# CONTEXTUAL BANDITS AND OPTIMISTICALLY UNIVERSAL LEARNING

BY MOÏSE BLANCHARD[*] STEVE HANNEKE[†], PATRICK JAILLET[‡]

[*]*Massachusetts Institute of Technology, moiseb@mit.edu*

[†]*Purdue University, steve.hanneke@gmail.com*

[‡]*Massachusetts Institute of Technology, jaillet@mit.edu*

We consider the contextual bandit problem on general action and context spaces, where the learner's rewards depend on their selected actions and an observable context. This generalizes the standard multi-armed bandit to the case where side information is available, e.g., patients' records or customers' history, which allows for personalized treatment. We focus on consistency—vanishing regret compared to the optimal policy—and show that for large classes of non-i.i.d. contexts, consistency can be achieved regardless of the time-invariant reward mechanism, a property known as *universal consistency*. Precisely, we first give necessary and sufficient conditions on the context-generating process for universal consistency to be possible. Second, we show that there always exists an algorithm that guarantees universal consistency whenever this is achievable, called an *optimistically universal* learning rule. Interestingly, for finite action spaces, learnable processes for universal learning are exactly the same as in the full-feedback setting of supervised learning, previously studied in the literature. In other words, learning can be performed with partial feedback without any generalization cost. The algorithms balance a trade-off between *generalization* (similar to structural risk minimization) and *personalization* (tailoring actions to specific contexts). Lastly, we consider the case of added continuity assumptions on rewards and show that these lead to universal consistency for significantly larger classes of data-generating processes.

**1. Introduction.** The contextual bandits setting is one of the core important problems in sequential statistical decision-making. Abstractly, in the contextual bandit setting, a learner (or decision maker) interacts with a reward mechanism iteratively. At each iteration, the learner observes a *context* (or covariate vector) $x \in \mathcal{X}$ and selects an *arm* (or action) $a \in \mathcal{A}$ to perform; it then receives a (potentially stochastic) reward depending on the context and selected action. For example, a store may serve a sequence of customers, and for each provide a list of product recommendations, and receive reward if the recommendation leads to a purchase. The key distinctions between the contextual bandit setting and standard supervised learning (or regression) are that (1) the learner's objective is to obtain a near-maximum average reward over time (rather than merely estimating the reward conditional means), and (2) the learner only observes the reward corresponding to the arm it chose. These aspects introduce a fundamental trade-off between *exploration* and *exploitation*: that is, while some arms may have high estimated reward values, other arms may have higher uncertainty in their rewards: in particular, uncertainty about whether they would yield an even higher reward, so that selecting that arm may provide information about the potential for higher future rewards.

1.1. *Universal Consistency.* In the contextual bandit setting, a learner is *consistent* if its average reward converges to the maximum-possible average reward obtained with an optimal policy. Naturally, one would aim for learning procedures that ensure consistency for a broad class of problem instances. In particular, we are interested in *universal consistency* which asks that a learning rule achieves consistency for any underlying reward mechanism and as a by-product, any optimal policy. The equivalent notion can be defined for the

full-information case: for a stream of data $(\mathbb{X}, \mathbb{Y}) = (X_t, Y_t)_{t \geq 1}$ of instances modeled as a stochastic process on $\mathcal{X} \times \mathcal{Y}$, a learning rule with predictions $\hat{Y}_t$ is consistent if it has vanishing excess error compared to any fixed measurable predictor function $f : \mathcal{X} \to \mathcal{Y}$, i.e., $\frac{1}{T} \sum_{t=1}^{T} \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \to 0 \ (a.s.)$. Then, an algorithm is universally consistent if it is consistent irrespective of the generating process for the values $\mathbb{Y}$ from the instances $\mathbb{X}$. In this standard full-feedback setting, there are many works establishing universal consistency, beginning with the seminal work of [39] who proved universal consistency for a broad family of *local average estimators*. Later works extended these results, to guarantee *strong* universal consistency (i.e., almost sure convergence), other categories of learning rules, more general conditions on the metric space $\mathcal{X}$, and more general loss functions [12, 17]. More recently, [23, 19, 11] gave minimal assumptions on the space $\mathcal{X}$ for universal consistency—essentially-separable metric spaces. All of these works were restricted to *i.i.d.* data $(X_t, Y_t)_{t \geq 1}$ sampled from a joint distribution on $\mathcal{X} \times \mathcal{Y}$. Some of these results aimed to relax the i.i.d. assumption by considering non-i.i.d. mixing, stationary ergodic data generating processes [29, 18, 17] or satisfying the law of large numbers [28, 14, 38].

1.2. *Optimistically universal learning.* In the present work we pursue a theory of universal consistency under provably-*minimal* assumptions on the sequence of contexts. This type of theory falls into a framework known as *optimistically universal learning*, introduced by [20], that can be succinctly summarized as "learning whenever learning is possible". The idea is to identify the minimal assumption on the data sequence sufficient for universal consistency to be possible. Such an assumption is then both necessary and sufficient, and therefore amounts to merely assuming that universally consistent learning is *possible*: aptly named the *optimist's assumption*. For any given process $\mathbb{X}$ satisfying this minimal assumption, by definition there must exist a universally consistent learning rule. However, the interesting question becomes whether the optimist's assumption alone is sufficient to guarantee universal consistency for some well-designed learning rule: that is, whether there exists a single learning rule that is universally consistent for *every* process $\mathbb{X}$ satisfying the optimist's assumption. Such a learning rule is said to be *optimistically universal*.

1.3. *Optimistically universal learning with full-feedback.* The first general analysis of optimistically universal learning and provably-minimal assumptions for universal consistency in the full-feedback setting was introduced by [20]. He provided general necessary and sufficient conditions for the existence of universally consistent learning rules for inductive learning—where one can only observe a finite amount of data $(X_t, Y_t)_{t \leq n}$ before committing to a prediction rule for the future steps $T \geq n$— and for a slight variation called self-adaptive learning—where the learner only observes a finite amount of values $(X_t, Y_t)_{t \leq n}$ but can continue to update its predictions from the testing observations $X_{n+1}, \ldots, X_T$, that is it continues to learn from test data. Interestingly, while there do not exist optimistically universal inductive learning rules, there do exist explicitly defined optimistically universal self-adaptive learning rules. That work focused mostly on the *noiseless* function learning setting where some unknown function $f^* : \mathcal{X} \to \mathcal{Y}$ defines the values exactly via $Y_t = f(X_t)$. It also left open the question of characterizing universal learning for the standard *online learning* framework, in which the learner can update its predictions from the complete available test and value data $(X_t, Y_t)_{t \leq T}$ [see also 21].

Addressing the online learning problem, in the noiseless setting, [6] provided a simpler characterization and algorithm for unbounded losses while [4, 5] provided a solution for the main case of interest of bounded losses. In particular, while the nearest neighbor algorithm may not be universally consistent even for i.i.d. data [9], for noiseless responses, a simple variant with restricted memory is optimistically universal [4]. For the generic case of

noisy responses, [22] showed that universal learning can be achieved even for arbitrarily dependent responses on large classes of processes. The complete characterization of universal online learning with noise was given in [7], showing that under mild conditions on the value space—including totally-bounded metric spaces—optimistically universal learning is possible for arbitrary or adversarial responses without generalizability cost compared to noiseless responses.

1.4. *Universal learning with partial feedback.* The contextual-bandit formulation was first introduced for one-armed bandits [43, 35] in a rather restricted setting. Since then, progress has been made in the literature investigating stochastic contextual bandits under *parametric* assumptions [42, 24, 13, 8, 1, 32]. In the *non-parametric* setting, significant advances have been made to obtain minimax guarantees under smoothness conditions (e.g. Lipschitz) and margin assumptions [26, 34, 36, 31] with recent refinements including [16, 33].

However, to the best of our knowledge, there are no prior works establishing universal consistency even under all i.i.d. data sequences, i.e., consistency in the non-parametric setting without further assumptions. As such, the present work is also the first to propose such results and corresponding universally consistent learning rules. Closest to this work is the result from [45] which shows that if rewards are continuous in the contexts, strong consistency can be achieved with familiar non-parametric methods, for Euclidean context spaces. This work significantly generalizes this result to unrestricted reward mechanisms, separable metric action and context spaces, and non-i.i.d. data.

Non-i.i.d. data has also been widely studied in the literature. Most relevant to our work are non-i.i.d. generating processes for contexts. Examples include customers' profile distribution, which may change depending on seasonal patterns, or the extension of clinical trials to new populations. In these cases, the distribution of contexts $x$ changes while the underlying conditional distribution remains unchanged, a phenomenon known as *covariate-shift*. Such formalism was adopted in works on domain adaptation for classification [40, 15, 2]. Moreover, several works have also considered distributional shifts in both contexts and responses for bandit problems, in both parametric [3, 27, 44, 10] and non-parametric settings [41].

1.5. *Summary of the present work.* In the present work we study optimistically universal learning in a partially-supervised setting: namely, standard contextual bandits [37, 25] with stationary reward functions. Precisely, there exists a time-invariant conditional probability distribution $P_{r|a,x}$ such that the reward $r_t$ at each iteration is sampled according to the distribution $P_{r|a=a_t,x=X_t}$ where $a_t$ (resp. $X_t$) denotes the selected action (resp. observed context) at time $t$, independently from the past history. We are interested in online learning, where the learner may observe all past rewards $r_{t'}$ and contexts $X_{t'}$, $t' < t$, when choosing its action $a_t$ given the context $X_t$. We aim to achieve average reward $\frac{1}{T}\sum_{t=1}^{T} r_t$ that is (almost surely) competitive with any fixed policy $\mathcal{X} \to \mathcal{A}$ as $T \to \infty$.

1.5.1. *Bounded unrestricted rewards.* We first focus on the classical assumption that rewards are bounded. We show there always exists an optimistically universal learning rule. Our approach to proving this is to first characterize which processes $\mathbb{X}$ admit universally consistent learning rules, and then use this characterization to inform the design and analysis of a learning rule, which will be universally consistent under every such process. However, this approach turns out to require three separate cases: namely, $\mathcal{A}$ finite, $\mathcal{A}$ countably infinite, and $\mathcal{A}$ uncountably infinite. Each of these cases gives rise to a different characterization of the set of processes $\mathbb{X}$ under which universally consistent learning is possible for contextual bandits, a fact which itself is of independent interest. Moreover, each of these sets of processes corresponds to known families of processes from the past literature on optimistically

universal learning. When $\mathcal{A}$ is finite, the set of processes admitting universal consistency for contextual bandits is equivalent to the family of processes admitting universally consistent online learning with full supervision: a family known as $\mathcal{C}_2$. While this fact appears natural, interestingly this is not the case when $\mathcal{A}$ is countably infinite. In that case, the set of processes admitting universal learning for contextual bandits is equivalent to the family of processes admitting universally consistent *inductive* learning with full supervision: a family known as $\mathcal{C}_1$, which is more restrictive than $\mathcal{C}_2$. Finally, when $\mathcal{A}$ is uncountably infinite, universal learning can never be achieved.

1.5.2. *Bounded rewards under continuity assumptions.*  For unrestricted rewards, although large classes of non-i.i.d. processes ($\mathcal{C}_1$ or $\mathcal{C}_2$) admit universal learning for countable action spaces, the answer for uncountable action spaces was very negative: universal consistency could never be achieved. However, we show that under continuity assumptions on the rewards, one can recover positive results for general action spaces. Further, in all cases, we provide optimistically universal learning rules. First, under the assumption that rewards are continuous, the characterization of processes admitting universal consistency now requires only two cases. If the action space is finite, the set of processes admitting universal learning remains unchanged and is $\mathcal{C}_2$. On the other hand, if the action space is infinite, this set becomes $\mathcal{C}_1$, irrespective of whether the action space was countably or uncountably infinite. Second, we consider a stronger assumption of uniform continuity on the rewards, in which the modulus of continuity of the expected reward in the actions $\bar{r}(\cdot, x)$ for $x \in \mathcal{X}$ are uniform over the context space $\mathcal{X}$. Under this assumption, universal learning under the more general set of processes $\mathcal{C}_2$ becomes possible for a significantly larger class of action spaces, namely totally-bounded action spaces. Otherwise, universal learning is achievable exactly on $\mathcal{C}_1$ processes.

1.5.3. *Unbounded rewards.*  Last, we consider the most general case of unbounded rewards. It is known that the family of processes admitting universal consistency with full supervision and *unbounded* losses is very restrictive. These are processes visiting only a finite number of distinct instances in $\mathcal{X}$, known as $\mathcal{C}_3$. For contextual bandits, in the standard case of unrestricted rewards, we show that there is a simple dichotomy: if the action space is countable then the set of processes admitting universal learning is still $\mathcal{C}_3$; however, if the action space is uncountably infinite, universal learning can never be achieved. Nevertheless, under continuity assumptions on the rewards, universal learning can always be achieved under $\mathcal{C}_3$ processes. Again, we give optimistically universal learning rules for all cases.

1.6. *Overview of probability-theoretic contributions.*  In this work, we make use of the conditions $\mathcal{C}_1$, $\mathcal{C}_2$, and $\mathcal{C}_3$ on stochastic processes from the universal learning literature to characterize the set of processes admitting universal learning. Along the way to establishing these results, another significant contribution of this work is establishing new equivalent characterizations of the families $\mathcal{C}_1$ and $\mathcal{C}_2$, crucial for the design of our optimistically universal algorithms. In particular, we establish a new connection between these two families: proving that $\mathcal{C}_2$ can essentially be characterized by processes that would be in $\mathcal{C}_1$ if we were to replace duplicate values in the sequence $\mathbb{X}$ by some default value $x_0$. As a result, $\mathcal{C}_2$ processes differ from $\mathcal{C}_1$ processes only through duplicates: if a process $\mathbb{X}$ is guaranteed to almost never visit exactly the same context (e.g. i.i.d. processes with density) the properties $\mathcal{C}_1$ and $\mathcal{C}_2$ are equivalent. This fact has further interesting implications, such as a new technique for the design of optimistically universal learning rules for online learning with full supervision; prior to this, only one approach was known to yield such learning rules, based on a modified nearest neighbor algorithm [4]. The new approach suggested in the present work is instead based on an explicit model selection technique, in the spirit of structural risk minimization, analogous to the optimistically universal self-adaptive learning technique developed by [20].

1.7. *Overview of algorithmic techniques.*   We present an overview of the optimistically universal learning rule for finite action sets, Algorithm 5, which encompasses the main algorithmic innovation in this work. We use the property that $\mathcal{C}_2$ processes without duplicates satisfy the $\mathcal{C}_1$ property (Proposition 3.2) to separate times into two classes: points not appearing often recently and points which have many duplicates recently.

1. For the points in the first category, which behave as $\mathcal{C}_1$ processes, we use an approach similar to structural risk minimization: we aim to achieve sublinear regret compared to a constructed countable set of policies that is empirically dense. To do so, we use a restarting technique introduced in [20]: we use classical bandit algorithms as a subroutine to achieve sublinear regret with respect to a fixed finite number of policies, and occasionally restart the bandit learner to gradually increase the number of competing policies considered.
2. For the points in the second category, we use a completely different strategy. Intuitively, these correspond to instances with many duplicates in the recent past, hence it is advantageous to assign each frequent instance an independent bandit learner. In particular, this specific bandit learner is tailored to that point's rewards only and completely disregards historical data from other points.

Interestingly, we can interpret the general strategy as balancing a tradeoff between *generalization* and *personalization*. The first strategy aims to find a policy that performs well at an aggregate level for points with few duplicates. On the other hand, the algorithm performs pure personalization for specific points that have many recent repetitions. This schematic presentation hides many details. In particular, to obtain vanishing excess error compared to the optimal policy, the algorithm needs to balance the generalization/personalization tradeoff carefully, to obtain the required generalization property. In effect, we allow for a cap $M$ of duplicates for each instance in the recent past to be treated with the generalization strategy, and adaptively increase this cap. To adaptively increase this cap, the algorithm occasionally uses "exploration" times to estimate the performance of each strategy, and decides to increase the cap based on these estimates. Last, in order to have decisions robust to non-stationarity in the sequence of contexts, the algorithm selects actions based on recent data: the learning procedure is broken down by periods that contain a given proportion of the past data, then this proportion adaptively decays to $0$.

1.8. *Outline of the paper.*   The remainder of the paper is organized as follows. After giving the definitions and main results in Section 2, we provide in Section 3 new characterizations of stochastic process classes as well as base algorithms, used to construct our learning rules. With these tools, we study optimistic learning with bounded rewards for finite (Section 4), countably infinite (Section 5), and uncountable (Section 6) action sets. We then show that universal learning can be achieved on larger classes of processes under continuity assumptions on the rewards in Section 7. Last, in Section 8 we consider the more restrictive case of unbounded rewards.

## 2. Preliminaries and main results.

2.1. *Formal setup and problem formulation.*    The goal of this paper is to study the general framework of contextual bandits in an online setting. Given a separable metrizable Borel context space $(\mathcal{X}, \mathcal{B})$ and a separable metrizable Borel action space $\mathcal{A}$, the learner interacts with the contextual bandit at each iteration $t \geq 1$ of the learning process in the following fashion. First, the learner observes a context $X_t \in \mathcal{X}$, then selects an action $\hat{a}_t \in \mathcal{A}$ based on the past history only. As a result of the action, the learner receives a reward $r_t$. We will suppose for the most part that the rewards are bounded $r_t \in [0, \bar{r}] = \mathcal{R}$ for some known $\bar{r} \geq 0$. Hence, except for Section 8 in which we consider unbounded rewards, we will take without loss of generality $\bar{r} = 1$. Crucially, the learning rule can only use the past history, which is defined formally as follows.

DEFINITION 2.1 (Learning rule).    A *learning rule* is a sequence $f. = (f_t)_{t \geq 1}$ of possibly randomized measurable functions $f_t : \mathcal{X}^{t-1} \times \mathcal{R}^{t-1} \times \mathcal{X} \to \mathcal{A}$. The action selected at time $t$ by the learning rule is $\hat{a}_t = f_t((X_s)_{s \leq t-1}, (r_s)_{s \leq t-1}, X_t)$.

We suppose that the contexts are generated from a stochastic process $\mathbb{X} = (X_t)_{t \in \mathbb{N}}$ on $\mathcal{X}$. Further, we assume that rewards are sampled from a distribution conditionally on the context and actions. Formally, we assume that there exists a time-invariant conditional distribution $P_{r|a,x}$ such that the rewards $(r_t)_{t \geq 1}$ are conditionally independent given their respective selected action $a_t$ and observed context $x_t$, and follow this conditional distribution. Hence, $(r_t \mid a_t, x_t)_{t \geq 1} \overset{iid.}{\sim} P_{r|a,x}$. To emphasize the conditional dependence of $r_t$ on the actions and context, we denote $r_t(a, x)$ (resp. $r_t(a)$) the reward at time $t$, had the selected action been $a \in \mathcal{A}$ and the observed context $x \in \mathcal{X}$ (resp. when the context at time $t$ is clear). Further, by abuse of notation, we will refer to a reward mechanism $r$ as a random variable $r \sim P_{r|a,x}$. For instance, we use the notation $\bar{r}(a, x) = \mathbb{E}[r \mid a, x]$ to denote the immediate expected reward for any $a \in \mathcal{A}$ and $x \in \mathcal{X}$. When we investigate unbounded rewards in Section 8, we will assume that the random variable $r(a, x)$ is integrable for any $(a, x) \in \mathcal{A} \times \mathcal{X}$. We investigate three settings for the reward mechanism $r$: unrestricted, continuous, and uniformly-continuous. For the two last settings, we suppose that $\mathcal{A}$ is a separable metric space with metric $d$. We formally define the two continuity assumptions below.

DEFINITION 2.2.    The reward mechanism $r$ is continuous if for any $x \in \mathcal{X}$, the immediate expected reward function $\bar{r}(\cdot, x) : \mathcal{A} \to [0, 1]$ is continuous.
The reward mechanism $r$ is uniformly-continuous if for any $\epsilon > 0$ there exists $\Delta(\epsilon) > 0$ with

$$\forall x \in \mathcal{X}, \forall a, a' \in \mathcal{A}, \quad d(a, a') \leq \Delta(\epsilon) \Rightarrow |\bar{r}(a, x) - \bar{r}(a', x)| \leq \epsilon.$$

Our goal is to design algorithms that intuitively converge to the optimal policy $\pi^* : \mathcal{X} \to \mathcal{A}$ that selects for any context $x \in \mathcal{X}$ an optimal arm in $\arg\max_{a \in \mathcal{A}} \bar{r}(a, x)$. Such an *optimal* policy $\pi^*$ is well-defined for finite $\mathcal{A}$; however, for infinite $\mathcal{A}$, this may no longer exist (e.g., if $\sup_{a \in \mathcal{A}} \bar{r}(a, x)$ is not attained). Thus, to be fully general, we instead ask that the regret of the algorithm be sublinear compared to *any* fixed measurable policy $\pi^* : \mathcal{X} \to \mathcal{A}$. We are then interested in learning rules that are consistent irrespective of the unknown reward mechanism $r$, i.e., which intuitively converge to the (near-)optimal policy for all reward mechanisms. We follow the definitions from the universal learning literature for general processes as introduced in [20].

DEFINITION 2.3 (Consistence and universal consistency). Let $\mathbb{X}$ be a stochastic process on $\mathcal{X}$, $r$ be a reward mechanism, and $f.$ be a learning rule. Denote by $(\hat{a}_t)_{t \geq 1}$ its selected actions. We say that $f.$ is *consistent* under $\mathbb{X}$ with rewards $r$ if for any measurable policy $\pi^* : \mathcal{X} \to \mathcal{A}$,

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} r_t(\pi^*(X_t)) - r_t(\hat{a}_t) \leq 0, \quad (a.s.).$$

We say that a learning rule is *universally consistent* if it is consistent under $\mathbb{X}$ for any reward mechanism $r$.

Unfortunately, universal consistency is not always achievable. For example, on $\mathcal{X} = \mathbb{N}$, under the process $\mathbb{X} = (t)_{t \geq 1}$, there does not exist any universally consistent learning rule even in the simplest framework of noiseless—realizable—online learning with full-feedback—when one observes not only the reward $r_t(\hat{a}_t)$ but the complete vector $(r_t(a))_{a \in \mathcal{A}}$ at step $t$ [20, 4]. Two natural questions then arise. First, when is universal consistency possible? And second, which algorithms are universally consistent for a large family of such stochastic processes? To this end, we introduce the notion of optimistically universal learning rules, that "learn whenever learning is possible".

DEFINITION 2.4 (Optimistically universal learning rule). Denote by $\mathcal{C}$ the set of processes $\mathbb{X}$ on $\mathcal{X}$ such that there exists a learning rule universally consistent under $\mathbb{X}$.

We say that a learning rule $f.$ is *optimistically universal* if it is universally consistent under any process $\mathbb{X} \in \mathcal{C}$.

Similarly, we define $\mathcal{C}^c$ (resp. $\mathcal{C}^{uc}$) the set of processes admitting universal learning under continuous (resp. uniformly-continuous) rewards, and define accordingly the notion of optimistically universal learning rule for continuous (resp. uniformly-continuous) rewards. In this paper, we answer the informal questions described above by 1. characterizing the set of learnable processes and 2. showing that there indeed exists and providing optimistically universal learning rules.

2.2. *Useful classes of stochastic processes.* In this subsection we present the key conditions arising in the characterizations of processes on $\mathcal{X}$ admitting universal learning. Let us first start with some notation. For any stochastic process $\mathbb{X} = (X_t)_{t \geq 1}$, we denote $\mathbb{X}_{\leq t} = (X_s)_{s \leq t}$ for any $t \geq 1$. We also introduce the empirical limsup frequency $\hat{\mu}_{\mathbb{X}}$ as follows,

$$\hat{\mu}_{\mathbb{X}}(A) = \limsup_{T \to \infty} \sum_{t=1}^{T} \mathbb{1}_A(X_t), \quad A \in \mathcal{B}.$$

The first condition asks that the set function $\mathbb{E}[\hat{\mu}_{\mathbb{X}}(\cdot)]$ forms a *continuous sub-measure*.

DEFINITION 2.5 (Condition 1 [20]). For every monotone sequence $\{A_k\}_{k=1}^{\infty}$ of measurable subsets of $\mathcal{X}$ with $A_k \downarrow \emptyset$,

$$\lim_{k \to \infty} \mathbb{E}[\hat{\mu}_{\mathbb{X}}(A_k)] = 0.$$

We define $\mathcal{C}_1$ as the set of processes $\mathbb{X}$ satisfying this condition.

