
The Minimax Regret of Sequential Probability Assignment, Contextual Shtarkov Sums, and Contextual Normalized Maximum Likelihood

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the fundamental problem of sequential probability assignment, also
2 known as online learning with logarithmic loss, with respect to an arbitrary, pos-
3 sibly nonparametric hypothesis class. Our goal is to obtain a complexity measure
4 for the hypothesis class that characterizes the minimax regret and to determine
5 a general, minimax optimal algorithm. Notably, the sequential ℓ_∞ entropy, ex-
6 tensively studied in the literature (Rakhlin and Sridharan, 2015, Bilodeau et al.,
7 2020, Wu et al., 2023), was shown to not characterize minimax risk in general. In-
8 spired by the seminal work of Shtarkov (1987) and Rakhlin, Sridharan, and Tewari
9 (2010), we introduce a novel complexity measure, the *contextual Shtarkov sum*,
10 corresponding to the Shtarkov sum after projection onto a multiary context tree,
11 and show that the worst case log contextual Shtarkov sum equals the minimax re-
12 gret. Using the contextual Shtarkov sum, we derive the minimax optimal strategy,
13 dubbed *contextual Normalized Maximum Likelihood* (cNML). Our results hold
14 for sequential experts, beyond binary labels, which are settings rarely considered
15 in prior work. To illustrate the utility of this characterization, we provide a short
16 proof of a new regret upper bound in terms of sequential ℓ_∞ entropy, unifying
17 and sharpening state-of-the-art bounds by Bilodeau et al. (2020) and Wu et al.,
18 (2023).

19 1 Introduction

20 Sequential probability assignment is a fundamental problem with connections to information theory
21 [Ris84; MF98; XB00], machine learning [CL06; Vov95; RST15; FKLMS18; Sha20], and portfolio
22 optimization [Kel56; Cov74; Cov91; CO96; Fed91]. In the original non-contextual setup, the learner
23 aims to assign probabilities to a series of labels, which are revealed sequentially. The goal is to offer
24 probabilistic forecasts over the label set such that the probability assigned to any observed sequence
25 is comparable to that assigned by the best model in any fixed class of models.

26 The celebrated work of Shtarkov [Sht87] characterized minimax regret for context-free sequential
27 probability assignment in terms of what is now known as the *Shtarkov sum*, and subsequently de-
28 scribed the minimax algorithm, *Normalized Maximum Likelihood* (NML). NML represents the ideal
29 probabilistic forecast in the sense of minimax regret, providing a benchmark for universal coding and
30 prediction strategies. While often not used directly due to its computational complexity, NML has
31 guided the design of practical algorithms and informed the development of efficient approximation
32 methods. The principles underlying NML have inspired advances in both information theory and
33 online learning, establishing fundamental limits and serving as critical benchmarks for performance
34 evaluation.

35 In this work, we study the problem of sequential probability assignment with contexts, which
 36 has been analyzed in recent works (e.g. [RS15; BFR20; WHGS23]) under the framework of on-
 37 line supervised learning formalized by Rakhlin, Sridharan, and Tewari [RST10]. In this setup,
 38 the problem is modeled as a T -round game between a *learner* and the *nature*: On each round
 39 $t = 1, \dots, T$, the learner observes a context x_t from nature and predicts a distribution \hat{p}_t over some
 40 finite label space \mathcal{Y} . Then nature reveals a label $y_t \in \mathcal{Y}$ and the learner incurs a logarithmic loss
 41 $\ell(\hat{p}_t, y_t) = -\log(\hat{p}_t(y_t))$, where $\hat{p}_t(y_t)$ is the probability assigned to label y_t by \hat{p}_t . The perfor-
 42 mance of the learner is measured by the *regret* with respect to a class \mathcal{F} of *experts*, defined as the
 43 difference between the total loss of the learner and that of the best expert in \mathcal{F} . The value of primary
 44 interest is the *minimax regret*, that is, the worst-case regret by the best learner over arbitrarily adap-
 45 tive data sequences. The minimax regret serves as a benchmark for all algorithms and as a target for
 46 studies of adaptivity. Our goal is to address several fundamental questions:

47 *Can we find a natural complexity measure of \mathcal{F} that characterizes the minimax regret, enabling us*
 48 *to analyze the minimax regret in new ways? And can we identify, in view of this complexity*
 49 *measure, a general, minimax optimal algorithm?*

50 Notably, the sequential covering number of \mathcal{F} , a well studied measure of complexity, has been
 51 shown not to characterize the minimax regret on its own [RS15; BFR20; WHGS23]. This fact
 52 distinguishes sequential probability assignment and log loss: while sequential covering numbers
 53 enable a tight analysis in online learning problems with convex Lipschitz losses, like absolute loss
 54 [RST15] and square loss [RS14a], they do not yield minimax rates for log loss on some classes.
 55 Tackling such classes evidently requires new techniques.

56 **Main contributions.**

- 57 1. We introduce a new complexity measure, which we call the *contextual Shtarkov sum*, that serves
 58 as a natural generalization of the Shtarkov sum from the context-free setting. We show that the
 59 minimax regret is characterized by the worst-case contextual Shtarkov sum.
- 60 2. We derive the minimax optimal algorithm, dubbed *contextual Normalized Maximum Likelihood*
 61 (cNML), using a data-dependent variant of the contextual Shtarkov sum, thereby generalizing
 62 NML from the context-free setting.
- 63 3. We apply contextual Shtarkov sums to the study of sequential entropy bounds on the minimax
 64 regret. Doing so, we provide a short proof of a new regret upper bound in terms of sequential
 65 entropy that unifies and even *improves* on state-of-the-arts bounds by [BFR20] and [WHGS23].
 66 Our results extend beyond the binary label setting studied by recent work to arbitrary finite label
 67 sets.

68 **Related work.** Sequential probability assignment has been studied extensively. Various aspects
 69 of this problem have been investigated, including sequences with and without side information
 70 (contexts), parametric and nonparametric hypothesis classes, and stochastic or adversarial data-
 71 generating mechanisms. This problem has a long history in the machine learning community, see
 72 [CL06, Ch. 9] and the references therein. In the information theory community, this problem is also
 73 known as universal prediction [MF98], where the regret is also referred to as redundancy with respect
 74 to a set of codes. This has been studied in both stochastic and adversarial settings [Fre96; Ris86;
 75 Ris96; Sht87; XB97; DS04; MF98; OS04; Sha06; Szp98], where the focus was primarily on the
 76 parametric classes of experts. Closely related topics include universal source coding [Kol65; Sol64;
 77 Fit66; Dav73], data compression with arithmetic coding [Ris76; RL81; ZL77; ZL78; FMG92], and
 78 the minimum description length (MDL) principle [Ris78; Ris84; Ris87; BRY98; Grü05; HY01].
 79 Lastly, this topic is intimately tied with sequential gambling and portfolio optimization, as pointed
 80 out by [Kel56; Cov74; Cov91; CO96; Fed91].

81 A classical result in context-free sequential probability assignment is that the minimax regret is
 82 equal to the log contextual Shtarkov sum [Sht87], and the minimax algorithm is the well-known
 83 Normalized Maximum Likelihood. When the set of contexts is known in advance to the forecaster,
 84 namely, a fixed design setting, the minimax regret is equivalent to the log Shtarkov sum of the
 85 function class when projected onto the known set of contexts [JSS21; WHGS23].

86 To handle rich hypothesis classes, [CL99; OH99] upper bounded the regret in terms of the (non-
 87 sequential) uniform covering number of the class. However, this complexity measure proved to be

88 insufficient for obtaining optimal rates. [RS15] improved regret upper bounds by proposing a se-
 89 quential covering measure. Thereafter, by utilizing the self-concordance property and the curvature
 90 of the log loss, [BFR20] further improved the upper bound in terms of the sequential covering num-
 91 ber for nonparametric Lipschitz classes, through a non-constructive proof. [WHGS23] proposed a
 92 Bayesian algorithmic approach in order to upper bound the regret using a global notion of sequential
 93 covering. Notably, both the global and local sequential covering numbers do not fully characterize
 94 the regret, and the algorithm in [WHGS23] is not minimax optimal.

95 Online learning with respect to arbitrary hypothesis classes and the zero-one loss, in the realizable
 96 case, is known to be characterized by the Littlestone dimension [Lit88]. The agnostic case was ad-
 97 dressed by [BPS09; ABDMNY21]. Understanding sequential complexities in online learning with
 98 Lipschitz losses was extensively studied by [RST10; RS14a; RS14b; RST15]. Note that the logarith-
 99 mic loss is neither Lipschitz nor bounded. Recently, [AHKKV23] characterized online regression
 100 in the realizable case, for any approximate pseudo-metric, such as the ℓ_p loss.

101 2 Preliminaries

102 **Notation.** For a positive integer K , let $[K] := \{1, 2, \dots, K\}$. For a finite set \mathcal{K} with $|\mathcal{K}| = K$, we
 103 use $\Delta(\mathcal{K})$ to denote the set of all distributions on \mathcal{K} . We may identify \mathcal{K} with $[K]$ (under arbitrary
 104 enumeration of elements in \mathcal{K}) and treat elements of $\Delta(\mathcal{K})$ as vectors in \mathbb{R}^K . For a vector $p \in \mathbb{R}^K$
 105 and $i \in [K]$, let $p(i)$ be the i -th coordinate of p . Let $\Delta^+(\mathcal{K}) = \{p \in \Delta(\mathcal{K}) : p(k) > 0, \forall k \in \mathcal{K}\}$.
 106 For a general finite sequence $(a_i)_{i=1}^N$, we will use $a_{n:m}$ to denote the sub-sequence (a_n, \dots, a_m) for
 107 any $n \leq m$ and the empty sequence for $n > m$. For any set \mathcal{A} , let $\mathcal{A}^* = \cup_{k \geq 0} \mathcal{A}^k$ be the set of all
 108 finite length sequences over \mathcal{A} .

109 **Sequential probability assignment and minimax regret.** Let \mathcal{X} be the context space and \mathcal{Y} be the
 110 finite label space. In each round $t \in [T]$ during the game of sequential probability assignment, the
 111 learner receives a context $x_t \in \mathcal{X}$ from nature and assigns a probability distribution $\hat{p}_t \in \Delta(\mathcal{Y})$ to
 112 the possible labels. Then nature reveals the true label $y_t \in \mathcal{Y}$ and the learner incurs a loss $\ell(\hat{p}_t, y_t) =$
 113 $-\log(\hat{p}_t(y_t))$. Throughout, the learner is required to predict nearly as well as the best expert from
 114 an expert class, which is modeled as an arbitrary hypothesis class $\mathcal{F} \subseteq \{(\mathcal{X} \times \mathcal{Y})^* \times \mathcal{X} \rightarrow \Delta(\mathcal{Y})\}$.
 115 More formally, the goal of the learner is make their *regret* with respect to \mathcal{F} ,

$$\mathcal{R}_T(\mathcal{F}; \hat{p}_{1:T}, x_{1:T}, y_{1:T}) = \sum_{t=1}^T \ell(\hat{p}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f(x_{1:t}, y_{1:t-1}), y_t),$$

116 as small as possible for all sequences \mathbf{x} and \mathbf{y} generated by nature, possibly in an adversarial manner.
 117 Here $f(x_{1:t}, y_{1:t-1}) \in \Delta(\mathcal{Y})$ can be understood as the prediction made by expert f at round t using
 118 past observations $(x_{1:t-1}, y_{1:t-1})$ as well as the fresh context x_t . The main focus is to study the
 119 *minimax regret* $\mathcal{R}_T(\mathcal{F})$, which can be written as the following extensive form

$$\mathcal{R}_T(\mathcal{F}) = \sup_{x_1} \inf_{\hat{p}_1} \sup_{y_1} \dots \sup_{x_T} \inf_{\hat{p}_T} \sup_{y_T} \mathcal{R}_T(\mathcal{F}; \hat{p}_{1:T}, x_{1:T}, y_{1:T}),$$

120 where $x_t \in \mathcal{X}$, $\hat{p}_t \in \Delta(\mathcal{Y})$ and $y_t \in \mathcal{Y}, \forall t \in [T]$.

121 **Remark 2.1 (Sequential vs non-sequential experts)** Experts f as mappings from $(\mathcal{X} \times \mathcal{Y})^* \times \mathcal{X}$ to
 122 $\Delta(\mathcal{Y})$ are sometimes called *fully sequential* experts [WHGS23] due to their ability to predict based
 123 on the past history. However, the literature (e.g. [RS15; BFR20; WHGS23]) often considers the
 124 more limited notion of *non-sequential* experts, modeled as $\mathcal{F} \subseteq \{\mathcal{X} \rightarrow \Delta(\mathcal{Y})\}$, reflecting the fact
 125 that prediction made by each expert f is simply $f(x_t)$ in each round t . In contrast, our results are
 126 more general as our novel techniques can be applied to the more flexible sequential experts.

127 **Multiary trees.** The complexity of online learning problems stems from the sequential and adaptive
 128 nature of the adversary, which we can capture with *multiary trees*. Formally, for a general space \mathcal{A}
 129 and a finite set \mathcal{K} , an \mathcal{A} -valued \mathcal{K} -ary tree \mathbf{v} of depth d is a sequence of mappings $v_t : \mathcal{K}^{t-1} \rightarrow \mathcal{A}$ for
 130 $t \in [d]$. A *path* in a depth- d tree is a sequence $\varepsilon = (\varepsilon_1, \dots, \varepsilon_d) \in \mathcal{K}^d$. We use the notation $\mathbf{v}_t(\varepsilon)$ to
 131 denote $\mathbf{v}_t(\varepsilon_1, \dots, \varepsilon_{t-1})$ for $t \in [d]$ and the boldface notation $\mathbf{v}(\varepsilon)$ to denote $(\mathbf{v}_1(\varepsilon), \dots, \mathbf{v}_d(\varepsilon)) \in$
 132 \mathcal{A}^d . Throughout we will only consider \mathcal{Y} -ary trees valued in either \mathcal{X} or $\Delta(\mathcal{Y})$, where the paths are
 133 denoted by the boldface \mathbf{y} . We refer to \mathcal{X} -valued trees as *context trees* and $\Delta(\mathcal{Y})$ -valued trees as
 134 *probabilistic trees*.

135 **Time-varying context sets.** So far we consider the context set \mathcal{X} to be constant over time. But
 136 all of our results can be extended easily to allow for time-varying context spaces. Details of this
 137 generalization can be found in Appendix C.

138 2.1 Prior work: the Shtarkov sum in context-free and fixed designs

139 Before introducing our complexity measure that characterizes $\mathcal{R}_T(\mathcal{F})$, we review some prior set-
 140 tings where the minimax regret can be characterized by the well-studied *Shtarkov sum*. First we
 141 introduce the notion of likelihood of a hypothesis f with respect to a context and label sequence,
 142 which plays a key role in defining complexity measures and optimal algorithms.

143 **Definition 2.2 (Likelihood)** For $f : (\mathcal{X} \times \mathcal{Y})^* \times \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ and length- d sequences $x_{1:d} \in$
 144 $\mathcal{X}^d, y_{1:d} \in \mathcal{Y}^d$, the likelihood $P_f(y_{1:d}|x_{1:d})$ is defined as

$$P_f(y_{1:d}|x_{1:d}) = \prod_{t=1}^d f(x_{1:t}, y_{1:t-1})(y_t),$$

145 where we use the compact notation $f(x_{1:t}, y_{1:t-1})$ for $f(x_1, y_1, \dots, x_{t-1}, y_{t-1}, x_t)$.

146 In the classical context-free setting where \mathcal{X} can be thought of as a singleton, any sequential expert
 147 f degenerates to a joint distribution over label sequences. Indeed, given any label sequence $y_{1:t-1}$,
 148 $f(y_{1:t-1}) \in \Delta(\mathcal{Y})$ can be interpreted as the conditional distribution f assigns to the next label y_t .

149 We use $P_f(y_{1:d}) = \prod_{t=1}^d f(y_{1:t-1})(y_t)$ to denote this distribution. Similarly, the learner's strategy
 150 is also specified by a joint distribution that is decomposed to a sequence of conditional distributions
 151 $\hat{p}_t = \hat{p}_t(\cdot|y_{1:t-1}) \in \Delta(\mathcal{Y})$. In this setup the minimax regret $\mathcal{R}_T(\mathcal{F})$ is characterized by the Shtarkov
 152 sum [Sht87].

