

# LEARNING TO ANSWER FROM CORRECT DEMONSTRATIONS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We study the problem of learning to generate an answer (or completion) to a question (or prompt), where there could be multiple correct answers, any one of which is acceptable at test time. Learning is based on demonstrations of some correct answer to each training question, as in Supervised Fine Tuning (SFT). We formalize the problem as offline imitation learning in contextual bandits, with demonstrations from some optimal policy, without explicitly observed rewards. Prior work assumes that the demonstrator belongs to a low-complexity policy class, which motivates maximum likelihood estimation (i.e., log-loss minimization). In contrast, we propose relying only on the reward model (specifying which answers are correct) being in a low-cardinality class, which we argue is a weaker assumption. We show that likelihood maximization methods can fail in this case, and instead devise an alternative novel approach that learns with sample complexity logarithmic in the cardinality of the reward class. Our approach and guarantees are robust and apply even when learning from arbitrary demonstrators and to the relaxed pass@ $k$  error setting. Our work motivates looking beyond likelihood maximization when learning from demonstrations.

## 1 INTRODUCTION

Many real-world problems involve generating answers to a question, where there may be many *equally good* responses. A math question can have millions of equally valid but differently written solutions, a coding task can admit many different perfectly working implementations, and a recommendation query can be satisfied by multiple items. The learner’s challenge is *not* to reproduce all correct responses, but to generate a *single good answer*.

This recurring structure can be formalized as *contextual bandits*. Each question corresponds to a context  $x \in \mathcal{X}$ , each candidate response corresponds to an action  $y \in \mathcal{Y}$ , and there is an unknown reward function  $r_*(x, y)$  (which we assume is binary for simplicity)<sup>1</sup>, indicating whether a response is good (equivalently correct when binary rewards). We consider the problem of learning from demonstrations given by *some optimal demonstrator*  $\pi_*$ , i.e., training set of  $x_i \sim \mathcal{D}, y_i \sim \pi_*(\cdot | x_i)$ , where  $\pi_*(\cdot | x)$  is supported on the set of optimal actions for the context  $x$ , given by

$$\sigma_*(x) := \{y \in \mathcal{Y} : r_*(x, y) = 1\}.$$

Note that  $|\sigma_*(x)|$  can be huge and there may be many (non-unique) optimal demonstrators. Our goal is to learn a predictor policy  $\hat{\pi}$  that produces a *good* action (response) to an unseen context (question) sampled from  $\mathcal{D}$ , i.e., it has low error of outputting actions that are incorrect, as captured by the following loss function.

$$L_{\mathcal{D}, \sigma_*}(\hat{\pi}) = \mathbb{E}_{x \sim \mathcal{D}, \hat{y} \sim \hat{\pi}(\cdot | x)} [\mathbb{1}\{\hat{y} \notin \sigma_*(x)\}] = \mathbb{E}_{x \sim \mathcal{D}, \hat{y} \sim \hat{\pi}(\cdot | x)} [\mathbb{1}\{r_*(x, \hat{y}) \neq 1\}]. \quad (1)$$

This is *offline imitation learning* (Rajaraman et al., 2020; Rashidinejad et al., 2021) in contextual bandits with an *optimal* demonstrator. In modern large language models (LLMs), this is exactly the problem addressed during the *supervised fine-tuning* (SFT) phase,<sup>2</sup> where the model is trained on

<sup>1</sup>We will return to the more general bounded reward functions later in the presentation.

<sup>2</sup>In LLMs, the actual response  $y$  consists of a sequence of tokens, and so it is a general Markov Decision Process (MDP). However, the transitions of this MDP are deterministic and rewards are generally sparse, and hence, contextual bandits provide with a useful abstraction for the purpose of our study. Indeed, this setup exactly corresponds to the prompt-completion formulation for LLMs (Ouyang et al., 2022; Rafailov et al., 2023; Huang et al., 2025; Wu et al., 2025).

054 curated datasets of prompt–completion pairs. Importantly, these demonstrations are *high-quality* but  
 055 *not exhaustive*.

056 An important remark about our objective (Eq. (1)) is that it is entirely reward-driven, focus-  
 057 ing solely on **validity**, which aligns with the *reward maximization* view of LLMs. We do  
 058 not require matching the distribution of  $\pi_*$  (see Appendix A for more discussion on this is-  
 059 sue). The typical situation we consider in our setup is during SFT, where the goal is to im-  
 060 prove the model’s performance on a specific task, and there are many ways to achieve good  
 061 performance—for example, producing gold-winning IMO solutions (reward = 1, i.e.,  $\sigma(\text{question}) =$   
 062  $\{\text{all completely correct solution texts to the problem}\}$ ). This set is unfathomably large, due both to  
 063 variation in solution approaches and details, and to differences in the text itself, down to word  
 064 choices, spacing, and punctuation. Rather than matching how a distribution over experts  $\pi_*$ , would  
 065 write their solution, producing any single correct solution is sufficient to win a gold medal at the  
 066 IMO.

067 **Policy Class Assumption and Likelihood Maximization.** A common approach when learning  
 068 from demonstrations is to assume that the demonstrator  $\pi_*$  belongs to a low-capacity policy class  
 069  $\Pi$ . This motivates maximum likelihood estimation (MLE), or equivalently log-loss minimization,  
 070 as a natural learning rule; see, e.g., Foster et al. (2024).<sup>3</sup> Indeed, in SFT for LLMs, one fits a model  
 071 by minimizing log-loss on prompt–completion pairs. Theoretically, prior work shows that under  
 072 the assumption  $\pi_* \in \Pi$ , MLE enjoys sharp  $O(\frac{\log |\Pi|}{m})$  convergence guarantees with  $m$  samples—  
 073 both in terms of the loss in Eq.(1) and in matching the distribution of the demonstrator  $\pi_*$  under  
 074 a certain notion of distance, ensuring both validity and coverage at the same time (Cohen et al.,  
 075 2024; Foster et al., 2024). However, for this guarantee to be meaningful,  $\log |\Pi|$  must be small. But  
 076 in learning from demonstrations, there may be many equally good demonstrations, and the policy  
 077 class modeling them can be extremely large (e.g., modeling the generative process of all optimal  
 078 demonstrators, such as graduate students producing solutions to math problems).

## 080 1.1 OUR CONTRIBUTIONS

081 In this work, departing from prior work, we instead propose an alternative to the above—we only  
 082 assume that the underlying *reward model class* has low cardinality (Section 2), i.e., the set-valued  
 083 function  $\sigma_* : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$  defining correct answers, comes from a low-cardinality model class  $\mathcal{S} \subseteq$   
 084  $(2^{\mathcal{Y}})^{\mathcal{X}}$  and the demonstrator  $\pi_*$  is supported on  $\sigma_*$ . This is a strictly weaker and more realistic  
 085 assumption: when creating QA datasets, often high-quality responses to prompts are hand-picked,  
 086 without any attempt to enumerate all valid responses or generate from them according to a fixed  
 087 specified apriori distribution. Thus, imposing any further assumptions on this generative sampling  
 088 may be overly restrictive, as done in the low-capacity policy class assumption.

- 090 • **Simple Failures of MLE Under the Low-Cardinality Reward Class (Section 3).** We  
 091 demonstrate that MLE, which is minimax optimal under the low-cardinality policy class  
 092 assumption and also enjoys convergence in distribution (Foster et al., 2024; Cohen et al.,  
 093 2024), *can fail to generalize* for low-cardinality *reward model classes* even in simple situ-  
 094 ations (Theorems 1 and 7).
- 095 • **New Learning Algorithms (Section 4).** We first present a simple voting-based learner that  
 096 ensures learnability for finite  $\mathcal{S}$  (Theorem 2). This already goes beyond MLE, though its  
 097 sample complexity can be as large as  $\Omega(|\mathcal{S}|)$  (Theorem 11). We then present a *simple* and  
 098 *novel* learner that succeeds with only  $O(\log |\mathcal{S}|)$  samples (Theorem 4), and thus optimal,  
 099 showing that exponential improvement in sample complexity is achievable with the right  
 100 inductive bias. This sample complexity has no dependence on  $|\mathcal{Y}|$  or  $|\sigma(x)|$  which can be  
 101 huge.
- 102 • **Extensions (Section 5).** (1) We generalize our optimal learner for general bounded (not  
 103 necessarily binary) reward model classes as long as the demonstrator is optimal (Sec-  
 104 tion 5.1). (2) We then also see how this algorithm is robust and can be generalized to  
 105 situations when the demonstrator  $\pi_*$  is not necessarily optimal, i.e., not supported on

106 <sup>3</sup>Even beyond imitation learning, where other forms of feedback may be available, it is routine to take the  
 107 view that  $\pi_* \in \Pi$  for some low-capacity  $\Pi$  (Yun et al., 2025; Zhan et al., 2023; Xie et al., 2024; Zhang et al.,  
 2025; Agarwal et al., 2025; Huang et al., 2024).

Rule \ Assump.	low-cardinality $\Pi$	low-cardinality $\mathcal{S}$
MLE	$\propto \log  \Pi $	May not learn
Our Learner	$\propto \log  \Pi $	$\propto \log  \mathcal{S} $

Table 1: Comparison of MLE and our learner under the low-cardinality assumptions on  $\Pi$  and  $\mathcal{S}$ . See Section 5 for the extensions of our optimal learner, for (1) general bounded reward classes (2) non-optimal demonstrator (3) pass@ $k$  error.

$\sigma_*$  (Theorem 5). In this case, one can still compete with the loss of the demonstrator  $L_{\mathcal{D}, \sigma_*}(\pi_*)$  up to a small constant blowup (e.g.,  $1.5 L_{\mathcal{D}, \sigma_*}(\pi_*)$ ). (3) Given recent interest in the relaxed notion of pass@ $k$  error (a relaxation of the objective in Eq. (1)), we show that  $\Theta(\log |\mathcal{S}| / \log k)$  samples are both sufficient for our learner and, in the worst case, necessary for any learner, yielding at most a  $\log k$  improvement in sample complexity from an information-theoretic perspective (Theorem 6).

We defer a broader discussion of our results in the context of other recent works to Section 6.

## 2 SETTING AND PROBLEM DEFINITION

We now formalize our setup. Let  $\mathcal{X}$  and  $\mathcal{Y}$  respectively be countable sets of all plausible contexts and actions. Recall the terminology introduced in Section 1 of a ground-truth support function  $\sigma_* : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$ , a marginal distribution  $\mathcal{D} \in \Delta(\mathcal{X})$  and the expert demonstrator  $\pi_* : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ . We use  $(x, y) \sim \mathcal{D} \times \pi$  to denote a joint distribution where  $x \sim \mathcal{D}$  and  $y \sim \pi(\cdot | x)$ .

**Two Function Approximations.** Note that our setup makes a distinction between the *support* of optimal answers  $\sigma_*$ , and the demonstrator’s policy  $\pi_*$ . Keeping the distinction in mind, there are two natural types of function approximations even in the realizable case:

- **Model class approximation.** There is a class of support functions  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  modeling optimal actions under different contexts. The demonstrator  $\pi_*$  is optimal and supported on some  $\sigma_* \in \mathcal{S}$  (i.e.,  $\text{supp } \pi_*(\cdot | x) \subseteq \sigma_*(x)$ ).
- **Policy class approximation.** There is a policy class  $\Pi \subseteq (\Delta(\mathcal{Y}))^{\mathcal{X}}$  such that the unknown  $\pi_* \in \Pi$ .<sup>4</sup> Note that this does not directly specify  $\sigma_*$ . It is natural to consider  $\sigma_*(x) = \sigma_{\pi_*}(x) = \text{supp } \pi_*(\cdot | x)$  as the ground-truth support function for evaluating the loss (Eq.(1)).

We note that while the two views are closely related, they differ significantly under *cardinality-based capacity control* on  $\Pi$  and  $\mathcal{S}$ . In particular, assuming small  $|\mathcal{S}|$  is much weaker, since the class of all optimal policies supported on it (defined below) can be huge, especially when  $|\sigma(x)|$  are large.

$$\Pi_{\mathcal{S}} := \bigcup_{\sigma \in \mathcal{S}} \Pi_{\sigma}, \text{ where } \Pi_{\sigma} := \{\text{Any } \pi : \mathcal{X} \rightarrow \Delta(\mathcal{Y}) \text{ s.t. } \text{supp } \pi(\cdot | x) \subseteq \sigma(x), \forall x \in \mathcal{X}\}. \quad (2)$$

On the other hand, for any policy class  $\Pi$ , the associated model class  $\mathcal{S}_{\Pi} := \bigcup_{\pi \in \Pi} \{\sigma_{\pi} \mid \sigma_{\pi}(x) = \text{supp } \pi(\cdot | x), \forall x \in \mathcal{X}\}$  has cardinality no larger than  $\Pi$  (i.e.,  $|\mathcal{S}_{\Pi}| \leq |\Pi|$ ).

In either cases, the goal is to find an  $\varepsilon$ -suboptimal policy  $\hat{\pi}(\mathcal{S}) : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ , for the loss in Eq.(1) from i.i.d. samples  $S = \{(x_i, y_i) \sim_{iid} (\mathcal{D} \times \pi_*) : i \in [m]\}$ , as formalized below.

**Definition 1** (Learning from Correct Demonstrations). *We say that the reward model class  $\mathcal{S}$  is learnable from correct demonstrations by an estimator  $\hat{\pi} : (\mathcal{X} \times \mathcal{Y})^* \rightarrow (\Delta(\mathcal{Y}))^{\mathcal{X}}$  with sample*

<sup>4</sup>We note that Foster et al. (2024) considered a more general setting than ours, aiming to compete with  $\pi_*$  (not necessarily optimal) under arbitrary bounded rewards, while assuming  $|\Pi| < \infty$ . To relax this assumption, we instead rely on the reward model class, focusing first on the special case where  $\pi_*$  is guaranteed to be optimal. We will return to the case when this may not happen in Section 5.

complexity  $m_{\mathcal{S}, \hat{\pi}} : (0, 1) \times (0, 1) \rightarrow \mathbb{N}$ , if for any  $\varepsilon, \delta \in (0, 1)$ , for any sample size  $m \geq m_{\mathcal{S}, \hat{\pi}}(\varepsilon, \delta)$ , for any  $\mathcal{D}, \sigma_*, \pi_*$ , we have

$$\mathbb{P}_{S \sim (\mathcal{D} \times \pi_*)^m} [L_{\mathcal{D}, \sigma_*}(\hat{\pi}(S)) \leq \varepsilon] \geq 1 - \delta.$$

The learnability for  $\Pi$  follows exactly the same definition after replacing  $\sigma_*(x) = \text{supp } \pi_*(\cdot | x)$  for measuring the loss.

### 3 WHY MLE IS GOOD FOR LOW-CARDINALITY $\Pi$ BUT FAILS FOR $\mathcal{S}$ ?

Given a policy class  $\Pi \subseteq (\Delta(\mathcal{Y}))^{\mathcal{X}}$ , the (conditional) Maximum Likelihood Estimator (MLE) (or the log-loss minimizer) is defined as

$$\text{MLE}_{\Pi}(S) = \arg \max_{\pi \in \Pi} \prod_{i=1}^m \pi(y_i | x_i) = \arg \min_{\pi \in \Pi} - \sum_{i=1}^m \log \pi(y_i | x_i). \quad (\text{MLE})$$

**Low-cardinality  $\Pi$ .** We first discuss the case when we assume  $\Pi$  has a low cardinality. In this case, MLE is actually *minimax optimal* with respect to  $|\Pi|$ . The proposition below formalizes  $O(\frac{\log |\Pi|}{m})$  convergence rate (in similar spirit to Cohen et al. (2024); Foster et al. (2021)) for our setup.

**Proposition 1.** For any finite  $\Pi \subseteq (\Delta(\mathcal{Y}))^{\mathcal{X}}$ , for any  $\mathcal{D}$  and  $\pi_* \in \Pi$ , with probability at least  $1 - \delta$  over  $S \sim (\mathcal{D} \times \pi_*)^m$ , any  $\hat{\pi}_{\text{mle}}(S) \in \text{MLE}_{\Pi}(S)$  enjoys the following guarantee:

$$L_{\mathcal{D}, \sigma_{\pi_*}}(\hat{\pi}_{\text{mle}}(S)) \leq D_{\text{H}}^2(\hat{\pi}_{\text{mle}}(S), \pi_*) \leq \frac{6 \log(2|\Pi|/\delta)}{m}.$$

Thus  $\Pi$  is learnable with  $\hat{\pi}_{\text{mle}}$  (cf. Definition 1) with sample complexity:  $m_{\Pi, \hat{\pi}_{\text{mle}}}(\varepsilon, \delta) = \frac{6}{\varepsilon} \log(2|\Pi|/\delta)$ .

In the above,  $D_{\text{H}}^2(\mathbb{P}, \mathbb{Q}) := \sum_{z \in \mathcal{Z}} (\sqrt{\mathbb{P}(z)} - \sqrt{\mathbb{Q}(z)})^2$  is the squared Hellinger distance between distributions  $\mathbb{P}$  and  $\mathbb{Q}$  over a discrete domain  $\mathcal{Z}$ . The proof (in Appendix C) uses the ideas from Foster et al. (2024); one can first use the standard guarantees in density estimation to establish the convergence in the squared Hellinger distance for  $D_{\text{H}}^2(\hat{\pi}_{\text{mle}}(S), \pi_*)$ , followed by controlling the loss (Eq. (1)) in terms of  $D_{\text{H}}^2(\hat{\pi}_{\text{mle}}(S), \pi_*)$ . However, this guarantee is only meaningful when  $\log |\Pi|$  is small. In the context of LLMs, this assumption may be somewhat unrealistic, since post-SFT models are not as perfect as this guarantee suggests—namely, matching the distribution of the demonstrator (see Section 6 for further discussion).

