
Optimal Mistake Bounds for Transductive Online Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We resolve a 30-year-old open problem concerning the power of unlabeled data in
2 online learning by tightly quantifying the gap between transductive and standard
3 online learning. In the standard setting, the optimal mistake bound is characterized
4 by the Littlestone dimension d of the concept class \mathcal{H} (Littlestone, 1987). We prove
5 that in the transductive setting, the mistake bound is at least $\Omega(\sqrt{d})$. This consti-
6 tutes an exponential improvement over previous lower bounds of $\Omega(\log \log(d))$,
7 $\Omega(\sqrt{\log(d)})$ and $\Omega(\log(d))$, due respectively to Ben-David, Kushilevitz, and Man-
8 sour (1995, 1997) and Hanneke, Moran, and Shafer (2023). We also show that this
9 lower bound is tight: for every d , there exists a class of Littlestone dimension d
10 with transductive mistake bound $O(\sqrt{d})$. Our upper bound also improves upon the
11 best known upper bound of $(2/3) \cdot d$ from Ben-David et al. (1997). These results
12 establish a quadratic gap between transductive and standard online learning, thereby
13 highlighting the benefit of advance access to the unlabeled instance sequence. This
14 contrasts with the PAC setting, where transductive and standard learning exhibit
15 similar sample complexities.

16 1 Introduction

17 The transductive model is a basic and well-studied framework in learning theory, dating back to
18 the early works of Vapnik. It has been investigated both in statistical and online settings, and is
19 motivated by the principle that to make good predictions on a specific set of test instances, one need
20 not construct a fully general classifier that performs well on the entire domain — including points
21 that may never actually appear. Rather, it may be sufficient to tailor predictions for a fixed, known set
22 of instances.

23 This perspective naturally connects to a broader question in learning theory: what is the value
24 of unlabeled data? In the transductive setting, the learner is given the sequence of unlabeled test
25 instances in advance and is then required to predict their labels one by one. Thus, the transductive
26 model can be viewed as a natural formalization of learning with unlabeled data: the test instances are
27 known in advance, but their labels are not. The central question is whether such prior access to the
28 unlabeled sequence can help reduce the number of prediction mistakes — compared to the standard
29 online model, where the instances arrive and are labeled one at a time.

30 Recall for instance that in the standard PAC¹ model of supervised learning, there are cases where
31 access to unlabeled data is not helpful. Indeed, the “hard population distributions” used to prove the

¹Probably Approximately Correct. For an exposition of the standard terminology and results mentioned in this paragraph see, e.g., Shalev-Shwartz and Ben-David (2014).

standard VC² lower bound are constructed by taking a fixed and known marginal distribution over a VC-shattered set. Namely, the cases that are hardest to learn in the PAC setting include ones where the learner knows the marginal distribution over the domain, and can therefore generate as much unlabeled data as it wishes. And yet, in those cases, access to unlabeled data provides no acceleration compared to an algorithm (like ERM³) that does not use unlabeled data.

Seeing as unlabeled data is often a lot easier to obtain than labeled data, there have been considerable efforts to understand when and to what extent can access to unlabeled data accelerate learning.⁴

In particular, it is natural to ask, for which plausible models of learning is access to unlabeled data beneficial? Online learning (Littlestone, 1987) is perhaps the model of learning that is most-extensively studied in learning theory after the PAC model and its variants. Therefore, the general question considered in this paper is:

Question 1. *Quantitatively, how much (if at all) is access to unlabeled data beneficial for learning in the online learning setting?*

This question is naturally instantiated by comparing *transductive* online learning — where the learner has advance access to the full sequence x_1, x_2, \dots, x_n of unlabeled instances — with *standard* online learning, where no such access is given. This perspective has also been adopted in prior work: for example, Kakade and Kalai (2005), Cesa-Bianchi and Shamir (2013), and Hoi, Sahoo, Lu, and Zhao (2021) (Section 7.3) all describe transductive online learning as a setting in which the learner has access to “unlabeled data”. We thus refine the question above as follows:

Question 2. *Quantitatively, how much (if at all) is learning in the transductive online learning setting easier than learning in the standard online learning setting? Specifically, how much is the optimal number of mistakes in the transductive setting smaller than in the standard setting?*

Addressing this question, our main result (Theorem 1.1) states that the optimal number of mistakes in the transductive setting (with access to unlabeled data) is at most quadratically smaller than in the standard setting (without unlabeled data). Furthermore, there are hypothesis classes for which a quadratic gap is achieved.

1.1 Setting: Standard vs. Transductive Online Learning

Standard online learning (Littlestone, 1987) is a zero-sum, perfect- and complete-information game played over n rounds between two players, a *learner* and an *adversary*. The game is played with respect to a domain set \mathcal{X} and a hypothesis class $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$ (consisting of functions $\mathcal{X} \rightarrow \{0, 1\}$), where n , \mathcal{X} and \mathcal{H} are fixed and known to both players. The game proceeds as follows.

For each round $t = 1, 2, \dots, n$:

- a. The adversary selects an instance $x_t \in \mathcal{X}$ and sends it to the learner.
- b. The learner selects a prediction $\hat{y}_t \in \{0, 1\}$ and sends it to the adversary.
- c. The adversary selects a label $y_t \in \{0, 1\}$ and sends it to the learner. The selected label must be *realizable*, meaning that $\exists h \in \mathcal{H} \forall i \in [t]: h(x_i) = y_i$.

Game 1: The standard online learning setting.

The number of mistakes for a learner L and an adversary A is $M_{\text{std}}(\mathcal{H}, n, L, A) = |\{t \in [n] : \hat{y}_t \neq y_t\}|$. We are interested in understanding the *optimal number of mistakes*, which is

$$M_{\text{std}}(\mathcal{H}) = \sup_{n \in \mathbb{N}} \sup_{A \in \mathcal{A}} \inf_{L \in \mathcal{L}} M_{\text{std}}(\mathcal{H}, n, L, A),$$

²Vapnik–Chervonenkis.

³Empirical Risk Minimization.

⁴The literature on semi-supervised learning is surveyed in Joachims (1999); Zhu (2005); Zhu and Goldberg (2009); Zhu (2010); Chapelle, Schölkopf, and Zien (2006). Theoretical works on the topic include Benedek and Itai (1991); Blum and Mitchell (1998); Ben-David, Lu, Pál, and Sotáková (2008); Balcan and Blum (2010); Darnstädt, Simon, and Szörényi (2013); Göpfert, Ben-David, Bousquet, Gelly, Tolstikhin, and Urner (2019).

65 where \mathcal{A} and \mathcal{L} are the set of all deterministic adversaries and learners, respectively.⁵
 66 The *transductive* online learning setting (Ben-David et al., 1995, 1997) is similar, except that the
 67 learner has access to the full sequence of unlabeled instances in advance. Namely,

The adversary selects a *sequence* $x_1, x_2, \dots, x_n \in \mathcal{X}$ and sends it to the learner.

For each round $t = 1, 2, \dots, n$:

- a. The learner selects a *prediction* $\hat{y}_t \in \{0, 1\}$ and sends it to the adversary.
- b. The adversary selects a *label* $y_t \in \{0, 1\}$ and sends it to the learner. The selected label must be *realizable*, meaning that $\exists h \in \mathcal{H} \forall i \in [t]: h(x_i) = y_i$.

Game 2: The transductive online learning setting.

68 The optimal number of mistakes for the transductive setting is defined exactly as before,

$$M_{\text{tr}}(\mathcal{H}, n, L, A) = |\{t \in [n] : \hat{y}_t \neq y_t\}|, \text{ and } M_{\text{tr}}(\mathcal{H}) = \sup_{n \in \mathbb{N}} \sup_{A \in \mathcal{A}} \inf_{L \in \mathcal{L}} M_{\text{tr}}(\mathcal{H}, n, L, A),$$

69 with the only difference between the standard quantity $M_{\text{std}}(\mathcal{H})$ and the transductive quantity $M_{\text{tr}}(\mathcal{H})$
 70 being in how the game is defined.

71 1.2 Main Result

72 Notice that for every hypothesis class \mathcal{H} , $M_{\text{tr}}(\mathcal{H}) \leq M_{\text{std}}(\mathcal{H})$. Indeed, in the transductive setting
 73 the adversary declares the sequence x at the start of the game. This reduces the number of mistakes
 74 because the transductive adversary is less powerful (it cannot adaptively alter the sequence mid-game),
 75 and also because the transductive learner is more powerful (it has more information).⁶

76 While for some classes $M_{\text{tr}}(\mathcal{H}) = M_{\text{std}}(\mathcal{H})$, we study the largest possible separation. The best
 77 previous lower bound on M_{tr} , due to Hanneke, Moran, and Shafer (2023), states that for every
 78 class \mathcal{H} ,

$$M_{\text{tr}}(\mathcal{H}) \geq \Omega(\log(d)),$$

79 where $d = M_{\text{std}}(\mathcal{H})$. In the other direction, Ben-David et al. (1997) constructed⁷ a class \mathcal{H} such that
 80 $M_{\text{std}}(\mathcal{H}) = d$ and

$$M_{\text{tr}}(\mathcal{H}) \leq \frac{2}{3}d.$$

81 This left an exponential gap between the best known lower and upper bounds on M_{tr} , namely $\Omega(\log d)$
 82 versus $\frac{2}{3}d$. Our main result closes this gap:

83 **Theorem 1.1** (Main result).

- 84 • For every hypothesis class $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$,

$$M_{\text{tr}}(\mathcal{H}) = \Omega(\sqrt{d}),$$

85 where $d = M_{\text{std}}(\mathcal{H})$.

- 86 • On the other hand, for every d there exists a hypothesis class \mathcal{H} with $M_{\text{std}}(\mathcal{H}) = d$ and

$$M_{\text{tr}}(\mathcal{H}) = O(\sqrt{d}).$$

87 This result is stated in considerably greater detail in Theorems B.1 and D.1.

⁵Because the adversary selects y_t *after* seeing \hat{y}_t , randomness is not beneficial for either party, and we assume without loss of generality that both the learner and the adversary are deterministic. As is common in learning theory, we avoid questions of computability and allow the learner and adversary to be any function. See Appendix A for formal definitions of \mathcal{A} and \mathcal{L} .

⁶One could also define an intermediate setting, where the adversary is less powerful because it must select the sequence at the start of the game and cannot change it during the gameplay, but the learner does not have more information because the adversary only reveals the instances in the sequence one at a time as in the standard setting. However, this intermediate setting would not model the learner having *access* to unlabeled data.

⁷Their class consists of all disjoint unions of $\Theta(d)$ functions from a specific constant-sized class.

88 1.3 Related Works

89 The notion of *transductive inference* as a more efficient alternative to *inductive inference* in statistical
 90 learning theory was introduced by Vapnik (1979, 2006); Gammernan, Vovk, and Vapnik (1998);
 91 Chapelle, Vapnik, and Weston (1999). The *online learning* setting is due to Littlestone (1987), who
 92 also proved that the optimal number of mistakes is characterized by the Littlestone dimension.

93 The *transductive online learning* setting studied in the current paper, was first defined by Ben-
 94 David, Kushilevitz, and Mansour (1995), who used the name *worst sequence off-line model*. Among
 95 other results, they showed a lower bound of $\Omega(\log \log(d))$ on the number of mistakes required to
 96 learn a class with Littlestone dimension d . The authors subsequently presented an exponentially
 97 stronger lower bound of $\Omega(\sqrt{\log(d)})$ in Ben-David, Kushilevitz, and Mansour (1997). However,
 98 understanding where the optimal number of mistakes is situated within the range $\left[\Omega(\sqrt{\log(d)}), d\right]$
 99 remained an open question.

100 Kakade and Kalai (2005) presented an oracle-efficient algorithm for the transductive online learning
 101 setting, and may have been the first to use that name. Their result was subsequently improved upon
 102 by Cesa-Bianchi and Shamir (2013).

103 The present work is most similar to that of Hanneke, Moran, and Shafer (2023) which, among other
 104 results, gave a quadratically-stronger mistake lower bound of $\Omega(\log(d))$ for classes with Littlestone
 105 dimension d in the transductive online setting. The proof of our lower bound utilizes some of their
 106 ideas, but yields a quantitative improvement by combining it with some new ideas.

107 Hanneke, Raman, Shaeiri, and Subedi (2024) studied a setting of *multi-class* transductive online
 108 learning where the number of possible labels is unbounded.

109 2 Technical Overview

110 In this section we explain some of the main ideas in our proofs. Formal definition appear in
 111 Appendix A. Full formal statements of the results, as well as detailed rigorous proofs, appear in
 112 Appendices B to D.

113 2.1 Paths in Trees

114 We make extensive use of the following notion. Given a perfect binary tree T_d of depth d , every
 115 function $f : T_d \rightarrow \{0, 1\}$ defines a unique *path* in the tree. The path is a sequence of nodes
 116 $\text{path}(f) = (x_{i_0}, x_{i_1}, \dots, x_{i_d})$, as explained in Figure 1c. See Appendix A for formal definitions.

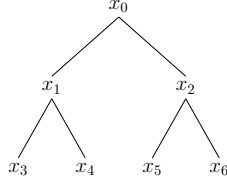
117 2.2 Proof Ideas for the Lower Bound

118 We start with an elementary observation about the adversary’s dilemma in the transductive online
 119 learning setting. Before round t of the game, the adversary selected a full sequence of instances
 120 $x_1, x_2, \dots, x_n \in \mathcal{X}$, and assigned some initial labels $y_1, y_2, \dots, y_{t-1} \in \{0, 1\}$. At the start of round
 121 t , the adversary must consider the *version space*,

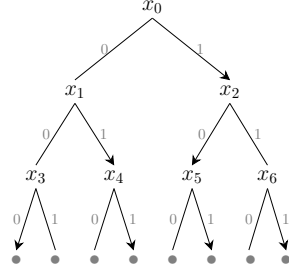
$$\mathcal{H}_t = \{h \in \mathcal{H} : (\forall i \in [t-1] : h(x_i) = y_i)\}.$$

122 If all $h \in \mathcal{H}_t$ assign $h(x_t) = b$ for some $b \in \{0, 1\}$, then the adversary has no choice but to assign
 123 the label $y_t = b$. Otherwise, the adversary can *force a mistake* at time t . Namely, after seeing the
 124 learner’s prediction \hat{y}_t , the adversary can assign $y_t = 1 - \hat{y}_t$, incrementing the number of learner
 125 mistakes by 1.

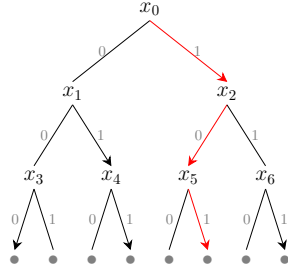
126 But “just because you can, doesn’t mean you should”. If the adversary is greedy and forces a mistake
 127 at time t , they may pay dearly for that later. As an extreme example, consider the case where there
 128 is a single $h_1 \in \mathcal{H}_t$ that assigns $h_1(x_t) = 1$, and all other functions $h \in \mathcal{H}_t$ assign $h(x_t) = 0$. If
 129 the adversary forces a mistake at time t , the version space at all subsequent times $s > t$ will be
 130 $\mathcal{H}_s = \{h_1\}$, and the adversary will be prevented from forcing any further mistakes.



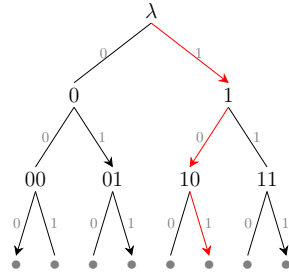
(a) A *perfect binary tree* of depth 2. Each *node* is labeled by an element of the domain \mathcal{X} . These labels need not be distinct (e.g., it is possible that $x_1 = x_6$). x_0 is the *root* of the tree, x_0, x_1 and x_2 are *internal nodes*, and x_3, \dots, x_6 are the *leaves*.



(b) A function $f : \mathcal{X} \rightarrow \{0, 1\}$ assigns a binary label to each node in the tree, represented here by edges with arrowhead tips. This figure depicts the function $f(x_i) = \mathbb{1}(i \notin \{2, 3\})$. (Note that the gray dots (•) in the figure are purely a pictorial detail. In this paper they are not considered nodes or leaves of the tree.)



(c) Every function $f : \mathcal{X} \rightarrow \{0, 1\}$ defines a *path* in the tree, which is a sequence $u_0, u_1, u_2, \dots, u_{d-1}$, where u_0 is the root, d is the depth of the tree, and for each $i \in [d-1]$, u_i is the b -child of u_{i-1} with $b = f(u_{i-1}) \in \{0, 1\}$. This figure shows that the function f from Figure 1b has $\text{path}(f) = (x_0, x_2, x_5)$, depicted in *red*. In particular, x_2 is ‘on-path’ for f , but x_6 is ‘off-path’ for f .



(d) In this paper we use a naming convention where, without loss of generality, we identify the domain elements x_i that are assigned to nodes with bit strings. The root is identified with the empty string λ , and for each pair of nodes u, v such that u is the b -child of v (for $b \in \{0, 1\}$), we have $u = v \circ b$, where ‘ \circ ’ denotes string concatenation. (Because the x_i ’s may not be distinct, a domain element may be identified with more than one bit string.)

Figure 1: Paths in trees.

131 A natural strategy for the adversary is therefore to be greedy up to a certain limit. Namely, at each
 132 time t the adversary computes the ratio⁸

$$r_t = \frac{|\{h \in \mathcal{H}_t : h(x_t) = 1\}|}{|\mathcal{H}_t|}.$$

133 If $r_t \in [\varepsilon, 1 - \varepsilon]$ for some parameter $\varepsilon > 0$ (“the version space is not too unbalanced”), then
 134 the adversary forces a mistake. Otherwise, the adversary assigns the majority label, i.e., $y_t =$
 135 $\mathbb{1}(r_t \geq 1/2)$. This ensures that the version space does not shrink too fast:

- 136 • If no mistake is forced, then $|\mathcal{H}_{t+1}| \geq (1 - \varepsilon) \cdot |\mathcal{H}_t|$, and
- 137 • If a mistake is forced, $|\mathcal{H}_{t+1}| \geq \varepsilon \cdot |\mathcal{H}_t|$.

138 In particular, at the end of the game, the version space \mathcal{H}_{n+1} is of size

$$|\mathcal{H}_{n+1}| \geq \varepsilon^M \cdot (1 - \varepsilon)^{n-M} \cdot |\mathcal{H}| = \varepsilon^M \cdot (1 - \varepsilon)^n \cdot 2^d, \quad (1)$$

139 where M is the number of mistakes that the adversary forces and n is the length of the sequence.

⁸For a class \mathcal{H} of Littlestone dimension d , the adversary will use only a subset of \mathcal{H} of cardinality 2^d that shatters a Littlestone tree of depth $d - 1$. So without loss of generality, we assume that \mathcal{H} has cardinality 2^d (in particular, \mathcal{H} is finite), and the ratio is well-defined.