For our purposes, we will need to extend this definition to extended stochastic processes which may take values on a subset of possibly random times $\mathcal{T} \subset \mathbb{N}$ instead of the complete set of times $\mathbb{N}$. Overloading the notation $\mathcal{C}_1$, we refer to the same condition $\mathcal{C}_1$ for extended stochastic processes which satisfy the equivalent condition.

DEFINITION 2.6 (Extended condition 1).   Given a possibly random set of times $\mathcal{T} \subset \mathbb{N}$, $\tilde{\mathbb{X}} = (X_t)_{t \in \mathcal{T}}$ satisfies the following condition: for every monotone sequence of measurable sets of $\mathcal{X}$ with $A_k \downarrow \emptyset$,

$$\lim_{k \to \infty} \mathbb{E} \left[ \limsup_{T \to \infty} \frac{1}{T} \sum_{t \leq T, t \in \mathcal{T}} \mathbb{1}_{A_k}(X_t) \right] = 0.$$

As an important remark, the set of $\mathcal{C}_1$ extended stochastic processes is larger than the processes $\tilde{\mathbb{X}} = (X_t)_{t \in \mathcal{T}}$ satisfying $(X_{t_i})_{i \geq 1} \in \mathcal{C}_1$ where $\mathcal{T} = \{t_1 \leq t_2 \leq \ldots\}$ is an enumeration of $\mathcal{T}$. For instance, on $\mathcal{X} = \mathbb{N}$, the process $(X_t = t)_{t \geq 1}$ does not belong to $\mathcal{C}_1$—the decreasing sequence $A_k = \{n \geq k\}$ disproves the condition. However, for any increasing sequence of times $t_k = \omega(k)$, the extended process $\tilde{\mathbb{X}} = (X_t)_{t \in \{t_k, k \geq 1\}}$ with $X_{t_k} = k$ for all $k \geq 1$, belongs to $\mathcal{C}_1$ because $|\{t_k \leq T, k \geq 1\}| = o(T)$.

We then introduce a weaker condition on stochastic processes which asks that the process visits a sublinear number of sets from any measurable partition of $\mathcal{X}$.

DEFINITION 2.7 (Condition 2 [20]).   For every sequence $\{A_k\}_{k=1}^{\infty}$ of disjoint measurable subsets of $\mathcal{X}$,

$$|\{k : \mathbb{X}_{\leq T} \cap A_k \neq \emptyset\}| = o(T) \text{ (a.s.)}.$$

Denote by $\mathcal{C}_2$ the set of all processes $\mathbb{X}$ satisfying this condition.

It is known [20] that $\mathcal{C}_1 \subset \mathcal{C}_2$ and that i.i.d. processes, stationary ergodic processes, stationary processes and processes satisfying the law of large numbers—for any $A \in \mathcal{B}$, the limit $\lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}_A(X_t)$ exists almost surely—belong to $\mathcal{C}_1$. Therefore, both $\mathcal{C}_1$ and $\mathcal{C}_2$ are very general classes of processes.

Last, we introduce a significantly stronger assumption asking that the process only visits a finite number of distinct points.

DEFINITION 2.8 (Condition 3 [20, 6]).

$$|\{x : \mathbb{X} \cap \{x\} \neq \emptyset\}| < \infty \text{ (a.s.)}.$$

Denote by $\mathcal{C}_3$ the set of all processes $\mathbb{X}$ satisfying this condition.

2.3. *Main results.*   We are now ready to present our main results. We show that the set of processes admitting universal learning $\mathcal{C}$ corresponds to one of the classes of processes $\mathcal{C}_3 \subset \mathcal{C}_1 \subset \mathcal{C}_2$ and depends *only* on the action set $\mathcal{A}$. A summary of the charcaterizations is provided in Table 1. In addition, we always provide optimistically universal learning rules for each case, which we construct in the next sections. In the main setting of bounded rewards, the relevant alternatives are whether $\mathcal{A}$ is finite, countably infinite, or uncountable.

THEOREM 2.9 (Unrestricted bounded rewards).   *Let $\mathcal{X}$ be a separable metrizable Borel context space and $\mathcal{A}$ an action space.*

- *If $\mathcal{A}$ is finite and $|\mathcal{A}| \geq 2$, then $\mathcal{C} = \mathcal{C}_2$.*
- *If $\mathcal{A}$ is infinite and countable, then $\mathcal{C} = \mathcal{C}_1$.*
- *If $\mathcal{A}$ is an uncountable separable metrizable Borel space, then $\mathcal{C} = \emptyset$.*

*In all cases there is an optimistically universal learning rule.*

TABLE 1
*Characterization of learnable instance processes for universal learning in contextual bandits depending on properties of the action space $\mathcal{A}$.*

| | **Unrestricted rewards** | **Continuous rewards** | **Uniformly-continuous rewards** |
|---|---|---|---|
| **Bounded rewards** | Finite: $\mathcal{C}_2$<br>Countably infinite: $\mathcal{C}_1$<br>Uncountable: $\emptyset$ | Finite: $\mathcal{C}_2$<br>Infinite: $\mathcal{C}_1$ | Totally-bounded: $\mathcal{C}_2$<br>Non-totally-bounded: $\mathcal{C}_1$ |
| **Unbounded rewards** | Countable: $\mathcal{C}_3$<br>Uncountable: $\emptyset$ | $\mathcal{C}_3$ | $\mathcal{C}_3$ |

We recall that $\mathbb{X} \in \mathcal{C}_2$ is necessary to achieve universal learning under $\mathbb{X}$ even in the simplest online learning setting with full-feedback and noiseless values [20, 4]. Therefore, Therorem 2.9 shows that universal consistence for contextual bandits is achievable for finite action sets at no extra generalizability cost. Unfortunately, in uncountable action spaces, universal consistence is not achievable. A natural question then becomes whether with additional mild assumptions on the rewards one can recover the large classes of processes $\mathcal{C}_1$ or $\mathcal{C}_2$ for universal learning. In particular, we assume that $(\mathcal{A}, d)$ is a separable metric space and first consider the case of *continuous* rewards. Under this first assumption, we show that we can achieve universal consistency on all $\mathcal{C}_1$ processes with an optimistically universal learning rule.

THEOREM 2.10 (Continuous bounded rewards).    *Let $\mathcal{X}$ be a separable metrizable Borel context space and $(\mathcal{A}, d)$ a separable metric action space.*

- *If $\mathcal{A}$ is finite and $|\mathcal{A}| \geq 2$, then $\mathcal{C}^c = \mathcal{C}_2$.*
- *If $\mathcal{A}$ is infinite, then $\mathcal{C}^c = \mathcal{C}_1$.*

*In all cases there is an optimistically universal learning rule for continuous rewards.*

As a result, under the the continuity assumption, one recovers the set of processes $\mathcal{C}_1$ for infinite action spaces. However, it is not sufficient to recover the largest set $\mathcal{C}_2$ which is necessary even in the noiseless full-feedback setting. To this ends, we consider the stronger assumption that rewards are uniformly-continuous and show that one can to recover the set of learnable processes $\mathcal{C}_2$ for totally-bounded action spaces.

THEOREM 2.11 (Uniformly-continuous bounded rewards).    *Let $\mathcal{X}$ be a separable metrizable Borel context space and $(\mathcal{A}, d)$ a separable metric action space.*

- *If $\mathcal{A}$ is totally-bounded and $|\mathcal{A}| \geq 2$, then $\mathcal{C}^{uc} = \mathcal{C}_2$.*
- *If $\mathcal{A}$ is non-totally-bounded, then $\mathcal{C}^{uc} = \mathcal{C}_1$.*

*In all cases there is an optimistically universal learning rule for uniformly-continuous rewards.*

Last, we investigate the more restrictive case of unbounded rewards in $\mathcal{R} = [0, \infty)$. [5] showed that even in the simplest noiseless and full-feedback online learning framework, for unbounded rewards, $\mathcal{C}_3$ is necessary for universal learning. We show that although it forms a restrictive class of processes, universal learning under $\mathcal{C}_3$ processes is still possible for contextual bandits. However, continuity or uniform continuity assumptions are not sufficient to enlarge this set of learnable processes.

THEOREM 2.12 (Unbounded rewards).    *Let $\mathcal{X}$ be a separable metrizable Borel context space and $(\mathcal{A}, d)$ a separable metric action space.*

- If $\mathcal{A}$ is countable, and $|\mathcal{A}| \geq 2$, then $\mathcal{C} = \mathcal{C}_3$. If $\mathcal{A}$ is uncountable, then $\mathcal{C} = \emptyset$.
- $\mathcal{C}^c = \mathcal{C}^{uc} = \mathcal{C}_3$.

*In all cases there is an optimistically universal learning rule for the corresponding rewards (unrestricted, continuous or uniformly-continuous).*

### 3. Base ingredients for the proofs and algorithms.

3.1. *Equivalent characterizations of stochastic process classes.* We give new characterizations of the classes $\mathcal{C}_1$ and $\mathcal{C}_2$ of independent interest.

We first show that for processes $\mathbb{X} \notin \mathcal{C}_1$, we can construct a measurable partition visited linearly by the process up to a known maximum number of duplicates in the instances for each set of the partition. This also characterizes $\mathcal{C}_1$.

LEMMA 3.1. *For any $\mathbb{X} \notin \mathcal{C}_1$, there exists a disjoint sequence $\{B_i\}_{i=1}^{\infty}$ of measurable subsets of $\mathcal{X}$ with $\bigcup_{i \in \mathbb{N}} B_i = \mathcal{X}$, and a sequence $N_i$ in $\mathbb{N}$ such that, letting $i_t$ be the unique $i \in \mathbb{N}$ with $X_t \in B_i$, with probability strictly greater than zero, it holds that*

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}] > 0.$$

*In fact, $\mathbb{X} \notin \mathcal{C}_1$ if and only if this holds.*

PROOF. Suppose $\mathbb{X} \notin \mathcal{C}_1$. By Lemma 14 of [20], there exists a disjoint sequence $\{B_i\}_{i=1}^{\infty}$ of measurable subsets of $\mathcal{X}$ such that, on an event $\mathcal{E}_0$ of probability strictly great than $0$, it holds that

$$\lim_{j \to \infty} \hat{\mu}_{\mathbb{X}} \left( \bigcup_{i \geq j} B_i \right) > 0.$$

Without loss of generality, we may suppose $B_1 = \mathcal{X} \setminus \bigcup_{i>1} B_i$ so that $\bigcup_{i \in \mathbb{N}} B_i = \mathcal{X}$. Define a random variable $\alpha$ as

$$\alpha = \lim_{j \to \infty} \hat{\mu}_{\mathbb{X}} \left( \bigcup_{i \geq j} B_i \right).$$

Inductively define sequences $T_k$, $J_k$ in $\mathbb{N}$ as follows. Let $T_0 = 0$ and $J_0 = 1$. For each $k \in \mathbb{N}$, suppose $T_{k-1}$ and $J_{k-1}$ are defined, elements of $\mathbb{N}$, and define $T_k$ and $J_k$ as follows. Note that, by definition of $\hat{\mu}_{\mathbb{X}}$, there exists an $\mathbb{X}$-dependent random variable $\tau_k \in \mathbb{N}$ with $\tau_k > T_{k-1}$ such that

$$\frac{1}{\tau_k} \left| \mathbb{X}_{\leq \tau_k} \cap \bigcup_{i \geq J_{k-1}} B_i \right| \geq (1/2) \hat{\mu}_{\mathbb{X}} \left( \bigcup_{i \geq J_{k-1}} B_i \right).$$

Moreover, by monotonicity of $\hat{\mu}_{\mathbb{X}}(\cdot)$, the right hand side is no smaller than $\alpha/2$. Let $T_k \in \mathbb{N}$ be any finite non-random value such that

$$\mathbb{P}(\tau_k > T_k) < \mathbb{P}(\mathcal{E}_0) 2^{-k-2}.$$

Next note that, since the sets $B_i$ are disjoint, there exists a finite $\mathbb{X}$-dependent random variable $j_k \in \mathbb{N}$ with $j_k > J_{k-1}$ such that

$$\mathbb{X}_{\leq T_k} \cap \bigcup_{i \geq j_k} B_i = \emptyset.$$

Let $J_k \in \mathbb{N}$ be any finite non-random value such that

$$\mathbb{P}(j_k > J_k) < \mathbb{P}(\mathcal{E}_0)2^{-k-2}.$$

In particular, on the event that $j_k \leq J_k$, it holds that

$$\mathbb{X}_{\leq T_k} \cap \bigcup_{i \geq J_k} B_i = \emptyset,$$

which implies that

$$\mathbb{X}_{\leq T_k} \cap \bigcup_{i \geq J_{k-1}} B_i = \mathbb{X}_{\leq T_k} \cap \bigcup_{J_{k-1} \leq i < J_k} B_i.$$

Thus, if both events $\tau_k \leq T_k$ and $j_k \leq J_k$ hold, it must be that

$$\frac{1}{\tau_k} \left| \mathbb{X}_{\leq \tau_k} \cap \bigcup_{J_{k-1} \leq i < J_k} B_i \right| \geq \alpha/2,$$

or equivalently,

(1)
$$\frac{1}{\tau_k} \sum_{t=1}^{\tau_k} \mathbb{1}[i_t \in \{J_{k-1} \leq i < J_k\}] \geq \alpha/2.$$

This completes the inductive definition of the sequences $T_k$ and $J_k$.

To specify the $N_i$ values, for each $k \in \mathbb{N}$ and $i \in \{J_{k-1}, \ldots, J_k - 1\}$, define $N_i = T_k$. Note that the event $\mathcal{E}_1 = \mathcal{E}_0 \cap \bigcap_{k \in \mathbb{N}} \{\tau_k \leq T_k\} \cap \{j_k \leq J_k\}$ has probability at least

$$\mathbb{P}(\mathcal{E}_0) - \sum_{k \in \mathbb{N}} \mathbb{P}(\mathcal{E}_0)2^{-k-1} = \mathbb{P}(\mathcal{E}_0)/2 > 0$$

by the union bound. On the event $\mathcal{E}_1$, (1) holds for every $k \in \mathbb{N}$. Since $\tau_k \leq T_k$ on $\mathcal{E}_1$, we also trivially have that every $i \in \{J_{k-1}, \ldots, J_k - 1\}$ and $t \in [\tau_k]$ satisfy $|\mathbb{X}_{<t} \cap B_i| < t \leq T_k = N_i$. Together, these facts imply that on $\mathcal{E}_1$, every $k \in \mathbb{N}$ satisfies

$$\frac{1}{\tau_k} \sum_{t=1}^{\tau_k} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}] \geq \alpha/2.$$

Since we also have $\alpha > 0$ on the event $\mathcal{E}_1$, and since $T_k$ is strictly increasing, and $\tau_k > T_{k-1}$ implies $\tau_k \to \infty$ as $k \to \infty$, altogether we have that on the event $\mathcal{E}_1$,

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}]$$

$$\geq \limsup_{k \to \infty} \frac{1}{\tau_k} \sum_{t=1}^{\tau_k} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}] \geq \alpha/2 > 0.$$

We establish the final claim that such a result is not possible for $\mathbb{X} \in \mathcal{C}_1$, as follows. Fix any $\mathbb{X} \in \mathcal{C}_1$. For any disjoint sequence $B_i$ of measurable subsets of $\mathcal{X}$, and any sequence $N_i \in \mathbb{N}$, define $C_n = \bigcup \{B_i : N_i > n\}$, and note that $C_n \downarrow \emptyset$. For every $n \in \mathbb{N}$, we have

(2)
$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}]$$

$$\leq \limsup_{T\to\infty} \frac{1}{T}\sum_{t=1}^{T}\left(\mathbb{1}[|\mathbb{X}_{<t}\cap B_{i_t}|<n]+\mathbb{1}[N_{i_t}>n]\right)$$

$$\leq \left(\limsup_{T\to\infty}\frac{1}{T}\sum_{t=1}^{T}\mathbb{1}[|\mathbb{X}_{<t}\cap B_{i_t}|<n]\right)+\hat{\mu}_{\mathbb{X}}(C_n).$$

For any $m\in\mathbb{N}$, any $t\geq m$ has $\mathbb{1}[|\mathbb{X}_{<t}\cap B_{i_t}|<n]\leq\mathbb{1}[|\mathbb{X}_{<m}\cap B_{i_t}|<n]$, so that

$$\limsup_{T\to\infty}\frac{1}{T}\sum_{t=1}^{T}\mathbb{1}[|\mathbb{X}_{<t}\cap B_{i_t}|<n]$$

$$\leq\limsup_{T\to\infty}\frac{m}{T}+\frac{1}{T}\sum_{t=1}^{T}\mathbb{1}[|\mathbb{X}_{<m}\cap B_{i_t}|<n]=\hat{\mu}_{\mathbb{X}}\left(\bigcup\{B_i:|\mathbb{X}_{<m}\cap B_i|<n\}\right).$$

Since the first expression above has no dependence on $m$, the conclusion remains valid in the limit of $m\to\infty$, so that

$$\limsup_{T\to\infty}\frac{1}{T}\sum_{t=1}^{T}\mathbb{1}[|\mathbb{X}_{<t}\cap B_{i_t}|<n]\leq\lim_{m\to\infty}\hat{\mu}_{\mathbb{X}}\left(\bigcup\{B_i:|\mathbb{X}_{<m}\cap B_i|<n\}\right),$$

which equals zero almost surely (by Lemmas 13 and 14 of [20]). Altogether, for any $n\in\mathbb{N}$, with probability one, (2) is at most $\hat{\mu}_{\mathbb{X}}(C_n)$. Again, since (2) has no dependence on $n$, this inequality remains valid in the limit as $n\to\infty$, so that with probability one, (2) is at most

$$\lim_{n\to\infty}\hat{\mu}_{\mathbb{X}}(C_n),$$

which equals zero almost surely (by Lemma 13 of [20]). The conclusion that (2) equals zero almost surely follows by the union bound. □

Next, we give a new characterization of $\mathcal{C}_2$ processes, which also provides motivation for our generalization of $\mathcal{C}_1$ to extended processes in Definition 2.6. This extension will be essential in our algorithms.

PROPOSITION 3.2.    *Let $\mathbb{X}$ be a stochastic process on $\mathcal{X}$, and define for any $M\geq 1$,*

$$\mathcal{T}^{\leq M}=\left\{t\geq 1:\sum_{t'\leq t}\mathbb{1}[X_{t'}=X_t]\leq M\right\},$$

*the set of times which are duplicates of index at most $M$. In particular, $\mathcal{T}^{\leq 1}$ is the set of times where we delete all duplicates. The following are equivalent.*

1. $\mathbb{X}\in\mathcal{C}_2$.
2. $(X_t)_{t\in\mathcal{T}^{\leq 1}}\in\mathcal{C}_1$.
3. *For all $M\geq 1$, $(X_t)_{t\in\mathcal{T}^{\leq M}}\in\mathcal{C}_1$.*

Essentially, the main difference between extended $\mathcal{C}_1$ and $\mathcal{C}_2$ processes lies in the multiple occurrences of instance points. In particular, if $\mathbb{X}$ never visits the same instance point twice almost surely, as is the case of i.i.d. process with densities, then $\mathbb{X}\in\mathcal{C}_1$ if and only if $\mathbb{X}\in\mathcal{C}_2$.

PROOF. We start by showing $(2) \Rightarrow (1)$. Suppose that a process $\mathbb{X}$ is not in $\mathcal{C}_2$. We aim to show that $\mathbb{X}$ disproves the second property. Because $\mathbb{X} \notin \mathcal{C}_2$, there exists a sequence of disjoint measurable sets $(B_i)_{i \geq 1}$, $\epsilon, \delta > 0$ such that with probability $\delta > 0$

$$\limsup_{T \to \infty} \frac{|\{i : \mathbb{X}_{\leq T} \cap B_i \neq \emptyset\}|}{T} \geq \epsilon.$$

Denote by $\mathcal{A}$ this event, and consider the sets $A_i = \bigcup_{j \geq i} B_j$ for $i \geq 1$. Now fix $i \geq 1$. For any $T \geq 1$, we have

$$\sum_{t \leq T, t \in \mathcal{T}^{\leq 1}} \mathbb{1}_{A_i}(X_t) = |A_i \cap \mathbb{X}_{\leq T}| \geq |\{j \geq i : B_j \cap \mathbb{X}_{\leq T} \neq \emptyset\}| \geq |\{j : \mathbb{X}_{\leq T} \cap B_j \neq \emptyset\}| - (i-1),$$

where in the first inequality we used the fact that the $B_j$ are disjoint for all $j \geq i$, but included within $A_i$. As a result, on the event $\mathcal{A}$ we have $\limsup_{T \to \infty} \frac{1}{T} \sum_{t \leq T, t \in \mathcal{T}^{\leq 1}} \mathbb{1}_{A_i}(X_t) \geq \epsilon$. Hence,

$$\mathbb{E}\left[\limsup_{T \to \infty} \frac{1}{T} \sum_{t \leq T, t \in \mathcal{T}^{\leq 1}} \mathbb{1}_{A_i}(X_t)\right] \geq \epsilon \mathbb{P}[\mathcal{A}] = \epsilon \delta.$$

This holds for all $i \geq 1$ but $A_i \downarrow \emptyset$, which shows that $\mathbb{X}$ does not satisfy property $(2)$.

To prove $(1) \Rightarrow (2)$, now suppose that property $(2)$ is not satisfied by $\mathbb{X}$. We aim to show that $\mathbb{X} \notin \mathcal{C}_2$. Then, there exists a sequence of measurable sets $A_i \downarrow \emptyset$, $\epsilon > 0$ and an increasing sequence of indices $(i_k)_{k \geq 1}$ such that for all $k \geq 1$

$$\mathbb{E}\left[\limsup_{T \to \infty} \frac{|A_{i_k} \cap \mathbb{X}_{\leq T}|}{T}\right] \geq \epsilon.$$

Because the sets $A_i$ are decreasing and the quantity within the expectation is increasing in the set $A$, this shows that for all $i \geq 1$, we have $\mathbb{E}\left[\limsup_{T \to \infty} \frac{|A_i \cap \mathbb{X}_{\leq T}|}{T}\right] \geq \epsilon$. Therefore, for any $i \geq 1$ because $\mathbb{E}\left[\limsup_{T \to \infty} \frac{|A_i \cap \mathbb{X}_{\leq T}|}{T}\right] \leq \mathbb{P}\left[\limsup_{T \to \infty} \frac{|A_i \cap \mathbb{X}_{\leq T}|}{T} \geq \frac{\epsilon}{2}\right] + \frac{\epsilon}{2}$ we obtain for all $i \geq 1$

$$\mathbb{P}\left[\limsup_{T \to \infty} \frac{|A_i \cap \mathbb{X}_{\leq T}|}{T} \geq \frac{\epsilon}{2}\right] \geq \frac{\epsilon}{2}.$$

Again, because the inner quantity is increasing in the set $A$, we obtain

$$\mathbb{P}\left[\limsup_{T \to \infty} \frac{|A_i \cap \mathbb{X}_{\leq T}|}{T} \geq \frac{\epsilon}{2}, \forall i \geq 1\right] = \lim_{I \to \infty} \mathbb{P}\left[\limsup_{T \to \infty} \frac{|A_i \cap \mathbb{X}_{\leq T}|}{T} \geq \frac{\epsilon}{2}, 1 \leq i \leq I\right]$$

$$= \lim_{I \to \infty} \mathbb{P}\left[\limsup_{T \to \infty} \frac{|A_I \cap \mathbb{X}_{\leq T}|}{T} \geq \frac{\epsilon}{2}\right]$$

$$\geq \frac{\epsilon}{2}.$$

We will denote by $\mathcal{H}$ this event in which for all $i \geq 1$, we have $\limsup_{T \to \infty} \frac{|A_i \cap \mathbb{X}_{\leq T}|}{T} \geq \frac{\epsilon}{2}$. Under the event $\mathcal{H}$, for any $i, t^0 \geq 1$, there always exists $t^1 > t^0$ such that $\frac{|A_i \cap \mathbb{X}_{\leq t^1}|}{t^1} \geq \frac{\epsilon}{4}$. We construct a sequence of times $(t_p)_{p \geq 1}$ and indices $(i_p)_{p \geq 1}$, $(u_p)_{p \geq 1}$ by induction as follows. We first pose $i_1 = t_0 = 0$. Now assume that for $p \geq 1$, the time $t_{p-1}$ and index $i_p$ are defined. Let $t_p > t_{p-1}$ such that

$$\mathbb{P}\left[\mathcal{H}^c \cup \bigcup_{t_{p-1} < t \leq t_p} \left\{\frac{|A_{i_p} \cap \mathbb{X}_{\leq t}|}{t} \geq \frac{\epsilon}{4}\right\}\right] \geq 1 - \frac{\epsilon}{2^{p+3}}.$$

This is also possible because $\mathcal{H} \subset \bigcup_{t > t_{p-1}} \left\{ \frac{|A_{i_p} \cap \mathbb{X}_{\leq t}|}{t} \geq \frac{\epsilon}{4} \right\}$. Last, let $i_{p+1} > i_p$ such that $\mathbb{P}[A_{i_{p+1}} \cap \mathbb{X}_{\leq t_p} \neq \emptyset] \leq \frac{\epsilon}{2^{p+3}}$ which is possible since $A_u \downarrow \emptyset$ as $u \to \infty$. We denote $\mathcal{E}_p$ this event. Then,

$$\mathbb{P}\left[ \mathcal{H}^c \cup \bigcup_{t_{p-1} < t \leq t_p} \left\{ \frac{|(A_{i_p} \setminus A_{i_{p+1}}) \cap \mathbb{X}_{\leq t}|}{t} \geq \frac{\epsilon}{4} \right\} \right]$$

$$\geq \mathbb{P}\left[ \mathcal{E}_p \cap \mathcal{H}^c \cup \bigcup_{t_{p-1} < t \leq t_p} \left\{ \frac{|A_{i_p} \cap \mathbb{X}_{\leq t}|}{t} \geq \frac{\epsilon}{4} \right\} \right] \geq 1 - \frac{\epsilon}{2^{p+2}}.$$

We denote $\mathcal{F}_p$ this event. This ends the recursive construction of times $t_p$ and indices $i_p$ for all $p \geq 1$. Note that by construction, $\mathbb{P}[\mathcal{F}_p^c] \leq \frac{\epsilon}{2^{p+2}}$. Hence, by union bound, the event $\mathcal{H} \cap \bigcap_{p \geq 1} \mathcal{F}_p$ has probability $\mathbb{P}[\mathcal{H} \cap \bigcap_{p \geq 1} \mathcal{F}_p] \geq \mathbb{P}[\mathcal{H}] - \frac{\epsilon}{4} \geq \frac{\epsilon}{4}$. For conciseness, denote $B_p = A_{i_p} \setminus A_{i_{p+1}}$. On the event $\mathcal{H} \cap \bigcap_{p \geq 1} \mathcal{F}_p$ we showed that for all $p \geq 1$, there exists $t_{p-1} < t \leq t_p$ such that $|B_p \cap \mathbb{X}_{\leq t}| \geq \frac{\epsilon}{4} t$, and $(B_p)_{p \geq 1}$ is a sequence of disjoint measurable sets.

