

Universal Regression with Adversarial Responses *

Moïse Blanchard
MIT
moiseb@mit.edu

Patrick Jaillet
MIT
jaillet@mit.edu

Abstract

We provide algorithms for regression with adversarial responses under large classes of non-i.i.d. instance sequences, on general separable metric spaces, with *provably minimal* assumptions. We also give characterizations of learnability in this regression context. We consider *universal consistency* which asks for strong consistency of a learner without restrictions on the value responses. Our analysis shows that such an objective is achievable for a significantly larger class of instance sequences than stationary processes, and unveils a fundamental dichotomy between value spaces: whether finite-horizon mean estimation is achievable or not. We further provide *optimistically universal* learning rules, i.e., such that if they fail to achieve universal consistency, any other algorithms will fail as well. For unbounded losses, we propose a mild integrability condition under which there exist algorithms for adversarial regression under large classes of non-i.i.d. instance sequences. In addition, our analysis also provides a learning rule for mean estimation in general metric spaces that is consistent under adversarial responses without any moment conditions on the sequence, a result of independent interest.

Keywords. Statistical learning theory, consistency, non-parametric estimation, generalization, stochastic processes, online learning, metric spaces

1 Introduction

1.1 Motivation and background

We study the classical statistical problem of metric-valued regression. Given an instance metric space $(\mathcal{X}, \rho_{\mathcal{X}})$ and a value metric space $(\mathcal{Y}, \rho_{\mathcal{Y}})$ with a loss ℓ , one observes instances in \mathcal{X} and aims to predict the corresponding values in \mathcal{Y} . The learning procedure follows an iterative process where successively, the learner is given an instance X_t and predicts the value Y_t based on the historical samples and the new instance. The learner’s goal is to minimize the loss of its predictions \hat{Y}_t compared to the true value Y_t . In particular, $\mathcal{Y} = \{0, 1\}$ (resp. $\mathcal{Y} = \{0, \dots, k\}$) with 0-1 loss corresponds to binary (resp. multiclass) classification while $\mathcal{Y} = \mathbb{R}$ corresponds to the classical regression setting. Motivated by the increase of new types of data in numerous data analysis applications— e.g., data lying on spherical spaces [Cha89; MJM00], manifolds [Shi+09; Dav+10; Fle13], Hilbert spaces [Zai+19], Hadamard spaces [LM21]—we will study the case where both instances and value spaces are general separable metric spaces. This general setting adopted in the recent literature on universal learning [Han21b; CK22; Bla22] includes and extends the specific classification and regression settings mentioned above. In this context, we model the stream of data as a general stochastic process $(\mathbb{X}, \mathbb{Y}) := (X_t, Y_t)_{t \geq 1}$, and are interested in *consistent* predictions that have vanishing average *excess* loss compared to any fixed measurable predictor functions $f : \mathcal{X} \rightarrow \mathcal{Y}$, i.e., $\frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \rightarrow 0$ (*a.s.*). Naturally, one would hope that the algorithm converges for a large class of value functions. Thus, we are interested in *universally consistent* learning rules that are consistent irrespective of the value process \mathbb{Y} .

*Accepted, Annals of Statistics, June 2023

The i.i.d. version of this problem where one assumes that the sequence (\mathbb{X}, \mathbb{Y}) is i.i.d. has been extensively studied. A classical result is that for binary classification in Euclidean spaces, k -nearest neighbor (kNN) with $k/\ln T \rightarrow \infty$ and $k/T \rightarrow 0$ is universally consistent under mild assumptions on the distribution of (X_1, Y_1) [Sto77; Dev+94; DGL13]. These results were then extended to a broader class of spaces [DGL13; Gyö+02] and more recently, [Han+21; GW21; CK22] provided universally consistent algorithms for any essentially separable metric space \mathcal{X} which are precisely those for which universal consistency is achievable for i.i.d. pairs $(X_t, Y_t)_{t \geq 1}$ of instances and responses. In parallel, a significant line of work aimed to obtain such results in non-i.i.d. settings, notably relaxations of the i.i.d. assumptions such as stationary ergodic processes [MYG96; GLM99; Gyö+02] or processes satisfying the law of large numbers [MKN99; GG09; SHS09].

1.2 Optimistic universal learning

In this work, we aim to understand which are the minimal assumptions on the data sequences for which universal consistency is still achievable. As such, we follow the *optimistic decision theory* [Han21a] which formalizes the paradigm of “learning whenever learning is possible”. Precisely, the *provably minimal* assumption for a given objective is that this task is achievable, or in other words that learning is possible. The goal then becomes to 1. characterize for which settings this objective is achievable and 2. if possible, provide learning rules that achieve this objective whenever it is achievable. These are called *optimistically universal* learning rules and enjoy the convenient property that if they failed the objective, any other algorithms would fail as well.

1.3 Related works in universal learning

This paradigm was recently used to study minimal assumptions for the noiseless (realizable) case where there exists an unknown underlying function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ such that $Y_t = f^*(X_t)$ [Han21a]. In this setting, the two questions described above were recently settled. For bounded losses, a simple variant of the nearest neighbor algorithm is optimistically universal [BC22; Bla22] and learnable processes are significantly larger than stationary processes. On the other hand, for unbounded losses, universal regression is extremely restrictive since the only learnable processes are those which visit a finite number of points almost surely [BCH22]. Yet, the general non-realizable setting was not characterized. As an initial result, for bounded losses, [Han22] proposed an algorithm that achieves universal consistency for a large class of processes \mathbb{X} , which intuitively asks that the sub-measure induced by empirical visits of the input sequence be continuous. There is however a significant gap between the proposed condition and the learnable processes in the bounded noiseless setting. [Han22] then left open the question of identifying the precise provably-minimal conditions to achieve consistency, and whether there exists an optimistically universal learning rule.

1.4 Adversarial responses and related works in learning with experts

The consistency results in [Han22] hold for arbitrary value processes \mathbb{Y} , arbitrarily correlated to the instance process \mathbb{X} . We consider the slightly more general *adversarial* responses and show that we can obtain the same results as for adversarial processes, without any generalizability cost. Formally, adversarial responses can not only arbitrarily depend on the instance sequence \mathbb{X} , but may also depend on past predictions and past randomness used by the learner. This is a non-trivial generalization for randomized algorithms—note that randomization is necessary to obtain guarantees for general online learning problems [BC12; Sli19]. There is a rich theory for arbitrary or adversarial responses \mathcal{Y} when the reference functions $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ are restricted to specific function classes \mathcal{F} . As a classical example, for the noiseless binary classification setting, there exist learning rules which guarantee a finite number of mistakes for arbitrary sequences \mathbb{X} , if and only if the class \mathcal{F} has finite Littlestone dimension [Lit88]. Other restrictions on the function class have been considered [CL06; BPS09; RST15]. Universal learning diverges from this line of work by imposing no restrictions on function classes—namely *all* measurable functions—but instead restricting instance processes \mathbb{X} to the optimistic set where universal consistency is achievable. Nevertheless, the algorithms we introduce

for adversarial responses use as subroutine the traditional exponentially weighted forecaster for learning with expert advice from the online learning literature, also known as the Hedge algorithm [LW94; Ces+97; FS97].

1.5 Contributions

In this paper, we provide answers to two fundamental questions in universal regression. First, we exactly characterize the set of processes we call *learnable*. These are instance processes \mathbb{X} for which universal learning is possible, i.e., consistency is achieved for every process $(X_t, Y_t)_{t \geq 1}$ with covariate sequence \mathbb{X} . Second, we provide optimistically universal learning rules, i.e., a unique algorithm that achieves universal consistency for all processes \mathbb{X} for which this is achievable by some learning rule. The specific answers to these questions depend on the value space and loss (\mathcal{Y}, ℓ) as detailed below.

1.5.1 Universal learning with empirically integrable responses

We introduce a mild moment-type assumption on the responses \mathbb{Y} , namely *empirical integrability*, that roughly asks that one can bound the tails of the empirical first moment of \mathbb{Y} . We then proceed to analyze the processes for which learning adversarial responses guaranteed to satisfy this assumption, is achievable. The answer depends on a property of the value space and loss (\mathcal{Y}, ℓ) which we denote F-TIME.

- If every ball $B_\ell(y, r)$ of (\mathcal{Y}, ℓ) satisfies the F-TIME property, the class of processes \mathbb{X} for which universal consistency under adversarial empirically integrable responses may be achieved is the so-called Sub-linear Measurable Visits (SMV) class. This coincides with the class of processes that admits universal learning for bounded losses in the realizable setting (noiseless responses) [Bla22]. In particular, this shows that for value spaces with bounded losses satisfying F-TIME, one can extend consistency results from the realizable setting to the adversarial one at no generalizability cost.
- Otherwise, the classes of processes \mathbb{X} for which one can achieve universal consistency for empirically integrable responses is a smaller class called Continuous Submeasure (CS). This is a condition that was already considered by [Han22], which showed that for bounded metric losses, one can achieve universal learning under CS processes. Our results show that whenever the F-TIME condition is not satisfied for bounded losses, CS is also a necessary condition for universal learning.