2.3.2 Construction of the Hypothesis Class

We now describe a construction of a hypothesis class that is easy to learn in the transductive setting, using the idea of a sparse encoding. Recall that a class \mathcal{H} has Littlestone dimension at least d (Definition A.6 in Appendix A) if there exists a Littlestone tree of depth $d - 1$ such that for every $b \in \{0, 1\}^d$ there exists $h = h_b \in \mathcal{H}$ such that the values on the path of h agree with b . More formally, $\forall i \in [d] : h(b_{\leq i}) = b_i$, and in particular $\text{path}(h) = (\lambda, b_{\leq 1}, b_{\leq 2}, b_{\leq 3}, \dots, b_{\leq d-1})$. Thus, when constructing a class that shatters a specific Littlestone tree of depth $d - 1$, we need to define 2^d functions $\{h_b : b \in \{0, 1\}^d\}$. For each function h_b , the on-path values of the function are fixed (fully determined by b), while for the remaining values there is complete freedom (for the nodes u that are off-path we may assign any values $h_b(u) \in \{0, 1\}$).

Perhaps the simplest way to construct a class of Littlestone dimension d is simply to assign all on-path values as required, and assign 0 to all other values. Namely, if u is a prefix of b then $h_b(u) = b_{|u|+1}$, and otherwise $h_b(u) = 0$. In a sense, this is the ‘minimal’ class of Littlestone dimension d for a specific Littlestone tree.¹¹

Observe that the ‘minimal’ class does not have the desired property of being easy to learn in the transductive setting.¹² However, a certain variation of the ‘minimal’ class that embeds a sparse encoding does satisfy the requirement. In this variation, on-path value of the function h_b are assigned as they must (as determined by b), while the off-path values are sampled independently using a biased coin, such that each of them is 0 with high probability, but has a small probability of being 1. The probability is chosen carefully so that the class satisfies some simple combinatorial properties, as described further in Section 2.3.6 and Lemma D.2.

2.3.3 Naïve Learning Strategy

We now explain in broad strokes how the probabilistic construction of the hypothesis class in Section 2.3.2 is useful for learning with few mistakes in the transductive setting.

Notice that when predicting labels for the ‘minimal’ class with nodes in breadth-first order, the learner knows at each step whether they are labeling an on-path or off-path node, because the learner has already seen the correct labels for all ancestors of the current node. For off-path nodes, the learner knows that the true label is 0, so it never makes mistakes on off-path nodes, but it also gains no new information when the true labels for off-path nodes are revealed. No risk, but no reward either. Instead, all the information about the true labeling function is revealed only at on-path nodes, where the adversary has complete freedom to assign labels and force mistakes. That’s why the adversary can force d mistakes.

For the randomly-chosen class, when predicting labels for off-path nodes, the learner may still safely predict a label of 0. But the reasoning for this is quite different. Conceptually, every off-path label is part of a sparse codeword that identifies the correct labeling function.¹³ Because the coin is biased, each bit of the codeword is easy to guess (it is likely to be 0), but every time that the adversary reveals that the true label for an off-path node is indeed 0, the learner gains a small (nonzero) amount of information about the true labeling function. Additionally, when the adversary selects an off-path label of 1, that reveals a lot of information about the true labeling function (such labels are rare in the hypothesis class), and therefore the adversary cannot force many off-path mistakes. Overall, the information about the true labeling function is ‘smeared’ throughout all labels of the tree (0s and 1s, on-path and off-path).¹⁴

¹¹More formally, this is a class with a minimal number of nodes labeled 1.

¹²The adversary can declare a sequence x consisting of all the nodes in the tree in breadth-first order, and then force d mistakes — one mistake in each layer (depth) of the tree. Specifically, regardless of how the adversary selects the labels, for each $i \in [d]$ there exists a node u_i at depth i that is on-path. When it is time for the learner to predict a label for this u_i , the learner knows that u_i is on-path because it has seen the correct labels for all the ancestors of u_i . However, the adversary has the freedom to extend the path arbitrarily to the left or to the right, and can therefore force a mistake on u_i .

¹³The coin-flips for off-path labels are all independent. For example, if X is a set of nodes all of which are off-path for a subset H of the hypothesis class, then the random variables $\{h(x) : h \in H, x \in X\}$ are i.i.d.

¹⁴Furthermore, the labels for most not-too-small subsets of the nodes reveal a lot of information about the correct labeling function — not just for a particular subset of nodes. These properties led us to code-name this construction while working on the paper as ‘everything everywhere all at once’ (in reference to a 2022 film of

Thus, the naïve general strategy for the learner when using the probabilistically-constructed class is to learn most of the information about the true labeling function by observing off-path labels. By the time the learner reaches an on-path node, it hopefully has already learned enough about the true labeling function in order to make a good prediction on that node.

However, making this general strategy work requires overcoming some very substantial obstacles:

1. Recall that in the transductive setting, the adversary can present the nodes of the tree in any order of its choosing — it does not have to present the tree in breadth-first order. The naïve strategy works only if the learner sees many off-path nodes before it sees most on-path nodes. But what happens if the adversary decides to present many on-path nodes near the beginning of the sequence? To handle this, the learner incorporates a strategy we call ‘danger zone minimization’, as described in Section 2.3.4.
2. Another, equally problematic, issue also arises from the fact that the sequence presented by the adversary might not be in breadth-first order. Recall that breadth-first order¹⁵ has the property that for every node u in the sequence, all the ancestors of u appear *before* u in the sequence. This means that by the time the learner needs to predict a label for u , the learner knows whether u is on-path or off-path for the true labeling function. But what happens if the adversary presents u before some of u ’s ancestors? Or omits some of u ’s ancestors from the sequence altogether? In this case the learner doesn’t know if u is on-path or off-path, and this presents a double hazard. One hazard is that the learner doesn’t know what label to predict for u — if u is off-path, the learner can simply predict 0, but if it is on-path it must do something more elaborate. The second hazard is that, after seeing the correct label for u , it is not clear what the learner can infer from it. If u is off-path, its label should be interpreted as part of a sparse encoding of the labeling function. But if u is on-path, the interpretation must be entirely different. To overcome this challenge, the learner incorporates a strategy we call ‘splitting experts’, described in Section 2.3.5.
3. Limiting off-path mistakes. Thanks to the coin’s bias, most off-path nodes have a true label of 0. Nonetheless, each function in the hypothesis class still has an expected number of $2^{\Omega(d)}$ off-path nodes labeled 1, so the learner can afford to misclassify only a vanishing fraction of them! To limit the number of mistakes, the learner extracts information from the sparse encoding and executes a ‘transition to Halving’ strategy, as described in Section 2.3.6.

2.3.4 Danger Zone Minimization

Utilizing information from the ‘sparse encoding’ of the off-path nodes to make good predictions on on-path nodes requires that the learner first see the true labels for many off-path nodes. Until that happens, the learner expects to make many mistakes on on-path nodes. However, whether a node is on-path or off-path is not fixed in advanced — the adversary may decide this adaptively, in response to the learners predictions.

Danger zone minimization is a strategy used by the learner, to force the adversary to assign few nodes in the beginning of the sequence as on-path (otherwise, if initial nodes are assigned to be on-path by the adversary, then the learner will make few mistakes on those nodes). This is analogous to the standard Halving algorithm (Algorithm 7), but instead of minimizing the cardinality of the set of consistent hypotheses (the ‘version space’), the learner minimizes a subset of the domain (the ‘danger zone’).

Concretely, at the beginning of the game the learner initializes a set $S = \{x_1, x_2, \dots, x_{t_{\max}}\}$ consisting of the first $t_{\max} = 2^{\Omega(\sqrt{d})}$ instances in the sequence x selected by the adversary. This set represents the ‘danger zone’ — nodes in the beginning of the sequence that have not been labeled yet, that *might* be on-path, and that are not ancestors of a previously-labeled on-path node.¹⁶ To predict a label for an instance x_i , the learner selects a label \hat{y}_i such that if \hat{y}_i is wrong, the danger zone will shrink by at least $1/3$. That is, for $b \in \{0, 1\}$, if the set S_b of b -descendants of x_i has cardinality $|S_b| \geq |S|/3$, the

that name). This is in contrast to the ‘minimal’ function, where the information is concentrated entirely on the function path. The asymmetry between the ‘minimal’ class and the probabilistic class is similar to that between the binary and one-hot encodings in Section 2.3.1 above.

¹⁵As well as depth-first order.

¹⁶If u is an ancestor of some on-path node v , and v is a b -descendant of u for $b \in \{0, 1\}$, then the true label for u must be b .

266 learner predicts $\hat{y}_i = b$. Then, if the adversary selects $y_i = 1 - b$, that implies that all b -descendants
 267 of x_i are off-path for the true labeling functions. Therefore, the learner removes all b -descendants of
 268 x_i from the danger zone, and the new cardinality is $|S \setminus S_b| \leq (2/3) \cdot |S|$. This guarantees that the
 269 learner can make at most $O(\log(t_{\max})) = O(\sqrt{d})$ such mistakes before the danger zone is empty.¹⁷

270 If neither S_0 nor S_1 have cardinality at least $|S|/3$, the learner predicts $\hat{y}_i = 0$. If $y_i = 1$ and x_i is
 271 on-path for the true labeling function, then the learner updates the danger zone to be $S_0 \cup S_1$,¹⁸ again
 272 shrinking the danger zone by a factor of at most $2/3$. Otherwise, if $y_i = 1$ and x_i is off-path, then
 273 it was an off-path node labeled 1 (which is rare), and the learner can afford to misclassify it (see
 274 Section 2.3.6).

275 2.3.5 Splitting Experts

276 The danger zone minimization strategy requires that the learner know whether the node u being
 277 classified is on-path or off-path for the true labeling function. However, if u appears in the sequence
 278 before some of its ancestors, the learner does not know this. To overcome this difficulty, the learner
 279 implements a variant of the standard *multiplicative weights algorithm* using *splitting experts*. This
 280 means that initially there is a single expert executing danger zone minimization. When a node u is
 281 reached for which danger zone minimization requires knowing whether u is on-path or off-path and
 282 that information is not yet evident, each expert is split into two experts, one of which continues the
 283 execution of danger zone minimization under the assumption that u is on-path, and the other under
 284 the opposite assumption. Thus, at each point in time, there exists precisely one expert for which all
 285 path-related assumptions are correct, and therefore that expert will make at most $O(\sqrt{d})$ mistakes.
 286 The multiplicative weights algorithm guarantees that the overall number of mistakes will be linear in
 287 the the number of mistakes of the best expert, i.e., $O(\sqrt{d})$.

288 2.3.6 Transition to Halving

289 The hypothesis class is engineered such that it satisfies the following property: there are at most
 290 $2^{O(\sqrt{d})}$ functions in the hypothesis class that agree with any set of $t_{\max} = 2^{\Omega(\sqrt{d})}$ labels, or that agree
 291 that a set of $\Theta(\sqrt{d})$ nodes are all off-path and labeled 1 (this follows from Lemma D.2).

292 Therefore, once the true labels for the first t_{\max} instances $x_1, x_2, \dots, x_{t_{\max}}$ have been revealed, or once
 293 $\Theta(\sqrt{d})$ off-path labels of 1 have been revealed (whichever happens first), the learner can *transition*
 294 *to halving*: stop doing danger zone minimization, and instead predict the labels for the remaining
 295 nodes using the standard Halving algorithm (Algorithm 7) on the subset of the hypothesis class that
 296 survived. Halving on $2^{O(\sqrt{d})}$ functions is guaranteed to make at most $O(\sqrt{d})$ mistakes (Fact E.1).

297 However, seeing as the learner lacks information on which nodes are off-path, it uses experts, and
 298 each expert maintains different path-related assumptions. Thus, each expert decides separately at
 299 which point to transition to Halving. The unique expert that makes only correct assumptions will
 300 transition ‘at the right time’. That expert will make at most $O(\sqrt{d})$ mistakes during danger zone
 301 minimization, and then at most $O(\sqrt{d})$ additional mistakes during halving.

302 This completes our overview of the upper bound.

303 3 Organization

304 Complete rigorous mathematical details are deferred to the supplementary materials. Formal defi-
 305 nitions appear in Appendix A. Formal statements and proofs for the lower bound and upper bound

¹⁷Once the danger zone is empty, the learner cannot make any further on-path mistakes within the prefix $x_1, x_2, \dots, x_{t_{\max}}$. And it will make at most $O(\sqrt{d})$ mistakes on the remaining nodes $x_{t_{\max}+1}, x_{t_{\max}+2}, \dots$, as explained in Section 2.3.6.

¹⁸Because on-path nodes must be either be descendants or ancestors of x_i , and the definition of the danger zone does not require that it contain ancestors of nodes that have been labeled.

appear in Appendix B and Appendix D, respectively. Optimal sequence length is discussed in Appendix C.

References

- Maria-Florina Balcan and Avrim Blum. A discriminative model for semi-supervised learning. *J. ACM*, 57(3):19:1–19:46, 2010. doi:10.1145/1706591.1706599. URL <https://doi.org/10.1145/1706591.1706599>.
- Shai Ben-David, Eyal Kushilevitz, and Yishay Mansour. Online learning versus offline learning. In Paul M. B. Vitányi, editor, *Computational Learning Theory, Second European Conference, EuroCOLT '95, Barcelona, Spain, March 13-15, 1995, Proceedings*, volume 904 of *Lecture Notes in Computer Science*, pages 38–52. Springer, 1995. doi:10.1007/3-540-59119-2_167. URL https://doi.org/10.1007/3-540-59119-2_167.
- Shai Ben-David, Eyal Kushilevitz, and Yishay Mansour. Online learning versus offline learning. *Mach. Learn.*, 29(1):45–63, 1997. doi:10.1023/A:1007465907571. URL <https://doi.org/10.1023/A:1007465907571>.
- Shai Ben-David, Tyler Lu, Dávid Pál, and Miroslava Sotáková. Learning low-density separators. *CoRR*, abs/0805.2891, 2008. URL <http://arxiv.org/abs/0805.2891>.
- Gyora M. Benedek and Alon Itai. Learnability with respect to fixed distributions. *Theor. Comput. Sci.*, 86(2):377–390, 1991. doi:10.1016/0304-3975(91)90026-X. URL [https://doi.org/10.1016/0304-3975\(91\)90026-X](https://doi.org/10.1016/0304-3975(91)90026-X).
- Avrim Blum and Tom M. Mitchell. Combining labeled and unlabeled data with co-training. In Peter L. Bartlett and Yishay Mansour, editors, *Proceedings of the Eleventh Annual Conference on Computational Learning Theory, COLT 1998, Madison, Wisconsin, USA, July 24-26, 1998*, pages 92–100. ACM, 1998. doi:10.1145/279943.279962. URL <https://doi.org/10.1145/279943.279962>.
- Olivier Bousquet, Steve Hanneke, Shay Moran, Ramon van Handel, and Amir Yehudayoff. A theory of universal learning. In Samir Khuller and Virginia Vassilevska Williams, editors, *STOC 2021: 53rd Annual ACM SIGACT Symposium on Theory of Computing, Virtual Event, Italy, June 21-25, 2021*, pages 532–541. ACM, 2021. doi:10.1145/3406325.3451087. URL <https://doi.org/10.1145/3406325.3451087>.
- Nicolò Cesa-Bianchi and Ohad Shamir. Efficient transductive online learning via randomized rounding. In Bernhard Schölkopf, Zhiyuan Luo, and Vladimir Vovk, editors, *Empirical Inference - Festschrift in Honor of Vladimir N. Vapnik*, pages 177–194. Springer, 2013. doi:10.1007/978-3-642-41136-6_16. URL https://doi.org/10.1007/978-3-642-41136-6_16.
- Olivier Chapelle, Vladimir N. Vapnik, and Jason Weston. Transductive inference for estimating values of functions. In Sara A. Solla, Todd K. Leen, and Klaus-Robert Müller, editors, *Advances in Neural Information Processing Systems 12, [NIPS Conference, Denver, Colorado, USA, November 29 - December 4, 1999]*, pages 421–427. The MIT Press, 1999. URL <http://papers.nips.cc/paper/1699-transductive-inference-for-estimating-values-of-functions>.
- Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien, editors. *Semi-Supervised Learning*. The MIT Press, 2006. ISBN 9780262033589. doi:10.7551/MITPRESS/9780262033589.001.0001. URL <https://doi.org/10.7551/mitpress/9780262033589.001.0001>.
- Malte Darnstädt, Hans Ulrich Simon, and Balázs Szörényi. Unlabeled data does provably help. In Natacha Portier and Thomas Wilke, editors, *30th International Symposium on Theoretical Aspects of Computer Science, STACS 2013, February 27 - March 2, 2013, Kiel, Germany*, volume 20 of *LIPICs*, pages 185–196. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 2013. doi:10.4230/LIPICs.STACS.2013.185. URL <https://doi.org/10.4230/LIPICs.STACS.2013.185>.
- Alexander Gammernan, Volodya Vovk, and Vladimir N. Vapnik. Learning by transduction. In Gregory F. Cooper and Serafin Moral, editors, *UAI 1998: Proceedings of the Fourteenth Conference on Uncertainty in Artificial Intelligence, University of Wisconsin Business School, Madison, Wisconsin, USA, July 24-26, 1998*, pages 148–155. Morgan Kaufmann, 1998. URL https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=243&proceeding_id=14.