Now for any $p \geq 1$, we will construct a countable partition of $B_p$ that separates all points falling in $B_p$ within time horizon $t_p$. Let $\delta_p > 0$ such that

$$\mathbb{P}\left[ \min_{u,v \leq t_p : X_u \neq X_v} \rho(X_u, X_v) \leq \delta_p \right] \leq \frac{\epsilon}{2^{p+3}}.$$

We denote by $\mathcal{G}_p$ the complementary of this event. Note that $\mathbb{P}[\bigcup_{p \geq 1} \mathcal{G}_p^c] \leq \frac{\epsilon}{8}$. As a result, the event $\mathcal{I} := \mathcal{H} \cap \bigcap_{p \geq 1} (\mathcal{F}_p \cap \mathcal{G}_p)$ has probability at least $\frac{\epsilon}{8}$. We will show that on this event, $\mathbb{X}$ disproves the $\mathcal{C}_2$ condition. Precisely, let $(x^i)_{i \geq 1}$ a dense sequence of $\mathcal{X}$. We will denote the balls of $\mathcal{X}$ by $B(x, r) = \{x' : \rho(x, x') < r\}$. Define the following partition of $\mathcal{X}$,

$$\mathcal{P}(\delta): \quad P_i(\delta) = B(x^i, \delta) \setminus \bigcup_{j < i} B(x^j, \delta).$$

Finally, for any $p, i \geq 1$, define $P_i^p := P_i(\delta_p) \cap B_p$. We can note that $\bigcup_{i \geq 1} P_i^p = B_p$. Further, the sets $(B_i^p)_{i, p \geq 1}$ are all disjoint, and form a countable sequence. However, on the event $\mathcal{I}$, for every $p \geq 1$, there exists a time $t_{p-1} < t \leq t_p$ such that $|B_p \cap \mathbb{X}_{\leq t}| \geq \frac{\epsilon}{4} t$. But because the event $\mathcal{G}_p$ is satisfied, all the points falling in $B_p$ within horizon $t \leq t_p$ are separated by at least $\delta_p$, hence fall in distinct sets $B_i^p$. As a result,

$$|\{i \geq 1 : P_i^p \cap \mathbb{X}_{\leq t} \neq \emptyset\}| \geq |B_p \cap \mathbb{X}_{\leq t}| \geq \frac{\epsilon}{4} t.$$

This shows that on the event $\mathcal{I}$, for every $p \geq 1$, there exists $t > t_{p-1}$ such that $|\{i, p \geq 1 : P_i^p \cap \mathbb{X}_{\leq t} \neq \emptyset\}| \geq \frac{\epsilon}{4} t$, and as a result

$$\limsup_{T \to \infty} \frac{|\{i, p \geq 1 : P_i^p \cap \mathbb{X}_{\leq T} \neq \emptyset\}|}{T} \geq \frac{\epsilon}{4}.$$

The fact that $\mathbb{P}[\mathcal{I}] \geq \frac{\epsilon}{8}$ ends the proof that $\mathbb{X} \notin \mathcal{C}_2$, and that the first proposition is equivalent to $\mathcal{C}_2$.

We now show the equivalence $(2) \Leftrightarrow (3)$. We clearly have $(3) \Rightarrow (2)$. Now suppose that $\mathbb{X}$ satisfies $(2)$. Let $M > 1$ and $A$ be a measurable set. Then, for any $T \geq 1$, we have

$$\frac{1}{T} \sum_{t \leq T, t \in \mathcal{T}^{\leq M}} \mathbb{1}_A(X_t) \leq M \frac{|A \cap \mathbb{X}_{\leq t}|}{T} = \frac{M}{T} \sum_{t \leq T, t \in \mathcal{T}^{\leq 1}} \mathbb{1}_A(X_t).$$

Because $(X_t)_{t \in \mathcal{T}^{\leq 1}} \in \mathcal{C}_1$, we obtain as a result $(X_t)_{t \in \mathcal{T}^{\leq M}} \in \mathcal{C}_1$ using the definition. This ends the proof of the proposition. $\qquad \square$

As a consequence of Proposition 3.2, we obtain new major insights on the noiseless full-feedback setting. In this setting, an online learning sequentially observes an instance $X_t \in \mathcal{X}$, predicts a value $\hat{Y}_t \in \mathcal{Y}$ then observes the true value $Y_t = f^*(X_t)$ for some unknown measurable function $f : \mathcal{X} \to \mathcal{Y}$. Similarly to the notion of universal consistency for contextual bandits, the goal is to find learning rules satisfying $\frac{1}{T}\sum_{t=1}^{T} \ell(Y_t, \hat{Y}_t) \to 0 \quad (a.s.)$, where $\ell$ is a given near-metric on $\mathcal{Y}$. For this setting, [20] gave a algorithm combining the Hedge algorithm and a "dense" countable family of measurable functions, universally consistent under $\mathcal{C}_1$ processes. [4] then gave a simple 1-nearest-neighbor-based algorithm 2C1NN and showed that in general separable Borel metrizable spaces [5], it is universally consistent under $\mathcal{C}_2$ processes, which are also necessary for universal learning [20]. Proposition 3.2 directly implies that combining the original algorithm from [20] on new instances $X_t$, i.e., on times $\mathcal{T}^{\leq 1}$, with memorization for previously observed instances also yields an optimistically universal learning rule. Unfortunately, such direct argument does not extend to a noisy setting [7] where the values $Y_t$ may not come from a fixed measurable function $f^*(X_t)$.

3.2. *Learning with experts algorithms.* We give the main ingredients that will be used as sub-routine in our algorithms. We start by recalling classical result on the regret of EXP3.

THEOREM 3.3 (Expected regret of EXP3 [8]). *If EXP3 is run with parameters $\eta_t = \sqrt{\frac{\ln K}{tK}}$ on a multi-armed bandit with $K$ arms, then the pseudo regret satisfies*

$$\max_{i=1,\ldots,k} \mathbb{E}\left[\sum_{t=1}^{T} r_i(t)\right] - \mathbb{E}\left[\sum_{t=1}^{T} r_{\hat{i}_t}(t)\right] \leq 2\sqrt{TK\ln K}.$$

We will also need an algorithm for adversarial multi-armed bandits that holds with high probability $1 - \delta$, with parameters that do not depend on the confidence $\delta$ nor the horizon $T$.

THEOREM 3.4 (High-probability regret of EXP3.IX [30]). *There exists an algorithm EXP3.IX for adversarial multi-armed bandit with $K \geq 2$ arms such that for any $\delta \in (0, 1)$ and $T \geq 1$,*

$$\max_{i\in[K]}\sum_{t=1}^{T}(r_t(a_i) - r_t(\hat{a}_t)) \leq 4\sqrt{KT\ln K} + \left(2\sqrt{\frac{KT}{\ln K}} + 1\right)\ln\frac{2}{\delta},$$

*with probability at least $1 - \delta$.*

Specifically we will always use a very simplified version of this result. There exists a universal constant $c > 0$ such that

$$\max_{i\in[K]}\sum_{t=1}^{T}(r_t(a_i) - r_t(\hat{a}_t)) \leq c\sqrt{KT\ln K}\ln\frac{1}{\delta},$$

with probability $1 - \delta$ for $\delta \leq \frac{1}{2}$. This has the following corollary which allows one to consider a countable family of experts asymptotically, based on an argument from [22, Corollary 4]. We use the same construction to design an algorithm EXPINF for learning with a countably infinite number of experts—the original proof extended the Hedge algorithm to infinite number of experts in the full-feedback setting. Precisely, we use an increasing sequence of times $(T_i)_{i\geq 1}$ such that the learning rule performs an independent EXP3.IX algorithm during each period $[T_i, T_{i+1})$. During this period, the EXP3.IX learner is run with $i$ arms consisting in the experts $E_k$ for $k \leq i$. To ease the computations, we choose $T_i = \sum_{j<i} j^3 = \frac{i^2(i+1)^2}{4}$, which yields the following bounds.

COROLLARY 3.5.    *There is an online learning rule* EXPINF *using bandit feedback such that for any countably infinite set of experts* $\{E_1, E_2, \ldots\}$ *(possibly randomized), for any* $T \geq 1$ *and* $0 < \delta \leq \frac{1}{2}$, *with probability at least* $1 - \delta$,

$$\max_{1 \leq i \leq T^{1/8}} \sum_{t=1}^{T} (r_t(E_{i,t}) - r_t(\hat{a}_t)) \leq cT^{3/4}\sqrt{\ln T} \ln \frac{T}{\delta}.$$

*where* $c > 0$ *is a universal constant. Further, with probability one on the learning and the experts, there exists* $\hat{T}$ *such that for any* $T \geq 1$,

$$\max_{1 \leq i \leq T^{1/8}} \sum_{t=1}^{T} (r_t(E_{i,t}) - r_t(\hat{a}_t)) \leq \hat{T} + cT^{3/4}\sqrt{\ln T} \ln T.$$

PROOF.   Denote by $(T_i = \sum_{j<i} j^3)_{i \geq 1}$ the restarting times used in the definition of EXPINF, and by $\hat{a}_t$ its selected action at time $t$. Theorem 3.4 implies that for any $i \geq 1$, with probability at least $0 < \delta < \frac{1}{2}$,

$$\max_{1 \leq j \leq i} \sum_{t=T_i}^{T_{i+1}-1} r_t(E_{j,t}) - r_t(\hat{a}_t) \leq c\sqrt{i(T_{i+1} - T_i)\ln i} \ln \frac{1}{\delta} = ci^2\sqrt{\ln i} \ln \frac{1}{\delta}.$$

Now fix $T \geq 1$ and $\delta > 0$. Let $i \geq 0$ such that $T_{i+1} \leq T < T_{i+2}$. Then summing the above equations gives that with probability at least $\delta$,

$$\max_{1 \leq j \leq T^{1/8}} \sum_{t=1}^{T} r_t(E_{j,t}) - r_t(\hat{a}_t) \leq T_{\lceil T^{1/8} \rceil} + (T - T_{i+1}) + \sum_{t=T_{\lceil T^{1/8} \rceil}}^{T_{i+1}-1} r_t(E_{i,t}) - r_t(\hat{a}_t)$$

$$\leq T_{\lceil T^{1/8} \rceil} + (i+1) + c\frac{i(i+1)(2i+1)}{6}\sqrt{\ln i} \ln \frac{i}{\delta}.$$

Now note that $i \sim \sqrt{2}T^{1/4}$ and $T_{\lceil T^{1/8} \rceil} \sim \frac{\sqrt{T}}{4}$ as $T \to \infty$. Therefore, there exists a universal constant $\tilde{c}$ such that for all $T \geq 1$, the right-hand term is upper bounded by $\tilde{c}T^{3/4}\sqrt{\ln T} \ln \frac{T}{\delta}$. This ends the proof of the first claim.

Now for any $T \geq 1$, using the probabilities of error $\delta_T = \frac{1}{T^2}$ which are summable, the Borel-Cantelli lemma implies that on an event of probability one, there exists $\hat{T}$ such that for any $T \geq \hat{T}$,

$$\max_{1 \leq j \leq T^{1/8}} \sum_{t=1}^{T} r_t(E_{j,t}) - r_t(\hat{a}_t) \leq \tilde{c}T^{3/4}\sqrt{\ln T} \ln(T^3) = 3\tilde{c}T^{3/4}\sqrt{\ln T} \ln T,$$

which ends the proof of the second claim by redefining the constant $c > 0$.    □

**4. Finite action space.**   In this section, we assume that the action space $\mathcal{A}$ is finite and we show that in this case, the set of processes $\mathbb{X}$ admitting universal learning is exactly $\mathcal{C}_2$. In other terms, we can recover the same processes which admit universal learning in the full-feedback setting.

We start by showing that the $\mathcal{C}_2$ condition is necessary for universal consistency, which is a direct consequence from its necessity in the full-feedback case [20].

THEOREM 4.1.    *If* $2 \leq |\mathcal{A}| < \infty$, $\mathbb{X} \in \mathcal{C}_2$ *is necessary for universal consistency, i.e.,* $\mathcal{C} \subset \mathcal{C}_2$.

PROOF. In the full-information feedback setting, [20, Theorem 37] showed that $\mathbb{X} \in \mathcal{C}_2$ is necessary for universal learning even for noiseless responses in binary classification. We will present a simple reduction from the full-feedback to the partial-feedback setting. Precisely, let $a_0, a_1 \in \mathcal{A}$ be two distinct actions. To any measurable function $f : \mathcal{X} \to \{0, 1\}$ we associate a deterministic reward function $r_f : \mathcal{X} \times \mathcal{A} \to [0, 1]$ as follows

$$r_f(x, a) = f(x)\mathbb{1}[a = a_1] + (1 - f(x))\mathbb{1}[a = a_0], \quad x \in \mathcal{X}, a \in \mathcal{A}.$$

Note that any action $a \in \mathcal{A} \setminus \{a_0, a_1\}$ always has reward 0. Now suppose that for a process $\mathbb{X}$ there exists an universally consistent learning rule $f.$ for contextual bandits. Then, we can consider the following learning rule for the complete-feedback setting, recursively defined as

$$\tilde{f}_t(\boldsymbol{x}_{\leq t-1}, \boldsymbol{y}_{\leq t-1}, x_t) = \mathbb{1}[f_t(\boldsymbol{x}_{\leq t-1}, (\mathbb{1}[\tilde{f}_i(\boldsymbol{x}_{\leq i-1}, \boldsymbol{y}_{\leq i-1}, x_i) = y_i])_{i \leq t-1}, x_t) = a_1].$$

for any $t \geq 1$, $\boldsymbol{x}_{\leq t} \in \mathcal{X}^{t-1}$ and $\boldsymbol{y}_{\leq t-1} \in \{0, 1\}^{t-1}$. We now shows that $\tilde{f}.$ is universally consistent for the noiseless full-feedback setting. For any measurable function $f : \mathcal{X} \to \{0, 1\}$, the learning rule $f.$ is consistent for the rewards $r_f$. In particular, if we denote by $\hat{a}_t$ the action selected by $f.$ at time $t$, using the measurable policy $\pi_f : x \in \mathcal{X} \mapsto a_0\mathbb{1}[f(x) = 0] + a_1\mathbb{1}[f(x) = 1] \in \mathcal{A}$ which always selects the best action we obtain

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} r_t(\pi_f(X_t)) - r_t(\hat{a}_t) = \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[\hat{a}_t \neq \pi_f(X_t)] \leq 0, \quad (a.s.).$$

Now consider the actions $\hat{a}_t$ selected under $\mathbb{X}$ and rewards $r_f$ and denote by $\tilde{Y}_t$ the prediction of $\tilde{f}.$ at time $t$ under $\mathbb{X}$ and values $Y_t = f(X_t)$ for $t \geq 1$. By construction, for any $t \geq 1$, we have $\mathbb{1}[\hat{a}_t \neq \pi_f(X_t)] \geq \mathbb{1}[\tilde{Y}_t \neq f(X_t)]$. Then, almost surely $\frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[\tilde{Y}_t \neq f(X_t)] \xrightarrow[n \to \infty]{} 0$. This shows that $\tilde{f}.$ is universally consistent for noiseless responses in binary classification, hence $\mathbb{X} \in \mathcal{C}_2$, which completes the proof. $\square$

We now present a learning rule for contextual bandits, which we will next show is universally consistent under any $\mathcal{C}_2$ process.

This learning rule at time $t$ has different behaviour depending on the number of occurrences of $X_t$ that were observed in the past. Precisely, for any time $t$, we compute a corresponding category $p$ such that the number of past occurrences of $X_t$ belongs in the interval $[4^p, 4^{p+1})$. The learning rule will behave completely separately on times from different categories. The formal definition is given by the function below

$$\text{CATEGORY}(t, \mathbb{X}_{\leq t}) = \left\lfloor \log_4 \left( \sum_{t' \leq t} \mathbb{1}[X_{t'} = X_t] \right) \right\rfloor.$$

For convenience we may write $\text{CATEGORY}(t)$ instead of $\text{CATEGORY}(t, \mathbb{X}_{\leq t})$. Further, for a given category $p$, the algorithm will proceed by periods $[T_p^q, T_p^{q+1})$ defined as follows. For any $p \geq 0$ and $q \geq p2^p$, we define the times $T_p^q = 2^k + \frac{i}{2^p} 2^k$, where $q = k2^p + i$ with $0 \leq i < 2^p$. Note that the sequence $(T_p^q)_q$ has an exponential behaviour with rate between $2^{-p-1}$ and $2^{-p}$. We will refer to $[T_p^q, T_p^{q+1})$ as the period $q$ for category $p$. We then define the function $\text{PERIOD}(t)$ which returns the index $q$ such that $T_p^q \leq t \leq T_p^{q+1}$ where $p$ is the category of $p$. Now let $(\pi^l)_{l \geq 1}$ be a sequence of measurable functions from $\mathcal{X}$ to $\mathcal{A}$ that are dense within measurable functions under $\mathcal{C}_1$ processes. Intuitively, the learning rule combines two strategies: a strategy 0 which applies a separate EXP3 algorithm to each distinct instance, and a strategy 1 which performs the best policy within a subset of the policies $(\pi^l)_{l \geq 1}$. In order to know which strategy to apply, the learning rule computes an estimate of the counterfactual

loss of strategy $i$, using classical importance sampling on some allocated exploration times for strategy $i$. On the exploitation times, the learning rule uses these estimates to perform the best strategy.

We first define the procedure ASSIGNPURPOSE which taking as input a time $t$ determines whether this time will be used for exploration of strategy 0 (output 0), strategy 1 (output 1), or exploitation (output 2). Intuitively, ASSIGNPURPOSE selects exploration times randomly with small probability while ensuring that times $t, t'$ from the same category $p$, period $q$, and that are duplicates $X_t = X_{t'}$ are assigned the same output, hence will serve for the same exploration or exploitation purpose. The algorithm is formally defined in Algorithm 1.

---

**Input:** time $t$, $\mathbb{X}_{\leq t}$, CATEGORY($t'$) for $t' \leq t$, ASSIGNPURPOSE($t'$) for $t' < t$.
**Output:** ASSIGNPURPOSE($t$) $\in \{0, 1, 2\}$.
$p = $ CATEGORY($t$); $q = $ PERIOD($t$)
**if** *exists $t' < t$ with* CATEGORY($t'$) $= p$; PERIOD($t'$) $= q$ *and* $X_t = X_{t'}$ **then**     // Not the first
  occurrence of $X_t$ in current period
    |   Return ASSIGNPURPOSE($t'$)
**else**                              // First occurrence of $X_t$ in current period
    |  $p_t = 1/(2t^{1/4})$
    |  $U_t \sim \mathcal{U}([0,1])$
    |  **if** $U_t \leq p_t$ **then** Return 0                   // Exploration for strategy 0
    |  **else if** $p_t < U_t \leq 2p_t$ **then** Return 1      // Exploration for strategy 1
    |  **else** Return 2                             // Exploitation
**end**

**Algorithm 1:** ASSIGNPURPOSE

---

Next, we define the subroutine EXPLORE($i; t$) that will be called on exploration times $t$ for strategy $i$. We first define it to estimate the performance of strategy 0. The subroutine updates an estimator $\hat{R}_p^0(q)$ of the loss that would be incurred by using strategy 0 for all times in category $p$ during period $q$. EXPLORE($0, \cdot$) is defined formally in Algorithm 2.

---

**Input:** time $t$, $\mathbb{X}_{\leq t}$, CATEGORY($t'$) for $t' \leq t$, rewards $\boldsymbol{r}_{<t}$, $\hat{R}_p^0(q)$ for $p \geq 0, q \geq p2^p$.
**Output:** Selects action $\hat{a}_t$ and updates $\hat{R}_p^0(q)$ for $p = $ CATEGORY($t$), $q = $ PERIOD($t$).
$p = $ CATEGORY($t$), $q = $ PERIOD($t$)
$S_t = \{t' < t : $ CATEGORY($t'$) $= p$, PERIOD($t'$) $= q, X_{t'} = X_t\}$
$\hat{a}_t = $ EXP3$_{\mathcal{A}}(\hat{\boldsymbol{a}}_{S_t}, \boldsymbol{r}_{S_t})$
Receive reward $r_t$
Let $t' = \min S_t$                                     // First occurrence of $X_t$
$\hat{R}_p^0(q) \leftarrow \hat{R}_p^0(q) + \frac{r_t}{p_{t'}}$                    // Update estimate $\hat{R}_p^0(q)$

**Algorithm 2:** EXPLORE($0; \cdot$)

---

Then, we define EXPLORE($1, \cdot$). It updates an estimator $\hat{R}_p^l(q)$ of the loss that would have been incurred using the policy $\pi^l$ for all times in category $p$ during period $q$, for all $l \geq 1$. Because there is an infinite number of such policies, they are introduced sequentially in the estimation process. EXPLORE is defined formally in Algorithm 3.

The estimates $\hat{R}_p^l(q)$ updated by EXPLORE are then used to select the strategy to perform on exploitation times. The learning rule that we will define acts separately on times from different categories: for any category $p \geq 0$, before starting phase $q$, the learning rule commits to performing strategy $\mathcal{P}_p(q) \in \{0, 1\}$, for times of that phase $q$ for category $p$. The

---

**Input:** time $t$, $\mathbb{X}_{\leq t}$, CATEGORY$(t')$ for $t' \leq t$, rewards $\boldsymbol{r}_{<t}$, $\hat{R}_p^l(q)$ for $l \geq 1, p \geq 0, q \geq p2^p$.