Also, in both cases, we give an optimistically universal learning rule, that is implicit for the first case—it uses as subroutine the learning rule for mean-estimation—and explicit for the second. These results resolve an open question from [Han22].

Intuitively, the property F-TIME asks that, for any fixed tolerance $\epsilon > 0$, there is a learning rule that solves the analogous prediction problem without covariates \mathbb{X} —*mean-estimation*—in finite time within the tolerance ϵ . This property is satisfied for “reasonable” value spaces, e.g., totally-bounded spaces or countably-many-classes classification (\mathbb{N}, ℓ_{01}) , but we also provide an explicit example of bounded metric space that does not satisfy this condition.

To motivate the introduction of the empirical integrability condition we show that a weaker moment-type assumption on responses—that $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) < \infty$ (*a.s.*) for some $y_0 \in \mathcal{Y}$ —is not sufficient to extend the results from the bounded loss case to unbounded losses, resolving an open question from [BCH22]. Further, empirical integrability is essentially necessary to obtain consistency results: it is automatically satisfied if the loss is bounded and for the i.i.d. setting it exactly asks that responses Y have finite first moment.

As a direct implication of this work, finite second moment $\mathbb{E}[Y^2]$ is sufficient to achieve consistency for stationary ergodic processes. This result relaxes the conditions of all past works to the best of our knowledge, which required finite fourth moment $\mathbb{E}[Y^4]$ [GO07].

1.5.2 Universal learning with unrestricted responses

For completeness, we also characterize the set of learnable processes without assuming empirical integrability on responses. Since the two notions coincide for bounded losses, we focus on unbounded losses. While there

always exists an optimistically universal learning rule, the precise class of universally learnable processes depends on an alternative involving the mean-estimation problem. Either mean-estimation on (\mathcal{Y}, ℓ) is impossible and universal learning is never achievable, or universal learning is achievable for processes that only visit a finite number of distinct points, a property called Finite Support (FS). Along the way, we show that mean-estimation with adversarial responses is always possible for metric losses, a result of independent interest.

1.6 Organization of the paper

After presenting the learning framework and definitions in Section 2, we describe in Section 3 our main results. Although these are stated for general value spaces under the empirical integrability constraint, the proofs build upon the bounded loss case. We follow this proof structure: in Section 4 we consider totally-bounded value spaces for which we can give explicit optimistically universal learning rules, in Section 5 we consider general bounded loss spaces. We then turn to unbounded and mean estimation in Section 6. Last, in Section 7 we introduce the empirical integrability and prove our general results for unbounded losses. We discuss open directions in Section 8.

2 Formal setup

We provide the necessary definitions, concepts and conditions.

2.1 Instance and value spaces

Consider a separable metric *instance space* $(\mathcal{X}, \rho_{\mathcal{X}})$ equipped with its Borel σ -algebra \mathcal{B} , and a separable metric *value space* $(\mathcal{Y}, \rho_{\mathcal{Y}})$ given with a loss ℓ . We recall that a metric space is *separable* if it contains a dense countable set. Unless mentioned otherwise, we suppose that the loss is a power of a metric, i.e., there exists $\alpha \geq 1$ such that the loss is $\ell = (\rho_{\mathcal{Y}})^{\alpha}$. As a remark, all of the results in this work can be generalized to *essentially separable* metric instance spaces, a condition introduced by [Han+21] which was shown to be the largest class of metric spaces for which learning possible. However, for the sake of exposition, we restrict ourselves to separable metric spaces. We denote $\bar{\ell} := \sup_{y_1, y_2 \in \mathcal{Y}} \ell(y_1, y_2)$. In the first Sections 4 and 5 of this work, we suppose that the loss ℓ is *bounded*, i.e., $\bar{\ell} < \infty$. The case of *unbounded* losses is addressed in the next Sections 6 and 7. We also introduce the notion of near-metrics for which we will provide some results. We say that ℓ is a near-metric on \mathcal{Y} if it is symmetric, satisfies $\ell(y, y) = 0$ for all $y \in \mathcal{Y}$, for any $y' \neq y \in \mathcal{Y}$ we have $\ell(y, y') > 0$, and it satisfies a relaxed triangle inequality $\ell(y_1, y_2) \leq c_{\ell}(\ell(y_1, y_3) + \ell(y_2, y_3))$ where c_{ℓ} is a finite constant.

2.2 Online learning on adversarial responses

We consider the *online learning* framework where at step $t \geq 1$, one observes a new instance $X_t \in \mathcal{X}$ and predicts a value $\hat{Y}_t \in \mathcal{Y}$ based on the past history $(X_u, Y_u)_{u \leq t-1}$ and the new instance X_t only. The learning rule may be randomized, where the private randomness used at each iteration t is drawn from a fixed probability space \mathcal{R} and independent of the data generation process used to generate Y_t .

Definition 2.1. An *online learning rule* is a sequence $f := \{f_t, R_t\}_{t \geq 1}$ of measurable functions $f_t : \mathcal{R} \times \mathcal{X}^{t-1} \times \mathcal{Y}^{t-1} \times \mathcal{X} \rightarrow \mathcal{Y}$ together with a distribution R_t on \mathcal{R} .

The prediction at time t of f is $f_t(r_t; (X_u)_{u \leq t-1}, (Y_u)_{u \leq t-1}, X_t)$ where $r_t \sim R_t$ is independent of the new value X_t and the past history $(X_u, Y_u)_{u \leq t}$. For simplicity, we may omit the internal randomness r_t and write directly $f_t : \mathcal{X}^{t-1} \times \mathcal{Y}^{t-1} \times \mathcal{X} \rightarrow \mathcal{Y}$. We are interested in general data-generating processes. To this means, a possible very general choice of instances and values are general stochastic processes $(\mathbb{X}, \mathbb{Y}) := \{(X_t, Y_t)\}_{t \geq 1}$ on the product space $\mathcal{X} \times \mathcal{Y}$. This corresponds to the arbitrarily dependent responses under instance processes \mathbb{X} [Han22]. In this work, we consider the slightly more general *adversarial responses* where the value Y_t is also allowed to depend on the past private randomness $(r_u)_{u \leq t-1}$ used by the learning rule f .

Definition 2.2. Let $\mathbb{X} = (X_t)_{t \geq 1}$ be a stochastic process on \mathcal{X} . An *adversarial response mechanism* on \mathbb{X} is a stochastic process $\{(\tilde{X}_t, \mathbf{Y}_t)\}_{t \geq 1}$ where $\tilde{X}_t \in \mathcal{X}$, $\mathbf{Y}_t = \mathbf{Y}_t(\cdot | \cdot)$ is a Markov kernel from \mathcal{R}^{t-1} to \mathcal{Y} , and $(\tilde{X}_t)_{t \geq 1}$ has same distribution as \mathbb{X} .

For a given learning rule f , having observed the sampled randomness $r_1, \dots, r_{t-1} \in \mathcal{R}$ used by the learning rule before time t , the target value at time t is $Y_t = \mathbf{Y}_t(r_1, \dots, r_{t-1})$. Again, for simplicity, we will refer to the adversarial response mechanism as \mathbb{Y} , which allows us to view the data generating process as a usual stochastic process on $\mathcal{X} \times \mathcal{Y}$. Of course, if the learning rule is *deterministic*, adversarial responses are equivalent to arbitrary dependent responses as in [Han22], but this is not necessarily the case for general *randomized* algorithms.

2.3 Empirically integrable responses

We introduce a novel assumption on the responses, namely *empirical integrability*.

Definition 2.3. A process $(Y_t)_{t \geq 1}$ is *empirically integrable* if there exists $y_0 \in \mathcal{Y}$ such that for any $\epsilon > 0$, almost surely there exists $M \geq 0$ for which

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M} \leq \epsilon.$$

Unless mentioned otherwise, we will focus on the case where responses satisfy this property. This is a mild assumption on the responses. Indeed, it is worth noting that this condition is always satisfied if the loss ℓ is bounded. Further, if for some $y_0 \in \mathcal{Y}$, $\ell(y_0, Y_t)$ admits moments of order $p > 1$, the empirical integrability condition is also satisfied.

2.4 Universal consistency

In this general setting, we are interested in online learning rules which achieve low long-run average loss compared to any fixed prediction function for general adversarial mechanisms. Given a learning rule f and an adversarial process (\mathbb{X}, \mathbb{Y}) , for any measurable function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$, we denote the long-run average excess loss as

$$\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f^*) := \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T (\ell(f_t(\mathbb{X}_{\leq t-1}, \mathbb{Y}_{\leq t-1}, X_t), Y_t) - \ell(f^*(X_t), Y_t)).$$

We can then define the notion of consistency which asks that the excess loss compared to any measurable function vanishes to zero.

Definition 2.4. Let (\mathbb{X}, \mathbb{Y}) be an adversarial process and f a learning rule. f is consistent under (\mathbb{X}, \mathbb{Y}) if for any measurable function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$, we have $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f^*) \leq 0$, (a.s.).