**Low-Cardinality  $\mathcal{S}$ : Simple Failures of MLE over Natural Policy Classes.** We now ask what can we say about the performance of MLE when only  $\mathcal{S}$  is finite. It is important to note that our problem (cf. Definition 1) is only specified in terms of the model class  $\mathcal{S}$ . To use likelihood maximization, we need to consider a policy class modeling the underlying generative process. In what follows, we describe two natural ways of defining the policy class based on the model class  $\mathcal{S}$ .

First, given the promise that  $\pi_*$  is supported on some  $\sigma_* \in \mathcal{S}$ , one natural choice is to consider the class of  $\Pi_{\mathcal{S}}$  from Eq. (2) of all policies supported on  $\mathcal{S}$ . Note that  $|\Pi_{\mathcal{S}}|$  may be infinite now making the guarantee in Proposition 1 vacuous. We shall see that MLE fails to produce even correct responses. It is simple to observe that, for any  $x$  observed in the training set,  $\text{MLE}_{\Pi_{\mathcal{S}}}(S)$  always outputs an action according to the empirical distribution of observed  $y_i$ 's for  $x_i = x$ . However, on unseen contexts, it may output any action  $y \in \sigma(x)$  for some consistent  $\sigma$ . Thus, this MLE on unseen contexts is exactly identical to consistency-based learner:

$$\hat{\pi}_{\text{Con}}(S)(x) = \hat{y} \text{ for some } \hat{y} \in \sigma(x) \text{ where } \sigma \in V(S) := \{\sigma \in \mathcal{S} : y_i \in \sigma(x_i), \forall (x_i, y_i) \in S\}. \quad (3)$$

While such a consistency-based learner is known to be optimal (for finite classes) in supervised learning problems like binary or multiclass classification, it fails for our problem and thus, MLE over  $\Pi_{\mathcal{S}}$  also fails. Consider  $\mathcal{S} = \{\sigma_0, \sigma_{01}\}$ , where  $\sigma_0(x) = \{0\}$  and  $\sigma_{01}(x) = \{0, 1\}$  for all  $x$ . If the true hypothesis is  $\sigma_* = \sigma_0$ , then all observed labels are 0. However,  $\sigma_{01}$  also remains consistent, and thus  $\text{MLE}_{\Pi_{\mathcal{S}}}(S)$  may output 1 at test time—*failing to generalize and go beyond memorization*. Thus, for an input distribution  $\mathcal{D}$  such that the *missing mass* (i.e., unobserved contexts) is arbitrarily close to 1, we get the following failure.

**Theorem 1** (Failure of MLE over  $\Pi_S$ ). *There exists  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  with  $|\mathcal{S}| = |\mathcal{Y}| = 2$  and  $\mathcal{X} = \mathbb{N}$  and a choice of  $(\sigma_*, \pi_*)$  such that for every sample size  $m$  and  $\gamma \in (0, 1)$ , there exists a marginal distribution  $\mathcal{D}$  such that for  $S \sim (\mathcal{D} \times \pi_*)^m$ , some  $\hat{\pi}_{\text{mle}}(S) \in \text{MLE}_{\Pi_S}(S)$  has the following guarantee:*

$$\mathbb{P}_{S \sim (\mathcal{D} \times \pi_*)^m} (L_{\mathcal{D}, \sigma_*}(\hat{\pi}_{\text{mle}}(S)) \geq 1 - \gamma) = 1.$$

In the above example, MLE over  $\Pi_S$  essentially *overfits*: it achieves zero error on the training data but fails to generalize to unseen data. A simple cause of it is just the fact that the induced policy class  $\Pi_S$  is too rich. As a remedy, we may consider a restricted policy class  $\Pi_{\text{unif}, \mathcal{S}}$  with size  $|\mathcal{S}|$ :

$$\Pi_{\text{unif}, \mathcal{S}} := \{\pi_{\text{unif}, \sigma} : \sigma \in \mathcal{S}\} \text{ where } \pi_{\text{unif}, \sigma}(\cdot | x) = \text{Unif}(\sigma(x)).$$

This is a natural candidate for a class with restricted capacity, where the learner only knows  $\pi_*$  is supported on some  $\sigma_* \in \mathcal{S}$ . However, the optimal demonstrator’s policy  $\pi_*$  need not coincide with the canonical choice  $\pi_{\text{unif}, \sigma_*}$ , although it is perfectly supported on  $\sigma_*$ . So  $\Pi_{\text{unif}, \mathcal{S}}$  is actually misspecified, i.e., it may not contain  $\pi_*$ . This mismatch suffices to make the MLE fail again on the restricted capacity class  $\Pi_{\text{unif}, \mathcal{S}}$ ; see Theorem 7 and its proof in Appendix D.

**Remark 3.1** (MLE achieves overlap). *An interesting property of MLE over the restricted class  $\Pi_{\text{unif}, \mathcal{S}} = \{\pi_{\text{unif}, \sigma} : \sigma \in \mathcal{S}\}$  is that it still achieves an overlap with ground-truth correct answers for any sample size  $m \geq \frac{1}{\varepsilon} (\log |\mathcal{S}| + \log(1/\delta))$ . For any  $\hat{\pi}_{\text{mle}}(S) \in \text{MLE}_{\Pi_{\text{unif}, \mathcal{S}}}(S)$ , we have*

$$\mathbb{P}_{S \sim (\mathcal{D} \times \pi_*)^m} (\mathbb{P}_{x \sim \mathcal{D}}[\text{supp } \hat{\pi}_{\text{mle}}(S)(\cdot | x) \cap \sigma_*(x) = \emptyset] \leq \varepsilon) \geq 1 - \delta.$$

*Thus, its predictions overlap with the ground-truth responses on all but an  $\varepsilon$ -fraction of inputs, though it may still generate answers outside the support with some non-trivial probability. See Appendix D.1 for further discussion on this.*

## 4 LEARNING ALGORITHMS

In this section, we present our main result: a sample-efficient learner (Section 4.2). En route to this, we first see how natural approaches, though sufficient to ensure learnability under finite  $\mathcal{S}$ , may have a disappointing linear dependence on  $|\mathcal{S}|$  in the sample complexity (Section 4.1).

### 4.1 WARM-UP

Our rule is simple: output from the common intersection of consistent hypotheses if it is non-empty, and otherwise output an arbitrary  $y$  within the support of some consistent  $\sigma$ . This suffices to ensure learnability.

**Input:** Sample  $S = \{(x_i, y_i) : i \in [m]\}$  and a finite model class  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$ .

- Let  $V(S) := \{\sigma \in \mathcal{S} : y_i \in \sigma(x_i), \forall (x_i, y_i) \in S\}$
- Return the predictor  $\text{COMMON-INTERSECTION}(S) = \hat{\pi}_{\text{CI}}(S) : \mathcal{X} \rightarrow \mathcal{Y}$  as follows:

$$\hat{\pi}_{\text{CI}}(S)(x) = \begin{cases} y \in \bigcap_{\sigma \in V(S)} \sigma(x), & \text{if } \bigcap_{\sigma \in V(S)} \sigma(x) \neq \emptyset; \\ \text{arbitrary } y & \text{otherwise.} \end{cases}$$

**Theorem 2.** *Any finite  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  is learnable (cf. Definition 1) using the rule  $\hat{\pi}_{\text{CI}}$  with sample complexity  $m_{\mathcal{S}, \hat{\pi}_{\text{CI}}}(\varepsilon, \delta) = \varepsilon^{-1} |\mathcal{S}| (\log |\mathcal{S}| + \log(1/\delta))$ .*

The proof is in Appendix E.2. The dependence on  $|\mathcal{S}|$  in Theorem 2 is  $O(|\mathcal{S}| \log |\mathcal{S}|)$ , in contrast to the logarithmic dependence in standard supervised learning. We show this dependence is tight in the worst case (up to a log factor) for this rule, and even for a seemingly stronger variant that outputs by majority vote over consistent hypotheses:  $\hat{\pi}_{\text{Maj}}(S)(x) = \arg \max_{y \in \mathcal{Y}} |\{\sigma \in V(S) : y \in \sigma(x)\}|$ . See Theorem 11 in Appendix E.3.

### 4.2 SAMPLE EFFICIENT LEARNER

We now design a sample-efficient learner that achieves optimal *logarithmic dependence* on  $|\mathcal{S}|$ —even under the finite cardinality assumption on  $\mathcal{S}$  and without assuming anything further about the

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

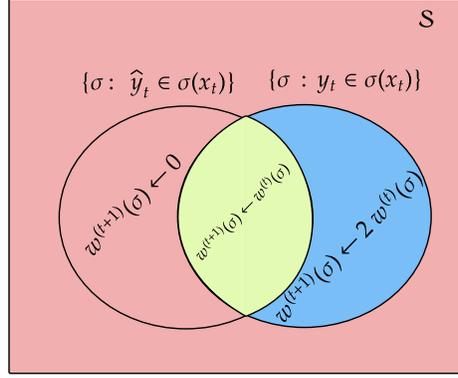


Figure 1: A visualization of the update rule of Algorithm 1 during  $t^{\text{th}}$  round. The weight of the hypotheses in the red, green, and blue regions are respectively set to zero, unchanged, and doubled.

optimal demonstrator’s policy, we obtain sample complexity  $\propto \log |\mathcal{S}|$ , independent of  $|\mathcal{X}|$ ,  $|\mathcal{Y}|$ , or  $\sup_{\sigma, x} |\sigma(x)|$ . To achieve this, we first turn our attention to the even more challenging online version of the problem. The statistical estimator will be designed by doing online-to-batch conversion.

**Online Version.** The adversary chooses  $\sigma_* \in \mathcal{S}$ . In each round  $t$ :

- The adversary chooses  $x_t \in \mathcal{X}$ . The learner predicts  $\hat{y}_t \in \mathcal{Y}$ .
- The adversary shows some  $y_t \in \sigma_*(x_t)$ .  
(Importantly, the feedback does not inform the learner whether  $\hat{y}_t$  was a mistake or not.)

The algorithm maintains weight function  $w^{(t)} : \mathcal{S} \rightarrow \mathbb{R}$  in each round.

---

**Algorithm 1** MISTAKE-UNAWARE-WEIGHT-UPDATE

---

**Input:** A finite support class  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$ .

- Initialize  $w^{(1)}(\sigma) = 1$  for all  $\sigma \in \mathcal{S}$ .
- In every round, receiving  $x_t$ :
  1. Output  $\hat{y}_t = \arg \max_{y \in \mathcal{Y}} \sum_{\sigma \in \mathcal{S}} w^{(t)}(\sigma) \mathbf{1}\{y \in \sigma(x_t)\}$ .
  2. On receiving  $y_t$ , update the weights

$$w^{(t+1)}(\sigma) \leftarrow \begin{cases} 0 & \text{for all } \sigma \in \mathcal{S} \text{ with } y_t \notin \sigma(x_t); \\ w^{(t)}(\sigma) & \text{for all } \sigma \in \mathcal{S} \text{ with both } y_t, \hat{y}_t \in \sigma(x_t); \\ 2w^{(t)}(\sigma) & \text{for all } \sigma \in \mathcal{S} \text{ with } y_t \in \sigma(x_t) \text{ and } \hat{y}_t \notin \sigma(x_t); \end{cases}$$


---

The key distinction from the standard *Majority* is in the final step of upweighting certain hypotheses: even without knowing whether  $\hat{y}_t$  was a mistake, the algorithm doubles the weight of all hypotheses that exclude  $\hat{y}_t$  but include the observed response  $y_t$  (see also Figure 1).

**Theorem 3** (Online Guarantee). *On any sequence  $((x_t, y_t))_{t \in \mathbb{N}}$  realizable by some  $\sigma_* \in \mathcal{S}$ , Algorithm 1 makes at most  $\log_2 |\mathcal{S}|$  mistakes.*

*Proof.* Letting  $W_{t+1} = \sum_{\sigma \in \mathcal{S}} w^{(t+1)}(\sigma)$  be the total weight in the system after completion of  $t$  rounds, we first note that the sequence  $\{W_t\}_t$  is non-increasing. This is because of the property of the algorithm that, during every round  $t$ , the weight added to the system is at most the weight eliminated from it. Formally,

$$W_{t+1} = \sum_{y_t \in \sigma(x_t) \text{ and } \hat{y}_t \notin \sigma(x_t)} 2w^{(t)}(\sigma) + \sum_{y_t, \hat{y}_t \in \sigma(x_t)} w^{(t)}(\sigma) \leq \sum_{y_t \in \sigma(x_t) \text{ or } \hat{y}_t \notin \sigma(x_t)} w^{(t)}(\sigma) \leq W_t,$$

where the first inequality follows from the property of the algorithm that it always chooses  $\hat{y}_t = \arg \max_{y \in \mathcal{Y}} \sum_{\sigma \in \mathcal{S}} w^{(t)}(\sigma) \mathbf{1}\{y \in \sigma(x_t)\}$  (see also Figure 1). Now if the algorithm made  $M$  mistake on a realizable sequence for some  $\sigma_* \in \mathcal{S}$  at the end some  $t$  number of rounds, then it must be that

$$w^{(t+1)}(\sigma_*) = 2^M \leq W_{t+1} \leq W_1 = |\mathcal{S}|, \text{ which implies } M \leq \log_2 |\mathcal{S}|.$$

□

We use the online-to-batch conversion of outputting a randomized predictor based on a random stopping time (for every test context).

**Input:** Sample  $S = \{(x_i, y_i) : i \in [m]\}$  and a finite model class  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$ .

- Run Algorithm 1 once over  $S$ , and record  $w^{(t)}$  before each round. Define

$$\hat{\pi}_t(x) = \arg \max_{y \in \mathcal{Y}} \sum_{\sigma \in \mathcal{S}} w^{(t)}(\sigma) \mathbf{1}\{y \in \sigma(x)\}. \quad (4)$$

- On a test  $x \in \mathcal{X}$ , sample  $I \sim \text{Unif}\{1, \dots, m\}$ , and return  $\hat{\pi}_{\text{o2b}}(S)(x) := \hat{\pi}_I(x)$ , i.e.,

$$\hat{\pi}_{\text{o2b}}(S)(x) = \frac{1}{m} \sum_{t=1}^m \hat{\pi}_t(x). \quad (5)$$

**Theorem 4** (Statistical Guarantee). *For any finite  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$ , the estimator  $\hat{\pi}_{\text{o2b}}$  in Eq. (5) achieves the following guarantee for any joint distribution  $(\mathcal{D} \times \pi_*)$  supported on  $\sigma_* \in \mathcal{S}$ :*

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}, \sigma_*}(\hat{\pi}_{\text{o2b}}(S))] \leq \frac{\log_2 |\mathcal{S}|}{m},$$

and, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$L_{\mathcal{D}, \sigma_*}(\hat{\pi}_{\text{o2b}}(S)) \leq \frac{1 + 2 \log |\mathcal{S}| + 12 \log \left( \frac{\log m}{\delta} \right)}{m}.$$

This implies that  $\mathcal{S}$  is learnable (cf. Definition 1) using the estimator  $\hat{\pi}_{\text{o2b}}$  with sample complexity  $m_{\mathcal{S}, \hat{\pi}_{\text{o2b}}}(\varepsilon, \delta) = O(\varepsilon^{-1}(\log |\mathcal{S}| + \log(1/\varepsilon\delta)))$ .

See Appendix E for the proof. It follows from concentration for martingale difference sequences. To obtain a sharper dependence on  $1/\varepsilon$  than in standard analyses (Cesa-Bianchi et al., 2004; Tewari & Kakade, 2008), we need a tighter control on the sum of martingale difference sequence, for which we apply Freedman’s inequality (Lemma 2) with some algebraic manipulations.

## 5 EXTENSIONS

We now provide generalization of our optimal learner to important variants of our main setting.

### 5.1 GENERAL BOUNDED REWARD CLASSES

Consider a finite reward class  $\mathcal{R}$ , consisting of functions  $r : (\mathcal{X} \times \mathcal{Y}) \rightarrow [0, 1]$ . The demonstrator  $\pi_*$  is still optimal with respect to some unknown  $r_* \in \mathcal{R}$ , and is supported on

$$\sigma_*(x) := \arg \max_{y \in \mathcal{Y}} r_*(x, y).$$

Both our optimal learner (and its pass@ $k$  variant, discussed in Section 5.3) can be implemented on the corresponding support class

$$\mathcal{S}_{\mathcal{R}} = \left\{ \sigma_r \mid \sigma_r(x) = \arg \max_{y \in \mathcal{Y}} r(x, y), \forall x \in \mathcal{X} \right\}.$$

Since  $\pi_*$  is optimal, the reward sub-optimality of our learner can be directly controlled by the loss function in Eq. (1). In particular, we obtain the guarantee that the value of our predictor policy  $V(\hat{\pi}) \geq V(\pi_*) - \varepsilon$  with high probability, at a sample complexity determined by the size of  $\mathcal{S}_{\mathcal{R}}$ , which satisfies  $|\mathcal{S}_{\mathcal{R}}| \leq |\mathcal{R}|$ .