**Output:** Selects action $\hat{a}_t$ and updates $\hat{R}_p^l(q)$ for $p = $ CATEGORY$(t)$, $q = $ PERIOD$(t)$.

$p = $ CATEGORY$(t)$, $q = $ PERIOD$(t)$, $k = \lfloor \log_2 t \rfloor$

$l_t = \mathcal{U}(\{1, \ldots, k\})$                // Uniform exploration

$\hat{a}_t = \pi^{l_t}(X_t)$

Receive reward $r_t$

Let $t' = \min\{s < t : $ CATEGORY$(s) = p,$ PERIOD$(s) = q, X_s = X_t\}$   // First occurrence of $X_t$

$\hat{R}_p^l(q) \leftarrow \hat{R}_p^l(q) + \frac{k}{p_{t'}} r_t \mathbb{1}[l = l_t], \quad 1 \leq l \leq k$      // Update estimate $\hat{R}_p^{l_t}(q)$

---

**Algorithm 3:** EXPLORE$(1; \cdot)$

choice of strategy $\mathcal{P}_p(q)$ is performed by a subroutine SELECTSTRATEGY which applies an $\eta_p = \mathcal{O}(2^{-p/2})$ average reward penalty for strategy 0 then select the strategy that obtained the highest adjusted estimated reward during the previous period. Last, if during the current period $q$, strategy 0 obtained the highest adjusted reward, we select this strategy for the future periods $q < q' \leq q + p2^p$. This ensures that if by mistake the rule selected $\mathcal{P}_p(q) = 1$, the loss incurred during this period is mitigated for the next strategy selection: the current performance until time $T^{q+1}$ is negligible up to a small average loss starting from time $T_p^{q+2^p+1}$. The construction of SELECTSTRATEGY is detailed in Algorithm 4.

---

**Input:** Category $p$, phase $q$, variable states $\hat{R}_p^l(t)$ for $t < T_p^{q+1}$

**Output:** Selects strategy $\mathcal{P}_p(r)$ for some future phases $r > q$.

$\eta_p = 10 \frac{\sqrt{|\mathcal{A}| \ln |\mathcal{A}|}}{2^{p/4}}, k = \lfloor \log_2 T_p^q \rfloor$

**if** $\mathcal{P}_p(q+1)$ *has not been defined yet* **then**

    **if** $\hat{R}_p^0(q) - \eta_p(T_p^{q+1} - T_p^q) \geq \max_{1 \leq l \leq k} \hat{R}_p^l(q)$ **then**

        $\mathcal{P}_p(q') = 0, \quad q < q' \leq q + p2^p$      // perform strategy 0 until current

        performance is negligible up to a $\mathcal{O}(2^{-p})$ average loss

    **else**

        $\mathcal{P}_p(q+1) = 1$

    **end**

**end**

---

**Algorithm 4:** SELECTSTRATEGY

We are now ready to define the learning rule for stochastic rewards. On exploration times, the learning rule calls the subroutine EXPLORE, and on exploitation times, the learning rule performs the corresponding strategy $\mathcal{P}_p(q)$ for times in category $p$ during phase $q$. The construction of the learning rule is detailed in Algorithm 5.

The main result of this section is that this learning rule is optimistically universal.

THEOREM 4.2. *Let $\mathcal{X}$ a metrizable separable Borel space and $\mathcal{A}$ a finite action set. Then, there exists an optimistically universal learning rule and the set of learnable processes is $\mathcal{C} = \mathcal{C}_2$.*

$\hat{R}_p^l = 0, l \geq 0, p \geq 0; \mathcal{P}_p(p2^{p+5}) = 0, p \geq 0$      // Initialization

**for** $t \geq 1$ **do**
   Observe context $X_t$
   $p = \text{CATEGORY}(t), q = \text{PERIOD}(t)$
   **if** $t < 2^{32p}$ **then**      // Initially perform strategy 0 without period restriction
     $S_t = \{t' < t : \text{CATEGORY}(t') = p, X_{t'} = X_t\}$
     $\hat{a}_t = \text{EXP3}_{\mathcal{A}}(\hat{\boldsymbol{a}}_{S_t}, \boldsymbol{r}_{S_t})$
   **else if** $i := \text{ASSIGNPURPOSE}(t) \leq 1$ **then**
     $\text{EXPLORE}(i; t)$
   **else**      // Perform strategy $\mathcal{P}_p(q)$
     **if** $\mathcal{P}_p(q) = 0$ **then**
       $S_t = \{t' < t : \text{CATEGORY}(t') = p, \text{PERIOD}(t') = q, X_{t'} = X_t\}$
       $\hat{a}_t = \text{EXP3}_{\mathcal{A}}(\hat{\boldsymbol{a}}_{S_t}, \boldsymbol{r}_{S_t})$
     **else**
       $k = \lfloor \log_2 T_p^q \rfloor$
       $S_t = \{t' < t : \text{CATEGORY}(t') = p, \text{PERIOD}(t') = q, \text{ASSIGNPURPOSE}(t') = 2\}$
       $l_t = \text{EXP3.IX}_{\{1,\dots,k\}}(\boldsymbol{l}_{S_t}, \boldsymbol{r}_{S_t})$      // Select policy $\pi^{l_t}$
       $\hat{a}_t = \pi^{l_t}(X_t)$
     **end**
     Receive reward $r_t$
   **end**
   $\mathcal{E} = \{(p', q') : q' \geq p'2^{p'+5}, t = T_{p'}^{q'+1} - 1\}$
   **for** $(p', q') \in \mathcal{E}$ **do**
     $\text{SELECTSTRATEGY}(p', q')$      // At the end of a phase $[T_{p'}^{q'}, T_{p'}^{q'-1})$, select strategy for future phases
   **end**
**end**

**Algorithm 5:** An optimistically universal learning rule for stochastic rewards

PROOF. We will denote by $\hat{a}_t$ the action selected by the learning rule at time $t$. For any $p \geq 0$, we define the set $\mathcal{T}_p$ of times in category $p$ as follows

$$\mathcal{T}_p = \left\{ t \geq 1 : 4^p \leq \sum_{t' \leq t} \mathbb{1}[X_{t'} = X_t] < 4^{p+1} \right\},$$

i.e. the set of times which correspond to duplicates of index in $[4^p, 4^{p+1})$. We also define

$$\mathcal{T}_p^{exp,i} = \{t \geq 2^{32p} : \text{ASSIGNPURPOSE}(t) = i\}, \quad i \in \{0, 1\},$$

$$\tilde{\mathcal{T}}_p = \{t \geq 2^{32p} : \text{ASSIGNPURPOSE}(t) = 2\},$$

the set of exploration times for strategy $i$ in category $p$, and exploitation times in category $p$, respectively. For convenience, we also define $\mathcal{T}_p(q) = \mathcal{T}_p \cap [T_p^q, T_p^{q+1})$ times in category $p$ and phase $q$. Last, we define $A_p(q) = |\mathcal{T}_p(q) \cap (\mathcal{T}_p^{exp,0} \cup \mathcal{T}_p^{exp,1})|$ the number of exploration times in period $q$ for category $p$.

Now fix a process $\mathbb{X} \in \mathcal{C}_2$ and let $r$ be a reward mechanism on $\mathcal{A} \times \mathcal{X}$. We recall the notation $\bar{r}(\cdot, \cdot) = \mathbb{E}[r(\cdot, \cdot)]$ for the average reward. We aim to show that $f.$ is consistent under $\mathbb{X}$ for the rewards given by $r$. We first define the policy $\pi^*$ given by

$$\pi^*(x) = \arg\max_{a \in \mathcal{A}} \bar{r}(a, x),$$

where ties are broken by the lexicographic rule. This function is measurable given that $\mathcal{A}$ is finite. Further, it is an optimal policy in the sense that for any measurable function $\pi : \mathcal{X} \to \mathcal{A}$ and any $x \in \mathcal{X}$, $\bar{r}(\pi(x), x) \leq \bar{r}(\pi^*(x), x)$.

For $p \geq 0$, we first analyze the reward estimates $\hat{R}_p^l(q)$ for $q \geq p2^{p+5}$ ($T_p^{p2^{p+5}} = 2^{32p}$) and $l \geq 0$. First note that the exploration times $\mathcal{T}_p^{exp,0}$ and $\mathcal{T}_p^{exp,1}$ were constructed precisely so that times corresponding to the same instance and within the same period, fall in the same set $\mathcal{T}_p^{exp,0}$, $\mathcal{T}_p^{exp,1}$, or $\tilde{\mathcal{T}}_p$. For simplicity, we will write $\mathcal{X}_p(q) = \{X_t, t \in \mathcal{T}_p(q)\}$ the set of visited instances during period $q$ of category $p$, and for $x \in \mathcal{X}_p(q)$ we denote $t_p(q; x) = \min\{t \in \mathcal{T}_p(q) : X_t = x\}$ the first time of occurrence of $x$ in period $q$. Then, we can write

$$\hat{R}_p^0(q) = \sum_{x \in \mathcal{X}_p(q)} \frac{\mathbb{1}[U_{t_p(q;x)} \leq p_{t_p(q;x)}]}{p_{t_p(q;x)}} \sum_{t \in \mathcal{T}_p(q), X_t = x} \tilde{r}_t$$

where $\tilde{r}_t$ is the reward at time $t$ that would have been obtained by performing strategy 0 during period $q$, i.e., assigning an independent EXP3 learner for each different instance in this period. We compare $\hat{R}_p^0(T)$ to the average reward obtained by the optimal policy $\pi^*$,

$$\bar{R}_p^*(q) := \sum_{t \in \mathcal{T}_p(q)} \bar{r}(\pi^*(X_t), X_t).$$

Observe that conditionally on $\mathbb{X}$, the terms in the sum of $\hat{R}_p^0(q)$ are independent. For any $x \in \mathcal{X}_p(q)$, let $\bar{R}_p^0(q; x) = \mathbb{E}[\sum_{t \in \mathcal{T}_p(q), X_t = x} \tilde{r}_t \mid \mathbb{X}]$, the average reward obtained by strategy 0 on the instance $x$. We will use the notation $N_p(q; x) = |\{t \in \mathcal{T}_p(q), X_t = x\}| \leq 4^{p+1}$ for the number of occurrences of the instance $x$ within $\mathcal{T}_p$. Note that

$$\left| \frac{\mathbb{1}[U_{t_p(q;x)} \leq p_{t_p(q;x)}]}{p_{t_p(q;x)}} \sum_{t \in \mathcal{T}_p(q), X_t = x} \tilde{r}_t \right| \leq \frac{N_p(q;x)}{p_{t_p(q;x)}} \leq 2^{2p+3}(T_p^{q+1})^{1/4},$$

and that $|\mathcal{X}_p(q)| \leq \frac{T_p^{q+1}}{2^{2p}}$ since by definition of $\mathcal{T}_p$ each instance has already occurred $4^p$ times. As a result, we can apply Hoeffding's inequality to obtain

$$\mathbb{P}\left[ \left| \hat{R}_p^0(q) - \sum_{x \in \mathcal{X}_p(q)} \bar{R}_p^0(q; x) \right| \leq (T_p^{q+1})^{\frac{7}{8}} \mid \mathbb{X} \right] \geq 1 - 2 \exp\left( -\frac{(T_p^{q+1})^{1/4}}{2^{2p+5}} \right) := 1 - 2p_1(p, q)$$

Now applying Theorem 3.3 to each pseudo-regret $\bar{R}_p^0(q; x)$ yields

$$\sum_{x \in \mathcal{X}_p(q)} \bar{R}_p^0(q; x)$$

$$\geq \sum_{x \in \mathcal{X}_p(q)} N_p(q; x) \left( \max_{a \in \mathcal{A}} \bar{r}(a, x) - 2\sqrt{\frac{|\mathcal{A}| \ln |\mathcal{A}|}{N_p(q; x)}} \right)$$

$$\geq \bar{R}_p^*(q) - 2\sqrt{\frac{|\mathcal{A}| \ln |\mathcal{A}|}{2^{p/2}}} (T_p^{q+1} - T_p^q) - 2\sqrt{|\mathcal{A}| \ln |\mathcal{A}|} \sum_{x \in \mathcal{X}_p(q), N_p(q;x) \leq 2^{p/2}} N_p(q; x)$$

$$\geq \bar{R}_p^*(q) - 2\frac{\sqrt{|\mathcal{A}| \ln |\mathcal{A}|}}{2^{p/4}} (T_p^{q+1} - T_p^q) - 2\sqrt{|\mathcal{A}| \ln |\mathcal{A}|} \frac{2^{p/2}}{4^p} T_p^{q+1}$$

$$\geq \bar{R}_p^*(q) - 6\frac{\sqrt{|\mathcal{A}| \ln |\mathcal{A}|}}{2^{p/4}} (T_p^{q+1} - T_p^q).$$

where in the third inequality, we used the fact that instances appearing in $\mathcal{T}_p$ before $T_p^{q+1}$ are visited at least $4^p$ times before horizon $T_p^{q+1}$, by construction of $\mathcal{T}_p$; and in the last inequality we used $2^{-p-1}T_p^{q+1} \leq T_p^{q+1} - T_p^q \leq 2^{-p}T_p^q$. Also, note that $\bar{R}_p^*(q) \geq \sum_{x\in\mathcal{X}_p(q)} \bar{R}_p^0(q;x)$. As a result, taking the expectation over $\mathbb{X}$, we obtain that with probability at least $1 - 2p_1(p,q)$,

$$(3) \qquad \left| \hat{R}_p^0(q) - \bar{R}_p^*(q) \right| \leq (T_p^{q+1})^{\frac{7}{8}} + 6\frac{\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}}{2^{p/4}}(T_p^{q+1} - T_p^q).$$

Now consider the quantity $\tilde{R}_p^0(q)$, the reward that would be obtained for exploitation times on period $q$ if strategy 0 was applied. We have

$$\tilde{R}_p^0(q) = \sum_{x\in\mathcal{X}_p(q)} \sum_{t\in\mathcal{T}_p(q),X_t=x,t\in\tilde{\mathcal{T}}_p} \tilde{r}_t$$

$$\geq \sum_{x\in\mathcal{X}_p(q)} \sum_{t\in\mathcal{T}_p(q),X_t=x} \tilde{r}_t - A_p(q).$$

Similarly as above, using Hoeffding's inequality, we have

$$\mathbb{P}\left[ \sum_{x\in\mathcal{X}_p(q)} \sum_{t\in\mathcal{T}_p(q),X_t=x} \tilde{r}_t \geq \sum_{x\in\mathcal{X}_p(q)} \bar{R}_p^0(q;x) - (T_p^{q+1})^{3/4} \right] \geq 1 - e^{-\frac{\sqrt{T_p^{q+1}}}{2^{2p+3}}} := 1 - p_2(p,q).$$

As a result, with probability $1 - p_2(p,q)$, we have

$$(4) \qquad \tilde{R}_p^0(q) \geq \bar{R}_p^*(q) - 6\frac{\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}}{2^{p/4}}(T_p^{q+1} - T_p^q) - (T_p^{q+1})^{3/4} - A_p(q).$$

We now turn to the estimates $\hat{R}_p^l(q)$ for $l \geq 1$. Note that the estimation of $R_p^l(q)$ only starts at time $2^l$. Hence, we can consider $k(q) = \lfloor \log_2 T_p^q \rfloor = \lfloor \frac{q}{2^p} \rfloor$ and observe that during period $q$, the only estimates $\hat{R}_p^l(q)$ that are considered are for $1 \leq l \leq k(q)$. Therefore, similarly as for the estimates $\hat{R}_p^0(q)$, we can write for $q \geq p2^{p+5}$ and $1 \leq l \leq k(q)$,

$$\hat{R}_p^l(q) = \sum_{x\in\mathcal{X}_p(q)} \frac{\mathbb{1}[p_{t_p(q;x)} < U_{t_p(q;x)} \leq 2p_{t_p(q;x)}]}{p_{t_p(q;x)}} \sum_{t\in\mathcal{T}_p(q)X_t=x} k(t)\mathbb{1}[l = l_t]r(\pi^l(x),x),$$

where $k(t)$ is the number of policies $\pi^l$ tested at time $t$, i.e. $k(t) = \lfloor \log_2 t \rfloor$. Conditionally on $\mathbb{X}$ and $\boldsymbol{U}$ we can apply Hoeffding's inequality to obtain

$$\mathbb{P}\left[ \left| \hat{R}_p^l(q) - \sum_{x\in\mathcal{X}_p(q)} \sum_{t\in\mathcal{T}_p(q),X_t=x} \frac{\mathbb{1}[p_{t_p(q;x)} < U_{t_p(q;x)} \leq 2p_{t_p(q;x)}]}{p_{t_p(q;x)}} \bar{r}(\pi^l(x),x) \right| \right.$$

$$\left. \leq (T_p^{q+1})^{7/8} \mid \mathbb{X}, \boldsymbol{U} \right]$$

$$\geq 1 - 2e^{-\frac{2(T_p^{q+1})^{7/4}}{(T_p^{q+1} - T_p^q)4(\log_2 T_p^{q+1})^2\sqrt{T_p^{q+1}}}} \geq 1 - 2e^{-\frac{2^p(T_p^{q+1})^{1/4}}{4(\log_2 T_p^{q+1})^2}} := 1 - 2p_3(p,q).$$

For convenience, let us denote by $\hat{R}_{p,bis}^l(q)$ the sum in the above inequality. We also define $\bar{R}_p^l(q) = \sum_{t\in\mathcal{T}_p(q)} \bar{r}(\pi^l(X_t),X_t)$ the expected reward of policy $l$ on period $q$. Now, similarly as before, we have

$$0 \leq \sum_{t\in\mathcal{T}_p(q),X_t=x} \frac{\mathbb{1}[p_{t_p(q;x)} < U_{t_p(q;x)} \leq 2p_{t_p(q;x)}]}{p_{t_p(q;x)}} \bar{r}(\pi^l(x),x) \leq \frac{N_p(q;x)}{p_{t_p(q;x)}} \leq 2^{2p+3}(T_p^{q+1})^{1/4}.$$

As a result, conditionally on $\mathbb{X}$, Hoeffding's inequality yields

$$\mathbb{P}[|\hat{R}^l_{p,bis}(q) - \bar{R}^l_p(q)| \leq (T_p^{q+1})^{7/8} \mid \mathbb{X}] \geq 1 - 2p_1(p,q).$$

Thus, with probability at least $1 - 2p_1(p,q) - 2p_3(p,q)$ we have

(5) $$|\hat{R}^l_p(q) - \bar{R}^l_p(q)| \leq (T_p^{q+1})^{7/8}.$$

Next, we consider the quantity $\tilde{R}^1_p(q)$, the reward that would have been obtained for exploitation times on period $q$ if strategy 1 was applied. Then, using Theorem 3.4, we have with probability at least $1 - e^{-(T_p^{q+1})^{1/4}} := 1 - p_4(p,q)$,

$$\max_{1 \leq l \leq k(q)} \sum_{t \in \mathcal{T}_p(q) \cap \tilde{\mathcal{T}}_p} r_t(\pi^l(X_t), X_t) - \tilde{R}^1_p(q) \leq c\sqrt{k(q) \ln k(q)(T_p^{q+1} - T_p^q)}(T_p^{q+1})^{1/4}$$

$$\leq c(T_p^{q+1})^{3/4} \ln T_p^{q+1}.$$

As a result, we have

$$\tilde{R}^1_p(q) \geq \max_{1 \leq l \leq k(q)} \sum_{t \in \mathcal{T}_p(q)} r_t(\pi^l(X_t), X_t) - c(T_p^{q+1})^{3/4} \ln T_p^{q+1} - A_p(q).$$

Now, by Hoeffding's inequality, for every $1 \leq l \leq k(q)$, with probability at least $1 - e^{-2^p \sqrt{T_p^{q+1}}} := 1 - p_5(p,q)$,

$$\sum_{t \in \mathcal{T}_p(q)} r_t(\pi^l(X_t), X_t) \geq \bar{R}^l_p(q) - (T_p^{q+1})^{3/4}.$$

Hence, with probability $1 - p_4(p,q) - k(q)p_5(p,q)$ we have

(6) $$\tilde{R}^1_p(q) \geq \max_{1 \leq l \leq k(q)} \bar{R}^l_p(q) - (T_p^{q+1})^{3/4} - c(T_p^{q+1})^{3/4} \ln T_p^{q+1} - A_p(q).$$

We will also need the quantity $\tilde{R}^1_p(q;T)$ for $T_p^q \leq T < T_p^{q+1}$ which is the reward that would have been obtained for exploitation times from $T_p^q$ to $T$. The exact same arguments as above show that with probability at least $1 - p_4(p,q) - k(q)p_5(p,q)$ we have

(7) $$\tilde{R}^1_p(q;T) \geq \max_{1 \leq l \leq k(q)} \sum_{t \in \mathcal{T}_p(q), t \leq T} \bar{r}(\pi^l(X_t), X_t) - (T_p^{q+1})^{3/4} - c(T_p^{q+1})^{3/4} \ln T_p^{q+1}$$

$$- A_p(q).$$

Last, we now bound the exploration terms $A_p(q)$ to show that exploration times are negligible. Writing $A_p(q) = \sum_{x \in \mathcal{X}_p(q)} \mathbb{1}[U_{t_p(q;x)} \leq 2p_{t_p(q;x)}]N_p(q;x)$, and because $\frac{N_p(q;x)}{t_p(q;x)^{1/4}} \leq 2^{2p+2}(T_p^{q+1})^{1/4}$, using Hoeffding's inequality we obtain that with probability at least $1 - e^{-\frac{(T_p^{q+1})^{1/4}}{2^{2p+3}}} := 1 - p_6(p,q)$,

(8) $$A_p(q) \leq \sum_{x \in \mathcal{X}_p(q)} \frac{N_p(q;x)}{t_p(q;x)^{1/4}} + (T_p^{q+1})^{7/8} \leq \frac{T_p^{q+1} - T_p^q}{(T_p^q)^{1/4}} + (T_p^{q+1})^{7/8} \leq 2(T_p^{q+1})^{7/8}.$$

Now recalling that $k(q) \leq \frac{q}{2^p}$, we have that

$$\sum_{p \geq 0} \sum_{q \geq p2^{p+5}} 2p_1(p,q) + p_2(p,q) + p_6(p,q) + k(q)(2p_1(p,q) + 2p_3(p,q))$$

$$+ (p_4(p,q) + k(q)p_5(p,q))(1 + T_p^{q+1} - T_p^q) < \infty.$$

As a result, the Borel-Cantelli lemma implies that on an event $\mathcal{E}$ of probability one, there exists $\hat{T}_1$ such that for any $p \geq 0$, $q \geq p2^{p+5}$ Eq (3), (4), (6) and (8) are satisfied, and (5) is satisfied for $q \geq l, p2^{p+5}$, and Eq (7) is satisfied for $T_p^q \leq T < T_p^{q+1}$.

We are now ready to prove the universal consistence of the learning rule. First, we pose $\epsilon_p = 2\frac{\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}}{2^{p/4}}$ and aim to show that the average error made by the learning rule on $\mathcal{T}_p$ is $\mathcal{O}(\epsilon_p)$ uniformly over time. Note in particular that $\sum_{p \geq 0} \epsilon_p < \infty$. For any $T \geq 1$, we define $\mathcal{R}_p(T) = \sum_{t \leq T, t \in \mathcal{T}_p} r_t$ the reward obtained by the learning rule, and $\bar{R}_p^*(T) = \sum_{t \leq T, t \in \mathcal{T}_p} \bar{r}(\pi^*(X_t), X_t)$ the reward obtained by the optimal policy. To do so, we first start by analyzing the regret on the first period $[1, 2^{32p})$ where there is no exploration and the learning rule uses EXP3.IX learners on each new instance. For $T < 2^{32p}$ let $\mathcal{X}_p(T) := \{X_t, t \in \mathcal{T}_p, t \leq T\}$. Note that $|\mathcal{X}_p(T)| \leq \frac{T}{4^p}$ by definition of $\mathcal{T}_p$. For $x \in \mathcal{X}_p(T)$, let $N_p(T; x) = |\{t \leq T, t \in \mathcal{T}_p, X_t = x\}| \leq 2^{2p+2}$ and $\bar{R}_p^0(T; x) := \mathbb{E}[\sum_{t \leq T, t \in \mathcal{T}_p, X_t = x} \tilde{r}_t \mid \mathbb{X}]$ where $\tilde{r}_t$ is the reward obtained if we used strategy 0. Now by Theorem 3.4, for every $x \in \mathcal{X}_p(T)$, with probability at least $1 - e^{-p^2 T^{1/2^7}}$, we have

$$\sum_{t \leq T, t \in \mathcal{T}_p(T), X_t = x} r_t(\pi^*(x), x) - \mathcal{R}_p(T) \leq cpT^{1/2^7}\sqrt{|\mathcal{A}|\ln|\mathcal{A}|N_p(T; x)}$$

As a result, with probability at least $1 - Te^{-p^2 T^{1/2^7}} := 1 - p_7(p, T)$,

$$\sum_{t \leq T, t \in \mathcal{T}_p} r_t(\pi^*(x), x) - r_t \leq \frac{2^p}{4^p}T + \sum_{x \in \mathcal{X}_p(T), N_p(T;x) \geq 2^p} \sum_{t \leq T, t \in \mathcal{T}_p, X_t = x} r_t(\pi^*(x), x) - r_t$$

$$\leq \frac{T}{2^p} + cp\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}T^{1/2^7} \sum_{x \in \mathcal{X}_p(T), N_p(T;x) \geq 2^p} \sqrt{N_p(T; x)}$$

$$\leq \frac{T}{2^p} + cp\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}T^{1/2^7}\frac{T}{2^{p/2}}$$

$$\leq \frac{T}{2^p} + \frac{c}{2}\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}T^{1-1/2^7}\log_2 T,$$

where in the last inequality, we used $2^{2p} \leq T < 2^{32p}$, thus $2^{p/2} \geq T^{1/64}$. Then, by Hoeffding's inequality, we have with probability $1 - e^{-2p^2\sqrt{T}} := 1 - p_8(p, T)$,

$$\sum_{t \leq T, t \in \mathcal{T}_p} r_t(\pi^*(x), x) \geq \bar{R}_p^*(T) - \frac{\log_2 T}{2}T^{3/4}.$$

Finally, with probability at least $1 - p_7(p, T) - p_8(p, T)$, we obtain

$$(9) \qquad \mathcal{R}_p(T) \geq \bar{R}_p^*(T) - \frac{1+c}{2}\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}T^{1-1/2^7}\log_2 T - \frac{T}{2^p}.$$

Noting that $\sum_{p \geq 0}\sum_{T \geq 1} p_7(p, T) + p_8(p, T) < \infty$, the Borel-Cantelli lemma implies that on an event $\mathcal{F}$ of probability one, there exists $\hat{T}_2$ such that for all $T \geq \hat{T}_2$, and $p \geq 0$ such that $T < 2^{32p}$, Eq (9) holds. We will now suppose that the event $\mathcal{E} \cap \mathcal{F}$ of probability one is met.