- Christina Göpfert, Shai Ben-David, Olivier Bousquet, Sylvain Gelly, Ilya O. Tolstikhin, and Ruth Urner. When can unlabeled data improve the learning rate? In Alina Beygelzimer and Daniel Hsu, editors, *Conference on Learning Theory, COLT 2019, 25-28 June 2019, Phoenix, AZ, USA*, volume 99 of *Proceedings of Machine Learning Research*, pages 1500–1518. PMLR, 2019. URL <http://proceedings.mlr.press/v99/gopfert19a.html>.
- Steve Hanneke, Shay Moran, and Jonathan Shafer. A trichotomy for transductive online learning. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine, editors, *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/3e32af2df2cd13dfbcbe6e8d38111068-Abstract-Conference.html.
- Steve Hanneke, Vinod Raman, Amirreza Shaeiri, and Unique Subedi. Multiclass transductive online learning. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang, editors, *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/6f244818d72b2a4be9b1225d1344e950-Abstract-Conference.html.
- Steven C. H. Hoi, Doyen Sahoo, Jing Lu, and Peilin Zhao. Online learning: A comprehensive survey. *Neurocomputing*, 459:249–289, 2021. doi:10.1016/J.NEUCOM.2021.04.112. URL <https://doi.org/10.1016/j.neucom.2021.04.112>.
- Thorsten Joachims. Transductive inference for text classification using support vector machines. In Ivan Bratko and Saso Dzeroski, editors, *Proceedings of the Sixteenth International Conference on Machine Learning (ICML 1999), Bled, Slovenia, June 27 - 30, 1999*, pages 200–209. Morgan Kaufmann, 1999.
- Sham M. Kakade and Adam Kalai. From batch to transductive online learning. In *Advances in Neural Information Processing Systems 18 [Neural Information Processing Systems, NIPS 2005, December 5-8, 2005, Vancouver, British Columbia, Canada]*, pages 611–618, 2005. URL <https://proceedings.neurips.cc/paper/2005/hash/17693c91d9204b7a7646284bb3adb603-Abstract.html>.
- Nick Littlestone. Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm. *Mach. Learn.*, 2(4):285–318, 1987. doi:10.1007/BF00116827. URL <https://doi.org/10.1007/BF00116827>.
- Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014. ISBN 978-1-10-705713-5. URL <http://www.cambridge.org/de/academic/subjects/computer-science/pattern-recognition-and-machine-learning/understanding-machine-learning-theory-algorithms>.
- Vladimir N. Vapnik. *Estimation of Dependencies Based on Empirical Data*. Nauka, Moscow, 1979. URL <https://www.ipu.ru/node/63854/publications>. In Russian.
- Vladimir N. Vapnik. *Estimation of Dependences Based on Empirical Data*. Springer, 2nd edition, 2006. ISBN 978-0-387-30865-4. doi:10.1007/0-387-34239-7. URL <https://doi.org/10.1007/0-387-34239-7>.
- Xiaojin Zhu. Semi-supervised learning literature survey. Technical report, Department of Computer Sciences, University of Wisconsin–Madison, 2005.
- Xiaojin Zhu. Semi-supervised learning. In Claude Sammut and Geoffrey I. Webb, editors, *Encyclopedia of Machine Learning*, pages 892–897. Springer, 2010. doi:10.1007/978-0-387-30164-8_749. URL https://doi.org/10.1007/978-0-387-30164-8_749.
- Xiaojin Zhu and Andrew B. Goldberg. *Introduction to Semi-Supervised Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2009. ISBN 978-3-031-00420-9. doi:10.2200/S00196ED1V01Y200906AIM006. URL <https://doi.org/10.2200/S00196ED1V01Y200906AIM006>.

410 Technical Appendices and Supplementary Material

411 A Preliminaries

412 A.1 Basic Notation

413 **Notation A.1.** $\mathbb{N} = \{1, 2, 3, \dots\}$, i.e., $0 \notin \mathbb{N}$. $\log(\cdot)$ and $\ln(\cdot)$ denote logarithm to base 2 and e ,
414 respectively.

415 **Notation A.2** (Sequences). Let \mathcal{X} be a set and $n, k \in \mathbb{N}$. For a sequence $x = (x_1, \dots, x_n) \in \mathcal{X}^n$,
416 we write $x_{\leq k}$ to denote the subsequence (x_1, \dots, x_k) . If $k \leq 0$ then $x_{\leq k}$ denotes the empty sequence,
417 which is also denoted by $\lambda = \mathcal{X}^0$. We use the notation $\mathcal{X}^{\leq n} = \cup_{i=0}^n \mathcal{X}^i$.

418 A.2 Standard Online Learning

419 Let \mathcal{X} be a set, and let $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$ be a collection of functions called a *hypothesis class*. A *learner*
420 *strategy* or simply *learner* for the standard online learning game (Game 1) is a function

$$L : \bigcup_{i=0}^{n-1} (\mathcal{X} \times \{0, 1\})^i \times \mathcal{X} \rightarrow \{0, 1\},$$

421 where $n \in \mathbb{N}$ is the number of rounds in the game. The set of all such learner strategies is denoted \mathcal{L}_n .
422 An *adversary strategy* or simply *adversary* for the standard online learning game is a pair of functions

$$A_{\text{instance}} : \bigcup_{i=0}^{n-1} (\mathcal{X} \times \{0, 1\} \times \{0, 1\})^i \rightarrow \mathcal{X}, \text{ and}$$

$$A_{\text{label}} : \bigcup_{i=1}^{n-1} (\mathcal{X} \times \{0, 1\} \times \{0, 1\})^i \times \{0, 1\} \rightarrow \{0, 1\}.$$

423 The set of all such adversary strategies is denoted \mathcal{A}_n .

424 Semantically, the interpretation of these strategies is that in each round $t \in [n]$ of Game 1, the
425 adversary selects an instance

$$x_t = A_{\text{instance}}(x_1, \hat{y}_1, y_1, \dots, x_{t-1}, \hat{y}_{t-1}, y_{t-1}) \in \mathcal{X},$$

426 then the learner makes a prediction

$$\hat{y}_t = L(x_1, y_1, \dots, x_{t-1}, y_{t-1}, x_t) \in \{0, 1\},$$

427 and finally, the adversary assigns a label

$$y_t = A_{\text{label}}(x_1, \hat{y}_1, y_1, \dots, x_{t-1}, \hat{y}_{t-1}, y_{t-1}, \hat{y}_t) \in \{0, 1\}.$$

428 The adversary's function A_{label} must satisfy *realizability*, meaning that there exists $h \in \mathcal{H}$ such that

$$\forall t \in [n] : y_t = h(x_t).$$

429 The number of mistakes is in a game with n rounds and hypothesis class \mathcal{H} between learner L and
430 adversary A is

$$M_{\text{std}}(\mathcal{H}, n, L, A) = |\{t \in [n] : \hat{y}_t \neq y_t\}|.$$

431 A.3 Transductive Online Learning

432 Given \mathcal{X} and \mathcal{H} as in Appendix A.2, a learner strategy for the *transductive online learning setting*
433 (Game 2) is a function

$$L : \mathcal{X}^n \times \bigcup_{i=0}^{n-1} \{0, 1\}^i \rightarrow \{0, 1\},$$

434 where $n \in \mathbb{N}$ is the number of rounds in the game. An adversary strategy consists of a sequence
435 $x \in \mathcal{X}^n$ and an *adversary labeling strategy*, which is a function

$$A : \left(\bigcup_{i=0}^{n-1} \{0, 1\}^{2i} \right) \times \{0, 1\} \rightarrow \{0, 1\}.$$

436 The sets of all such learner and adversary strategies are denoted \mathcal{L}_n and \mathcal{A}_n respectively.
 437 Semantically, the interpretation of these strategies is that at the start of Game 2, the adversary selects
 438 the sequence x . Then, in each round $t \in [n]$, the learner makes a prediction

$$\hat{y}_t = L(x, y_1, \dots, y_{t-1}) \in \{0, 1\},$$

439 and then the adversary assigns a label

$$y_t = A(\hat{y}_1, y_1, \dots, \hat{y}_{t-1}, y_{t-1}, \hat{y}_t) \in \{0, 1\}.$$

440 Exactly as in Appendix A.2, the adversary's function A must satisfy realizability, namely,

$$\exists h \in \mathcal{H} \forall t \in [n] : y_t = h(x_t),$$

441 and the number of mistakes is in a game with sequence length n and hypothesis class \mathcal{H} between
 442 learner L and adversary A is

$$M_{\text{tr}}(\mathcal{H}, n, L, A) = |\{t \in [n] : \hat{y}_t \neq y_t\}|.$$

443 A.4 Mistake Bounds

444 In this paper, we study *optimal mistake bounds*, or the *optimal number of mistakes*, which is the value
 445 of Games 1 and 2. For $M \in \{M_{\text{std}}, M_{\text{tr}}\}$, the optimal number of mistakes in a game with hypothesis
 446 class \mathcal{H} and sequence length n is,

$$M(\mathcal{H}, n) = \sup_{A \in \mathcal{A}_n} \inf_{L \in \mathcal{L}_n} M(\mathcal{H}, n, L, A).$$

447 The optimal number of mistakes for hypothesis class \mathcal{H} is

$$M(\mathcal{H}) = \sup_{n \in \mathbb{N}} M(\mathcal{H}, n).$$

448 **Remark A.3.** As is common in learning theory literature, in both Game 1 and Game 2, we take
 449 the sets \mathcal{L}_n and \mathcal{A}_n to be the sets of all (deterministic) functions. In this paper, we do not consider
 450 randomized strategies. By allowing arbitrary functions, we ignore issues relating to computability.

451 A.5 Trees

452 **Definition A.4** (Notation for binary trees). *Let $d \in \mathbb{N} \cup \{0\}$. A perfect binary tree of depth d is a*
 453 *collection of $2^{d+1} - 1$ nodes, which we identify with the collection of binary strings*

$$T_d = \{ \{0, 1\}^k : k \in \{0, 1, 2, \dots, d\} \}.$$

454 *The empty string, denoted $\lambda = \{0, 1\}^0$, is a member of T_d and is called the root of the tree. Every*
 455 *string $u \in \{0, 1\}^d$ is called a leaf. The depth of a node $u \in T_d$, denoted $|u|$, is the length of u as a*
 456 *string, namely, the integer k such that $u \in \{0, 1\}^k$.*

457 *For two nodes $u, v \in T_d$, we say that u is a parent of v , and that v is a child of u , if $v = u \circ 0$ or*
 458 *$v = u \circ 1$, where \circ denotes string concatenation. More fully, for $b \in \{0, 1\}$, we say that v is a b -child*
 459 *of u if $v = u \circ b$.*

460 *Recursively, we define that u is an ancestor of v and that v is a descendant of u , and write $u \preceq v$, if*
 461 *one of the following holds:*

462 $\bullet v = u \text{ function } \in T_d : u \preceq w \preceq v.$

463 *For $b \in \{0, 1\}$, we say that v is a b -descendant of u , denoted $u \preceq_b v$, if v is a descendant of the*
 464 *b -child of u .*

465 **Definition A.5** (Paths in a binary tree). *Let $d, k \in \mathbb{N}$, $k \leq d$. Let $u \in \{0, 1\}^k$ be a node in T_d .*
 466 *The path to u is the unique sequence $\text{path}(u) = (u_0, u_1, u_2, \dots, u_k)$ such that $u_0 = \lambda$ is the root,*
 467 *$u_k = u$, and u_i is a child of u_{i-1} for all $i \in [k]$.*

468 *Let $f : T_d \rightarrow \{0, 1\}$ be a function. The path of f is the unique sequence $\text{path}(f) =$
 469 *$(u_0, u_1, u_2, \dots, u_d)$ such that $u_0 = \lambda$ is the root, and for each $i \in [d]$, $u_i = u_{i-1} \circ f(u_{i-1})$.*
 470 *Namely, u_i is the $f(u_{i-1})$ -child of u_{i-1} .**

471 *For a hypothesis class $\mathcal{H} \subseteq \{0, 1\}^{T_d}$, we define $\text{path}(\mathcal{H}) = \text{path}(f_{\text{maj}})$, where $f_{\text{maj}} : T_d \rightarrow \{0, 1\}$*
 472 *is given by*

$$f_{\text{maj}}(u) = \mathbb{1} \left(\frac{|\{h \in \mathcal{H} : h(u) = 1\}|}{|\mathcal{H}|} \geq \frac{1}{2} \right).$$

473 A.6 Littlestone Dimension

474 **Definition A.6** (Littlestone, 1987). Let \mathcal{X} be a set, let $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$, and let $d \in \mathbb{N} \cup \{0\}$. We say
 475 that \mathcal{H} shatters the binary tree T_d if there exists a mapping $T_d \rightarrow \mathcal{X}$ given by $u \mapsto x_u$ such that for
 476 every $u \in \{0, 1\}^{d+1}$ there exists $h_u \in \mathcal{H}$ such that

$$\forall i \in [d + 1] : h(x_{u_{\leq i-1}}) = u_i.$$

477 The Littlestone dimension of \mathcal{H} , denoted $\text{LD}(\mathcal{H})$, is the supremum over all $d \in \mathbb{N}$ such that there
 478 exists a Littlestone tree of depth $d - 1$ that is shattered by \mathcal{H} .

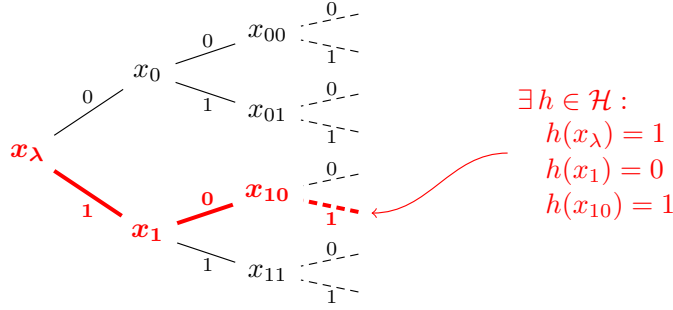


Figure 2: A shattered Littlestone tree of depth 2. The empty sequence is denoted by λ .

(Source: Bousquet et al., 2021)

479 **Theorem A.7** (Littlestone, 1987). Let \mathcal{X} be a set and let $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$ such that $d = \text{LD}(\mathcal{H}) < \infty$.
 480 Then there exists a strategy for the learner that guarantees that the learner will make at most d
 481 mistakes in the standard (non-transductive) online learning setting, regardless of the adversary's
 482 strategy and of the number n of instances to be labeled. Furthermore, there exists an adversary that
 483 forces every learner to make at least $\min\{n, d\}$ mistakes.

484 B Lower Bound

485 B.1 Statement

486 Our $\Omega(\sqrt{d})$ lower bound states the following.

487 **Theorem B.1** (Lower bound). There exists a constant $d_0 \geq 0$ as follows. Let $d \in \mathbb{N}$, $d \geq d_0$, let \mathcal{X}
 488 be a set, and let $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$ be a hypothesis class with $\text{LD}(\mathcal{H}) = d$. Then for every $k \in [d]$ there
 489 exist a sequence $x \in \mathcal{X}^n$ of length $n = O(k \cdot 2^{\sqrt{k}})$ such that for every learning rule L there exists
 490 $h \in \mathcal{H}$ such that

$$M_{\text{tr}}(L, x, h) \geq \sqrt{k}. \quad (3)$$

491 Moreover, this is witnessed by a simple adaptive labeling strategy for the adversary (as in Algorithm 1)
 492 that forces every learning rule to make at least \sqrt{k} mistakes on the (same) sequence x . Furthermore,
 493 for any integer $n \in \mathbb{N}$,

$$M_{\text{tr}}(\mathcal{H}, n) \geq \min \left\{ \sqrt{d}, \lfloor \log(n) \rfloor \right\}. \quad (4)$$

494 See Section 2.2 for a general overview of this result and the main proof ideas. In the following
 495 subsections we prove Theorem B.1. Algorithm 1 gives an explicit construction of the adversary that
 496 witnesses the lower bound, using Algorithm 2 as a subroutine. We start with presenting some initial
 497 observations about the behavior of these algorithms in Appendix B.2.

Assumptions:

- $d \in \mathbb{N}, \varepsilon = 2^{-\sqrt{d}/2}$.
- $T = T_d$ is a perfect binary tree of depth d .
- $\mathcal{H} \subseteq \{0, 1\}^T$ is a class that shatters T .

TRANSDUCTIVEADVERSARY (\mathcal{H}):

```

 $(x_1, x_2, \dots, x_n) \leftarrow \text{CONSTRUCTSEQUENCE}(\mathcal{H})$  ▷ See Algorithm 2.
send  $(x_1, x_2, \dots, x_n)$  to learner
 $\mathcal{H}_0 \leftarrow \mathcal{H}$ 
for  $t \in [n]$ :
    receive  $\hat{y}_t$  from learner
     $r_t \leftarrow \frac{|\{h \in \mathcal{H}_{t-1} : h(x_t) = 1\}|}{|\mathcal{H}_{t-1}|}$ 
     $y_{\text{maj}} \leftarrow \mathbb{1}(r_t \geq 1/2)$ 
     $y_t \leftarrow \begin{cases} y_{\text{maj}} & r_t \notin [\varepsilon, 1 - \varepsilon] \\ 1 - \hat{y}_t & \text{otherwise} \end{cases}$ 
    send  $y_t$  to learner
     $\mathcal{H}_t \leftarrow \{h \in \mathcal{H}_{t-1} : h(x_t) = y_t\}$ 

```

Algorithm 1: The strategy for the adversary that achieves the lower bound in Theorem B.1. Note that while the construction of the sequence x is not entirely trivial, the adversary's strategy for labeling this sequence is very simple.

498 **B.2 Analysis of the Adversary**

499 **Claim B.2.** Let $d \in \mathbb{N}$, let $M = \sqrt{d}/10$, and let $\mathcal{H} \subseteq \{0, 1\}^{T_d}$ be a hypothesis class. Consider an
500 execution of $\text{CONSTRUCTSEQUENCE}(\mathcal{H})$ as in Algorithm 2 that produces a sequence x_1, x_2, \dots, x_t .
501 Then:

- 502 (a) For all $i \in [t]$, $\text{path}(x_i)$ is a subsequence of x_0, x_1, \dots, x_i .
- 503 (b) The length t of the sequence satisfies $t < n_d$, where $n_d = (d + 1) \cdot 2^{M+1}$.

504 *Proof.*

- 505 (a) Fix $i \in [t]$. It suffices to show that for all $u \in T_d$, if $u \preceq x_i$ then $u \in (x_1, x_2, \dots, x_i)$.
506 Proceed by induction on i . For the base case $i = 1$, the claim holds because $x_1 = \lambda$.

507 For the induction step, assume the claim holds for $i \in [t - 1]$. Let $u \preceq x_{i+1}$, we prove that
508 $u \in (x_1, x_2, \dots, x_{i+1})$. Assume $x_{i+1} \neq \lambda$ (otherwise, there is nothing to prove).

509 Because x_{i+1} appears in the sequence x , it must have been added to \mathcal{Q} before it was added to
510 x . The only place where items that are not λ are added to \mathcal{Q} is in the line $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{x_t \circ y\}$.
511 Namely, there exist an index $j \in [i]$ and a bit $y \in \{0, 1\}$ such that $x_{i+1} = x_j \circ y$ (note
512 that $j < i + 1$ because x_j was added to the sequence before x_{i+1}). If $x_j = u$ we are
513 done. Otherwise, note that x_j is the parent of x_{i+1} , and therefore $u \preceq x_j$. By the induction
514 hypothesis, $u \in (x_1, x_2, \dots, x_j)$. This concludes the proof.

- 515 (b) Items are added to the the sequence x only if they were previously added to \mathcal{Q} . Therefore,
516 the length of the sequence x is $t = U + 1$, where U is the number of times that the line
517 “ $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{x_t \circ y\}$ ” was executed (we add 1 because λ is added to \mathcal{Q} elsewhere).

