostic-to-realizable*
 362 *reductions in the context of private learning.*

363 **References**

- 364 Noga Alon, Roi Livni, Maryanthe Malliaris, and Shay Moran. Private PAC learning implies
365 finite littlestone dimension. In Moses Charikar and Edith Cohen, editors, *Proceedings of*
366 *the 51st Annual ACM SIGACT Symposium on Theory of Computing, STOC 2019, Phoenix,*
367 *AZ, USA, June 23-26, 2019*, pages 852–860. ACM, 2019. doi: 10.1145/3313276.3316312.
368 URL <https://doi.org/10.1145/3313276.3316312>.
- 369 Noga Alon, Amos Beimel, Shay Moran, and Uri Stemmer. Closure properties for private
370 classification and online prediction. In *COLT*, volume 125 of *Proceedings of Machine*
371 *Learning Research*, pages 119–152. PMLR, 2020.
- 372 Raef Bassily, Abhradeep Guha Thakurta, and Om Dipakbhai Thakkar. Model-agnostic
373 private learning. In *NeurIPS*, pages 7102–7112, 2018.
- 374 Amos Beimel, Kobbi Nissim, and Uri Stemmer. Characterizing the sample complexity of
375 private learners. In *ITCS*, pages 97–110. ACM, 2013a.
- 376 Amos Beimel, Kobbi Nissim, and Uri Stemmer. Private learning and sanitization: Pure vs.
377 approximate differential privacy. In *APPROX-RANDOM*, pages 363–378, 2013b.
- 378 Amos Beimel, Kobbi Nissim, and Uri Stemmer. Learning privately with labeled and
379 unlabeled examples. *Algorithmica*, 83(1):177–215, 2021.
- 380 Gavin Brown, Mark Bun, Vitaly Feldman, Adam D. Smith, and Kunal Talwar. When is
381 memorization of irrelevant training data necessary for high-accuracy learning? In *STOC*,
382 pages 123–132. ACM, 2021.
- 383 Mark Bun and Mark Zhandry. Order-revealing encryption and the hardness of private
384 learning. In *TCC (A1)*, volume 9562 of *Lecture Notes in Computer Science*, pages 176–206.
385 Springer, 2016.
- 386 Mark Bun, Kobbi Nissim, Uri Stemmer, and Salil P. Vadhan. Differentially private release
387 and learning of threshold functions. In *FOCS*, pages 634–649, 2015.
- 388 Mark Bun, Thomas Steinke, and Jonathan Ullman. Make up your mind: The price of online
389 queries in differential privacy. *CoRR*, abs/1604.04618, 2016. URL <http://arxiv.org/abs/1604.04618>.
- 391 Edith Cohen, Xin Lyu, Jelani Nelson, Tamás Sarlós, and Uri Stemmer. \tilde{O} ptimal differentially
392 private learning of thresholds and quasi-concave optimization. *CoRR*, abs/2211.06387,
393 2022. doi: 10.48550/arXiv.2211.06387. URL [https://doi.org/10.48550/arXiv.2211.](https://doi.org/10.48550/arXiv.2211.06387)
394 [06387](https://doi.org/10.48550/arXiv.2211.06387).
- 395 Yuval Dagan and Vitaly Feldman. PAC learning with stable and private predictions. In
396 *COLT*, volume 125 of *Proceedings of Machine Learning Research*, pages 1389–1410. PMLR,
397 2020.
- 398 Cynthia Dwork and Vitaly Feldman. Privacy-preserving prediction. In Sébastien Bubeck,
399 Vianney Perchet, and Philippe Rigollet, editors, *Conference On Learning Theory, COLT*
400 *2018, Stockholm, Sweden, 6-9 July 2018*, volume 75 of *Proceedings of Machine Learning*
401 *Research*, pages 1693–1702. PMLR, 2018. URL [http://proceedings.mlr.press/v75/](http://proceedings.mlr.press/v75/dwork18a.html)
402 [dwork18a.html](http://proceedings.mlr.press/v75/dwork18a.html).
- 403 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to
404 sensitivity in private data analysis. In *TCC*, pages 265–284, 2006.
- 405 Cynthia Dwork, Guy N. Rothblum, and Salil P. Vadhan. Boosting and differential privacy.
406 In *FOCS*, pages 51–60, 2010.
- 407 Vitaly Feldman and David Xiao. Sample complexity bounds on differentially private
408 learning via communication complexity. *SIAM J. Comput.*, 44(6):1740–1764, 2015. doi:
409 10.1137/140991844. URL <http://dx.doi.org/10.1137/140991844>.

- 410 Haim Kaplan, Katrina Ligett, Yishay Mansour, Moni Naor, and Uri Stemmer. Privately
411 learning thresholds: Closing the exponential gap. In *COLT*, volume 125 of *Proceedings of*
412 *Machine Learning Research*, pages 2263–2285. PMLR, 2020.
- 413 Shiva Prasad Kasiviswanathan, Homin K. Lee, Kobbi Nissim, Sofya Raskhodnikova, and
414 Adam Smith. What can we learn privately? *SIAM J. Comput.*, 40(3):793–826, 2011.
- 415 Michael J. Kearns, Mallesh M. Pai, Ryan M. Rogers, Aaron Roth, and Jonathan R. Ullman.
416 Robust mediators in large games. *CoRR*, abs/1512.02698, 2015.
- 417 Frank McSherry and Kunal Talwar. Mechanism design via differential privacy. In *FOCS*,
418 pages 94–103. IEEE, Oct 20–23 2007.
- 419 Anupama Nandi and Raef Bassily. Privately answering classification queries in the agnostic
420 PAC model. In *ALT*, volume 117 of *Proceedings of Machine Learning Research*, pages
421 687–703. PMLR, 2020.
- 422 L. G. Valiant. A theory of the learnable. *Commun. ACM*, 27(11):1134–1142, November 1984.
423 ISSN 0001-0782. doi: 10.1145/1968.1972. URL [http://doi.acm.org/10.1145/1968.](http://doi.acm.org/10.1145/1968.1972)
424 [1972](http://doi.acm.org/10.1145/1968.1972).
- 425 Vladimir N. Vapnik and Alexey Y. Chervonenkis. On the uniform convergence of relative
426 frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16
427 (2):264–280, 1971.

428 **A Additional Preliminaries from PAC Learning**

429 It is well known that a sample of size $\Theta(\text{VC}(C))$ is necessary and sufficient for the PAC
 430 learning of a concept class C , where the Vapnik-Chervonenkis (VC) dimension of a class C
 431 is defined as follows:

432 **Definition A.1** (VC-Dimension [Vapnik and Chervonenkis, 1971]). Let C be a concept class
 433 over a domain X . For a set $B = \{b_1, \dots, b_\ell\} \subseteq X$ of ℓ points, let $\Pi_C(B) = \{(c(b_1), \dots, c(b_\ell)) : c \in C\}$
 434 be the set of all dichotomies that are realized by C on B . We say that the set $B \subseteq X$ is *shattered*
 435 by C if C realizes all possible dichotomies over B , in which case we have $|\Pi_C(B)| = 2^{|B|}$.