Next we consider the case of $T \geq 2^{32p}$, and let $q_0 \geq p2^{p+5}$ such that $T_p^{q_0} \leq T < T_p^{q_0+1}$. Then, consider

$$\mathcal{S}_p^0 := \left\{ p2^{p+5} \leq q < q_0 : \hat{R}_p^0(q) - \eta_p(T_p^{q+1} - T_p^q) \geq \max_{1 \leq l \leq k(q)} \hat{R}_p^k(q) \right\},$$

the set of phases where the learning rule estimated that strategy 0 performed better than strategy 1. Next, let $\mathcal{P}_p^i = \{p2^{p+5} \leq q < q_0 : \mathcal{P}_p(q) = i\}$ the set of phases where the learning

rule performed strategy $i$ for $i \in \{0,1\}$. An important observation is that for two phases $q_1 < q_2 \in \mathcal{S}_p^0 \cap \mathcal{P}_p^1$, if strategy 1 should not have been performed, then $q_2 > q_1 + p2^p$. In particular, we have $T_p^{q_1} \leq 2^{-p} T_p^{q_2}$, hence $T_p^{q_1+1} - T_p^{q_1} \leq 2^{-p}(T_p^{q_2+1} - T_p^{q_2})$. This allows to dissipate the errors made during phases where the algorithm performs strategy 1 by mistake. Precisely, using a descending induction we obtain

$$\sum_{q \in \mathcal{S}_p^0 \cap \mathcal{P}_p^1} T_p^{q+1} - T_p^q \leq \frac{T_p^{q_0} - T_p^{q_0-1}}{1 - 2^{-p}} \leq 2 \cdot 2^{-p} T_p^{q_0} \leq 2^{-p+1} T \leq 2\epsilon_p T.$$

On all other phases $\mathcal{P}_p^0 \cup (\mathcal{P}_p^1 \setminus \mathcal{S}_p^0)$, the performance of the learning rule is close to having performed strategy 0 on all phases. Indeed, using Eq (6) we obtain

$$\sum_{q \in (\mathcal{P}_p^1 \setminus \mathcal{S}_p^0)} \tilde{R}_p^1(q) \geq \sum_{q \in (\mathcal{P}_p^1 \setminus \mathcal{S}_p^0)} \max_{l=1,\ldots,k(q)} \bar{R}_p^l(q) - (T_p^{q+1})^{3/4} - c(T_p^{q+1})^{3/4} \ln T_p^{q+1} - A_p(q)$$

$$\geq \sum_{q \in (\mathcal{P}_p^1 \setminus \mathcal{S}_p^0)} \max_{l=1,\ldots,k(q)} \hat{R}_p^l(q) - \sum_{q < q_0} \left( 4(T_p^{q+1})^{7/8} + c(T_p^{q+1})^{3/4} \ln T_p^{q+1} \right)$$

$$\geq \sum_{q \in (\mathcal{P}_p^1 \setminus \mathcal{S}_p^0)} \hat{R}_p^0(q) - \eta_p T_p^{q_0} - 4(4 + c \ln T_p^{q_0}) T^{15/16}$$

$$\geq \sum_{q \in (\mathcal{P}_p^1 \setminus \mathcal{S}_p^0)} \bar{R}_p^*(q) - \eta_p T - 3\epsilon_p T - 4(5 + c \ln T) T^{15/16}.$$

In the second inequality, we used Eq (5) and in the third inequality, we used the definition of $\mathcal{S}_p^0$ and the identities $\sum_{q \leq q_0} (T_p^q)^{7/8} \leq (T_p^{q_0})^{7/8} \frac{2^p}{1-2^{-7/8}} \leq 2^{p+2}(T_p^{q_0})^{7/8} \leq 4T^{15/16}$. In the last inequality, we used Eq (3). Next, using Eq (4) we have directly

$$\sum_{q \in \mathcal{P}_p^0} \tilde{R}_p^0(q) \geq \sum_{q \in \mathcal{P}_p^0} \bar{R}_p^*(q) - 3\epsilon_p T - 3 \cdot 4T^{15/16}.$$

Combining the two above inequalities and observing that $\eta_p = 5\epsilon_p$ gives

$$\sum_{2^{32p} \leq t < T_p^{q_0}, t \in \mathcal{T}_p} r_t \geq \sum_{q \in \mathcal{P}_p^0} \tilde{R}_p^0(q) + \sum_{q \in \mathcal{P}_1 \setminus \mathcal{S}^0} \tilde{R}_p^1(q)$$

$$\geq \sum_{q \in \mathcal{P}_p^0 \cup (\mathcal{P}_p^1 \setminus \mathcal{S}_p^0)} \bar{R}_p^*(q) - 11\epsilon_p T - (32 + 4c \ln T) T^{15/16}$$

$$\geq \sum_{p2^{p+5} \leq q < q_0} \bar{R}_p^*(q) - \sum_{q \in \mathcal{S}_p^0 \cap \mathcal{P}_p^1} (T_p^{q+1} - T_p^q) - 11\epsilon_p T - (32 + 4c \ln T) T^{15/16}$$

$$\geq \sum_{p2^{p+5} \leq q < q_0} \bar{R}_p^*(q) - 13\epsilon_p T - (32 + 4c \ln T) T^{15/16}.$$

Now recalling the former estimate of $\mathcal{R}_p(T)$ for $T < 2^{32p}$, we obtain

$$\mathcal{R}_p(T) \geq \mathcal{R}_p(2^{32p} - 1) + \sum_{2^{32p} \leq t < T_p^{q_0}, t \in \mathcal{T}_p} r_t$$

$$\geq \bar{R}_p^*(T) - 2\frac{T}{2^p} - \frac{1+c}{2}\sqrt{|\mathcal{A}| \ln |\mathcal{A}|} T^{1-1/2^7} \log_2 T - 13\epsilon_p T - (32 + 4c \ln T) T^{15/16}$$

$$\geq \bar{R}_p^*(T) - \frac{1+c}{2}\sqrt{|\mathcal{A}| \ln |\mathcal{A}|} T^{1-1/2^7} \log_2 T - (32 + 4c \ln T) T^{15/16} - 15\epsilon_p T$$

where the term $\frac{T}{2^p}$ comes from the fact that $T - (T_p^{q_0} - 1) \le T_p^{q_0+1} - T_p^{q_0} \le \frac{T}{2^p}$. Now note that if $t \in \mathcal{T}_p$, there were at least $4^p$ duplicates, hence $t \ge 4^p$. As a result, we can always suppose without loss of generality that $T \ge 4^p$. Combining with the case $T < 2^{32p}$, we obtain that for all $T \ge \max(\hat{T}_1, \hat{T}_2)$, $p \ge 0$ with $t \ge 4^p$,

$$(10) \qquad \mathcal{R}_p(T) \ge \bar{R}_p^*(T) - (33 + 5c)\sqrt{|\mathcal{A}| \ln |\mathcal{A}|} T^{1-1/2^7} \log_2 T - 15\epsilon_p T.$$

This ends the proof that on times $\mathcal{T}_p$, the learning rule has an average error at most $\mathcal{O}(\epsilon_p)$ on the event $\mathcal{E} \cap \mathcal{F}$. Because $\sum_{p \ge 0} \epsilon_p < \infty$, we can afford to converge on each set $\mathcal{T}_p$ to the optimal policy independently.

Precisely, we aim to show that

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^T \bar{r}(\pi^*(X_t), X_t) - r_t \le 0, \quad (a.s.).$$

Fix $0 < \epsilon \le 1, \delta > 0$ and let $p_0$ such that $\sum_{p \ge p_0} \epsilon_p < \frac{\epsilon}{15}$. Because $\mathbb{X} \in \mathcal{C}_2$, by Proposition 3.2, $\mathbb{X}^{\le 4^{p_0}} \in \mathcal{C}_1$. As a result, because the sequence of policies $(\pi^l)_l$ is dense under $\mathcal{C}_1$ processes, there exists $l_0 \ge 1$ such that

$$\mathbb{E}\left[ \limsup_{T \to \infty} \frac{1}{T} \sum_{t \le T, t \in \mathcal{T}^{\le 4^{p_0}}} \mathbb{1}[\pi^*(X_t) \ne \pi^{l_0}(X_t)] \right] \le \frac{\epsilon\delta}{2^{2p_0+2} p_0}.$$

Then, by the dominated convergence theorem, there exists $T_0$ such that

$$\mathbb{E}\left[ \sup_{T \ge T_0} \frac{1}{T} \sum_{t \le T, t \in \mathcal{T}^{\le 4^{p_0}}} \mathbb{1}[\pi^*(X_t) \ne \pi^{l_0}(X_t)] \right] \le \frac{\epsilon\delta}{2^{2p_0+1} p_0}.$$

In particular, on an event $\mathcal{B}_\delta$ of probability at least $1 - \delta$, the Markov inequality yields that for all $T \ge T_0$,

$$\sum_{t \le T, t \in \mathcal{T}^{\le 4^{p_0}}} \mathbb{1}[\pi^*(X_t) \ne \pi^{l_0}(X_t)] \le \frac{\epsilon}{2^{2p_0+1} p_0} T.$$

In particular, the above equation holds if we replace $\mathcal{T}^{\le 4^{p_0}}$ by $\mathcal{T}_p$ for any $p < p_0$. Now suppose that the event $\mathcal{E} \cap \mathcal{F} \cap \mathcal{B}_\delta$ of probability at least $1 - \delta$ is met. For any $p < p_0$ and $q \ge p2^{p+5}$ such that $T_p^q \ge \hat{T} := \max(\hat{T}_1, \hat{T}_2, 2^{l_0}, 2^{32p_0})$, because $T_p^q \ge 2^{l_0}$, we have

$$\max_{1 \le l \le k(q)} \hat{R}_p^k(q) \ge \hat{R}_p^{l_0}(q)$$

$$\ge \bar{R}_p^{l_0}(q) - (T_p^{q+1})^{7/8}$$

$$\ge \bar{R}_p^*(q) - (T_p^{q+1})^{7/8} - \sum_{t \in \mathcal{T}_p(q)} \mathbb{1}[\pi^*(X_t) \ne \pi^{l_0}(X_t)]$$

$$\ge \hat{R}_p^0(q) - 2(T_p^{q+1})^{7/8} - 3\epsilon_p(T_p^{q+1} - T_p^q) - 2^{-2p-1} T_p^{q+1}$$

$$\ge \hat{R}_p^0(q) - 2(T_p^{q+1})^{7/8} - 4\epsilon_p(T_p^{q+1} - T_p^q).$$

where in the second inequality we used Eq (5) and in the fourth we used Eq (3). In the last inequality, we used $2^{-p-1} T_p^{q+1} \le T_p^{q+1} - T_p^q$. Now let $T_1$ such that $2T^{7/8} < \frac{\epsilon_p}{2^{p+1}} T$ for any $T \ge T_1$. Then, for any $p < p_0$ and $q \ge p2^{p+5}$ such that $T_p^q \ge \tilde{T} := \max(\hat{T}, T_1)$, we have

$$\max_{1 \le l \le k(q)} \hat{R}_p^k(q) > \hat{R}_p^0(q) - 5\epsilon_p(T_p^{q+1} - T_p^q),$$

which implies $\mathcal{P}_p(q+1) = 1$ since $\eta_p = 5\epsilon_p$ if $\mathcal{P}_p(q+1)$ was not already defined. In other terms, starting from time $2^{p_0}\tilde{T}$, the learning rule always chooses strategy 1 for categories $p < p_0$. We now bound the error of the learning rule on $\mathcal{T}_p$ for $p < p_0$. Let $\tilde{q}$ such that $T_p^{\tilde{q}-1} \leq 2^{p_0}\tilde{T} < T_p^{\hat{q}}$. For any $T \geq 2^{p_0}\tilde{T}$ and $q(T)$ such that $T_p^{q(T)} \leq T < T_p^{q(T)+1}$, we can write

$$
\mathcal{R}_p(T) - \bar{R}_p^*(T) \geq \sum_{\tilde{q}<q<q(T)} (\tilde{R}_p^1(q) - \bar{R}_p^*(q)) + \tilde{R}_p^1(q(T),T) - \sum_{t\in\mathcal{T}_p(q),t\leq T} \bar{r}(\pi^*(X_t),X_t)
$$

$$
- 2^{p_0}\tilde{T} - \sum_{q<q(T)} A_p(q)
$$

$$
\geq \sum_{\tilde{q}<q<q(T)} (R_p^{l_0}(q) - \bar{R}_p^*(q)) - \sum_{t\in\mathcal{T}_p(q),t\leq T} \mathbb{1}[\pi^*(X_t) \neq \pi^{l_0}(X_t)] - 2^{p_0}\tilde{T}
$$

$$
- \sum_{q\leq q(T)} (2A_p(q) + (T_p^{q+1})^{3/4} + c(T_p^{q+1})^{3/4}\ln T_p^{q+1})
$$

$$
\geq - \sum_{t\leq T,t\in\mathcal{T}_p} \mathbb{1}[\pi^*(X_t) \neq \pi^{l_0}(X_t)] - 2^{p_0}\tilde{T} - 4(3+c)(T_p^{q(T)+1})^{15/16}\ln T_p^{q(T)+1}
$$

$$
\geq -2^{p_0}\tilde{T} - 16(3+c)T^{15/16}\ln T - \frac{\epsilon}{2^{p_0}p_0}T.
$$

where in the second inequality we applied Eq (6) and Eq (7), and in the third inequality, we used the identity $\sum_{q\leq q(T)}(T_p^{q+1} - T_p^q)^{3/4} \leq 4(T_p^{q(T)+1})^{7/8}$ proved earlier. As a result, we can write

$$
\sum_{p<p_0} \bar{R}_p^*(T) - \mathcal{R}_p(T) \leq p_0 2^{p_0}\tilde{T} + 16p_0(3+c)T^{15/16}\ln T + \epsilon T.
$$

Now because the events $\mathcal{E}, \mathcal{F}$ are met, using Eq (10), we also have for $T \geq \tilde{T}$

$$
\sum_{p\geq p_0} \bar{R}_p^*(T) - \mathcal{R}_p(T) = \sum_{p_0\leq p<\log_4 T} \bar{R}_p^*(T) - \mathcal{R}_p(T)
$$

$$
\leq (17+3c)\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}T^{1-1/2^7}(\log_2 T)^2 + 15\sum_{p\geq p_0} \epsilon_p \cdot T
$$

$$
\leq (17+3c)\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}T^{1-1/2^7}(\log_2 T)^2 + \epsilon T
$$

Summing the two above inequalities gives

$$
\sum_{t=1}^T \bar{r}(\pi^*(X_t),X_t) - r_t \leq p_0 2^{p_0}\tilde{T} + 16p_0(3+c)T^{15/16}\ln T
$$

$$
+ (17+3c)\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}T^{1-1/2^7}(\log_2 T)^2 + 2\epsilon T.
$$

As a result, on the event $\mathcal{E} \cap \mathcal{F} \cap \mathcal{G} \cap \mathcal{B}_\delta$ of probability at least $1 - \delta$, we have

$$
\limsup_{T\to\infty} \frac{1}{T} \sum_{t=1}^T \bar{r}(\pi^*(X_t),X_t) - r_t \leq 2\epsilon.
$$

Because this holds for any $\delta > 0$ and $0 < \epsilon < 1$, this shows that almost surely, we have $\limsup_{T\to\infty} \frac{1}{T}\sum_{t=1}^T \bar{r}(\pi^*(X_t),X_t) - r_t \leq 0$. We denote by $\mathcal{C}$ this event. We now formally

show that the learning rule is universally consistent. Let $\pi : \mathcal{X} \to \mathcal{A}$ be a measurable function. First, by the Hoeffding inequality, we have for $T \geq 1$,

$$\mathbb{P}\left[\left|\sum_{t=1}^{T} r_t(\pi(X_t), X_t) - \bar{r}_t(\pi(X_t), X_t)\right| \leq T^{3/4}\right] 1 - e^{-2\sqrt{T}}.$$

As a result, the Borel-Cantelli lemma implies that on an event $\mathcal{H}$ of probability one, there exists $\hat{T}_4$ such that for all $T \geq \hat{T}_4$, $|\sum_{t=1}^{T} r_t(\pi(X_t), X_t) - \bar{r}_t(\pi(X_t), X_t)| \leq T^{3/4}$. Then, on $\mathcal{C} \cap \mathcal{H}$ of probability one, for any $T \geq \hat{T}_4$ we have

$$\sum_{t=1}^{T} r(\pi(X_t), X_t) - r_t \leq \sum_{t=1}^{T} \bar{r}(\pi(X_t), X_t) - r_t + T^{3/4}$$

$$\leq \sum_{t=1}^{T} \bar{r}(\pi^*(X_t), X_t) - r_t + T^{3/4}.$$

Thus, $\limsup_{T \to \infty} \sum_{t=1}^{T} \bar{r}(\pi(X_t), X_t) - r_t \leq 0$. This ends the proof that the learning rule is universally consistent under any $\mathcal{C}_2$ process. Now recall that $\mathcal{C}_2$ is a necessary condition for universal learning by Theorem 4.1. Hence, the set of learnable processes is exactly $\mathcal{C} = \mathcal{C}_2$ and the learning rule is optimistically universal. □

## 5. Countably infinite action space.
We next turn to the case where the action space is infinite $|\mathcal{A}| = \infty$ but countable. The goal of this section is to show that the set of processes admitting universal learning now becomes $\mathcal{C}_1$. This contrasts with the full-feedback setting where universal learning is optimistically achievable under $\mathcal{C}_2$ processes when a property F-TiME on the value space $(\mathcal{Y}, \ell)$ is satisfied [7]. Intuitively, this asks that mean-estimation is possible in finite time for any prescribed error tolerance. Of interest to the discussion of this section with countable number of actions, [7] showed that countably-infinite classification $(\mathcal{Y}, \ell) = (\mathbb{N}, \ell_{01})$ satisfies the F-TiME property and, their learning rule is universally consistent under $\mathcal{C}_2$ processes even under noisy and adversarial responses.

For countable action sets, there is a simple optimistically universal learning rule defined as follows. From [20, Lemma 24], because $\mathcal{A}$ is countable, the $0 - 1$ loss on $\mathcal{A}$ is a separable metric, thus, there exists a countable set $\Pi$ of measurable policies $\pi : \mathcal{X} \to \mathcal{A}$ such that for every $\mathbb{X} \in \mathcal{C}_1$, for every measurable $\pi^\star : \mathcal{X} \to \mathcal{A}$,

$$\mathbb{E}\left[\inf_{\pi \in \Pi} \hat{\mu}_{\mathbb{X}}(\{x : \pi(x) \neq \pi^\star(x)\})\right] \leq \inf_{\pi \in \Pi} \mathbb{E}\left[\hat{\mu}_{\mathbb{X}}(\{x : \pi(x) \neq \pi^\star(x)\})\right] = 0,$$

which implies in particular that almost surely, $\inf_{\pi \in \Pi} \hat{\mu}_{\mathbb{X}}(\{x : \pi(x) \neq \pi^\star(x)\}) = 0$. Enumerate $\Pi = \{\pi_1, \pi_2, \ldots\}$. For any $\mathbb{X}$, we consider the countable set of experts $\{E_1, E_2, \ldots\}$ such that $E_{i,t} = \pi_i(X_t)$. Our learning rule then applies EXPINF from Corollary 3.5 with this family of experts.

THEOREM 5.1. *Let $\mathcal{X}$ be a separable Borel-metrizable space and $\mathcal{A}$ a countable infinite action set. Then, there is an optimistically universal learning rule and the set of learnable processes is $\mathcal{C} = \mathcal{C}_1$.*

PROOF. We start by showing that the learning rule defined above is universally consistent on any $\mathbb{X} \in \mathcal{C}_1$ process. This proof is essentially identical to that of [22, Theorem 1]. Indeed,

denoting by $\hat{a}_t$ the action selected by the learning rule at time $t$, Corollary 3.5 implies that on an event $\mathcal{E}$ of probability one, for any $\pi \in \Pi$, we have

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} r_t(\pi(X_t)) - r_t(\hat{a}_t) \leq 0.$$

Now fix a measurable policy $\pi^\star : \mathcal{X} \to \mathcal{A}$. For any $\pi \in \Pi$, because the rewards lie in $[0, 1]$, on $\mathcal{E}$,

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} r_t(\pi^*(X_t)) - r_t(\hat{a}_t)$$

$$\leq \hat{\mu}_{\mathbb{X}}(\{x : \pi(x) \neq \pi^*(x)\}) + \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} r_t(\pi(X_t)) - r_t(\hat{a}_t)$$

$$\leq \hat{\mu}_{\mathbb{X}}(\{x : \pi(x) \neq \pi^*(x)\}).$$

Also, by construction of the countable set $\Pi$, on an event $\mathcal{F}$ of probability one, we have $\inf_{\pi \in \Pi} \hat{\mu}_{\mathbb{X}}(\{x : \pi(x) \neq \pi^\star(x)\}) = 0$. Thus, on $\mathcal{E} \cap \mathcal{F}$, the above inequality shows that $\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} r_t(\pi^*(X_t)) - r_t(\hat{a}_t) \leq 0$. Hence, the learning rule is universally consistent under $\mathcal{C}_1$ processes with adversarial responses.

Next, we show that the condition $\mathbb{X} \in \mathcal{C}_1$ is necessary for the existence of a universally consistent learning rule, even for function learning. Let $\mathbb{X}$ be any process with $\mathbb{X} \notin \mathcal{C}_1$. By Lemma 3.1, there exists a sequence $\{B_i\}_{i=1}^{\infty}$ of disjoint measurable subsets of $\mathcal{X}$ with $\bigcup_{i \in \mathbb{N}} B_i = \mathcal{X}$, and a sequence $\{N_i\}_{i=1}^{\infty}$ in $\mathbb{N}$ such that, on a $\sigma(\mathcal{X})$-measurable event $\mathcal{E}_0$ of probability strictly great than zero,

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}] > 0,$$

where $i_t$ is the unique $i \in \mathbb{N}$ with $X_t \in B_i$.