## 5.2 LEARNING FROM ARBITRARY DEMONSTRATOR

We now discuss what happens if  $\pi_*$  is not necessarily optimal (again consider only binary rewards for simplicity). Note that in our Algorithm 1, we assign zero weight to all hypotheses that do not contain  $y_t$ —here we critically relied on the promise that the demonstrator is optimal. However, this strategy may be detrimental when it is not the case. In order to find an estimator, one needs to minimize a certain notion of *regret* of outputting bad responses with respect to any potential  $\sigma \in \mathcal{S}$  in the online case, simply from demonstrations. To do so, we again output  $\hat{y}_t$  that is most voted, but upon receiving  $y_t$ , do a “softer” update:

$$w^{(t+1)}(\sigma) \leftarrow w^{(t)}(\sigma) \cdot \alpha^{\mathbb{1}[\hat{y}_t \notin \sigma(x_t)]} \cdot \beta^{-\mathbb{1}[y_t \notin \sigma(x_t)]}, \text{ for some } \alpha, \beta \geq 1.$$

Then  $(\alpha, \beta)$  are set so that, it ensures a certain *regret minimization* in the online case. Applying the standard online-to-batch conversion yields the following guarantee:

**Theorem 5** (Learning from arbitrary demonstrator). *For any finite class  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$ , there is an estimator  $\hat{\pi}$  such that for any unknown joint distribution  $\mathcal{D} \times \pi_*$  (where  $\pi_*$  is not necessarily optimal) and reference  $\sigma_* \in \mathcal{S}$ , with probability  $1 - \delta$  over  $S \sim (\mathcal{D} \times \pi_*)^m$ , we have*

$$L_{\mathcal{D}, \sigma_*}(\hat{\pi}(S)) \leq 1.41 L_{\mathcal{D}, \sigma_*}(\pi_*) + \frac{1 + 6 \log_2 |\mathcal{S}| + 40 \log \left( \frac{2 \log 6m}{\delta} \right)}{m}.$$

See Appendix F for the proof and details of the algorithm.

## 5.3 pass@k ERROR MINIMIZATION

In modern practice, pass@k accuracy is often used as a benchmark (e.g., Chen et al. (2021); Orlanski et al. (2022); Dalal et al. (2025)). This relaxes the original goal by allowing a (stochastic) policy  $\hat{\mu} : \mathcal{X} \rightarrow \Delta(\mathcal{Y}^k)$  to output  $k$  answers. The pass@k loss is defined as

$$L_{\mathcal{D}, \sigma_*}(\hat{\mu}) = \mathbb{E}_{x \sim \mathcal{D}}, \mathbb{E}_{\mathbf{y}=(y^{(1)}, \dots, y^{(k)}) \sim \hat{\mu}(\cdot|x)} \left[ \mathbb{1}\{y^{(i)} \notin \sigma_*(x); \forall i \in [k]\} \right]. \quad (6)$$

The policy  $\hat{\mu}$  allows for any joint distribution over the set of  $k$  responses, which also includes adaptive sampling and need not be a product distribution (i.e., repeated independent sampling from a policy  $\hat{\pi} : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ ).

Our goal is to understand how the parameter  $k$  impacts sample complexity. We show that, under this relaxation, the worst-case improvement is only a  $\log k$  factor in terms of the cardinality parameter. Again we return to the case where the demonstrator  $\pi_*$  is optimal. We extend our Algorithm 1 to this setting: the learner outputs a greedily picked set of responses  $\{\hat{y}_t^{(1)}, \dots, \hat{y}_t^{(k)}\}$  that covers most “uncovered mass” each time. Upon receiving  $y_t$ , the update is

$$w^{(t+1)}(\sigma) \leftarrow (k+1) w^{(t)}(\sigma),$$

for all hypotheses that contain  $y_t$  but neither of  $\{\hat{y}_t^{(1)}, \dots, \hat{y}_t^{(k)}\}$ . (The other updates are same.)

This leads to a sharp mistake bound of  $\log_{k+1} |\mathcal{S}|$  in the online case, which then translates to a sample complexity in the statistical case, by the same online-to-batch conversion. See Appendix G for the details and formalization of the following informal result.

**Theorem 6** (Informal: pass@k error). *The minimax mistake bound in the online case as well as sample complexity bound in the statistical case are  $\Theta(\log_k |\mathcal{S}|)$  for the family of finite classes of fixed size, when learning from optimal demonstrator.*

## 6 DISCUSSION, BROADER CONCLUSIONS, AND FUTURE WORK

We studied the problem of learning to answer from correct demonstrations (aka imitation learning in contextual bandits from optimal demonstrator). Our goal was to move away from the common assumption that the demonstrator lies in a low-cardinality policy class  $\Pi$ , where MLE guarantees strong convergence to distribution of the demonstrator, under a certain notion of distance, thereby it can match the demonstrator also in performance. However, supervised fine-tuned models trained

via log-loss are not perfect (Ji et al., 2023; OpenAI, 2023)—they often fail to produce even correct responses. Even theoretically, the above results are in contrast to the alternative perspective of (Kalai & Vempala, 2024; Kalai et al., 2025), which argues that calibrated models (a coverage related notion) must necessarily hallucinate (i.e., incur high loss under our objective). The conclusion is that these guarantees, under the assumption  $\pi_* \in \Pi$  with a low-capacity  $\Pi$  can be considered too crude to capture the current situations with LLMs (see also Section 6.2 of Foster et al. (2024)). Instead, we proposed a low-cardinality *reward model* class, which we argued is strictly weaker and more realistic assumption for current LLMs. This shift lets us expose simple failure cases of likelihood maximization and motivates new estimators that go beyond it in order to learn.

**Broader Implications.** Our new learner highlights that demonstrations available in SFT can carry much more information about the reward structure than what likelihood maximization—the standard approach for LLM fine-tuning—can extract. Another implication of our analysis concerns *hallucinations*. Our explanation is new: they arise naturally in the *prompted* question–answering scenario when learning is carried out via likelihood maximization. This departs from other recent theoretical perspectives (Kalai & Vempala, 2024; Kalavasis et al., 2025; Kalai et al., 2025). The closest to our work is (Kalai et al., 2025, Section 3.2), which matches our contextual bandit setting and similarly shows that hallucinations arise for statistical reasons. However, our failure cases are simpler and tied directly to MLE.

While their work attempts to bound hallucination rates in terms of a certain notion of “calibration”—encouraged by log-loss objective—our work goes further and asks: if *validity* and *coverage* are in tension, can we at least ensure correctness (reward maximization), when coverage may not be the primary goal anyway. We show that indeed it can be possible within our framework. There may, of course, be limited situations where coverage is important for fairness-related reasons, but even in such cases, the model should perhaps give a more comprehensive response indicating its uncertainty (e.g., (Kirchhof et al., 2025)). Thus the focus should perhaps be on validity/goodness/usefulness, with appropriate notion of utility, rather than attempting to learn the more difficult task of matching the distribution, cf. the quote from Vladimir Vapnik:

“When solving a problem of interest, do not solve a more general problem as an intermediate step.”  
—Vladimir Vapnik

**Continuous Classes.** This motivates the question what classes  $\mathcal{S}$  should we use. In practice, we generally use parametric function classes, and thus a natural question is whether we can handle natural parametric continuous classes. For example, given a parameterized function class  $h_w : (\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}$  assigning real-valued scores, we can consider the associated support function  $\sigma_w : x \mapsto \{y : h_w(x, y) \geq 0\}$ . Natural choices for  $h_w$  include linear classes  $h_w(x, y) = \langle w, \phi(x, y) \rangle$ , low-rank score matrices  $W \in \mathbb{R}^{|\mathcal{X}| \times |\mathcal{Y}|}$  with  $h_w(x, y) = W_{x,y}$ , or even transformers  $f_w(x, y) \in \mathbb{R}$  assigning scores to  $(x, y)$ . Do the gaps we observed also manifest in these natural classes? What is the sample complexity of offline imitation learning for them? Can we characterize the complexity of the support class  $\mathcal{S}_{\mathcal{F}} = \{\sigma_f : x \mapsto \{y : f(x, y) \geq 0\} \mid f \in \mathcal{F}\}$  in terms of abstract properties of  $\mathcal{F}$ , both for general classes and for simple parametric ones? What algorithms achieve this? One option is discretization and running our algorithm, yielding a log-cardinality bound statistically, but the algorithm’s memory is exponential in the number of parameters and thus impractical. It is an important question to see if there are exact, approximate, or heuristic methods for learning with continuous classes that scale more favorably with the number of parameters (i.e., model size), and whether the inductive biases about the reward structure can be directly incorporated when learning from demonstrations, without modeling the policy of the demonstrator.

**Technical Open Questions.** Even for finite classes, important technical and conceptual gaps remain in our understanding. Our main learning rule in Section 4.2 is *randomized* and *improper* (see definition in Appendix A), whereas the simpler rule COMMON-INTERSECTION in Section 4.1, which is deterministic and proper, faces an  $\Omega(|\mathcal{S}|)$  barrier. Is this gap fundamental? Can we design simpler rules even for finite classes that achieves  $O(\log |\mathcal{S}|)$  sample complexity? Moreover, when the demonstrator is suboptimal, can we compete with the demonstrator’s loss  $L_{\mathcal{D}, \sigma_*}(\pi_*)$  without incurring blow-up, either with our rule or any other?

## REFERENCES

- 486  
487  
488 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic  
489 bandits. *Advances in neural information processing systems*, 24, 2011.
- 490 Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert Schapire. Taming  
491 the monster: A fast and simple algorithm for contextual bandits. In *International conference on*  
492 *machine learning*, pp. 1638–1646. PMLR, 2014.
- 493 Alekh Agarwal, Christoph Dann, and Teodor V Marinov. Design considerations in offline  
494 preference-based rl. *arXiv preprint arXiv:2502.06861*, 2025.
- 495  
496 DeepSeek AI. Deepseek models overview. <https://www.deepseek.com>, 2025.
- 497 Anthropic. Claude 3 model family. <https://www.anthropic.com/news/claude-3>,  
498 2024.
- 499  
500 Léon Bottou, Jonas Peters, Joaquin Quiñero-Candela, Denis X Charles, D Max Chickering, Elon  
501 Portugaly, Dipankar Ray, Patrice Simard, and Ed Snelson. Counterfactual reasoning and learning  
502 systems. *arXiv preprint arXiv:1209.2355*, 2012.
- 503 Florian Le Bronnec, Alexandre Verine, Benjamin Negrevergne, Yann Chevaleyre, and Alexandre  
504 Allauzen. Exploring precision and recall to assess the quality and diversity of llms. *arXiv preprint*  
505 *arXiv:2402.10693*, 2024.
- 506  
507 Nicolo Cesa-Bianchi, Alex Conconi, and Claudio Gentile. On the generalization ability of on-line  
508 learning algorithms. *IEEE Transactions on Information Theory*, 50(9):2050–2057, 2004.
- 509 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared  
510 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Sandhini Puri,  
511 Gretchen Krueger, Mikhail Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan,  
512 Scott Gray, Nick Ryder, Michael Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian,  
513 Clemens Winter, Phil Tillet, Felipe Such, Dave Cummings, Matthias Plappert, Fotios Chantzis,  
514 Elizabeth Barnes, Ariel Herbert-Voss, William H. Guss, Alex Nichol, Igor Paino, Nikolas Tezak,  
515 Jie Tang, Igor Babuschkin, Suchir Balaji, Szymon Jain, William Saunders, Jesse Hsu, Ryan Crow-  
516 der, Anjali Srinivasan, Andrew Ho, Tom Gray, Nicholas Ryder, Dario Amodei, Ilya Sutskever,  
517 and Wojciech Zaremba. Evaluating large language models trained on code. *arXiv preprint*  
518 *arXiv:2107.03374*, 2021.
- 519 Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff func-  
520 tions. In *Proceedings of the fourteenth international conference on artificial intelligence and*  
521 *statistics*, pp. 208–214. JMLR Workshop and Conference Proceedings, 2011.
- 522 Lee Cohen, Yishay Mansour, Shay Moran, and Han Shao. Probably approximately precision and  
523 recall learning. *arXiv preprint arXiv:2411.13029*, 2024.
- 524  
525 Uri Dalal, Meirav Segal, Zvika Ben-Haim, Dan Lahav, and Omer Nevo. Leveraging llm inconsis-  
526 tency to boost pass@k performance. *arXiv preprint arXiv:2505.12938*, 2025.
- 527 Xingyu Dang, Christina Baek, Kaiyue Wen, Zico Kolter, and Aditi Raghunathan. Weight ensem-  
528 bling improves reasoning in language models. *arXiv preprint arXiv:2504.10478*, 2025. URL  
529 <https://arxiv.org/abs/2504.10478>.
- 530  
531 Google DeepMind. Introducing gemini: Our largest and most capable ai model.  
532 <https://blog.google/technology/ai/google-gemini-ai/>, 2023.
- 533 Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning. *arXiv*  
534 *preprint arXiv:1103.4601*, 2011.
- 535  
536 Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of  
537 interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.
- 538  
539 Dylan J Foster, Adam Block, and Dipendra Misra. Is behavior cloning all you need? understanding  
horizon in imitation learning. *Advances in Neural Information Processing Systems*, 37:120602–  
120666, 2024.

- 540 Audrey Huang, Wenhao Zhan, Tengyang Xie, Jason D Lee, Wen Sun, Akshay Krishnamurthy, and  
541 Dylan J Foster. Correcting the mythos of kl-regularization: Direct alignment without overopti-  
542 mization via chi-squared preference optimization. *arXiv preprint arXiv:2407.13399*, 2024.
- 543 Audrey Huang, Adam Block, Qinghua Liu, Nan Jiang, Akshay Krishnamurthy, and Dylan J Foster.  
544 Is best-of-n the best of them? coverage, scaling, and optimality in inference-time alignment. *arXiv*  
545 *preprint arXiv:2503.21878*, 2025.
- 546  
547 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Eric Ishii, Yejin Bang, An-  
548 drea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM*  
549 *Computing Surveys*, 55(12):1–38, 2023.
- 550 Adam Tauman Kalai and Santosh S Vempala. Calibrated language models must hallucinate. In  
551 *Proceedings of the 56th Annual ACM Symposium on Theory of Computing*, pp. 160–171, 2024.
- 552  
553 Adam Tauman Kalai, Ofir Nachum, Santosh S. Vempala, and Edwin Zhang. Why language models  
554 hallucinate. Technical report, OpenAI, September 2025.
- 555 Alkis Kalavasis, Anay Mehrotra, and Grigoris Velegkas. On the limits of language generation:  
556 Trade-offs between hallucination and mode-collapse. In *Proceedings of the 57th Annual ACM*  
557 *Symposium on Theory of Computing*, pp. 1732–1743, 2025.
- 558  
559 Michael Kirchhof, Luca Füger, Adam Goliński, Eeshan Gunesh Dhekane, Arno Blaas, and Sinead  
560 Williamson. Self-reflective uncertainties: Do llms know their internal answer distribution? *arXiv*  
561 *preprint arXiv:2505.20295*, 2025.
- 562 John Langford and Tong Zhang. The epoch-greedy algorithm for multi-armed bandits with side  
563 information. *Advances in neural information processing systems*, 20, 2007.
- 564  
565 Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 566 Gen Li, Laixi Shi, Yuxin Chen, Yuantao Gu, and Yuejie Chi. Breaking the sample complexity barrier  
567 to regret-optimal model-free reinforcement learning. *Advances in Neural Information Processing*  
568 *Systems*, 34:17762–17776, 2021.
- 569  
570 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- 571 Gabriel Orlanski, Seonhye Yang, and Michael Healy. Evaluating how fine-tuning on bimodal data  
572 affects code generation. *arXiv preprint arXiv:2211.07842*, 2022.
- 573  
574 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong  
575 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to fol-  
576 low instructions with human feedback. *Advances in neural information processing systems*, 35:  
577 27730–27744, 2022.
- 578 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea  
579 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*  
580 *in neural information processing systems*, 36:53728–53741, 2023.
- 581 Nived Rajaraman, Lin Yang, Jiantao Jiao, and Kannan Ramchandran. Toward the fundamental limits  
582 of imitation learning. *Advances in Neural Information Processing Systems*, 33:2914–2924, 2020.
- 583  
584 Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. Bridging offline rein-  
585 forcement learning and imitation learning: A tale of pessimism. *Advances in Neural Information*  
586 *Processing Systems*, 34:11702–11716, 2021.
- 587 Shai Shalev-Shwartz and Shai Ben-David. *Understanding Machine Learning: From Theory to*  
588 *Algorithms*. Cambridge University Press, 2014.
- 589  
590 Yuda Song, Julia Kempe, and Remi Munos. Outcome-based exploration for llm reasoning. *arXiv*  
591 *preprint arXiv:2509.06941*, 2025. doi: 10.48550/arXiv.2509.06941.
- 592  
593 Adith Swaminathan and Thorsten Joachims. Counterfactual risk minimization: Learning from  
logged bandit feedback. In *International conference on machine learning*, pp. 814–823. PMLR,  
2015.