436 The VC dimension of the class C , denoted $\text{VC}(C)$, is the cardinality of the largest set $B \subseteq X$
 437 shattered by C .

438 **Theorem A.2** (VC bound). *Let C be a concept class over a domain X . For $\alpha, \beta < 1/2$, there
 439 exists an (α, β, n) -PAC learner for C , where $n = \frac{8\text{VC}(C)\log(\frac{13}{\alpha}) + 4\log(\frac{2}{\beta})}{\alpha}$.*

440 **B Proof of Theorem 3.3**

441 The proof of Theorem 3.3 follows from algorithms HypothesisLearner, AccuracyBoost
 442 and claims B.1, B.2, all described below.

443 In Algorithm HypothesisLearner we assume that the everlasting differentially private
 444 prediction interface \mathcal{A} was fed with n i.i.d. samples taken from some (unknown) distribution
 445 \mathcal{D} and labeled by an unknown concept $c \in C$. Assuming the sequence of hypotheses $\{h_r\}_{r \geq 0}$
 446 produced by \mathcal{A} satisfies

$$\forall r \text{ error}_{\mathcal{D}}(c, h_r) \leq \alpha \quad (1)$$

447 we use it to construct – with constant probability – a hypothesis h with error bounded by
 448 $O(\alpha)$.

Algorithm HypothesisLearner

Parameters: $0 < \beta \leq 1/8$, $R = \lceil |X| \log(|X|) \log(1/\beta) \rceil$

Input: hypothesis sequence $\{h_r\}_{r \geq 0}$

1. for all $x \in X$ let $L_x = \emptyset$
 2. for $r = 0, 1, 2, \dots, R$
 - (a) select x uniformly at random from X and let $L_x = L_x \cup \{h_r(x)\}$
 3. if $L_x = \emptyset$ for some $x \in X$ then fail, output an arbitrary hypothesis, and halt
 /* $\Pr[\exists x \text{ such that } L_x = \emptyset] \leq |X|(1 - \frac{1}{|X|})^R \approx |X|e^{-R/|X|} = \beta$ */
 4. for all $x \in X$ let r_x be sampled uniformly at random from L_x
 5. construct the hypothesis h , where $h(x) = r_x$
-

449 **Claim B.1.** *If executed on a hypothesis sequence satisfying Equation 1 then with probability at
 450 least 3/4 Algorithm HypothesisLearner outputs a hypothesis h satisfying $\text{error}_{\mathcal{D}}(c, h) \leq 8\alpha$.*

451 *Proof.* Having $\mathcal{D}, c \in C$ fixed, and given a hypothesis h , we define $e_h(x)$ to be 1 if $h(x) \neq c(x)$
 452 and 0 otherwise. Thus, we can write $\text{error}_{\mathcal{D}}(c, h) = \mathbb{E}_{x \sim \mathcal{D}}[e_h(x)]$.

453 Observe that when Algorithm HypothesisLearner does not fail, r_x (and hence $h(x)$) is cho-
 454 sen with equal probability among $(h_1(x), h_2(x), \dots, h_R(x))$ and hence $\mathbb{E}_{\theta}[e_h(x)] = \mathbb{E}_{i \in [R]}[e_{h_i}(x)]$
 455 where θ denotes the randomness of HypothesisLearner. We get:

$$\begin{aligned} \mathbb{E}_{\theta}[\text{error}_{\mathcal{D}}(c, h)] &= \mathbb{E}_{\theta}[\mathbb{E}_{x \sim \mathcal{D}}[e_h(x)]] = \mathbb{E}_{x \sim \mathcal{D}}[\mathbb{E}_{\theta}[e_h(x)]] \\ &= \mathbb{E}_{x \sim \mathcal{D}}[\mathbb{E}_{i \in [R]}[e_{h_i}(x)]] = \mathbb{E}_{i \in [R]}[\mathbb{E}_{x \sim \mathcal{D}}[e_{h_i}(x)]] \\ &\leq \mathbb{E}_{i \sim \mathcal{R}}[\alpha] = \alpha. \end{aligned}$$

456 By Markov inequality, we have $\Pr_{\theta}[\text{error}_{\mathcal{D}}(c, h) \geq 8\alpha] \leq 1/8$. The claim follows noting that
 457 Algorithm HypothesisLearner fails with probability at most $\beta \leq 1/8$. \square

458 The second part of the transformation is Algorithm AccuracyBoost that applies Algorithm
 459 HypothesisLearner $O(\log(1/\beta))$ times to obtain with high probability a hypothesis with
 460 $O(\alpha)$ error.

Algorithm AccuracyBoost

Parameters: $\beta, R = 104 \ln \frac{1}{\beta}$

Input: R labeled samples with n examples each (S_1, \dots, S_R) where $S_i \in (X \times \{0, 1\})^n$

1. for $i = 1, 2 \dots R$
 - (a) execute $\mathcal{A}(S_i)$ to obtain a hypothesis sequence $\{h_r^i\}_{r \geq 0}$
 - (b) execute Algorithm WeakHypothesisLearner on $\{h_r^i\}_{r \geq 0}$ to obtain hypothesis h^i
 2. construct the hypothesis \hat{h} , where $\hat{h}(x) = \text{maj}(h^1(x), \dots, h^R(x))$.
-

461 **Claim B.2.** *With probability $1 - \beta$, Algorithm AccuracyBoost output a 24α -good hypothesis*
 462 *over distribution \mathcal{D} .*

Proof. Define B_i to be the event where the sequence of hypotheses $\{h_r^i\}_{r \geq 0}$ produced in Step 1a of AccuracyBoost does not satisfy Equation 1. We have,

$$\Pr[\text{error}_{\mathcal{D}}(c, h_i) > 8\alpha] \leq \Pr[B] + (1 - \Pr[B]) \cdot \Pr[\text{error}_{\mathcal{D}}(c, h) > 8\alpha] \leq \beta + 1/4 < 3/8.$$

463 Hence, by the Chernoff bound, when $R \geq 104 \ln \frac{1}{\beta}$, we have at least $7R/8$ hypotheses are
 464 8α -good over distribution \mathcal{D} . Consider the worst case, in which $R/8$ hypotheses always
 465 output wrong labels. To output a wrong label of x , we require at least $3R/8$ hypotheses to
 466 output wrong labels. Thus \hat{h} is 24α -good over distribution \mathcal{D} . \square

467 C Tools from Prior Works

468 C.1 Algorithm LabelBoost [Beimel et al., 2021]

Algorithm LabelBoost [Beimel et al., 2021]

Parameters: A concept class C .

Input: A partially labeled database $S \circ T \in (X \times \{0, 1, \perp\})^*$.

% We assume that the first portion of the database (denoted S) contains labeled examples.
 The algorithm outputs a similar database where both S and T are (re)labeled.