Next, we define the function $f^\star$. Enumerate $\mathcal{A} = \{a_1, a_2, \ldots\}$, and for each $i \in \mathbb{N}$, let $A_i = \{a_1, \ldots, a_{2N_i}\}$. For each $i \in \mathbb{N}$, let $a_i^\star$ be an element of $A_i$. Denote by $\bar{a} = \{a_i^\star\}_{i \in \mathbb{N}}$. Then for each $i \in \mathbb{N}$ and each $x \in B_i$, define $f_{\bar{a}}^\star(x, a) = \mathbb{1}[a = a_i^\star]$. Also define $\mathbf{a}_i^\star$ as $\mathrm{Uniform}(A_i)$ (independent over $i$ and all independent of $\mathbb{X}$ and the randomness of the learning rule), and $\bar{\mathbf{a}} = \{\mathbf{a}_i^\star\}_{i \in \mathbb{N}}$. Then for any learning rule $\hat{f}_t$, denoting by $\hat{a}_t$ its actions when $f^\star = f_{\bar{\mathbf{a}}}^\star$ is as constructed above, we have

$$\sup_{\bar{a}} \mathbb{E}\left[ \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \left( \sup_{a \in \mathcal{A}} r_t(a) - r_t(\hat{a}_t) \right) \middle| \bar{\mathbf{a}} = \bar{a} \right]$$

$$= \sup_{\bar{a}} \mathbb{E}\left[ \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[\hat{a}_t \neq \mathbf{a}_{i_t}] \middle| \bar{\mathbf{a}} = \bar{a} \right]$$

$$\geq \mathbb{E}\left[ \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[\hat{a}_t \neq \mathbf{a}_{i_t}] \right]$$

$$\geq \mathbb{E}\left[ \mathbb{1}_{\mathcal{E}_0} \cdot \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}] \mathbb{1}[\hat{a}_t \neq \mathbf{a}_{i_t}] \right].$$

By the law of total expectation, this last expression above equals

$$\mathbb{E}\left[\mathbb{1}_{\mathcal{E}_0} \cdot \mathbb{E}\left[\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}] \mathbb{1}[\hat{a}_t \neq \mathbf{a}_{i_t}] \Big| \mathbb{X}, \hat{f}_\cdot \right]\right],$$

where conditioning on $\hat{f}_\cdot$ indicates we condition on the independent randomness of the learning rule. Since the average is bounded for any fixed $T$, Fatou's lemma, together with the fact that $\mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}]$ is $\sigma(\mathbb{X})$-measurable, imply the expression above is at least as large as

(11) $$\mathbb{E}\left[\mathbb{1}_{\mathcal{E}_0} \cdot \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}] \mathbb{P}\left(\hat{a}_t \neq \mathbf{a}_{i_t} \Big| \mathbb{X}, \hat{f}_\cdot \right)\right].$$

Let $\hat{N}_t = |\mathbb{X}_{\leq t} \cap B_{i_t}|$ and $\hat{A}_t = \{\hat{a}_{t'} : t' \leq t, i_{t'} = i_t\}$. Note that, conditioned on $\hat{f}_\cdot$ and $\mathbb{X}$, the probability that $\mathbf{a}_{i_t} \in \hat{A}_t$ is at most $\hat{N}_t \frac{1}{|A_{i_t}|} = \frac{\hat{N}_t}{2N_{i_t}}$. In particular, this implies that if $\hat{N}_t \leq N_{i_t}$, the conditional probability (given $\hat{f}_\cdot$ and $\mathbb{X}$) that $\hat{a}_t \neq \mathbf{a}_{i_t}$ is at least $1 - \frac{\hat{N}_t}{2N_{i_t}} \geq \frac{1}{2}$. Thus, (11) is no smaller than

(12) $$\mathbb{E}\left[\mathbb{1}_{\mathcal{E}_0} \cdot \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}] \cdot \frac{1}{2}\right].$$

By definition of the event $\mathcal{E}_0$, there is a nonzero probability that

$$\mathbb{1}_{\mathcal{E}_0} \cdot \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}] > 0,$$

and since the quantity on the left hand size is non-negative, this further implies the expectation in (12) is also strictly greater than zero.

Altogether, this implies there exists a (non-random) choice of $\bar{a}$ such that, choosing $f^\star = f_{\bar{a}}^\star$, the actions $\hat{a}_t$ made by the learning rule $\hat{f}_t$ satisfy

$$\mathbb{E}\left[\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \left(\sup_{a \in \mathcal{A}} r_t(a) - r_t(\hat{a})\right)\right] > 0,$$

and since the quantity in the expectation is non-negative, this further implies that for this choice of $f^\star$, with non-zero probability,

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \left(\sup_{a \in \mathcal{A}} r_t(a) - r_t(\hat{a})\right) > 0.$$

Thus, $\hat{f}_t$ is not universally consistent for function learning. Since this holds for any choice of learning rule $\hat{f}_\cdot$, this completes the proof. $\qquad\square$

**6. Uncountable action space.** We next consider the case of uncountable action spaces. In this section, we assume that $\mathcal{A}$ is an uncountable separable Borel metrizable space. In this case, we will show that universal consistency is impossible even in the simplest setting where rewards are a deterministic, i.e., $r_t(a) = f^*(X_t, a)$ for some unknown measurable function $f^* : \mathcal{X} \times \mathcal{A} \to [0, 1]$. The argument is based on a dichotomy depending whether there exists a non-atomic probability measure $\mu$ on $\mathcal{A}$, i.e., such that for all $a \in \mathcal{A}$, we have $\mu(\{x\}) = 0$. If this is not the case, we will need the following simple result which states that any stochastic process $\mathbb{X}$ takes values in a countable set $Supp(\mathbb{X})$ almost surely.

LEMMA 6.1. *Let $\mathcal{X}$ a metrizable separable Borel space such that there does not exist a non-atomic probability measure on $\mathcal{X}$. Then, for any random variable $X$ on $\mathcal{X}$ there exists a countable set $Supp(X) \subset \mathcal{X}$ such that almost surely, $X \in Supp(X)$. Similarly, for any stochastic process $\mathbb{X}$ on $\mathcal{X}$ there exists a countable set $Supp(\mathbb{X}) \subset \mathcal{X}$ such that almost surely $\forall t \geq 1, X_t \in Supp(\mathbb{X})$.*

PROOF. Fix $\mathcal{X}$ such a space and let $X$ be a random variable on $\mathcal{X}$. Let $Supp(X) = \{x \in \mathcal{X} : \mathbb{P}[X = x] > 0\}$. Suppose by contradiction that $\mathbb{P}[X \notin Supp(X)] > 0$ and denote $\mathcal{E}$ the corresponding event. Because $\mathbb{P}[\mathcal{E}] > 0$ we can consider a random variable $Y \sim X|\mathcal{E}$. For instance take $(X_i)_{i \geq 1}$ an i.i.d. process following the distribution of $X$, fix $x_0 \in \mathcal{X}$ a fixed arbitrary instance, and pose

$$Y = \begin{cases} X_{\hat{k}} & \text{if } \{i \geq 1 : X_i \notin Supp(X)\} \neq \emptyset, \quad \hat{k} = \min\{i \geq 1 : X_i \notin Supp(X)\}, \\ x_0 & \text{otherwise.} \end{cases}$$

Because the first time $k$ such that $X_k \notin Supp(X)$ is a geometric variable $\mathcal{G}(1 - \mathbb{P}[\mathcal{E}])$, the event $\mathcal{F} = \{\exists i \geq 1 : X_i \notin Supp(X)\}$ has probability one. We now show that $Y$ is non-atomic. First observe that $Y \notin Supp(X)$. Then, if $x \in \mathcal{X} \notin Supp(X)$, we have

$$\mathbb{P}[Y = x] = \mathbb{P}[\{Y = x\} \cap \mathcal{F}] = \mathbb{P}[\{X_{\hat{k}} = x\} \cap \mathcal{F}] \leq \mathbb{P}\left[\bigcup_{i \geq 1}\{X_i = x\}\right] \leq \sum_{i \geq 1} \mathbb{P}[X_i = x] = 0.$$

where in the first equality we used the fact that $\mathbb{P}[\mathcal{F}^c] = 0$. Therefore $Y$ is non-atomic which contradicts the hypothesis on $\mathcal{X}$. As a result, almost surely $X \in Supp(X)$. It now suffices to check that $Supp(X)$ is countable, which is guaranteed by the identity $1 = \mathbb{P}[X \in Supp(X)] = \sum_{x \in Supp(X)} \mathbb{P}[X = x]$, since each term of the sum is positive. This ends the proof of the first claim.

Now let $\mathbb{X}$ be a stochastic process on $\mathcal{X}$ and define $Supp(\mathbb{X}) = \bigcup_{t \geq 1} Supp(X_t)$. Then $Supp(\mathbb{X})$ is countable as countable union of countable sets and

$$\mathbb{P}[\exists t \geq 1 : X_t \notin Supp(\mathbb{X})] \leq \sum_{t \geq 1} \mathbb{P}[X_t \notin Supp(\mathbb{X})] \leq \sum_{t \geq 1} \mathbb{P}[X_t \notin Supp(X_t)] = 0.$$

This ends the proof of the lemma. $\square$

We are now ready to show that no process admits universal learning when the action set is uncountable.

THEOREM 6.2. *If $\mathcal{A}$ is an uncountable separable Borel metrizable space, then there does not exist any $\mathbb{X}$ admitting universal consistency for measurable function learning.*

PROOF. Fix a learning rule $f.$ and for any $a^* \in \mathcal{A}$, we define the reward function $f_{a^*}^*(x, a) = \mathbb{1}[a = a^*]$ for $x \in \mathcal{X}, a \in \mathcal{A}$. We also define the policy $\pi_{a^*} : x \in \mathcal{X} \mapsto a^* \in \mathcal{A}$. We first consider the case where there exists a non-atomic probability measure $\mu$ on $\mathcal{A}$. Then, for any $t \geq 1$, and consider the case where $a^*$ is sampled from the distribution $\mu$ independently from the process $\mathbb{X}$ and the randomness of the learning rule. Then we have

$$\mathbb{P}_{a^* \sim \mu}[f_t(\mathbb{X}_{<t}, (0)_{<t}, X_t) = a^*] = \mathbb{E}_{\mathbb{X}, f_t}[\mathbb{P}_{a^* \sim \mu}(f_t(\mathbb{X}_{<t}, (0)_{<t}, X_t) = a^*)] = 0.$$

Denote by $\mathcal{E}_t$ this event. Then, by the union bound, $\mathbb{P}[\bigcap_{t \geq 1} \mathcal{E}_t] = 1$. The law of total probability implies that there exists a deterministic choice of $a^*$ such that

$$\mathbb{P}[\forall t \geq 1, f_t(\mathbb{X}_{<t}, (0)_{<t}, X_t) \neq a^*] = 1,$$

where the probability is taken over $\mathbb{X}$ and the randomness of the learning rule.

Now suppose that there does not exist non-atomic probability measures on $\mathcal{A}$. From Lemma 6.1, for any probability measure $\mu$ on $\mathcal{A}$, we can construct a countable set $Supp(\mu) \subset \mathcal{X}$ such that $\mu(Supp(\mu)) = 1$. Now consider the set

$$S = \bigcup_{t \geq 1} Supp(f_t(\mathbb{X}_{\leq t-1}, (0)_{\leq t-1}, X_t)).$$

Then, $S$ is countable as the union of countable sets. Since $\mathcal{A}$ is uncountable, let $a^* \in \mathcal{A} \setminus S$. By construction, on an event of probability one, for all $t \geq 1$, we have $f_t(\mathbb{X}_{\leq t-1}, (0)_{\leq t-1}, X_t) \neq a^*$.

In both cases, we found an action $a^* \in \mathcal{A}$ such that on an event $\mathcal{E}$ of probability one over $\mathbb{X}$ and the randomness of the learning rule, having received 0 reward in the past history at time step $t$, the learning rule does not select $a^*$, hence receives reward 0 at time $t$ as well. Thus, by induction, denoting by $\hat{a}_t$ the action selected by the learning rule at time $t$ for reward $f_{a^*}^*$, we have $\mathcal{E} \subset \{\forall t \geq 1, \hat{a}_t \neq a^*\}$. Thus, on $\mathcal{E}$,

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} r_t(\pi^*(X_t)) - r_t(\hat{a}_t) = 1.$$

Because $\mathcal{E}$ has probability one, this shows that $f.$ is not universally consistent. $\square$

## 7. Universal learning under continuity assumptions.
In Section 6 we showed that for general uncountable separable metric actions spaces, without further assumptions on the rewards, one cannot achieve universal consistency. The goal of this section is to show that adding mild continuity assumptions on the rewards enables to significantly enlarge the set of processes admitting universal learning.

7.1. *Continuous rewards.* In this section, we suppose that the rewards are continuous as defined in Definition 2.2, and show that universal consistency on $\mathcal{C}_1$ processes is still achievable. For bounded separable metric action spaces $(\tilde{\mathcal{A}}, \tilde{d})$, [20] showed that there is countable set of measurable policies $\Pi$ such that for any measurable $\pi^* : \mathcal{X} \to \tilde{\mathcal{A}}$ and $\mathbb{X} \in \mathcal{C}_1$,

$$\inf_{\pi \in \Pi} \mathbb{E}\left[ \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \tilde{d}(\pi^*(X_t), \pi(X_t)) \right] = 0.$$

In general, the action space $(\mathcal{A}, d)$ is unbounded, however, $(\mathcal{A}, d \wedge 1)$ is a separable bounded metric space on which we can apply the above result. This provides a countable set of measurable policies $\Pi$ such that for any measurable $\pi^* : \mathcal{X} \to \mathcal{A}$ and $\mathbb{X} \in \mathcal{C}_1$,

$$\inf_{\pi \in \Pi} \mathbb{E}\left[ \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \tilde{d}(\pi^*(X_t), \pi(X_t)) \wedge 1 \right] = 0.$$

From this observation, we can get the following lemma.

LEMMA 7.1. *Let $\mathcal{X}$ be a separable metrizable Borel space and $(\mathcal{A}, d)$ be a separable metric space. For any measurable function $\pi^* : \mathcal{X} \to \mathcal{A}$, on an event of probability one, for all $i \geq 1$, there exists $\pi^i \in \Pi$ such that*

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[d(\pi^*(X_t), \pi^i(X_t)) \geq 2^{-i}] \leq 2^{-i},$$

*for all $i \geq 1$, $\frac{1}{T} \sum_{t \leq T} r_t(\pi^i(X_t)) - \bar{r}_t(\pi^i(X_t)) \to 0$ and similarly for $\pi^*$.*

PROOF. By construction of the countable set of policies $\Pi$, for any $i \geq 1$, there exists $\pi^i \in \Pi$ such that

$$\mathbb{E}\left[\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} d(\pi^*(X_t), \pi^i(X_t)) \wedge 1\right] \leq 2^{-3i}.$$

Then, Markov's inequality implies that with probability at least $1 - 2^{-i}$.

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} d(\pi^*(X_t), \pi^i(X_t)) \wedge 1 \leq 2^{-2i}.$$

Applying Markov's inequality a second time, we obtain

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[d(\pi^*(X_t), \pi^i(X_t)) \geq 2^{-i}] \leq 2^i \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} d(\pi^*(X_t), \pi^i(X_t)) \wedge 1 \leq 2^{-i}.$$

The Borel-Cantelli lemma implies that on an event $\mathcal{E}$ of probability one, for $i$ sufficiently large, there exists $\pi^i \in \Pi$ with $\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[d(\pi^*(X_t), \pi^i(X_t)) \geq 2^{-i}] \leq 2^{-i}$. Clearly, this implies that this is the case for all $i \geq 1$. For any $i \geq 1$, Azuma's inequality implies that with probability at least $1 - 4e^{-2i\sqrt{T}}$, we have

$$\left|\sum_{t=1}^{T} r_t(\pi^i(X_t)) - \bar{r}_t(\pi^i(X_t))\right|, \left|\sum_{t=1}^{T} r_t(\pi^*(X_t)) - \bar{r}_t(\pi^*(X_t))\right| \leq 2iT^{3/4}.$$

Because $\sum_{T \geq 1} \sum_{i \geq 1} e^{-2i\sqrt{T}} < \infty$, the Borel-Cantelli lemma implies that on an event $\mathcal{F}$ of probability one, for all $i \geq 1$, $\frac{1}{T} \sum_{t \leq T} r_t(\pi^i(X_t)) - \bar{r}_t(\pi^i(X_t)) \to 0$ and similarly for $\pi^*$. Therefore, on the event $\mathcal{E} \cap \mathcal{F}$ of probability one, all events are satisfied, which ends the proof of the lemma. $\square$

Using Lemma 7.1, we will show that the EXPINF algorithm over the set of policies $\Pi$ is optimistically universal for continuous rewards.

THEOREM 7.2. *Let $(\mathcal{A}, d)$ be an infinite separable metric space. Then, EXPINF is optimistically univesal for continuous rewards and the set of learnable processes for continuous rewards is $\mathcal{C}^c = \mathcal{C}_1$.*

PROOF. We start by showing that EXPINF is universally consistent under continuous rewards under $\mathcal{C}_1$ processes. Let $\mathbb{X} \in \mathcal{C}_1$ and continuous rewards $(r_t)_t$ and let $\pi^* : \mathcal{X} \to \mathcal{A}$ be measurable policy. We denote $\mathcal{E}$ the event on which the guarantee for EXPINF of Corollary 3.5 holds. For convenience, we also note $\hat{a}_t$ the action selected by the learning rule at time $t$. For any $x \in \mathcal{X}$, and $\epsilon > 0$, we define

$$\Delta_\epsilon(x) = \sup_{a \in \mathcal{A}: d(a, \pi^*(x)) \leq \epsilon} |\bar{r}(a, x) - \bar{r}(\pi^*(x), x)|.$$

Next, fix $\delta > 0$, and for any $\epsilon > 0$, let $A(\epsilon, \delta) = \{x \in \mathcal{X} : \Delta_\epsilon(x) \geq \delta\}$. Note that for any $x \in \mathcal{X}$, by continuity of $\bar{r}(\cdot, x)$, for any $\delta > 0$, $\bigcap_{\epsilon > 0} A(\epsilon, \delta) = \emptyset$. By Lemma 7.1, on an event $\mathcal{F}$ of probability one, for any $i \geq 1$, there exists $\pi^i \in \Pi$ such that

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} \mathbb{1}[d(\pi^*(X_t), \pi^i(X_t) \geq 2^{-i}] \leq 2^{-i},$$

$\frac{1}{T}\sum_{t\leq T} r_t(\pi^i(X_t) - \bar{r}_t(\pi^i(X_t)) \to 0$ and similarly for $\pi^*$. As a result, on $\mathcal{F}$, for any $i \geq 1$,

$$\limsup_{T\to\infty} \frac{1}{T}\sum_{t\leq T} r_t(\pi^*(X_t), X_t) - r_t(\pi^i(X_t), X_t)$$

$$\leq \hat{\mu}_{\mathbb{X}}(A(2^{-i}, \delta)) + 2^{-i} + \limsup_{T\to\infty} \frac{1}{T} \sum_{\substack{t\leq T, \\ d(\pi^*(X_t), \pi^i(X_t))<2^{-i} \\ \Delta_{2^{-i}}(X_t)<\delta}} \bar{r}_t(\pi^*(X_t), X_t) - \bar{r}_t(\pi^i(X_t), X_t)$$

$$\leq \hat{\mu}_{\mathbb{X}}(A(2^{-i}, \delta)) + 2^{-i} + \delta.$$

Because $\mathbb{X} \in \mathcal{C}_1$ and $A(2^{-i}, \delta) \downarrow \emptyset$, on an event $\mathcal{G}(\delta)$ of probability one, we have that $\hat{\mu}_{\mathbb{X}}(A(2^{-i}, \delta)) \xrightarrow[i\to\infty]{} 0$. Last, let $\delta_j = 2^{-j}$ for any $j \geq 0$. On the event $\mathcal{E} \cap \mathcal{F} \cap \bigcap_{j\geq 0} \mathcal{G}(\delta_j)$ of probability one, combining Corollary 3.5 together with the above inequality implies that for any $j \geq 0$,

$$\limsup_{T\to\infty} \frac{1}{T}\sum_{t\leq T} r_t(\pi^*(X_t), X_t) - r_t(\hat{a}_t, X_t) \leq \delta_j.$$

Thus, $\limsup_{T\to\infty} \frac{1}{T}\sum_{t\leq T} r_t(\pi^*(X_t), X_t) - r_t(\hat{a}_t, X_t) \leq 0$ $(a.s.)$, which shows that EXPINF is universally consistent under $\mathbb{X}$ for stationary rewards. This ends the proof of the theorem.

We now show that $\mathcal{C}_1$ is necessary for universal consistency. The proof is analogous to that of Theorem 5.1 in which we proved that for unrestricted rewards on countably infinite action sets, $\mathcal{C}_1$ is necessary for universal learning. Suppose that $\mathbb{X} \in \mathcal{C}_1$ and let $f.$ be a learning rule. Using the same arguments, there exist a partition of $\mathcal{X}$ in measurable sets $\{B_i\}_{i\geq 1}$ and a sequence $\{N_i\}_{i\geq 1}$ of integers such that with non-zero probability,

$$\limsup_{T\to\infty} \frac{1}{T}\sum_{t=1}^{T} \mathbb{1}[|\mathbb{X}_{<t} \cap B_{i_t}| < N_{i_t}] > 0,$$

where $i_t$ is the index such that $X_t \in B_i$. As in the original proof, let $\{a_i, i \geq 1\}$ be a sequence of distinct actions and let $A_i = \{a_1, \ldots, a_{2N_i}\}$ for $i \geq 1$. We also define $\epsilon_i = \min_{a\neq a'\in A_i} d(a, a')$ the minimum distance within $A_i$ actions. For any sequence $\bar{a} = \{a_i^*\}_{i\in\mathbb{N}}$ where $a_i^* \in A_i$ for $i \geq 1$, we define a deterministic reward $r_{\bar{a}}^*$ with

$$r_{\bar{a}}^*(a, x) = \max\left(1 - \frac{2d(a, a_i^*)}{\epsilon_i}, 0\right),$$

for any $x \in B_i$, which defines a proper measurable continuous reward. We also define the rewards $\tilde{r}_{\bar{a}}^*(a, x) = \mathbb{1}[a = a_i^*]$ for $a \in \mathcal{A}$ and $x \in B_i$. We now define the learning rule $\tilde{f}.$ which at each step $t$ computes the action $\hat{a}$ chosen by the learning rule $f.$, selects the action $\tilde{a}_t := \arg\min_{a'\in A_i} d(\hat{a}, a')$ where $i \geq 1$ is the unique index with $X_t \in B_i$, receives a reward $r_t$, then reports the reward $\max\left(1 - \frac{2d(\hat{a}, \tilde{a})}{\epsilon_i}, 0\right)$, which will be then used by $f.$ for future action selections. Note that on $B_i$, the rewards $r_{\bar{a}}^*$ were defined so that they are identically zero outside of the balls $B_d(a, \epsilon_i)$ for $a \in A_i$. These are disjoint, so the report of reward given by $\tilde{f}.$ to its internal run of $f.$ coincides exactly with what $f.$ would have received by selecting action $\hat{a}$ instead of $\tilde{a}$. Further, one can observe that selecting one of the nearest element within $A_i$ always increases the reward because the balls $B_d(a, \epsilon_i)$ for $a \in A_i$ are disjoint. Therefore, $\tilde{f}.$ always receives higher reward than $f.$ at any step. Now observe that $\tilde{f}.$ always observes a reward in $\{0, 1\}$. Hence, for any choice of $\bar{a}$, at any step $t$, $\tilde{f}_t$ has the same rewards on $r_{\bar{a}}^*$ as

it would have obtained on the rewards $\tilde{r}_{\bar{a}}^*$. Therefore,

$$\limsup_{T\to\infty} \frac{1}{T} \sum_{t=1}^{T} \left( \sup_{a\in\mathcal{A}} r_{\bar{a},t}^*(a) - r_{\bar{a},t}^*(\hat{a}_t) \right) \geq \limsup_{T\to\infty} \frac{1}{T} \sum_{t=1}^{T} \left( \sup_{a\in\mathcal{A}} r_{\bar{a},t}^*(a) - r_{\bar{a},t}^*(\tilde{a}_t) \right)$$

$$= \limsup_{T\to\infty} \frac{1}{T} \sum_{t=1}^{T} \left( \sup_{a\in\mathcal{A}} \tilde{r}_{\bar{a},t}^*(a) - \tilde{r}_{\bar{a},t}^*(\tilde{a}_t) \right),$$

where $\hat{a}_t$ (resp. $\tilde{a}_t$) denotes the action selected by $f_.$ (resp. $\tilde{f}_.$) at time $t$. However, the proof of Theorem 5.1 precisely shows that there exists a choice of $\bar{a}$ such that with non-zero probability, $\limsup_{T\to\infty} \frac{1}{T} \sum_{t=1}^{T} \left( \sup_{a\in\mathcal{A}} \tilde{r}_{\bar{a},t}^*(a) - \tilde{r}_{\bar{a},t}^*(\tilde{a}_t) \right) > 0$. Now observe that the measurable function $\pi(x) = a_i^*$ where $x \in B_i$ always selects the best action. This show that $f_.$ is not consistent on rewards $r_{\bar{a}}^*$, hence not universally consistent. This shows that $\mathbb{X} \notin \mathcal{C}^c$ and completes the proof of the theorem. □

### 7.2. Uniformly-continuous rewards.

In the last section, we showed that adding a continuity constraint on the rewards allowed to learn $\mathcal{C}_1$ processes even when the action space $\mathcal{A}$ is infinite. Unfortunately, this additional assumption on the rewards is not sufficient to obtain universal consistency on the more general class of processes $\mathcal{C}_2$. In this section, we strengthen the assumptions on the rewards and suppose that they are uniformly-continuous in the actions as per Definition 2.2.