518 Consider a function f that maps a node u in the sequence x to the value of the index b' at
 519 the time that u was added to \mathcal{Q} . Namely, if $u = x_t \circ y$ and the index b' had some value β
 520 when the line “ $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{x_t \circ y\}$ ” was executed, then $f(u) = \beta$.

521 Notice that “ $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{x_t \circ y\}$ ” is executed only if the condition $x_t \in \text{path}(\mathcal{H}_{b'})$ is satisfied
 522 in the previous line. By Item (a), the line “ $\mathcal{H}_{b'} \leftarrow \{h \in \mathcal{H}_b : h(x_t) = y\}$ ”, and the fact that
 523 items in \mathcal{Q} are processed in lexicographic order, it follows that there exists some assignment
 524 of labels $\ell : \text{path}(x_t) \rightarrow \{0, 1\}$ such that

$$\mathcal{H}_{b'} \subseteq \{h \in \mathcal{H} : (\forall u \in \text{path}(x_t) : h(u) = \ell(u))\}.$$

Assumptions:

- $d \in \mathbb{N}$, $M = \sqrt{d}/10$, $\varepsilon = 2^{-\sqrt{d}/2}$.
- $T = T_d$ is a perfect binary tree of depth d .
- λ , the empty string, is the root of T .
- $\mathcal{H} \subseteq \{0, 1\}^T$ is a class that shatters T .

CONSTRUCTSEQUENCE(\mathcal{H}):

```

 $\mathcal{H}_\lambda \leftarrow \mathcal{H}$ 
 $\mathbb{H}_0 \leftarrow \{\mathcal{H}_\lambda\}$  ▷ A set of classes indexed by bit strings.
 $\mathcal{Q} \leftarrow \{\lambda\}$  ▷ A set of nodes to be processed.
 $t \leftarrow 0$ 
while  $|\mathcal{Q}| > 0$ :
   $t \leftarrow t + 1$ 
   $x_t \leftarrow \min_{\leq_{\text{lex}}} \mathcal{Q}$  ▷ Pop the lexicographically first string from  $\mathcal{Q}$  and add
   $\mathcal{Q} \leftarrow \mathcal{Q} \setminus \{x_t\}$  it to the output sequence.
   $\mathbb{H}_t \leftarrow \emptyset$ 
  for  $\mathcal{H}_b \in \mathbb{H}_{t-1}$ :
     $r \leftarrow \frac{|\{h \in \mathcal{H}_b : h(x_t) = 1\}|}{|\mathcal{H}_b|}$ 
     $\mathcal{Y} \leftarrow \begin{cases} \{0, 1\} & (r \in [\varepsilon, 1 - \varepsilon]) \wedge (|b| \leq M) \\ \{\mathbb{1}(r \geq 1/2)\} & \text{otherwise} \end{cases}$  ▷
    Adversary will force mistakes on the first  $M$  balanced nodes.
    for  $y \in \mathcal{Y}$ :
       $b' \leftarrow \begin{cases} b & |\mathcal{Y}| = 1 \\ b \circ y & |\mathcal{Y}| = 2 \end{cases}$  ▷ Restrict class to agree with  $y$ . If splitting the class in two to force a mistake then create new indices.
       $\mathcal{H}_{b'} \leftarrow \{h \in \mathcal{H}_b : h(x_t) = y\}$ 
       $\mathbb{H}_t \leftarrow \mathbb{H}_t \cup \{\mathcal{H}_{b'}\}$ 
      if  $x_t \in \text{path}(\mathcal{H}_{b'}) \wedge |x_t| < d$ : ▷ If  $x_t$  is on-path for  $\mathcal{H}_{b'}$  and it has a  $y$ -child, add that child to  $\mathcal{Q}$ .
         $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{x_t \circ y\}$ 
return  $(x_1, x_2, \dots, x_t)$ 

```

Algorithm 2: A subroutine of Algorithm 1 for selecting the sequence x .

Consequently, $x_t \in \text{path}(\mathcal{G})$ for any class \mathcal{G} that is a restriction of $\mathcal{H}_{b'}$; in particular, because the only way that $\mathcal{H}_{b'}$ might be modified later during the execution of Algorithm 2 is by restricting it further, it follows that $x_t \in \text{path}(\mathcal{H}_{b'})$ when the line “ $\mathcal{Q} \leftarrow \mathcal{Q} \cup \{x_t \circ y\}$ ” is executed and in all subsequent times.

However, $|\text{path}(\mathcal{G})| = d + 1$ for any class $\mathcal{G} \subseteq \{0, 1\}^{T_d}$. This implies that f maps at most $(d + 1)$ nodes to each bit string.¹⁹ In other words, for any bit string b , the size of the preimage satisfies $|f^{-1}(b)| \leq d + 1$.

Thus,

$$\begin{aligned}
t &= 1 + |\{2, 3, \dots, t\}| \\
&= 1 + \sum_{\substack{b \in \{0, 1\}^k \\ k \leq M}} |\{i \in \{2, 3, \dots, t\} : f(x_i) = b\}| \\
&= 1 + \sum_{\substack{b \in \{0, 1\}^k \\ k \leq M}} |f^{-1}(b)| \\
&\leq 1 + \sum_{\substack{b \in \{0, 1\}^k \\ k \leq M}} (d + 1) \\
&\leq 1 + (d + 1) \cdot (2^{M+1} - 1). \\
&< (d + 1) \cdot 2^{M+1},
\end{aligned}$$

as desired. \square

Claim B.3. Let $d \in \mathbb{N}$, and let $\mathcal{H} \subseteq \{0, 1\}^{T_d}$ be a hypothesis class. Consider an execution of TRANSDUCTIVEADVERSARY(\mathcal{H}) as in Algorithm 1. Let

$$\mathcal{H}_0, \mathcal{H}_1, \dots, \mathcal{H}_n$$

be the sequence of hypothesis classes created by TRANSDUCTIVEADVERSARY, and let

$$\mathbb{H}_0, \mathbb{H}_1, \dots, \mathbb{H}_n$$

be the sequence of collections created by the subroutine CONSTRUCTSEQUENCE (Algorithm 2).

Then

$$\forall t \in \{0, 1, \dots, n\} : \mathcal{H}_t \in \mathbb{H}_t.$$

Proof. Proceed by induction on $t \in \{0, 1, \dots, n\}$. The base case $t = 0$ is satisfied, because $\mathcal{H}_0 = \mathcal{H} \in \{\mathcal{H}\} = \mathbb{H}_0$. For the induction step, assume that $\mathcal{H}_{i-1} \in \mathbb{H}_{i-1}$ for some $i \in [n]$. We prove that $\mathcal{H}_i \in \mathbb{H}_i$.

Let y_i be the label assigned to x_i by TRANSDUCTIVEADVERSARY. Then

$$\mathcal{H}_i = \{h \in \mathcal{H}_{i-1} : h(x_i) = y_i\}.$$

Consider the iteration of the while loop in CONSTRUCTSEQUENCE that starts with $t \leftarrow i$. By the induction hypothesis, $\mathcal{H}_{i-1} \in \mathbb{H}_{i-1}$. Therefore, in this iteration of the while loop, there will be an iteration of the “for $\mathcal{H}_b \in \mathbb{H}_{t-1}$ ” loop where $\mathcal{H}_b = \mathcal{H}_{i-1}$. In that iteration, $y_i \in \mathcal{Y}$ by construction of y_i and \mathcal{Y} . Therefore, in the iteration of the “for $y \in \mathcal{Y}$ ” loop in which $y = y_i$,

$$\mathcal{H}_{b'} = \{h \in \mathcal{H}_b : h(x_i) = y\} = \{h \in \mathcal{H}_{i-1} : h(x_i) = y_i\} = \mathcal{H}_i.$$

The class $\mathcal{H}_{b'}$ is then added to $\mathbb{H}_i = \mathbb{H}_t$ in the line “ $\mathbb{H}_t \leftarrow \mathbb{H}_t \cup \{\mathcal{H}_{b'}\}$ ”. Furthermore, no class is ever removed from \mathbb{H}_t . So $\mathcal{H}_i \in \mathbb{H}_i$, as desired. \square

¹⁹Note that $x_t \circ y$ could be added to \mathcal{Q} twice: once with $y = 0$ and a second time with $y = 1$. However, when this happens, the index b is replaced by two indices of the form $b' = b \circ y$; for each such b' , there just is a single node of the form $x_t \circ y$ such that $f(x_t \circ y) = b'$.

Claim B.4. Let $d \in \mathbb{N}$, let $k \in \{0, 1, \dots, d\}$, and let $\mathcal{H} \subseteq \{0, 1\}^{T_d}$ be a hypothesis class. Consider an execution of `TRANSDUCTIVEADVERSARY` (\mathcal{H}) as in Algorithm 1 where the adversary constructs a sequence of nodes $x_1, x_2, \dots, x_n \in T_d$ and a sequence of classes $\mathcal{H}_0, \mathcal{H}_1, \dots, \mathcal{H}_n \subseteq \{0, 1\}^{T_d}$. Then there exists $i \in [n]$ such that

1. $|x_i| = k$, and
2. $x_i \in \text{path}(\mathcal{H}_{i-1})$.

Proof. Proceed by induction on k . For the base case $k = 0$, notice that $x_1 = \lambda$, $|\lambda| = 0$, and $\lambda \in \text{path}(\mathcal{H}_0)$.

For the induction step, assume the claim holds for some $k \in \{0, 1, \dots, d-1\}$, and let $i_k \in [n]$ such that $|x_{i_k}| = k$ and $x_{i_k} \in \text{path}(\mathcal{H}_{i_k-1})$; we prove that the claim holds for $k+1$ as well.

Consider the iteration of the while loop in `CONSTRUCTSEQUENCE` in which x_{i_k} is added to the sequence (i.e., the iteration starting with $t \leftarrow i_k$). By Claim B.3, $\mathcal{H}_{i_k-1} \in \mathbb{H}_{i_k-1}$. Hence, within this iteration of the while loop, there is an iteration of the “for $\mathcal{H}_b \in \mathbb{H}_{t-1}$ ” loop such that $\mathcal{H}_b = \mathcal{H}_{i_k-1}$. By construction, the set \mathcal{Y} always contains the label predicted by the adversary, so $y_{i_k} \in \mathcal{Y}$. Because $x_{i_k} \in \text{path}(\mathcal{H}_{i_k-1}) = \text{path}(\mathcal{H}_b)$, it follows that $x_{i_k} \in \text{path}(\mathcal{H}_{b'})$. Seeing as $|x_{i_k}| < d$, in the last line of the iteration of the “for $y \in \mathcal{Y}$ ” loop with $y = y_{i_k}$, the node $x_{i_{k+1}} := x_{i_k} \circ y_{i_k}$ is added to \mathcal{Q} . This guarantees that $x_{i_{k+1}}$ will eventually be popped from \mathcal{Q} and added to the sequence returned by `CONSTRUCTSEQUENCE`. Once a node has been added to the sequence, it is never removed.

Notice that $|x_{i_{k+1}}| = |x_{i_k}| + 1 = k+1$, satisfying Item 1. Therefore, it remains to show Item 2, namely, to show that $x_{i_{k+1}} \in \text{path}(\mathcal{H}_{i_{k+1}-1})$.

Let $(\lambda = u_0, u_1, u_2, \dots, u_k = x_{i_k}) = \text{path}(x_{i_k})$ be the root-to-node path to x_{i_k} . By Item (a) in Claim B.2, $\text{path}(x_{i_k})$ is a subsequence of $x_{\leq i_k}$. By the construction of the version space \mathcal{H}_{i_k} , \mathcal{H}_{i_k} is restricted to agree with the root-to-node path to x_{i_k} , as well as with the label y_{i_k} .²⁰ So \mathcal{H}_{i_k} is restricted to agree with the root-to-node path to $x_{i_k} \circ y_{i_k} = x_{i_{k+1}}$. This implies that,

$$\forall \mathcal{G} \subseteq \mathcal{H}_{i_k} : x_{i_{k+1}} \in \text{path}(\mathcal{G}).$$

In particular, because $\mathcal{H}_{i_{k+1}-1} \subseteq \mathcal{H}_{i_k}$, it follows that $x_{i_{k+1}} \in \text{path}(\mathcal{H}_{i_{k+1}-1})$, as desired. \square

B.3 Proof

Claim B.5. Let $d \in \mathbb{N}$, $d \geq 800$, and let $\mathcal{H} \subseteq \{0, 1\}^{T_d}$ with cardinality $|\mathcal{H}| = 2^{d+1}$ be a hypothesis class that shatters T_d . Consider an execution of `TRANSDUCTIVEADVERSARY` (\mathcal{H}) as in Algorithm 1. Let x_1, \dots, x_n and y_1, \dots, y_n be nodes and labels selected by the adversary, let $\mathcal{H}_0, \mathcal{H}_1, \dots, \mathcal{H}_n$ be the hypothesis classes defined by the adversary, and let $\hat{y}_1, \dots, \hat{y}_n$ be the labels selected by the learner.

Let

$$S = \{s_0, s_1, s_2, \dots\} = \{t \in [n] : r_t \in [\varepsilon, 1 - \varepsilon]\}$$

be the set of indices where the adversary forces a mistake.

For each $k \in \{0, 1, 2, \dots, \lfloor \sqrt{d} \rfloor\}$, define

$$s_{\max}(k) = \max \left\{ t \in [n] : |x_t| \leq k\sqrt{d} \right\}.$$

Then $|S| \geq \lfloor \sqrt{d} \rfloor + 1$, and for all $k \in \{0, 1, 2, \dots, \lfloor \sqrt{d} \rfloor\}$:

1. $s_k \leq s_{\max}(k)$, and
2. $|\mathcal{H}_{s_k}| \geq 2^{d-k\sqrt{d}}$.

²⁰Formally, there exists a mapping $j : \text{path}(x_{i_k}) \rightarrow [i_k]$ such that for each $u \in \text{path}(x_{i_k})$, $u = x_{j(u)}$, and $\mathcal{H}_{i_k} \subseteq \{h \in \mathcal{H} : (\forall u \in \text{path}(x_{i_k}) : h(u) = y_{j(u)})\}$.

586 *Proof.* Proceed by induction on k . For the base case $k = 0$, $s_{\max}(0) = 1$ because $x_1 = \lambda$ is the
 587 root, and all other nodes in the sequence are strictly deeper than the root. Because \mathcal{H} shatters T_d and
 588 $|\mathcal{H}| = 2^{d+1}$, \mathcal{H} is perfectly balanced on the root λ , meaning that for all $v \in \{0, 1\}$, if we denote

$$\mathcal{H}_{(v)} = \{h \in \mathcal{H} : h(\lambda) = v\},$$

589 then

$$|\mathcal{H}_{(v)}| = \frac{1}{2} \cdot 2^{d+1} = 2^d. \quad (5)$$

590 So $r_1 = 1/2$, and therefore the adversary forces a mistake on index $s_0 = s_{\max}(0) = 1$, satisfying
 591 Item 1. Furthermore, Eq. (5) implies,

$$|\mathcal{H}_{s_0}| = |\mathcal{H}_1| = |\mathcal{H}_{(y_1)}| = 2^d = 2^{d-k\sqrt{d}},$$

592 satisfying Item 2.

593 For the induction step, assume for some $k \in \{0, 1, 2, \dots, \lfloor \sqrt{d} \rfloor - 1\}$ that $|S| \geq k + 1$ and Items 1
 594 and 2 hold for k . We show that $|S| \geq k + 2$ and Items 1 and 2 hold for $k + 1$ as well. To
 595 establish $|S| \geq k + 1$ and Item 1 for $k + 1$, it suffices to show that there exists an index t such that
 596 $s_k < t \leq s_{\max}(k + 1)$ and $r_t \in [\varepsilon, 1 - \varepsilon]$. Assume for contradiction that no such index exists. Then,

$$\begin{aligned} |\mathcal{H}_{s_{\max}(k+1)}| &\geq |\mathcal{H}_{s_k}| \cdot \prod_{t=s_k+1}^{s_{\max}(k+1)} \max\{r_t, 1 - r_t\} && (r_t \notin [\varepsilon, 1 - \varepsilon] \text{ implies } y_t = y_{\text{maj}}, \\ &&& \text{so larger subclass survives}) \\ &\geq 2^{d-k\sqrt{d}} \cdot \prod_{t=s_k+1}^{s_{\max}(k+1)} \max\{r_t, 1 - r_t\} && (\text{Induction hypothesis Item 2}) \\ &\geq 2^{d-k\sqrt{d}} \cdot (1 - \varepsilon)^{n_d} && (r_t \notin [\varepsilon, 1 - \varepsilon]; \text{Item (b) in Claim B.2}) \\ &\geq 2^{d-k\sqrt{d}-1}, \end{aligned} \quad (6)$$

597 where the last inequality holds because

$$(1 - \varepsilon)^{n_d} = \left(1 - 2^{-\sqrt{d}/2}\right)^{(d+1) \cdot 2^{\sqrt{d}/10+1}} \geq \frac{1}{2} \quad (7)$$

598 for our choice of $d \geq 800$.

599 On the other hand, by Claim B.4, there exists $t \in [s_{\max}(k + 1)]$ such that $|x_t| = (k + 1)\sqrt{d}$ and
 600 $x_t \in \text{path}(\mathcal{H}_{t-1})$. By Item (a) in Claim B.2 and the construction of \mathcal{H}_t , the class \mathcal{H}_t agrees with the
 601 root-to-node path to x_t . Namely, there exists an assignment of labels $\ell : \text{path}(x_t) \rightarrow \{0, 1\}$ such
 602 that

$$\mathcal{H}_t \subseteq \{h \in \mathcal{H} : (\forall u \in \text{path}(x_t) : h(u) = \ell(u))\}.$$

603 Seeing as $|\text{path}(x_t)| = |x_t| + 1 = (k + 1)\sqrt{d} + 1$, this implies that

$$|\mathcal{H}_{s_{\max}(k+1)}| \leq |\mathcal{H}_t| \leq 2^{d+1} 2^{-(k+1)\sqrt{d}-1} = 2^{d-(k+1)\sqrt{d}}. \quad (8)$$

604 Combining Eqs. (6) and (8) gives

$$2^{d-k\sqrt{d}-1} \leq |\mathcal{H}_{s_{\max}(k+1)}| \leq 2^{d-(k+1)\sqrt{d}}.$$

605 Namely,

$$d - k\sqrt{d} - 1 \leq d - (k + 1)\sqrt{d} \implies d \leq 1,$$

606 which is a contradiction to our choice of $d \geq 800$. This implies that $|S| \geq k + 2$ and Item 1 holds for
 607 $k + 1$, namely, the index $s_{k+1} \in S$ satisfies $s_{k+1} \leq s_{\max}(k + 1)$.