1. Initialize $H = \emptyset$.
 2. Let $P = \{p_1, \dots, p_\ell\}$ be the set of all points $p \in X$ appearing at least once in $S \circ T$.
 Let $\Pi_C(P) = \{(c(p_1), \dots, c(p_\ell)) : c \in C\}$ be the set of all dichotomies generated by C on P .
 3. For every $(z_1, \dots, z_\ell) \in \Pi_C(P)$, add to H an arbitrary concept $c \in C$ s.t. $c(p_i) = z_i$ for every $1 \leq i \leq \ell$.
 4. Choose $h \in H$ using the exponential mechanism with privacy parameter $\epsilon=1$, solution set H , and the database S .
 5. (Re)label $S \circ T$ using h , and denote the resulting database $(S \circ T)^h$, that is, if $S \circ T = (x_i, y_i)_{i=1}^t$ then $(S \circ T)^h = (x_i, y'_i)_{i=1}^t$ where $y'_i = h(x_i)$.
 6. Output $(S \circ T)^h$.
-

469 **Lemma C.1** (privacy of Algorithm LabelBoost [Beimel et al., 2021]). *Let \mathcal{A} be an (ϵ, δ) -*
 470 *differentially private algorithm operating on partially labeled databases. Construct an algorithm*
 471 *\mathcal{B} that on input a partially labeled database $S \circ T \in (X \times \{0, 1, \perp\})^*$ applies \mathcal{A} on the outcome of*
 472 *LabelBoos($S \circ T$). Then, \mathcal{B} is $(\epsilon + 3, 4\epsilon\delta)$ -differentially private.*

473 Consider an execution of `LabelBoost` on a database $S \circ T$, and assume that the examples in
474 S are labeled by some target concept $c \in C$. Recall that for every possible labeling \vec{z} of the
475 elements in S and in T , algorithm `LabelBoost` adds to H a hypothesis from C that agrees
476 with \vec{z} . In particular, H contains a hypothesis that agrees with the target concept c on S
477 (and on T). That is, $\exists f \in H$ s.t. $\text{error}_S(f) = 0$. Hence, the exponential mechanism (on Step 4)
478 chooses (w.h.p.) a hypothesis $h \in H$ s.t. $\text{error}_S(h)$ is small, provided that $|S|$ is roughly $\log|H|$,
479 which is roughly $\text{VC}(C) \cdot \log(|S| + |T|)$ by Sauer’s lemma. So, algorithm `LabelBoost` takes an
480 input database where only a small portion of it is labeled, and returns a similar database in
481 which the labeled portion grows exponentially.

Lemma C.2 (utility of Algorithm `LabelBoost` [Beimel et al., 2021]). *Fix α and β , and let $S \circ T$ be s.t. S is labeled by some target concept $c \in C$, and s.t.*

$$|T| \leq \frac{\beta}{e} \text{VC}(C) \exp\left(\frac{\alpha|S|}{2\text{VC}(C)}\right) - |S|.$$

482 *Consider the execution of `LabelBoost` on $S \circ T$, and let h denote the hypothesis chosen on Step 4.*
483 *With probability at least $(1 - \beta)$ we have that $\text{error}_S(h) \leq \alpha$.*

484 C.2 Algorithm `BetweenThresholds` [Bun et al., 2016]

Algorithm `BetweenThresholds` [Bun et al., 2016]

Input: Database $S \in X^n$.

Parameters: $\varepsilon, t_\ell, t_u \in (0, 1)$ and $n, k \in \mathbb{N}$.

1. Sample $\mu \sim \text{Lap}(2/\varepsilon n)$ and initialize noisy thresholds $\hat{t}_\ell = t_\ell + \mu$ and $\hat{t}_u = t_u - \mu$.
 2. For $j = 1, 2, \dots, k$:
 - (a) Receive query $q_j : X^n \rightarrow [0, 1]$.
 - (b) Set $c_j = q_j(S) + v_j$ where $v_j \sim \text{Lap}(6/\varepsilon n)$.
 - (c) If $c_j < \hat{t}_\ell$, output L and continue.
 - (d) If $c_j > \hat{t}_u$, output R and continue.
 - (e) If $c_j \in [\hat{t}_\ell, \hat{t}_u]$, output \top and halt.
-

485 **Lemma C.3** (Privacy for `BetweenThresholds` [Bun et al., 2016]). *Let $\varepsilon, \delta \in (0, 1)$ and $n \in \mathbb{N}$.*
486 *Then algorithm `BetweenThresholds` is (ε, δ) -differentially private for any adaptively-chosen*
487 *sequence of queries as long as the gap between the thresholds t_ℓ, t_u satisfies*

$$t_u - t_\ell \geq \frac{12}{\varepsilon n} (\log(10/\varepsilon) + \log(1/\delta) + 1).$$

488 **Lemma C.4** (Accuracy for `BetweenThresholds` [Bun et al., 2016]). *Let $\alpha, \beta, \varepsilon, t_\ell, t_u \in (0, 1)$*
489 *and $n, k \in \mathbb{N}$ satisfy*

$$n \geq \frac{8}{\alpha \varepsilon} (\log(k + 1) + \log(1/\beta)).$$

490 *Then, for any input $x \in X^n$ and any adaptively-chosen sequence of queries q_1, q_2, \dots, q_k , the*
491 *answers $a_1, a_2, \dots, a_{\leq k}$ produced by `BetweenThresholds` on input x satisfy the following with*
492 *probability at least $1 - \beta$. For any $j \in [k]$ such that a_j is returned before `BetweenThresholds`*
493 *halts,*

- 494 • $a_j = \text{L} \implies q_j(x) \leq t_\ell + \alpha,$
- 495 • $a_j = \text{R} \implies q_j(x) \geq t_u - \alpha,$ and
- 496 • $a_j = \top \implies t_\ell - \alpha \leq q_j(x) \leq t_u + \alpha.$

497 **Observation 2.** *Using standard composition theorems for differential privacy (see, e.g., Dwork*
498 *et al. [2010]), we can assume that algorithm `BetweenThresholds` takes another parameter c ,*
499 *and halts after c times of outputting \top . In this case, the algorithm satisfies $(\varepsilon', 2c\delta)$ -differential*
500 *privacy, for $\varepsilon' = \sqrt{2c \ln(\frac{1}{c\delta})} \varepsilon + c\varepsilon(e^\varepsilon - 1)$.*

501 **D Some Technical Facts**

502 We refer to the execution of steps 3a-3g of algorithm `GenericBBL` as a *phase* of the algorithm,
 503 indexed by $i = 1, 2, 3, \dots$

504 The original `BetweenThresholds` needs to halt when it outputs \top . In `GenericBBL`, we toler-
 505 ance it to halt at most c_i times in the phase i . We prove `BetweenThresholds` in `GenericBBL`
 506 is $(1, \delta)$ -differentially private.

507 **Claim D.1.** For $\delta < 1$, Mechanism `BetweenThresholds` used in step 3c in the i -th iteration, is
 508 $(1, \delta)$ -differentially private.