For example, if (\mathbb{X}, \mathbb{Y}) is an i.i.d. process on $\mathcal{X} \times \mathcal{Y}$ following a distribution μ where μ has a finite first-order moment, achieving consistency is equivalent to reaching the optimal risk $R^* := \inf_{f^*} \mathbb{E}_{(X, Y) \sim \mu} [\ell(f^*(X), Y)]$, where the infimum is taken over all measurable functions $f^* : \mathcal{X} \rightarrow \mathcal{Y}$. As introduced in [Han21a; Han22], consistency against all measurable function is the natural extension of consistency for i.i.d. processes (\mathbb{X}, \mathbb{Y}) to non-i.i.d. settings. The goal of universal learning is to design learning rules that are consistent for any adversarial process \mathbb{Y} that is empirically integrable.

Definition 2.5. Let \mathbb{X} be a stochastic process on \mathcal{X} and f a learning rule. f is *universally consistent* under \mathbb{X} for empirically integrable adversarial responses if for any adversarial process $(\tilde{\mathbb{X}}, \mathbb{Y})$ with $\tilde{\mathbb{X}} \sim \mathbb{X}$ and such that \mathbb{Y} is empirically integrable, f is consistent.

2.5 Optimistic universal learning

Given this regression setup, we define SOLAR (Strong universal Online Learning with Adversarial Responses) as the set of processes \mathbb{X} for which universal consistency with adversarial responses is *achievable*,

$$\text{SOLAR} = \{\mathbb{X} : \exists f. \text{ universally consistent learning rule under } \mathbb{X} \\ \text{for empirically integrable adversarial responses}\}.$$

Note that this learning rule is allowed to depend on the process \mathbb{X} . Similarly, in the realizable (noiseless) setting, one can define the set SOUL (Strong Online Universal Learning) of processes for which there exists a learning rule that is universally consistent for realizable responses when the loss is bounded (and hence, the empirical integrability condition is always satisfied). Of course, $\text{SOLAR} \subset \text{SOUL}$. We are then interested in learning rules that would achieve universal consistency whenever possible.

Definition 2.6. A learning rule f is *optimistically universal* for adversarial regression with empirically integrable responses if it is universally consistent under all $\mathbb{X} \in \text{SOLAR}$ for adversarial empirically integrable responses.

Similarly, we say that a learning rule is optimistically universal for noiseless regression if it is universally consistent under all $\mathbb{X} \in \text{SOUL}$ for noiseless responses when the loss is bounded. In this general framework, the main interests of optimistic learning are 1. identifying the set of learnable processes with adversarial responses SOLAR, 2. determining whether there exists an optimistically universal learning rule, and 3. constructing one if it exists.

Remark 2.7. Except for Section 6.1 in which we assume that the loss is a metric $\alpha = 1$, one can generalize our results to any symmetric and discernible losses ℓ such that for any $0 < \epsilon \leq 1$, there exists a constant c_ϵ such that for all $y_1, y_2, y_3 \in \mathcal{Y}$, $\ell(y_1, y_2) \leq (1 + \epsilon)\ell(y_1, y_3) + c_\epsilon\ell(y_2, y_3)$. Without loss of generality, we can further assume that c_ϵ is non-increasing in ϵ . This is a stronger assumption than having a near-metric ℓ , for which we also give some results in Sections 4 and 7.

3 Main results

We introduce some conditions on stochastic processes. For any process \mathbb{X} on \mathcal{X} , given any measurable set $A \in \mathcal{B}$ of \mathcal{X} , let $\hat{\mu}_{\mathbb{X}}(A) := \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{1}_A(X_t)$. We consider the condition CS (Continuous Sub-measure) defined as follows.

Condition CS: For every decreasing sequence $\{A_k\}_{k=1}^\infty$ of measurable sets in \mathcal{X} with $A_k \downarrow \emptyset$, $\mathbb{E}[\hat{\mu}_{\mathbb{X}}(A_k)] \xrightarrow[k \rightarrow \infty]{} 0$.

It is known that this condition is equivalent to $\mathbb{E}[\hat{\mu}_{\mathbb{X}}(\cdot)]$ being a continuous sub-measure [Han21a], hence the adopted name CS. Importantly, CS processes contain in particular i.i.d., stationary ergodic or stationary processes. We now introduce a second condition SMV (Sublinear Measurable Visits) which asks that for any partition, the process \mathbb{X} visits a sublinear number of sets of the partition.

Condition SMV: For every disjoint sequence $\{A_k\}_{k=1}^\infty$ of measurable sets of \mathcal{X} with $\bigcup_{k=1}^\infty A_k = \mathcal{X}$, (every countable measurable partition),

$$|\{k \geq 1 : A_k \cap \mathbb{X}_{\leq T} \neq \emptyset\}| = o(T), \quad (a.s.).$$

This condition is significantly weaker and allows to consider a larger family of processes $\text{CS} \subset \text{SMV}$, with $\text{CS} \subsetneq \text{SMV}$ whenever \mathcal{X} is infinite [Han21a]. Note that these sets depend on the instance space $(\mathcal{X}, \rho_{\mathcal{X}})$. This dependence is omitted for simplicity. We first consider bounded losses. In the *noiseless* case, where there exists some unknown measurable function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ such that the stochastic process \mathbb{Y} is given as $Y_t = f^*(X_t)$ for all $t \geq 1$, [Bla22] showed that learnable processes are exactly $\text{SOUL} = \text{SMV}$ for bounded

losses. [Bla22] also introduced a learning rule 2-Capped-1-Nearest-Neighbor (2C1NN), variant of the classical 1NN algorithm, which is *optimistically universal* in the noiseless case for bounded losses. Interestingly, we show that this same learning rule is universally consistent for unbounded losses in the noiseless setting with empirically integrable responses.

Theorem 3.1. *Let (\mathcal{Y}, ℓ) be a separable near-metric space. Then, 2C1NN is optimistically universal in the noiseless setting with empirically integrable responses, i.e., for all processes $\mathbb{X} \in \text{SMV}$ and for all measurable target functions $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ such that $(f^*(X_t))_{t \geq 1}$ is empirically integrable, $\mathcal{L}_{(\mathbb{X}, (f^*(X_t))_{t \geq 1})}(\text{2C1NN}, f^*) = 0$ (a.s.).*

In general, one has $\text{SOLAR} \subset \text{SMV}$. It was posed as a question whether we could recover the complete set SMV for learning under adversarial—or arbitrary—processes [Han22].

Question [Han22]: For bounded losses, does there exist an online learning rule that is universally consistent for arbitrary responses under all processes $\mathbb{X} \in \text{SMV}$ (= SOUL)?

We answer this question with an alternative. Depending on the bounded value space (\mathcal{Y}, ℓ) , either $\text{SOLAR} = \text{SMV}$ or $\text{SOLAR} = \text{CS}$, but in both cases there exists an optimistically universal learning rule. We now introduce the property F-TiME (Finite-Time Mean Estimation) on the value space (\mathcal{Y}, ℓ) which characterizes this alternative.

Property F-TiME: *For any $\eta > 0$, there exists a horizon time $T_\eta \geq 1$, an online learning rule $g_{\leq T_\eta}$ such that for any $\mathbf{y} := (y_t)_{t=1}^{T_\eta}$ of values in \mathcal{Y} and any value $y \in \mathcal{Y}$, we have*

$$\frac{1}{T_\eta} \mathbb{E} \left[\sum_{t=1}^{T_\eta} \ell(g_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t) \right] \leq \eta.$$

We are now ready to state our main results for bounded value spaces. The first result shows that if the value space satisfies the above property locally, we can universally learn all the processes in SOUL even under adversarial responses.

Theorem 3.2. *Suppose that any ball of (\mathcal{Y}, ℓ) , $B_\ell(y, r)$ satisfies F-TiME. Then, $\text{SOLAR} = \text{SMV}$ and there exists an optimistically universal learning rule f for adversarial regression with empirically integrable responses., i.e., such that for any stochastic process (\mathbb{X}, \mathbb{Y}) on $\mathcal{X} \times \mathcal{Y}$ with $\mathbb{X} \in \text{SMV}$ and \mathbb{Y} empirically integrable, for any measurable function $f : \mathcal{X} \rightarrow \mathcal{Y}$ we have $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f^*) \leq 0$, (a.s.).*

F-TiME defines a non-trivial alternative, and an explicit construction of a non-F-TiME bounded metric space $(\mathcal{Y}, \rho_{\mathcal{Y}})$ is given in Section 5.1 with $\mathcal{Y} = \mathbb{N}$. Nevertheless, F-TiME is satisfied by a large class of spaces, e.g., any totally-bounded metric space and countable classification $(\mathcal{Y}, \ell) = (\mathbb{N}, \ell_{01})$ satisfy F-TiME. Hence, we can universally learn all SOUL processes with adversarial responses, for countable classification (the empirical integrability condition is automatically satisfied because the loss is bounded). If F-TiME is not satisfied locally, we have the following result which shows that learning under CS is still possible but universal learning beyond CS processes cannot be achieved.