608 Item 2 follows by a calculation similar to Eq. (6). Seeing as s_{k+1} is the first index after s_k where a
 609 mistake is forced,

$$\begin{aligned} |\mathcal{H}_{s_{k+1}}| &\geq \varepsilon \cdot (1 - \varepsilon)^{s_{k+1}-s_k} \cdot |\mathcal{H}_{s_k}| \\ &\geq \varepsilon \cdot (1 - \varepsilon)^{n_d} \cdot |\mathcal{H}_{s_k}| \end{aligned} \quad (\text{Item (b) in Claim B.2})$$

$$\begin{aligned}
&\geq \varepsilon \cdot (1 - \varepsilon)^{n_d} \cdot 2^{d-k\sqrt{d}} && \text{(Induction hypothesis Item 2)} \\
&\geq 2^{-\sqrt{d}/2} \cdot 2^{d-k\sqrt{d}-1} && \text{(Eq. (7) and choice of } \varepsilon) \\
&\geq 2^{d-(k+1)\sqrt{d}},
\end{aligned}$$

as desired. \square

Finally, we complete the proof of the lower bound.

Proof of Theorem B.1. Fix $d_0 = 800$ and assume $d \geq d_0$. Fix $k \leq d$. Seeing as $\text{LD}(\mathcal{H}) = d$, \mathcal{H} shatters the tree T_k . By replacing \mathcal{H} with a suitable subset of \mathcal{H} as necessary, assume without loss of generality that $\mathcal{H} \subseteq \{0, 1\}^{T_k}$ and $|\mathcal{H}| = 2^{k+1}$.

Consider an execution of $\text{TRANSDUCTIVEADVERSARY}(\mathcal{H})$ as in Algorithm 1. By Item (b) in Claim B.2, Algorithm 1 constructs a sequence $x = (x_1, x_2, \dots, x_n)$ of length $n < n_k$ for

$$n_k = (d + 1) \cdot 2^{\sqrt{d}/10+1} = O(d \cdot 2^{\sqrt{d}}).$$

By Claim B.5, for every learning rule, the adversary forces at least $\lfloor \sqrt{k} \rfloor + 1 \geq \sqrt{k}$ mistakes on the sequence x , establishing Eq. (3) in Theorem B.1 and the “moreover” sentence that follows it.

Fix a length $n \in \mathbb{N}$. Let k be the largest integer such that $2^{\lceil \sqrt{k} \rceil} \leq n$ and $k \leq d$. By Eq. (3), there exists some sequence on which the adversary can force every learning rule to make at least \sqrt{k} mistakes. By Theorem C.2, this implies that there exists a sequence of length $2^{\lceil \sqrt{k} \rceil} \leq n$ on which the adversary can force every learning rule to make at least $\sqrt{k} = \min \{ \sqrt{d}, \lfloor \log(n) \rfloor \}$ mistakes. Namely,

$$M_{\text{tr}}(\mathcal{H}, n) \geq \min \{ \sqrt{d}, \lfloor \log(n) \rfloor \},$$

as in Eq. (4). \square

C Sequence Length

In this section, we show that if there exists a sequence on which the adversary can force M mistakes, then a sequence of length $2^M - 1$ is sufficient, and this upper bound is tight for some classes.²¹

Definition C.1 (Minimal sequence). *Let \mathcal{X} be a set, let $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$ be a class, and let $M \in \mathbb{N}$.*

The minimal sequence length for forcing M mistakes for the class \mathcal{H} , denoted $\text{MinLen}(\mathcal{H}, M)$ is

$$\text{MinLen}(\mathcal{H}, M) = \inf \{ n \in \mathbb{N} : (\exists x \in \mathcal{X}^n : M_{\text{tr}}(\mathcal{H}, x) \geq M) \}.$$

In words, $\text{MinLen}(\mathcal{H}, M)$ is the smallest integer n for which there exists a sequence of length n on which the adversary can force at least n mistakes; if no such sequence exists, then $\text{MinLen}(\mathcal{H}, M) = \infty$.

Theorem C.2 (Minimal sequence bound). *Let \mathcal{X} be a set, and fix $M \in \mathbb{N}$. Then for any class $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$, if $\text{MinLen}(\mathcal{H}, M) < \infty$ then*

$$\text{MinLen}(\mathcal{H}, M) \leq 2^M - 1.$$

Furthermore, there exists a class $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$ for which $\text{MinLen}(\mathcal{H}, M) = 2^M - 1$.

Theorem C.2 is a corollary of the tree rank characterization of M_{tr} from Ben-David et al. (1997). For completeness, we present a direct proof of Theorem C.2 that does not directly invoke that characterization. Roughly, given an adversary A_0 that forces every learner to make at least M mistakes on a (possibly long) sequence x , we apply two modifications to obtain new adversaries

$$A_0 \rightsquigarrow A_1 \rightsquigarrow A_2.$$

A_1 forces M mistakes and has a specific structure that we call ‘rigidity’, but it still uses the same (possibly long) sequence x . Capitalizing on the rigid structure, A_2 selects a subsequence of x of length at most $2^M - 1$, and forces M mistakes on that subsequence.

²¹Of course, there also exist classes for which a shorter sequence is sufficient. For instance, if the class shatters (in the VC sense) a subset of the domain of cardinality M , then a sequence of length M suffices.

643 C.1 Rigid Adversary

644 **Definition C.3** (Rigid adversary). *Let $n \in \mathbb{N}$, let \mathcal{X} be a set, and let*

$$A : \left(\bigcup_{k=0}^{n-1} \{0, 1\}^{2k} \right) \times \{0, 1\} \rightarrow \{0, 1\}$$

645 *be an adversary strategy for some fixed sequence $x \in \mathcal{X}^n$. We say that A is rigid if there exists a*
 646 *function*

$$f : \bigcup_{k=0}^{n-1} \{0, 1\}^k \rightarrow \{0, 1, \star\}$$

647 *such that for all $k \in \{0, 1, \dots, n-1\}$ and all $y, \hat{y} \in \{0, 1\}^k$,*

$$A(\hat{y}_1, y_1, \dots, \hat{y}_k, y_k, \hat{y}_{k+1}) = \begin{cases} f(y_1, \dots, y_k) & f(y_1, \dots, y_k) \in \{0, 1\} \\ 1 - \hat{y}_{k+1} & f(y_1, \dots, y_k) = \star \end{cases}.$$

648 Note that if an adversary is rigid, then the function f that witnesses this is uniquely determined.

649 **Claim C.4** (Rigid adversary exists). *Let $n, M \in \mathbb{N}$, let \mathcal{X} be a set, let $x \in \mathcal{X}^n$, and let $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$*
 650 *be a class. Let A be an adversary strategy that forces every learner to make at least M mistakes on x .*
 651 *Then there exists an adversary strategy A^* such that:*

- 652 1. *A^* forces every learner to make at least M mistakes on x and A^* is rigid.*
- 653 2. *Let f be the function that witnesses the rigidity of A^* . Then for every $y \in \{0, 1\}^n$, the*
 654 *sequence*

$$f(y_{\leq 0}), f(y_{\leq 1}), f(y_{\leq 2}), \dots, f(y),$$

655 *has at least M members equal to \star .*

656 *Proof of Claim C.4.* For Item 1, consider the adversary strategy A^* that simulates an execution of A ,
 657 as in Algorithm 3. In broad strokes, A^* functions as a middle-man between the learner and A . As
 658 the learner makes a sequence of predictions $\tilde{y} \in \{0, 1\}^n$, the adversary A^* generates a sequence of
 659 (possibly different) predictions $\hat{y} \in \{0, 1\}^n$, and sends those to the adversary A . Adversary A sees
 660 only the predictions \hat{y} , and assigns labels $y \in \{0, 1\}^n$, which are relayed back to the learner by A^*
 661 with no modifications.

662 First, observe that A^* satisfies the realizability requirement. Indeed, A^* simulates an execution of A
 663 such that the sequence of labels y_1, \dots, y_n sent by A^* to the learner is exactly the sequence of labels
 664 selected by A . Seeing as A is realizable, every sequence of labels selected by A is realizable, and
 665 therefore every sequence of labels selected by A^* must be realizable as well.

666 Second, observe that A^* forces every learner to make at least M mistakes. To see this, notice that in
 667 Algorithm 3,

$$\sum_{t \in [n]} \mathbb{1}(\tilde{y}_t \neq y_t) \geq M. \tag{9}$$

668 Indeed, A forces every learner to make at least M mistakes, and in particular this applies to a learner
 669 that makes predictions \tilde{y} as in the simulation. Furthermore, observe that A^* only alters the predictions
 670 it receives from the learner in cases when it selects a label that is accepted by A , namely,

$$\forall t \in [n] : \tilde{y}_t \neq \hat{y}_t \implies \tilde{y}_t = y_t. \tag{10}$$

671 Therefore, if $E = \{t \in [n] : \tilde{y}_t = \hat{y}_t\}$, then

$$\begin{aligned} \sum_{t \in [n]} \mathbb{1}(\tilde{y}_t \neq y_t) &= \sum_{t \in E} \mathbb{1}(\tilde{y}_t \neq y_t) + \sum_{t \in [n] \setminus E} \mathbb{1}(\tilde{y}_t \neq y_t) \\ &= \sum_{t \in E} \mathbb{1}(\tilde{y}_t \neq y_t) + 0 && \text{(By Eq. (10))} \\ &= \sum_{t \in E} \mathbb{1}(\hat{y}_t \neq y_t) && \text{(Definition of } E\text{)} \end{aligned}$$

Assumptions:

- $n \in \mathbb{N}$, \mathcal{X} is a set, $x \in \mathcal{X}^n$ is a fixed sequence of instances.
- $A : \left(\bigcup_{k=0}^{n-1} \{0, 1\}^{2k} \right) \times \{0, 1\} \rightarrow \{0, 1\}$ is an adversary labeling strategy for x .

RIGIDADVERSARY:

```

send  $x_1, \dots, x_n$  to the learner
for  $t = 1, 2, \dots, n$ :
    receive prediction  $\hat{y}_t$  from learner
    if  $A(\tilde{y}_1, y_1, \dots, \tilde{y}_{t-1}, y_{t-1}, 0) = 0$ :
         $\tilde{y}_t \leftarrow 0$ 
    else if  $A(\tilde{y}_1, y_1, \dots, \tilde{y}_{t-1}, y_{t-1}, 1) = 1$ :
         $\tilde{y}_t \leftarrow 1$ 
    else:
         $\tilde{y}_t \leftarrow \hat{y}_t$ 
    send prediction  $\tilde{y}_t$  to  $A$ 
    receive label  $y_t$  from  $A$ 
    send label  $y_t$  to learner

```

Algorithm 3: Construction of a rigid adversary, by simulating a given adversary A .

$$\leq \sum_{t \in [n]} \mathbb{1}(\hat{y}_t \neq y_t). \quad (11)$$

672 Combining Eqs. (9) and (11) implies that A forces at least M mistakes.

673 Third, we show that A^* is rigid. We claim that there exists a function $g : \{0, 1\}^{\leq n-1} \rightarrow \{0, 1\}^{\leq n-1}$
 674 such that for every $t \in \{0, 1, 2, \dots, n-1\}$,

$$(\tilde{y}_1, \dots, \tilde{y}_t) = g(y_1, \dots, y_t).$$

675 Proceed by induction on t . For the base case $t = 0$ there is nothing to prove. For the induction step,
 676 we assume the claim holds for some $t = k < n - 1$, and show that it holds for $t = k + 1$. From
 677 Algorithm 3, \tilde{y}_{k+1} satisfies

$$\tilde{y}_{k+1} = \begin{cases} 0 & A(\tilde{y}_1, y_1, \dots, \tilde{y}_k, y_k, 0) = 0 \\ 1 & A(\tilde{y}_1, y_1, \dots, \tilde{y}_k, y_k, 0) = A(\tilde{y}_1, y_1, \dots, \tilde{y}_k, y_k, 1) = 1 \\ 1 - y_{k+1} & \text{otherwise} \end{cases} \quad (12)$$

678 The first two cases in Eq. (12) are immediate from Algorithm 3, and the remaining case occurs when
 679 A forces a mistake at time $k + 1$, namely, when A selects $y_{k+1} = 1 - \tilde{y}_{k+1}$. Thus, \tilde{y}_{k+1} is a function
 680 of $y_{\leq k+1}$ and $\tilde{y}_{\leq k}$. By the induction hypothesis, $\tilde{y}_{\leq k} = g(y_{\leq k})$, so \tilde{y}_{k+1} is simply a function of
 681 $y_{\leq k+1}$. This establishes the existence of the desired function g .

682 Hence, A^* is rigid, as witnessed by the function

$$f(y_1, \dots, y_k) = \begin{cases} 0 & A(\tilde{y}_1, y_1, \dots, \tilde{y}_k, y_k, 0) = 0 \\ 1 & A(\tilde{y}_1, y_1, \dots, \tilde{y}_k, y_k, 0) = A(\tilde{y}_1, y_1, \dots, \tilde{y}_k, y_k, 1) = 1 \\ \star & \text{otherwise} \end{cases},$$

683 where f is a well-defined function because $\tilde{y}_{\leq k} = g(y_{\leq k})$.

684 We have seen that A^* is a valid (realizable) adversary that forces every learner to make at least M
 685 mistakes, and it is rigid. This concludes the proof of Item 1.

686 Finally, For Item 2, note that $\tilde{y}_t \neq y_t$ only if A forces a mistake at time t in the sense that A selects
 687 $y_t = 1 - b$ for any prediction $b \in \{0, 1\}$ provided at time t . If A forces a mistake at time t , then A^*

688 forces a mistake at time t as well. Therefore, if $\tilde{y}_t \neq y_t$, then $f(y_{<t}) = \star$, namely, \tilde{y}_t makes mistakes
 689 only when the value of f is \star . By Eq. (9), \tilde{y}_t makes at least M mistakes throughout the game, so
 690 there must be at least M rounds where f outputs \star , as desired. \square

691 C.2 Essential Indices

692 **Definition C.5.** Let $n, M \in \mathbb{N}$, let \mathcal{X} be a set, let $x \in \mathcal{X}^n$, and let $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$ be a class. Let A
 693 be a rigid adversary strategy witnessed by function f . We say that an index $t \in [n]$ is essential for
 694 A for forcing M mistakes on x if there exists a sequence $y \in \{0, 1\}^{t-1}$ such that $f(y) = \star$ and the
 695 sequence

$$f(y_{\leq 0}), f(y_{\leq 1}), f(y_{\leq 2}), \dots, f(y_{\leq t-1})$$

696 contains at most $M - 1$ members equal to \star .

697 **Claim C.6.** Let $n, M \in \mathbb{N}$, let \mathcal{X} be a set, let $x \in \mathcal{X}^n$, and let $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$ be a class. Let A be
 698 a rigid adversary strategy. Then $[n]$ contains at most $2^M - 1$ indices that are essential for A for
 699 forcing M mistakes on x .

700 *Proof.* For each essential index $t \in [n]$, there exists a label sequence $y \in \{0, 1\}^{t-1}$ that witnesses
 701 that t is essential, as in Definition C.5. Each label sequence y is a witness for at most one index (the
 702 index $|y| + 1$), so it suffices to show that the set $Y \subseteq \{0, 1\}^{\leq n-1}$ of all witness label sequences is of
 703 cardinality at most $2^M - 1$.

704 Think of Y as a collection of nodes in the binary tree T_{n-1} (Definition A.4). By Definition C.5, if
 705 $y \in Y$, then the collection of all ancestors of y in Y has cardinality

$$|\{y_{\leq i} : i \in \{0, 1, 2, \dots, |y| - 1\}\} \cap Y| \leq M - 1.$$

706 Namely, Y is a subtree of depth at most $d = M - 1$ in the binary tree T_{n-1} .²² Hence, the number of
 707 nodes in Y is at most

$$2^{d+1} - 1 = 2^M - 1,$$

708 as desired. \square

²²The depth of a subtree is s if the longest root-to-node path contains $s + 1$ nodes from the subtree.

Assumptions:

- $n, M \in \mathbb{N}$, \mathcal{X} is a set, $x \in \mathcal{X}^n$ is a fixed sequence of instances.
- $A_1 : \left(\bigcup_{k=0}^{n-1} \{0, 1\}^{2^k} \right) \times \{0, 1\} \rightarrow \{0, 1\}$ is a rigid adversary labeling strategy for x that forces every learner to make at least M mistakes on the sequence x , and satisfies Items 1 and 2 in Claim C.4.
- $I = \{i_1, i_2, \dots, i_k\} \subseteq [n]$ is the set of indices that are essential for A for forcing M mistakes on x , and $i_1 \leq i_2 \leq \dots \leq i_k$. By Claim C.6, $k \leq 2^M - 1$.

MINIMALADVERSARY:

```

send  $x_{i_1}, x_{i_2}, \dots, x_{i_k}$  to the learner
for  $t = 1, 2, \dots, n$ :
  if  $t \in I$ :
    receive prediction  $\hat{y}_t$  from learner
    send prediction  $\hat{y}_t$  to  $A_1$ 
    receive label  $y_t$  from  $A_1$ 
    send label  $y_t$  to learner
  else:
    send prediction  $\hat{y}_t = 0$  to  $A_1$ 
    receive label  $y_t$  from  $A_1$ 

```

Algorithm 4: Construction of an adversary that forces M mistakes using a sequence x of length at most $2^M - 1$. In the proof of Theorem C.2, this adversary is A_2 . Internally, it simulates a rigid adversary A_1 .

710 *Proof of Theorem C.2.* If $\text{MinLen}(\mathcal{H}, M) < \infty$, then there exist a sequence $x \in \mathcal{X}^n$, and an adver-
 711 sary A_0 that forces every learner to make at least M mistakes on x . By Claim C.4, there exists a rigid
 712 adversary A_1 that causes every learner to make at least M mistakes on x ,²³ and also satisfies Item 2
 713 in Claim C.4. Let f be the function that witnesses the rigidity of A_1 . By Claim C.6, the set $I \subseteq [n]$
 714 of indices that are essential for A_1 for forcing M mistakes on x has cardinality $k = |I| \leq 2^M - 1$.

715 Algorithm 4 defines a new adversary, A_2 , which forces every learner to make at least M mistakes on
 716 a sequence of length k . A_2 is realizable, because A_1 is realizable.²⁴

717 To see that adversary A_2 forces every learner to make at least M mistakes, let y_1, \dots, y_n be the
 718 sequence of labels assigned by A_2 . Seeing as A_2 assigns the same labels as A_1 , and A_1 satisfies
 719 Item 2 in Claim C.4, it follows that there are at least M indices $j \in [n]$ such that $f(y_{\leq j-1}) = \star$. Fix
 720 $J \subseteq [n]$ to be the first M such indices. Then $J \subseteq I$, namely, all the indices in J are essential for A_1
 721 for forcing M mistakes on x (Definition C.5).