153 **Proposition 2.3 ([Sht87])** *In the context-free setting, for any hypothesis class \mathcal{F} and horizon T , the*
 154 *Shtarkov sum $S_T(\mathcal{F})$ is defined as*

$$S_T(\mathcal{F}) = \sum_{y_{1:T} \in \mathcal{Y}^T} \sup_{f \in \mathcal{F}} P_f(y_{1:T}).$$

155 *Moreover, the minimax regret is given by $\mathcal{R}_T(\mathcal{F}) = \log S_T(\mathcal{F})$, and the unique minimax optimal*
 156 *strategy is the normalized maximum likelihood (NML) distribution given by*

$$p_{nml}(y_{1:T}) = \frac{\sup_{f \in \mathcal{F}} P_f(y_{1:T})}{\sum_{y'_{1:T} \in \mathcal{Y}^T} \sup_{f \in \mathcal{F}} P_f(y'_{1:T})}, \quad \forall y_{1:T} \in \mathcal{Y}^T.$$

157 To go beyond this classical context-agnostic setting and incorporate contextual information, prior
 158 work (e.g. [JSS21]) also considered an easier problem than the aforementioned sequential probabili-
 159 ty assignment, by forcing nature to reveal the context sequence $x_{1:T}$ to the learner at the start of the
 160 game. This is known as the *fixed design* setting or *transductive online learning* [WHGS23], where
 161 the goal is to characterize the so-called *fixed design maximal* minimax regret

$$\mathcal{R}_T^{\text{FD}}(\mathcal{F}) := \sup_{x_{1:T} \in \mathcal{X}^T} \inf_{\hat{p}_1} \sup_{y_1} \dots \inf_{\hat{p}_T} \sup_{y_T} \mathcal{R}_T(\mathcal{F}; \hat{p}_{1:T}, x_{1:T}, y_{1:T}).$$

162 It is straightforward to see that after projecting on $x_{1:T}$, the hypothesis class \mathcal{F} again collapses to a
 163 set of joint distributions over \mathcal{Y}^T specified by the likelihood function in Definition 2.2. Moreover,
 164 this set of distributions can be accessed by the learner from the start, so the fixed design setting can
 165 be essentially reduced to the context-free setting. To be more specific, for any $f \in \mathcal{F}$, it induces an
 166 expert in the context-free setting after being projected on $x_{1:T}$, which is denoted by $f|_{x_{1:T}}$ and

$$f|_{x_{1:T}}(y_{1:t-1}) := f(x_{1:t}, y_{1:t-1}) \in \Delta(\mathcal{Y}), \forall t \in [T], y_{1:t-1} \in \mathcal{Y}^{t-1},$$

167 and let $\mathcal{F}|_{x_{1:T}} := \{f|_{x_{1:T}} : f \in \mathcal{F}\}$. Then given any predetermined $x_{1:T}$, the learner is equivalently
 168 competing with $\mathcal{F}|_{x_{1:T}}$ in the context-free setting. With the following natural variant of the Shtarkov
 169 sum, we can easily characterize $\mathcal{R}_T^{\text{FD}}(\mathcal{F})$.

170 **Definition 2.4 (Conditional Shtarkov sum)** Given a context sequence $x_{1:T} \in \mathcal{X}^T$, the Shtarkov
 171 sum of \mathcal{F} conditioned on $x_{1:T}$ is

$$S_T(\mathcal{F}|x_{1:T}) := \sum_{y_{1:T} \in \mathcal{Y}^T} \sup_{f \in \mathcal{F}} P_f(y_{1:T}|x_{1:T}).$$

172 In fact, $S_T(\mathcal{F}|x_{1:T})$ is just the Shtarkov sum of the projected class $\mathcal{F}|_{x_{1:T}}$ in the context-free setting.
 173 The following result characterizes the fixed-design setting:

174 **Proposition 2.5 (Minimax regret, fixed design [JSS21])** *In the fixed design setting, for any hy-*
 175 *pothesis class \mathcal{F} and horizon T , the fixed design maximal minimax regret is*

$$\mathcal{R}_T^{\text{FD}}(\mathcal{F}) = \sup_{x_{1:T} \in \mathcal{X}^T} \log S_T(\mathcal{F}|x_{1:T}),$$

176 *and, given any context sequence $x_{1:T}$, the minimax optimal response is NML with respect to $\mathcal{F}|_{x_{1:T}}$.*

177 3 Minimax regret via contextual Shtarkov sum

178 Now we state one of our main results about the characterization of the minimax regret of sequential
 179 probability assignment. First we introduce the key concept of *contextual Shtarkov sum*, which is a
 180 natural generalization of Shtarkov sum in the context-free setting.

181 **Definition 3.1 (Contextual Shtarkov sum)** The *contextual Shtarkov sum* $S_T(\mathcal{F}|\mathbf{x})$ of a hypothesis
 182 class \mathcal{F} on a given context tree \mathbf{x} of depth T is defined as

$$S_T(\mathcal{F}|\mathbf{x}) := \sum_{\mathbf{y} \in \mathcal{Y}^T} \sup_{f \in \mathcal{F}} P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})).$$

183 Just like the conditional Shtarkov sum, the contextual Shtarkov sum $S_T(\mathcal{F}|\mathbf{x})$ can be interpreted as
 184 the Shtarkov sum of the projected class $\mathcal{F}|_{\mathbf{x}} := \{f|_{\mathbf{x}} : f \in \mathcal{F}\}$ where $f|_{\mathbf{x}}$ is the induced context-free
 185 expert specified by

$$f|_{\mathbf{x}}(\mathbf{y}) = P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})), \forall \mathbf{y} \in \mathcal{Y}^T,$$

186 for any depth- T context tree \mathbf{x} . Next we show that the minimax regret $\mathcal{R}_T(\mathcal{F})$ is characterized by
 187 the worst-case contextual Shtarkov sum:

188 **Theorem 3.2 (Main result: minimax regret)** *For any hypothesis class $\mathcal{F} \subseteq \{(\mathcal{X} \times \mathcal{Y})^* \times \mathcal{X} \rightarrow$*
 189 *$\Delta(\mathcal{Y})\}$ and horizon T ,*

$$\mathcal{R}_T(\mathcal{F}) = \sup_{\mathbf{x}} \log S_T(\mathcal{F}|\mathbf{x}),$$

190 *where the supremum is taken over all context trees \mathbf{x} (i.e., \mathbf{x} is \mathcal{X} -valued) of depth T .*

191 Since any context sequence $x_{1:T}$ can be thought as a special context tree \mathbf{x} that is constant in
 192 each level $t \in [T]$ (i.e., $\mathbf{x}_t(\mathbf{y}) = x_t, \forall \mathbf{y}$), we can find that the supremum over context trees in
 193 Theorem 3.2 strictly subsumes the supremum over context sequences in Proposition 2.5. Thus we
 194 can see the separation between $\mathcal{R}_T(\mathcal{F})$ and $\mathcal{R}_T^{\text{FD}}(\mathcal{F})$ is clearly exhibited.

195 The proof of Theorem 3.2 is provided in Appendix A but we give a brief sketch of it here.

196 **Proof sketch.** The proof starts from swapping the pairs of inf and sup (after randomizing the labels
 197 revealed by the nature) in the extensive formulation of $\mathcal{R}_T(\mathcal{F})$ to move to the *dual game*, where the
 198 learner predicts *after* seeing the action of the nature. Trivially the value of this swapped game is a
 199 lower bound for $\mathcal{R}_T(\mathcal{F})$, and after rearranging we get that

$$\text{the value of the swapped game} = \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}}[\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})] \leq \mathcal{R}_T(\mathcal{F}),$$

200 where the supremum is taken over all context trees \mathbf{x} and probabilistic trees \mathbf{p} , of depth T . Also
 201 $\mathbb{E}_{\mathbf{y} \sim \mathbf{p}}$ means the nested conditional expectations $\mathbb{E}_{y_1 \sim \mathbf{p}_1(\mathbf{y})} \mathbb{E}_{y_2 \sim \mathbf{p}_2(\mathbf{y})} \cdots \mathbb{E}_{y_T \sim \mathbf{p}_T(\mathbf{y})}$.

202 Similar to the proof of Lemma 6 in [BFR20] for the binary label setting, we apply the minimax
 203 theorem with a tweak that we devise to handle multiary labels to derive that

$$\mathcal{R}_T(\mathcal{F}) = \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}}[\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})] \quad (1)$$

204 under some mild regularity condition for \mathcal{F} . A key observation is that the supremum over depth- T
 205 probabilistic trees \mathbf{p} is equivalent to the supremum over joint distributions P over \mathcal{Y}^T . Based on this
 206 observation and a few algebraic manipulations, we can re-write $\sup_{\mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}}[\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})]$
 207 as

$$\sup_{P \in \Delta(\mathcal{Y}^T)} H(P) + \mathbb{E}_{\mathbf{y} \sim P} \left[\sup_{f \in \mathcal{F}} \log P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})) \right]$$

208 given any context tree \mathbf{x} , and the value of this maximization problem can be easily computed to be
 209 $\log \left(\sum_{\mathbf{y}} \sup_{f \in \mathcal{F}} P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})) \right) = \log S_T(\mathcal{F}|\mathbf{x})$. Thus,

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}) &= \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}}[\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})] \\ &= \sup_{\mathbf{x}} \sup_{P \in \Delta(\mathcal{Y}^T)} H(P) + \mathbb{E}_{\mathbf{y} \sim P} \left[\sup_{f \in \mathcal{F}} \log P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})) \right] = \sup_{\mathbf{x}} \log S_T(\mathcal{F}|\mathbf{x}). \end{aligned}$$

210 However, Eq. (1) is not guaranteed when there is no assumed regularity condition for \mathcal{F} . To get
 211 away from this, prior works would have to add a particular hypothesis to the class \mathcal{F} such that the
 212 enlarged class allows for the minimax swap [RS15; BFR20]. Nevertheless, even adding a mere
 213 hypothesis may lead to suboptimal analysis for some classes \mathcal{F} , say when $\mathcal{R}_T(\mathcal{F})$ is of constant
 214 order. To completely get rid of any regularity assumption and obtain a unified characterization of
 215 the minimax regret for arbitrary class \mathcal{F} , we provide a novel argument as follows. For an arbitrary
 216 class \mathcal{F} , we study a smooth truncated version of it, denoted by \mathcal{F}^δ for any level $\delta \in (0, 1/2)$, such
 217 that \mathcal{F}^δ always validates the use of the minimax theorem and hence $\mathcal{R}_T(\mathcal{F}^\delta) = \sup_{\mathbf{x}} \log S_T(\mathcal{F}^\delta|\mathbf{x})$.
 218 Then we give a series of refined analysis comparing the minimax regrets and contextual Shtarkov
 219 sums of \mathcal{F} and \mathcal{F}^δ that yields

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}) &\leq \mathcal{R}_T(\mathcal{F}^\delta) + T \log(1 + |\mathcal{Y}|\delta) = \sup_{\mathbf{x}} \log S_T(\mathcal{F}^\delta|\mathbf{x}) + T \log(1 + |\mathcal{Y}|\delta) \\ &\leq \log \left(\sup_{\mathbf{x}} S_T(\mathcal{F}|\mathbf{x}) + \delta \cdot C(T, |\mathcal{Y}|) \right) + T \log(1 + |\mathcal{Y}|\delta), \end{aligned}$$

220 where $C(T, |\mathcal{Y}|) < \infty$ is a positive constant that only depends on T and $|\mathcal{Y}|$. Sending $\delta \rightarrow 0^+$
 221 will conclude that $\mathcal{R}_T(\mathcal{F}) \leq \sup_{\mathbf{x}} \log S_T(\mathcal{F}|\mathbf{x})$, which finishes the whole proof as we already have
 222 $\mathcal{R}_T(\mathcal{F}) \geq \sup_{\mathbf{x}} \log S_T(\mathcal{F}|\mathbf{x})$ from the start.

223 3.1 Applications: an improved regret upper bound in terms of sequential entropy

224 To illustrate the utility of our characterization in Theorem 3.2, we walk through some examples
 225 where we are able to recover and *sharpen* existing regret upper bounds with relatively short proofs
 226 via contextual Shtarkov sum. As a start, we provide a short proof in Appendix A.5 of the classical
 227 regret bound for a finite hypothesis class.

228 **Proposition 3.3 (Finite classes)** For any finite hypothesis class \mathcal{F} and horizon T , $\mathcal{R}_T(\mathcal{F}) \leq$
 229 $\log |\mathcal{F}|$.

230 Let us go back to the binary label setting with non-sequential experts, that is, $\mathcal{Y} = \{0, 1\}$ and
 231 $\mathcal{F} \subseteq [0, 1]^{\mathcal{X}}$, and $f(x) \in [0, 1]$ is interpreted as the probability assigned to label 1 by this expert f .
 232 We will show a regret bound that *outperforms* the state-of-the-art ones in [BFR20; WHGS23] with
 233 a surprisingly simple proof. To proceed, we need the following notation. Given a context tree \mathbf{x} of
 234 depth T , let $\mathcal{F} \circ \mathbf{x} = \{f \circ \mathbf{x} : f \in \mathcal{F}\}$, where $f \circ \mathbf{x}$ is the $[0, 1]$ -valued tree such that

$$(f \circ \mathbf{x})_t(\mathbf{y}) = f(\mathbf{x}_t(\mathbf{y})), \forall \mathbf{y} \in \mathcal{Y}^T.$$

235 Next we introduce the definitions of sequential ℓ_∞ covers and entropy.

236 **Definition 3.4 (Sequential ℓ_∞ cover and entropy)** Given a hypothesis class $\mathcal{F} \subseteq [0, 1]^{\mathcal{X}}$ and a
 237 context tree \mathbf{x} of depth T , we say a collection of \mathbb{R} -valued trees $V_{\mathbf{x}, \alpha}$ is a sequential cover of $\mathcal{F} \circ \mathbf{x}$
 238 at scale $\alpha > 0$ if for any $f \in \mathcal{F}$, $\mathbf{y} \in \mathcal{Y}^T$, there exists some $v \in V_{\mathbf{x}, \alpha}$ such that

$$|f(\mathbf{x}_t(\mathbf{y})) - v_t(\mathbf{y})| \leq \alpha, \forall t \in [T].$$

239 Let the sequential ℓ_∞ covering number $\mathcal{N}_\infty(\mathcal{F} \circ \mathbf{x}, \alpha, T)$ be the size of the smallest such cover. The
 240 sequential ℓ_∞ entropy of \mathcal{F} at scale α and depth T is defined as the logarithm of the worst-case
 241 sequential covering number:

$$\mathcal{H}_\infty(\mathcal{F}, \alpha, T) := \sup_{\mathbf{x}} \log \mathcal{N}_\infty(\mathcal{F} \circ \mathbf{x}, \alpha, T).$$

242

243 **Definition 3.5 (Global sequential ℓ_∞ cover and entropy)** Given a hypothesis class $\mathcal{F} \subseteq [0, 1]^\mathcal{X}$,
 244 we say a collection of mappings $\mathcal{G}_\alpha \subseteq [0, 1]^{\mathcal{X}^*}$ is a *global sequential cover* of \mathcal{F} at scale $\alpha > 0$ and
 245 depth T if for any $f \in \mathcal{F}, x_{1:T} \in \mathcal{X}^T$, there exists some $g \in \mathcal{G}_\alpha$ such that

$$|f(x_t) - g(x_{1:t})| \leq \alpha, \forall t \in [T].$$

246 Let the *global sequential ℓ_∞ covering number* $\mathcal{N}_G(\mathcal{F}, \alpha, T)$ be the size of the smallest such cover.
 247 The *global sequential ℓ_∞ entropy* of \mathcal{F} at scale α and depth T is defined as

$$\mathcal{H}_G(\mathcal{F}, \alpha, T) := \log \mathcal{N}_G(\mathcal{F}, \alpha, T).$$

248

249 **Proposition 3.6 ([BFR20; WHGS23])** For any $\mathcal{F} \subseteq [0, 1]^\mathcal{X}$ and horizon T ,

$$\mathcal{R}_T(\mathcal{F}) \leq \min \left\{ \underbrace{\inf_{\alpha > 0} \{4T\alpha + c\mathcal{H}_\infty(\mathcal{F}, \alpha, T)\}}_{\text{[BFR20]}}, \underbrace{\inf_{\alpha > 0} \{T \log(1 + 2\alpha) + \mathcal{H}_G(\mathcal{F}, \alpha, T)\}}_{\text{[WHGS23]}} \right\},$$

250 where $c = \frac{2 - \log(2)}{\log(3) - \log(2)} \in (3, 4)$.