Proof. Let ϵ'_i, δ'_i be as in Step 3c. Since $e^{\epsilon'_i} - 1 < 2\epsilon'_i$ for $0 < \epsilon'_i < 1$, we have

$$\sqrt{2c_i \ln\left(\frac{1}{c_i \delta'_i}\right) \cdot \epsilon'_i + c_i \epsilon'_i (e^{\epsilon'_i} - 1)} \leq \sqrt{2c_i \ln\left(\frac{2}{\delta}\right) \cdot \epsilon'_i + 2c_i \epsilon_i'^2} = \frac{\sqrt{2}}{3} + \frac{2}{9 \ln\left(\frac{2}{\delta}\right)} \leq 1.$$

509 The proof is concluded by using observation 2. □

510 In Claim D.2- D.5, we prove that with high probability, `BetweenThresholds` in step 3d halts
 511 within $64\alpha_i$ times. We prove it by 4 steps:

- 512 1. prove that with high probability, most hypothesis in step 3b have high accuracy
 513 (Claim D.2).
- 514 2. prove that if most hypothesis in step 3b have high accuracy, then with high probability,
 515 the queries in `BetweenThresholds` are closed to 0 or 1 (Claim D.3).
- 516 3. prove that if the queries in `BetweenThresholds` are closed to 0 or 1, then
 517 `BetweenThresholds` in step 3d will outputs L or R with high probability (Claim D.4).
- 518 4. prove that if `BetweenThresholds` outputs L or R , then every single phase fails with low
 519 probability (Claim D.5).

520 **Claim D.2.** If $\beta_i \leq 1/32$ and $T_i \geq 96 \ln \frac{1}{\alpha_i}$, then with probability $1 - \alpha_i$, $\frac{15T_i}{16}$ hypotheses in step 3b
 521 are α_i -good with respect to g_i , where g_i is the concept of S_i .

Proof. By the VC bound (Theorem A.2), for each $t \in [T_i]$, we have

$$\Pr[\text{error}_{\mathcal{D}}(f_t, g_i) \leq \alpha_i] \geq 1 - \beta_i.$$

522 By Chernoff bound, if $T_i \geq \frac{16+256\beta_i}{(1-16\beta_i)^2} \ln \frac{1}{\alpha_i}$, then with probability $1 - \alpha_i$, we have $\frac{15T_i}{16}$ hypothe-
 523 ses have $\text{error}_{\mathcal{D}}(f_t, g_i) \leq \alpha_i$. When $\beta_i \leq 1/32$, it is sufficient to set $T_i \geq 96 \ln \frac{1}{\alpha_i}$. □

524 **Claim D.3.** If $\alpha_i \leq 1/16$ and $\frac{15T_i}{16}$ hypotheses in step 3b are α_i -good with respect to g_i , where g_i
 525 is the concept of S_i , then $\Pr_{x \sim \mathcal{D}}[|q(x) - \frac{1}{2}| \leq \frac{3}{8}] \leq 15\alpha_i$.

526 *Proof.* W.l.o.g. assume $g_i(x) = 1$, where g_i is the concept of S_i , so it is sufficient to prove
 527 $\Pr_{x \sim \mathcal{D}}[q(x) \leq \frac{7}{8}] \leq 8\alpha_i$. Consider the worst case that $\frac{T_i}{16}$ "bad" hypotheses output 0. In that
 528 case, $q(x) \leq \frac{7}{8}$ when $\frac{T_i}{16}$ of α_i -good hypotheses output 0. So that with probability $15\alpha_i$, we
 529 have $q(x) \leq \frac{7}{8}$. (see Figure 2)

530 □

531 **Claim D.4.** Let $t_u < 1/2 + 1/8$ and $t_\ell > 1/2 - 1/8$. For a query q such that $q(S) > 7/8$ (similarly,
 532 for $q(S) < 1/8$), Algorithm `BetweenThresholds` outputs R (similarly, L) with probability at least

533 $1 - \exp\left(-\frac{T_i}{144 \sqrt{c_i \ln\left(\frac{2}{\delta}\right)}}\right).$

534 *Proof.* Wlog assume $q(S) > 7/8$, it is sufficient to show

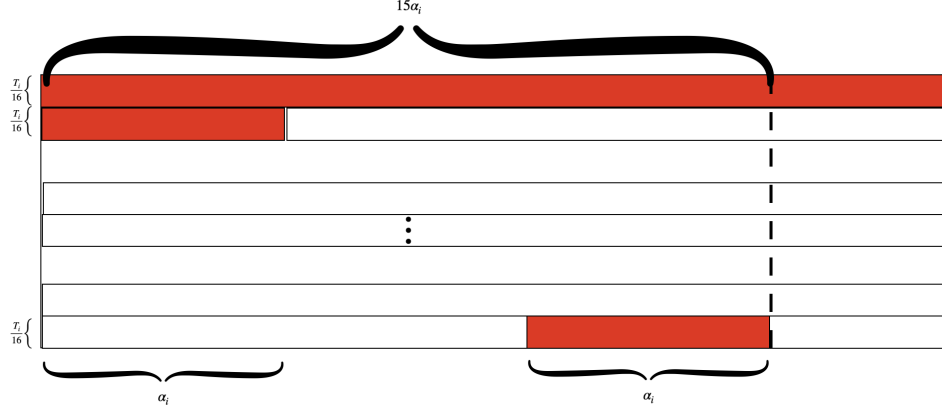


Figure 2: The horizontal represents the input point. The vertical represents the hypothesis. The red parts represent the incorrect prediction. We let $\frac{T_i}{16}$ hypothesis predict all labels as 0. To let $q(x) \leq \frac{7}{8}$, there must exist $\frac{T_i}{16}$ hypothesis output 0. In the worst case, at most $15\alpha_i$ of points are labeled as 0.

$$\begin{aligned}
\Pr[\text{BetweenThreshold outputs } R] &= \Pr[q(S) + \text{Lap}(6/\varepsilon' T_i) > t_u + \text{Lap}(2/\varepsilon' T_i)] \\
&> \Pr[\text{Lap}(6/\varepsilon' T_i) > -1/8] \cdot \Pr[\text{Lap}(2/\varepsilon' T_i) < 1/8] \\
&= \left(1 - \frac{1}{2} \exp\left(-\frac{T_i}{144\sqrt{c_i \ln(\frac{2}{\delta})}}\right)\right) \cdot \left(1 - \frac{1}{2} \exp\left(-\frac{T_i}{48\sqrt{c_i \ln(\frac{2}{\delta})}}\right)\right) \\
&> 1 - \exp\left(-\frac{T_i}{144\sqrt{c_i \ln(\frac{2}{\delta})}}\right).
\end{aligned}$$

535

□

536 **Claim D.5.** For any phase i , *BetweenThresholds* outputs \top at most $64\alpha_i R_i$ times with proba-
537 bility at most β_i .