722 Therefore, for each $j \in J$, A_2 includes the instance x_j in the sequence of length k sent to the learner.
 723 Then, in round j of the n rounds simulated by A_2 :

- 724 • The learner makes a prediction $\hat{y}_j \in \{0, 1\}$ corresponding to instance x_j .

²³This is Item 1 in Claim C.4.

²⁴The argument for realizability is the same as in the proof of Claim C.4.

725 • Adversary A_2 sends prediction \hat{y}_j to adversary A_1 . Because $f(y_{\leq j-1}) = \star$, adversary A_1
 726 assigns the label $y_j = 1 - \hat{y}_j$. Adversary A_2 then sends that label y_j to the learner. So the
 727 learner makes a mistake on x_j .

728 Hence, the learner makes at least $|J| = M$ mistakes, as desired. \square

729 D Upper Bound

730 D.1 Statement

731 The following result states that the lower bound of Theorem B.1 is tight for some classes.

732 **Theorem D.1** (Upper bound, and separation between standard and transductive online learning).
 733 *For every integer $d \geq 23$, there exists a hypothesis class $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$ with a domain \mathcal{X} of size
 734 $|\mathcal{X}| = 2^d - 1$ such that $\text{LD}(\mathcal{H}) = d$ and the following two conditions hold for all $n \in \mathbb{N}$:*

- 735 1. $M_{\text{tr}}(\mathcal{H}, n) \leq 48 \cdot \sqrt{d}$.
- 736 2. $M_{\text{std}}(\mathcal{H}, n) = \min \{n, d\}$.

737 D.2 Hypothesis Class

738 In this section we construct the hypothesis class for Theorem D.1.

739 **Lemma D.2.** *Let $d \in \mathbb{N}$, $d \geq 22$. Let T_d be a perfect binary tree of depth d , as in Definition A.4.
 740 Then there exists a collection of functions $\mathcal{H} \subseteq \{0, 1\}^{T_d}$ such that $\text{LD}(\mathcal{H}) = d + 1$ and the following
 741 two conditions hold for all $H \subseteq \mathcal{H}$ and all $X \subseteq T_d$:*

- 742 1. *If $\forall h \in H \forall x \in X : x \notin \text{path}(h) \wedge h(x) = 0$, then $\min \{|H|, |X|\} < 2^{2\sqrt{d}}$.*
- 743 2. *If $\forall h \in H \forall x \in X : x \notin \text{path}(h) \wedge h(x) = 1$, then $|H| < 2^{2\sqrt{d}}$ or $|X| < 3\sqrt{d}$.*

744 The proof employs the probabilistic method, showing that a hypothesis class sampled randomly from
 745 a suitable distribution has the desired properties with very high probability.

746 *Proof.* Let \mathcal{P} be a probability distribution over hypothesis classes. Formally, $\mathcal{P} \in$
 747 $\Delta\left(\left(\{0, 1\}^{T_d}\right)^{2^{d+1}}\right)$ is a distribution over vectors of hypotheses. Each vector $\mathcal{H} \in \text{supp}(\mathcal{P})$ consists
 748 of 2^{d+1} hypotheses,

$$\mathcal{H} = (h_b)_{b \in \{0, 1\}^{d+1}},$$

749 where for each $b \in \{0, 1\}^{d+1}$, hypothesis h_b is a function $h_b : T_d \rightarrow \{0, 1\}$ sampled independently
 750 as following:

- 751 • For each $i \in [d] \cup \{0\}$: $h_b(b_{\leq i}) = b_{i+1}$. (In particular, with probability 1, $\text{path}(h_b) =$
 752 $(b_{\leq 0}, b_{\leq 1}, \dots, b_{\leq d})$, each entry in the vector \mathcal{H} is unique, and \mathcal{H} shatters T_d .)
- 753 • For each $x \in T_d \setminus \text{path}(h_b)$, the bit $h_b(x) \in \{0, 1\}$ is sampled $\text{Ber}(2^{-\sqrt{d}})$ independently of
 754 all other bits in \mathcal{H} , i.e., $\mathbb{P}[h_b(x) = 1] = \mathbb{P}[h_b(x) = 1 \mid \{h_{b'}\}_{b' \neq b}, \{h_b(x')\}_{x' \neq x}] = 2^{-\sqrt{d}}$.

755 Fix $B \subseteq \{0, 1\}^{d+1}$ and $X \subseteq T_d$, and let $E(B, X, y)$ denote the event

$$\{\forall b \in B \forall x \in X : x \notin \text{path}(h_b) \wedge h_b(x) = y\}. \quad (13)$$

756 Seeing as each off-path label $h_b(x) \in \{0, 1\}$ is sampled independently,

$$\begin{aligned} \mathbb{P}_{\mathcal{H} \sim \mathcal{P}}[E(B, X, 0)] &= \prod_{(b, x) \in B \times X} \mathbb{P}_{\mathcal{H} \sim \mathcal{P}}[x \notin \text{path}(h_b) \wedge h_b(x) = 0] \\ &\leq (1 - 2^{-\sqrt{d}})^{|B \times X|}. \end{aligned} \quad (14)$$

757 Hence,

$$\begin{aligned}
& \mathbb{P}_{\mathcal{H} \sim \mathcal{P}} \left[\exists B \subseteq \{0, 1\}^{d+1} \exists X \subseteq T_d : E(B, X, 0) \wedge \min \{|B|, |X|\} \geq 2^{2\sqrt{d}} \right] \\
& \leq \binom{|\{0, 1\}^{d+1}|}{\lceil 2^{2\sqrt{d}} \rceil} \binom{|T_d|}{\lceil 2^{2\sqrt{d}} \rceil} (1 - 2^{-\sqrt{d}})^{2^{4\sqrt{d}}} && \text{(union bound, Eq. (14))} \\
& < 2 \binom{2^{d+1}}{2^{2\sqrt{d}} + 1} (1 - 2^{-\sqrt{d}})^{2^{4\sqrt{d}}} \\
& < 2 \cdot 2^{(d+1) \cdot (2^{2\sqrt{d}} + 1)} \cdot e^{-2^{-\sqrt{d}} \cdot 2^{4\sqrt{d}}} && \left(\binom{n}{k} < n^k \text{ for } k \geq e; 1 + x \leq e^x \text{ for } x \in \mathbb{R} \right) \\
& < 2^{(d+2)2^{2\sqrt{d}}} \cdot 2^{-2^{-\sqrt{d}} \cdot 2^{4\sqrt{d}}} && (2^{2\sqrt{d}} \geq d + 2 \text{ for } d \geq 1) \\
& = 2^{2^{2\sqrt{d}} \cdot (d+2-2^{2\sqrt{d}})} \\
& < 2^{-2^{2\sqrt{d}}} && (d + 2 - 2^{2\sqrt{d}} < -1 \text{ for } d \geq 22) \\
& && (15)
\end{aligned}$$

758 Similarly,

$$\mathbb{P}_{\mathcal{H} \sim \mathcal{P}} [\forall b \in B \forall x \in X : x \notin \text{path}(h_b) \wedge h_b(x) = 1] \leq 2^{-\sqrt{d} \cdot |B \times X|}, \quad (16)$$

759 so

$$\begin{aligned}
& \mathbb{P}_{\mathcal{H} \sim \mathcal{P}} \left[\exists B \subseteq \{0, 1\}^{d+1} \exists X \subseteq T_d : E(B, X, 0) \wedge |H| \geq 2^{2\sqrt{d}} \wedge |X| \geq 3\sqrt{d} \right] \\
& \leq \binom{|\{0, 1\}^{d+1}|}{\lceil 2^{2\sqrt{d}} \rceil} \binom{|T_d|}{\lceil 3\sqrt{d} \rceil} \cdot 2^{-\sqrt{d} \cdot 2^{2\sqrt{d}} \cdot 3\sqrt{d}} && \text{(union bound, Eq. (16))} \\
& < 2 \binom{2^{d+1}}{2^{2\sqrt{d}} + 1} \cdot 2^{-3d \cdot 2^{2\sqrt{d}}} \\
& < 2^{2d \cdot 2^{2\sqrt{d}}} \cdot 2^{-3d \cdot 2^{2\sqrt{d}}} && \text{(for } d \geq 2) \\
& < 2^{-d2^{2\sqrt{d}}}. && (17)
\end{aligned}$$

760 Applying a union bound to Eqs. (15) and (17) gives

$$\mathbb{P}_{\mathcal{H} \sim \mathcal{P}} [\mathcal{H} \text{ satisfies Items 1 and 2}] \geq 1 - 2^{-2^{2\sqrt{d}}} - 2^{-d2^{2\sqrt{d}}} \geq 1 - 10^{-100}.$$

761 In particular, there exists a collection \mathcal{H} that satisfies Items 1 and 2. Furthermore, this collection
762 has $\text{LD}(\mathcal{H}) = d + 1$ (namely, $\text{LD}(\mathcal{H}) \geq d + 1$ because it shatters T_d ; and $\text{LD}(\mathcal{H}) \leq d + 1$ because
763 $|\mathcal{H}| = 2^{d+1}$). \square

764 D.3 Algorithm

765 In this section we describe Algorithms 5, 6a and 6b, which together constitute the learning algorithm
766 that achieves the $O(\sqrt{d})$ mistake upper bound in the transductive setting, as in Theorem D.1. See
767 Section 2.3 for a general overview of these algorithms.

768 D.3.1 How Experts Work

769 We start with some preliminary remarks about experts in Algorithms 5, 6a and 6b.

770 **Experts.** A tuple $e = (S, u, H)$ defines an expert that can make predictions using the procedure
771 `EXPERT.PREDICT(e, \cdot)`. The tuple e reflects two kinds of information:

- 772 1. *Knowledge.* Information that the expert *knows* with certainty. Specifically, this reflects the
773 labels y_1, y_2, \dots sent by the adversary so far. All experts see the labels sent by the adversary,
774 so this knowledge is the same for all experts.

2. *Assumptions.* At certain times, experts make *assumptions* about things that are not known for certain. Specifically, experts assume that certain nodes x are on-path ($x \in \text{path}(h)$) or off-path ($x \notin \text{path}(h)$) with respect to the correct labeling function $h : T_d \rightarrow \{0, 1\}$. Assumptions are simply guesses that may be wrong, and therefore when an expert needs to make such an assumption, it splits into two experts (as described below), with one expert assuming $x \in \text{path}(h)$, and the other expert assuming $x \notin \text{path}(h)$. This ensures that there always exists an expert for which all assumptions are correct.

In greater detail, the contents of the state tuple $e = (S, u, H)$ represents the knowledge and assumptions of the expert as follows:

- $u \in T_d$ – This single node encodes everything the expert knows and assumes about which of the nodes labeled so far are on-path. Observe that if $v_1, v_2, \dots, v_k \in T_d$ are nodes that are assumed to be on-path (and all these assumptions are consistent), then these k assumptions can be represented succinctly by assigning $u = v_{i^*}$ where v_{i^*} is the deepest node among v_1, v_2, \dots, v_k . Therefore, u simply holds the deepest node in the tree that is known or assumed to be on-path. At the start of the algorithm, this value is initialized to be $u = \lambda$, because the root is known to be on-path regardless of the target function.
- $S \subseteq T_d$ – the ‘danger zone’, as described in Section 2.3.4. This is a collection that contains all nodes in the prefix $x_{\leq t_{\max}} = (x_1, x_2, \dots, x_{t_{\max}})$ of the sequence to be classified that have not been labeled yet and *might* be on-path for the true labeling function h given what the expert knows and assumes so far. However, S is not required to contain ancestors of nodes that are assumed to be on-path. Initially, S equals the prefix $x_{\leq t_{\max}}$. As information accumulates, nodes that cannot be on-path are removed from S . For instance, if $x_i \in T_d$ is assigned label $y_i \in \{0, 1\}$ by the adversary, then any $(1 - y_i)$ -descendant of x_i (including x_i itself) may safely be removed from S .
- $H \subseteq \{0, 1\}^{T_d}$ – the version space of the experts, i.e., the collection of all functions that could be the correct labeling function given everything that the expert knows and assumes. Initially, H contains all functions in \mathcal{H} . As information accumulates, some functions are ruled out. Specifically, a function h can be removed from H for two reasons: (i) the adversary assigns a label $y \neq h(x)$ to some node $x \in T_d$; (ii) the expert makes an assumption that some $x \in T_d$ is on-path for the correct labeling function but $x \notin \text{path}(h)$, or vice versa, the expert assumes that x is off-path for the correct labeling function but $x \in \text{path}(h)$.

Updates and splits. An expert can be modified using the procedure $\text{EXPERT.UPDATE}(e, \cdot, \cdot)$. This procedure either returns a single modified tuple (S, u, H) (in the first two return statements in the procedure), in which case we think of the expert as being *updated*; or alternatively, the procedure returns two tuples $e_{\in} = (S_{\in}, u_{\in}, H_{\in})$ and $e_{\notin} = (S_{\notin}, u_{\notin}, H_{\notin})$ (in the third return statement), in which case we think of the expert as being *split* into two experts. e_{\in} corresponds to adding an assumption that the most recently presented node x_t is on-path for the correct labeling function, and e_{\notin} corresponds to adding the opposite assumption.

Ancestry. At the end of each iteration of the outer ‘for’ loop in Algorithm 5, for each expert $e \in E_{t+1}$ there exists a unique *ancestry* sequence $\text{ancestry}(e) = (e_1, e_2, \dots, e_{t+1})$ such that $e_1 = (\{x_1, \dots, x_{t_{\max}}\}, \lambda, \mathcal{H})$ is the initial single expert that was created before the start of the outer ‘for’ loop, $e_{t+1} = e$ is the latest version of the expert, and for each $i \in [t]$, the expert e_{i+1} is either equal to e_i , or it was created by executing $\text{EXPERT.UPDATE}(e_i, \cdot, \cdot)$.²⁵

²⁵Note that in this paper, we use genealogical metaphors in two distinct contexts that should not be confused. First, as is customary, we use “child”, “parent”, “ancestor” and “descendant” to describe relations between nodes in the binary tree T_d , which constitutes the domain of our hypothesis class. Separately from that, we use “ancestor” and “descendant” to describe relations between experts.

This overlap in terminology can partially be excused by the fact that the history of experts also forms a binary tree. Indeed, initially there is a single expert (the root of the tree), and experts can split into two, corresponding to a node having two children as in a binary tree. Seeing as experts cannot merge, the expert history corresponds precisely to a binary tree. (However, the domain T_d is a *perfect* binary tree, whereas the binary tree corresponding to expert genealogy need not be balanced).

To reduce confusion, we use $\text{path}(\cdot)$ only for nodes in T_d , and $\text{ancestry}(\cdot)$ only for experts, even though these operators are mathematically equivalent (however, $\text{path}(\cdot)$ is defined not only for nodes in T_d but also for functions $T_d \rightarrow \{0, 1\}$).

Assumptions:

- $d, n \in \mathbb{N}$, λ is the empty string.
- $\mathcal{H} \subseteq \{0, 1\}^{T_d}$ is the class that exists by Lemma D.2.
- $x_1, x_2, \dots, x_n \in T_d$ are points to be classified.

TRANSDUCTIVELEARNER($\mathcal{H}, d, (x_1, x_2, \dots, x_n)$):

$t \leftarrow 0, t_{\max} \leftarrow 2^{4\sqrt{d}}$
 $e \leftarrow (\{x_1, \dots, x_{t_{\max}}\}, \lambda, \mathcal{H})$ \triangleright The initial expert. An expert is defined by a 3-tuple.
 $w(e) \leftarrow 1$ \triangleright Assign the initial expert a weight of 1.
 $E_1 \leftarrow \{e\}$ $\triangleright E_t$ is the set of experts used for predicting \hat{y}_t .
 $E_2, \dots, E_n, E_{n+1} \leftarrow \emptyset$

for $t \leftarrow 1, 2, \dots, n$:

$\hat{y}_t \leftarrow \mathbb{1} \left(\sum_{e \in E_t} w(e) \cdot \text{EXPERT.PREDICT}(e, x_t) \geq \frac{1}{2} \right)$ \triangleright

A weighted majority, using
Algorithm 6a.

send prediction \hat{y}_t to adversary

receive correct label $y_t \in \{0, 1\}$ from adversary

for $e \in E_t$: \triangleright Update the experts.

if $\text{EXPERT.PREDICT}(e, x_t) = y_t$:

$E_{t+1} \leftarrow E_{t+1} \cup \{e\}$ \triangleright If expert e made a correct prediction,
 keep it and its weight unchanged.

else:

$U \leftarrow \text{EXPERT.UPDATE}(e, x_t, y_t)$ \triangleright If e made a mistake, update e using
 Algorithm 6b. This might cause e to
 be split into two experts.

for $e' \in U$:

$E_{t+1} \leftarrow E_{t+1} \cup \{e'\}$ \triangleright Add updated expert(s) to E_{t+1} .
 $w(e') \leftarrow w(e)/(2 \cdot |U|)$ \triangleright When e makes a mistake, its weight
 is decreased by a factor of 2 and then
 split equally between its descendants.

Algorithm 5: A transductive online learning algorithm that makes at most $O(\sqrt{d})$ mistakes. It is a variant of the multiplicative weights algorithm that employs splitting experts. Namely, we start with a single expert, and when an expert makes a mistake it may split into two experts. The behavior of the experts is defined in Algorithms 6a and 6b.

Assumptions:

- $d \in \mathbb{N}, x \in T_d$.
- $e = (S, u, H)$ is a tuple that defines an expert:
 - $S \subseteq T_d$ – a collection of nodes that could be on-path for the true labeling function given what the expert knows and assumes.
 - $u \in T_d$ – the deepest node known or assumed to be on-path by the expert.
 - $H \subseteq \{0, 1\}^{T_d}$ – the collection of all functions that could be the correct labeling function given what the expert knows and assumes.

EXPERT.PREDICT(e, x):

```

( $S, u, H$ )  $\leftarrow e$                                 ▷ Unpack the state that defines the expert.

if  $|H| \leq 2^{2\sqrt{d}}$ :
    return HALVING.PREDICT( $H, x$ )                    ▷ Once  $H$  becomes small enough, simulate the
                                                    Halving algorithm (Algorithm 7). [Case I]
if  $x \preceq u$ :
    return  $b \in \{0, 1\}$  such that  $x \preceq_b u$           ▷  $u$  is assumed to be on-path. If  $u$  is a  $b$ -
                                                    descendant of  $x$ , then the correct label for
                                                     $x$  must be  $b$ . [Case II]

return  $\mathbb{1}\left(\frac{|\{x' \in S : x \preceq_1 x'\}|}{|S|} \geq \frac{1}{3}\right)$   ▷ If there is  $b \in \{0, 1\}$  such that at least a 1/3
                                                    of suspected on-path nodes are  $b$ -descendants
                                                    of  $x$ , then output  $b$ . Otherwise (when at least
                                                    2/3 of  $S$  are non-descendants of  $x$ ), output 0.
                                                    [Cases III to VI]
```

Algorithm 6a: A subroutine of Algorithm 5 that defines how an expert makes predictions.