251 It is easy to show that $\mathcal{H}_\infty(\mathcal{F}, \alpha, T) \leq \mathcal{H}_G(\mathcal{F}, \alpha, T)$, but, in general, the two bounds in Propo-
 252 sition 3.6 are incomparable due to constants and different dependence on α (more discussions on
 253 these bounds are deferred to Appendix C). Starting from the contextual Shtarkov sum, we are able
 254 to derive a bound that combines the best of these two bounds :

255 **Theorem 3.7 (Main result: sequential entropy bound)** For any $\mathcal{F} \subseteq [0, 1]^\mathcal{X}$ and horizon T ,

$$\mathcal{R}_T(\mathcal{F}) \leq \inf_{\alpha > 0} \left\{ T \log(1 + 2\alpha) + \mathcal{H}_\infty(\mathcal{F}, \alpha, T) \right\}.$$

256

257 **Proof of Theorem 3.7** For any scale $\alpha > 0$ and depth- T context tree \mathbf{x} , let $V_{\mathbf{x}, \alpha}$ be a sequential
 258 cover of $\mathcal{F} \circ \mathbf{x}$ at scale α with size $\mathcal{N}_\infty(\mathcal{F} \circ \mathbf{x}, \alpha, T)$. We can always assume $V_{\mathbf{x}, \alpha}$ to be $[0, 1]$ -valued
 259 without loss of generality because otherwise we can just truncate it without violating its coverage
 260 guarantee. Define the smoothed covering set $\tilde{V}_{\mathbf{x}, \alpha} = \left\{ \tilde{v} : \forall t \in [T], \tilde{v}_t(\cdot) = \frac{v_t(\cdot) + \alpha}{1 + 2\alpha}, v \in V_{\mathbf{x}, \alpha} \right\}$,
 261 inspired by [BFR23; WHGS23]. Then for any $f \in \mathcal{F}, \mathbf{y} \in \mathcal{Y}^T$, there exists some $v \in V_{\mathbf{x}, \alpha}$ such
 262 that $|f(\mathbf{x}_t(\mathbf{y})) - v_t(\mathbf{y})| \leq \alpha, \forall t \in [T]$ and hence \tilde{v} satisfies

$$\frac{f(\mathbf{x}_t(\mathbf{y}))}{\tilde{v}_t(\mathbf{y})} \leq 1 + 2\alpha, \quad \frac{1 - f(\mathbf{x}_t(\mathbf{y}))}{1 - \tilde{v}_t(\mathbf{y})} \leq 1 + 2\alpha.$$

263 Hence

$$P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})) = \prod_{t=1}^T f(\mathbf{x}_t(\mathbf{y}))^{y_t} (1 - f(\mathbf{x}_t(\mathbf{y})))^{1-y_t} \leq (1 + 2\alpha)^T \prod_{t=1}^T \tilde{v}_t(\mathbf{y})^{y_t} (1 - \tilde{v}_t(\mathbf{y}))^{1-y_t},$$

264 and

$$\begin{aligned} \sum_{\mathbf{y}} \sup_{f \in \mathcal{F}} P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})) &\leq (1 + 2\alpha)^T \sum_{\mathbf{y}} \sup_{\tilde{v} \in \tilde{V}_{\mathbf{x}, \alpha}} \prod_{t=1}^T \tilde{v}_t(\mathbf{y})^{y_t} (1 - \tilde{v}_t(\mathbf{y}))^{1-y_t} \\ &\leq (1 + 2\alpha)^T \sum_{\tilde{v} \in \tilde{V}_{\mathbf{x}, \alpha}} \sum_{\mathbf{y}} \prod_{t=1}^T \tilde{v}_t(\mathbf{y})^{y_t} (1 - \tilde{v}_t(\mathbf{y}))^{1-y_t} = (1 + 2\alpha)^T |\tilde{V}_{\mathbf{x}, \alpha}|, \end{aligned}$$

265 where the last equality is derived by treating \tilde{v} as sequential experts and applying Lemma D.1.
 266 Finally,

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}) &= \sup_{\mathbf{x}} \log \left(\sum_{\mathbf{y}} \sup_{f \in \mathcal{F}} P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})) \right) \\ &\leq \sup_{\mathbf{x}} \log \left((1 + 2\alpha)^T |\tilde{V}_{\mathbf{x}, \alpha}| \right) \\ &= \sup_{\mathbf{x}} \log \left((1 + 2\alpha)^T |V_{\mathbf{x}, \alpha}| \right) = T \log(1 + 2\alpha) + \mathcal{H}_\infty(\mathcal{F}, \alpha, T). \end{aligned}$$

267 Since our choice of α is arbitrary, we conclude that

$$\mathcal{R}_T(\mathcal{F}) \leq \inf_{\alpha > 0} \left\{ T \log(1 + 2\alpha) + \mathcal{H}_\infty(\mathcal{F}, \alpha, T) \right\}.$$

268

■

269 3.2 Contextual Shtarkov sums through martingales

270 We can relate our characterization of the minimax regret to the more extensively studied *sequential*
 271 *Rademacher complexity*, which arises in online learning problems with hypothesis class $\mathcal{F} \subseteq [0, 1]^{\mathcal{X}}$
 272 and bounded convex losses like absolute loss. Specifically, the (conditional) sequential Rademacher
 273 complexity [RST15] is defined by

$$\mathfrak{R}_T(\mathcal{F}; \mathbf{x}) := \mathbb{E}_\varepsilon \left[\sup_{f \in \mathcal{F}} \sum_{t=1}^T \varepsilon_t f(\mathbf{x}_t(\varepsilon)) \right],$$

274 where \mathbf{x} is a depth- T binary context tree and $\varepsilon = (\varepsilon_1, \dots, \varepsilon_T) \in \{\pm 1\}^T$ is a sequence of i.i.d.
 275 Rademacher random variables. A notable feature of $\mathfrak{R}_T(\mathcal{F}; \mathbf{x})$ is that it is the expected supremum
 276 of the sum of a martingale differences, i.e., for any f , $\mathbb{E}[\varepsilon_t f(\mathbf{x}_t(\varepsilon)) | \varepsilon_1, \dots, \varepsilon_{t-1}] = 0$. Likewise,
 277 $S_T(\mathcal{F} | \mathbf{x})$ also admits a martingale interpretation. To see this, let $\mathcal{F} \subseteq \{(\mathcal{X} \times \mathcal{Y})^* \times \mathcal{X} \rightarrow \Delta(\mathcal{Y})\}$
 278 and rewrite $S_T(\mathcal{F} | \mathbf{x})$ for any context tree \mathbf{x} :

$$S_T(\mathcal{F} | \mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}^T} \sup_{f \in \mathcal{F}} P_f(\mathbf{y} | \mathbf{x}(\mathbf{y})) = \mathbb{E}_{\mathbf{y}} \left[\sup_{f \in \mathcal{F}} \prod_{t=1}^T (|\mathcal{Y}| \cdot f(\mathbf{x}_{1:t}(\mathbf{y}), y_{1:t-1})(y_t)) \right],$$

279 where $\mathbf{y} = (y_1, \dots, y_T)$ is a sequence of i.i.d. variables following the uniform distribution over \mathcal{Y} .
 280 It is easy to check that $\mathbb{E}[|\mathcal{Y}| \cdot f(\mathbf{x}_{1:t}(\mathbf{y}), y_{1:t-1})(y_t) | y_1, \dots, y_{t-1}] = 1$, and thus

$$\left\{ \prod_{s=1}^t (|\mathcal{Y}| \cdot f(\mathbf{x}_{1:s}(\mathbf{y}), y_{1:s-1})(y_s)) \right\}_{t \in [T]}$$

281 is a martingale with respect to filtration $\mathcal{F}_t = \sigma(y_1, \dots, y_t), t \in [T]$. It would be of independent
 282 interest to study the contextual Shtarkov sums more quantitatively by developing new tools for such
 283 product-type martingales.

284 4 Contextual NML, the minimax optimal algorithm

285 So far we have settled the minimax regret of sequential probability assignment in a nonconstruc-
 286 tive way. Now we switch to the algorithmic lens to study the optimal strategy that achieves the
 287 minimax regret. Remarkably, we show that the minimax optimal algorithm can be described by a
 288 data-dependent variant of the contextual Shtarkov sum, which is named contextual Shtarkov sum
 289 *with prefix*.

290 **Definition 4.1 (Contextual Shtarkov sum with prefix)** Given sequences $x_{1:t} \in \mathcal{X}^t, y_{1:t} \in \mathcal{Y}^t, t \in$
 291 $[T]$ and a context tree \mathbf{x} of depth $T - t$, the contextual Shtarkov sum $S_T^{x_{1:t}, y_{1:t}}(\mathcal{F} | \mathbf{x})$ of \mathcal{F} on \mathbf{x} with
 292 prefix $x_{1:t}, y_{1:t}$ is defined as

$$S_T^{x_{1:t}, y_{1:t}}(\mathcal{F} | \mathbf{x}) = \sum_{\mathbf{y} \in \mathcal{Y}^{T-t}} \sup_{f \in \mathcal{F}} P_f(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y})).$$

293 Now we present our prediction strategy, *contextual normalized maximum likelihood* (cNML),
 294 which is summarized in Algorithm 1. In each round t , with $x_{1:t}, y_{1:t-1}$ as past observations,
 295 the learner first checks whether $\sup_{f \in \mathcal{F}} P_f(y_{1:t-1} | x_{1:t-1}) > 0$ since if that is not the case and
 296 $\sup_{f \in \mathcal{F}} P_f(y_{1:t-1} | x_{1:t-1}) = 0$, the cumulative losses of all experts in \mathcal{F} have already blown up to
 297 $+\infty$ and the learner only needs to predict any $\hat{p} \in \Delta^+(\mathcal{Y})$ in all remaining rounds. On the other
 298 hand, if $\sup_{f \in \mathcal{F}} P_f(y_{1:t-1} | x_{1:t-1}) > 0$, then

$$\max_{y \in \mathcal{Y}} \sup_{\mathbf{x}} S_f^{x_{1:t}, (y_{1:t-1}, y)}(\mathcal{F} | \mathbf{x}) > 0$$

Algorithm 1 Contextual Normalized Maximum Likelihood (cNML)

Input: Hypothesis class \mathcal{F} , horizon T

For $t = 1, 2, \dots, T$ **do**

1. Observe context $x_t \in \mathcal{X}$
2. If $\sup_{f \in \mathcal{F}} P_f(y_{1:t-1} | x_{1:t-1}) > 0$, predict $\hat{p}_t \in \Delta(\mathcal{Y})$ with

$$\hat{p}_t(y) = \frac{\sup_{\mathbf{x}} S_T^{x_{1:t}, (y_{1:t-1}, y)}(\mathcal{F} | \mathbf{x})}{\sum_{y' \in \mathcal{Y}} \sup_{\mathbf{x}} S_T^{x_{1:t}, (y_{1:t-1}, y')}(\mathcal{F} | \mathbf{x})}, \forall y \in \mathcal{Y}, \quad (2)$$

and otherwise set \hat{p}_t to be an arbitrary member of $\Delta^+(\mathcal{Y})$

3. Receive label $y_t \in \mathcal{Y}$

End for

299 and the \hat{p}_t given by Eq. (2) is indeed a valid member of $\Delta(\mathcal{Y})$ (shown in Appendix B) and is used as
300 the learner's prediction. The following theorem shows that cNML is the minimax optimal algorithm,
301 with proof deferred to Appendix B.

302 **Theorem 4.2 (Main result: optimal algorithm)** *The contextual normalized maximum likelihood*
303 *strategy (Algorithm 1) is minimax optimal.*

304 To see that cNML is reduced to NML in the context-free setting, it suffices to consider the case
305 where $\sup_{f \in \mathcal{F}} P_f(y_{1:T}) > 0$ since otherwise NML will simply assign 0 probability on this se-
306 quence $y_{1:T}$ and during the actual round-wise implementation of NML, it also predicts an arbitrary
307 element from $\Delta^+(\mathcal{Y})$ in those rounds t where $\sup_{f \in \mathcal{F}} P_f(y_{1:t-1}) = 0$. Now for any $y_{1:T}$ such that
308 $\sup_{f \in \mathcal{F}} P_f(y_{1:T}) > 0$, the prediction by cNML in each round t is

$$\hat{p}_t(y) = \frac{\sum_{\mathbf{y} \in \mathcal{Y}^{T-t}} \sup_{f \in \mathcal{F}} P_f(y_{1:t-1}, y, \mathbf{y})}{\sum_{\mathbf{y}' \in \mathcal{Y}^{T-t+1}} \sup_{f \in \mathcal{F}} P_f(y_{1:t-1}, \mathbf{y}')}, \forall y \in \mathcal{Y}$$

309 which can be summarized into a joint density over $y_{1:T}$ by

$$\hat{p}(y_{1:T}) = \frac{\sup_{f \in \mathcal{F}} P_f(y_{1:T})}{\sum_{y'_{1:T} \in \mathcal{Y}^T} \sup_{f \in \mathcal{F}} P_f(y'_{1:T})}.$$

310 Recall that this is exactly the NML prediction $p_{nml}(y_{1:T})$.

311 5 Discussions

312 In this paper, we characterize the minimax regret and the optimal prediction strategy for sequen-
313 tial probability assignment, generalizing the classical results in the context-free setting. Moreover,
314 our results are general enough to subsume the setting of multiary labels and sequential hypothesis
315 classes, which has not been sufficiently explored before. Remarkably, our characterization holds
316 for arbitrary hypothesis classes that may not admit the regularity assumptions implicitly required by
317 prior works (e.g. [RST15; BFR20]).

318 For future works, it would be interesting to study the minimax regret of specific classes more quan-
319 titatively using our contextual Shtarkov sums. It is also intriguing to consider the setting of infinite
320 labels. Although most of our arguments would go through under sufficient regularity conditions, a
321 more systematic study is needed. On the practical side, it is important to develop algorithms that are
322 more computationally efficient than cNML and with provable guarantees.

323 References

324 [ABDMNY21] N. Alon, O. Ben-Eliezer, Y. Dagan, S. Moran, M. Naor, and E. Yogev. "Adversarial
325 laws of large numbers and optimal regret in online classification". In: *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*. 2021,
326 pp. 447–455 (cit. on p. 3).
327