We start by giving necessary conditions for uniformly-continuous rewards. To do so, we will need the following simple reduction, showing that some necessary conditions provided in the unrestricted rewards case can be used in the uniformly-continuous setting as well.

LEMMA 7.3. *Let $(\mathcal{A}, d)$ be a separable metric space. Let $S \subset \mathcal{A}$ such that we have $\min_{a,a'\in S} d(a, a') > 0$. Then, $\mathcal{C}^{uc}(\mathcal{A}) \subset \mathcal{C}(S)$.*

PROOF. Intuitively, we restrict the problem on $\mathcal{A}$ to the actions $S$. Formally, let $\eta = \frac{1}{3} \min_{a,a'\in S} d(a, a')$ and observe that any reward function $r : S \to [0, \bar{r}]$ can be extended to a uniformly-continuous function $F(r) : \mathcal{X} \to \mathcal{A}$ as follows.

$$F(r)(a) = \max\left( 0, \max_{a'\in S} r(a') - d(a, a')\frac{\bar{r}}{\eta} \right), \quad a \in \mathcal{A}.$$

Note that this function is $\frac{\bar{r}}{\eta}$−Lipschitz, hence uniformly-continuous—in the case where rewards are stochastic, we can still apply this transformation at the realization-level. Further, the sets $B_d(a', \eta)$ for $a' \in S$ are all disjoint by triangular inequality. Thus, for all $a' \in S$, we have $F(r)(a') = r(a')$. We now describe the reduction from uniformly-continuous rewards on $\mathcal{A}$ to unrestricted rewards on $S$. Let $\mathbb{X} \in \mathcal{C}(\mathcal{A})$ and we denote by $\hat{a}_t$ the action selected at time $t$ by an universally consistent learner $f_.$ under $\mathbb{X}$ for uniformly-continuous rewards on $\mathcal{A}$. We now construct a learning rule for unrestricted rewards on $S$. First, for $a \in \mathcal{A}$, denote by $NN_S(a) = \operatorname{argmin}_{a'\in S} d(a, a')$ the index of the nearest neighbor of $a$ in $S$ where ties are broken arbitrarily, e.g., by lexicographic order (necessarily, $S$ is countable because $\mathcal{A}$ is separable). We consider the learning rule which selects the actions $NN_S(\hat{a}_t)$, i.e.,

$$f_t^S(\boldsymbol{x}_{\leq t-1}, \boldsymbol{r}_{\leq t-1}, x_t) = NN_S(f_t(\boldsymbol{x}_{\leq t-1}, \boldsymbol{r}_{\leq t-1}, x_t))$$

for all $x_{\leq t} \in \mathcal{X}^t$ and $r_{\leq t-1} \in [0, \bar{r}]^{t-1}$. We aim to show that $f_.^S$ is universally consistent under $\mathbb{X}$ for unrestricted rewards on $S$. Fix any reward mechanism $r$ on the action space $S$. We consider the reward mechanism $\tilde{r}$ on the action space $\mathcal{A}$ as follows,

$$\tilde{r}_t(a, x) = F(r(\cdot \mid x))(a),$$

for any $a \in \mathcal{A}$. Note that the mechanism $\tilde{r}$ only depends on the nearest neighbor of selected actions. Denote $\tilde{a}_t$ the corresponding selected action. Observe that by construction of the functional $F$, for any $t \geq 1$, $\tilde{r}_t(\tilde{a}_t) \geq \tilde{r}_t(\hat{a}_t)$. Thus, by monotonicity, $f^S_.$ is also consistent on reward mechanism $\tilde{r}$. Now note that $\tilde{f}_.$ only selects actions within $S$ and receives the same rewards that would have been observed by running the learning rule on reward mechanism $r$. As a result, $f^S_.$ is also consistent for reward $r$. This ends the proof that it is universally consistent under $\mathbb{X}$ and hence $\mathbb{X} \in \mathcal{C}(S)$. This ends the proof of the proposition.

$\square$

As a direct consequence of Lemma 7.3 and the results from previous sections, we can use the necessary conditions from the unrestricted reward setting by changing the terms "finite action set" (resp. "countably infinite action set") into "totally-bounded action set" (resp. "non-totally-bounded action set").

COROLLARY 7.4.  *Let $\mathcal{A}$ be a non-totally-bounded metric space. Then, $\mathcal{C}^{uc} \subset \mathcal{C}_1$. Let $\mathcal{A}$ be a totally-bounded metric space with $|\mathcal{A}| > 2$. Then, $\mathcal{C}^{uc} \subset \mathcal{C}_2$.*

We now turn to sufficient conditions and show that we can recover the results from the unrestricted case as well. For non-totally-bounded value spaces, the EXPINF learning rule from Theorem 7.2 is already universally consistent under $\mathcal{C}_1$ processes, which is a necessary condition by Corollary 7.4. As a result, imposing the uniformly-continuous assumption on the rewards does not improve the set of learnable processes.

THEOREM 7.5.  *Let $\mathcal{X}$ be a separable Borel metrizable space and $\mathcal{A}$ a non-totally-bounded metric space. Then, $\mathcal{C}^{uc} = \mathcal{C}_1$.*

Next, we consider totally-bounded actions spaces and generalize the learning rule for stochastic rewards in finite action spaces. Recall that this learning rule associates to each time a category $p = \text{CATEGORY}(t)$, based on the number of previous occurrences of $X_t$, and works separately on each category. Within each category, the algorithm balances between two strategies: strategy 0 which uses independent EXP3 learners for each distinct instance, and strategy 1 which performs EXPINF. We adapt the algorithm in the following way. First, the EXP3 learners from strategy 0 search for the best action within $\mathcal{A}(\delta_p)$, an $\delta_p-$net of $\mathcal{A}$ where $\delta_p$ will be defined carefully. Note that since $\mathcal{A}$ is possibly infinite, restricting strategy 0 to finite action sets is necessary. However, we aim for arbitrary precision, hence we will have $\delta_p \to 0$ as $p \to \infty$. Second, for strategy 1, we use the countable set of functions $\Pi$ defined as for the EXPINF algorithm in Theorem 7.2.

THEOREM 7.6.  *Let $\mathcal{A}$ be a totally-bounded metric space. Then, there exists an optimistically universal learning rule for stationary and uniformly-continuous rewards, and learnable processes are $\mathcal{C}^{uc} = \mathcal{C}_2$.*

PROOF. We first define the new learning rule. CATEGORY and ASSIGNPURPOSE are left unchanged. We will use the countable set of policies $\Pi = \{\pi^l, l \geq 1\}$ as in the continuous case in Lemma 7.1, for $\text{EXPLORE}(1; \cdot)$, and Algorithm 5. Further, in $\text{EXPLORE}(0; \cdot)$ and Algorithm 5, $\text{EXP3}_{\mathcal{A}}$ is replaced by $\text{EXP3}_{\mathcal{A}(\delta_p)}$. Finally, in SELECTSTRATEGY, $\eta_p = 10\frac{\sqrt{|\mathcal{A}| \ln |\mathcal{A}|}}{2^{p/4}}$ is replaced by $\eta_p = 10\frac{\sqrt{|\mathcal{A}(\delta_p)| \ln |\mathcal{A}(\delta_p)|}}{2^{p/4}}$, where we will define $\delta_p$ shortly. In the original proof of the universal consistence of the algorithm, we showed that the average error of the learning rule on category $p$, $\mathcal{T}_p$ is $\mathcal{O}(\tilde{\epsilon}_p)$ where $\tilde{\epsilon}_p = 2\frac{\sqrt{|\mathcal{A}| \ln |\mathcal{A}|}}{2^{p/4}}$. Similarly, we now define $\epsilon_p := 2\frac{\sqrt{|\mathcal{A}(\delta_p)| \ln |\mathcal{A}(\delta_p)|}}{2^{p/4}}$. A key feature of the proof is that since we had $\sum_p \tilde{\epsilon}_p < \infty$, the

learner can afford to converge on each set $\mathcal{T}_p$ separately. We mimic this behavior by choosing $\delta_p$ such that $\sum_p \epsilon_p < \infty$. Precisely, we pose

$$\delta_p = \min\{2^{-i} : |\mathcal{A}(2^{-i})|\ln|\mathcal{A}(2^{-i})| \leq 2^{p/4}\}.$$

As a result, we obtain directly $\epsilon_p \leq 2^{-1-p/8}$ which is summable, and $\delta_p \to 0$ as $p \to \infty$.

We now show that this learning rule is universally consistent under processes $\mathbb{X} \in \mathcal{C}_2$ by adapting the proof of Theorem 4.2. Fix $r$ a reward mechanism. For every $\epsilon > 0$, there exists $\Delta(\epsilon)$ such that

$$\forall x \in \mathcal{X}, \forall a, a' \in \mathcal{A}, \quad d(a, a') \leq \Delta(\epsilon) \Rightarrow |\bar{r}(a, x) - \bar{r}(a', x)| \leq \epsilon.$$

For every $\delta > 0$, we will also define $\epsilon(\delta) = 2\inf\{\epsilon > 0 : \Delta(\epsilon) \geq \delta\}$. By uniform-continuity, $\epsilon(\delta) \to 0$ as $\delta \to 0$ and because of the factor 2, we have

$$\forall x \in \mathcal{X}, \forall a, a' \in \mathcal{A}, \quad d(a, a') \leq \delta \Rightarrow |\bar{r}(a, x) - \bar{r}(a', x)| \leq \epsilon(\delta).$$

Now observe that in the original proof, the probabilistic bounds $p_i(p, q)$ for $1 \leq i \leq 8$ do not depend on the cardinality of the action set. Therefore, on the same event $\mathcal{E} \cap \mathcal{F}$ of probability one, Eq (3), (4), (5), (6), (7) and (8) hold starting from some time $\hat{T}$, for the intended values of $p, q, T$. The only difference, however, is that in strategy 0, we perform EXP3 over the restricted action set $\mathcal{A}(\delta_p)$. As a result, for any $x \in \mathcal{X}$, we have

$$\max_{a \in \mathcal{A}(\delta_p)} \bar{r}(a, x) \geq \max_{a \in \mathcal{A}} \bar{r}(a, x) - \epsilon(\delta_p).$$

As a result, Eq (3) should be replaced with

$$\hat{R}_p^0(q) \geq \bar{R}_p^*(q) - (T_p^{q+1})^{\frac{7}{8}} - 6\frac{\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}}{2^{p/4}}(T_p^{q+1} - T_p^q) - \epsilon(\delta_p)|\mathcal{T}_p(q)|$$

$$\hat{R}_p^0(q) \leq \bar{R}_p^*(q) - (T_p^{q+1})^{\frac{7}{8}} - 6\frac{\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}}{2^{p/4}}(T_p^{q+1} - T_p^q).$$

Note that the additional term $\epsilon(\delta_p)|\mathcal{T}_p(q)|$ is not present in the upper bound because searching over $\mathcal{A}$ (in $\bar{R}_p^*(q)$) is always better than searching over $\mathcal{A}(\delta_p)$ (in $\hat{R}_p^0(q)$). Similarly, Eq (4) should be replaced with

$$\tilde{R}_p^0(q) \geq \bar{R}_p^*(q) - 6\frac{\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}}{2^{p/4}}(T_p^{q+1} - T_p^q) - (T_p^{q+1})^{3/4} - A_p(q) - \epsilon(\delta_p)|\mathcal{T}_p(q)|.$$

Similarly, the adapted Eq (9) becomes

$$\mathcal{R}_p(T) \geq \bar{R}_p^*(T) - \frac{1+c}{2}\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}T^{1-1/2^7}\log_2 T - \frac{T}{2^p} - \epsilon(\delta_p)|\mathcal{T}_p \cap \{t \leq T\}|.$$

Furthering the same bounds, Eq (10) becomes

$$\mathcal{R}_p(T) \geq \bar{R}_p^*(T) - (33 + 5c)\sqrt{|\mathcal{A}|\ln|\mathcal{A}|}T^{1-1/2^7}\log_2 T - 15\epsilon_p T - 2\epsilon(\delta_p)|\mathcal{T}_p \cap \{t \leq T\}|.$$

We are now ready to prove universal consistence of our learning rule. Fix $0 < \epsilon < 1$, and as in the original proof, let $p_0$ such that $\sum_{p \geq p_0} \epsilon_p < \frac{\epsilon}{15}$, because $\sum_p \epsilon_p < \infty$. Again, we have $\mathbb{X}^{\leq 4^{p_0}} \in \mathcal{C}_1$ and as a result, we can apply Lemma 7.1. As a result, on an event $\mathcal{H}$ of probability one, for all $\epsilon > 0$, there exists $i(\epsilon) \geq 1$ such that $2^{-i(\epsilon)} \leq \Delta(\epsilon), \epsilon$ and $\pi^{i(\epsilon)} \in \Pi$ such that

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t \leq T, t \in \mathcal{T}^{\leq 4^{p_0}}} \bar{r}_t(\pi^*(X_t)) - \bar{r}_t(\pi^i(X_t))$$

$$\le \limsup_{T\to\infty} \frac{1}{T} \sum_{t\le T, t\in\mathcal{T}^{\le 4^{p_0}}} \mathbb{1}[d(\pi^*(X_t), \pi^{i(\epsilon)}(X_t)) \ge 2^{-i(\epsilon)}]$$

$$+ \limsup_{T\to\infty} \frac{1}{T} \sum_{t\le T, t\in\mathcal{T}^{\le 4^{p_0}}} (\bar{r}_t(\pi^*(X_t)) - \bar{r}_t(\pi^{i(\epsilon)}(X_t))) \mathbb{1}[d(\pi^*(X_t), \pi^{i(\epsilon)}(X_t)) \le \Delta(\epsilon)]$$

$$\le 2^{-i(\epsilon)} + \epsilon \le 2\epsilon,$$

where $\pi^*$ denotes the optimal policy. We define the events $\mathcal{E}, \mathcal{F}$ as in the original proof. In the rest of the proof, we will now suppose that the event $\mathcal{E} \cap \mathcal{F} \cap \mathcal{H}$ of probability one is satisfied. On this event, because the parameter $\epsilon > 0$ was arbitrary in the above derivations, gthere exists $l_0 \ge 1$ (random index) such that

$$\limsup_{T\to\infty} \frac{1}{T} \sum_{t\le T, t\in\mathcal{T}^{\le 4^{p_0}}} \bar{r}_t(\pi^*(X_t)) - \bar{r}_t(\pi^{l_0}(X_t)) \le \frac{\epsilon}{2^{2p_0+2}}.$$

Following the same arguments as in the original proof, for $p < p_0$, and $T_p^q$ sufficiently large, we need to adapt the following estimates.

$$\max_{1\le l\le k(q)} \hat{R}_p^k(q) \ge \hat{R}_p^{l_0}(q)$$

$$\ge \bar{R}_p^{l_0}(q) - (T_p^{q+1})^{7/8}$$

$$\ge \bar{R}_p^*(q) - (T_p^{q+1})^{7/8} - \sum_{t\in\mathcal{T}_p(q)} (r_t(\pi^*(X_t)) - \bar{r}_t(\pi^{l_0}(X_t)))$$

$$\ge \hat{R}_p^0(q) - 2(T_p^{q+1})^{7/8} - 3\epsilon_p(T_p^{q+1} - T_p^q) - \sum_{t\in\mathcal{T}_p(q)} r_t(\pi^*(X_t)) - \bar{r}_t(\pi^{l_0}(X_t)).$$

Then, observe that

$$\limsup_{q\to\infty} \frac{2(T_p^{q+1})^{7/8} + 3\epsilon_p(T_p^{q+1} - T_p^q) + \sum_{t\in\mathcal{T}_p(q)} r_t(\pi^*(X_t)) - \bar{r}_t(\pi^{l_0}(X_t))}{T_p^{q+1} - T_p^q} \le 4\epsilon_p < \eta_p.$$

Thus, as in the original proof, starting from some time $\tilde{T}$, the learning rule always chooses strategy 1 over strategy 0 for all categories $p \le p_0$.

We continue the same arguments to obtain for $p < p_0$ and $T \ge 2^{p_0}\tilde{T}$,

$$\mathcal{R}_p(T) - \bar{R}_p^*(T) \ge -2^{p_0}\tilde{T} - 16(3+c)T^{15/16}\ln T - \sum_{t\le T, t\in\mathcal{T}_p} \bar{r}_t(\pi^*(X_t)) - \bar{r}_t(\pi^{l_0}(X_t)),$$

which yields

$$\sum_{p<p_0} \bar{R}_p^*(T) - \mathcal{R}_p(T) \le p_0 2^{p_0}\tilde{T} + 16p_0(3+c)T^{15/16}\ln T + \sum_{t\le T} \bar{r}_t(\pi^*(X_t)) - \bar{r}_t(\pi^{l_0}(X_t)).$$

Noting that $\limsup_{T\to\infty} \frac{1}{T} \sum_{t\le T} \bar{r}_t(\pi^*(X_t)) - \bar{r}_t(\pi^{l_0}(X_t)) \le \epsilon$, from there, the same arguments show that the learning rule is universally consistent.                                      □

As a summary, with the uniform-continuity assumption we could generalize all results from the unrestricted rewards case with the corresponding totally-bounded/non-totally-bounded dichotomy on action spaces.

**8. Unbounded rewards.** In this section, we allow for unbounded rewards $\mathcal{R} = [0, \infty)$ and start with the unrestricted rewards setting—no continuity assumption. Recall that in this setting, we assume that for any context $x \in \mathcal{X}$ and action $a \in \mathcal{A}$, the random variable $r(a, x)$ is integrable so that the immediate expected reward is well defined.

When $\mathcal{A}$ is uncountable, we showed that even for bounded rewards, no process $\mathbb{X}$ admits universal learning. Therefore, we will focus on the case when $\mathcal{A}$ is finite or countably infinite, and show that $\mathcal{C}_3$ determines whether universal consistency is possible. Moreover, a simple variant of EXPINF suffices for optimistically universal learning as follows. Enumerate $\mathcal{A} = \{a_1, a_2, \ldots, a_{|\mathcal{A}|}\}$ (or $\mathcal{A} = \{a_1, a_2, \ldots\}$ for countably infinite $\mathcal{A}$) and for any observed instance $x \in \mathcal{X}$, we run an independent EXPINF where the experts of the sequence are the constant policies equal to $a_i$ for $1 \leq i \leq |\mathcal{A}|$, i.e., the expert $E_i$ always selects action $a_i$.

THEOREM 8.1. *Let $\mathcal{A}$ be a countable action set with $|\mathcal{A}| \geq 2$. Then, there is an optimistically universal learning rule and the set of learnable processes admitting universal is $\mathcal{C}_3$.*

The fact that $\mathcal{C}_3$ characterizes universal learning was already the case in the noiseless full-feedback setting [6], hence Theorem 8.1 shows that for unrestricted rewards, we can achieve universal learning in the partial feedback setting without generalization cost.

PROOF. First, even in the full-information feedback setting, $\mathbb{X} \in \mathcal{C}_3$ is known to be necessary for universal consistency [6]. A fortiori in the bandit setting, this condition is still necessary $\mathcal{C} \subset \mathcal{C}_3$.

We now show that the learning rule defined above is universally consistent under $\mathcal{C}_3$ processes. For simplicity, we denote by $\hat{a}_t$ the action selected be the learning rule at time $t$. Fix $\mathbb{X} \in \mathcal{C}_3$ and define $S = \{x \in \mathcal{X} : \mathbb{X} \cap \{x\} \neq \emptyset\}$ the support of the process. By definition of $\mathcal{C}_3$, almost surely, $|S| < \infty$. We denote by $\mathcal{E}$ this event of probability one. Next, for any $x \in S$, we define $\mathcal{T}(x) = \{t : X_t = x\}$ and let $\tilde{S} = \{x \in S : |\mathcal{T}(x)| = \infty\}$ the set of points which are visited an infinite number of times. Recall that the learning rule performs an independent EXPINF subroutine on the times $\mathcal{T}(x)$ for all $x \in S$. As a result, by Corollary 3.5, for any $x \in \tilde{S}$, with probability one, for all $a \in \mathcal{A}$,

$$\limsup_{T \to \infty} \frac{1}{|\mathcal{T}(x) \cap \{t \leq T\}|} \sum_{t \in \mathcal{T}(x), t \leq T} r_t(a) - r_t(\hat{a}_t) \leq 0.$$

Now observe that $\tilde{S}$ is countable. Hence, by the union bound, on an event $\mathcal{F}$ of probability one, for all $x \in \tilde{S}$ and $a \in \mathcal{A}$, we have

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t \leq T, t \in \mathcal{T}(x)} r_t(a) - r_t(\hat{a}_t) \leq \limsup_{T \to \infty} \frac{1}{|\mathcal{T}(x) \cap \{t \leq T\}|} \sum_{t \leq T, t \in \mathcal{T}(x)} r_t(a) - r_t(\hat{a}_t) \leq 0.$$

In the rest of the proof, we suppose that $\mathcal{E} \cap \mathcal{F}$ is met. On $\mathcal{E}$, there exists $\hat{T} = 1 + \max\{t : X_t = x, x \in S \setminus \tilde{S}\}$ such that for any $T \geq \hat{T}$, we have $X_t \in \tilde{S}$. Then, for any policy $\pi^* : \mathcal{X} \to \mathcal{A}$, and $T \geq 1$, we have

$$\sum_{t=1}^{T} r_t(\pi^*(X_t)) - r_t(\hat{a}_t) \leq \sum_{t \leq \hat{T}} r_t(\pi^*(X_t)) + \sum_{x \in \tilde{S}} \sum_{t \leq T, t \in \mathcal{T}(x)} r_t(a) - r_t(\hat{a}_t).$$

As a result, because $\mathcal{F}$ is met,

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} r_t(\pi^*(X_t)) - r_t(\hat{a}_t) \leq \sum_{x \in \tilde{S}} \limsup_{T \to \infty} \frac{1}{T} \sum_{t \leq T, t \in \mathcal{T}(x)} r_t(a) - r_t(\hat{a}_t) \leq 0.$$

using the fact that $\mathbb{P}[\mathcal{E} \cap \mathcal{F}] = 1$, we proved that the learning rule is universally consistent under any $\mathcal{C}_3$ process. This ends the proof of the theorem. $\square$

The last remaining question is whether this very restrictive set of processes $\mathcal{C}_3$ can be improved under the continuity and uniform-continuity assumptions from Definition 2.2.

Unfortunately, we show that this is not the case for continuous rewards, however, the continuity assumption allows to achieve universal consistence on $\mathcal{C}_3$ processes even on uncountable action spaces. Recall that by Theorem 6.2, universal consistency was not achievable for uncountable spaces in the unrestricted reward case.