Theorem 3.3. *Suppose that there exists a ball $B_\ell(y, r)$ of (\mathcal{Y}, ℓ) that does not satisfy F-TiME. Then, $\text{SOLAR} = \text{CS}$ and there exists an optimistically universal learning rule f for adversarial regression with empirically integrable responses., i.e., such that for any stochastic process (\mathbb{X}, \mathbb{Y}) on $(\mathcal{X}, \mathcal{Y})$ with $\mathbb{X} \in \text{CS}$ and \mathbb{Y} empirically integrable, then, for any measurable function $f : \mathcal{X} \rightarrow \mathcal{Y}$ we have $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f^*) \leq 0$, (a.s.).*

For metric losses $\ell = \rho_{\mathcal{Y}}$, it was already known [Han22] that universal learning under adversarial responses under all processes in CS is achievable by some learning rule. Hence, Theorem 3.3 implies that this learning rule is automatically optimistically universal for adversarial regression for all metric value spaces with bounded loss which do not satisfy F-TiME. However, our result is stronger in that consistency holds for any power of a metric loss $\ell = \rho_{\mathcal{Y}}^\alpha$, $\alpha \geq 1$ and unbounded value spaces.

Remark 3.4. As a direct consequence of Theorems 3.2 and 3.3, for stationary ergodic processes, finite second moment of the values $\mathbb{E}[Y^2] < \infty$ suffices for consistency, in agreement with the known results for the i.i.d. setting. This relaxes the fourth-moment conditions $\mathbb{E}[Y^4] < \infty$ proposed in the literature [GO07].

We now consider removing the empirical integrability assumption. As mentioned above, for bounded losses this assumption is automatically satisfied, hence Theorems 3.2 and 3.3 apply directly, with a simplified alternative: whether (\mathcal{Y}, ℓ) satisfies F-TIME.

Corollary 3.5. *Suppose that ℓ is bounded.*

- *If (\mathcal{Y}, ℓ) satisfies F-TIME. Then, $\text{SOLAR} = \text{SMV}(= \text{SOUL})$.*
- *If (\mathcal{Y}, ℓ) does not satisfy F-TIME. Then, $\text{SOLAR} = \text{CS}$.*

Further, an optimistically universal learning rule for adversarial regression always exists, i.e., achieving universal consistency with adversarial responses under any $\mathbb{X} \in \text{SOLAR}$.

It remains to analyze the case of unbounded losses without empirical integrability assumption on the responses. To avoid confusions, we denote by SOLAR-U the set of processes that admit universal learning with adversarial (unrestricted) responses. Unfortunately, even in the noiseless setting, universal learning is extremely restrictive in that case. Specifically, the set of universally learnable processes SOUL for noiseless responses is reduced to the set FS (Finite Support) of processes that visit a finite number of different points almost surely [BCH22].

Condition FS: The process \mathbb{X} satisfies $|\{x \in \mathcal{X} : \{x\} \cap \mathbb{X} \neq \emptyset\}| < \infty$ (a.s.).

We show that in the adversarial setting we still have $\text{SOLAR-U} = \text{FS}$ when ℓ is a metric: we can solve the fundamental problem of mean estimation where one sequentially makes predictions of a sequence \mathbb{Y} of values in (\mathcal{Y}, ℓ) and aims to have a better long-run average loss than any fixed value. If responses \mathbb{Y} are i.i.d. this is the Fréchet means estimation problem [EJ20; Sch22; Jaf22]. Our main result on mean estimation holds in general spaces and is of independent interest.

Theorem 3.6. *Let (\mathcal{Y}, ℓ) be a separable metric space. There exists an online learning rule f that is universally consistent for adversarial mean estimation, i.e., for any adversarial process \mathbb{Y} on \mathcal{Y} , almost surely, for all $y \in \mathcal{Y}$,*

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T (\ell(f_t(\mathbb{Y}_{\leq t-1}), Y_t) - \ell(y, Y_t)) \leq 0.$$

Further, we show that for powers of metric we may have $\text{SOLAR-U} = \emptyset$. Specifically, for real-valued regression with Euclidean norm and loss $|\cdot|^\alpha$ and $\alpha > 1$, adversarial regression or mean estimation are not achievable. We then show that we have an alternative: either mean estimation with adversarial responses is achievable, $\text{SOLAR-U} = \text{FS}$ and we have an optimistically universal learning rule; or mean estimation is not achievable and $\text{SOLAR-U} = \emptyset$. Thus, even in the best case scenario for unbounded losses, $\text{SOLAR-U} = \text{FS}$, which is already extremely restrictive. [BCH22] asked whether imposing moment conditions on the responses would allow recovering the large set SMV as learnable processes instead. Specifically, they formulated the following question.

Question [BCH22]: For unbounded losses ℓ , does there exist an online learning rule f which is consistent under every $\mathbb{X} \in \text{SMV}$, for every measurable function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ such that there exists $y_0 \in \mathcal{Y}$ with $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, f^*(X_t)) < \infty$ (a.s.), i.e., such that we have $\mathcal{L}_{\mathbb{X}}(f, f^*) = 0$ (a.s.)?

We answer negatively to this question. Under this first-moment condition, universal learning under all SMV processes is not achievable even in this noiseless case. We show the stronger statement that noiseless universal learning under all processes having pointwise convergent relative frequencies—which are included

Table 1: Characterization of learnable instance processes in universal consistency (ME = Mean Estimation).

Learning setting	Bounded loss	Unbounded loss	Unbounded loss with empirically integrable responses
Noiseless responses	SOUL = SMV [Bla22]	SOUL = FS [BCH22]	Identical to bounded loss [This paper]
Adversarial (or arbitrary) responses	$\text{SOLAR} \supset \text{CS}$ (metric loss) [Han22] Does (\mathcal{Y}, ℓ) satisfy F-TIME? $\left\{ \begin{array}{l} \text{Yes} \quad \text{SOLAR} = \text{SMV} \\ \text{No} \quad \text{SOLAR} = \text{CS} \end{array} \right.$ [This paper]	Is ME achievable? $\left\{ \begin{array}{l} \text{Yes} \quad \text{SOLAR-U} = \text{FS} \\ \text{No} \quad \text{SOLAR-U} = \emptyset \end{array} \right.$ [This paper]	Identical to bounded loss [This paper]

Table 2: Proposed learning rules for universal consistency (ME = Mean Estimation and EI = Empirical Integrability).²

Learning setting	Loss (and response/setting constraints)	Learning rule	Guarantees for which processes \mathbb{X} ?	Optimist. universal?	Reference
I.i.d. responses	Finite or countable class., 01-loss	OptiNet	i.i.d.	No	[Han+21]
	Real-valued regression + integrable	Proto-NN	i.i.d.	No	[GW21]
	Metric loss + integrable	MedNet	i.i.d.	No	[CK22]
Noiseless responses (realizable)	Bounded loss	2C1NN	SMV	Yes	[Bla22]
	Unbounded loss	Memorization	FS	Yes	[BCH22]
	Unbounded + EI	2C1NN	SMV	Yes	[This paper]
Adversarial (or arbitrary) responses	Bounded loss + metric loss	Hedge-variant	CS	Not always	[Han22]
	Bounded loss + F-TIME	$(1 + \delta)$ C1NN-hedged	SMV	Yes	[This paper]
	Bounded loss + not F-TIME	Hedge-variant 2	CS	Yes	[This paper]
	Unbounded loss + ME	ME-algorithm	FS	Yes	[This paper]
	Unbounded loss + not ME	N/A	\emptyset	N/A	[This paper]
	Unbounded + EI + local F-TIME	EI- $(1 + \delta)$ C1NN-hedged	SMV	Yes	[This paper]
Unbounded + EI + not local F-TIME	EI-Hedge-variant	CS	Yes	[This paper]	

in CS—is not achievable. However, under the empirical integrability condition introduced above we are able to recover all positive results from bounded losses.

Tables 1 and 2 summarize known results in the literature and our contributions. As a reminder, $\text{FS} \subset \text{CS} \subset \text{SMV}$ in general, and $\text{FS} \subsetneq \text{CS} \subsetneq \text{SMV}$ whenever \mathcal{X} is infinite [Han21a].

4 An optimistically universal learning rule for totally-bounded value spaces

We start our analysis of universal learning under adversarial responses with *totally-bounded* value spaces, for which we can give simple and explicit algorithms. Hence, we suppose in this section that the value space (\mathcal{Y}, ℓ) is totally-bounded, i.e., for any $\epsilon > 0$ there exists a finite ϵ -net \mathcal{Y}_ϵ of \mathcal{Y} such that for any $y \in \mathcal{Y}$, there exists $y' \in \mathcal{Y}_\epsilon$ with $\ell(y, y') < \epsilon$. In particular, a totally-bounded space is necessarily bounded and separable. The goal of this section is to show that for such value spaces, adversarial universal regression is achievable for all processes in SMV as in the noiseless setting (the empirical integrability assumption is automatically satisfied in this context). Further, we explicitly construct an optimistically universal learning rule for adversarial responses.

We recall that in the noiseless setting, the 2C1NN learning rule achieves universal consistency for all

²In our paper, an algorithm is optimistically universal if it is universally consistent for all processes under which universal learning is possible in the considered setting. OptiNet, Proto-NN, and MedNet are optimistically universal in another sense, their guarantees hold in all metric spaces for which universal learning with i.i.d. pairs of instances and responses is achievable: *essentially separable* spaces $(\mathcal{X}, \rho_{\mathcal{X}})$ [Han+21]. Our learning rules also enjoy this second optimistic property.