- 328 [AHKKV23] I. Attias, S. Hanneke, A. Kalavasis, A. Karbasi, and G. Velegkas. “Optimal learners for realizable regression: Pac learning and online learning”. In: *Advances in Neural Information Processing Systems* 36 (2023) (cit. on p. 3).
- 329
- 330
- 331 [BRY98] A. Barron, J. Rissanen, and B. Yu. “The minimum description length principle in coding and modeling”. In: *IEEE transactions on information theory* 44.6 (1998), pp. 2743–2760 (cit. on p. 2).
- 332
- 333
- 334 [BPS09] S. Ben-David, D. Pál, and S. Shalev-Shwartz. “Agnostic Online Learning.” In: *COLT*. Vol. 3. 2009, p. 1 (cit. on p. 3).
- 335
- 336 [BFR20] B. Bilodeau, D. Foster, and D. Roy. “Tight bounds on minimax regret under logarithmic loss via self-concordance”. In: *International Conference on Machine Learning*. PMLR. 2020, pp. 919–929 (cit. on pp. 2, 3, 5–7, 9, 12, 22).
- 337
- 338
- 339 [BFR23] B. Bilodeau, D. J. Foster, and D. M. Roy. “Minimax rates for conditional density estimation via empirical entropy”. In: *The Annals of Statistics* 51.2 (2023), pp. 762–790 (cit. on pp. 7, 12).
- 340
- 341
- 342 [CL06] N. Cesa-Bianchi and G. Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006 (cit. on pp. 1, 2).
- 343
- 344 [CL99] N. Cesa-Bianchi and G. Lugosi. “Minimax regret under log loss for general classes of experts”. In: *Proceedings of the Twelfth annual conference on computational learning theory*. 1999, pp. 12–18 (cit. on p. 2).
- 345
- 346
- 347 [Cov74] T. M. Cover. “Universal gambling schemes and the complexity measures of Kolmogorov and Chaitin”. In: *Technical Report, no. 12* (1974) (cit. on pp. 1, 2).
- 348
- 349 [Cov91] T. M. Cover. “Universal portfolios”. In: *Mathematical finance* 1.1 (1991), pp. 1–29 (cit. on pp. 1, 2).
- 350
- 351 [CO96] T. M. Cover and E. Ordentlich. “Universal portfolios with side information”. In: *IEEE Transactions on Information Theory* 42.2 (1996), pp. 348–363 (cit. on pp. 1, 2).
- 352
- 353
- 354 [Dav73] L. D. Davisson. “Universal noiseless coding”. In: *IEEE Trans. Inf. Theory* 19 (1973), pp. 783–795 (cit. on p. 2).
- 355
- 356 [DS04] M. Drmota and W. Szpankowski. “Precise minimax redundancy and regret”. In: *IEEE Transactions on Information Theory* 50.11 (2004), pp. 2686–2707 (cit. on p. 2).
- 357
- 358
- 359 [Fed91] M. Feder. “Gambling using a finite state machine”. In: *IEEE Transactions on Information Theory* 37.5 (1991), pp. 1459–1465 (cit. on pp. 1, 2).
- 360
- 361 [FMG92] M. Feder, N. Merhav, and M. Gutman. “Universal prediction of individual sequences”. In: *IEEE transactions on Information Theory* 38.4 (1992), pp. 1258–1270 (cit. on p. 2).
- 362
- 363
- 364 [Fit66] B. Fitingof. “Optimal encoding with unknown and variable message statistics”. In: *Probl. Inform. Transm.* 2 (1966), pp. 3–11 (cit. on p. 2).
- 365
- 366 [FKLMS18] D. J. Foster, S. Kale, H. Luo, M. Mohri, and K. Sridharan. “Logistic regression: The importance of being improper”. In: *Conference on learning theory*. PMLR. 2018, pp. 167–208 (cit. on p. 1).
- 367
- 368
- 369 [Fre96] Y. Freund. “Predicting a binary sequence almost as well as the optimal biased coin”. In: *Proceedings of the ninth annual conference on Computational learning theory*. 1996, pp. 89–98 (cit. on p. 2).
- 370
- 371
- 372 [Grü05] P. Grünwald. “Minimum description length tutorial”. In: (2005) (cit. on p. 2).
- 373
- 374 [HY01] M. H. Hansen and B. Yu. “Model selection and the principle of minimum description length”. In: *Journal of the American Statistical Association* 96.454 (2001), pp. 746–774 (cit. on p. 2).
- 375
- 376 [JSS21] P. Jacquet, G. Shamir, and W. Szpankowski. “Precise minimax regret for logistic regression with categorical feature values”. In: *Algorithmic Learning Theory*. PMLR. 2021, pp. 755–771 (cit. on pp. 2, 4, 5).
- 377
- 378
- 379 [Kel56] J. L. Kelly. “A new interpretation of information rate”. In: *the bell system technical journal* 35.4 (1956), pp. 917–926 (cit. on pp. 1, 2).
- 380
- 381 [Kol65] A. N. Kolmogorov. “Three approaches to the definition of the concept “quantity of information””. In: *Problemy peredachi informatsii* 1.1 (1965), pp. 3–11 (cit. on p. 2).
- 382
- 383

- 384 [Lit88] N. Littlestone. “Learning quickly when irrelevant attributes abound: A new linear-
385 threshold algorithm”. In: *Machine learning* 2 (1988), pp. 285–318 (cit. on p. 3).
- 386 [MF98] N. Merhav and M. Feder. “Universal prediction”. In: *IEEE Transactions on Infor-
387 mation Theory* 44.6 (1998), pp. 2124–2147 (cit. on pp. 1, 2).
- 388 [MG22] J. Mourtada and S. Gaïffas. “An improper estimator with optimal excess risk in
389 misspecified density estimation and logistic regression”. In: *Journal of Machine
390 Learning Research* 23.31 (2022), pp. 1–49 (cit. on p. 19).
- 391 [OH99] M. Opper and D. Haussler. “Worst case prediction over sequences under log
392 loss”. In: *The Mathematics of Information Coding, Extraction and Distribution*.
393 Springer, 1999, pp. 81–90 (cit. on p. 2).
- 394 [OS04] A. Orlitsky and N. P. Santhanam. “Speaking of infinity [iid strings]”. In: *IEEE
395 Transactions on Information Theory* 50.10 (2004), pp. 2215–2230 (cit. on p. 2).
- 396 [RS14a] A. Rakhlin and K. Sridharan. “Online non-parametric regression”. In: *Conference
397 on Learning Theory*. PMLR, 2014, pp. 1232–1264 (cit. on pp. 2, 3).
- 398 [RS14b] A. Rakhlin and K. Sridharan. “Statistical Learning and Sequential Prediction”. In:
399 2014 (cit. on p. 3).
- 400 [RS15] A. Rakhlin and K. Sridharan. “Sequential probability assignment with binary al-
401 phabets and large classes of experts”. In: *arXiv preprint arXiv:1501.07340* (2015)
402 (cit. on pp. 2, 3, 6, 22).
- 403 [RST10] A. Rakhlin, K. Sridharan, and A. Tewari. “Online learning: Random averages,
404 combinatorial parameters, and learnability”. In: *Advances in Neural Information
405 Processing Systems* 23 (2010) (cit. on pp. 2, 3).
- 406 [RST15] A. Rakhlin, K. Sridharan, and A. Tewari. “Sequential complexities and uniform
407 martingale laws of large numbers”. In: *Probability theory and related fields* 161
408 (2015), pp. 111–153 (cit. on pp. 1–3, 8, 9, 12).
- 409 [Ris78] J. Rissanen. “Modeling by shortest data description”. In: *Automatica* 14.5 (1978),
410 pp. 465–471 (cit. on p. 2).
- 411 [Ris84] J. Rissanen. “Universal coding, information, prediction, and estimation”. In: *IEEE
412 Transactions on Information theory* 30.4 (1984), pp. 629–636 (cit. on pp. 1, 2).
- 413 [Ris86] J. Rissanen. “Complexity of strings in the class of Markov sources”. In: *IEEE
414 Transactions on Information Theory* 32.4 (1986), pp. 526–532 (cit. on p. 2).
- 415 [Ris87] J. Rissanen. “Stochastic complexity”. In: *Journal of the Royal Statistical Society:
416 Series B (Methodological)* 49.3 (1987), pp. 223–239 (cit. on p. 2).
- 417 [RL81] J. Rissanen and G. Langdon. “Universal modeling and coding”. In: *IEEE Trans-
418 actions on Information Theory* 27.1 (1981), pp. 12–23 (cit. on p. 2).
- 419 [Ris76] J. J. Rissanen. “Generalized Kraft inequality and arithmetic coding”. In: *IBM Jour-
420 nal of research and development* 20.3 (1976), pp. 198–203 (cit. on p. 2).
- 421 [Ris96] J. J. Rissanen. “Fisher information and stochastic complexity”. In: *IEEE transac-
422 tions on information theory* 42.1 (1996), pp. 40–47 (cit. on p. 2).
- 423 [Sha06] G. I. Shamir. “On the MDL principle for iid sources with large alphabets”. In:
424 *IEEE transactions on information theory* 52.5 (2006), pp. 1939–1955 (cit. on p. 2).
- 425 [Sha20] G. I. Shamir. “Logistic regression regret: What’s the catch?” In: *Conference on
426 Learning Theory*. PMLR, 2020, pp. 3296–3319 (cit. on p. 1).
- 427 [Sht87] Y. M. Shtarkov. “Universal Sequential Coding of Single Messages”. In: *Problems
428 of Information Transmission* 23 (3 1987), pp. 3–17 (cit. on pp. 1, 2, 4).
- 429 [Sio58] M. Sion. “On general minimax theorems”. In: *Pacific Journal of Mathematics* 8
430 (1958), pp. 171–176 (cit. on p. 13).
- 431 [Sol64] R. Solmonoff. “A formal theory of inductive inference. I”. In: *II Information and
432 Control* 7 (1964), pp. 224–254 (cit. on p. 2).
- 433 [Szp98] W. Szpankowski. “On asymptotics of certain recurrences arising in universal cod-
434 ing”. In: *PROBLEMS OF INFORMATION TRANSMISSION C/C OF PROBLEMY
435 PEREDACHI INFORMATSII* 34 (1998), pp. 142–146 (cit. on p. 2).
- 436 [Vov95] V. G. Vovk. “A game of prediction with expert advice”. In: *Proceedings of the
437 eighth annual conference on Computational learning theory*. 1995, pp. 51–60 (cit.
438 on p. 1).

439 [WHGS23] C. Wu, M. Heidari, A. Grama, and W. Szpankowski. “Regret Bounds for Log-loss
440 via Bayesian Algorithms”. In: *IEEE Transactions on Information Theory* (2023)
441 (cit. on pp. 2–4, 6, 7, 12, 22, 23).

442 [XB97] Q. Xie and A. R. Barron. “Minimax redundancy for the class of memoryless
443 sources”. In: *IEEE Transactions on Information Theory* 43.2 (1997), pp. 646–657
444 (cit. on p. 2).

445 [XB00] Q. Xie and A. R. Barron. “Asymptotic minimax regret for data compression, gam-
446 bling, and prediction”. In: *IEEE Transactions on Information Theory* 46.2 (2000),
447 pp. 431–445 (cit. on p. 1).

448 [ZL77] J. Ziv and A. Lempel. “A universal algorithm for sequential data compression”. In:
449 *IEEE Transactions on information theory* 23.3 (1977), pp. 337–343 (cit. on p. 2).

450 [ZL78] J. Ziv and A. Lempel. “Compression of individual sequences via variable-rate cod-
451 ing”. In: *IEEE transactions on Information Theory* 24.5 (1978), pp. 530–536 (cit.
452 on p. 2).

453 A Proofs for Section 3

454 **Notations.** When the context and label sequences $x_{1:T}, y_{1:T}$ are clear from the context, we may
455 use f_t to denote the probability vector $f(x_{1:t}, y_{1:t-1}) \in \Delta(\mathcal{Y})$ produced by hypothesis f at time t
456 for notational convenience. We also adopt the notation for repeated operators in [RST15; BFR20],
457 denoting $\text{Opt}_1 \cdots \text{Opt}_T[\cdots]$ by $\left\langle \left\langle \text{Opt}_t \right\rangle \right\rangle_{t=1}^T [\cdots]$. For any discrete distribution P and discrete
458 random variables X, Y , let $H(P)$ be the entropy of P and $H(X|Y)$ be the conditional entropy of
459 X given Y .

460 A.1 Minimax swap

461 As standard in online learning literature, we will first move to a dual game after applying a minimax
462 swap at each round of the game. Under mild assumptions, the value of the original game coincides
463 with the that of the swapped game. More specifically, we have:

464 **Lemma A.1** *Whenever \mathcal{F} satisfies that for every sequence $x_{1:T} \in \mathcal{X}^T, y_{1:T} \in \mathcal{Y}^T$,*

$$\inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f(x_{1:t}, y_{1:t-1}), y_t) < \infty, \quad (3)$$

465 *we have that*

$$\mathcal{R}_T(\mathcal{F}) = \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}}[\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})], \quad (4)$$

466 *where the supremum is taken over all \mathcal{X} -valued \mathcal{Y} -ary trees \mathbf{x} and $\Delta(\mathcal{Y})$ -valued \mathcal{Y} -ary trees \mathbf{p} , of
467 depth T . Also $\mathbb{E}_{\mathbf{y} \sim \mathbf{p}}$ means the nested conditional expectations $\mathbb{E}_{y_1 \sim \mathbf{p}_1(\mathbf{y})} \mathbb{E}_{y_2 \sim \mathbf{p}_2(\mathbf{y})} \cdots \mathbb{E}_{y_T \sim \mathbf{p}_T(\mathbf{y})}$.*

468 To deal with the unboundedness of log loss in the proof, we introduce the following truncation
469 method inspired by [BFR23; WHGS23], generalizing the one in [BFR20] which was specific to
470 binary labels.

471 **Definition A.2 (Smooth truncation)** The general smooth truncation map $\tau_\delta : \Delta(\mathcal{Y}) \rightarrow \Delta(\mathcal{Y})$ is
472 defined such that for all $p \in \Delta(\mathcal{Y})$ and $y \in \mathcal{Y}$,

$$\tau_\delta(p)(y) = \frac{p(y) + \delta}{1 + |\mathcal{Y}|\delta},$$

473 given threshold $\delta \in (0, 1/2)$.

474 It is easy to check that $\tau_\delta(p)$ is indeed a valid member in $\Delta(\mathcal{Y})$ and $\tau_\delta(p)(y) \in [\delta/(1 + |\mathcal{Y}|\delta), (1 +$
475 $\delta)/(1 + |\mathcal{Y}|\delta)]$. Moreover, it is not hard to verify that $\tau_\delta(\Delta(\mathcal{Y})) = \{p \in \Delta(\mathcal{Y}) : p(y) \in [\delta/(1 +$
476 $|\mathcal{Y}|\delta), (1 + \delta)/(1 + |\mathcal{Y}|\delta)], \forall y \in \mathcal{Y}\}$. We will use $\Delta^\delta(\mathcal{Y})$ to denote this image set $\tau_\delta(\Delta(\mathcal{Y}))$.

477 **Proof of Lemma A.1** Fix $\delta \in (0, 1/2)$. By restricting the learner's prediction \hat{p}_t to $\Delta^\delta(\mathcal{Y})$, we get
 478 an upper bound on $\mathcal{R}_T(\mathcal{F})$:

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}) &\leq \left\langle \left\langle \sup_{x_t} \inf_{\hat{p}_t \in \Delta^\delta(\mathcal{Y})} \sup_{y_t} \right\rangle \right\rangle_{t=1}^T \left[\sum_{t=1}^T \ell(\hat{p}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right] \\ &= \left\langle \left\langle \sup_{x_t} \inf_{\hat{p}_t \in \Delta^\delta(\mathcal{Y})} \sup_{y_t} \right\rangle \right\rangle_{t=1}^{T-1} \sup_{x_T} \inf_{\hat{p}_T \in \Delta^\delta(\mathcal{Y})} \sup_{p_T} \mathbb{E}_{y_T \sim p_T} \left[\sum_{t=1}^T \ell(\hat{p}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right]. \end{aligned}$$

479 Now we can apply Sion's minimax theorem [Sio58] to the function

$$A(\hat{p}_T, p_T) = \mathbb{E}_{y_T \sim p_T} \left[\sum_{t=1}^T \ell(\hat{p}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right]$$

480 to derive that

$$\inf_{\hat{p}_T \in \Delta^\delta(\mathcal{Y})} \sup_{p_T \in \Delta(\mathcal{Y})} A(\hat{p}_T, p_T) = \sup_{p_T \in \Delta(\mathcal{Y})} \inf_{\hat{p}_T \in \Delta^\delta(\mathcal{Y})} A(\hat{p}_T, p_T).$$

481 This is because:

- 482 1. $A(\hat{p}_T, p_T)$ is convex and continuous in $\hat{p}_T \in \Delta^\delta(\mathcal{Y})$ and
- 483 2. $A(\hat{p}_T, p_T)$ is concave and continuous in $p_T \in \Delta(\mathcal{Y})$, which is further due to that $A(\hat{p}_T, p_T)$
 484 is linear in p_T and is bounded given Eq. (3).