538 *Proof.* For a single query, if $t_u < 1/2 + 1/8$ and $q(S) > 7/8$ (similarly, $t_\ell > 1/2 - 1/8$
539 and $q(S) < 1/8$), by Claim D.4, *BetweenThresholds* outputs \top with probability at

540 most $\exp\left(-\frac{T_i}{144\sqrt{c_i \ln(\frac{2}{\delta})}}\right) = \exp\left(-\frac{T_i}{144\sqrt{64\alpha_i R_i \ln(\frac{2}{\delta})}}\right) < \alpha_i$. Combine Claim D.2 and D.3,

541 *BetweenThresholds* outputs \top with probability at most $32\alpha_i$. By the Chernoff bound

542 and $R_i \geq \frac{3\ln(\frac{1}{\beta_i})}{\alpha_i}$, *BetweenThresholds* outputs \top more than $64\alpha_i R_i$ times with probability

543 at most β_i . □

544 In step 3f, *GenericBBL* takes a random subset of size $\lambda_{i+1} T_{i+1}$ from \hat{D}'_i . We show that the
545 size of \hat{D}'_i is at least $\lambda_{i+1} T_{i+1}$.

546 **Claim D.6.** When $\varepsilon \leq 1$, for any $i \geq 1$, we always have $|\hat{D}'_i| \geq \lambda_{i+1} T_{i+1}$.

547 *Proof.* Let $m = 3 + \exp(\varepsilon + 4) < 200$. By the step 3c, step 3e and step 3f, $|\hat{D}_j| = \frac{\varepsilon |D_j|}{m} = \frac{25600 |S_j|}{m} \geq$
548 $128 |S_j| = 128 \lambda_j T_j$. Then it is sufficient to verify $128 \lambda_j T_j \geq \lambda_{j+1} T_{j+1}$

We can verify that

$$4\lambda_j = 4 \cdot \frac{8\text{VC}(C)\log(\frac{13}{\alpha_i}) + 4\log(\frac{2}{\beta_i})}{\alpha_i} = 4 \cdot \frac{8\text{VC}(C)(\log(\frac{13}{\alpha_{j+1}}) - 1) + 4(\log(\frac{2}{\beta_{j+1}}) - 1)}{2\alpha_{j+1}} \geq \lambda_{j+1}$$

and

$$32T_j = \frac{32\tau \cdot \lambda_i \cdot \log(\frac{1}{\delta}) \cdot \log^2(\frac{\lambda_i}{\varepsilon\alpha_i\beta_i\delta})}{\alpha_i\varepsilon} \geq \frac{32\tau \cdot \lambda_i \cdot \log(\frac{1}{\delta}) \cdot \log^2(\frac{\lambda_{i+1}}{16\varepsilon\alpha_{i+1}\beta_{i+1}\delta})}{8\alpha_{i+1}\varepsilon} \geq \lambda_{j+1}T_{j+1}.$$

549 The last inequality holds because $\lambda_j \geq 4$ and $\alpha_j, \beta_j \leq 1/2$. \square

550 To apply the privacy and accuracy of *LabelBoost* and *BetweenThresholds*, the sizes of the
551 databases need to satisfy the inequalities in lemma C.2, C.3 and C.4. We verify that in each
552 phase, the sizes of the databases always satisfy the requirement.

Claim D.7. *Let $\alpha, \beta, \delta < 1/16$, $\varepsilon \leq 1$, and $\text{VC}(C) \geq 1$. Then for any $i \geq 1$, we have*

$$T_i \geq \frac{8}{\alpha_i\varepsilon'} (\log(|D_i| + 1) + \log(1/\beta_i)).$$

553 *Proof.* By claim D.6 and step 3c, $|D_i| = \frac{25600|S_i|}{\varepsilon} = \frac{25600\lambda_i T_i}{\varepsilon}$. Since

$$\begin{aligned} \frac{8}{\alpha_i\varepsilon'} (\log(|D_i| + 1) + \log(1/\beta_i)) &= \frac{24\sqrt{64\alpha_i|D_i|\ln(\frac{2}{\delta})}}{\sqrt{2}\alpha_i} \cdot (\log(|D_i| + 1) + \log(1/\beta_i)) \\ &= O\left(\sqrt{\frac{\lambda_i T_i \log(\frac{1}{\delta})}{\alpha_i\varepsilon}} \left(\log\left(\frac{\lambda_i T_i}{\varepsilon\beta_i}\right)\right)\right) \\ &= O\left(\sqrt{\frac{\lambda_i T_i \log(\frac{1}{\delta})}{\alpha_i\varepsilon}} \cdot \log\left(\frac{\lambda_i \log(\frac{1}{\delta})}{\alpha_i\beta_i\varepsilon}\right)\right), \end{aligned}$$

554 and $T_i = \frac{\tau \cdot \lambda_i \cdot \log(\frac{1}{\delta}) \cdot \log^2(\frac{\lambda_i}{\varepsilon\alpha_i\beta_i\delta})}{\alpha_i\varepsilon}$, where $\tau \geq 1.1 * 10^{10}$, the inequality always holds. \square

555 **Claim D.8.** *When $\varepsilon \leq 1$, for any $i \geq 1$, we have $|\hat{D}_i| \leq \frac{\beta_i}{\varepsilon} \text{VC}(C) \exp\left(\frac{\alpha_i|\hat{S}_i|}{2\text{VC}(C)}\right) - |\hat{S}_i|$.*

Proof. By claim D.6, step 3c and step 3f,

$$|\hat{D}_i| = \frac{\varepsilon|D_i|}{m} = O(\lambda_i T_i) = O\left(\text{VC}(C)\log^2(\text{VC}(C)) \cdot \text{poly}\left(\frac{1}{\alpha_i}, \log\left(\frac{1}{\beta_i}\right), \frac{1}{\varepsilon}, \log\left(\frac{1}{\delta}\right)\right)\right)$$

556 and

$$\begin{aligned} |\hat{S}_i| &= \frac{\varepsilon|S_i|}{m} \\ &= O(\varepsilon\lambda_i T_i) = O(\lambda_i T_i) \\ &= O\left(\text{VC}(C)\log^2(\text{VC}(C)) \cdot \text{poly}\left(\frac{1}{\alpha_i}, \log\left(\frac{1}{\beta_i}\right), \frac{1}{\varepsilon}, \log\left(\frac{1}{\delta}\right)\right)\right). \end{aligned} \tag{2}$$

Note that

$$\frac{\beta_i}{\varepsilon} \text{VC}(C) \exp\left(\frac{\alpha_i|\hat{S}_i|}{2\text{VC}(C)}\right) = \Omega\left(\text{VC}^2(C) \cdot \exp\left(\text{poly}\left(\frac{1}{\alpha_i}, \log\left(\frac{1}{\beta_i}\right), \frac{1}{\varepsilon}, \log\left(\frac{1}{\delta}\right)\right)\right)\right),$$

557 for $T_i = \frac{\tau \cdot \lambda_i \cdot \log(\frac{1}{\delta}) \cdot \log^2(\frac{\lambda_i}{\varepsilon\alpha_i\beta_i\delta})}{\alpha_i\varepsilon}$, the inequality holds when $\tau \geq 1$. \square

558 **Claim D.9.** For every $i \geq 1$, we have

$$t_u - t_\ell \geq \frac{12}{\varepsilon'_i T_i} \left(\log(10/\varepsilon'_i) + \log(1/\delta'_i) + 1 \right).$$