SMV processes [Bla22]. At each iteration t , This rule performs the nearest neighbor rule over a restricted dataset instead of the complete history $\mathbb{X}_{\leq t-1}$. The dataset is updated by keeping track of the number of times each point X_u was used as nearest neighbor. This number is then capped at 2 by deleting from the current dataset any point which has been used twice as representative. Unfortunately, this learning rule is not optimistically universal for adversarial responses. More generally, [CK22] noted that any learning rule which only outputs observed historical values cannot be consistent, even in the simplest case of $\mathcal{X} = \{0\}$ and i.i.d. responses \mathbb{Y} . For instance, take $\mathcal{Y} = \bar{B}(0, 1)$ the closed ball of radius 1 in the plane \mathbb{R}^2 with the euclidean loss, consider the points $A, B, C \in \mathcal{Y}$ representing the equilateral triangle $e^{2ik\pi/3}$ for $k = 0, 1, 2$, and let \mathbb{Y} be an i.i.d. process following the distribution which visits A, B or C with probability $\frac{1}{3}$. Predictions within observed values, i.e., A, B or C , incur an average loss of $\frac{2}{3}\sqrt{3} > 1$ where 1 is the loss obtained with the fixed value $(0, 0)$.

To construct an optimistically universal learning rule for adversarial responses, we first generalize a result from [Bla22]. Instead of the 2C1NN learning rule, we use $(1 + \delta)$ C1NN rules for $\delta > 0$ arbitrarily small. Similarly as in 2C1NN, each new input X_t is associated to a representative $\phi(t)$ used for the prediction $\hat{Y}_t = Y_{\phi(t)}$. In the $(1 + \delta)$ C1NN rule, each point is used as a representative at most twice with probability δ and at most once with probability $1 - \delta$. In order to have this behavior irrespective of the process \mathbb{X} , which can be thought of been chosen by a (limited) adversary within the SOUL processes, the information of whether a point can allow for 1 or 2 children is only revealed when necessary. Specifically, at any step $t \geq 1$, the algorithm initiates a search for a representative $\phi(t)$. It successively tries to use the nearest neighbor of X_t within the current dataset and uses it as a representative if allowed by the maximum number of children that this point can have. However, the information whether a potential representative u can have at most 1 or 2 children is revealed only when u already has one child.

- If u allows for 2 children, it will be used as final representative $\phi(t)$.
- Otherwise, u is deleted from the dataset and the search for a representative continues.

The rule is formally described in Algorithm 1, where $\bar{y} \in \mathcal{Y}$ is an arbitrary value, and the maximum number of children that a point X_t can have is represented by $1 + U_t$. In this formulation, all Bernoulli $\mathcal{B}(\delta)$ samples are drawn independently of the past history. Note that if $\delta = 1$, the $(1 + \delta)$ C1NN learning rule coincides with the 2C1NN rule of [Bla22].

Theorem 4.1. *Fix $\delta > 0$. For any separable Borel space $(\mathcal{X}, \mathcal{B})$ and any separable near-metric output setting (\mathcal{Y}, ℓ) with bounded loss, in the noiseless setting, $(1 + \delta)$ C1NN is optimistically universal.*

We now construct our algorithm. This learning rule uses a collection of algorithms f^ϵ which each yield an asymptotic error at most a constant factor from $\epsilon^{\frac{1}{\alpha+1}}$. Now fix $\epsilon > 0$ and let \mathcal{Y}_ϵ be a finite ϵ -net of \mathcal{Y} for ℓ . Recall that we denote by $\bar{\ell}$ the supremum loss. We pose

$$T_\epsilon := \left\lceil \frac{\bar{\ell}^2 \ln |\mathcal{Y}_\epsilon|}{2\epsilon^2} \right\rceil \quad \text{and} \quad \delta_\epsilon := \frac{\epsilon}{2T_\epsilon}.$$

The quantity T_ϵ will be the horizon window used by our learning rule to make its prediction using the $(1 + \delta_\epsilon)$ C1NN learning rule. Precisely, let ϕ be the representative function from the $(1 + \delta_\epsilon)$ C1NN learning rule. Note that this representative function $\phi(t)$ is defined only for times t where a new instance X_t is revealed, otherwise $(1 + \delta_\epsilon)$ C1NN uses simple memorization $\hat{Y}_t = Y_u$. For simplicity, we will denote by $\mathcal{N} = \{t : \forall u < t, X_u \neq X_t\}$ these times of new instances. For $t \in \mathcal{N}$, we denote by $d(t)$ the depth of time t within the graph constructed by $(1 + \delta_\epsilon)$ C1NN, and define the horizon $L_t = d(t) \bmod T_\epsilon$. Intuitively, the learning rule f^ϵ performs the classical Hedge algorithm [CL06] on clusters of times that are close within the graph ϕ . Precisely, we define the equivalence relation between times as follows:

$$t_1 \stackrel{\phi}{\sim} t_2 \iff \begin{cases} \phi^{L_{u_1}}(u_1) = \phi^{L_{u_2}}(u_2) & \text{and } |\{u < t_i : X_u = X_{t_i}\}| \leq \frac{T_\epsilon}{\epsilon}, i = 1, 2 \\ \text{or} \\ X_{t_1} = X_{t_2} & \text{and } |\{u < t_i : X_u = X_{t_i}\}| > \frac{T_\epsilon}{\epsilon}, i = 1, 2, \end{cases}$$

Input: Historical samples $(X_t, Y_t)_{t < T}$ and new input point X_T
Output: Predictions $\hat{Y}_t = (1 + \delta)C1NN_t(\mathbf{X}_{<t}, \mathbf{Y}_{<t}, X_t)$ for $t \leq T$

```

 $\hat{Y}_1 := \bar{y}$  // Arbitrary prediction at  $t = 1$ 
 $\mathcal{D}_2 \leftarrow \{1\}; n_1 \leftarrow 0;$  // Initialisation
for  $t = 2, \dots, T$  do
  if exists  $u < t$  such that  $X_u = X_t$  then
     $\hat{Y}_t := Y_u$ 
  else
     $continue \leftarrow True$  // Begin search for available representative  $\phi(t)$ 
    while  $continue$  do
       $\phi(t) \leftarrow \min \{l \in \arg \min_{u \in \mathcal{D}_t} \rho_{\mathcal{X}}(X_t, X_u)\}$ 
      if  $n_{\phi(t)} = 0$  then // Candidate representative has no children
         $\mathcal{D}_{t+1} \leftarrow \mathcal{D}_t \cup \{t\}$ 
         $continue \leftarrow False$ 
      else // Candidate representative has one child
         $U_{\phi(t)} \sim \mathcal{B}(\delta)$ 
        if  $U_{\phi(t)} = 0$  then
           $\mathcal{D}_t \leftarrow \mathcal{D}_t \setminus \{\phi(t)\}$ 
        else
           $\mathcal{D}_{t+1} \leftarrow (\mathcal{D}_t \setminus \{\phi(t)\}) \cup \{t\}$ 
           $continue \leftarrow False$ 
        end
      end
    end
     $\hat{Y}_t := Y_{\phi(t)}$ 
     $n_{\phi(t)} \leftarrow n_{\phi(t)} + 1$ 
     $n_t \leftarrow 0$ 
  end

```

Algorithm 1: The $(1 + \delta)C1NN$ learning rule

where $u_i = \min\{u : X_u = X_{t_i}\}$ is the first occurrence of the considered instance point X_{t_i} . Hence, multiple occurrences of the same instance value fall in the same cluster and for new instance points times $t \in \mathcal{N}$, all times of a given cluster share the same ancestor up to generation at most $T_\epsilon - 1$. Additionally, a cluster is dedicated to instance points that have a significant number of duplicates. To make its prediction at time t , f^ϵ performs the Hedge algorithm based on values observed on its current cluster $\{u \leq t : u \stackrel{\phi}{\sim} t\}$. Let $\eta_\epsilon := \sqrt{\frac{8 \ln |\mathcal{Y}_\epsilon|}{\epsilon^2 T_\epsilon}}$ and define the losses $L_y^t = \sum_{u < t: u \stackrel{\phi}{\sim} t} \ell(Y_u, y)$. The learning rule $f_t^\epsilon(\mathbb{X}_{\leq t-1}, \mathbb{Y}_{\leq t-1}, X_t)$ outputs a random value in \mathcal{Y}_ϵ independently from the past history with

$$\mathbb{P}(\hat{Y}_t(\epsilon) = y) = \frac{e^{-\eta_\epsilon L_y^t}}{\sum_{z \in \mathcal{Y}_\epsilon} e^{-\eta_\epsilon L_z^t}}, \quad y \in \mathcal{Y}_\epsilon,$$

where, for simplicity, we denoted $\hat{Y}_t(\epsilon)$ the prediction given by the learning rule f^ϵ at time t .