485 Hence

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}) &\leq \left\langle \left\langle \sup_{x_t} \inf_{\hat{p}_t \in \Delta^\delta(\mathcal{Y})} \sup_{y_t} \right\rangle \right\rangle_{t=1}^{T-1} \sup_{x_T} \sup_{p_T} \inf_{\hat{p}_T \in \Delta^\delta(\mathcal{Y})} \mathbb{E}_{y_T \sim p_T} \left[\sum_{t=1}^T \ell(\hat{p}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right] \\ &= \left\langle \left\langle \sup_{x_t} \inf_{\hat{p}_t \in \Delta^\delta(\mathcal{Y})} \sup_{y_t} \right\rangle \right\rangle_{t=1}^{T-2} \sup_{x_{T-1}} \inf_{\hat{p}_{T-1} \in \Delta^\delta(\mathcal{Y})} \sup_{p_{T-1}} \mathbb{E}_{y_{T-1} \sim p_{T-1}} \\ &\quad \left[\sum_{t=1}^{T-1} \ell(\hat{p}_t, y_t) + \sup_{x_T} \sup_{p_T} \left[\inf_{\hat{p}_T \in \Delta^\delta(\mathcal{Y})} \mathbb{E}_{y_T \sim p_T} \ell(\hat{p}_T, y_T) - \mathbb{E}_{y_T \sim p_T} \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right] \right]. \end{aligned}$$

486 Again the order of $\inf_{\hat{p}_{T-1} \in \Delta^\delta(\mathcal{Y})}$ and $\sup_{p_{T-1} \in \Delta(\mathcal{Y})}$ with respect to

$$B(\hat{p}_{T-1}, p_{T-1}) = \mathbb{E}_{y_{T-1} \sim p_{T-1}} \left[\sum_{t=1}^{T-1} \ell(\hat{p}_t, y_t) + \sup_{x_T} \sup_{p_T} \left[\inf_{\hat{p}_T \in \Delta^\delta(\mathcal{Y})} \mathbb{E}_{y_T \sim p_T} \ell(\hat{p}_T, y_T) - \mathbb{E}_{y_T \sim p_T} \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right] \right]$$

487 can be swapped due to the same reason as above, leading to

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}) &\leq \left\langle \left\langle \sup_{x_t} \inf_{\hat{p}_t \in \Delta^\delta(\mathcal{Y})} \sup_{y_t} \right\rangle \right\rangle_{t=1}^{T-2} \sup_{x_{T-1}} \sup_{p_{T-1}} \inf_{\hat{p}_{T-1} \in \Delta^\delta(\mathcal{Y})} \mathbb{E}_{y_{T-1} \sim p_{T-1}} \\ &\quad \left[\sum_{t=1}^{T-1} \ell(\hat{p}_t, y_t) + \sup_{x_T} \sup_{p_T} \left[\inf_{\hat{p}_T \in \Delta^\delta(\mathcal{Y})} \mathbb{E}_{y_T \sim p_T} \ell(\hat{p}_T, y_T) - \mathbb{E}_{y_T \sim p_T} \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right] \right] \\ &= \left\langle \left\langle \sup_{x_t} \inf_{\hat{p}_t \in \Delta^\delta(\mathcal{Y})} \sup_{y_t} \right\rangle \right\rangle_{t=1}^{T-3} \sup_{x_{T-2}} \inf_{\hat{p}_{T-2} \in \Delta^\delta(\mathcal{Y})} \sup_{p_{T-2}} \mathbb{E}_{y_{T-2} \sim p_{T-2}} \\ &\quad \left\{ \sum_{t=1}^{T-2} \ell(\hat{p}_t, y_t) + \sup_{x_{T-1}} \sup_{p_{T-1}} \left[\inf_{\hat{p}_{T-1} \in \Delta^\delta(\mathcal{Y})} \mathbb{E}_{y_{T-1} \sim p_{T-1}} \ell(\hat{p}_{T-1}, y_{T-1}) \right. \right. \\ &\quad \left. \left. + \mathbb{E}_{y_{T-1} \sim p_{T-1}} \sup_{x_T} \sup_{p_T} \left[\inf_{\hat{p}_T \in \Delta^\delta(\mathcal{Y})} \mathbb{E}_{y_T \sim p_T} \ell(\hat{p}_T, y_T) - \mathbb{E}_{y_T \sim p_T} \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right] \right] \right\}. \end{aligned}$$

488 Repeating this procedure through all T rounds yields

$$\mathcal{R}_T(\mathcal{F}) \leq \left\langle \left\langle \sup_{x_t} \sup_{p_t} \mathbb{E}_{y_t \sim p_t} \right\rangle \right\rangle_{t=1}^T \sup_{f \in \mathcal{F}} \left[\sum_{t=1}^T \inf_{\hat{p}_t \in \Delta^\delta(\mathcal{Y})} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{p}_t, y_t)] - \ell(f_t, y_t) \right].$$

489 By Lemma A.7, we know that we do not lose too much by restricting learner's prediction to $\Delta^\delta(\mathcal{Y})$:

$$\mathcal{R}_T(\mathcal{F}) \leq \left\langle \left\langle \sup_{x_t} \sup_{p_t} \mathbb{E}_{y_t \sim p_t} \right\rangle \right\rangle_{t=1}^T \sup_{f \in \mathcal{F}} \left[\sum_{t=1}^T \inf_{\hat{p}_t} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{p}_t, y_t)] - \ell(f_t, y_t) \right] + |\mathcal{Y}| \delta T.$$

490 Sending $\delta \rightarrow 0^+$ on the RHS of the above inequality, we get

$$\mathcal{R}_T(\mathcal{F}) \leq \left\langle \left\langle \sup_{x_t} \sup_{p_t} \mathbb{E}_{y_t \sim p_t} \right\rangle \right\rangle_{t=1}^T \sup_{f \in \mathcal{F}} \left[\sum_{t=1}^T \inf_{\hat{p}_t} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{p}_t, y_t)] - \ell(f_t, y_t) \right].$$

491 It is easy to see that on the RHS of the above inequality, the inner infimum over $\hat{p}_t \in \Delta(\mathcal{Y})$ is
492 achieved at $\hat{p}_t = p_t$ due to the nature of log loss. So

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}) &\leq \left\langle \left\langle \sup_{x_t} \sup_{p_t} \mathbb{E}_{y_t \sim p_t} \right\rangle \right\rangle_{t=1}^T \sup_{f \in \mathcal{F}} \left[\sum_{t=1}^T \mathbb{E}_{y_t \sim p_t} [\ell(p_t, y_t)] - \ell(f_t, y_t) \right] \\ &= \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} [\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})], \end{aligned}$$

493 where in the last equality we use the compact notation of trees to further simplify our expression
494 and this concludes the proof. ■

495 **Lemma A.3** For any hypothesis class \mathcal{F} and horizon T ,

$$\sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} [\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})] = \sup_{\mathbf{x}} \log S_T(\mathcal{F}|\mathbf{x}). \quad (5)$$

496 It is implied that whenever \mathcal{F} satisfies Eq. (4), we have

$$\mathcal{R}_T(\mathcal{F}) = \sup_{\mathbf{x}} \log S_T(\mathcal{F}|\mathbf{x}).$$

497

498 **Proof of Lemma A.3** First we can see that the outcome sequence $y_{1:T}$ generated under any tree
499 \mathbf{p} is the same thing as $y_{1:T}$ generated by its associated joint distribution over \mathcal{Y}^T , and vice versa.
500 So we can replace the supremum over trees \mathbf{p} in the LHS of Eq. (5) by the supremum over joint
501 distributions P over \mathcal{Y}^T . Hence,

$$\begin{aligned} \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} [\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})] &= \sup_{\mathbf{x}, P} \mathbb{E}_{\mathbf{y} \sim P} [\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})] \\ &= \sup_{\mathbf{x}, P} \mathbb{E}_{\mathbf{y} \sim P} \left[\sum_{t=1}^T \ell(P_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right], \end{aligned}$$

502 where P_t denotes the conditional distribution $P_t(\cdot|y_{1:t-1}) \in \Delta(\mathcal{Y})$ of y_t under P given $y_{1:t-1}$.

503 Now fix the context tree \mathbf{x} and distribution P . Then we can see that $\mathbb{E}_{\mathbf{y} \sim P} [\ell(P_t, y_t)] =$
504 $H(y_t|y_{1:t-1})$. So $\mathbb{E}_{\mathbf{y} \sim P} [\sum_{t=1}^T \ell(P_t, y_t)] = \sum_{t=1}^T H(y_t|y_{1:t-1}) = H(P)$. Further notice that

$$\inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) = \inf_{f \in \mathcal{F}} (-\log P_f(y_{1:T}|x_{1:T})) = -\sup_{f \in \mathcal{F}} \log P_f(y_{1:T}|x_{1:T}).$$

505 So naturally we define the map $F_{\mathbf{x}} : \mathcal{Y}^T \rightarrow \mathbb{R} \cup \{-\infty\}$ by

$$F_{\mathbf{x}}(\mathbf{y}) = \sup_{f \in \mathcal{F}} \log P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})),$$

506 and then we see that

$$\mathbb{E}_{\mathbf{y} \sim P} \left[\sum_{t=1}^T \ell(P_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right] = H(P) + \mathbb{E}_{\mathbf{y} \sim P} [F_{\mathbf{x}}(\mathbf{y})].$$

507 For any given tree \mathbf{x} , the optimization problem

$$\sup_{P \in \Delta(\mathcal{Y}^T)} H(P) + \mathbb{E}_{\mathbf{y} \sim P} [F_{\mathbf{x}}(\mathbf{y})]$$

508 is actually a maximization problem, for which the optimal P^* is given by

$$P^*(\mathbf{y}) = \frac{\exp(F_{\mathbf{x}}(\mathbf{y}))}{\sum_{\mathbf{y}'} \exp(F_{\mathbf{x}}(\mathbf{y}'))} = \frac{\sup_{f \in \mathcal{F}} P_f(\mathbf{y}|\mathbf{x}(\mathbf{y}))}{\sum_{\mathbf{y}'} \sup_{f \in \mathcal{F}} P_f(\mathbf{y}'|\mathbf{x}(\mathbf{y}'))}, \forall \mathbf{y} \in \mathcal{Y}^T.$$

509 Note that the above formula for P^* is also valid when $F_{\mathbf{x}}(\mathbf{y}) = -\infty$ for some \mathbf{y} , since P^* should
 510 be supported on $\{\mathbf{y} \in \mathcal{Y}^T : F_{\mathbf{x}}(\mathbf{y}) > -\infty\}$, and $F_{\mathbf{x}}(\mathbf{y})$ cannot be $-\infty$ for all \mathbf{y} due to Lemma D.1.
 511 The associated value of this maximization problem is

$$\log \left(\sum_{\mathbf{y}} \exp(F_{\mathbf{x}}(\mathbf{y})) \right) = \log \left(\sum_{\mathbf{y}} \sup_{f \in \mathcal{F}} P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})) \right).$$

512 Therefore,

$$\begin{aligned} \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} [\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})] &= \sup_{\mathbf{x}, P} \mathbb{E}_{\mathbf{y} \sim P} \left[\sum_{t=1}^T \ell(P_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right] \\ &= \sup_{\mathbf{x}} \sup_P \left\{ H(P) + \mathbb{E}_{\mathbf{y} \sim P} [F_{\mathbf{x}}(\mathbf{y})] \right\} \\ &= \sup_{\mathbf{x}} \log \left(\sum_{\mathbf{y}} \sup_{f \in \mathcal{F}} P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})) \right) \\ &= \sup_{\mathbf{x}} \log S_T(\mathcal{F}|\mathbf{x}). \end{aligned}$$

513 ■

514 In the proof of Lemma A.1, if we do not restrict the learner's prediction and simply swap the order
 515 of inf and sup to produce an inequality at each time t , we will reach the following folklore result.

516 **Lemma A.4** For any hypothesis class \mathcal{F} and horizon T ,

$$\mathcal{R}_T(\mathcal{F}) \geq \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} [\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})]. \quad (6)$$

517

518 **Proof of Lemma A.4** To get Eq. (6), we simply need to reverse the order of sup and inf at each time
 519 in the extensive formulation of minimax regret and produce an inequality:

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}) &= \sup_{x_1} \inf_{\hat{p}_1} \sup_{y_1} \cdots \sup_{x_T} \inf_{\hat{p}_T} \sup_{y_T} \mathcal{R}_T(\mathcal{F}; \hat{p}_{1:T}, x_{1:T}, y_{1:T}) \\ &= \left\langle \left\langle \sup_{x_t} \inf_{\hat{p}_t} \sup_{y_t} \right\rangle \right\rangle_{t=1}^T \left[\sum_{t=1}^T \ell(\hat{p}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f(x_{1:t}, y_{1:t-1}), y_t) \right] \\ &= \left\langle \left\langle \sup_{x_t} \inf_{\hat{p}_t} \sup_{y_t} \right\rangle \right\rangle_{t=1}^{T-1} \sup_{x_T} \inf_{\hat{p}_T} \sup_{p_T} \mathbb{E}_{y_T \sim p_T} \left[\sum_{t=1}^T \ell(\hat{p}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right] \\ &\geq \left\langle \left\langle \sup_{x_t} \inf_{\hat{p}_t} \sup_{y_t} \right\rangle \right\rangle_{t=1}^{T-1} \sup_{x_T} \sup_{p_T} \inf_{\hat{p}_T} \mathbb{E}_{y_T \sim p_T} \left[\sum_{t=1}^T \ell(\hat{p}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right] \\ &= \left\langle \left\langle \sup_{x_t} \inf_{\hat{p}_t} \sup_{y_t} \right\rangle \right\rangle_{t=1}^{T-1} \left[\sum_{t=1}^{T-1} \ell(\hat{p}_t, y_t) + \sup_{x_T} \sup_{p_T} \left[\inf_{\hat{p}_T} \mathbb{E}_{y_T \sim p_T} \ell(\hat{p}_T, y_T) - \mathbb{E}_{y_T \sim p_T} \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \right] \right]. \end{aligned}$$

520 Iterating the argument and rearranging terms as above, we will get that

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}) &\geq \left\langle \left\langle \sup_{x_t} \sup_{p_t} \mathbb{E}_{y_t \sim p_t} \right\rangle \right\rangle_{t=1}^T \sup_{f \in \mathcal{F}} \left[\sum_{t=1}^T \inf_{\hat{p}_t} \mathbb{E}_{y_t \sim p_t} [\ell(\hat{p}_t, y_t)] - \ell(f_t, y_t) \right] \\ &= \left\langle \left\langle \sup_{x_t} \sup_{p_t} \mathbb{E}_{y_t \sim p_t} \right\rangle \right\rangle_{t=1}^T \sup_{f \in \mathcal{F}} \left[\sum_{t=1}^T \mathbb{E}_{y_t \sim p_t} [\ell(p_t, y_t)] - \ell(f_t, y_t) \right] \\ &= \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} [\mathcal{R}_T(\mathcal{F}; \mathbf{p}(\mathbf{y}), \mathbf{x}(\mathbf{y}), \mathbf{y})]. \end{aligned}$$

521 ■

522 **A.2 Smooth truncated hypothesis class**

523 To remove the reliance on Eq. (3), we introduce a smooth truncated version of \mathcal{F} that always satisfies
 524 Eq. (3) and study its minimax regret as well as contextual Shtarkov sums, compared to those of the
 525 untruncated class \mathcal{F} . To be more specific, we will apply the smooth truncation map to hypotheses:
 526 for any $\delta \in (0, 1/2)$ and $f : (\mathcal{X} \times \mathcal{Y})^* \times \mathcal{X} \rightarrow \Delta(\mathcal{Y})$, we use f^δ to denote its smooth truncated
 527 counterpart $\tau_\delta \circ f$; for any hypothesis class \mathcal{F} , we use \mathcal{F}^δ to denote the corresponding smooth
 528 truncated class $\tau_\delta \circ \mathcal{F} = \{\tau_\delta \circ f : f \in \mathcal{F}\}$. It is easy to verify that any smooth truncated class \mathcal{F}^δ
 529 satisfies Eq. (3) and hence

$$\mathcal{R}_T(\mathcal{F}^\delta) = \sup_{\mathbf{x}} \log S_T(\mathcal{F}^\delta | \mathbf{x}).$$

530 Next we control the effect of truncation on the minimax regret.

531 **Lemma A.5** For any \mathcal{F}, T and $\delta \in (0, 1/2)$,

$$\mathcal{R}_T(\mathcal{F}) \leq \mathcal{R}_T(\mathcal{F}^\delta) + T \cdot \log(1 + |\mathcal{Y}|\delta).$$

532

533 **Proof of Lemma A.5** Fix threshold $\delta \in (0, 1/2)$ and hypothesis f . By Lemma A.7, for any given
 534 sequences $x_{1:T}, y_{1:T}$, there is

$$\sum_{t=1}^T \ell(f^\delta(x_{1:t}, y_{1:t-1}), y_t) - \sum_{t=1}^T \ell(f(x_{1:t}, y_{1:t-1}), y_t) \leq T \cdot \log(1 + |\mathcal{Y}|\delta). \quad (7)$$

535 Then, for any sequence of predictions $\hat{p}_{1:T}$,

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}; \hat{p}_{1:T}, x_{1:T}, y_{1:T}) &= \sum_{t=1}^T \ell(\hat{p}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f_t, y_t) \\ &\leq \sum_{t=1}^T \ell(\hat{p}_t, y_t) - \inf_{f^\delta \in \mathcal{F}^\delta} \sum_{t=1}^T \ell(f_t^\delta, y_t) + T \cdot \log(1 + |\mathcal{Y}|\delta) \\ &= \mathcal{R}_T(\mathcal{F}^\delta; \hat{p}_{1:T}, x_{1:T}, y_{1:T}) + T \cdot \log(1 + |\mathcal{Y}|\delta), \end{aligned}$$

536 which concludes the proof. ■

537 **Lemma A.6** There exists a constant $M(T) < \infty$ that only depends on T such that for any $f, x_{1:T} \in$
 538 $\mathcal{X}^T, y_{1:T} \in \mathcal{Y}^T$ and $\delta \in (0, 1/2)$,

$$P_{f^\delta}(y_{1:T} | x_{1:T}) \leq P_f(y_{1:T} | x_{1:T}) + \delta \cdot M(T).$$