Assumptions:

- d, x, e, S, u, H – as in Algorithm 6a.
- y – the correct label for x , as selected by the adversary.

EXPERT.UPDATE(e, x, y):

$(S, u, H) \leftarrow e$	▷ Unpack the state that defines the expert.
$H \leftarrow \text{HALVING.UPDATE}(H, x, y)$	▷ Update the version space, as in the Halving algorithm (Algorithm 7).
if $ H \leq 2^{2\sqrt{d}}$:	▷ If the version space is small, we just simulate the Halving algorithm, so the update is complete. [Case III]
return $\{(S, u, H)\}$	
for $b \in \{0, 1\}$:	
$S_b \leftarrow \{x' \in S : x \preceq_b x'\}$	▷ Set of suspected on-path nodes that are b -descendant of x .
if $ S_{(1-y)} / S \geq 1/3$:	
$S' \leftarrow S \setminus S_{(1-y)}$	▷ At least 1/3 of suspected on-path nodes were b -
return $\{(S', u, H)\}$	descendants of x , and therefore the expert predicted label $\hat{y} = b$. But the correct label was $y = 1 - b$. Remove all b -descendants of x from S . [Case IV]
else:	
$S_{\notin} \leftarrow S; u_{\notin} \leftarrow u$	▷ Split e in two. First, construct e_{\notin} to be an updated version of e after adding the assumption that $x \notin \text{path}(h)$ for the correct labeling function h .
$H_{\notin} = \{h \in H : x \notin \text{path}(h)\}$	
$e_{\notin} \leftarrow (S_{\notin}, u_{\notin}, H_{\notin})$	
$S_{\in} \leftarrow S_0 \cup S_1$	▷ Next, construct e_{\in} to be an updated version of e adding the assumption $x \in \text{path}(h)$. S_{\in} contains only nodes that are descendants of x .
$u_{\in} \leftarrow u$	
if $u_{\in} \preceq x$:	▷ u_{\in} represents updating the prior assumption that u is on path by adding that x is also on path.
$u_{\in} \leftarrow x$	
$H_{\in} = \{h \in H : x \in \text{path}(h)\}$	▷ H_{\in} is obtained by updating the version space to include only function where x is on path.
$e_{\in} \leftarrow (S_{\in}, u_{\in}, H_{\in})$	
return $\{e_{\notin}, e_{\in}\}$	▷ [Cases V and VI]

Algorithm 6b: A subroutine of Algorithm 5 that defines how an expert is updated (and possibly split into two) when it makes a mistake.

Assumptions:

- \mathcal{X} a set, $k \in \mathbb{N}$.
- $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$.
- $x, x_1, \dots, x_k \in \mathcal{X}, y \in \{0, 1\}$.

HALVING($\mathcal{H}, (x_1, x_2, \dots, x_k)$):

$\mathcal{H}_1 \leftarrow \mathcal{H}$

for $i \in [k]$:

$\hat{y}_i \leftarrow \text{HALVING.PREDICT}(\mathcal{H}, x_i)$

send prediction \hat{y}_i to adversary

receive correct label $y_i \in \{0, 1\}$ from adversary

$\mathcal{H}_{i+1} \leftarrow \text{HALVING.UPDATE}(\mathcal{H}_i, x_i, y_i)$

HALVING.PREDICT(\mathcal{H}, x):

return $\mathbb{1}(\sum_{h \in \mathcal{H}} h(x) \geq \frac{1}{2})$

HALVING.UPDATE(\mathcal{H}, x, y):

return $\{h \in \mathcal{H} : h(x) = y\}$

Algorithm 7: This is the well-known halving algorithm. The experts in Algorithms 6a and 6b simulate this algorithm once their version space becomes small enough.

819 D.4 Analysis

820 In this section we prove our main result, Theorem D.1.

821 D.4.1 Assumption-Consistent Expert

822 Occasionally, when an expert is updated, it makes an assumption about whether the most-recently
 823 presented node x_t is on-path or off-path with respect to the true labeling function h . In these
 824 updates, the expert is split into two: one expert assumes that $x_t \in \text{path}(h)$, and the other assumes
 825 $x_t \notin \text{path}(h)$. Clearly, by splitting into two in this manner, we preserve the invariant that the set of
 826 experts always contains a ‘vindicated’ expert e^* such that all the assumptions made by e^* are correct.
 827 This simple observation is made formal in the following definition and claim.

828 **Definition D.3** (Assumption consistency). *For an expert $e \in E_{t+1}$ with $\text{ancestry}(e) =$
 829 $(e_1, e_2, \dots, e_{t+1})$, and an index $i \in [t]$, we say that the $i \rightarrow (i+1)$ update of e was assump-
 830 tion-consistent with a function $h : T_d \rightarrow \{0, 1\}$ if one of the following conditions hold:*

- 831 • $e_{i+1} = e_i$; or
- 832 • e_{i+1} was the single expert returned by $\text{EXPERT.UPDATE}(e_i, x_i, y_i)$; or
- 833 • $\text{EXPERT.UPDATE}(e_i, x_i, y_i)$ returned two experts $(S_{\in}, u_{\in}, H_{\in})$ and $(S_{\notin}, u_{\notin}, H_{\notin})$ (as in
 834 the third return statement of EXPERT.UPDATE), and furthermore,

$$e_{i+1} = \begin{cases} (S_{\in}, u_{\in}, H_{\in}) & x_i \in \text{path}(h) \\ (S_{\notin}, u_{\notin}, H_{\notin}) & x_i \notin \text{path}(h). \end{cases} \quad (18)$$

835 We say that an expert $e \in E_{t+1}$ is assumption-consistent with h if for all $i \in [t]$, the $i \rightarrow (i+1)$
 836 update of e was assumption-consistent with h .

837 **Claim D.4** (Existence of assumption-consistent expert). *Let $d, n, t \in \mathbb{N}$, $t \leq n$, let $\mathcal{H} \subseteq \{0, 1\}^{T_d}$, let*
 838 *$x_1, \dots, x_n \in T_d$, and let $h : T_d \rightarrow \{0, 1\}$. Consider an execution of*

$$\text{TRANSDUCTIVELEARNER}(\mathcal{H}, d, (x_1, x_2, \dots, x_n))$$

839 *as in Algorithm 5. Then, at the end of the t -th iteration of the outer ‘for’ loop in TRANSDUCTIVE-*
 840 *LEARNER, there exists a unique expert $e_{t+1}^* \in E_{t+1}$ that is assumption-consistent with h .*

841 *Proof.* We prove by induction that, for all $s \in [t + 1]$, E_s contains a unique expert that is assumption-
 842 consistent with h . The base case $s = 1$ is clear, because E_1 contains only a single expert that was
 843 never modified. For the induction step, let e_s^* be the unique assumption-consistent expert in E_s , and
 844 consider the $s \rightarrow (s + 1)$ update. Notice that by Definition D.3,

- 845 • For all $e \in E_s \setminus \{e_s^*\}$, every expert $e' \in E_{s+1}$ such that $e' = e$ or e' was created by
 846 executing $\text{EXPERT.UPDATE}(e_s, x_s, y_s)$ is not assumption-consistent with h ; and
- 847 • Either $\text{EXPERT.UPDATE}(e_s^*, x_s, y_s)$ is not executed and $e_s^* \in E_{s+1}$ (e_s^* is added
 848 to E_{s+1} without any modification), or precisely one of the experts created by
 849 $\text{EXPERT.UPDATE}(e_s^*, x_s, y_s)$ and added to E_{s+1} is assumption-consistent with h .

850 Seeing as the $s \rightarrow (s + 1)$ update executes at most $\text{EXPERT.UPDATE}(e, x_s, y_s)$ once for each $e \in E_s$,
 851 it follows that E_{s+1} contains precisely one expert that is assumption-consistent with h . \square

852 An expert $e = (S, u, H)$ that is assumption-consistent with the correct labeling function enjoys two
 853 simple properties. The first property is that the node u in the expert encodes correct information
 854 about which previously seen nodes are on-path for the correct labeling function.

855 The second property is that the set S contains all future nodes that are on-path for the correct labeling
 856 function and are also deeper in the tree than all nodes assumed to be on-path so far. These two
 857 properties are formalized in the following claim.

858 **Claim D.5** (Properties of assumption-consistent expert). *Let $d, n, t \in \mathbb{N}$, $t \leq n + 1$, let $\mathcal{H} \subseteq \{0, 1\}^{T_d}$,*
 859 *let $x_1, \dots, x_n \in T_d$. Consider an execution of*

$$\text{TRANSDUCTIVELEARNER}(\mathcal{H}, d, (x_1, x_2, \dots, x_n))$$

860 *as in Algorithm 5. Assume that the adversary selects labels $y_1, y_2, \dots, y_n \in \{0, 1\}$ that are consistent*
 861 *with some function $h : T_d \rightarrow \{0, 1\}$. Let $e_t^* = (S_t^*, u_t^*, H_t^*) \in E_t$ be the unique expert in E_t that is*
 862 *assumption-consistent with h .²⁶ Then the following two properties hold:*

- 863 1. $u_t^* \in \text{path}(h)$.
- 864 2. $\{x \in \{x_t, x_{t+1}, \dots, x_{t_{\max}}\} : x \in \text{path}(h) \wedge x \not\in u_t^*\} \subseteq S_t^*$.

865 *Proof of Claim D.5.* The proof proceeds by induction on t . For the base case $t = 1$, E_1 contains
 866 a single expert $e_1^* = (S_1^*, u_1^*, H_1^*)$ where $u_1^* = \lambda$ is the root of T_d . Indeed, $\lambda \in \text{path}(h)$ for
 867 any function $h : T_d \rightarrow \{0, 1\}$. This establishes the base case for Item 1. Additionally, $S_1^* =$
 868 $\{x_1, x_2, \dots, x_{t_{\max}}\}$, satisfying the base case for Item 2.

869 For the induction step, we assume that the claim holds for some integer $t = i$, and show that it
 870 holds for $t = i + 1$ as well. First, we establish Item 1. If $e_{i+1}^* = e_i^*$, then the claim is immediate.
 871 Otherwise, by Definition D.3 and the first two return statements in EXPERT.UPDATE , either
 872 $e_{i+1}^* = (S_{i+1}^*, u_{i+1}^*, H_{i+1}^*)$ has $u_{i+1}^* = u_i^* \in \text{path}(h)$, in which case the claim is immediate, or else
 873 e_{i+1}^* satisfies Eq. (18), namely,

$$e_{i+1}^* = \begin{cases} (S_{\in}, u_{\in}, H_{\in}) & x_i \in \text{path}(h) \\ (S_{\notin}, u_{\notin}, H_{\notin}) & x_i \notin \text{path}(h). \end{cases}$$

874 As defined in EXPERT.UPDATE , u_{\in} is equal either to u_i^* or to x_i , so if $x_i \in \text{path}(h)$ then

$$u_{i+1}^* = u_{\in} \in \{u_i^*, x_i\} \subseteq \text{path}(h).$$

²⁶Recall that e_t^* exists by Claim D.4.

On the other hand, if $x_i \notin \text{path}(h)$ then we get $u_{i+1}^* = u_{\notin} = u_i^* \in \text{path}(h)$. We see that in all cases, $u_{i+1}^* \in \text{path}(h)$ as desired. This concludes the proof of Item 1.

For Item 2, again, if $e_{i+1}^* = e_i^*$, then the claim is immediate. Otherwise, consider the various ways in which u_{i+1}^* and S_{i+1}^* can be assigned by EXPERT.UPDATE. In the first return statement, $u_{i+1}^* = u_i^*$ and $S_{i+1}^* = S_i^*$, and the claim is immediate.

The second return statement assigns $u_{i+1}^* = u_i^*$ and $S_{i+1}^* = S_i^* \setminus S_{1-y_i}$, where S_{1-y_i} is the set of $(1 - y_i)$ -descendants of x_i (including x_i itself). Notice that regardless of whether x_i is on-path for the correct labeling function h or not, none of the $(1 - y_i)$ -descendants of x_i (except possibly x_i itself) can be on-path for h , because h assigns a label y_i to x_i . And seeing as Item 2 only requires that S_{i+1}^* contain nodes from $\{x_{i+1}, x_{i+2}, \dots, x_{t_{\max}}\}$, it is also safe to remove x_i . Therefore, removing S_{1-y_i} preserves Item 2.

For the third return statement, there are two possibilities. The first possibility is that $u_{i+1}^* = u_{\notin} = u_i^*$ and $S_{i+1}^* = S_{\notin} = S_i^*$, in which case the claim is immediate. The second possibility assigns $u_{i+1}^* = u_{\in}$, and $S_{i+1}^* = S_{\in} = S_0 \cup S_1$, namely, S_{i+1}^* is constructed by removing the non-descendants of x_i from S_i^* . By Eq. (18), this happens when $x_i \in \text{path}(h)$, so all non-descendants of x_i or either off-path for h , or they are ancestors of x_i . Seeing as $x_i \in \text{path}(h)$ and $u_i^* \in \text{path}(h)$, and u_{\in} is the deeper node between these two, any node that is an ancestor of x_i is also an ancestor of $u_{i+1}^* = u_{\in}$. Thus, all the nodes removed or either off-path for h , or they are ancestors of u_{i+1}^* , satisfying Item 2. (Similarly, any node that is an ancestor of u_i^* is also an ancestor of u_{i+1}^* , so we do not need to add any new nodes to S_{i+1}^* that are not included in S_i^* .)

We see that in all cases, Item 2 is preserved, as desired. \square

D.4.2 Transition to Halving

Claim D.6. Let $d, n, t \in \mathbb{N}$, $d \geq 16$, let $\mathcal{H} \subseteq \{0, 1\}^{T_d}$, and let $x_1, \dots, x_n \in T_d$. Consider an execution of

$$\text{TRANSDUCTIVELEARNER}(\mathcal{H}, (x_1, x_2, \dots, x_n))$$

as in Algorithm 5. Let $t > t_{\max} = 2^{4\sqrt{d}}$ and let $e = (S, u, H) \in E_t$ be an expert. Then

$$|H| \leq 2^{2\sqrt{d}}.$$

Proof of Claim D.6. Assume for contradiction that $|H| > 2^{2\sqrt{d}}$. Let $H' \subseteq H$ be an arbitrary subset of size $2^{2\sqrt{d}} + 1$. Let

$$P = \cup_{h \in H'} \text{path}(h).$$

Seeing as each root-to-leaf path contains $d + 1$ nodes,

$$|P| \leq |H'| \cdot (d + 1) \leq (2^{2\sqrt{d}} + 1) \cdot (d + 1) \leq d2^{2\sqrt{d}+1}. \quad (19)$$

Let y_1, y_2, \dots, y_t be the labels provided by the adversary in the first t_{\max} iterations. The line in EXPERT.UPDATE constructing H using HALVING.UPDATE(H, x, y) ensures that

$$\forall h \in H \forall i \in [t_{\max}] : h(x_i) = y_i. \quad (20)$$

Consider two cases:

• **Case I.** $\sum_{i=1}^{t_{\max}} y_i \leq t_{\max}/2$. Then the set

$$X_0 = \{x_i : i \in [t_{\max}] \wedge y_i = 0\}$$

has cardinality $|X_0| \geq t_{\max}/2$. Let $X'_0 = X_0 \setminus P$. By Eq. (19),

$$|X'_0| \geq \frac{t_{\max}}{2} - d2^{2\sqrt{d}+1} = 2^{4\sqrt{d}} - d2^{2\sqrt{d}+1}. \quad (21)$$

From the choice of X'_0 , the inclusion $H' \subseteq H$, and Eq. (20),

$$\forall h \in H' \forall x \in X'_0 : x \notin \text{path}(h) \wedge h(x) = 0. \quad (22)$$

909 Seeing as $|H'| > 2^{2\sqrt{d}}$, Eq. (22) and Item 1 from Lemma D.2 imply that

$$|X'_0| \leq 2^{2\sqrt{d}}. \quad (23)$$

910 Combining Eqs. (21) and (23) yields

$$\begin{aligned} 2^{2\sqrt{d}} &\geq |X'_0| \geq 2^{4\sqrt{d}} - d2^{2\sqrt{d}+1} \\ &\geq 2^{4\sqrt{d}-1} \end{aligned} \quad (d \geq 16),$$

911 which is a contradiction.

912 • **Case II.** $\sum_{i=1}^{t_{\max}} y_i > t_{\max}/2$. A similar argument gives a contradiction by defining

$$X_1 = \{x_i : i \in [t_{\max}] \wedge y_i = 1\}, \text{ and } X'_1 = X_1 \setminus P.$$

913 As before,

$$|X'_1| \geq \frac{t_{\max}}{2} - d2^{2\sqrt{d}+1} \geq 2^{4\sqrt{d}} - d2^{2\sqrt{d}+1}. \quad (24)$$

914 for all $d \in \mathbb{N}$. However, $|H'| > 2^{2\sqrt{d}}$ and Item 2 imply that

$$|X'_1| < 3\sqrt{d}, \quad (25)$$

915 which is a contradiction. \square

916 D.4.3 Performance of Best Expert

917 **Claim D.7** (Existence of expert with large weight). *Let $d, n \in \mathbb{N}$, $d \geq 16$, let $\mathcal{H} \subseteq \{0, 1\}^{T_d}$, and let*
 918 *$x_1, \dots, x_n \in T_d$. Consider an execution of*

$$\text{TRANSDUCTIVELEARNER}(\mathcal{H}, (x_1, x_2, \dots, x_n))$$

919 *as in Algorithm 5. Then, at the end of the execution, there exists $e \in E_{n+1}$ such that*

$$w(e) \geq 2^{-48\sqrt{d}}.$$

920 Note that the lower bound in Claim D.7 does not depend on n .

921 *Proof.* Fix a hypothesis $h \in \mathcal{H}$ such that $h(x_t) = y_t$ for all $t \in [n]$ (such an h exists because the
 922 adversary must always select a realizable label).

923 By Claim D.4, there exists $e_{n+1}^* \in E_{n+1}$ that is assumption-consistent with h . Let $\text{ancestry}(e_{n+1}^*) =$
 924 $(e_1^*, e_2^*, \dots, e_{n+1}^*)$. We argue that this ancestry sequence makes few mistakes. Specifically, for each
 925 $t \in [n]$, let $\hat{y}_t^* = \text{EXPERT.PREDICT}(e_t^*, x_t)$. We claim that

$$m = \sum_{t=1}^n \mathbb{1}(\hat{y}_t^* \neq y_t) \leq 24\sqrt{d}.$$

926 Indeed, let $B = \{t \in [n] : \hat{y}_t^* \neq y_t\}$ be the set of m indices where a mistake was made. For
 927 each $t \in B$, let $e_t^* = (S, u, H)$, and note that each $t \in B$ has a corresponding execution of
 928 $\text{EXPERT.PREDICT}(e_t^*, x_t)$, and an execution of $\text{EXPERT.UPDATE}(e_t^*, x_t, y_t)$ that produces e_{t+1}^* . We
 929 partition the indices in B into six cases (six sets), and bound the number of indices that fall in each.