Having constructed the learning rules f^ϵ , we are now ready to define our final learning rule f_\cdot . Let $\epsilon_i = 2^{-i}$ for all $i \geq 0$. Intuitively, it aims to select the best prediction within the rules f^{ϵ_i} . If there were a finite number of such predictors, we could directly use the algorithms for learning with experts from the literature [CL06]. Instead, we introduce these predictors one at a time: at step $t \geq 1$ we only consider the indices $I_t := \{i \leq \ln t\}$. We then compute an estimate $\hat{L}_{t-1, i}$ of the loss incurred by each predictor f^{ϵ_i} for $i \in I_t$ and select a random index \hat{i}_t independent from the past history from an exponentially-weighted distribution based on the estimates $\hat{L}_{t-1, i}$. The final output of our learning rule is $\hat{Y}_t := \hat{Y}_t(\epsilon_{\hat{i}_t})$. The complete algorithm is formally described in Algorithm 3. The following lemma quantifies the loss of the rule f_\cdot compared to the best rule f^{ϵ_i} .

Input: Historical samples $(X_t, Y_t)_{t < T}$ and new input point X_T ,
Representatives $\phi_\epsilon(\cdot)$ and depths $d_\epsilon(\cdot)$ constructed iteratively within $(1 + \delta_\epsilon)\text{C1NN}$.

Output: Predictions $\hat{Y}_t(\epsilon) = f_t^\epsilon(\mathbf{X}_{<t}, \mathbf{Y}_{<t}, X_t)$ for $t \leq T$

\mathcal{Y}_ϵ an ϵ -net of \mathcal{Y}

$T_\epsilon := \left\lceil \frac{\bar{\ell}^2 \ln |\mathcal{Y}_\epsilon|}{2\epsilon^2} \right\rceil$, $\eta_\epsilon := \sqrt{\frac{8 \ln |\mathcal{Y}_\epsilon|}{\ell^2 T_\epsilon}}$

for $t = 1, \dots, T$ **do**

$L_y^t = \sum_{u < t: u \in \mathcal{Y}_\epsilon} \ell(Y_u, y), \quad y \in \mathcal{Y}_\epsilon$	// Losses on the cluster given by ϕ_ϵ
$p^t(y) = \frac{\exp(-\eta_\epsilon L_y^t)}{\sum_{z \in \mathcal{Y}_\epsilon} \exp(-\eta_\epsilon L_z^t)}, \quad y \in \mathcal{Y}_\epsilon$	
$\hat{Y}_t \sim p^t$	

end

Algorithm 2: The f^ϵ learning rule

Input: Historical samples $(X_t, Y_t)_{t < T}$ and new input point X_T ,
Predictions $\hat{Y}_t(\epsilon_i)$ from the learning rules $f^{:\epsilon_i}$.

Output: Predictions \hat{Y}_t for $t \leq T$

$w_{0,0} = 1, t_i := \lceil e^i \rceil, \quad i \geq 0$

$I_t = \{i \leq \ln t\}, \eta_t = \sqrt{\frac{\ln t}{t}}, \quad t \geq 1$

for $t = 1, \dots, T$ **do**

$L_{t-1,i} := \sum_{s=t_i}^{t-1} \ell(\hat{Y}_s(\epsilon_i), Y_s), \hat{L}_{t-1,i} := \sum_{s=t_i}^{t-1} \hat{\ell}_s, \quad i \in I_t$	
$w_{t-1,i} = e^{\eta_t(\hat{L}_{t-1,i} - L_{t-1,i})}$	
$p_t(i) = \frac{w_{t-1,i}}{\sum_{j \in I_t} w_{t-1,j}}$	
$\hat{i}_t \sim p_t(\cdot)$	// model selection
$\hat{Y}_t = \hat{Y}_t(\epsilon_{\hat{i}_t})$	
$\hat{\ell}_t := \frac{\sum_{i \in I_t} w_{t-1,i} \ell(\hat{Y}_t(\epsilon_i), Y_t)}{\sum_{i \in I_t} w_{t-1,i}}$	

end

Algorithm 3: An optimistically universal learning rule for totally bounded spaces

Lemma 4.2. *Almost surely, there exists $\hat{t} \geq 0$ such that*

$$\forall t \geq \hat{t}, \forall i \in I_t, \quad \sum_{s=t_i}^t \ell(\hat{Y}_t, Y_t) \leq \sum_{s=t_i}^t \ell(\hat{Y}_t(\epsilon_i), Y_t) + (2 + \bar{\ell} + \bar{\ell}^2) \sqrt{t \ln t}.$$

We are now ready to show that Algorithm 3 is universally consistent under SMV processes.

Theorem 4.3. *Suppose that (\mathcal{Y}, ℓ) is totally-bounded. There exists an online learning rule f which is universally consistent for adversarial responses under any process $\mathbb{X} \in \text{SMV}(= \text{SOUL})$, i.e., for any process (\mathbb{X}, \mathbb{Y}) on $(\mathcal{X}, \mathcal{Y})$ with adversarial response, such that $\mathbb{X} \in \text{SMV}$, then for any measurable function $f : \mathcal{X} \rightarrow \mathcal{Y}$, we have $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f) \leq 0$, (a.s.).*

Proof sketch. First observe that Lemma 4.2 allows us to combine predictors f^ϵ : if individually they perform well, Algorithm 3 achieves the best long-term average excess loss among them. We then proceed to show that f^ϵ has low average error in the long run. First, $(1 + \delta_\epsilon)\text{C1NN}$ is universally consistent on SMV processes in the noiseless setting by Theorem 4.1. This intuitively shows that for noiseless functions, the value at time $\phi_\epsilon(t)$ provides a good representative for the value at time t . Extrapolating this argument, we show that if two times are close (for the graph metric) within the graph formed by ϕ_ϵ , they will have close

values for any fixed function in the long run. As a result, times in the same cluster defined by ϕ^ϵ share similar values in the long run. The f^ϵ rule precisely aims to learn the best predictor by cluster using the classical Hedge algorithm. Because it can only ensure low regret compared to a finite number of options, we use ϵ -nets of the value space \mathcal{Y} . The reason why we need to have $(1 + \delta_\epsilon)$ C1NN instead of the known 2C1NN algorithm is that for a given time T , we need to ensure low excess error even though some clusters might not be completed. Because the tree formed by ϕ_ϵ resembles a $(1 + \delta_\epsilon)$ -branching process, the fraction of times which belong to unfinished clusters is only a small fraction ϵT of the T times, hence does not affect the average long-term excess error significantly. Altogether, we show that f^ϵ has $\mathcal{O}(\epsilon^{\frac{1}{\alpha+1}})$ long-term average excess error compared to any fixed function for any SMV process, which ends the proof.

As a result, $\text{SMV} \subset \text{SOLAR}$ for totally-bounded value spaces. Recalling that for bounded values $\text{SMV} = \text{SOUL}$ [Bla22], i.e., processes $\mathbb{X} \notin \text{SMV}$ are not universally learnable even in the noiseless setting, we have $\text{SOLAR} \subset \text{SMV}$. Thus we obtain a complete characterization of the processes which admit universal learning with adversarial responses: $\text{SOLAR} = \text{SMV}$. Further, the proposed learning rule is optimistically universal for adversarial regression.

Corollary 4.4. *Suppose that (\mathcal{Y}, ℓ) is totally-bounded. Then, $\text{SOLAR} = \text{SMV}$, and there exists an optimistically universal learning rule for adversarial regression, i.e., which achieves universal consistency with adversarial responses under any process $\mathbb{X} \in \text{SOLAR}$.*

This is a first step towards the more general Theorem 5.5. Indeed, one can note that F-TIME is satisfied by any totally-bounded value space: given a fixed error tolerance $\eta > 0$, consider a finite $\frac{\eta}{2}$ -net $\mathcal{Y}_{\eta/2}$ of \mathcal{Y} . Because this is a finite set, we can perform the classical Hedge algorithm [CL06] to have $\Theta(\sqrt{T \ln |\mathcal{Y}_{\eta/2}|})$ regret compared to the best fixed value of $\mathcal{Y}_{\eta/2}$. For example, if $\alpha = 1$, posing $T_\eta = \Theta(\frac{4}{\eta^2} \ln |\mathcal{Y}_{\eta/2}|)$ enables to have a regret at most $\frac{\eta}{2} T_\eta$ compared to any fixed value of $\mathcal{Y}_{\eta/2}$, hence regret at most ηT_η compared to any value of \mathcal{Y} . This achieves F-TIME, taking a deterministic time $\tau_\eta := T_\eta$.

5 Characterization of learnable processes for bounded losses

While Section 4 focused on totally-bounded value spaces, the goal of this section is to give a full characterization of the set SOLAR of processes for which adversarial regression is achievable and provide optimistically universal algorithms, for any *bounded* value space.

5.1 Negative result for non-totally-bounded spaces

Although for all bounded value spaces (\mathcal{Y}, ℓ) , noiseless universal learning is achievable on all $\text{SMV} (= \text{SOUL})$ processes, this is not the case for adversarial regression in non-totally-bounded spaces. We show in this section that extending Corollary 4.4 to any bounded value space is impossible: the set of learnable processes for adversarial regression may be reduced to CS only, instead of SMV.