539

540 **Proof of Lemma A.6** Fix threshold $\delta \in (0, 1/2)$, hypothesis f and sequences $x_{1:T}, y_{1:T}$. Then

$$\begin{aligned} P_{f^\delta}(y_{1:T} | x_{1:T}) &= \prod_{t=1}^T f_t^\delta(y_t) = \prod_t \left(\frac{f_t(y_t) + \delta}{1 + |\mathcal{Y}|\delta} \right) \\ &\leq \prod_t (f_t(y_t) + \delta) \\ &= \prod_t f_t(y_t) + \delta \cdot \sum_t \prod_{t' \neq t} f_{t'}(y_{t'}) + \dots + \delta^T \\ &\leq \prod_t f_t(y_t) + \delta \cdot M(T) \\ &= P_f(y_{1:T} | y_{1:T}) + \delta \cdot M(T), \end{aligned}$$

541 where we can set $M(T) = T + \binom{T}{2} + \binom{T}{3} + \dots + \binom{T}{T}$ since $f_t(y_t)$'s are bounded by 1. ■

542 **A.3 Putting together**

543 Now we are fully prepared to finish the proof of Theorem 3.2, our main result in Section 3.

544 **Proof of Theorem 3.2** By Lemma A.6, we have that for any context tree \mathbf{x} of depth T ,

$$\sum_{\mathbf{y} \in \mathcal{Y}^T} \sup_{f^\delta \in \mathcal{F}^\delta} P_{f^\delta}(\mathbf{y} | \mathbf{x}(\mathbf{y})) \leq \sum_{\mathbf{y} \in \mathcal{Y}^T} \sup_{f \in \mathcal{F}} P_f(\mathbf{y} | \mathbf{x}(\mathbf{y})) + \delta \cdot M(T) \cdot |\mathcal{Y}|^T.$$

545 Thus

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}^\delta) &= \sup_{\mathbf{x}} \log S_T(\mathcal{F}^\delta | \mathbf{x}) \\ &= \sup_{\mathbf{x}} \log \left(\sum_{\mathbf{y} \in \mathcal{Y}^T} \sup_{f^\delta \in \mathcal{F}^\delta} P_{f^\delta}(\mathbf{y} | \mathbf{x}(\mathbf{y})) \right) \\ &\leq \sup_{\mathbf{x}} \log \left(\sum_{\mathbf{y} \in \mathcal{Y}^T} \sup_{f \in \mathcal{F}} P_f(\mathbf{y} | \mathbf{x}(\mathbf{y})) + \delta \cdot M(T) \cdot |\mathcal{Y}|^T \right) \\ &= \log \left(\sup_{\mathbf{x}} \sum_{\mathbf{y} \in \mathcal{Y}^T} \sup_{f \in \mathcal{F}} P_f(\mathbf{y} | \mathbf{x}(\mathbf{y})) + \delta \cdot M(T) \cdot |\mathcal{Y}|^T \right). \end{aligned}$$

546 Together with Lemma A.5, we get that for any $\delta \in (0, 1/2)$,

$$\mathcal{R}_T(\mathcal{F}) \leq \log \left(\sup_{\mathbf{x}} \sum_{\mathbf{y} \in \mathcal{Y}^T} \sup_{f \in \mathcal{F}} P_f(\mathbf{y} | \mathbf{x}(\mathbf{y})) + \delta \cdot M(T) \cdot |\mathcal{Y}|^T \right) + T \cdot \log(1 + |\mathcal{Y}|\delta). \quad (8)$$

547 After sending $\delta \rightarrow 0^+$ on the RHS of Eq. (8),

$$\mathcal{R}_T(\mathcal{F}) \leq \log \left(\sup_{\mathbf{x}} \sum_{\mathbf{y} \in \mathcal{Y}^T} \sup_{f \in \mathcal{F}} P_f(\mathbf{y} | \mathbf{x}(\mathbf{y})) \right) = \sup_{\mathbf{x}} \log S_T(\mathcal{F} | \mathbf{x}).$$

548 Recall that we have $\mathcal{R}_T(\mathcal{F}) \geq \sup_{\mathbf{x}} \log S_T(\mathcal{F} | \mathbf{x})$ from Lemma A.4 and Lemma A.3. So finally,

$$\mathcal{R}_T(\mathcal{F}) = \sup_{\mathbf{x}} \log S_T(\mathcal{F} | \mathbf{x}).$$

549 ■

550 **A.4 Additional proofs**

551 **Lemma A.7** For any $p \in \Delta(\mathcal{Y})$ and $\delta \in (0, 1/2)$,

$$\ell(\tau_\delta(p), y) \leq \ell(p, y) + \log(1 + |\mathcal{Y}|\delta) \leq \ell(p, y) + |\mathcal{Y}|\delta, \forall y \in \mathcal{Y}.$$

552

553 **Proof of Lemma A.7** By direct computation, for any $y \in \mathcal{Y}$,

$$\begin{aligned} \ell(\tau_\delta(p), y) - \ell(p, y) &= \log \left(\frac{p(y)}{p(y) + \delta} \cdot (1 + |\mathcal{Y}|\delta) \right) \\ &\leq \log(1 + |\mathcal{Y}|\delta) \\ &\leq |\mathcal{Y}|\delta. \end{aligned}$$

554 ■

555 **A.5 Proof of Proposition 3.3**

556 Starting from Theorem 3.2 that $\mathcal{R}_T(\mathcal{F}) = \sup_{\mathbf{x}} \log S_T(\mathcal{F}|\mathbf{x})$, we have

$$\begin{aligned} \mathcal{R}_T(\mathcal{F}) &= \sup_{\mathbf{x}} \log \left(\sum_{\mathbf{y}} \sup_{f \in \mathcal{F}} P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})) \right) \\ &\leq \sup_{\mathbf{x}} \log \left(\sum_{\mathbf{y}} \sum_{f \in \mathcal{F}} P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})) \right) \\ &= \sup_{\mathbf{x}} \log \left(\sum_{f \in \mathcal{F}} \sum_{\mathbf{y}} P_f(\mathbf{y}|\mathbf{x}(\mathbf{y})) \right) = \log |\mathcal{F}|, \end{aligned}$$

557 where the last equality is due to Lemma D.1.

558 **B Proofs for Section 4**

559 **Notations.** Again we may use f_t to denote the probability vector $f(x_{1:t}, y_{1:t-1}) \in \Delta(\mathcal{Y})$ produced
 560 by hypothesis f at time t when the context and label sequences $x_{1:T}, y_{1:T}$ are clear from the context.
 561 For a context tree \mathbf{x} of depth $T-t$ and a path $\mathbf{y} \in \mathcal{Y}^{T-t}$, we re-index $\mathbf{x}(\mathbf{y})$ as $(\mathbf{x}_{t+1}(\mathbf{y}), \dots, \mathbf{x}_T(\mathbf{y}))$
 562 whenever it takes the last $T-t$ entries of the entire context sequence. And we do the same for the
 563 probabilistic tree \mathbf{p} as well. That is, whenever $\mathbf{y} = (y_{t+1}, \dots, y_T) \in \mathcal{Y}^{T-t}$ takes the last $T-t$
 564 entries of the whole label sequence and $\mathbf{y} \sim \mathbf{p}$, then we will denote this label generating process by
 565 $y_{t+1} \sim \mathbf{p}_{t+1}(\mathbf{y}), \dots, y_T \sim \mathbf{p}_T(\mathbf{y})$.

566 **B.1 Proof of Theorem 4.2**

567 **Proof of Theorem 4.2** Recall that the minimax regret is

$$\mathcal{R}_T(\mathcal{F}) = \left\langle \left\langle \sup_{x_t} \inf_{\hat{p}_t} \sup_{y_t} \right\rangle \right\rangle_{t=1}^T \left[\sum_{t=1}^T \ell(\hat{p}_t, y_t) - \inf_{f \in \mathcal{F}} \sum_{t=1}^T \ell(f(x_{1:t}, y_{1:t-1}), y_t) \right].$$

568 Through this extensive form of the minimax regret, we know that given $x_{1:t}, y_{1:t-1}$, the minimax
 569 prediction \hat{p}_t^* at round t is the one that minimizes the following expression over all $\hat{p}_t \in \Delta(\mathcal{Y})$:

$$\sup_{y_t} \left\langle \left\langle \sup_{x_s} \inf_{\hat{p}_s} \sup_{y_s} \right\rangle \right\rangle_{s=t+1}^T \left[\sum_{s=t}^T \ell(\hat{p}_s, y_s) - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f(x_{1:s}, y_{1:s-1}), y_s) \right]. \quad (9)$$

570 Define

$$G(\mathcal{F}, x_{1:t}, y_{1:t}) = \left\langle \left\langle \sup_{x_s} \inf_{\hat{p}_s} \sup_{y_s} \right\rangle \right\rangle_{s=t+1}^T \left[\sum_{s=t+1}^T \ell(\hat{p}_s, y_s) - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f(x_{1:s}, y_{1:s-1}), y_s) \right],$$

571 and now

$$\hat{p}_t^* = \operatorname{argmin}_{\hat{p}_t \in \Delta(\mathcal{Y})} \sup_{y_t} \left\{ \ell(\hat{p}_t, y_t) + G(\mathcal{F}, x_{1:t}, y_{1:t}) \right\}.$$

572 The crux of the proof is to show the following:

573 **Lemma B.1** For any hypothesis class \mathcal{F} and sequences $x_{1:t} \in \mathcal{X}^t, y_{1:t} \in \mathcal{Y}^t$,

$$G(\mathcal{F}, x_{1:t}, y_{1:t}) = \sup_{\mathbf{x}} \log S_T^{x_{1:t}, y_{1:t}}(\mathcal{F}|\mathbf{x}).$$

574 The proof of Lemma B.1 is done by essentially following the same strategy in Appendix A since
 575 $G(\mathcal{F}, x_{1:t}, y_{1:t})$ admits a similar extensive form with the minimax regret $\mathcal{R}_T(\mathcal{F})$. For completeness
 576 we provide its proof in Appendix B.2. Given Lemma B.1, we have

$$\begin{aligned} \hat{p}_t^* &= \operatorname{argmin}_{\hat{p}_t \in \Delta(\mathcal{Y})} \sup_{y_t} \left\{ \ell(\hat{p}_t, y_t) + \sup_{\mathbf{x}} \log S_T^{x_{1:t}, y_{1:t}}(\mathcal{F}|\mathbf{x}) \right\} \\ &= \operatorname{argmin}_{\hat{p}_t \in \Delta(\mathcal{Y})} \sup_{y_t} \log \left(\frac{\sup_{\mathbf{x}} S_T^{x_{1:t}, y_{1:t}}(\mathcal{F}|\mathbf{x})}{\hat{p}_t(y_t)} \right). \end{aligned}$$

577 We apply the following result to solve the above program:

578 **Lemma B.2** [MG22, Lemma 15] Let $g : \mathcal{Y} \rightarrow [0, +\infty]$ be a measurable function such that
 579 $\int_{\mathcal{Y}} g(y) d\mu \in (0, +\infty)$. Then,

$$\inf_p \sup_{y \in \mathcal{Y}} \log \frac{g(y)}{p(y)} = \log \left(\int_{\mathcal{Y}} g(y) \mu(dy) \right), \quad (10)$$

580 where the infimum in Eq. (10) spans over all probability densities $p : \mathcal{Y} \rightarrow [0, +\infty)$ with respect to
 581 μ , and the infimum is reached at

$$p^* = \frac{g}{\int_{\mathcal{Y}} g(y) d\mu}.$$

582 Letting $g(y) = \sup_{\mathbf{x}} S_T^{x_{1:t}, (y_{1:t-1}, y)}(\mathcal{F}|\mathbf{x}) \in [0, 1]$ and μ be the counting measure on the finite space
 583 \mathcal{Y} , we can apply Lemma B.2 whenever not all $g(y)$'s are 0. In this case, we solve that

$$\hat{p}_t^*(y) = \frac{g(y)}{\sum_{y' \in \mathcal{Y}} g(y')} = \frac{\sup_{\mathbf{x}} S_T^{x_{1:t}, (y_{1:t-1}, y)}(\mathcal{F}|\mathbf{x})}{\sum_{y' \in \mathcal{Y}} \sup_{\mathbf{x}} S_T^{x_{1:t}, (y_{1:t-1}, y')}(\mathcal{F}|\mathbf{x})}, \forall y \in \mathcal{Y}.$$

584 On the other hand, if $g(y) = 0, \forall y \in \mathcal{Y}$, then any \hat{p}_t such that $\hat{p}_t(y) > 0, \forall y \in \mathcal{Y}$, is an minimax
 585 optimal prediction. Moreover, it implies that $P_f(y_{1:t-1} | x_{1:t-1}) = 0, \forall f \in \mathcal{F}$. This is because for
 586 arbitrary context tree \mathbf{x} ,

$$\begin{aligned} 0 &= \sum_{y_t} \sum_{\mathbf{y} \in \mathcal{Y}^{T-t}} P_f(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y})) \\ &= \sum_{y_t} P_f(y_{1:t} | x_{1:t}) \\ &= P_f(y_{1:t-1} | x_{1:t-1}). \end{aligned}$$

587 So the cumulative loss for each expert f up to round $t - 1$ already blows up to $+\infty$ and
 588 the learner only needs to predict an arbitrary $\hat{p} \in \Delta^+(\mathcal{Y})$ in all remaining rounds to achieve
 589 $\mathcal{R}_T(\mathcal{F}; \hat{p}_{1:T}, x_{1:T}, y_{1:T}) = -\infty$.

590 Overall, we can see that the minimax optimal prediction $\hat{p}_t^* \in \Delta(\mathcal{Y})$ at round t given $x_{1:t}, y_{1:t-1}$ is

$$\hat{p}_t^*(y) = \frac{\sup_{\mathbf{x}} S_T^{x_{1:t}, (y_{1:t-1}, y)}(\mathcal{F}|\mathbf{x})}{\sum_{y' \in \mathcal{Y}} \sup_{\mathbf{x}} S_T^{x_{1:t}, (y_{1:t-1}, y')}(\mathcal{F}|\mathbf{x})}, \forall y \in \mathcal{Y},$$

591 if there exists $y \in \mathcal{Y}$ such that $\sup_{\mathbf{x}} S_T^{x_{1:t}, (y_{1:t-1}, y)}(\mathcal{F}|\mathbf{x}) > 0$. Otherwise, select \hat{p}_t^* to be an arbitrary
 592 element in $\Delta^+(\mathcal{Y})$ (and so do all remaining rounds). ■

593 B.2 Auxiliary lemmas

594 Recall that for any hypothesis class \mathcal{F} and sequences $x_{1:t} \in \mathcal{X}^t, y_{1:t} \in \mathcal{Y}^t$,

$$\begin{aligned} G(\mathcal{F}, x_{1:t}, y_{1:t}) &= \left\langle \left\langle \sup_{x_s} \inf_{\hat{p}_s} \sup_{y_s} \right\rangle \right\rangle_{s=t+1}^T \left[\sum_{s=t+1}^T \ell(\hat{p}_s, y_s) - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f(x_{1:s}, y_{1:s-1}), y_s) \right] \\ &= \left\langle \left\langle \sup_{x_s} \inf_{\hat{p}_s} \sup_{p_s} \mathbb{E}_{y_s \sim p_s} \right\rangle \right\rangle_{s=t+1}^T \left[\sum_{s=t+1}^T \ell(\hat{p}_s, y_s) - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f(x_{1:s}, y_{1:s-1}), y_s) \right]. \end{aligned}$$

595 To prove Lemma B.1, we need the following lemmas.

596 **Lemma B.3** For any hypothesis class \mathcal{F} and sequences $x_{1:t} \in \mathcal{X}^t, y_{1:t} \in \mathcal{Y}^t$,

$$G(\mathcal{F}, x_{1:t}, y_{1:t}) \geq \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} \left[\sum_{s=t+1}^T \mathbb{E}_{y_s \sim \mathbf{p}_s(\mathbf{y})} [\ell(\mathbf{p}_s(\mathbf{y}), y_s)] - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right]. \quad (11)$$

597 And whenever for every $x_{t+1:T} \in \mathcal{X}^{T-t}, y_{t+1:T} \in \mathcal{Y}^{T-t}$, it holds

$$\inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f(x_{1:s}, y_{1:s-1}), y_s) < \infty, \quad (12)$$

598 then

$$G(\mathcal{F}, x_{1:t}, y_{1:t}) = \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} \left[\sum_{s=t+1}^T \mathbb{E}_{y_s \sim \mathbf{p}_s(\mathbf{y})} [\ell(\mathbf{p}_s(\mathbf{y}), y_s)] - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right]. \quad (13)$$

599

600 **Proof of Lemma B.3** First we see that similar to the proof of Lemma A.4, we can reverse every pair
601 of sup over p_s and inf over \hat{p}_s in the extensive formulation of $G(\mathcal{F}, x_{1:t}, y_{1:t})$ and rearrange terms
602 to obtain

$$G(\mathcal{F}, x_{1:t}, y_{1:t}) \geq \left\langle \left\langle \sup_{x_s} \sup_{p_s} \mathbb{E}_{y_s \sim p_s} \right\rangle \right\rangle_{s=t+1}^T \left[\sum_{s=t+1}^T \inf_{\hat{p}_s} \mathbb{E}_{y_s \sim p_s} [\ell(\hat{p}_s, y_s)] - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right],$$

603 and again due to the nature of log loss,

$$\begin{aligned} G(\mathcal{F}, x_{1:t}, y_{1:t}) &\geq \left\langle \left\langle \sup_{x_s} \sup_{p_s} \mathbb{E}_{y_s \sim p_s} \right\rangle \right\rangle_{s=t+1}^T \left[\sum_{s=t+1}^T \mathbb{E}_{y_s \sim p_s} [\ell(p_s, y_s)] - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right] \\ &= \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} \left[\sum_{s=t+1}^T \mathbb{E}_{y_s \sim \mathbf{p}_s(\mathbf{y})} [\ell(\mathbf{p}_s(\mathbf{y}), y_s)] - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right], \end{aligned}$$

604 where in the last step we compress the expression using trees (of depth $T-t$) and Eq. (11) is proved.