- 594 Ambuj Tewari and Sham Kakade. Online-to-batch conversions. Lecture notes,  
595 CMSC 35900: Learning Theory, Toyota Technological Institute at Chicago, 2008.  
596 <https://home.ttic.edu/~tewari/lectures/lecture13.pdf>.  
597
- 598 Hugo Touvron et al. Llama: Open and efficient foundation language models. *arXiv preprint*  
599 *arXiv:2302.13971*, 2023.
- 600 Michael Woodroffe. A one-armed bandit problem with a concomitant variable. *Journal of the*  
601 *American Statistical Association*, 74(368):799–806, 1979.  
602
- 603 Fang Wu, Weihao Xuan, Ximing Lu, Zaid Harchaoui, and Yejin Choi. The invisible leash:  
604 Why rlvr may not escape its origin. *arXiv preprint arXiv:2507.14843*, 2025. URL  
605 <https://arxiv.org/abs/2507.14843>.
- 606 Tengyang Xie, Dylan J Foster, Akshay Krishnamurthy, Corby Rosset, Ahmed Awadallah, and  
607 Alexander Rakhlin. Exploratory preference optimization: Harnessing implicit q\*-approximation  
608 for sample-efficient rlhf. *arXiv preprint arXiv:2405.21046*, 2024.  
609
- 610 Yuhong Yang and Dan Zhu. Randomized allocation with nonparametric estimation for a multi-  
611 armed bandit problem with covariates. *The Annals of Statistics*, 30(1):100–121, 2002.
- 612 Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song,  
613 and Gao Huang. Does reinforcement learning really incentivize reasoning capacity  
614 in llms beyond the base model? *arXiv preprint arXiv:2504.13837*, 2025. URL  
615 <https://arxiv.org/abs/2504.13837>.
- 616 Jihun Yun, Juno Kim, Jongho Park, Junhyuck Kim, Jongha Jon Ryu, Jaewoong Cho, and Kwang-  
617 Sung Jun. Alignment as distribution learning: Your preference model is explicitly a language  
618 model. *arXiv preprint arXiv:2506.01523*, 2025.  
619
- 620 Wenhao Zhan, Masatoshi Uehara, Nathan Kallus, Jason D Lee, and Wen Sun. Provable offline  
621 preference-based reinforcement learning. *arXiv preprint arXiv:2305.14816*, 2023.
- 622 Yudi Zhang, Lu Wang, Meng Fang, Yali Du, Chenghua Huang, Jun Wang, Qingwei Lin, Mykola  
623 Pechenizkiy, Dongmei Zhang, Saravan Rajmohan, et al. Distill not only data but also rewards:  
624 Can smaller language models surpass larger ones?, 2025. URL <https://arxiv.org/abs/2502.19557>,  
625 2025.  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

648	APPENDIX TABLE OF CONTENTS	
649		
650		
651	<b>A</b>	<b>Miscellaneous Deferred Discussions</b> <span style="float: right;"><b>13</b></span>
652		
653	<b>B</b>	<b>Technical Preliminary Lemmas</b> <span style="float: right;"><b>14</b></span>
654		
655	<b>C</b>	<b>Proof of Proposition 1: MLE is Good for Low-Cardinality <math>\Pi</math></b> <span style="float: right;"><b>15</b></span>
656		
657	<b>D</b>	<b>Simple MLE Failures for Low-Cardinality <math>\mathcal{S}</math></b> <span style="float: right;"><b>16</b></span>
658		
659	D.1	Overlap of MLE . . . . . 18
660		
661	<b>E</b>	<b>Proofs from Section 4.2</b> <span style="float: right;"><b>19</b></span>
662		
663	E.1	Online-to-Batch Analysis via Freedman’s Inequality . . . . . 19
664		
665	E.2	Upper Bounds for Common Intersection and Majority . . . . . 19
666		
667	E.3	Lower Bounds for Common Intersection and Majority . . . . . 20
668		
669	<b>F</b>	<b>Algorithm and Analysis for Learning from Arbitrary Demonstrator</b> <span style="float: right;"><b>22</b></span>
670		
671	F.1	Online Learner with Softer Update . . . . . 22
672		
673	F.2	Online-to-Batch Conversion and Statistical Guarantee . . . . . 24
674		
675	<b>G</b>	<b>Algorithms and Proofs for <math>\text{pass}@k</math>-Error</b> <span style="float: right;"><b>25</b></span>
676		
677	G.1	Statistical Upper Bound . . . . . 27
678		
679	G.2	Lower Bounds for Online and Statistical Settings for $k$ -pass Error . . . . . 28

## A MISCELLANEOUS DEFERRED DISCUSSIONS

**Validity vs Coverage.** For LLMs, we discuss several important reasons why our objective in Eq. (1) is important and worth studying in its own right, despite not demanding coverage or matching the distribution of the demonstrator. (1) First, it mirrors real-world usage of deployed LLMs (OpenAI, 2023; Anthropic, 2024; DeepMind, 2023; AI, 2025; Touvron et al., 2023), where the feedback signal is based solely on the quality of the single output shown to the user. The models are never directly evaluated using distributional distances, nor do we believe such evaluation is even feasible. (2) The objective may tolerate mode collapse, which might seem limiting to those who view coverage as important for fairness-related reasons. Even in such situations, the model should perhaps indicate its uncertainty (Kirchhof et al., 2025). And there are many other situations—typically during SFT—where coverage is not the primary goal. (3) Achieving both “validity” (precision) and “coverage” (roughly recall or calibration) is challenging and often in tension; we have both empirical (Bronnec et al., 2024) and theoretical (Kalai & Vempala, 2024; Kalai et al., 2025) evidence of this. In such situations, it is important to at least target the validity. (4) Finally, given the extensive alignment and post-training that modern LLMs undergo, which are completely reward-driven, it is unclear whether they preserve any coverage guarantees over correct responses (Song et al., 2025; Dang et al., 2025; Yue et al., 2025; Wu et al., 2025).

**Proper vs Improper Learning.** The model class  $\mathcal{S}$  consists of set-valued functions  $\sigma$ , whereas the prediction is a single label. Thus, we define a notion of proper learning that is natural for our problem.

**Definition 2** (Proper Learning). *We call a learning rule  $\hat{\pi} : (\mathcal{X} \times \mathcal{Y})^* \rightarrow (\Delta(\mathcal{Y}))^{\mathcal{X}}$  proper if for any  $S \in (\mathcal{X} \times \mathcal{Y})^*$ , the policy  $\hat{\pi}(S)$  is supported on  $\sigma$  for some  $\sigma \in \mathcal{S}$ , i.e.  $\text{supp}(\hat{\pi}(S)(\cdot | x)) \subseteq \sigma(x)$  for all  $x \in \mathcal{X}$ .*

Our optimal learner is randomized and improper. On the other hand, the rule COMMON-INTERSECTION is *deterministic* and also *proper*, after making a choice when outputting an arbitrary  $y$ —in the case where the common intersection is empty, we output a  $y$  that always belongs to  $\sigma(x)$  for some fixed  $\sigma \in V(S)$ . The proofs of Theorems 2 and 9 still hold because we always treat an arbitrary  $y$  as a mistake anyway.

**Related Work in Learning in contextual bandits.** Contextual bandits have a long history in statistics. Early work includes Woodroffe (1979); Yang & Zhu (2002). The term contextual bandits was popularized by the paper Langford & Zhang (2007). A major line of work focuses on-line learning in contextual bandits (Chu et al., 2011; Abbasi-Yadkori et al., 2011; Agarwal et al., 2014); we refer interested readers to the recent book Lattimore & Szepesvári (2020) for more details. A parallel line of work focuses on the offline setting, where one learns from logged bandit feedback without active exploration. Influential contributions in this direction include Dudík et al. (2011); Bottou et al. (2012); Swaminathan & Joachims (2015). Most related to our work is imitation learning in contextual bandits, where the learner has access only to demonstrations from an expert. This perspective has been explored in recent theoretical studies (Rajaraman et al., 2020; Rashidinejad et al., 2021; Foster et al., 2024). They predominantly work under the low-capacity policy class assumption, which is different from our low-capacity reward class assumption.

## B TECHNICAL PRELIMINARY LEMMAS

We start by a technical lemma about the one-sided change of measure bound on an expectation of a bounded function in terms of the Hellinger distance (e.g. Lemma A.11 from Foster et al. (2021)). We will use the exact variant from (Foster et al., 2024, Lemma 3.11).

**Lemma 1** (Change-of-measure bound via Hellinger distance, Foster et al. (2024)). *Let  $(\mathcal{Z}, \mathcal{F})$  be a measurable space and let  $\mathbb{P}, \mathbb{Q}$  be probability measures on it. For every measurable function  $h : \mathcal{Z} \rightarrow \mathbb{R}$ :*

$$|\mathbb{E}_{\mathbb{P}}[h] - \mathbb{E}_{\mathbb{Q}}[h]| \leq \sqrt{\frac{\mathbb{E}_{\mathbb{P}}[h^2] + \mathbb{E}_{\mathbb{Q}}[h^2]}{2}} D_{\text{H}}(\mathbb{P}, \mathbb{Q}). \quad (7)$$

*In particular for  $h : \mathcal{Z} \rightarrow [0, R]$ ,*

$$\mathbb{E}_{\mathbb{P}}[h] \leq 2\mathbb{E}_{\mathbb{Q}}[h] + R D_{\text{H}}^2(\mathbb{P}, \mathbb{Q}). \quad (8)$$

We now specify Freedman’s inequality that provides us with a non-asymptotic bound on the sum of martingale difference sequence.

**Lemma 2** (Freedman’s inequality, Theorem 3 from Li et al. (2021)). *Consider a filtration  $\mathcal{F}_0 \subset \mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots$ , and write  $\mathbb{E}_i[\cdot] := \mathbb{E}[\cdot | \mathcal{F}_i]$ . Let*

$$Y_m = \sum_{i=1}^m X_i$$

*where  $(X_i)$  is a real-valued scalar sequence satisfying:*

$$|X_i| \leq R, \quad \mathbb{E}_{i-1}[X_i] = 0 \quad \text{for all } i \geq 1,$$

*for some constant  $R < \infty$ . Define the predictable variance process*

$$W_m := \sum_{i=1}^m \mathbb{E}_{i-1}[X_i^2],$$

*and assume deterministically that  $W_m \leq \sigma^2$  for some constant  $\sigma^2 < \infty$ . Then for any integer  $n \geq 1$ , with probability at least  $1 - \delta$ ,*

$$|Y_m| \leq \sqrt{8 \max\left\{W_m, \frac{\sigma^2}{2^n}\right\} \log\left(\frac{2n}{\delta}\right)} + \frac{4}{3}R \log\left(\frac{2n}{\delta}\right).$$

**Maximum Likelihood Estimation for Distribution Learning.** We now state guarantee for the maximum likelihood estimator (MLE) for density estimation, exactly similar to (Foster et al., 2024, Section B.4). Given a class of candidate densities  $\mathcal{G}$  and i.i.d. samples  $z_1, \dots, z_m \sim g_*$  (possibly not in  $\mathcal{G}$ ), we define the empirical negative log-likelihood (log-loss) of  $g \in \mathcal{G}$  as

$$L_{\log}(g) = - \sum_{i=1}^m \log g(z_i).$$

The maximum likelihood estimator is then

$$\hat{g}_{\text{mle}} \in \arg \min_{g \in \mathcal{G}} L_{\log}(g). \quad (9)$$

**Definition 3** (Log-loss covering number). *For a class  $\mathcal{G} \subseteq \Delta(\mathcal{Z})$ , we say that a subset  $\mathcal{G}' \subseteq \Delta(\mathcal{Z})$  is an  $\varepsilon$ -cover with respect to the log-loss if for all  $g \in \mathcal{G}$  there exists  $g' \in \mathcal{G}'$  such that  $\sup_{z \in \mathcal{Z}} \log(g(z)/g'(z)) \leq \varepsilon$ . We denote the size of the smallest such cover by  $\mathcal{N}_{\log}(\mathcal{G}, \varepsilon)$ .*

We have the following property of MLE’s convergence in the squared Hellinger distance with high probability.

**Proposition 2.** *With probability  $1 - \delta$  over  $m$  i.i.d. samples from any  $g_* \in \mathcal{G}$ ,*

$$D_{\text{H}}^2(g_*, \hat{g}_{\text{mle}}) \leq \inf_{\varepsilon > 0} \left\{ \frac{6 \log(2 \mathcal{N}_{\log}(\mathcal{G}, \varepsilon)/\delta)}{m} + 4\varepsilon \right\} + 2 \inf_{g \in \mathcal{G}} \log(1 + D_{\chi^2}(g_* \| g)) + 2\varepsilon_{\text{opt}}.$$

*In particular, if  $\mathcal{G}$  is finite and  $\varepsilon_{\text{opt}} = 0$ , the maximum likelihood estimator satisfies*

$$D_{\text{H}}^2(g_*, \hat{g}_{\text{mle}}) \leq \frac{6 \log(2|\mathcal{G}|/\delta)}{m} + 2 \inf_{g \in \mathcal{G}} \log(1 + D_{\chi^2}(g_* \| g)).$$

*Note that the term  $\inf_{g \in \mathcal{G}} \log(1 + D_{\chi^2}(g_* \| g))$  corresponds to the misspecification error, and is zero if  $g_* \in \mathcal{G}$ .*

We note that the proof of (Foster et al., 2024, Proposition B.1) contains a couple of minor typographical errors. Namely, in Eq.(20) therein, the authors aim to compare  $\tilde{g}$  and  $\hat{g}$ , but ended up comparing  $\tilde{g}$  and  $g_*$ . A similar mistake is repeated a couple of more times later in the proof without affecting the correctness of the argument.

## C PROOF OF PROPOSITION 1: MLE IS GOOD FOR LOW-CARDINALITY $\Pi$

In this, we present the proof of Proposition 1. As discussed in Section 3, the argument builds on a result from density estimation (Proposition 1) showing that MLE achieves convergence in squared Hellinger distance. We begin by providing intuition for why, even in this setting, MLE can obtain a guarantee expressed explicitly in terms of  $\log |\Pi|$ , without any dependence on  $|\mathcal{X}|$ ,  $|\mathcal{Y}|$ , or the support size  $\sup_{\sigma, x} |\sigma(x)|$ , in a special case of the problem.

**Intuition with the special cases of  $\Pi$ :** We provide a more transparent and direct proof for the special case when for every  $\pi \in \Pi$ ,  $x \in \mathcal{X}$ , the conditional density  $\pi(\cdot | x)$  puts a uniform distribution over exactly  $s$  members of  $\mathcal{Y}$  for some large but finite integer  $s$ . First, observe that in this special case we have a *dichotomy*; any hypothesis that does not contradict the data has the same likelihood as any other, so any  $\pi \in \Pi$  that does not contradict with the data is MLE. For the unknown  $\mathcal{D} \times \pi_*$ , we now consider any  $\pi$  such that

$$L_{\mathcal{D}, \sigma_*}(\pi) = \mathbb{P}_{x \sim \mathcal{D}, \hat{y} \sim \pi(\cdot | x)} (\hat{y} \notin \sigma_{\pi_*}(x)) > \varepsilon.$$

Then, due to the symmetry of the loss in the special case where each  $\pi \in \Pi$  puts a uniform distribution on exactly  $s$  responses, we have

$$\mathbb{P}_{x \sim \mathcal{D}, \hat{y} \sim \pi(\cdot | x)} (\hat{y} \notin \sigma_{\pi_*}(x)) = \mathbb{P}_{x \sim \mathcal{D}, y \sim \pi_*(\cdot | x)} (y \notin \sigma_{\pi}(x)) > \varepsilon,$$

where the key fact used is the ability to change the order of randomness between  $\hat{y} \sim \pi(\cdot | x)$  and  $y \sim \pi_*(\cdot | x)$ .

This shows that when we sample  $(x, y) \sim \mathcal{D} \times \pi_*$ , the probability that  $(x, y)$  does not fall in the support  $\sigma_\pi(x)$  exceeds  $\varepsilon$ . Hence, for any fixed  $\pi \in \Pi$ , after  $m$  i.i.d. draws

$$\mathbb{P}_S(\pi \text{ survives}) \leq (1 - \varepsilon)^m \leq e^{-\varepsilon m}.$$

Therefore, by a standard union bound,

$$\mathbb{P}_S(\exists \text{ bad } \pi \in \Pi \text{ that survives}) \leq |\Pi| e^{-\varepsilon m}.$$

The proposition follows by choosing  $m \geq m_{\Pi, \hat{\pi}_{\text{MLE}}}(\varepsilon, \delta) = \frac{\log |\Pi| + \log(1/\delta)}{\varepsilon}$ .