Proof. By step 3c, $t_u - t_\ell = 2\alpha_i$. Then we have

$$\begin{aligned} \frac{6}{\alpha_i \varepsilon'_i T_i} \left(\log(10/\varepsilon'_i) + \log(1/\delta'_i) + 1 \right) &= \frac{6\sqrt{64\alpha_i R_i \ln(\frac{2}{\delta})}}{\alpha_i T_i} \left(\log(10/\varepsilon'_i) + \log(1/\delta'_i) + 1 \right) \\ &= 6\sqrt{\frac{1638400 \ln(\frac{2}{\delta}) \lambda_i}{\alpha_i T_i}} \left(\log(10/\varepsilon'_i) + \log(1/\delta'_i) + 1 \right) \\ &= 6\sqrt{\frac{1638400 \ln(\frac{2}{\delta})}{\tau \log(\frac{1}{\delta}) \log^2(\frac{\lambda_i}{\varepsilon \alpha_i \beta_i \delta})}} \left(\log(10/\varepsilon'_i) + \log(1/\delta'_i) + 1 \right) \\ &= O(1), \end{aligned}$$

559 the inequality holds when $\tau > 10^{10}$. □

560 E Accuracy of Algorithm GenericBBL – proof of Theorem 5.2

561 We refer to the execution of steps 3a-3g of algorithm GenericBBL as a *phase* of the algorithm,
562 indexed by $i = 1, 2, 3, \dots$

563 We give some technical facts in Appendix D. In Claim E.1, we show that in each phase,
564 samples are labeled with high accuracy. In Claim E.2, we prove that algorithm GenericBBL
565 fails with low probability. In Claim E.4, we prove that algorithm GenericBBL predict the
566 labels with high accuracy.

Claim E.1. When Algorithm GenericBBL does not fail on phases 1 to i , then for phase $i + 1$ we have

$$\Pr \left[\exists g_{i+1} \in C \text{ s.t. } \text{error}_{S_{i+1}}(g_{i+1}) = 0 \text{ and } \text{error}_{\mathcal{D}}(g_{i+1}, c) \leq \sum_{j=1}^{i+1} \alpha_j \right] \geq 1 - 2 \sum_{j=0}^{i+1} \beta_j.$$

567 *Proof.* The proof is by induction on i . The base case for $i = 1$ is trivial, with $g_1 = c$. Assume
568 the claim holds for all $j \leq i$. By the properties of LabelBoost (Lemma C.2) and Claim D.8,
569 with probability at least $1 - \beta_{i+1}$ we have that S_{i+1} is labeled by a hypothesis $g_{i+1} \in C$ s.t.
570 $\text{error}_{S_i}(g_i, g_{i+1}) \leq \alpha_{i+1}$. Observe that the points in S_i (without their labels) are chosen i.i.d.
571 from \mathcal{D} , and hence, By Theorem A.2 (VC bounds) and $|S_i| \geq 128\lambda_i \geq \lambda_{i+1}$, with probability at
572 least $1 - \beta_{i+1}$ we have that $\text{error}_{\mathcal{D}}(g_i, g_{i+1}) \leq \alpha_{i+1}$. Hence, with probability $1 - 2\beta_{i+1}$, we have
573 $\text{error}_{\mathcal{D}}(g_i, g_{i+1}) \leq \alpha_{i+1}$. Finally, by the triangle inequality, $\text{error}_{\mathcal{D}}(g_{i+1}, c) \leq \sum_{j=1}^{i+1} \alpha_j$, except
574 with probability $2 \sum_{j=1}^{i+1} \beta_j$ □

575 Define the following good event.

576 **Event E_1 :** Algorithm GenericBBL never fails on the execution of
BetweenThresholds in step 3(d)iv.

577 **Claim E.2.** Event E_1 occurs with probability at least $1 - \beta$.

Proof. Using to union bound and Claim D.5,

$$\Pr[\text{Event } E_1 \text{ occurs}] \geq 1 - \beta.$$

578 □

579 Combining claims E.1 and E.2, we get:

Claim E.3. Let \mathcal{D} be an underlying distribution and let $c \in C$ be a target concept. Then

$$\Pr[\forall i \exists g_i \in C \text{ s.t. } \text{error}_{S_i}(g_i) = 0 \text{ and } \text{error}_{\mathcal{D}}(g_i, c) \leq \alpha] \geq 1 - 3\beta.$$

580 **Notations.** Consider the i th phase of Algorithm `GenericBBL`, and focus on the j -th iteration of Step 3. Fix all of the randomness in `BetweenThresholds`. Now observe that the output on step 3(d)iii is a deterministic function of the input $x_{i,j}$. This defines a hypothesis which we denote as $h_{i,j}$.

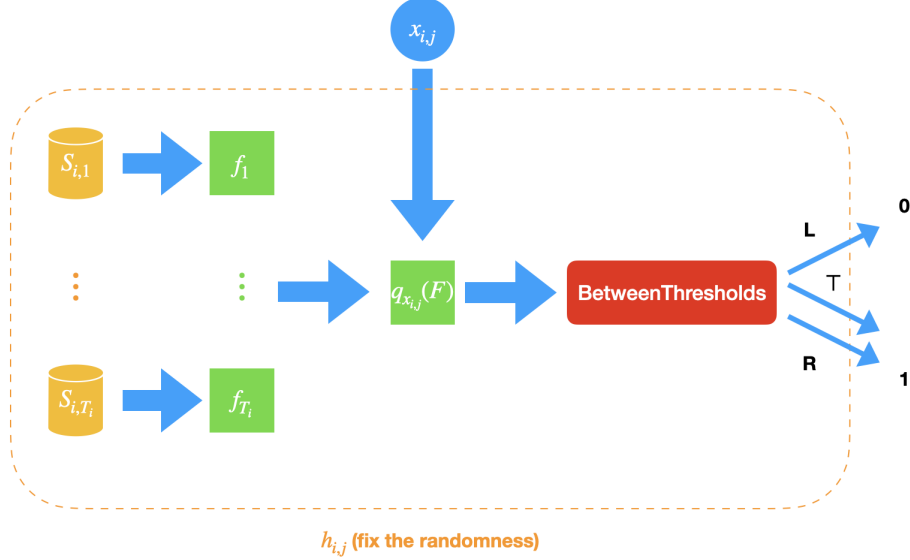


Figure 3: Hypothesis $h_{i,j}$

584 **Claim E.4.** For $\beta < 1/16$, with probability at least $1 - 4\beta$, all of the hypotheses defined above are
585 6α -good w.r.t. \mathcal{D} and c .

586 *Proof.* In the phase i , by Claim E.3, with probability at least $1 - 3\beta$ we have that S_i is labeled
587 by a hypothesis $g_i \in C$ satisfying $\text{error}_{\mathcal{D}}(g_i, c) \leq \alpha$. We continue with the analysis assuming
588 that this is the case.

589 On step 3a of the i th phase we divide S_i into T_i subsamples of size λ_i each, identify
590 a consistent hypothesis $f_t \in C$ for every subsample $S_{i,t}$, and denote $F_i = (f_1, \dots, f_{T_i})$. By
591 Theorem A.2 (VC bounds), every hypothesis in F_i satisfies $\text{error}_{\mathcal{D}}(f_t, g_i) \leq \alpha$ with probability
592 $3/4$, in which case, by the triangle inequality we have that $\text{error}_{\mathcal{D}}(f_t, c) \leq 2\alpha$.