Theorem 5.1. *Let $(\mathcal{X}, \mathcal{B})$ a separable Borel metrizable space. There exists a separable metric value space (\mathcal{Y}, ℓ) with bounded loss such that the following holds: for any process $\mathbb{X} \notin \text{CS}$, universal learning under \mathbb{X} for arbitrary responses is not achievable. Precisely, for any learning rule f_\cdot , there exists a process \mathbb{Y} on \mathcal{Y} , a measurable function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ and $\epsilon > 0$ such that with non-zero probability $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f_\cdot, f^*) \geq \epsilon$.*

In the proof, we explicitly construct a bounded metric space that does not satisfy F-TIME. More precisely, we choose $\mathcal{Y} = \mathbb{N} = \{i \geq 0\}$ and a specific metric loss ℓ with values in $\{0, \frac{1}{2}, 1\}$. For any $k \geq 1$, we pose $n_k := 2k(k-1) + 2^k - 1$ and define the sets

$$I_k := \{n_k, n_k + 1, \dots, n_k + 4k - 1\} \quad \text{and} \quad J_k := \{n_k + 4k, n_k + 4k + 1, \dots, n_{k+1} - 1\}.$$

These sets are constructed so that $|I_k| = 4k$, $|J_k| = 2^k$ for all $k \geq 1$, and together with $\{0\}$, they form a partition of \mathbb{N} . We now construct the loss ℓ . We pose $\ell(i, j) = \mathbb{1}_{i=j}$ for all $i, j \in \mathbb{N}$ unless there is $k \geq 1$

such that $(i, j) \in I_k \times J_k$ or $(j, i) \in I_k \times J_k$. It now remains to define the loss $\ell(i, j)$ for all $i \in I_k$ and $j \in J_k$. Note that for any $j \in J_k$, we have that $j - n_k - 4k \in \{0, \dots, 2^k - 1\}$. Hence we will use their binary representation which we write as $j - n_k - 4k = \{b_j^{k-1} \dots b_j^1 b_j^0\}_2 = \sum_{u=0}^{k-1} b_j^u 2^u$ where $b_j^0, b_j^1, \dots, b_j^{k-1} \in \{0, 1\}$ are binary digits. Finally, we pose

$$\begin{aligned} \ell(n_k + 4u, j) &= \ell(n_k + 4u + 1, j) = \frac{1 + b_j^u}{2}, \\ \ell(n_k + 4u + 2, j) &= \ell(n_k + 4u + 3, j) = \frac{2 - b_j^u}{2}, \end{aligned}$$

for all $u \in \{0, 1, \dots, k-1\}$ and $j \in J_k$.

Proof sketch. This value space does not belong to F-TIME because for any algorithm and horizon time k , there is a sequence of length k of elements in I_k with $y_u = n_k + 4(u-1) + 2b_u + c_u$ for $1 \leq u \leq k$ and $b_u, c_u \in \{0, 1\}$, such that the algorithm incurs an average excess loss $\frac{1}{4}$ per iteration compared to some fixed element of J_k . To find such a sequence, we sample randomly and independently Bernoulli variables $b_u, c_u \sim \mathcal{B}(\frac{1}{2})$. In hindsight, the best predictor of the sequence is $n_k + 4k + j$, where $j = b_1 \dots b_k$ in binary representation. However, the algorithm only observes these bits in an online fashion: at time t it incurs an excess loss cost if it guesses an element of I_k because it has probability at most $\frac{1}{4}$ of finding y_t . And if it predicts an element of J_k , it cannot know in advance the correct t -th bit to choose in their binary representation.

We then proceed to show that for this space $\text{SOLAR} = \text{CS} \subsetneq \text{SOUL}$. To do so, we show that for processes $\mathbb{X} \notin \text{CS}$ there exists a sequence of disjoint measurable sets $\{B_p\}_{p \geq 1}$ and increasing times $(t_p)_{p \geq 1}$ and $\epsilon > 0$ such that with non-zero probability,

$$\forall p \geq 1, \quad \mathbb{X}_{\leq t_{p-1}} \cap B_p = \emptyset \text{ and } \exists t_{p-1} < t \leq t_p : \frac{1}{t} \sum_{t'=1}^t \mathbb{1}_{B_p}(X_{t'}) \geq \epsilon.$$

On this event, an online algorithm does not receive any information for instances in B_p before time t_{p-1} . We then construct responses by $(t_{p-1}, t_p]$. During this period and for contexts in B_p , we choose the same difficult-to-predict sequence of values as above for $k = t_p - t_{p-1}$. On the other hand, because the sets B_p are disjoint, there exists a measurable function f^* that selects the best action in hindsight for each set B_p . Intuitively, within horizon t_p , the algorithm cannot gather enough information to achieve lower average excess error than $\frac{\epsilon}{4}$ compared to f^* , which shows that it is not universally consistent.

Although learning beyond CS is impossible in this case, there still exists an optimistically universal learning rule for adversarial responses. Indeed, the main result of [Han22] shows that for any bounded value space, there exists a learning rule which is consistent under all CS processes for arbitrary responses (when ℓ is a metric, i.e., $\alpha = 1$).

Theorem 5.2 ([Han22]). *Suppose that (\mathcal{Y}, ℓ) is metric and ℓ is bounded. Then, there exists an online learning rule f which is universally consistent for arbitrary responses under any process $\mathbb{X} \in \text{CS}$, i.e., such that for any stochastic process (\mathbb{X}, \mathbb{Y}) on $(\mathcal{X}, \mathcal{Y})$ with $\mathbb{X} \in \text{CS}$, then for any measurable function $f : \mathcal{X} \rightarrow \mathcal{Y}$, we have $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f) \leq 0$, (a.s.).*

The proof of this theorem given in [Han22] extends to adversarial responses. However, we defer the argument because we will later prove Theorem 3.3 which also holds for any loss $\ell = \rho_{\mathcal{Y}}^\alpha$ for $\alpha \geq 1$ and unbounded losses in Section 7. This shows that for any separable metric space $(\mathcal{X}, \rho_{\mathcal{X}})$, there exists a metric value space for which the learning rule proposed in [Han22] was already optimistically universal.

5.2 Adversarial regression for classification with a countable number of classes

Although we showed in the last section that adversarial regression under all SMV processes is not achievable for some non-totally-bounded spaces, we will show that there exist non-totally-bounded value spaces for which

we can recover SOLAR = SMV. Precisely, we consider the case of classification with countable number of classes (\mathbb{N}, ℓ_{01}) , with 0 – 1 loss $\ell_{01}(i, j) = \mathbb{1}_{i \neq j}$. The goal of this section is to prove that in this case, we can learn arbitrary responses under any SOUL process. The main difficulty with non-totally-bounded classification is that we cannot apply traditional online learning tools because ϵ -nets may be infinite. Hence, we first show a result that allows us to perform online learning with an infinite number of experts in the context of countable classification.

Lemma 5.3. *Let $t_0 \geq 1$. There exists an online learning rule f such that for any sequence $\mathbf{y} := (y_i)_{i \geq 1}^T$ of values in \mathbb{N} , we have that for $T \geq t_0$*

$$\sum_{t=1}^T \mathbb{E}[\ell_{01}(f_t(\mathbf{y}_{\leq t-1}), y_t)] \leq \min_{y \in \mathbb{N}} \sum_{t=1}^T \ell_{01}(y, y_t) + 1 + \ln 2 \sqrt{\frac{t_0}{2 \ln t_0}} + \sqrt{\frac{\ln t_0}{2t_0}}(t_0 + T),$$

and with probability $1 - \delta$,

$$\sum_{t=1}^T \mathbb{E}[\mathbb{1}_{f_t(\mathbf{y}_{\leq t-1})=y_t}] \geq \max_{y \in \mathbb{N}} \sum_{t=1}^T \mathbb{1}_{y=y_t} - 1 - \ln 2 \sqrt{\frac{t_0}{2 \ln t_0}} - \sqrt{\frac{\ln t_0}{2t_0}}(t_0 + T) - \sqrt{2T \ln \frac{1}{\delta}}.$$

Proof sketch. We adapt the classical Hedge algorithm, which in its standard form can only ensure sublinear regret compared to a fixed set of values. Instead, we only consider a small subset of candidate values that is enlarged occasionally with previously observed values $y \in \mathbb{Y}_{\leq t}$. This formalizes the intuition that even though there are a priori an infinite number of candidate values (\mathbb{N}), it is reasonable to only focus on values with high frequency in the observed sequence $\mathbb{Y}_{\leq t}$: if the next value y_{t+1} is not in this set, the algorithm incurs a loss 1, which would also be incurred by the best fixed predictor until time $t+1$ in hindsight.

We can therefore adapt the learning rules f^ϵ from Section 4 by replacing the Hedge algorithm with the algorithm from Lemma 5.3. Further adapting parameters, we obtain our main result for countable classification.

Theorem 5.4. *Let $(\mathcal{X}, \mathcal{B})$ be a separable Borel metrizable space. There exists an online learning rule f which is universally consistent for adversarial responses under any process $\mathbb{X} \in \text{SMV}$ for countable classification, i.e., such that for any adversarial process (\mathbb{X}, \mathbb{Y}) on $(\mathcal{X}, \mathbb{N})$ with $\mathbb{X} \in \text{SMV}$, for any measurable function $f^* : \mathcal{X} \rightarrow \mathbb{N}$, we have that $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f^*) \leq 0$, (a.s.).*

5.3 A complete characterization of universal regression on bounded spaces

The last two Sections 5.1 and 5.2 gave examples of non-totally-bounded value spaces for which we obtain respectively SOLAR = CS or SOLAR = SMV. In this section, we prove that there is an underlying alternative, defined by F-TIME, which enables us to precisely characterize the set SOLAR of learnable processes for adversarial regression.