**Proof for any general  $\Pi$ :** Consider any unknown but fixed marginal distribution  $\mathcal{D} \in \Delta(\mathcal{X})$ . For any conditional law  $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$ , let  $\mathbb{P}_{(\mathcal{D}, \pi)}$  denote the joint law over  $(\mathcal{X} \times \mathcal{Y})$  given by the marginal distribution  $\mathcal{D}$  and the conditional law  $\pi(\cdot | x)$ . First observe that for any  $S \in (\mathcal{X} \times \mathcal{Y})^*$ , the joint law  $\mathbb{P}_{(\mathcal{D}, \hat{\pi}_{\text{MLE}}(S))}$  is the MLE of among all joint distribution  $\{\mathbb{P}_{(\mathcal{D}, \pi)} : \pi \in \Pi\}$ . Using Proposition 2, for  $S \sim (\mathcal{D} \times \pi_*)^m$

$$\mathbb{P}_S \left( D_{\text{H}}^2(\mathbb{P}_{(\mathcal{D}, \pi_*)}, \mathbb{P}_{\mathcal{D}, \hat{\pi}_{\text{MLE}}(S)}) \leq \frac{6 \log(2|\Pi|/\delta)}{m} \right) \geq 1 - \delta. \quad (10)$$

Now let  $\sigma_* : \mathcal{X} \rightarrow 2^{\mathcal{Y}}$  be the associated support set valued function of valid responses for  $\sigma_*(x) = \text{supp}(\pi_*(\cdot | x))$ . Let us define a function  $\text{err} : (\mathcal{X} \times \mathcal{Y}) \rightarrow \{0, 1\}$  as

$$\text{err}(x, y) = \begin{cases} 1 & \text{if } y \notin \sigma_*(x), \\ 0 & \text{otherwise.} \end{cases}$$

Then using Lemma 1, we have for any conditional law  $\pi : \mathcal{X} \rightarrow \Delta(\mathcal{Y})$

$$L_{\mathcal{D}, \sigma_*}(\pi) = \mathbb{E}_{\mathbb{P}_{(\mathcal{D}, \pi)}}[\text{err}] \leq D_{\text{H}}^2(\mathbb{P}_{(\mathcal{D}, \pi_*)}, \mathbb{P}_{(\mathcal{D}, \pi)}),$$

where we used the fact that  $L_{\mathcal{D}, \sigma_*}(\pi_*) = \mathbb{E}_{\mathbb{P}_{\mathcal{D}, \pi_*}}[\text{err}] = 0$  and that  $\text{err}$  is a bounded function in  $[0, 1]$ . Combining this with (10), we obtain that with probability at least  $1 - \delta$  over  $S \sim (\mathcal{D} \times \pi_*)^m$ ,

$$L_{\mathcal{D}, \sigma_*}(\hat{\pi}_{\text{MLE}}(S)) \leq \frac{6 \log(2|\Pi|/\delta)}{m}.$$

## D SIMPLE MLE FAILURES FOR LOW-CARDINALITY $\mathcal{S}$

We first show that there is a simple instance of a support class, where some MLE over the entire class  $\Pi_{\mathcal{S}} := \bigcup_{\sigma} \Pi_{\sigma}$  fails.

*Proof of Theorem 1.* Fix any  $\gamma \in (0, 1)$ . Let  $\mathcal{Y} = \{0, 1\}$  and  $\mathcal{X} = \mathbb{N}$ . Define a support class  $\mathcal{S} = \{\sigma_0, \sigma_{01}\} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  by

$$\sigma_0(x) = \{0\} \quad \text{and} \quad \sigma_{01}(x) = \{0, 1\} \quad \forall x \in \mathcal{X}.$$

Choose the ground-truth support  $\sigma_* = \sigma_0$  and let the data-generating conditional be the point mass  $\pi_*(\cdot | x) = \delta_0(\cdot)$  for all  $x$ ; thus every observed label equals 0. Let  $\Pi_{\mathcal{S}}$  be the class of all policies supported on  $\mathcal{S}$  (see (2)). Now fix a sample size  $m \in \mathbb{N}$ . Set

$$q := \left\lceil \frac{m}{\gamma} \right\rceil,$$

and define the marginal  $\mathcal{D}$  to be the uniform distribution on  $[q] = \{1, 2, \dots, q\}$ , i.e.,  $\mathcal{D}(\{x\}) = 1/q$  for  $x \in [q]$  and 0 otherwise.

For any dataset  $S = \{(x_i, y_i)\}_{i=1}^m \sim (\mathcal{D} \times \pi_*)^m$ , write  $S_{\text{dis}} := \{x_i : i \in [m]\}$  for the set of distinct unlabeled inputs in  $S$  (so  $|S_{\text{dis}}| \leq m$ ). Consider the predictor  $\hat{\pi}$  defined by

$$\hat{\pi}(\cdot | x) = \begin{cases} \delta_0(\cdot), & x \in S_{\text{dis}}, \\ \delta_1(\cdot), & x \notin S_{\text{dis}}. \end{cases} \quad (11)$$

We claim that  $\hat{\pi} \in \text{MLE}_{\Pi_{\mathcal{S}}}(S)$ . Indeed, the log-likelihood is

$$\ell_{\log}(\pi; S) = \sum_{i=1}^m \log \pi(0 | x_i) = \sum_{x \in S_{\text{dis}}} N_x(S) \log \pi(0 | x),$$

where  $N_x(S) := |\{i : x_i = x\}|$ . This expression depends only on the values  $\pi(0 | x)$  for  $x \in S_{\text{dis}}$  and is maximized by setting  $\pi(0 | x) = 1$  for every  $x \in S_{\text{dis}}$ . For  $x \notin S_{\text{dis}}$  the likelihood does not constrain  $\pi(\cdot | x)$ , so any choice (for tie-breaking) is a valid maximizer; in particular, (11) yields a valid MLE in  $\Pi_{\mathcal{S}}$ .

We next evaluate its population loss against the support  $\sigma_*$ . Since  $\sigma_*(x) = \{0\}$  for all  $x$ ,

$$L_{\mathcal{D}, \sigma_*}(\hat{\pi}) = \mathbb{P}_{x \sim \mathcal{D}, \hat{y} \sim \hat{\pi}(\cdot | x)}(\hat{y} \notin \sigma_*(x)) = \mathbb{P}_{x \sim \mathcal{D}}(x \notin S_{\text{dis}}) = 1 - \frac{|S_{\text{dis}}|}{q}.$$

Using  $|S_{\text{dis}}| \leq m$  and  $q \geq m/\gamma$ ,

$$L_{\mathcal{D}, \sigma_*}(\hat{\pi}) \geq 1 - \frac{m}{q} \geq 1 - \gamma.$$

The bound holds deterministically for every sample  $S$ , hence

$$\mathbb{P}_{S \sim (\mathcal{D} \times \pi_*)^m}(L_{\mathcal{D}, \sigma_*}(\hat{\pi}) \geq 1 - \gamma) = 1.$$

□

We now show that the attempt to restrict the capacity of the class via another natural choice of considering  $\Pi_{\text{unif}, \mathcal{S}} = \bigcup_{\sigma \in \mathcal{S}} \{\pi_{\text{unif}, \sigma}\}$  also does not work, when the expert demonstrations  $\pi_*$  does not necessarily follow the distribution  $\pi_{\text{unif}, \sigma_*}$  while still showing optimal demonstrations.

**Theorem 7 (MLE Failure 2).** Fix  $\gamma \in (0, 1)$ . There exists  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  with  $|\mathcal{S}| = 2$ ,  $|\mathcal{X}| = 1$ ,  $|\mathcal{Y}| = 2\lceil 1/\gamma \rceil$ , such that for some choice of  $\mathcal{D} \times \pi_*$  for some  $\pi_*$  supported on  $\mathcal{S}$  such that, for every sample size  $m$ , for  $S \sim (\mathcal{D} \times \pi_*)^m$ , there is a unique  $\hat{\pi}_{\text{mle}}(S) \in \text{MLE}_{\Pi_{\text{unif}, \mathcal{S}}}(S)$  that suffers from the following loss:

$$\mathbb{P}_{S \sim (\mathcal{D} \times \pi_*)^m}(L_{\mathcal{D}, \sigma_*}(\hat{\pi}_{\text{mle}}(S)) \geq 1 - \gamma) = 1.$$

*Proof of Theorem 7.* Fix  $\gamma \in (0, 1)$  and let  $s := \lceil 1/\gamma \rceil$ . Take  $\mathcal{X} = \{x\}$  and

$$\mathcal{Y} = \{y^*\} \cup \{a_1, \dots, a_{s-1}\} \cup \{b_1, \dots, b_s\},$$

so  $|\mathcal{Y}| = 1 + (s-1) + s = 2s = 2\lceil 1/\gamma \rceil$ . Define  $\sigma_1, \sigma_2 \in (2^{\mathcal{Y}})^{\mathcal{X}}$  by

$$\sigma_1(x) = \{y^*, a_1, \dots, a_{s-1}\} \quad (\text{size } s), \quad \sigma_2(x) = \{y^*, b_1, \dots, b_s\} \quad (\text{size } s+1),$$

so  $\sigma_1(x) \cap \sigma_2(x) = \{y^*\}$  and they are otherwise disjoint. Let  $\mathcal{S} = \{\sigma_1, \sigma_2\}$  and

$$\Pi_{\text{unif}, \mathcal{S}} := \{\pi_{\text{unif}, \sigma} : \sigma \in \mathcal{S}\}, \quad \pi_{\text{unif}, \sigma}(y | x) = \begin{cases} \frac{1}{|\sigma(x)|}, & y \in \sigma(x), \\ 0, & \text{otherwise.} \end{cases}$$

Set  $\mathcal{D}$  to be the point mass at  $x$  and choose the ground-truth support  $\sigma_* = \sigma_2$  with data-generating conditional  $\pi_* = \delta_{y^*}$  (always emit  $y^*$ ). For any  $m$ , every dataset  $S \sim (\mathcal{D} \times \pi_*)^m$  equals  $\{(x, y^*)\}^m$ .

It is simple to see that  $\pi_{\text{unif}, \sigma_1} \in \text{MLE}_{\Pi_{\text{unif}, \mathcal{S}}}(S)$  is the unique maximum likelihood estimator. This is because

$$\prod_{i=1}^m \pi_{\text{unif}, \sigma_1}(y_i | x_i) = \left(\frac{1}{s}\right)^m, \quad \prod_{i=1}^m \pi_{\text{unif}, \sigma_2}(y_i | x_i) = \left(\frac{1}{s+1}\right)^m.$$

However, with  $\sigma_*(x) = \sigma_2(x)$ , the estimator  $\hat{\pi}_{\text{mle}}(S) = \pi_{\text{unif}, \sigma_1}$  has the error

$$L_{\mathcal{D}, \sigma_*}(\hat{\pi}_{\text{mle}}(S)) = \mathbb{P}_{\hat{y} \sim \pi_{\text{unif}, \sigma_1}(\cdot | x)}(\hat{y} \notin \sigma_2(x)) = 1 - \pi_{\text{unif}, \sigma_1}(y^* | x) = 1 - \frac{1}{s} = 1 - \frac{1}{\lceil 1/\gamma \rceil} \geq 1 - \gamma.$$

All bounds are deterministic given  $S$ , hence

$$\mathbb{P}_{S \sim (\mathcal{D} \times \pi_*)^m}(L_{\mathcal{D}, \sigma_*}(\hat{\pi}_{\text{mle}}(S)) \geq 1 - \gamma) = 1,$$

for every  $m$ , completing the proof. □

## D.1 OVERLAP OF MLE

We now show that MLE over the restricted capacity class  $\Pi_{\text{unif},\mathcal{S}}$  attains a *non-trivial overlap* at the statistical limit of Theorem 4, though its failure to directly optimize the objective of interest (cf. Section 3). Note that it is trivial to achieve overlap, by just using a policy which puts mass on one good response from each support function class. However, this would produce the responses outside the support with overwhelming probability. Thus, for the purpose for overlap to be meaningful, we call it a *non-trivial* one if the output policy is only supported on one of  $\sigma \in \mathcal{S}$  (some notion *properness*). The MLE over the class  $\Pi_{\text{unif},\mathcal{S}} = \bigcup_{\sigma \in \mathcal{S}} \{\pi_{\text{unif},\sigma}\}$  will be of this form and achieves a non-trivial overlap.

**Theorem 8.** *For any finite class  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  and an unknown joint realizable distribution  $\mathcal{D} \times \pi_*$ , where  $\pi_*$  is supported on some  $\sigma_* \in \mathcal{S}$ , for any estimator  $\hat{\pi}_{\text{mle}}(S) \in \text{MLE}_{\Pi_{\text{unif},\mathcal{S}}}(S)$ , we have the following guarantee: for any sample size  $m \geq \varepsilon^{-1} (\log |\mathcal{S}| + \log(1/\delta))$ , we have*

$$\mathbb{P}_{S \sim (\mathcal{D} \times \pi_*)^m} (\mathbb{P}_{x \sim \mathcal{D}} (\text{supp}(\hat{\pi}_{\text{mle}}(S)(\cdot | x)) \cap \sigma_*(x) = \emptyset) \leq \varepsilon) \geq 1 - \delta.$$

*Proof of Theorem 8.* The proof is simple. First note that  $\pi_{\text{unif},\sigma_*}$  has non-zero likelihood. Therefore, any policy in the set  $\text{MLE}_{\Pi_{\text{unif},\mathcal{S}}}(S)$  must have non-zero likelihood. Thus, for any  $\sigma$  for which  $\pi_{\text{unif},\sigma} \in \text{MLE}_{\Pi_{\text{unif},\mathcal{S}}}(S)$ , we must have that  $\sigma \in V(S) := \{\sigma \in \mathcal{S} : y_i \in \sigma(x_i) \forall (x_i, y_i) \in S\}$ . Therefore, in order to establish

$$\mathbb{P}_{S \sim (\mathcal{D} \times \pi_*)^m} (\mathbb{P}_{x \sim \mathcal{D}} (\text{supp}(\hat{\pi}_{\text{mle}}(S)(x)) \cap \sigma_*(x) = \emptyset) \leq \varepsilon) \geq 1 - \delta,$$

it suffices to establish

$$\mathbb{P}_S (\forall \sigma \in V(S) : \mathbb{P}_{x \sim \mathcal{D}} (\sigma(x) \cap \sigma_*(x) = \emptyset) \leq \varepsilon) \geq 1 - \delta. \quad (12)$$

Consider any bad  $\sigma \in \mathcal{S}$  such that  $\mathbb{P}_{x \sim \mathcal{D}} (\sigma(x) \cap \sigma_*(x) = \emptyset) > \varepsilon$ . With each draw  $(x_i, y_i) \sim (\mathcal{D} \times \pi_*)$ , we have that  $\sigma$  gets knocked-out of version space with probability at least  $\varepsilon$ , i.e.  $\mathbb{P}_{(x_i, y_i) \sim (\mathcal{D} \times \pi_*)} (y_i \notin \sigma(x_i)) > \varepsilon$ . Therefore, for any fixed  $\sigma$ , after sample  $S \sim (\mathcal{D} \times \pi_*)^m$

$$\mathbb{P}_S (\sigma \in V(S)) \leq (1 - \varepsilon)^m \leq e^{-\varepsilon m}.$$

Therefore, by a standard union bound,

$$\mathbb{P}_S (\exists \text{ bad } \sigma \in V(S) \text{ that survives}) \leq |\mathcal{S}| e^{-\varepsilon m} \leq |\mathcal{S}| 2^{-\varepsilon m}.$$

The theorem follows by noting that the  $|\mathcal{S}| 2^{-\varepsilon m} \leq \delta$  for any  $m \geq \frac{\log |\mathcal{S}| + \log(1/\delta)}{\varepsilon}$ .  $\square$

**Remark D.1** (Comparison with multi-class classification). *Note that for multiclass classification when  $\mathcal{S} \subseteq \mathcal{Y}^{\mathcal{X}}$  (i.e. all  $|\sigma(x)| = 1$ , the guarantee captured in (12) is enough to ensure learnability by just outputting a single predictor from  $\hat{\sigma} \in V(S)$  (i.e. consistent / ERM). This happens because the overlap implies that that labels are the same and so no error. However, for our problem despite this overlap, it is unclear how to output a single label so that it belongs to the support of  $\sigma_*$ . What would be sufficient for our problem is the following guarantee, where the quantifier  $\forall \sigma \in V(S)$  is taken inside the randomness of test point sampling  $x \sim \mathcal{D}$ :*

$$\mathbb{P}_S (\mathbb{P}_{x \sim \mathcal{D}} (\forall \sigma \in V(S) : \sigma(x) \cap \sigma_*(x) = \emptyset) \leq \varepsilon) \geq 1 - \delta.$$

*However, we know that this provably requires the sample size where there is  $\Omega(|\mathcal{S}|)$  dependence on cardinality—see the lower bound for COMMON-INTERSECTION estimator (Theorem 11).*

It may be possible to turn this into a predictor that directly starts to produce good responses, depending on the overlap among hypotheses and other types of feedback available in post-training (e.g., whether a generated response is good or not). This overlap can be captured by a parameter that reflects the need for repeated sampling and the number of feedback that must be queried, which in turn allows for a more quantitative understanding of how many feedbacks are required to guarantee performance in terms of this parameter. For example, this parameter would be maximum in the case of multi-class classification (Remark D.1) and no additional feedback is required. However, we leave it open to formulate an interesting setup that enables a study of both types of feedbacks together for our problem, and we do not attempt to investigate this any further.

## 972 E PROOFS FROM SECTION 4.2

973 We first show our proof of Theorem 4 and then return to the proofs for the rules  
974 COMMON-INTERSECTION and MAJORITY.  
975

### 976 E.1 ONLINE-TO-BATCH ANALYSIS VIA FREEDMAN’S INEQUALITY

977 We have an online learner MISTAKE-UNAWARE-WEIGHT-UPDATE (Algorithm 1) that makes at  
978 most  $\log_2 |\mathcal{S}|$  mistakes (Theorem 3). We now show, how using online-to-batch conversion via the  
979 estimator  $\widehat{\pi}_{\text{o2b}}$  (Eq. (5)), we can enjoy a similar sample complexity. The main difference from  
980 standard online-to-batch analysis is the use of Freedman’s inequality (Lemma 2), which yields a  
981 better dependence on  $1/\varepsilon$ .  
982  
983

984 *Proof of Theorem 4.* Let  $\ell_t = \mathbb{1}\{\widehat{\pi}_t(x_t) \notin \sigma_*(x_t)\}$ . Because  $\widehat{\pi}_t$  is a deterministic function of  
985  $S_{<t} = \{(x_i, y_i) : i < t\}$ , we have

$$986 \mathbb{E}[\ell_t \mid S_{<t}] = L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t).$$

987 Hence

$$988 \mathbb{E}_S [L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_{\text{o2b}}(S))] = \mathbb{E}_S \left[ \frac{1}{m} \sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t) \right] = \mathbb{E}_S \left[ \frac{1}{m} \sum_{t=1}^m \ell_t \right] \leq \frac{\log_2 |\mathcal{S}|}{m},$$

989 where in the last inequality we used Theorem 3, which guarantees  $\sum_{t=1}^m \ell_t \leq \log_2 |\mathcal{S}|$ .