593 Set $T_i \geq \frac{512(1-4\beta_i)\ln(\frac{1}{\beta_i})}{(1-64\beta_i)^2}$, using Chernoff bound, it holds that for at least $15T_i/16$ of the hy-
594 potheses in F_i have error $\text{error}_{\mathcal{D}}(f_t, g_i) \leq 2\alpha$ with probability at least $1 - \beta_i$. These hypotheses
595 have $\text{error}_{\mathcal{D}}(f_t, c) \leq 3\alpha$.

596 Let $m : X \rightarrow \{0, 1\}$ defined as $m(x) = \text{maj}_{f_t \in F_i}(f_t(x))$. For m to err on a point x (w.r.t. the target
597 concept c), it must be that at least $7/16$ -fraction of the 3α -good hypotheses in \hat{F}_i err on x .
598 Consider the worst case in Figure 4, we have $\text{error}_{\mathcal{D}}(m, c) \leq 6\alpha$

599 By Lemma C.4 and Claim D.7, with probability at least $1 - \beta_i$, all of the hypotheses defined
600 during the i th iteration satisfy this condition, and are hence 6α -good w.r.t. c and \mathcal{D} . By the
601 union bound, with probability $1 - 4\beta$, all the hypotheses are 6α -good. \square

602 E.1 Privacy analysis – proof of Claim 5.3

603 Fix $t \in \mathbb{N}$ and the adversary \mathcal{B} . We need to show that $\text{View}_{\mathcal{B},t}^0$ and $\text{View}_{\mathcal{B},t}^1$ (defined in
604 Figure 1) are (ϵ, δ) -indistinguishable. We will consider separately the case where the
605 executions differ in the training phase (Claim E.5) and the case where the difference occurs
606 during the prediction phase (Claim E.6).

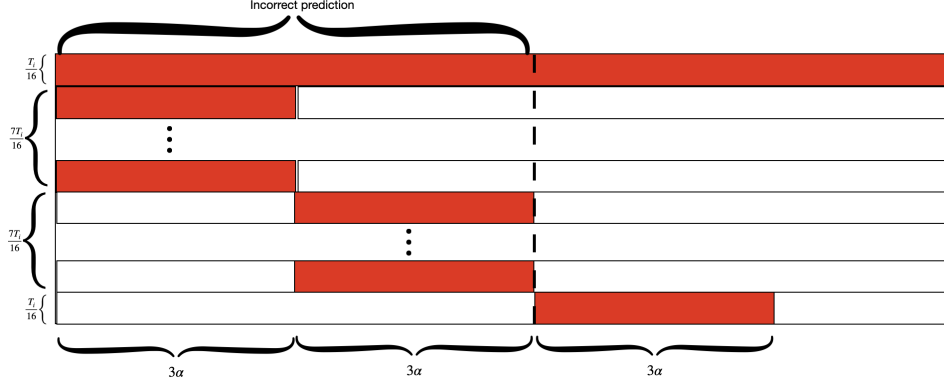


Figure 4: The horizontal represents the input point. The vertical represents the hypothesis. The red parts represent the incorrect prediction. We let $\frac{T_i}{16}$ hypothesis predict all labels incorrectly. To output an incorrect label, there must exist $\frac{7T_i}{16}$ hypothesis output the incorrect label. In the worst case, at most 6α of points are incorrectly classified.

607 **Privacy of the initial training set S .** Let $S^0, S^1 \in (X \times \{0, 1\})^n$ be neighboring datasets of
 608 labeled examples and let $\text{View}_{\mathcal{B}, t}^0$ and $\text{View}_{\mathcal{B}, t}^1$ be as in Figure 1 where $((x_1^0, y_1^0), \dots, (x_n^0, y_n^0)) =$
 609 S^0 and $((x_1^1, y_1^1), \dots, (x_n^1, y_n^1)) = S^1$.

610 **Claim E.5.** For all adversaries \mathcal{B} , for all $t > 0$, and for any two neighbouring database S^0 and S^1
 611 selected by \mathcal{B} , $\text{View}_{\mathcal{B}, t}^0$ and $\text{View}_{\mathcal{B}, t}^1$ are (ϵ, δ) -indistinguishable.

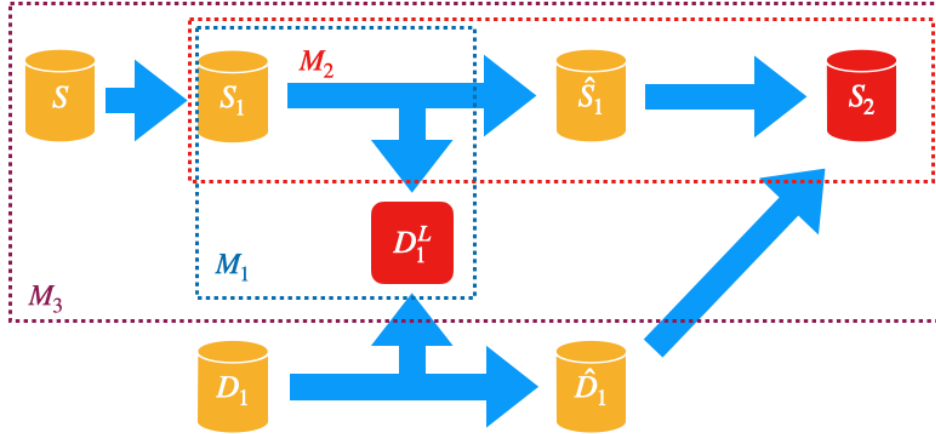


Figure 5: Privacy of the labeled sample S

612 *Proof.* Let $R'_1 = \min(t, R_1)$. Note that $\text{View}_{\mathcal{B}, R'_1}^b$ is a prefix of $\text{View}_{\mathcal{B}, t}^b$ which includes the
 613 labels Algorithm GenericBBL produces in Step 3(d)iii for the R'_1 first unlabeled points
 614 selected by \mathcal{B} . Let S_2^b be the result of the first application of algorithm LabelBoost in Step 3f
 615 of GenericBBL (if $t < R_1$ we set S_2^b as \perp). The creation of these random variables is depicted
 616 in Figure 5, where D_1^L denotes the labels Algorithm GenericBBL produces for the unlabeled
 617 points D_1 .

618 Observe that $\text{View}_{\mathcal{B},t}^b$ results from a post-processing (jointly by the adversary \mathcal{B} and Algo-
619 rithm `GenericBBL`) of the random variable $(\text{View}_{\mathcal{B},R'_1}^b, S_2^b)$, and hence it suffices to show that
620 $(\text{View}_{\mathcal{B},R'_1}^0, S_2^0)$ and $(\text{View}_{\mathcal{B},R'_1}^1, S_2^1)$ are (ϵ, δ) -indistinguishable.