THEOREM 8.2. *Let $\mathcal{X}$ be a separable metrizable Borel space and $(\mathcal{A}, d)$ be a separable metric space with $|\mathcal{A}| \geq 2$. Then, there is an optimistically universal learning rule for continuous unbounded rewards and the set of learnable processes for universal learning with continuous unbounded rewards is $\mathcal{C}_3$.*

PROOF. In the case of countable action set $\mathcal{A}$ with $|\mathcal{A}| \geq 2$, Theorem 8.1 already showed that $\mathcal{C}_3$ is sufficient for universal learning under continuous unbounded rewards. Therefore, it remains to show that in the case of uncountable action space, $\mathcal{C}_3$ is still sufficient for universal learning. More precisely, we will show that the same learning rule which assigns a distinct EXPINF learner to each distinct instance of $\mathbb{X}$ as defined in Theorem 8.1 is still universally consistent under $\mathcal{C}_3$ processes. The only difference is that we run the learners EXPINF on a dense sequence of actions $(a_i)_{i \geq 1}$ of the complete action set $\mathcal{A}$ which may be uncountable. Let $\mathbb{X} \in \mathcal{C}_3$. We use the same notations as in the original proof of Theorem 8.1 for the support $S = \{x \in \mathcal{X} : \mathbb{X} \cap \{x\} \neq \emptyset\}$, the event $\mathcal{E} = \{|S| < \infty\}$, $\mathcal{T}(x) = \{t : X_t = x\}$ for $x \in S$ and $\tilde{S} = \{x \in S : |\mathcal{T}(x)| = \infty\}$. By Corollary 3.5, for any $x \in \tilde{S}$, with probability one, for all $i \geq 1$, we have now

$$\limsup_{T \to \infty} \frac{1}{|\mathcal{T}(x) \cap \{t \leq T\}|} \sum_{t \in \mathcal{T}(x), t \leq T} r_t(a_i) - r_t(\hat{a}_t) \leq 0.$$

Let $a \in \mathcal{A}$ and $\epsilon > 0$, because $(a_i)_{i \geq 1}$ is dense in $\mathcal{A}$ and the immediate reward is continuous, there exists $i(\epsilon)$ such that $|\bar{r}(a_{i(\epsilon)}) - \bar{r}(a)| \leq \epsilon$. Now observe that by the union bound, for any $x \in \tilde{S}$, with probability one, by the law of large numbers one has for all $i \geq 1$,

$$\frac{1}{|\mathcal{T}(x) \cap \{t \leq T\}|} \sum_{t \in \mathcal{T}(x), t \leq T} r_t(a_i) \xrightarrow[T \to \infty]{} \bar{r}_t(a_i),$$

and similarly for $a$. As a result, for any $x \in \tilde{S}$, with probability one, for any $\epsilon > 0$,

$$\limsup_{T \to \infty} \frac{1}{|\mathcal{T}(x) \cap \{t \leq T\}|} \sum_{t \in \mathcal{T}(x), t \leq T} r_t(a) - r_t(\hat{a}_t)$$

$$\leq \bar{r}(a) - \bar{r}(a_{i(\epsilon)}) + \limsup_{T \to \infty} \frac{1}{|\mathcal{T}(x) \cap \{t \leq T\}|} \sum_{t \in \mathcal{T}(x), t \leq T} r_t(a_{i(\epsilon)}) - r_t(\hat{a}_t)$$

$$\leq \epsilon.$$

As a result, we showed that for any $x \in \tilde{S}$, and any $a \in \mathcal{A}$, with probability one,

$$\limsup_{T \to \infty} \frac{1}{|\mathcal{T}(x) \cap \{t \leq T\}|} \sum_{t \in \mathcal{T}(x), t \leq T} r_t(a) - r_t(\hat{a}_t) \leq 0.$$

Now fix $\pi^* : \mathcal{X} \to \mathcal{A}$ a measurable policy. Because $\tilde{S}$ is countable, by the union bound, on an event $\mathcal{F}$ of probability one, for all $x \in \tilde{S}$, we have

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t \leq T, t \in \mathcal{T}(x)} r_t(\pi^*(x)) - r_t(\hat{a}_t)$$

$$\leq \limsup_{T \to \infty} \frac{1}{|\mathcal{T}(x) \cap \{t \leq T\}|} \sum_{t \leq T, t \in \mathcal{T}(x)} r_t(\pi^*(x)) - r_t(\hat{a}_t) \leq 0.$$

Then, the same arguments as in the original proof show that on $\mathcal{E} \cap \mathcal{F}$, for any $T \geq 1$, one has

$$\limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} r_t(\pi^*(X_t)) - r_t(\hat{a}_t) \leq \sum_{x \in \tilde{S}} \limsup_{T \to \infty} \frac{1}{T} \sum_{t \leq T, t \in \mathcal{T}(x)} r_t(a) - r_t(\hat{a}_t) \leq 0.$$

Thus, the learning rule is universally consistent under $\mathcal{C}_3$ processes.

We now show that $\mathbb{X} \in \mathcal{C}_3$ is still necessary for universal learning with continuous rewards. For the unrestricted reward case, this was a direct consequence of a result of [6], which we now adapt for continuous rewards. First, for any $\mathbb{X} \notin \mathcal{C}_3$, they show that there exists a disjoint measurable partition $\{B_i\}_{i=1}^{\infty}$ such that with non-zero probability, $|\{i : \mathbb{X} \cap B_i \neq \emptyset\}| = \infty$ on an event $\mathcal{E}_0$. Then, they constructed a sequence of times $T_i$ for $i \geq 1$ such that on an event $\mathcal{E}$ of probability one, for sufficiently large indices $i$, $\tau_i := \min\{0\} \cup \{t : X_t \in B_i\} \leq T_i$. Now fix two distinct actions $a_0, a_1 \in \mathcal{A}$, let $\epsilon = \frac{d(a_0, a_1)}{3}$ and fix a learning rule $f.$. We denote by $\hat{a}_t$ its selected action at time $t$. Consider the following rewards

$$(13) \qquad r^{\boldsymbol{U}}(a, x) = \max\left(0, T_i\left(1 - \frac{d(a, a_{U_j})}{\epsilon}\right)\right), \quad x \in B_j,$$

for any binary sequence $\boldsymbol{U}$. Now suppose that they were sampled from an i.i.d. sequence of Bernouillis $\mathcal{B}(\frac{1}{2})$, independent of the process $\mathbb{X}$ and the randomness of the learning rule. Now observe that for any $i \geq 1$ such that $\tau_i \leq T_i$, with probability at least $\frac{1}{2}$ independently of the past, we have $\hat{a}_{\tau_i} \notin B(a_{U_j}, \epsilon)$, which implies $\max_{a \in \mathcal{A}} r_{\tau_i}^{\boldsymbol{U}}(a) - r_{\tau_i}^{\boldsymbol{U}}(\hat{a}_{\tau_i}) \geq T_i$. From there, the same arguments as in the original proof show that with probability one, this event occurs infinitely often and $\mathcal{E}$ is met, which by the law of total probability implies that there exists a deterministic choice of values for $\boldsymbol{U} = (U_j)_{j \geq 1}$ such that on the corresponding deterministic (hence stationary) rewards, the learning rule is not consistent on $\mathcal{E}_0 \cap \mathcal{E}$ which has non-zero probability. This shows that $\mathbb{X}$ does not admit universal learning even in the simplest case of deterministic continuous rewards. $\square$

Last, we investigate the case of uniformly-continuous unrestricted rewards. Unfortunately, the uniform continuity assumption over the immediate expected rewards does not provide any advantage over the continuity assumption.

PROPOSITION 8.3. *Let $\mathcal{X}$ be a separable metrizable Borel space and $\mathcal{A}$ be a separable metric space with $|\mathcal{A}| \geq 2$. Then, the set of learnable processes for universal learning with uniformly-continuous unbounded rewards is $\mathcal{C}_3$.*

PROOF. It suffices to show that the $\mathcal{C}_3$ condition is still necessary for universal learning under uniformly-continuous rewards since the sufficiency is guaranteed by Theorem 8.2. We adapt the proof of the necessity of $\mathcal{C}_3$ in the continuous unbounded reward case. Let $\mathbb{X} \notin \mathcal{C}_3$ and suppose that there exists an universally consistent learning rule $f.$ under $\mathbb{X}$ for uniformly-continuous unbounded rewards. We use the same notations as in the proof of Theorem 8.2. We

now define a sequence $(M_i)_{i\geq 1}$ recursively such that $M_1 = 2T_1$ and for any $i \geq 1$, $M_{i+1} = 2T_{i+1} + 4T_{i+1}\sum_{j\leq i} M_j$. Then, consider the following stochastic rewards

$$r(a,x) = \begin{cases} M_i\left(1 + \frac{d(a,a_0)\wedge d(a_0,a_1)}{d(a_0,a_1)}\right) & \text{w.p. } \frac{1}{2}, \\ M_i\left(1 - \frac{d(a,a_0)\wedge d(a_0,a_1)}{d(a_0,a_1)}\right) & \text{w.p. } \frac{1}{2}. \end{cases} \quad x \in B_i, i \geq 1.$$

These rewards are uniformly-continuous because for any $x \in \mathcal{X}$, the expected immediate reward is $\bar{r}(a,x) = 0$ for all $a \in \mathcal{A}$. Now for $u \in \{0,1\}$, define the constant policy $\pi^u : x \in \mathcal{X} \mapsto a_u \in \mathcal{A}$. Denote by $\hat{a}_t$ the action selected by the learning rule at time $t$. Because it is consistent under the rewards mechanism given by $r$, using $\pi^0$, $\pi^1$ and the union bound, we have that almost surely, for any $u \in \{0,1\}$,

$$(14) \qquad \limsup_{T\to\infty} \frac{1}{T}\sum_{t=1}^{T} r_t(a_u, X_t) - r_t(\hat{a}_t, X_t) \leq 0.$$

Now recall that on the event $\mathcal{E}_0$ of non-zero probability, we have $|\{i : \mathbb{X} \cap B_i \neq \emptyset\}| = \infty$. In other terms, $|\{i : \tau_i > 0\}| = \infty$. We then define the random sequence of indices $(i_k)_{k\geq 1}$ such that on $\mathcal{E}_0^c$, $i_k = 0$ for all $k \geq 1$ and on $\mathcal{E}_0$, the indices are defined recursively such that $i_1 = \operatorname{argmin}_{i\geq 1, \tau_i>0} \tau_i$ and for $k \geq 1$, we have $i_{k+1} = \operatorname{argmin}_{i>i_k, \tau_i>0} \tau_i$. The argmin are well defined because on $\mathcal{E}_0$, all the times $\tau_i$ for $i \in \{j \geq 1 : \tau_j > 0\}$ are distinct. As a result, by construction of the recursion, on $\mathcal{E}_0$, the sequence $(i_k)_{k\geq 1}$ is an increasing sequence of times and for all $k \geq 1$, we have

$$\{i : \mathbb{X}_{<\tau_{i_k}} \cap B_i \neq \emptyset\} = \{i : 0 < \tau_i < \tau_{i_k}\} \subset \{1 \leq i < i_k\}.$$

Now recall that on the event $\mathcal{E}$ of probability one, there exists $\hat{i} \geq 1$ such that for any $i \geq \hat{i}$, we have $\tau_i := \min\{0\} \cup \{t : X_t \in B_i\} \leq T_i$. Therefore, on $\mathcal{E}_0 \cap \mathcal{E}$, letting $\hat{k} = \min\{k : i_k \geq \hat{i}\}$, we have that for $k \geq \hat{k}$, and $u \in \{0,1\}$

$$\sum_{t=1}^{\tau_{i_k}-1} r_t(a_u, X_t) - r_t(\hat{a}_t, X_t) \geq \sum_{i:\mathbb{X}_{<\tau_{i_k}}\cap B_i\neq\emptyset} \sum_{t<\tau_{i_k}, X_t\in B_i} (-2M_i)$$

$$\geq -2\sum_{i<i_k} T_{i_k} M_i$$

$$\geq -\frac{M_{i_k}}{2} + T_{i_k}.$$

Now observe that on the event $\mathcal{E}_0 \cap \mathcal{E}$ which has non-zero probability, if $d(\hat{a}_{\tau_{i_k}}, a_0) \geq \frac{d(a_0,a_1)}{2}$ and the reward on $B_{i_k}$ at time $\tau_{i_k}$ is in its negative alternative, i.e., $r(a,x) = M_i\left(1 - \frac{d(a,a_0)\wedge d(a_0,a_1)}{d(a_0,a_1)}\right)$, we have

$$\frac{1}{\tau_{i_k}}\sum_{t=1}^{\tau_{i_k}} r_t(a_0, X_t) - r_t(\hat{a}_t, X_t) \geq \frac{1}{\tau_{i_k}}\left(\frac{M_{i_k}}{2} - \frac{M_{i_k}}{2} + T_{i_k}\right) \geq 1.$$

Now by construction, the negative alternative occurs with probability $\frac{1}{2}$, independently from the past history and the complete process $\mathbb{X}$. As a result, for any $k \geq 1$, we have

$$(15) \quad \mathbb{P}\left[\frac{1}{\tau_{i_k}}\sum_{t=1}^{\tau_{i_k}} r_t(a_0, X_t) - r_t(\hat{a}_t, X_t) \geq 1 \mid \mathcal{E}_0, \mathcal{E}, k \geq \hat{k}, d(\hat{a}_{\tau_{i_k}}, a_0) \geq \frac{d(a_0,a_1)}{2}\right] \geq \frac{1}{2}.$$

Similarly, one can check that on the event $\mathcal{E}_0 \cap \mathcal{E}$, if $d(\hat{a}_{\tau_{i_k}}, a_0) < \frac{d(a_0, a_1)}{2}$ and the reward on $B_{i_k}$ at time $\tau_{i_k}$ is in its positive alternative, we have

$$\frac{1}{\tau_{i_k}} \sum_{t=1}^{\tau_{i_k}} r_t(a_1, X_t) - r_t(\hat{a}_t, X_t) \geq \frac{1}{\tau_{i_k}} \left( \frac{M_i}{2} - \frac{M_{i_k}}{2} + T_{i_k} \right) \geq 1.$$

As a result, the same arguments as above give

$$(16) \quad \mathbb{P}\left[ \frac{1}{\tau_{i_k}} \sum_{t=1}^{\tau_{i_k}} r_t(a_1, X_t) - r_t(\hat{a}_t, X_t) \geq 1 \mid \mathcal{E}_0, \mathcal{E}, k \geq \hat{k}, d(\hat{a}_{\tau_{i_k}}, a_0) < \frac{d(a_0, a_1)}{2} \right] \geq \frac{1}{2}.$$

Finally, define for any $T \geq 1$ the event

$$\mathcal{F}_T = \left\{ \frac{1}{T} \sum_{t=1}^{T} r_t(a_0, X_t) - r_t(\hat{a}_t, X_t) \geq 1 \right\} \cup \left\{ \frac{1}{T} \sum_{t=1}^{T} r_t(a_1, X_t) - r_t(\hat{a}_t, X_t) \geq 1 \right\}.$$

We obtain for any $k \geq 1$,

$$\mathbb{P}[\mathcal{F}_{\tau_{i_k}} \mid \mathcal{E}_0, \mathcal{E}, k \geq \hat{k}]$$

$$\geq \mathbb{P}\left[ \mathcal{F}_{\tau_{i_k}} \mid \mathcal{E}_0, \mathcal{E}, k \geq \hat{k}, d(\hat{a}_{\tau_{i_k}}, a_0) \geq \frac{d(a_0, a_1)}{2} \right] \mathbb{P}\left[ d(\hat{a}_{\tau_{i_k}}, a_0) \geq \frac{d(a_0, a_1)}{2} \mid \mathcal{E}_0, \mathcal{E}, k \geq \hat{k} \right]$$

$$+ \mathbb{P}\left[ \mathcal{F}_{\tau_{i_k}} \mid \mathcal{E}_0, \mathcal{E}, k \geq \hat{k}, d(\hat{a}_{\tau_{i_k}}, a_0) < \frac{d(a_0, a_1)}{2} \right] \mathbb{P}\left[ d(\hat{a}_{\tau_{i_k}}, a_0) < \frac{d(a_0, a_1)}{2} \mid \mathcal{E}_0, \mathcal{E}, k \geq \hat{k} \right]$$

$$\geq \frac{1}{2} \mathbb{P}\left[ d(\hat{a}_{\tau_{i_k}}, a_0) \geq \frac{d(a_0, a_1)}{2} \mid \mathcal{E}_0, \mathcal{E}, k \geq \hat{k} \right] + \frac{1}{2} \mathbb{P}\left[ d(\hat{a}_{\tau_{i_k}}, a_0) < \frac{d(a_0, a_1)}{2} \mid \mathcal{E}_0, \mathcal{E}, k \geq \hat{k} \right]$$

$$= \frac{1}{2},$$

where in the second inequality we used Eq (15) and Eq (16). As a result, using Fatou's lemma

$$\mathbb{P}[\mathcal{F}_{\tau_{i_k}} \text{ occurs for infinitely many } k \geq 1 \mid \mathcal{E}_0, \mathcal{E}] \geq \limsup_{k \geq 1} \mathbb{P}[\mathcal{F}_{\tau_{i_k}} \mid \mathcal{E}_0, \mathcal{E}]$$

$$\geq \frac{1}{2} \limsup_{k \geq 1} \mathbb{P}[k \geq \hat{k} \mid \mathcal{E}_0, \mathcal{E}] = \frac{1}{2},$$

where in the last inequality, we used the dominated convergence theorem given that on the event $\mathcal{E}$, $\hat{k} < \infty$. As a result, we showed that

$$\mathbb{P}\left[ \exists u \in \{0, 1\}, \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} r_t(a_u, X_t) - r_t(\hat{a}_t, X_t) \geq 1 \mid \mathcal{E}_0, \mathcal{E} \right] \geq \frac{1}{2}.$$

However, because $\mathbb{P}[\mathcal{E} \cap \mathcal{E}_0] = \mathbb{P}[\mathcal{E}_0] > 0$, Eq (14) shows that

$$\mathbb{P}\left[ \forall u \in \{0, 1\}, \limsup_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} r_t(a_u, X_t) - r_t(\hat{a}_t, X_t) \geq 1 \mid \mathcal{E}_0, \mathcal{E} \right] = 1,$$

which contradicts the previous inequality. This shows that the learning rule was not consistent under the rewards $(r_t)_t$, hence not universally consistent under $\mathbb{X}$. This shows that $\mathcal{C}_3$ is necessary for universal learning and completes the proof. $\square$

## REFERENCES

[1] AUER, P. and CHIANG, C.-K. (2016). An algorithm with nearly optimal pseudo-regret for both stochastic and adversarial bandits. In *Conference on Learning Theory* 116–120. PMLR.

[2] BEN-DAVID, S. and URNER, R. (2012). On the hardness of domain adaptation and the utility of unlabeled target samples. In *International Conference on Algorithmic Learning Theory* 139–153. Springer.

[3] BESBES, O., GUR, Y. and ZEEVI, A. (2014). Stochastic multi-armed-bandit problem with non-stationary rewards. *Advances in neural information processing systems* **27**.

[4] BLANCHARD, M. (2022). Universal online learning: An optimistically universal learning rule. In *Conference on Learning Theory* 479–495. PMLR.

[5] BLANCHARD, M. and COSSON, R. (2022). Universal Online Learning with Bounded Loss: Reduction to Binary Classification. In *Conference on Learning Theory* 479–495. PMLR.

[6] BLANCHARD, M., COSSON, R. and HANNEKE, S. (2022). Universal Online Learning with Unbounded Losses: Memory Is All You Need. In *International Conference on Algorithmic Learning Theory* 107–127. PMLR.

[7] BLANCHARD, M. and JAILLET, P. (2022). Universal Regression with Adversarial Responses. *arXiv preprint arXiv:2203.05067*.

[8] BUBECK, S., CESA-BIANCHI, N. et al. (2012). Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning* **5** 1–122.

[9] CÉROU, F. and GUYADER, A. (2006). Nearest neighbor classification in infinite dimension. *ESAIM: Probability and Statistics* **10** 340–355.

[10] CHEN, Y., LEE, C.-W., LUO, H. and WEI, C.-Y. (2019). A new algorithm for non-stationary contextual bandits: Efficient, optimal and parameter-free. In *Conference on Learning Theory* 696–726. PMLR.

[11] COHEN, D. T. and KONTOROVICH, A. (2022). Learning with Metric Losses. In *Proceedings of $35^{th}$ Conference on Learning Theory*.

[12] DEVROYE, L., GYÖRFI, L. and LUGOSI, G. (1996). *A Probabilistic Theory of Pattern Recognition*. Springer-Verlag New York.

[13] GOLDENSHLUGER, A. and ZEEVI, A. (2009). Woodroofe's one-armed bandit problem revisited. *The Annals of Applied Probability* **19** 1603–1633.

[14] GRAY, R. M. (2009). *Probability, Random Processes, and Ergodic Properties*, second ed. Springer.

[15] GRETTON, A., SMOLA, A., HUANG, J., SCHMITTFULL, M., BORGWARDT, K. and SCHÖLKOPF, B. (2009). Covariate shift by kernel mean matching. *Dataset shift in machine learning* **3** 5.

[16] GUAN, M. and JIANG, H. (2018). Nonparametric stochastic contextual bandits. In *Proceedings of the AAAI Conference on Artificial Intelligence* **32**.

[17] GYÖRFI, L., KOHLER, M., ZAK, A. K. and WALK, H. (2002). *A Distribution-Free Theory of Nonparametric Regression*. Springer-Verlag New York.

[18] GYÖRFI, L., LUGOSI, G. and MORVAI, G. (1999). A Simple Randomized Algorithm for Sequential Prediction of Ergodic Time Series. *IEEE Transactions on Information Theory* **45** 2642–2650.

[19] GYÖRFI, L. and WEISS, R. (2021). Universal consistency and rates of convergence of multiclass prototype algorithms in metric spaces. *Journal of Machine Learning Research* **22** 1–25.

[20] HANNEKE, S. (2021). Learning Whenever Learning Is Possible: Universal Learning under General Stochastic Processes. *Journal of Machine Learning Research* **22** 1–116.

[21] HANNEKE, S. (2021). Open Problem: Is There an Online Learning Algorithm That Learns Whenever Online Learning Is Possible? In *Proceedings of the $34^{th}$ Conference on Learning Theory*.

[22] HANNEKE, S. (2022). Universally Consistent Online Learning with Arbitrarily Dependent Responses. In *Proceedings of the $33^{rd}$ International Conference on Algorithmic Learning Theory*.

[23] HANNEKE, S., KONTOROVICH, A., SABATO, S. and WEISS, R. (2021). Universal Bayes Consistency in Metric Spaces. *The Annals of Statistics* **To appear**.

[24] LANGFORD, J. and ZHANG, T. (2007). The epoch-greedy algorithm for multi-armed bandits with side information. *Advances in neural information processing systems* **20**.

[25] LATTIMORE, T. and SZEPESVÁRI, C. (2020). *Bandit algorithms*. Cambridge University Press.

[26] LU, T., PÁL, D. and PÁL, M. (2009). Showing relevant ads via context multi-armed bandits. In *Proceedings of AISTATS*.

[27] LUO, H., WEI, C.-Y., AGARWAL, A. and LANGFORD, J. (2018). Efficient contextual bandits in non-stationary worlds. In *Conference On Learning Theory* 1739–1776. PMLR.

[28] MORVAI, G., KULKARNI, S. R. and NOBEL, A. B. (1999). Regression Estimation from an Individual Stable Sequence. *Statistics* **33** 99–118.

[29] MORVAI, G., YAKOWITZ, S. and GYÖRFI, L. (1996). Nonparametric Inference for Ergodic, Stationary Time Series. *The Annals of Statistics* **24** 370–379.

[30] NEU, G. (2015). Explore no more: Improved high-probability regret bounds for non-stochastic bandits. *Advances in Neural Information Processing Systems* **28**.

[31] PERCHET, V. and RIGOLLET, P. (2013). The multi-armed bandit problem with covariates. *The Annals of Statistics* **41** 693–721.

[32] RAKHLIN, A. and SRIDHARAN, K. (2016). Bistro: An efficient relaxation-based method for contextual bandits. In *International Conference on Machine Learning* 1977–1985. PMLR.

[33] REEVE, H., MELLOR, J. and BROWN, G. (2018). The k-nearest neighbour ucb algorithm for multi-armed bandits with covariates. In *Algorithmic Learning Theory* 725–752. PMLR.

[34] RIGOLLET, P. and ZEEVI, A. (2010). Nonparametric bandits with covariates. *arXiv preprint arXiv:1003.1630*.

[35] SARKAR, J. (1991). One-armed bandit problems with covariates. *The Annals of Statistics* 1978–2002.

[36] SLIVKINS, A. (2011). Contextual bandits with similarity information. In *Proceedings of the 24th annual Conference On Learning Theory* 679–702. JMLR Workshop and Conference Proceedings.

[37] SLIVKINS, A. et al. (2019). Introduction to multi-armed bandits. *Foundations and Trends® in Machine Learning* **12** 1–286.

[38] STEINWART, I., HUSH, D. and SCOVEL, C. (2009). Learning from Dependent Observations. *Journal of Multivariate Analysis* **100** 175–194.

[39] STONE, C. J. (1977). Consistent Nonparametric Regression. *The Annals of Statistics* **5** 595–620.

[40] SUGIYAMA, M., NAKAJIMA, S., KASHIMA, H., BUENAU, P. and KAWANABE, M. (2007). Direct importance estimation with model selection and its application to covariate shift adaptation. *Advances in neural information processing systems* **20**.

[41] SUK, J. and KPOTUFE, S. (2021). Self-Tuning Bandits over Unknown Covariate-Shifts. In *Proceedings of the 32nd International Conference on Algorithmic Learning Theory*.

[42] WANG, C.-C., KULKARNI, S. R. and POOR, H. V. (2005). Bandit problems with side observations. *IEEE Transactions on Automatic Control* **50** 338–355.

[43] WOODROOFE, M. (1979). A one-armed bandit problem with a concomitant variable. *Journal of the American Statistical Association* **74** 799–806.

[44] WU, Q., IYER, N. and WANG, H. (2018). Learning contextual bandits in a non-stationary environment. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval* 495–504.

[45] YANG, Y. and ZHU, D. (2002). Randomized allocation with nonparametric estimation for a multi-armed bandit problem with covariates. *The Annals of Statistics* **30** 100–121.