990 For the high-probability statement, define the martingale differences

$$991 Z_t := L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t) - \ell_t, \quad \text{where } |Z_t| \leq 1 \text{ almost surely.}$$

992 Then  $\mathbb{E}[Z_t \mid S_{<t}] = 0$ , and

$$993 \mathbb{E}[Z_t^2 \mid S_{<t}] = \mathbb{E}[(L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t) - \ell_t)^2 \mid S_{<t}] = \text{Var}(\ell_t \mid S_{<t}) = L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t)(1 - L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t)) \leq L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t).$$

994 And taking  $W_m = \sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t)$  and  $\sigma^2 = m$  suffices, thus, using Lemma 2 with  $n = \log m$   
995 inequality gives us with probability  $1 - \delta$

$$996 \sum_{t=1}^m Z_t \leq \sqrt{8 \left( 1 + \sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t) \right) \log \left( \frac{\log m}{\delta} \right) + \frac{4}{3} \log \left( \frac{\log m}{\delta} \right)}$$

$$997 \leq \frac{1}{2} \left( 1 + \sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t) \right) + 4 \log \left( \frac{\log m}{\delta} \right) + \frac{4}{3} \log \left( \frac{\log m}{\delta} \right) \quad (\text{GM} \leq \text{AM})$$

998 Substituting  $Z_t$  and rearranging terms,

$$999 \sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t) \leq 1 + 2 \sum_{t=1}^m \ell_t + \frac{32}{3} \log \left( \frac{\log m}{\delta} \right)$$

1000 Finally noting that  $L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_{\text{o2b}}(S)) = \frac{1}{m} \sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_t)$  and that  $\sum_{t=1}^m \ell_t \leq \log_2 |\mathcal{S}|$  (by Theo-  
1001 rem 3), we obtain that with probability  $1 - \delta$ ,

$$1002 L_{\mathcal{D}, \sigma_*}(\widehat{\pi}_{\text{o2b}}(S)) \leq \frac{1 + 2 \log_2 |\mathcal{S}| + 12 \log \left( \frac{\log m}{\delta} \right)}{m}$$

1003  $\square$

### 1004 E.2 UPPER BOUNDS FOR COMMON INTERSECTION AND MAJORITY

1005 We start by analyzing the COMMON-INTERSECTION rule in the more difficult online setting which  
1006 helps for the intuition for the statistical setting.

1007 **Theorem 9** (Online Guarantee for COMMON-INTERSECTION). *On any sequence  $((x_t, y_t))_{t \in \mathbb{N}}$  re-  
1008 alizable by some  $\sigma_* \in \mathcal{S}$ , the rule COMMON-INTERSECTION (applied to the sequence seen so far)  
1009 makes at most  $|\mathcal{S}| - 1$  mistakes.*

*Proof of Theorem 9.* Consider any round  $t$  in which there was a mistake made by the rule. It must be that the set of consistent hypothesis  $V_t$  in that round, it must be that there was no *common intersection* in that round, i.e.,  $\bigcap_{\sigma \in V_t} \sigma(x_t) = \emptyset$ . That means even though we would not know whether we made a mistake in that round, observing  $y_t$  will eliminate at least one hypothesis from the version space (i.e.  $|V_{t+1}| \leq |V_t| - 1$ ). Therefore, the rule cannot make  $|V_1| - 1 = |\mathcal{S}| - 1$  mistakes on any realizable sequence.  $\square$

We now analyze the performance of this rule in the statistical version.

*Proof of Theorem 2.* Partition the  $m$  examples into  $K := |\mathcal{S}|$  consecutive blocks  $B_1, \dots, B_K$ , each of length  $n \geq \frac{1}{\varepsilon} (\log |\mathcal{S}| + \log(1/\delta))$ . Let  $V_t$  denote the version space just before block  $B_t$  begins. I.e. define the restricted dataset  $S_t = B_1 \cup \dots \cup B_{t-1}$  and

$$V_t = \{\sigma \in \mathcal{S} : y_i \in \sigma(x_i) \forall (x_i, y_i) \in S_t\}$$

with  $V_1 = \mathcal{S}$ . Define the region of  $x$ , where we do not have a common intersection among  $V_t$ .

$$A_t := \left\{ x \in \mathcal{X} : \bigcap_{\sigma \in V_t} \sigma(x) = \emptyset \right\}.$$

Note that  $A_{t+1} \subseteq A_t$  for all  $t \in [K]$  because  $V_{t+1} \subseteq V_t$ . Moreover, we never make an error when outputting from common-intersection region. Now, using these facts, we have

$$\begin{aligned} \mathbb{P}_{\mathcal{S}}(L_{\mathcal{D}, \sigma_*}(\hat{\pi}_{\text{CI}}(S)) > \varepsilon) &\leq \mathbb{P}_{\mathcal{S}}(\mathbb{P}_{x \sim \mathcal{D}}(A_{K+1}) > \varepsilon) \\ &= \mathbb{P}_{\mathcal{S}}(V_{K+1} \neq \emptyset \cap \mathbb{P}_{x \sim \mathcal{D}}(A_{K+1}) > \varepsilon) \quad (V_{K+1} \neq \emptyset \text{ always happens}) \\ &\leq \mathbb{P}_{\mathcal{S}}(\exists t \in [K], V_{t+1} = V_t \cap \mathbb{P}_{x \sim \mathcal{D}}(A_{K+1}) > \varepsilon) \\ &\leq \mathbb{P}_{\mathcal{S}}(\exists t \in [K], V_{t+1} = V_t \cap \mathbb{P}_{x \sim \mathcal{D}}(A_t) > \varepsilon) \\ &\leq \sum_{t=1}^K \mathbb{P}_{B_t}(\exists t \in [K], V_{t+1} = V_t \mid \mathbb{P}_{x \sim \mathcal{D}}(A_t) > \varepsilon) \\ &\leq \sum_{t=1}^K (1 - \varepsilon)^{|B_t|} = K(1 - \varepsilon)^n \\ &\leq |\mathcal{S}| 2^{-\varepsilon n} \leq |\mathcal{S}| \cdot \frac{\delta}{|\mathcal{S}|} = \delta. \end{aligned}$$

The above calculation formalizes the following argument. There are two cases to consider:

**Case 1.** If  $\mathbb{P}_{x \sim \mathcal{D}}(A_t) > \varepsilon$ , then the probability that no  $x \in A_t$  appears in block  $B_t$  is at most  $(1 - \varepsilon)^n \leq e^{-\varepsilon n} \leq 2^{-\varepsilon n} \leq \delta/|\mathcal{S}|$ . Otherwise, some  $x \in A_t$  appears; since  $\bigcap_{\sigma \in V_t} \sigma(x) = \emptyset$ , the observed label  $y \in \sigma_*(x)$  excludes at least one  $\sigma \in V_t$ , so  $|V_{t+1}| \leq |V_t| - 1$ .

**Case 2.** If  $\mathbb{P}_{x \sim \mathcal{D}}(A_t) \leq \varepsilon$ , then on  $A_t^c$  the intersection is nonempty, and because  $\sigma_* \in V_t$  it follows that the COMMON-INTERSECTION prediction is always correct there. Moreover, since  $V_{t+1} \subseteq V_t$ , the intersections can only grow and hence  $A_{t+1} \subseteq A_t$ . Therefore, once Case 2 holds, the final error remains below  $\varepsilon$ .

$$L_{\mathcal{D}, \sigma_*}(\hat{\pi}_{\text{CI}}(S)) \leq \mathbb{P}_{\mathcal{D}}(A_t) < \varepsilon.$$

Putting these together, with probability at least  $1 - K \cdot (\delta/|\mathcal{S}|) \geq 1 - \delta$ , each block in Case 1 eliminates at least one hypothesis, and there are at most  $K = |\mathcal{S}|$  such eliminations are even possible. Hence either Case 2 occurs in some block (giving final error  $< \varepsilon$ ), or Case 1 occurs in all  $K$  blocks, which is not possible in the realizable setting, arriving at a contradiction.  $\square$

### E.3 LOWER BOUNDS FOR COMMON INTERSECTION AND MAJORITY

For the lower bound, we will show the lower bound on a stronger rule which outputs according to a majority vote among the consistent hypothesis formalized below.

**Input:** Sample  $S = \{(x_i, y_i) : i \in [m]\}$  and a finite support hypothesis class  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$ .

- Let  $V(S) := \{\sigma \in \mathcal{S} : y_i \in \sigma(x_i), \forall (x_i, y_i) \in S\}$
- Return the predictor  $\text{MAJORITY}(S) = \hat{\pi}_{\text{Maj}}(S) : \mathcal{X} \rightarrow \mathcal{Y}$  defined as follows:

$$\hat{\pi}_{\text{Maj}}(S)(x) = \arg \max_{y \in \mathcal{Y}} |\{\sigma \in V(S) : y \in \sigma(x)\}|$$

**Corollary 1.** *Note that  $\hat{\pi}_{\text{Maj}}$  always outputs from the common-intersection whenever it is non-empty, and thus, it enjoys the same online as well as the statistical guarantees as in Theorem 9 and Theorem 2 respectively.*

The lower bounds will hold for the following instance of the of the class.

**Description of the class.** Fix  $d \in \mathbb{N}$ . Let

$$\mathcal{Y} = \{0, 1\}, \quad q := \lfloor \frac{d-1}{2} \rfloor, \quad \mathcal{X} := \{1, 2, \dots, q\}.$$

We define a hypothesis class  $\mathcal{S} = \{\sigma_1, \sigma_2, \dots, \sigma_d\} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  as follows.

- **Distinguished hypothesis.** Set  $\sigma_1(x) = \{1\}$  for every  $x \in \mathcal{X}$ . This will serve as the ground-truth hypothesis.
- **Adversarial hypotheses.** For each  $i \geq 2$ , require that  $0 \in \sigma_i(x)$  for all  $x \in \mathcal{X}$ . Moreover, for each  $t \in \{1, \dots, q\}$  we designate a *pair* of hypotheses,  $\sigma_{2t}, \sigma_{2t+1}$ , that both exclude label 1 at coordinate  $t$ :

$$1 \notin \sigma_{2t}(t), \quad 1 \notin \sigma_{2t+1}(t).$$

For all other coordinates  $x \neq t$ , these hypotheses include both labels, e.g.

$$\sigma_{2t}(x) = \sigma_{2t+1}(x) = \{0, 1\} \quad \text{for } x \neq t.$$

If  $(d-1)$  is odd, then there is one remaining index pairing. In that case, define  $\sigma_d(x) = \{0, 1\}$  for all  $x \in \mathcal{X}$ .

Thus  $\mathcal{S}$  has size exactly  $d$ , uses  $2q \leq d-1$  adversarial hypotheses to plant two “anti-1” voters at each coordinate  $t \in \mathcal{X}$ , and possibly one additional “neutral” hypothesis if  $d-1$  is odd. We are now ready to show the online lower bound.

**Theorem 10** (Online Lower Bounds for COMMON-INTERSECTION and MAJORITY). *For every  $d$ , there exists a hypothesis class  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  with  $|\mathcal{S}| = d$ ,  $|\mathcal{X}| = \lfloor (d-1)/2 \rfloor$  and  $|\mathcal{Y}| = 2$  such that both the rules make  $|\mathcal{X}|$  mistakes (i.e. a mistake on every round).*

*Proof of Theorem 10.* Note that it suffices to show the lower bound of simply the MAJORITY rule, which also implies the lower bound on COMMON-INTERSECTION. Consider the hypothesis class  $\mathcal{S}$  constructed above, and let the ground truth be  $\sigma_* = \sigma_1$ . Present the sequence of instances  $x_t = t$  for  $t = 1, \dots, q = |\mathcal{X}|$ . Then  $y_t = 1$  for all  $t$  under  $\sigma_*$ .

At each round  $t$ , the version space  $V_{t-1}$  contains  $\sigma_1$  together with all adversarial hypotheses that have not yet been eliminated. By construction, every adversarial hypothesis other than  $\sigma_1$  always includes 0, while at coordinate  $t$  at least two of them exclude 1. Hence

$$N_0(x_t; V_{t-1}) = |V_{t-1}| - 1 \quad \text{and} \quad N_1(x_t; V_{t-1}) \leq |V_{t-1}| - 2,$$

so the majority rule predicts 0 (which is an error according to  $\sigma_* = \sigma_1$ ) and errs, therefore the rule makes an error on every round, completing the proof.  $\square$

**Remark E.1.** *Note that the rule COMMON-INTERSECTION and MAJORITY respectively recover the textbook rules Consistent and Halving in the standard realizable online classification Shalev-Shwartz & Ben-David (2014). However, both the rules have a mistake bound of  $\Omega(|\mathcal{S}|)$  in our setup even when the labels are binary in the worst case (cf. Theorem 10). This is in sharp contrast with the standard classification where Halving enjoys  $\log_2 |\mathcal{H}|$  mistake bound. This failure is due to the set-valued nature of the support functions.*

We now show the statistical lower bound in a similar spirit.

**Theorem 11** (Statistical Lower Bounds for COMMON-INTERSECTION and MAJORITY). *For every  $d$ , there exists a problem instance  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  with  $|\mathcal{S}| = d$ ,  $|\mathcal{X}| = \lfloor (d-1)/2 \rfloor$ ,  $|\mathcal{Y}| = 2$  and some choice of realizable joint distribution  $(\mathcal{D} \times \pi_*)$  where  $\pi_*$  is supported on  $\sigma_* \in \mathcal{S}$  such that for any sample size  $m \leq |\mathcal{X}|/2$ , letting  $\hat{\pi}(S)$  to be either COMMON-INTERSECTION or MAJORITY have the following guarantee:*

$$\mathbb{P}_{S \sim \mathcal{D}^m} (L_{\mathcal{D}, \sigma_*}(\hat{\pi}(S)) \geq 1/2) = 1.$$

*Proof of Theorem 11.* It suffices to prove the claim for MAJORITY again; the bound for the other follows since MAJORITY is just a special instantiation of COMMON-INTERSECTION.

Consider the same class  $\mathcal{S}$  constructed above with  $|\mathcal{X}| = q$  and ground truth  $\sigma_* = \sigma_1$  (so the realizable label is always 1). Let  $\mathcal{D}$  be the uniform distribution on  $\mathcal{X}$ . Take  $\pi_*$  to be the only conditional distribution supported on  $\sigma_*$ , so the joint  $(\mathcal{D} \times \pi_*)$  is realizable.

Fix any sample size  $m \leq q/2$  and draw  $S \sim \mathcal{D}^m$ . Let  $S_{\text{unseen}} \subseteq \mathcal{X}$  be the set of coordinates *unseen* in  $S$ ; then  $|S_{\text{unseen}}| \geq q - m \geq q/2$ . Let  $V(S) \subseteq \mathcal{S}$  be the version space of hypotheses consistent with  $S$  (with respect to the labels of  $\sigma_*$  label, which are always 1).

Again, by construction, for each  $t \in S_{\text{unseen}}$  there are two designated adversarial hypotheses in  $V(S)$ . At such a point  $t \in S_{\text{unseen}}$ :

- Every hypothesis in  $V(S)$  includes label 0, except  $\sigma_*$ , so

$$N_0(t; V(S)) = |V(S)| - 1.$$

- Every hypothesis in  $V(S)$  includes label 1, except  $\sigma_{2t}, \sigma_{2t+1}$ , so

$$N_1(t; V(S)) \leq |V(S)| - 2.$$

Thus  $N_0(t; V(S)) > N_1(t; V(S))$ , and the majority rule outputs 0 (which is an error according to  $\sigma_*$ ). Thus,  $\sigma_* = \sigma_1$  is  $y_t = 1$ , MAJORITY errs on every unseen  $t \in S_{\text{unseen}}$ .

With  $\mathcal{D}$  uniform on  $\mathcal{X}$ ,

$$L_{\mathcal{D}, \sigma_*}(\hat{\pi}(S)) \geq \mathbb{P}_{x \sim \mathcal{D}}[x \in S_{\text{unseen}}] = \frac{|S_{\text{unseen}}|}{q} \geq 1 - \frac{m}{q} \geq \frac{1}{2}.$$

Since this lower bound holds for every realization of  $S$  with  $m \leq q/2$ , we have  $\mathbb{P}_{S \sim (\mathcal{D} \times \pi_*)^m} (L_{\mathcal{D}, \sigma_*}(\hat{\pi}(S)) \geq \frac{1}{2}) = 1$ . This proves the theorem.  $\square$

Because  $d = (2k)^q$ , we have  $q = \log_{2k} d$ , so the bound implies  $m = \Omega(\log_k d)$  under  $\mathcal{D} = \text{Unif}(\mathcal{X})$ .

**Remark E.2.** *Both online (Theorem 15) and statistical lower bounds (Theorem 15) for pass@k essentially demonstrate that one cannot do better than memorization below  $\Omega(\log_k d)$  barrier in the worst-case, even for the special case of the problem of multiclass classification  $\mathcal{S} \subseteq \mathcal{Y}^{\mathcal{X}}$  which is isomorphic to  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  with  $|\sigma(x)| = 1$ .*

## F ALGORITHM AND ANALYSIS FOR LEARNING FROM ARBITRARY DEMONSTRATOR

We will first begin by online learner which makes a softer update, and then show how the online-to-batch conversion enjoys the guarantee given in Theorem 5.