621 We follow the processes creating $\text{View}_{\mathcal{B},t}^b$ and S_2^b in Figure 5: (i) The mechanism M_1 cor-
622 responds to the loop in Step 3d of `GenericBBL` where labels are produced for the adver-
623 sarially chosen points D_1^b . By application of Lemma C.3, M_1 is $(1, \delta)$ -differentially private.
624 (ii) The mechanism M_2 , corresponds to the subsampling of \hat{S}_1^b from S_1^b and the applica-
625 tion of procedure `LabelBoost` on the subsample in Step 3f of `GenericBBL` resulting in
626 S_2^b . By application of Claim 2.7 and Lemma C.1, M_2 is $(\epsilon, 0)$ -differentially private. Thus
627 (M_1, M_2) is $(\epsilon + 1, \delta)$ -differentially private. (iii) The mechanism M_3 with input of S^b and
628 output $(D_1^{b,L}, S_2^b) = (\text{View}_{\mathcal{B},R'_1}^b, S_2^b)$ applies (M_1, M_2) on the sub-sample S_1^b obtained from
629 S^b in Step 2 of `GenericBBL`. By application of Claim 2.7 M_3 is $(\epsilon, \frac{4\epsilon\delta}{3+\exp(\epsilon+1)})$ -differentially
630 private. Since $\frac{4\epsilon\delta}{3+\exp(\epsilon+1)} \leq \delta$ for any ϵ , hence $(\text{View}_{\mathcal{B},R'_1}^0, S_2^0)$ and $(\text{View}_{\mathcal{B},R'_1}^1, S_2^1)$ are (ϵ, δ) -
631 indistinguishable \square

632 **Privacy of the unlabeled points D .** Let $D^0, D^1 \in X^t$ be neighboring datasets of unla-
633 beled examples and let $\text{View}_{\mathcal{B},t}^0$ and $\text{View}_{\mathcal{B},t}^1$ be as in Figure 1 where $(x_1^0, \dots, x_t^0) = D^0$ and
634 $(x_1^1, \dots, x_t^1) = D^1$.

635 **Claim E.6.** For all adversaries \mathcal{B} , for all $t > 0$, and for any two neighbouring databases D^0 and
636 D^1 selected by \mathcal{B} , $\text{View}_{\mathcal{B},t}^0$ and $\text{View}_{\mathcal{B},t}^1$ are (ϵ, δ) -indistinguishable.

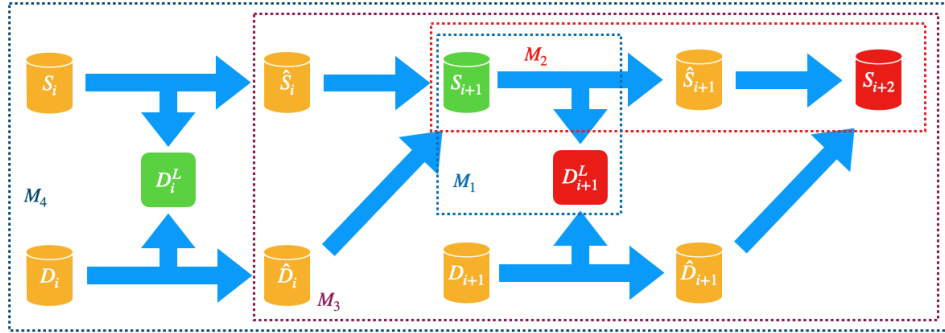


Figure 6: Privacy leakage of D_i

637 *Proof.* Let $D_1^0, D_2^0, \dots, D_k^0$ and $D_1^1, D_2^1, \dots, D_k^1$ be the set of unlabeled databases in step 3e of
638 `GenericBBL`. Without loss of generality, we assume D_i^0 and D_i^1 differ on one entry. When
639 $i = k$, $\text{View}_{\mathcal{B},t}^0 = \text{View}_{\mathcal{B},t}^1$ because all selected hypothesis are the same. When $i < k$, let
640 $R' = \min(\sum_{j=1}^{i+1} R_j, t)$.

641 Similar to the analysis if Claim E.5, $\text{View}_{\mathcal{B},t}^b$ results from a post-processing of the ran-
642 dom variable $(\text{View}_{\mathcal{B},R'}^b, S_{i+2}^b)$ (if $t < \sum_{j=1}^{i+1} R_j$ we set S_{i+2}^b as \perp). Note that $\text{View}_{\mathcal{B},R'_1}^b =$
643 $(D_1^{b,L}, \dots, D_i^{b,L^*}, D_{i+1}^{b,L})$, and $(D_1^{b,L}, \dots, D_{i-1}^{b,L}, D_i^{b,L^*})$ follow the same distribution for $b \in \{0, 1\}$,
644 where D_i^{b,L^*} is the labels of points in D_i^b expect the different point. So that it suffices to show
645 that $(D_{i+1}^{0,L}, S_2^0)$ and $(D_{i+1}^{1,L}, S_2^1)$ are (ϵ, δ) -indistinguishable.

646 We follow the processes creating $D_{i+1}^{b,L}$ and S_{i+2}^b in Figure 6: (i) The mechanism M_1 corre-
 647 sponds to the loop in Step 3d of `GenericBBL` where labels are produced for the adversarially
 648 chosen points D_{i+1}^b . By application of Lemma C.3, M_1 is $(1, \delta)$ -differentially private. (ii)
 649 The mechanism M_2 , corresponds to the subsampling of \hat{S}_{i+1}^b from S_{i+1}^b and the applica-
 650 tion of procedure `LabelBoost` on the subsample in Step 3f of `GenericBBL` resulting in
 651 S_{i+2}^b . By application of Claim 2.7 and Lemma C.1, M_2 is $(\varepsilon, 0)$ -differentially private. Thus
 652 (M_1, M_2) is $(\varepsilon + 1, \delta)$ -differentially private. (iii) The mechanism M_3 with input of \hat{D}_i^b and
 653 output $(D_{i+1}^{b,L}, S_{i+2}^b)$ applies (M_2, M_3) on S_{i+1} , which is generated from \hat{D}_i^b and in Step 3f
 654 of `GenericBBL`. By application of Claim C.1, M_3 is $(\varepsilon + 4, 4\varepsilon\delta)$ -differentially private. (iv)
 655 The mechanism M_4 , corresponds to the subsampling \hat{D}_i^b from D_i^b and the application of
 656 M_4 on \hat{D}_i^b . By application of Claim 2.7, M_4 is $(\varepsilon, \frac{16\varepsilon\delta}{3+\exp(\varepsilon+4)})$ -differentially private. Since
 657 $\frac{16\varepsilon\varepsilon}{3+\exp(\varepsilon+4)} \leq 1$ for any ε , $(D_{i+1}^{0,L}, S_2^0)$ and $(D_{i+1}^{1,L}, S_2^1)$ are (ε, δ) -indistinguishable. \square

658 **Remark E.7.** *The above proofs work on the adversarially selected D because: (i) Lemma C.3*
 659 *works on the adaptively selected queries. (We treat the hypothesis class F_i as the database, the*
 660 *unlabelled points $x_{i,\ell}$ as the query parameters.) (ii) `LabelBoost` generates labels by applying one*
 661 *private hypothesis on points. The labels are differentially private by post-processing.*