605 To show that the minimax swap is valid under Eq. (12), we follow the same strategy as in the proof
606 of Lemma A.1 by restricting the learner's prediction \hat{p}_s to $\Delta^\delta(\mathcal{Y})$ for any threshold $\delta \in (0, 1/2)$
607 which yields

$$\begin{aligned} G(\mathcal{F}, x_{1:t}, y_{1:t}) &\leq \left\langle \left\langle \sup_{x_s} \sup_{p_s} \mathbb{E}_{y_s \sim p_s} \right\rangle \right\rangle_{s=t+1}^T \left[\sum_{s=t+1}^T \inf_{\hat{p}_s \in \Delta^\delta(\mathcal{Y})} \mathbb{E}_{y_s \sim p_s} [\ell(\hat{p}_s, y_s)] - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right] \\ &\leq \left\langle \left\langle \sup_{x_s} \sup_{p_s} \mathbb{E}_{y_s \sim p_s} \right\rangle \right\rangle_{s=t+1}^T \left[\sum_{s=t+1}^T \inf_{\hat{p}_s} \mathbb{E}_{y_s \sim p_s} [\ell(\hat{p}_s, y_s)] - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right] + |\mathcal{Y}| \delta T. \end{aligned}$$

608 So Eq. (13) is proved by sending $\delta \rightarrow 0^+$ on the RHS of the last inequality and the established
609 Eq. (11). ■

610 **Lemma B.4** For any hypothesis class \mathcal{F} and sequences $x_{1:t} \in \mathcal{X}^t, y_{1:t} \in \mathcal{Y}^t$,

$$\sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} \left[\sum_{s=t+1}^T \mathbb{E}_{y_s \sim \mathbf{p}_s(\mathbf{y})} [\ell(\mathbf{p}_s(\mathbf{y}), y_s)] - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right] = \sup_{\mathbf{x}} \log S_T^{x_{1:t}, y_{1:t}}(\mathcal{F} | \mathbf{x}).$$

611

612 **Proof of Lemma B.4** The proof follows that of Lemma A.3. By replacing the probabilistic tree \mathbf{p}
613 by the joint distribution $P \in \Delta(\mathcal{Y}^{T-t})$, we get

$$\begin{aligned} &\sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} \left[\sum_{s=t+1}^T \mathbb{E}_{y_s \sim \mathbf{p}_s(\mathbf{y})} [\ell(\mathbf{p}_s(\mathbf{y}), y_s)] - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right] \\ &= \sup_{\mathbf{x}, P} \mathbb{E}_{\mathbf{y} \sim P} \left[\sum_{s=t+1}^T \ell(P_s, y_s) - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right] \\ &= \sup_{\mathbf{x}} \sup_{P \in \Delta(\mathcal{Y}^{T-t})} H(P) + \mathbb{E}_{\mathbf{y} \sim P} \left[\sup_{f \in \mathcal{F}} \log P_f(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y})) \right]. \end{aligned}$$

614 Similarly, for any fixed \mathbf{x} , define the map $F_{\mathbf{x}}^{x_{1:t}, y_{1:t}} : \mathcal{Y}^{T-t} \rightarrow \mathbb{R} \cup \{-\infty\}$ by

$$F_{\mathbf{x}}^{x_{1:t}, y_{1:t}}(\mathbf{y}) = \sup_{f \in \mathcal{F}} \log P_f(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y})),$$

615 and now we solve

$$\sup_{P \in \Delta(\mathcal{Y}^{T-t})} H(P) + \mathbb{E}_{\mathbf{y} \sim P} [F_{\mathbf{x}}^{x_{1:t}, y_{1:t}}(\mathbf{y})].$$

616 If there exists some $\mathbf{y} \in \mathcal{Y}^{T-t}$ such that $F_{\mathbf{x}}^{x_{1:t}, y_{1:t}}(\mathbf{y}) > -\infty$, then the optimal P^* is given by

$$P^*(\mathbf{y}) = \frac{\exp(F_{\mathbf{x}}^{x_{1:t}, y_{1:t}}(\mathbf{y}))}{\sum_{\mathbf{y}'} \exp(F_{\mathbf{x}}^{x_{1:t}, y_{1:t}}(\mathbf{y}'))} = \frac{\sup_{f \in \mathcal{F}} P_f(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y}))}{\sum_{\mathbf{y}'} \sup_{f \in \mathcal{F}} P_f(y_{1:t}, \mathbf{y}' | x_{1:t}, \mathbf{x}(\mathbf{y}'))}, \forall \mathbf{y} \in \mathcal{Y}^{T-t},$$

617 and then

$$\begin{aligned} & \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} \left[\sum_{s=t+1}^T \mathbb{E}_{y_s \sim \mathbf{p}_s(\mathbf{y})} [\ell(\mathbf{p}_s(\mathbf{y}), y_s)] - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right] \\ &= \sup_{\mathbf{x}} \sup_{P \in \Delta(\mathcal{Y}^{T-t})} H(P) + \mathbb{E}_{\mathbf{y} \sim P} [F_{\mathbf{x}}^{x_{1:t}, y_{1:t}}(\mathbf{y})] \\ &= \sup_{\mathbf{x}} \log \left(\sum_{\mathbf{y}} \sup_{f \in \mathcal{F}} P_f(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y})) \right) \\ &= \sup_{\mathbf{x}} \log S_T^{x_{1:t}, y_{1:t}}(\mathcal{F} | \mathbf{x}). \end{aligned}$$

618 However, if $F_{\mathbf{x}}^{x_{1:t}, y_{1:t}}(\mathbf{y}) = -\infty$ for all \mathbf{y} , then it implies that for any context tree \mathbf{x} , path \mathbf{y} , and
619 $f \in \mathcal{F}$, $P_f(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y})) = 0$ and hence,

$$\begin{aligned} & \sup_{\mathbf{x}, \mathbf{p}} \mathbb{E}_{\mathbf{y} \sim \mathbf{p}} \left[\sum_{s=t+1}^T \mathbb{E}_{y_s \sim \mathbf{p}_s(\mathbf{y})} [\ell(\mathbf{p}_s(\mathbf{y}), y_s)] - \inf_{f \in \mathcal{F}} \sum_{s=1}^T \ell(f_s, y_s) \right] \\ &= \sup_{\mathbf{x}} \sup_{P \in \Delta(\mathcal{Y}^{T-t})} H(P) + \mathbb{E}_{\mathbf{y} \sim P} [F_{\mathbf{x}}^{x_{1:t}, y_{1:t}}(\mathbf{y})] \\ &= -\infty \\ &= \sup_{\mathbf{x}} \log \left(\sum_{\mathbf{y}} \sup_{f \in \mathcal{F}} P_f(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y})) \right) \\ &= \sup_{\mathbf{x}} \log S_T^{x_{1:t}, y_{1:t}}(\mathcal{F} | \mathbf{x}), \end{aligned}$$

620 which finishes our proof. ■

621 Now we are able to prove the key result Lemma B.1.

622 **Proof of Lemma B.1** Fix any hypothesis class \mathcal{F} and sequences $x_{1:t} \in \mathcal{X}^t, y_{1:t} \in \mathcal{Y}^t$. First we
623 know

$$G(\mathcal{F}, x_{1:t}, y_{1:t}) \geq \sup_{\mathbf{x}} \log S_T^{x_{1:t}, y_{1:t}}(\mathcal{F} | \mathbf{x})$$

624 due to Eq. (11) and Lemma B.4. For the other direction, let us fix any threshold value $\delta \in (0, 1/2)$
625 and then

$$\begin{aligned} G(\mathcal{F}, x_{1:t}, y_{1:t}) &\leq G(\mathcal{F}^\delta, x_{1:t}, y_{1:t}) + T \cdot \log(1 + |\mathcal{Y}|\delta) \\ &= \sup_{\mathbf{x}} \log S_T^{x_{1:t}, y_{1:t}}(\mathcal{F}^\delta | \mathbf{x}) + T \cdot \log(1 + |\mathcal{Y}|\delta) \\ &= \sup_{\mathbf{x}} \log \left(\sum_{\mathbf{y} \in \mathcal{Y}^{T-t}} \sup_{f \in \mathcal{F}^\delta} P_{f^\delta}(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y})) \right) + T \cdot \log(1 + |\mathcal{Y}|\delta) \\ &\leq \sup_{\mathbf{x}} \log \left(\sum_{\mathbf{y} \in \mathcal{Y}^{T-t}} \sup_{f \in \mathcal{F}} P_f(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y})) + \delta \cdot M(T) \cdot |\mathcal{Y}|^T \right) + T \cdot \log(1 + |\mathcal{Y}|\delta) \\ &= \log \left(\sup_{\mathbf{x}} \sum_{\mathbf{y} \in \mathcal{Y}^{T-t}} \sup_{f \in \mathcal{F}} P_f(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y})) + \delta \cdot M(T) \cdot |\mathcal{Y}|^T \right) + T \cdot \log(1 + |\mathcal{Y}|\delta), \end{aligned}$$

626 where we have applied Lemma A.7, Lemma B.3, Lemma B.4, and Lemma A.6 accordingly. Simi-
627 larly, we send $\delta \rightarrow 0^+$ on the RHS of the last inequality and get

$$G(\mathcal{F}, x_{1:t}, y_{1:t}) \leq \sup_{\mathbf{x}} \log \left(\sum_{\mathbf{y} \in \mathcal{Y}^{T-t}} \sup_{f \in \mathcal{F}} P_f(y_{1:t}, \mathbf{y} | x_{1:t}, \mathbf{x}(\mathbf{y})) \right) = \sup_{\mathbf{x}} \log S_T^{x_{1:t}, y_{1:t}}(\mathcal{F} | \mathbf{x}),$$

628 which concludes the proof. ■

629 C Additional discussions

630 C.1 On the time-variant context space

631 In this section we generalize our analysis to the setting where the context space can evolve over time.
 632 We model time-varying context sets by a sequence of maps $\mathcal{X}_t : \mathcal{X}^{t-1} \times \mathcal{Y}^{t-1} \rightarrow 2^{\mathcal{X}}$, $t \in [T]$ as in
 633 [RS15; BFR20]. In each round t , instead of picking any context from \mathcal{X} , the nature is now required
 634 to only choose x_t from $\mathcal{X}_t(x_{1:t-1}, y_{1:t-1}) \subseteq \mathcal{X}$. Then the minimax regret with respect to $(\mathcal{X}_t)_{t \in [T]}$
 635 is rewritten as

$$\mathcal{R}_T(\mathcal{F}) = \left\langle \left\langle \sup_{x_t \in \mathcal{X}_t(x_{1:t-1}, y_{1:t-1})} \inf_{\hat{p}_t} \sup_{y_t} \right\rangle \right\rangle_{t=1}^T \mathcal{R}_T(\mathcal{F}; \hat{p}_{1:T}, x_{1:T}, y_{1:T}).$$

636 A context tree \mathbf{x} is *consistent* with respect to $(\mathcal{X}_t)_{t \in [T]}$ if for all $t \in [T]$ and $\mathbf{y} \in \mathcal{Y}^T$, $\mathbf{x}_t(\mathbf{y}) \in$
 637 $\mathcal{X}_t(x_{1:t-1}, y_{1:t-1})$. Then our results in Section 3 and Section 4 can be generalized simply by replac-
 638 ing the supremum over all context trees (of depth- T) by the supremum over all consistent context
 639 trees. For example, we will have

$$\mathcal{R}_T(\mathcal{F}) = \sup_{\mathbf{x}: \mathbf{x} \text{ is consistent}} \log S_T(\mathcal{F} | \mathbf{x}).$$

640 C.2 On the global and non-global sequential cover

641 Now we go back to consider the usual setting of binary label and constant experts, i.e., $\mathcal{Y} = \{0, 1\}$
 642 and $\mathcal{F} \subseteq [0, 1]^{\mathcal{X}}$. As mentioned in Section 3, previous works [BFR20; WHGS23] provided re-
 643 gret upper bounds based on ℓ_∞ sequential entropy. More specifically, both of their bounds are in
 644 the form of $O(\inf_{\alpha > 0} \{\alpha T + \mathcal{H}(\mathcal{F}, \alpha, T)\})$, with $\mathcal{H}(\mathcal{F}, \alpha, T)$ being either the non-global entropy
 645 $\mathcal{H}_\infty(\mathcal{F}, \alpha, T)$ or the global entropy $\mathcal{H}_G(\mathcal{F}, \alpha, T)$. It is then natural to ask which one of these
 646 two bounds is tighter. Although it is straightforward to prove that $\mathcal{H}_\infty(\mathcal{F}, \alpha, T)$ is no larger than
 647 $\mathcal{H}_G(\mathcal{F}, \alpha, T)$, the gap between them is at most a polylog factor, as shown below. The proof of
 648 $\mathcal{H}_\infty(\mathcal{F}, \alpha, T) \leq \mathcal{H}_G(\mathcal{F}, \alpha, T)$ is also included for completeness.

649 **Proposition C.1** *For any scale $\alpha > 0$, we have*

$$\mathcal{H}_\infty(\mathcal{F}, \alpha, T) \geq \min\{T, \sup_{\alpha' > \alpha} \text{sfat}_{2\alpha'}(\mathcal{F})\} \cdot \log(2).$$

650 *Therefore, together with $\mathcal{H}_\infty(\mathcal{F}, \alpha, T) \leq \mathcal{H}_G(\mathcal{F}, \alpha, T)$ and the folklore $\mathcal{H}_G(\mathcal{F}, \alpha, T) \leq$
 651 $O(\text{sfat}_\alpha(\mathcal{F}) \log(T/\alpha))$, we conclude that the regret upper bounds $O(\inf_{\alpha > 0} \{\alpha T +$
 652 $\mathcal{H}(\mathcal{F}, \alpha, T)\})$, $\mathcal{H} \in \{\mathcal{H}_\infty, \mathcal{H}_G\}$, differ by at most a polylog factor.*

653

654 **Proof of Proposition C.1** Fix any $\alpha' > \alpha > 0$ and let $d_{\alpha'}$ denote $\min\{T, \text{sfat}_{2\alpha'}(\mathcal{F})\}$. Then
 655 there exists a context tree \mathbf{x} and a witness tree \mathbf{s} , both of depth $d_{\alpha'}$, satisfying that for any path
 656 $\mathbf{y} \in \{0, 1\}^{d_{\alpha'}}$, there exists an $f \in \mathcal{F}$ such that

$$\forall t \in [d_{\alpha'}], (2y_t - 1) \cdot (f(x_t(\mathbf{y})) - s_t(\mathbf{y})) \geq \alpha' > \alpha. \quad (14)$$

657 Let $V_{\mathbf{x}, \alpha}$ be an arbitrary sequential ℓ_∞ covering of \mathcal{F} on \mathbf{x} . Now we select a path \mathbf{y} and a sequence
 658 of subsets $V_{\mathbf{x}, \alpha}^{(t)} \subseteq V_{\mathbf{x}, \alpha}$, $t \in [d_{\alpha'}]$ in the following recursive way. Define $V_{\mathbf{x}, \alpha}^{(0)} = V_{\mathbf{x}, \alpha}$. For
 659 each $t \in [d_{\alpha'}]$, choose $y_t \in \{0, 1\}$ such that $2y_t - 1 \in \{-1, +1\}$ is the minority among all
 660 $\text{sgn}(v_t(y_{1:t-1}) - s_t(y_{1:t-1}))$, $v \in V_{\mathbf{x}, \alpha}^{(t-1)}$ (ignoring those of 0's). Finally update $V_{\mathbf{x}, \alpha}^{(t)} = \{v \in$
 661 $V_{\mathbf{x}, \alpha}^{(t-1)} : \text{sgn}(v_t(y_{1:t-1}) - s_t(y_{1:t-1})) = 2y_t - 1\}$.