### F.1 ONLINE LEARNER WITH SOFTER UPDATE

Our goal is to bound the **regret** against the best hypothesis in hindsight in a certain sense. Let's define two types of mistakes over  $T$  rounds for any given hypothesis  $\sigma \in \mathcal{S}$ :

**Algorithm 2** Softer MISTAKE-UNAWARE-WEIGHT-UPDATE

**Input:** A finite support class  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  and the penalty parameters  $\alpha, \beta > 1$ .

- Initialize  $w^{(1)}(\sigma) = 1$  for all  $\sigma \in \mathcal{S}$ .
- In every round  $t = 1, \dots, T$ , receiving  $x_t$ :
  1. Output

$$\hat{y}_t = \arg \max_{y \in \mathcal{Y}} \sum_{\sigma \in \mathcal{S}} w^{(t)}(\sigma) \mathbf{1}\{y \in \sigma(x_t)\}.$$

2. On receiving  $y_t$ , update the weights:

$$w^{(t+1)}(\sigma) \leftarrow w^{(t)}(\sigma) \cdot \alpha^{\mathbf{1}\{\hat{y}_t \notin \sigma(x_t)\}} \cdot \beta^{-\mathbf{1}\{y_t \notin \sigma(x_t)\}} \quad \forall \sigma \in \mathcal{S}.$$

Here,  $\alpha > 1$  and  $\beta > 1$  are penalty parameters. Note that this update rule ensures that after  $T$  rounds, the weight of any hypothesis  $\sigma$  is exactly  $w^{(T+1)}(\sigma) = \alpha^{M_T^{\text{alg}}(\sigma)} \beta^{-M_T(\sigma)}$ .

- **Algorithm's Mistakes relative to  $\sigma$ :** This is the number of times our algorithm's prediction  $\hat{y}_t$  was not in  $\sigma$ 's valid set absolutely, irrespective of what  $y_t$  was.

$$M_T^{\text{alg}}(\sigma) := \sum_{t=1}^T \mathbf{1}[\hat{y}_t \notin \sigma(x_t)]$$

- **Hypothesis  $\sigma$ 's Mistakes:** This is the number of times the true label  $y_t$  was not in  $\sigma$ 's valid set.

$$M_T(\sigma) := \sum_{t=1}^T \mathbf{1}[y_t \notin \sigma(x_t)]$$

The total number of mistakes made by our algorithm with respect to  $\sigma$  is  $M_T^{\text{alg}}(\sigma)$ , and the number of mistakes made by  $\sigma$  itself is  $M_T(\sigma)$ . Note that in some rounds, both types of mistake could occur. We now analyze the algorithm with the first simple lemma that the total weight of the hypotheses in the system is again non-increasing under a minor condition on  $(\alpha, \beta)$ .

**Lemma 3** (Non-increasing Total Weight). *As long as  $\alpha, \beta > 1$  satisfies  $\alpha \leq 2 - 1/\beta$ , the total weight  $\{W_t\}_t$  in the system is non-increasing.*

*Proof of Lemma 3.* Let  $W_t = \sum_{\sigma \in \mathcal{S}} w^{(t)}(\sigma)$  be the total weight at the start of round  $t$ . Define the following disjoint sets of hypotheses:

$$\begin{aligned} S_1 &= \{\sigma \in \mathcal{S} \mid y_t \in \sigma(x_t), \hat{y}_t \in \sigma(x_t)\}, \\ S_2 &= \{\sigma \in \mathcal{S} \mid y_t \in \sigma(x_t), \hat{y}_t \notin \sigma(x_t)\}, \\ S_3 &= \{\sigma \in \mathcal{S} \mid y_t \notin \sigma(x_t), \hat{y}_t \in \sigma(x_t)\}, \\ S_4 &= \{\sigma \in \mathcal{S} \mid y_t \notin \sigma(x_t), \hat{y}_t \notin \sigma(x_t)\}. \end{aligned}$$

By the update rule

$$w^{(t+1)}(\sigma) = w^{(t)}(\sigma) \alpha^{\mathbf{1}\{\hat{y}_t \notin \sigma(x_t)\}} \beta^{-\mathbf{1}\{y_t \notin \sigma(x_t)\}},$$

we have:

$$w^{(t+1)}(\sigma) = w^{(t)}(\sigma) \begin{cases} 1, & \sigma \in S_1, \\ \alpha, & \sigma \in S_2, \\ 1/\beta, & \sigma \in S_3, \\ \alpha/\beta, & \sigma \in S_4. \end{cases}$$

Therefore:

$$W_{t+1} = W_t(S_1) + \alpha W_t(S_2) + \frac{1}{\beta} W_t(S_3) + \frac{\alpha}{\beta} W_t(S_4).$$

The change in total weight is:

$$W_t - W_{t+1} = (1 - \alpha)W_t(S_2) + \left(1 - \frac{1}{\beta}\right)W_t(S_3) + \left(1 - \frac{\alpha}{\beta}\right)W_t(S_4).$$

1242 By the prediction rule:  
1243

$$1244 \quad W_t(y_t) \leq W_t(\hat{y}_t) \quad \Rightarrow \quad W_t(S_2) \leq W_t(S_3).$$

1245  
1246 Thus:

$$1247 \quad W_t - W_{t+1} \geq (1 - \alpha)W_t(S_2) + \left(1 - \frac{1}{\beta}\right)W_t(S_2) + \left(1 - \frac{\alpha}{\beta}\right)W_t(S_4).$$

1248  
1249 A sufficient condition for  $W_{t+1} \leq W_t$  is:

$$1250 \quad 1 - \alpha + 1 - \frac{1}{\beta} \geq 0 \quad \text{and} \quad 1 - \frac{\alpha}{\beta} \geq 0,$$

1251  
1252 which simplifies to:

$$1253 \quad \alpha \leq 2 - \frac{1}{\beta} \quad \text{and} \quad \alpha \leq \beta.$$

1254  
1255 Since  $\alpha \leq \beta$  is implied by the first condition in relevant regimes, which completes the proof.  $\square$   
1256

1257  
1258 **Note that the realizable case is simply  $\beta = \infty$  and  $\alpha = 2$ .**  
1259

1260  
1261 **Final Regret Bound** We are now ready to bound a certain type of regret of the algorithm.

1262 **Theorem 12.** *Let  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$  be any finite model class, and let  $\{w^{(t)}(\sigma)\}$  be the weights generated  
1263 by Algorithm 2 when run with  $\alpha, \beta > 1$  satisfying  $\alpha \leq 2 - 1/\beta$ . Then for any  $\sigma \in \mathcal{S}$ , we have*  
1264

$$1265 \quad M_T^{\text{alg}}(\sigma) \leq \frac{\log |\mathcal{S}|}{\log(\alpha)} + M_T(\sigma) \cdot \frac{\log(\beta)}{\log(\alpha)}.$$

1266  
1267 *Proof.* Let  $W_t = \sum_{\sigma \in \mathcal{S}} w^{(t)}(\sigma)$  be the total weight at round  $t$ . By the update rule of Algorithm 2,  
1268 we have for any  $\sigma \in \mathcal{S}$ :

$$1269 \quad w^{(T+1)}(\sigma) = \alpha^{M_T^{\text{alg}}(\sigma)} \beta^{-M_T(\sigma)}.$$

1270  
1271 Since  $\alpha, \beta > 1$  satisfy  $\alpha \leq 2 - 1/\beta$ , Lemma 3 ensures that  $\{W_t\}$  is non-increasing, hence:

$$1272 \quad W_{T+1} \leq W_1 = |\mathcal{S}|.$$

1273  
1274 Therefore,

$$1275 \quad \alpha^{M_T^{\text{alg}}(\sigma)} \beta^{-M_T(\sigma)} = w^{(T+1)}(\sigma) \leq W_{T+1} \leq |\mathcal{S}|,$$

1276  
1277 and taking logarithms yields:

$$1278 \quad M_T^{\text{alg}}(\sigma) \log(\alpha) - M_T(\sigma) \log(\beta) \leq \log |\mathcal{S}|,$$

1279  
1280 which rearranges to:

$$1281 \quad M_T^{\text{alg}}(\sigma) \leq \frac{\log |\mathcal{S}|}{\log(\alpha)} + M_T(\sigma) \cdot \frac{\log(\beta)}{\log(\alpha)}.$$

1282  
1283  $\square$   
1284

## 1285 F.2 ONLINE-TO-BATCH CONVERSION AND STATISTICAL GUARANTEE

1286  
1287 Our learning rule for the statistical setting is again the same as Eq. (5), but by executing the softer  
1288 MISTAKE-UNAWARE-WEIGHT-UPDATE (Algorithm 2) with choice of  $(\alpha, \beta) = (4/3, 3/2)$ . We  
1289 again denote this predictor by  $\hat{\pi}_{\text{o2b}}$ . We show that Theorem 5 holds when we use this as an estimator.

1290 *Proof of Theorem 5.* Let  $\kappa := \log \beta / \log \alpha \leq 1.41$  and consider any fixed but unknown reference  
1291  $\sigma_*$  and the demonstrator  $\pi_*$  (which may not be optimal). Define

$$1292 \quad Z_t := L_{\mathcal{D}, \sigma_*}(\hat{\pi}_t) - \mathbb{1}\{\hat{\pi}_t(x_t) \notin \sigma_*(x_t)\} + \kappa (L_{\mathcal{D}, \sigma_*}(\pi_*) - \mathbb{1}\{y_t \notin \sigma_*(x_t)\}),$$

1293  
1294 Because  $\hat{\pi}_t$  is deterministic given  $S_{<t} = \{(x_i, y_i) : i < t\}$ , we have

$$1295 \quad \mathbb{E}[Z_t \mid S_{<t}] = L_{\mathcal{D}, \sigma_*}(\hat{\pi}_t) - L_{\mathcal{D}, \sigma_*}(\hat{\pi}_t) + \kappa (L_{\mathcal{D}, \sigma_*}(\pi_*) - L_{\mathcal{D}, \sigma_*}(\pi_*)) = 0,$$

so  $\{Z_t\}_{t=1}^m$  forms a martingale difference sequence. Moreover,  $|Z_t| \leq 1 + \kappa \leq 3$  almost surely.

We have

$$\mathbb{E}[Z_t | S_{<t}] = 0, \quad \text{Var}(Z_t | S_{<t}) \leq 2L_{\mathcal{D},\sigma_*}(\hat{\pi}_t) + 2\kappa^2 L_{\mathcal{D},\sigma_*}(\pi_*) \leq 2L_{\mathcal{D},\sigma_*}(\hat{\pi}_t) + 4L_{\mathcal{D},\sigma_*}(\pi_*).$$

Applying Freedman's inequality (Lemma 2) with  $W_m = 2\sum_{t=1}^m L_{\mathcal{D},\sigma_*}(\hat{\pi}_t) + 4\sum_{t=1}^m L_{\mathcal{D},\sigma}(\pi_*)$  and  $\sigma^2 = 6m$  yields that with probability  $1 - \delta$ :

$$\sum_{t=1}^m Z_t \leq \sqrt{8(1 + W_m) \log\left(\frac{2 \log 6m}{\delta}\right)} + \frac{4}{3}(1 + \kappa) \log\left(\frac{2 \log 6m}{\delta}\right).$$

By AM–GM inequality:

$$\sum_{t=1}^m Z_t \leq \frac{(1 + W_m)}{8} + 16 \log\left(\frac{2 \log 6m}{\delta}\right) + 4 \log\left(\frac{2 \log 6m}{\delta}\right).$$

Substituting  $Z_t$  and rearranging gives us:

$$\sum_{t=1}^m L_{\mathcal{D},\sigma_*}(\hat{\pi}_t) \leq \kappa L_{\mathcal{D},\sigma_*}(\pi_*) + 1 + 2 \left( \sum_{t=1}^m \mathbf{1}\{\hat{\pi}_t(x_t) \notin \sigma_*(x_t)\} - \kappa \sum_{t=1}^m \mathbf{1}\{y_t \notin \sigma_*(x_t)\} \right) + 40 \log\left(\frac{2 \log 6m}{\delta}\right),$$

Finally, using Theorem 12, and using the fact that

$$L_{\mathcal{D},\sigma_*}(\hat{\pi}_{\text{o2b}}(S)) = \frac{1}{m} \sum_{t=1}^m L_{\mathcal{D},\sigma_*}(\hat{\pi}_t),$$

we obtain that, with probability at least  $1 - \delta$ ,

$$L_{\mathcal{D},\sigma_*}(\hat{\pi}_{\text{o2b}}(S)) \leq \kappa L_{\mathcal{D},\sigma_*}(\pi_*) + \frac{1 + 6 \log_2 |S| + 40 \log\left(\frac{2 \log 6m}{\delta}\right)}{m}.$$

□

## G ALGORITHMS AND PROOFS FOR $\text{pass@}k$ -ERROR

In this, we provide our guarantees for  $\text{pass@}k$  error (and formalize Theorem 6). We first start by describing an online learner.

---

### Algorithm 3 Online $\text{pass@}k$ rule with greedy selection and mistake unaware updates

---

**Input:** A finite model class  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$ , and parameter  $k \in \mathbb{N}$ .

- Initialize  $V_1 = \mathcal{S}$  and  $w^{(1)}(\sigma) = 1$  for all  $\sigma \in \mathcal{S}$ .
- In every round, receiving  $x_t$ :
  1. For each  $y \in \mathcal{Y}$ , form the slice  $A_y^t = \{\sigma \in V_t : y \in \sigma(x_t)\}$ .
  2. (*Greedy top- $k$  selection*). Let  $\mathcal{Y}_0 = \emptyset$ . For  $i = 1, 2, \dots, k$  set

$$\hat{y}_t^{(i)} \in \arg \max_{y \in \mathcal{Y} \setminus \mathcal{Y}_{i-1}} w^{(t)} \left( A_y^t \setminus \bigcup_{z \in \mathcal{Y}_{i-1}} A_z^t \right), \quad \mathcal{Y}_i \leftarrow \mathcal{Y}_{i-1} \cup \{\hat{y}_t^{(i)}\}.$$

(Break ties arbitrarily.) Define  $U_t := \bigcup_{i=1}^k A_{\hat{y}_t^{(i)}}^t$ .

3. Output the  $k$  labels  $\hat{y}_t^{(1)}, \dots, \hat{y}_t^{(k)}$ .
4. (*Weight update*). Upon receiving label  $y_t$ :

$$w^{(t+1)}(\sigma) \leftarrow \begin{cases} 0, & \text{for all } \sigma \notin A_{y_t}^t; \\ w^{(t)}(\sigma) & \text{for all } \sigma \in A_{y_t}^t \cap U_t; \\ (k+1)w^{(t)}(\sigma), & \text{for all } \sigma \in A_{y_t}^t \setminus U_t. \end{cases}$$

So the version space:  $V_{t+1} \leftarrow A_{y_t}^t$ .

---

**Theorem 13** (Online pass@ $k$  guarantee). *On any sequence  $((x_t, y_t))_{t \in \mathbb{N}}$  realizable by some  $\sigma_* \in \mathcal{S}$ , Algorithm 3 makes at most  $\log_{k+1} |\mathcal{S}|$  mistakes (i.e., rounds with  $\{\hat{y}_t^{(1)}, \dots, \hat{y}_t^{(k)}\} \cap \sigma_*(x_t) = \emptyset$ ).*

*Proof.* Let  $V_t$  be the version space at the start of round  $t$ , and as in the algorithm write  $A_{y_t}^t = \{\sigma \in V_t : y \in \sigma(x_t)\}$  and  $U_t := \bigcup_{i=1}^k A_{\hat{y}_t^{(i)}}^t$ . Define the potential  $W_t := w^{(t)}(V_t) = \sum_{\sigma \in V_t} w^{(t)}(\sigma)$ , then we again have  $\{W_t\}_t$  is non-increasing.

$$\begin{aligned} W_{t+1} &= (k+1)w^{(t)}(A_{y_t}^t \setminus U_t) + w^{(t)}(A_{y_t}^t) = kw^{(t)}(A_{y_t}^t \setminus U_t) + w^{(t)}(A_{y_t}^t \setminus U_t) + w^{(t)}(A_{y_t}^t \cap U_t) \\ &\leq w^{(t)}(U_t \setminus A_{y_t}^t) + w^{(t)}(A_{y_t}^t \setminus U_t) + w^{(t)}(A_{y_t}^t \cap U_t) \\ &= w^{(t)}(U_t \cup A_{y_t}^t) \leq W_t, \end{aligned}$$

where in the first inequality arises from the following key relation:

$$w^{(t)}(U_t \setminus A_{y_t}^t) \geq kw^{(t)}(A_{y_t}^t \setminus U_t). \quad (13)$$

We defer its proof to the end of this section.