930 • **Case I.** *The execution of $\text{EXPERT.PREDICT}(e_t^*, x_t)$ exited via the first return statement in that*
 931 *procedure. This happens once $|H| \leq 2^{2\sqrt{d}}$, and from that point on, the expert and*
 932 *all subsequent experts in the ancestry are exactly simulating the HALVING algorithm*
 933 *(Algorithm 7) in both predictions and updates. Hence, by Fact E.1, B contains at most*
 934 *$m_I = 2\sqrt{d}$ such indices.*

935 • **Case II.** *The execution of $\text{EXPERT.PREDICT}(e_t^*, x_t)$ exited via the second return statement in*
 936 *that procedure. In particular $x \preceq u$, and the predicted label was $\hat{y}_t^* = b \in \{0, 1\}$ such*
 937 *that $x_t \preceq_b u$. Because e_t^* is assumption-consistent with h , Item 1 in Claim D.5 implies*
 938 *that $u \in \text{path}(h)$. Namely, we see that u is a b -descendant of x_t and $u \in \text{path}(h)$. It*
 939 *follows that $\hat{y}_t^* = b = h(x_t) = y_t$. So no mistakes are made in Case II, and the number*
 940 *of indices $t \in B$ that belong to Case II is simply $m_{II} = 0$.*

941 In the remaining cases, we assume that $\text{EXPERT.PREDICT}(e_t^*, x_t)$ exited via the third return statement
 942 in that procedure, so the prediction was

$$\hat{y}_t^* = \mathbb{1}\left(\frac{|S_1|}{|S|} \geq \frac{1}{3}\right), \quad (26)$$

943 where $S_1 = \{x' \in S : x_t \preccurlyeq_1 x'\}$. These cases are as follows.

944 • **Case III.** *The execution of $\text{EXPERT.UPDATE}(e_t^*, x_t, y_t)$ exited via the first return statement in*
 945 *that procedure.* Namely, after the update, the resulting expert e_{t+1}^* has $|H| \leq 2^{2\sqrt{d}}$.
 946 However, because we are not in Case I, at the beginning of the iteration expert e_t^* had
 947 $|H| > 2^{2\sqrt{d}}$. Seeing as the cardinality of H decreases monotonically throughout the
 948 ancestry e_1^*, \dots, e_{n+1}^* , this type of mistake can happen at most $m_{\text{III}} = 1$ times.

949 • **Case IV.** *The execution of $\text{EXPERT.UPDATE}(e_t^*, x_t, y_t)$ exited via the second return statement*
 950 *in that procedure.* In this case, $|S_{(1-y_t)}|/|S| \geq 1/3$, and $e_{t+1}^* = (\bar{S}', u, \bar{H})$ with
 951 $S' = S \setminus S_{1-y_t}$. So $|S'| \leq 2|S|/3$. Namely, the update causes the cardinality of the
 952 set S to be multiplied by a factor of at most $2/3$ and it strictly decreases. Seeing as the
 953 initial cardinality is t_{\max} , and cardinalities are integers, the number of times this can
 954 happen is at most

$$m_{\text{IV}} = \frac{\log(t_{\max})}{\log(3/2)} + 1 = \frac{4\sqrt{d}}{\log(3/2)} + 1. \quad (27)$$

955 In the remaining cases, we assume that the execution of $\text{EXPERT.UPDATE}(e_t^*, x_t, y_t)$ exited via the
 956 third return statement in that procedure. This implies that

$$|S_{\hat{y}_t^*}|/|S| < 1/3 \quad (28)$$

957 Combining this with Eq. (26), it follows $\hat{y}_t^* = 0$ and therefore $y_t = 1$. The remaining cases are as
 958 follows.

959 • **Case V.** $x_t \in \text{path}(h)$. Let $e_t^* = (S, u, H)$. Seeing as $|H| > 2^{2\sqrt{d}}$ (because we are not in
 960 Case I), Claim D.6 (with the assumption $d \geq 16$) implies that $t \leq t_{\max}$. By Item 2
 961 of Claim D.5, the facts $x_t \not\preccurlyeq u$ (we are not in Case II) and $x_t \in \text{path}(h)$ imply that
 962 $x_t \in S$. In particular, S is not empty.

963 Because the $t \rightarrow (t+1)$ update of e_{t+1}^* was assumption-consistent with h , Eq. (18)
 964 implies that $e_{t+1}^* = (S_{\in}, u_{\in}, H_{\in})$, with $S_{\in} = S_0 \cup S_1$. Observe that

- 965 • $|S_0|/|S| < 1/3$ (plugging $\hat{y}_t^* = 0$ into Eq. (28)); and
- 966 • $|S_1|/|S| < 1/3$ (because otherwise, by Eq. (26), the prediction would have been
 967 $\hat{y}_t^* = 1$).

968 Therefore,

$$|S_{\in}| \leq |S_0| + |S_1| \leq 2|S|/3. \quad (29)$$

969 As in Case IV, combining Eq. (29) and the fact that S is not empty imply an upper
 970 bound m_{V} on the number of times Case V can happen, with the bound being the same
 971 number $m_{\text{V}} = m_{\text{IV}}$ as in Eq. (27).

972 • **Case VI.** $x_t \notin \text{path}(h)$. So (x_t, y_t) is a pair such that $x_t \notin \text{path}(h)$ and $y_t = 1$. Assume for
 973 contradiction that this type of mistake can happen strictly more than

$$m_{\text{VI}} = 3\sqrt{d}$$

974 times. Let $t_1, t_2, \dots, t_{m_{\text{VI}}}$ be the indices of the first m_{VI} iterations of the outer ‘for’
 975 loop of $\text{TRANSDUCTIVELEARNER}$ in which this type of mistake happened. Note that
 976 if at the end of iteration $t_{m_{\text{VI}}}$, we had expert $e_{t_{m_{\text{VI}}}+1}^* = (S_{t_{m_{\text{VI}}}+1}, u_{t_{m_{\text{VI}}}+1}, H_{t_{m_{\text{VI}}}+1})$
 977 such that $|H_{t_{m_{\text{VI}}}+1}| \leq 2^{2\sqrt{d}}$, then from that point onwards, the expert would be
 978 simulating the halving algorithm, and in particular, it would not make any further
 979 mistake of the type in Case VI (all subsequent mistakes would belong to Case I). Hence,

980 by the assumption that strictly more than m_{VI} mistakes were made, it follows that
 981 $|H_{t_{m_{\text{VI}}+1}}| > 2^{2\sqrt{d}}$. Let

$$H^* = \{h' \in \mathcal{H} : (\forall t \in [m_{\text{VI}}] : h'(x_t) = 1 \wedge x_t \notin \text{path}(h'))\}.$$

982 Because $e_{t_{m_{\text{VI}}+1}}^*$ is assumption-consistent with h , and from the construction of $H_{t_{m_{\text{VI}}+1}}$
 983 using H_{\in} and H_{\notin} in EXPERT.UPDATE, it follows that $H_{t_{m_{\text{VI}}+1}} \subseteq H^*$. So there exist
 984 collections $H^* \subseteq \mathcal{H}$ and $X = \{x_t : t \in [m_{\text{VI}}]\} \subseteq T_d$ such that

- 985 • $|H^*| \geq |H_{t_{m_{\text{VI}}+1}}| > 2^{2\sqrt{d}}$,
- 986 • $|X| = m_{\text{VI}} = 3\sqrt{d}$,
- 987 • $\forall h' \in H^* \forall x \in X : h'(x) = 1$.
- 988 • $\forall h' \in H^* \forall x \in X : x \notin \text{path}(h')$.

989 This is a contradiction to the choice of \mathcal{H} , specifically, to Item 2 in Lemma D.2.

990 Thus, combining the analysis of all cases, we see that the number of mistakes made by the
 991 ancestry(e_{n+1}^*) is at most

$$\begin{aligned} m &\leq m_{\text{I}} + m_{\text{II}} + m_{\text{III}} + m_{\text{IV}} + m_{\text{V}} + m_{\text{VI}} \\ &\leq 2\sqrt{d} + 0 + 1 + \left(\frac{4\sqrt{d}}{\log(3/2)} + 1 \right) + \left(\frac{4\sqrt{d}}{\log(3/2)} + 1 \right) + 3\sqrt{d} \\ &\leq 24\sqrt{d}. \end{aligned}$$

992 The weights satisfy

$$w(e_{t+1}^*) \begin{cases} = w(e_t^*) & \hat{y}_t^* = y_t \\ \geq \frac{1}{4} \cdot w(e_t^*) & \hat{y}_t^* \neq y_t. \end{cases}$$

993 This implies that $w(e_{n+1}^*) \geq w(e_1^*) \cdot \prod_{t=1}^n 4^{-\mathbb{1}(\hat{y}_t \neq y_t)} = w(e_1^*) \cdot 4^{-m} \geq 4^{-24\sqrt{d}} = 2^{-48\sqrt{d}}$, as
 994 desired. \square

995 D.4.4 Multiplicative Weights Mistake Bound

996 **Claim D.8** (Mistake bound for multiplicative weights). *Let $d, n \in \mathbb{N}$, let $\alpha > 0$, let $\mathcal{H} \subseteq \{0, 1\}^{T_d}$,
 997 and let $x_1, \dots, x_n \in T_d$. Consider an execution of*

$$\text{TRANSDUCTIVELEARNER}(\mathcal{H}, (x_1, x_2, \dots, x_n))$$

998 *as in Algorithm 5. Assume that at the end of the execution, there exists $e^* \in E_{n+1}$ such that*

$$w(e^*) \geq 2^{-\alpha}.$$

999 *Then TRANSDUCTIVELEARNER makes at most α mistakes.*

1000 *Proof of Claim D.8.* For all $i \in [n+1]$, let $w(E_i) = \sum_{e \in E_i} w(e)$. For each $i \in [n]$, if $\hat{y}_i \neq y_i$, then
 1001 $w(E_{i+1}) \leq w(E_i)/2$. Hence, if TRANSDUCTIVELEARNER makes m mistakes, then by induction

$$w(E_{n+1}) \leq w(E_1) \cdot \prod_{t=1}^n 2^{-\mathbb{1}(\hat{y}_t \neq y_t)} = 2^{-m} \cdot w(E_1).$$

1002 So

$$2^{-\alpha} \leq w(e^*) \leq \sum_{e \in E_{n+1}} w(e) = w(E_{n+1}) \leq 2^{-m} \cdot w(E_1) = 2^{-m}.$$

1003 We conclude that

$$m \leq \alpha,$$

1004 as desired. \square

1005 D.5 Proof

1006 *Proof of Theorem D.1.* Fix an integer $d \geq 23$. Let $\mathcal{H} \subseteq \{0, 1\}^{T_{d-1}}$ be the class constructed by
 1007 invoking Lemma D.2 for the integer $d - 1 \geq 22$. We argue that this class satisfies the requirements
 1008 of Theorem D.1.

1009 By construction, \mathcal{H} is a class of Littlestone dimension precisely d . By Theorem A.7, this implies the
 1010 equality in Item 2.

1011 We now show the upper bound in Item 1. We argue that TRANSDUCTIVELERARNER (Algorithm 5)
 1012 satisfies this upper bound. By Claim D.7, at the end of the execution of TRANSDUCTIVELERARNER
 1013 there exists an expert $e \in E_{n+1}$ such that $w(e) \geq 2^{-48\sqrt{d}}$. By Claim D.8, this implies that the
 1014 number of mistakes made by TRANSDUCTIVELERARNER is at most $48\sqrt{d}$, as desired. \square

1015 E Halving

1016 **Fact E.1.** Let \mathcal{X} be a set, and let $\mathcal{H} \subseteq \{0, 1\}^{\mathcal{X}}$ be a hypothesis class. Then for all $n \in \mathbb{N}$, all
 1017 sequences $x \in \mathcal{X}^n$, and all realizable adversaries, HALVING (Algorithm 7) makes at most $\log(|\mathcal{H}|)$
 1018 mistakes in the transductive online learning (Game 2).²⁷ Namely,

$$\sup_{n \in \mathbb{N}} \sup_{A \in \mathcal{A}_n} M_{\text{tr}}(\mathcal{H}, n, \text{HALVING}, A) \leq \log(|\mathcal{H}|).$$

²⁷With the suitable syntactic modification, it also makes at most $\log(|\mathcal{H}|)$ mistakes in the standard online learning (Game 1).

1019 **NeurIPS Paper Checklist**

1020 **1. Claims**

1021 Question: Do the main claims made in the abstract and introduction accurately reflect the
1022 paper’s contributions and scope?

1023 Answer: [\[Yes\]](#)

1024 Justification: The main claims made in the abstract and introduction accurately reflect the
1025 paper’s contributions and scope.

1026 Guidelines:

- 1027 • The answer NA means that the abstract and introduction do not include the claims
1028 made in the paper.
- 1029 • The abstract and/or introduction should clearly state the claims made, including the
1030 contributions made in the paper and important assumptions and limitations. A No or
1031 NA answer to this question will not be perceived well by the reviewers.
- 1032 • The claims made should match theoretical and experimental results, and reflect how
1033 much the results can be expected to generalize to other settings.
- 1034 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1035 are not attained by the paper.

1036 **2. Limitations**

1037 Question: Does the paper discuss the limitations of the work performed by the authors?

1038 Answer: [\[NA\]](#)

1039 Justification: Purely rigorous mathematical results. We explain precisely what our proofs
1040 imply (and therefore also what they do not imply).

1041 Guidelines:

- 1042 • The answer NA means that the paper has no limitation while the answer No means that
1043 the paper has limitations, but those are not discussed in the paper.
- 1044 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1045 • The paper should point out any strong assumptions and how robust the results are to
1046 violations of these assumptions (e.g., independence assumptions, noiseless settings,
1047 model well-specification, asymptotic approximations only holding locally). The authors
1048 should reflect on how these assumptions might be violated in practice and what the
1049 implications would be.
- 1050 • The authors should reflect on the scope of the claims made, e.g., if the approach was
1051 only tested on a few datasets or with a few runs. In general, empirical results often
1052 depend on implicit assumptions, which should be articulated.
- 1053 • The authors should reflect on the factors that influence the performance of the approach.
1054 For example, a facial recognition algorithm may perform poorly when image resolution
1055 is low or images are taken in low lighting. Or a speech-to-text system might not be
1056 used reliably to provide closed captions for online lectures because it fails to handle
1057 technical jargon.
- 1058 • The authors should discuss the computational efficiency of the proposed algorithms
1059 and how they scale with dataset size.
- 1060 • If applicable, the authors should discuss possible limitations of their approach to
1061 address problems of privacy and fairness.
- 1062 • While the authors might fear that complete honesty about limitations might be used by
1063 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
1064 limitations that aren’t acknowledged in the paper. The authors should use their best
1065 judgment and recognize that individual actions in favor of transparency play an impor-
1066 tant role in developing norms that preserve the integrity of the community. Reviewers
1067 will be specifically instructed to not penalize honesty concerning limitations.

1068 **3. Theory assumptions and proofs**

1069 Question: For each theoretical result, does the paper provide the full set of assumptions and
1070 a complete (and correct) proof?

Answer: [Yes]

Justification: For each theoretical result, the paper provides the full set of assumptions and a complete (and correct) proof.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [NA]

Justification: The paper has no experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

1124 Question: Does the paper provide open access to the data and code, with sufficient instruc-
 1125 tions to faithfully reproduce the main experimental results, as described in supplemental
 1126 material?

1127 Answer: [NA]

1128 Justification: The paper does not include experiments requiring code.

1129 Guidelines:

- 1130 • The answer NA means that paper does not include experiments requiring code.
- 1131 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
 1132 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1133 • While we encourage the release of code and data, we understand that this might not be
 1134 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
 1135 including code, unless this is central to the contribution (e.g., for a new open-source
 1136 benchmark).
- 1137 • The instructions should contain the exact command and environment needed to run to
 1138 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
 1139 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1140 • The authors should provide instructions on data access and preparation, including how
 1141 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1142 • The authors should provide scripts to reproduce all experimental results for the new
 1143 proposed method and baselines. If only a subset of experiments are reproducible, they
 1144 should state which ones are omitted from the script and why.
- 1145 • At submission time, to preserve anonymity, the authors should release anonymized
 1146 versions (if applicable).
- 1147 • Providing as much information as possible in supplemental material (appended to the
 1148 paper) is recommended, but including URLs to data and code is permitted.

1149 **6. Experimental setting/details**

1150 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
 1151 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
 1152 results?

1153 Answer: [NA]

1154 Justification: The paper does not include experiments.

1155 Guidelines:

- 1156 • The answer NA means that the paper does not include experiments.
- 1157 • The experimental setting should be presented in the core of the paper to a level of detail
 1158 that is necessary to appreciate the results and make sense of them.
- 1159 • The full details can be provided either with the code, in appendix, or as supplemental
 1160 material.

1161 **7. Experiment statistical significance**

1162 Question: Does the paper report error bars suitably and correctly defined or other appropriate
 1163 information about the statistical significance of the experiments?

1164 Answer: [NA]

1165 Justification: The paper does not include experiments.

1166 Guidelines:

- 1167 • The answer NA means that the paper does not include experiments.
- 1168 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
 1169 dence intervals, or statistical significance tests, at least for the experiments that support
 1170 the main claims of the paper.
- 1171 • The factors of variability that the error bars are capturing should be clearly stated (for
 1172 example, train/test split, initialization, random drawing of some parameter, or overall
 1173 run with given experimental conditions).
- 1174 • The method for calculating the error bars should be explained (closed form formula,
 1175 call to a library function, bootstrap, etc.)

- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: The paper does not include experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: The research conducted in the paper conforms, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [No]

Justification: The work is purely theoretical with no immediate direct societal impacts foreseeable.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: The paper does not use existing assets.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

1331 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1332 non-standard component of the core methods in this research? Note that if the LLM is used
1333 only for writing, editing, or formatting purposes and does not impact the core methodology,
1334 scientific rigorousness, or originality of the research, declaration is not required.

1335 Answer: [NA]

1336 Justification: The core method development in this research does not involve LLMS as any
1337 important, original, or non-standard components.

1338 Guidelines:

- 1339 • The answer NA means that the core method development in this research does not
1340 involve LLMs as any important, original, or non-standard components.
- 1341 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
1342 for what should or should not be described.