662 First we argue that, if there is any time $t' \in [d_{\alpha'}]$ such that $V_{\mathbf{x}, \alpha}^{(t'-1)} \neq \emptyset$, $V_{\mathbf{x}, \alpha}^{(t')} = \emptyset$, then $V_{\mathbf{x}, \alpha}$
 663 is not a valid cover of \mathcal{F} on \mathbf{x} . Otherwise, recall we have selected $y_1, \dots, y_{t'-1}$. Now pick an
 664 arbitrary $y_{t'} \in \{0, 1\}$. By Eq. (14) we can find some $f \in \mathcal{F}$ such that $(2y_{t'} - 1) \cdot (f(x_{t'}(y_{1:t'-1})) -$
 665 $s_{t'}(y_{1:t'-1})) > \alpha$, $\forall t \in [t']$. Since $V_{\mathbf{x}, \alpha}$ is a covering at scale α , there is $v \in V_{\mathbf{x}, \alpha}$ such that $|v_t(\mathbf{y}) -$
 666 $f(x_t(\mathbf{y}))| \leq \alpha$, $\forall t \in [t']$. This implies that $\text{sgn}(f(x_t(\mathbf{y})) - s_t(\mathbf{y})) = \text{sgn}(v_t(\mathbf{y}) - s_t(\mathbf{y})) =$
 667 $2y_t - 1$, $\forall t \in [t']$. So we can always find some member of $V_{\mathbf{x}, \alpha}^{(t'-1)}$ to match the minority sign of
 668 $v_{t'}(y_{1:t'-1}) - s_{t'}(y_{1:t'-1})$, $v \in V_{\mathbf{x}, \alpha}^{(t'-1)}$, which means that $V_{\mathbf{x}, \alpha}^{(t')} \neq \emptyset$ and yields a contradiction.

669 Now we know that $|V_{\mathbf{x},\alpha}^{(t)}| \geq 1, \forall t \in [d_{\alpha'}]$. By design $|V_{\mathbf{x},\alpha}^{(t)}| \leq |V_{\mathbf{x},\alpha}^{(t-1)}|/2, \forall t \in [d_{\alpha'}]$, so we
670 must have $|V_{\mathbf{x},\alpha}| = |V_{\mathbf{x},\alpha}^{(0)}| \geq 2^{d_{\alpha'}}$. As the choice of covering is arbitrary, the covering number
671 $\mathcal{N}_{\infty}(\mathcal{F} \circ \mathbf{x}, \alpha, d_{\alpha'})$ is also lower bounded by $2^{d_{\alpha'}}$ and hence $\mathcal{H}_{\infty}(\mathcal{F}, \alpha, d_{\alpha'}) \geq d_{\alpha'} \cdot \log(2)$. If
672 $\sup_{\alpha' > \alpha} \text{sfat}_{2\alpha'}(\mathcal{F}) \leq T$, then we get that

$$\mathcal{H}_{\infty}(\mathcal{F}, \alpha, T) \geq \sup_{\alpha' > \alpha} \mathcal{H}_{\infty}(\mathcal{F}, \alpha, \text{sfat}_{2\alpha'}(\mathcal{F})) \geq \sup_{\alpha' > \alpha} \text{sfat}_{2\alpha'}(\mathcal{F}) \cdot \log(2).$$

673 If there is some $\alpha' > \alpha$ such that $\text{sfat}_{2\alpha'}(\mathcal{F}) \geq T$, then

$$\mathcal{H}_{\infty}(\mathcal{F}, \alpha, T) = \mathcal{H}_{\infty}(\mathcal{F}, \alpha, d_{\alpha'}) \geq T \cdot \log(2).$$

674 Combining these two cases together, we have

$$\mathcal{H}_{\infty}(\mathcal{F}, \alpha, T) \geq \min\{T, \sup_{\alpha' > \alpha} \text{sfat}_{2\alpha'}(\mathcal{F})\} \cdot \log(2).$$

675 ■

676

677 **Proposition C.2** Let \mathcal{G}_{α} be a global sequential α -covering of \mathcal{F} as defined in [WHGS23]. Then for
678 any context tree \mathbf{x} , there exists a sequential cover $V_{\mathbf{x},\alpha}$ of $\mathcal{F} \circ \mathbf{x}$ at scale α with $|V_{\mathbf{x},\alpha}| \leq |\mathcal{G}_{\alpha}|$. This
679 implies that $\mathcal{H}_{\infty}(\mathcal{F}, \alpha, T) \leq \log |\mathcal{G}_{\alpha}|$.

680

681 **Proof of Proposition C.2** Fix arbitrary context tree \mathbf{x} . For any $g \in \mathcal{G}_{\alpha}$, define the $[0, 1]$ -valued tree
682 v^g by $v_t^g(\mathbf{y}) = g(x_{1:t}(\mathbf{y}))$, $\forall t \in [T], \mathbf{y} \in \mathcal{Y}^T$. Now let $V_{\mathbf{x},\alpha} = \{v^g : g \in \mathcal{G}_{\alpha}\}$ and we will show that
683 $V_{\mathbf{x},\alpha}$ is indeed a sequential cover of $\mathcal{F} \circ \mathbf{x}$ at scale α .

684 For any $f \in \mathcal{F}$ and $\mathbf{y} \in \mathcal{Y}^T$, tree \mathbf{x} yields a length- T sequence $x_{1:T}(\mathbf{y})$ and by definition of the
685 global sequential covering, there exists $g \in \mathcal{G}_{\alpha}$ such that

$$|f(x_t(\mathbf{y})) - g(x_{1:t}(\mathbf{y}))| \leq \alpha, \forall t \in [T].$$

686 So by our construction of $V_{\mathbf{x},\alpha}$, $v^g \in V_{\mathbf{x},\alpha}$ holds

$$|f(x_t(\mathbf{y})) - v_t^g(\mathbf{y})| = |f(x_t(\mathbf{y})) - g(x_{1:t}(\mathbf{y}))| \leq \alpha, \forall t \in [T],$$

687 which yields our claim after observing $|V_{\mathbf{x},\alpha}| \leq |\mathcal{G}_{\alpha}|$. ■

688 D Additional proofs

689 **Lemma D.1** For any \mathcal{X} -valued \mathcal{Y} -ary context tree \mathbf{x} of depth T , and $f : (\mathcal{X} \times \mathcal{Y})^* \times \mathcal{X} \rightarrow \Delta(\mathcal{Y})$,
690 we have

$$\sum_{\mathbf{y} \in \mathcal{Y}^T} P_f(\mathbf{y} | \mathbf{x}(\mathbf{y})) = 1, \tag{15}$$

691 where we recall that $\mathbf{x}(\mathbf{y})$ denotes the context sequence $(\mathbf{x}_1(\mathbf{y}), \dots, \mathbf{x}_T(\mathbf{y}))$.

692

693 **Proof of Lemma D.1** This is done by induction on the depth T . The key observation is that for any
694 label sequence \mathbf{y} , $\mathbf{x}_t(\mathbf{y}) = \mathbf{x}_t(y_1, \dots, y_{t-1})$ only depends on the first $t-1$ labels. For $T=1$, any
695 context tree \mathbf{x} is represented by its root node $\mathbf{x}_1(\cdot) = \mathbf{x}_1 \in \mathcal{X}$ and hence

$$\sum_{y_1} P_f(y_1 | \mathbf{x}_1) = \sum_{y_1} f(\mathbf{x}_1)(y_1) = 1.$$

696 Suppose Eq. (15) holds for all context trees \mathbf{x} of depth $T \leq d$ and all sequential functions f . Now
697 given any context tree $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_{d+1})$ of depth $T = d+1$, we denote its depth d subtree

698 $(\mathbf{x}_1, \dots, \mathbf{x}_d)$ by $\mathbf{x}_{[d]}$. Then

$$\begin{aligned}
 \sum_{\mathbf{y} \in \mathcal{Y}^{d+1}} P_f(\mathbf{y} | \mathbf{x}(\mathbf{y})) &= \sum_{y_{1:d}} \sum_{y_{d+1}} P_f(y_{1:d+1} | \mathbf{x}_1, \mathbf{x}_2(y_1), \dots, \mathbf{x}_{d+1}(y_{1:d})) \\
 &= \sum_{y_{1:d}} \sum_{y_{d+1}} P_f(y_{1:d} | \mathbf{x}_1, \dots, \mathbf{x}_d(y_{1:d-1})) \cdot f(\mathbf{x}_1, \dots, \mathbf{x}_{d+1}(y_{1:d}), y_{1:d})(y_{d+1}) \\
 &= \sum_{y_{1:d}} P_f(y_{1:d} | \mathbf{x}_1, \dots, \mathbf{x}_d(y_{1:d-1})) \sum_{y_{d+1}} f(\mathbf{x}_1, \dots, \mathbf{x}_{d+1}(y_{1:d}), y_{1:d})(y_{d+1}) \\
 &= \sum_{y_{1:d}} P_f(y_{1:d} | \mathbf{x}_1, \dots, \mathbf{x}_d(y_{1:d-1})) \\
 &= \sum_{\mathbf{y} \in \mathcal{Y}^d} P_f(\mathbf{y} | \mathbf{x}_{[d]}(\mathbf{y})) = 1,
 \end{aligned}$$

699 where the last step is due to induction. We are done. ■

700 **NeurIPS Paper Checklist**

701 The checklist is designed to encourage best practices for responsible machine learning research, ad-
702 dressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
703 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should
704 follow the references and precede the (optional) supplemental material. The checklist does NOT
705 count towards the page limit.

706 Please read the checklist guidelines carefully for information on how to answer these questions. For
707 each question in the checklist:

- 708 • You should answer [Yes] , [No] , or [NA]
- 709 • [NA] means either that the question is Not Applicable for that particular paper or the
710 relevant information is Not Available.
- 711 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

712 **The checklist answers are an integral part of your paper submission.** They are visible to the
713 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it
714 (after eventual revisions) with the final version of your paper, and its final version will be published
715 with the paper.

716 The reviewers of your paper will be asked to use the checklist as one of the factors in their evalu-
717 ation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]
718 " provided a proper justification is given (e.g., "error bars are not reported because it would be too
719 computationally expensive" or "we were unable to find the license for the dataset we used"). In
720 general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased
721 in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your
722 best judgment and write a justification to elaborate. All supporting evidence can appear either in the
723 main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question,
724 in the justification please point to the section(s) where related material for the question can be found.

725 **IMPORTANT**, please:

- 726 • **Delete this instruction block, but keep the section heading “NeurIPS paper checklist”,**
- 727 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 728 • **Do not modify the questions and only use the provided macros for your answers.**

729 **1. Claims**

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

732 Answer: [Yes]

733 Justification: Abstract summarizes theorems we have proven.

734 Guidelines:

- 735 • The answer NA means that the abstract and introduction do not include the claims
736 made in the paper.
- 737 • The abstract and/or introduction should clearly state the claims made, including the
738 contributions made in the paper and important assumptions and limitations. A No or
739 NA answer to this question will not be perceived well by the reviewers.
- 740 • The claims made should match theoretical and experimental results, and reflect how
741 much the results can be expected to generalize to other settings.
- 742 • It is fine to include aspirational goals as motivation as long as it is clear that these
743 goals are not attained by the paper.

744 **2. Limitations**

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

746 Answer: [Yes]

747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799

Justification: We discuss limitations in the Discussion section.

Guidelines:

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

3. Theory Assumptions and Proofs

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

Answer: [Yes]

Justification: We don't see how to justify this without machine checkable proofs, which we have not provided.

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: There are no experiments.

Guidelines:

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

831 5. Open access to data and code

832 Question: Does the paper provide open access to the data and code, with sufficient instruc-

833 tions to faithfully reproduce the main experimental results, as described in supplemental

834 material?

835 Answer: [NA]

836 Justification: There is no data or code.

837 Guidelines:

- 838 • The answer NA means that paper does not include experiments requiring code.
- 839 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/](https://nips.cc/public/guides/CodeSubmissionPolicy)
- 840 [public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 841 • While we encourage the release of code and data, we understand that this might not
- 842 be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
- 843 including code, unless this is central to the contribution (e.g., for a new open-source
- 844 benchmark).
- 845 • The instructions should contain the exact command and environment needed to run to
- 846 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 847 • The authors should provide instructions on data access and preparation, including how
- 848 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 849 • The authors should provide scripts to reproduce all experimental results for the new
- 850 proposed method and baselines. If only a subset of experiments are reproducible, they
- 851 should state which ones are omitted from the script and why.
- 852 • At submission time, to preserve anonymity, the authors should release anonymized
- 853 versions (if applicable).
- 854

- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

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

Answer: [NA]

Justification: There are no experiments.

Guidelines:

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

7. Experiment Statistical Significance

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

Answer: [NA]

Justification: There are no experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, 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: There are no 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.

- 905
- 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).
- 906
- 907
- 908
- 909

9. Code Of Ethics

910

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

913 Answer: [Yes]

914 Justification: We have read the code and do not see any violation. Our work relates to the
915 mathematical foundations of a basic task in ML.

916 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).
- 917
- 918
- 919
- 920
- 921

10. Broader Impacts

922

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

925 Answer: [Yes]

926 Justification: The paper presents a mathematical characterization of the limits of probabilistic
927 forecasting, and a (meta)algorithm that achieves these limits. For any class of interest,
928 there remains significant work to realize that algorithm in an efficient way. As such, our
929 impact is most directly on the theoretical community, who might then have direct societal
930 impact by producing an minimax optimal algorithm. As such, our societal impact may be
931 great, but it will always be quite indirect.

932 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).
- 933
- 934
- 935
- 936
- 937
- 938
- 939
- 940
- 941
- 942
- 943
- 944
- 945
- 946
- 947
- 948
- 949
- 950
- 951
- 952
- 953
- 954

11. Safeguards

955

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

959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008

Answer: [NA]

Justification: We are not releasing models or data.

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: We use no such 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: No new assets are introduced.

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

1009 Question: For crowdsourcing experiments and research with human subjects, does the pa-
1010 per include the full text of instructions given to participants and screenshots, if applicable,
1011 as well as details about compensation (if any)?

1012 Answer: [NA]

1013 Justification: No such experiments were performed.

1014 Guidelines:

- 1015 • The answer NA means that the paper does not involve crowdsourcing nor research
1016 with human subjects.
- 1017 • Including this information in the supplemental material is fine, but if the main contri-
1018 bution of the paper involves human subjects, then as much detail as possible should
1019 be included in the main paper.
- 1020 • According to the NeurIPS Code of Ethics, workers involved in data collection, cura-
1021 tion, or other labor should be paid at least the minimum wage in the country of the
1022 data collector.

1023 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 1024 Subjects

1025 Question: Does the paper describe potential risks incurred by study participants, whether
1026 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1027 approvals (or an equivalent approval/review based on the requirements of your country or
1028 institution) were obtained?

1029 Answer: [NA]

1030 Justification: Paper does not involve crowdsourcing or research with human subjects.

1031 Guidelines:

- 1032 • The answer NA means that the paper does not involve crowdsourcing nor research
1033 with human subjects.
- 1034 • Depending on the country in which research is conducted, IRB approval (or equiva-
1035 lent) may be required for any human subjects research. If you obtained IRB approval,
1036 you should clearly state this in the paper.
- 1037 • We recognize that the procedures for this may vary significantly between institutions
1038 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1039 guidelines for their institution.
- 1040 • For initial submissions, do not include any information that would break anonymity
1041 (if applicable), such as the institution conducting the review.