Now suppose the algorithm makes  $M$  pass@ $k$  mistakes by the end of round  $t$ . On each mistake round we must have  $\sigma_* \in A_{y_t}^t \setminus U_t$ , so its weight is multiplied by  $(k+1)$ . Therefore

$$w^{(t+1)}(\sigma_*) = (k+1)^M \leq W_{t+1} \leq W_1 = |\mathcal{S}|,$$

which yields  $M \leq \log_{k+1} |\mathcal{S}|$ .  $\square$

**Proof of the key inequality (13).** To simplify the notation, we remove the time index  $t$  and write  $U_t$  as  $U$ , and  $A_{y_t}^t$  as  $A_{y_t}$ . Define the uncovered mass in  $A$  after selecting first  $i$  labels greedily as:

$$a_i := w^{(t)}\left(A_{y_t} \setminus \bigcup_{z \in \mathcal{Y}_i} A_z\right) \quad \text{so that} \quad a_0 = w^{(t)}(A_{y_t}), \quad a_k = w^{(t)}(A_{y_t} \setminus U),$$

and  $a_0 \geq a_1 \geq \dots \geq a_k$ . Correspondingly, define

$$s_i := a_{i-1} - a_i = w^{(t)}\left((A_{\hat{y}_t^{(i)}} \cap A_{y_t}) \setminus \bigcup_{z \in \mathcal{Y}_{i-1}} A_z\right).$$

In addition, we define the uncovered weight for which  $\hat{y}_t^{(i)}$  got picked as

$$m_i := w^{(t)}\left(A_{\hat{y}_t^{(i)}} \setminus \bigcup_{z \in \mathcal{Y}_{i-1}} A_z\right).$$

By the design of the greedy selections, we have

$$m_i \geq a_{i-1} \quad \text{for all } i \in [k]. \quad (14)$$

With these definitions and relations in place, we are ready to prove the desired inequality. Basic set inclusion / exclusion tells us that

$$\begin{aligned} w^{(t)}(U \setminus A_{y_t}) &= \sum_{i=1}^k w^{(t)}(A_{\hat{y}_t^{(i)}} \setminus A_{y_t} \cup A_{\hat{y}_t^{(1)}} \cdots \cup A_{\hat{y}_t^{(i-1)}}) \\ &= \sum_{i=1}^k \left\{ w^{(t)}(A_{\hat{y}_t^{(i)}} \setminus \bigcup_{z \in \mathcal{Y}_{i-1}} A_z) - w^{(t)}\left((A_{\hat{y}_t^{(i)}} \cap A_{y_t}) \setminus \bigcup_{z \in \mathcal{Y}_{i-1}} A_z\right) \right\} \\ &= \sum_{i=1}^k (m_i - s_i), \end{aligned}$$

where the last equality uses the definitions of  $m_i$  and  $s_i$ . Using (14), we obtain

$$m_i - s_i \geq a_{i-1} - (a_{i-1} - a_i) = a_i,$$

which allows us to further lower bound  $w^{(t)}(U \setminus A)$  as

$$w^{(t)}(U \setminus A_{y_t}) \geq \sum_{i=1}^k a_i \geq k a_k.$$

Here we use the fact that  $(a_i)_i$  is non-increasing. This proves the claim.

## G.1 STATISTICAL UPPER BOUND

**Input:** Sample  $S = \{(x_i, y_i) : i \in [m]\}$  and a finite model class  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$ .

- Run Algorithm 3 once over  $S$ , recording  $(V_t, w^{(t)})$  at the *start* of each round  $t \in [m]$ .
- Find the deterministic predictor  $\hat{\mu}_t : \mathcal{X} \rightarrow \mathcal{Y}^k$  used by the online algorithm from the snapshot. I.e. for any  $x \in \mathcal{X}$  and any  $t \in [m]$ , define slices (with respect to  $V_t$ )

$$A_y^t(x) := \{\sigma \in V_t : y \in \sigma(x)\}.$$

And greedily pick top  $k$  labels according to the rule described in Algorithm 3.  
Let  $\mathcal{Y}_0(x) = \emptyset$ . For  $i = 1, \dots, k$  set

$$\hat{y}_t^{(i)}(x) \in \arg \max_{y \in \mathcal{Y} \setminus \mathcal{Y}_{i-1}(x)} w^{(t)} \left( A_y^t(x) \setminus \bigcup_{z \in \mathcal{Y}_{i-1}(x)} A_z^t(x) \right)$$

$$\mathcal{Y}_i(x) \leftarrow \mathcal{Y}_{i-1}(x) \cup \{\hat{y}_t^{(i)}(x)\},$$

breaking ties by arbitrary fixed rule.

- Then the deterministic predictor  $\hat{\mu}_t : \mathcal{X} \rightarrow \mathcal{Y}^k$  is given by:

$$\hat{\mu}_t(x) := (\hat{y}_t^{(1)}(x), \dots, \hat{y}_t^{(k)}(x)).$$

- **Final batch predictor.** On a test input  $x$ , draw  $I \sim \text{Unif}\{1, \dots, m\}$  and output

$$\hat{\mu}_{\text{o2b}}(S)(x) := \hat{\mu}_I(x). \quad (15)$$

(Equivalently:  $\hat{\mu}_{\text{o2b}}$  is the uniform mixture over  $\{\hat{\mu}_t\}_{t=1}^m$  for every test point.)

Below is the statistical guarantee for this estimator in similar spirit to Theorem 4.

**Theorem 14** (Statistical Guarantee for pass@k). *For any finite model class  $\mathcal{S}$ , the estimator  $\hat{\mu}_{\text{o2b}}$  in Eq. (15) achieves the following guarantee for any joint distribution  $(\mathcal{D} \times \pi_*)$  supported on  $\sigma_* \in \mathcal{S}$ :*

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}, \sigma_*}(\hat{\mu}_{\text{o2b}}(S))] \leq \frac{\log_{k+1} |\mathcal{S}|}{m},$$

and, for any  $\delta \in (0, 1)$ , with probability at least  $1 - \delta$ ,

$$L_{\mathcal{D}, \sigma_*}(\hat{\mu}_{\text{o2b}}(S)) \leq \frac{1 + 2 \log_{k+1} |\mathcal{S}| + 12 \log\left(\frac{\log m}{\delta}\right)}{m}.$$

This implies that  $\mathcal{S}$  is learnable (cf. Definition 1) using the estimator  $\hat{\mu}_{\text{o2b}}$  with sample complexity

$$m_{\mathcal{S}, \hat{\mu}_{\text{o2b}}}(\varepsilon, \delta) = O\left(\frac{1}{\varepsilon} \left(\log_{k+1} |\mathcal{S}| + \log\left(\frac{1}{\varepsilon \delta}\right)\right)\right).$$

*Proof of Theorem 14.* The proof is exactly similar to that of Theorem 4 and given for completeness.

Let  $\ell_t = \mathbb{1}\{\hat{y}_t^{(i)}(x_t) \notin \sigma_*(x_t), \forall \hat{y}_t^{(i)}(x_t) \in \hat{\mu}_t(x_t)\}$ . Because  $\hat{\mu}_t$  is a deterministic function of  $S_{<t} = \{(x_i, y_i) : i < t\}$ , we have

$$\mathbb{E}[\ell_t \mid S_{<t}] = L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t).$$

Hence

$$\mathbb{E}_S [L_{\mathcal{D}, \sigma_*}(\hat{\mu}_{\text{o2b}}(S))] = \mathbb{E}_S \left[ \frac{1}{m} \sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t) \right] = \mathbb{E}_S \left[ \frac{1}{m} \sum_{t=1}^m \ell_t \right] \leq \frac{\log_{k+1} |\mathcal{S}|}{m},$$

where in the last inequality we used Theorem 13, which guarantees  $\sum_{t=1}^m \ell_t \leq \log_{k+1} |\mathcal{S}|$ .

For the high-probability statement, define the martingale differences

$$Z_t := L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t) - \ell_t, \quad \text{where } |Z_t| \leq 1 \text{ almost surely.}$$

Then  $\mathbb{E}[Z_t \mid S_{<t}] = 0$ , and

$$\mathbb{E}[Z_t^2 \mid S_{<t}] = \mathbb{E}[(L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t) - \ell_t)^2 \mid S_{<t}] = \text{Var}(\ell_t \mid S_{<t}) = L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t)(1 - L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t)) \leq L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t).$$

Taking  $W_m = \sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t)$  and  $\sigma^2 = m$  suffices; thus, using Lemma 2 with  $n = \log m$  gives, with probability  $1 - \delta$ ,

$$\begin{aligned} \sum_{t=1}^m Z_t &\leq \sqrt{8 \left(1 + \sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t)\right) \log\left(\frac{\log m}{\delta}\right) + \frac{4}{3} \log\left(\frac{\log m}{\delta}\right)} \\ &\leq \frac{1}{2} \left(1 + \sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t)\right) + 4 \log\left(\frac{\log m}{\delta}\right) + \frac{4}{3} \log\left(\frac{\log m}{\delta}\right) \quad (\text{GM} \leq \text{AM}) \end{aligned}$$

Substituting  $Z_t$  and rearranging terms,

$$\sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t) \leq 1 + 2 \sum_{t=1}^m \ell_t + \frac{32}{3} \log\left(\frac{\log m}{\delta}\right).$$

Finally, noting that  $L_{\mathcal{D}, \sigma_*}(\hat{\mu}_{\text{o2b}}(S)) = \frac{1}{m} \sum_{t=1}^m L_{\mathcal{D}, \sigma_*}(\hat{\mu}_t)$  and that  $\sum_{t=1}^m \ell_t \leq \log_{k+1} |\mathcal{S}|$  (by Theorem 13), we obtain that with probability  $1 - \delta$ ,

$$L_{\mathcal{D}, \sigma_*}(\hat{\mu}_{\text{o2b}}(S)) \leq \frac{1 + 2 \log_{k+1} |\mathcal{S}| + 12 \log\left(\frac{\log m}{\delta}\right)}{m}.$$

□

## G.2 LOWER BOUNDS FOR ONLINE AND STATISTICAL SETTINGS FOR $k$ -PASS ERROR

We next provide a lower bound that, information-theoretically, this dependence cannot be improved and we only gain a factor of  $1/\log k$  in sample complexity as well as mistake bound in the worst-case.

**Theorem 15** (Online  $\Omega(\log_{k+1} |\mathcal{S}|)$  pass@ $k$  mistake bound even for multiclass classification). *Fix integers  $k \geq 1$  and  $d \geq 2$ . There exists a problem instance  $\mathcal{S} \subseteq \mathcal{Y}^{\mathcal{X}}$  with  $|\mathcal{S}| \leq d$ ,  $|\mathcal{Y}| = k+1$ ,  $|\mathcal{X}| = \lfloor \log_{k+1} d \rfloor$  such that for any deterministic online learning algorithm that outputs at most  $k$  labels, there exists a sequence  $(x_t, y_t)_{t \in [|\mathcal{X}|]}$  realizable by some  $\sigma_* \in \mathcal{S}$  such that it makes mistake on every round.*

Note that our instance is an instance of multiclass classification problem  $\sigma : \mathcal{X} \rightarrow \mathcal{Y}$ . This is isomorphic to an instance  $\mathcal{S} \subseteq (2^{\mathcal{Y}})^{\mathcal{X}}$ , where  $|\sigma(x)| = 1$  for all  $x \in \mathcal{X}, \sigma \in \mathcal{S}$ .

*Proof of Theorem 15.* Let  $m := \lfloor \log_{k+1} d \rfloor$  and take  $\mathcal{X} = \{1, \dots, m\}$  and  $\mathcal{Y} = \{1, \dots, k+1\}$ . Consider the full product class  $\mathcal{S} = \mathcal{Y}^{\mathcal{X}}$ , which has size  $(k+1)^m \leq d$ . First of all, observe that in any round in which  $y_t$  does not belong to the list of  $(\hat{y}_t^{(1)}, \dots, \hat{y}_t^{(k)})$ , the mistake is made because we are in the multiclass classification setting.

For rounds  $t \in [m]$ , present a fresh coordinate  $x_t = t$ . Since  $|\mathcal{Y}| = k+1$ , there exists a label  $y_t \in \mathcal{Y}$  that the learner failed to output in the set;  $y_t \neq \hat{y}_t^{(i)}$  for all  $i \in [k]$ . Reveal this  $y_t$ . This forces a mistake on every round. Moreover, this sequence is realizable since  $\mathcal{S} = \mathcal{Y}^{\mathcal{X}}$  contains all functions from  $\mathcal{X}$  to  $\mathcal{Y}$ . □

**Theorem 16** (Statistical lower bound of  $\Omega(\log_k |\mathcal{S}|)$ ). *Fix integers  $k \geq 1$  and  $q \geq 1$ . Let  $\mathcal{X} = \{1, \dots, q\}$ ,  $\mathcal{Y} = \{1, \dots, 2k\}$ , and take the hypothesis class  $\mathcal{S} = \mathcal{Y}^{\mathcal{X}}$  (all multiclass functions), so its cardinality is  $d := |\mathcal{S}| = (2k)^q$ . Let  $\mathcal{D}$  be the uniform distribution on  $\mathcal{X}$ . Then for any estimators  $\hat{\mu} : (\mathcal{X} \times \mathcal{Y})^* \rightarrow \Delta(\mathcal{Y}^k)^{\mathcal{X}}$*

$$\inf_{\hat{\mu}} \sup_{\sigma \in \mathcal{S}} \mathbb{E}_{S \sim (\mathcal{D} \times \sigma)^m} \mathbb{P}_{x \sim \mathcal{D}} \mathbb{P}_{\hat{y}(x) \sim \hat{\mu}(\cdot|x)} [\sigma(x) \notin \hat{y}(x)] \geq \frac{1}{2} \left(1 - \frac{1}{q}\right)^m.$$

*In particular, to ensure expected error at most  $0 < \varepsilon < \frac{1}{2}$  for all  $\sigma \in \mathcal{S}$ , one needs*

$$m \geq \frac{\ln(1/(2\varepsilon))}{-\ln(1-1/q)} \geq q \ln\left(\frac{1}{2\varepsilon}\right).$$

1512 *Proof.* Fix any (possibly randomized) estimator  $\hat{\mu}$ . Let  $S = \{(x_i, y_i)\}_{i=1}^m$  be the training sample  
 1513 drawn i.i.d. from  $(\mathcal{D} \times \sigma)$  for  $\sigma \sim \text{Unif}(\mathcal{S})$ , and let  $U_S = \{x_i : 1 \leq i \leq m\} \subseteq \mathcal{X}$  be the set of  
 1514 distinct inputs seen in  $S$ . Draw  $x \sim \mathcal{D}$  independently of  $S$  and then  $\hat{\mathbf{y}}(x) \sim \hat{\mu}(S)(\cdot | x)$ .

1515 On any  $x \notin U_S$ , under the prior where  $\sigma$  is uniform over  $\mathcal{S}$ , for any (possibly randomized)  $k$ -list  
 1516  $\hat{\mathbf{y}}(x) \sim \hat{\mu}(\cdot | x)$ ,

$$1517 \mathbb{P}_\sigma[\sigma(x) \in \hat{\mathbf{y}}(x) | S, x \notin U_S] = \mathbb{E} \left[ \frac{\# \text{ of distinct labels in } \hat{\mathbf{y}}(x)}{|\mathcal{Y}|} \mid S, x \notin U_S \right] \leq \frac{k}{2k} = \frac{1}{2},$$

1518 so  $\mathbb{P}_\sigma[\sigma(x) \notin \hat{\mathbf{y}}(x) | S, x \notin U_S] \geq \frac{1}{2}$ . (Allowing duplicates in  $\hat{\mathbf{y}}(x)$  cannot decrease this probabili-  
 1519 ty.)

1520 If  $x \in U_S$ , the learner can always include the observed label and incur zero error on that  $x$ . There-  
 1521 fore, for any estimator  $\hat{\mu}$ ,

$$1522 \mathbb{P}_{\sigma, x, \hat{\mathbf{y}}(x) \sim \hat{\mu}(\cdot | x)}[\sigma(x) \notin \hat{\mathbf{y}}(x) | S] \geq \frac{1}{2} \cdot \mathbb{P}[x \notin U_S].$$

1523 Taking expectation over  $S$  and using  $\mathcal{D} = \text{Unif}(\mathcal{X})$  yields

$$1524 \mathbb{E}_S \mathbb{P}_{\sigma, x, \hat{\mathbf{y}}(x) \sim \hat{\mu}(\cdot | x)}[\sigma(x) \notin \hat{\mathbf{y}}(x)] \geq \frac{1}{2} \mathbb{E}_S[1 - |U_S|/q] = \frac{1}{2} \left(1 - \frac{1}{q}\right)^m,$$

1525 since  $\mathbb{E}[|U_S|] = q(1 - (1 - \frac{1}{q})^m)$ . Finally, by minimax principle

$$1526 \inf_{\hat{\mu}} \sup_{\sigma \in \mathcal{S}} \mathbb{E}_S \mathbb{P}_{x \sim \mathcal{D}} \mathbb{P}_{\hat{\mathbf{y}}(x) \sim \hat{\mu}(\cdot | x)}[\sigma(x) \notin \hat{\mathbf{y}}(x)] \geq \inf_{\hat{\mu}} \mathbb{E}_{\sigma \sim \text{Unif}(\mathcal{S})} \mathbb{E}_S \mathbb{P}_{x \sim \mathcal{D}} \mathbb{P}_{\hat{\mathbf{y}}(x) \sim \hat{\mu}(\cdot | x)}[\sigma(x) \notin \hat{\mathbf{y}}(x)]$$

$$1527 \geq \frac{1}{2} \left(1 - \frac{1}{q}\right)^m.$$

1528 For the sample-complexity bound, solve  $\frac{1}{2}(1 - \frac{1}{q})^m \leq \varepsilon$  for  $m$  and use  $-\ln(1 - 1/q) \leq 1/q$ .  $\square$