When F-TIME is satisfied by the value space, similarly to the case of countable classification, we recover SOLAR = SMV and there exists an optimistically universal rule. The corresponding algorithm follows the same general structure as the learning rule provided in Section 4 for totally-bounded-spaces, however, the learning rules f^ϵ need to be significantly modified. First, the Hedge algorithm should be replaced by the learning rule $g_{\leq t_\epsilon}$ provided by the F-TIME property. Second, as the horizon time t_ϵ of this learning rule is bounded, the clusters of points on which it is applied have to be adapted: we cannot simply use clusters by distance in the graph defined by the $(1 + \delta_\epsilon)\text{C1NN}$ algorithm. Instead, we construct clusters of smaller size t_ϵ among these larger graph-based clusters.

More precisely, we take the horizon time t_ϵ and the learning rule $g_{\leq t_\epsilon}^\epsilon$ satisfying the condition imposed by the assumption on (\mathcal{Y}, ℓ) . Then, let $T_\epsilon = \lceil \frac{t_\epsilon}{\epsilon} \rceil$. Similarly as before, we then define $\delta_\epsilon := \frac{\epsilon}{2T_\epsilon}$ and let ϕ be the representative function from the $(1 + \delta_\epsilon)\text{C1NN}$ learning rule. Then, we introduce the same equivalence relation between times $\overset{\phi}{\sim}$, which induces clusters of times. We define a sequence of i.i.d. copies $g^{\epsilon, t}$ of the

Input: Historical samples $(X_t, Y_t)_{t < T}$ and new input point X_T ,
Learning rule for finite-time mean estimation $g_{\leq t_\epsilon}^\epsilon$, $T_\epsilon = \lceil \frac{t_\epsilon}{\epsilon} \rceil$, $\delta_\epsilon := \frac{\epsilon}{2T_\epsilon}$.
Representatives $\phi_\epsilon(\cdot)$ constructed iteratively within $(1 + \delta_\epsilon)$ C1NN.
Output: Predictions $\hat{Y}_t(\epsilon) = f_t^\epsilon(\mathbf{X}_{<t}, \mathbf{Y}_{<t}, X_t)$ for $t \leq T$
for $t = 1, \dots, T$ **do**
 $\mathcal{C}(t) = \{u < t : u \stackrel{\phi_\epsilon}{\sim} t\}$
 if $\mathcal{C}(t) = \emptyset$ **then** $L_t = 0$ and initialize learner $g^{\epsilon, t}$;
 else
 $\psi(t) = \max \mathcal{C}(t)$
 if $L_{\psi(t)} < t_\epsilon - 1$ **then** $L_t = L_{\psi(t)} + 1$;
 else $L_t = 0$ and initialize learner $g^{\epsilon, t}$;
 end
 $\hat{Y}_t = g_{L_t+1}^{\epsilon, \psi^{L_t}(t)} \left(\{y_{\psi^{L_t+1-u}(t)}\}_{u=1}^{L_t} \right)$
end

Algorithm 4: The modified f^ϵ learning rule for value spaces (\mathcal{Y}, ℓ) satisfying F-TiME. When initializing a learner $g^{\epsilon, t}$ for finite-time mean estimation, its internal randomness is sampled independently from the past.

learning rule g^ϵ for all $t \geq 1$. This means that the randomness used within these learning rules is i.i.d, and the copy $g^{\epsilon, t}$ should be sampled only at time t , independently of the past history. Predictions are then made by blocks of size t_ϵ within the same cluster: at time t , let $u_1 < \dots < u_{L_t} < t$ be the elements of the current block. If the block does not contain t_ϵ elements yet, we use $g_{L_t+1}^{\epsilon, u_1}$ for the prediction at time t . Otherwise, we start a new block and use $g_1^{\epsilon, t}$. Hence, letting $\psi(t) = \max \mathcal{C}(t)$ be the last time in the same cluster as t (as defined by ϕ_ϵ) and L_t the size of the current block of t without counting t , we now define the learning rule f^ϵ such that for any sequence \mathbf{x}, \mathbf{y} ,

$$f_t^\epsilon(\mathbf{x}_{\leq t-1}, \mathbf{y}_{\leq t-1}, x_t) := g_{L_t+1}^{\epsilon, \psi^{L_t}(t)} \left(\{y_{\psi^{L_t+1-u}(t)}\}_{u=1}^{L_t} \right).$$

The complete learning rule is given in Algorithm 4. The learning rules f^ϵ are then combined into a single learning rule as in the original algorithm for totally-bounded spaces, following the same procedure given in Algorithm 3. We then show that it is universally consistent under SMV processes using same arguments as for Theorem 4.3.

Theorem 5.5. *Suppose that ℓ is bounded and (\mathcal{Y}, ℓ) satisfies F-TiME. Then, SOLAR = SMV(= SOUL) and there exists an optimistically universal learning rule for adversarial regression, i.e., which achieves universal consistency with adversarial responses under any process $\mathbb{X} \in \text{SMV}$.*

We are now interested in value spaces (\mathcal{Y}, ℓ) which do not satisfy F-TiME. We will show that in this case, SOLAR is reduced to the processes CS. We first introduce a second property on value spaces as follows.

Property 2: *For any $\eta > 0$, there exists a horizon time $T_\eta \geq 1$ and an online learning rule $g_{\leq \tau}$ where τ is a random time with $1 \leq \tau \leq T_\eta$ such that for any $\mathbf{y} := (y_t)_{t=1}^{T_\eta}$ of values in \mathcal{Y} and any value $y \in \mathcal{Y}$, we have*

$$\mathbb{E} \left[\frac{1}{\tau} \sum_{t=1}^{\tau} (\ell(g_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t)) \right] \leq \eta.$$

Remark 5.6. The random time τ may depend on the possible randomness of the learning rule g , but it does not depend on any of the values y_1, y_2, \dots on which the learning rule g may be tested. Intuitively, the learning rule uses some randomness which is first privately sampled and may be used by τ . This randomness is never explicitly revealed to the adversary choosing the values \mathbf{y} , but only implicitly through the realizations of the predictions.

Lemma 5.7. *Property F-TIME is equivalent to Property 2.*

Using this second property, we can then show that when F-TIME is not satisfied, universal consistency outside CS under adversarial responses is not achievable. In the proof, we only use stochastic processes (\mathbb{X}, \mathbb{Y}) , hence the same result holds if we only considered universal consistency under arbitrary responses.

Theorem 5.8. *Suppose that ℓ is bounded and (\mathcal{Y}, ℓ) does not satisfy F-TIME. Then, $\text{SOLAR} = \text{CS}$ and there exists an optimistically universal learning rule for adversarial regression, i.e., which achieves universal consistency with adversarial responses under any process $\mathbb{X} \in \text{CS}$.*

Proof sketch. First, from Theorem 5.2 we already have $\text{CS} \subset \text{SOLAR}$. The main difficulty is to prove that one cannot universally learn any process $\mathbb{X} \notin \text{CS}$. To do so, we re-use the property derived in the proof of Theorem 5.1 that for non-CS processes, one can find a disjoint sequence of sets $\{B_p\}_{p \geq 1}$, an increasing times $(t_p)_{p \geq 1}$ and $\epsilon > 0$ such that with non-zero probability for all $p \geq 1$, the process \mathbb{X} never visits B_p before time t_{p-1} and at some point between times $t_{p-1} + 1$ and t_p , the set B_p has been visited a proportion ϵ of times. Now (\mathcal{Y}, ℓ) does not satisfy F-TIME, hence does not satisfy Property 2 by Lemma 5.7 for some constant $\eta > 0$. Then, for $p \geq 1$, during period $(t_{p-1}, t_p]$, we define the values $\mathbb{Y}_{t_{p-1} < \cdot \leq t_p}$ when the instance process visits B_p as a sequence $\mathbf{y}_{t_{p-1} < \cdot \leq t_p}$ such that the algorithm has average excess loss at least η whenever \mathbb{X} visits B_p , compared to a fixed value $y_p^* \in \mathcal{Y}$. We note that the randomized version of F-TIME given by Lemma 5.7 is important because we do not know in advance when, between t_{p-1} and t_p , B_p has been visited a fraction ϵ of times: potentially, this time is random and there is a huge gap (exponential or more) between t_{p-1} and t_p . On the constructed stochastic process \mathbb{Y} , the algorithm does not have vanishing average excess loss compared to the function equal to y_p^* on B_p . This proves that no algorithm is universally consistent on \mathbb{X} .

This completes the proof of Corollary 3.5 and closes our study of universal learning with adversarial responses for bounded value spaces. Notably, there always exists an optimistically universal learning rule, however, this rule highly depends on the value space.

- If (\mathcal{Y}, ℓ) satisfies F-TIME, we can learn all $\text{SMV} = \text{SOUL}$ processes. The proposed learning rule of Theorem 5.5 is *implicit* in general. Indeed, to construct it one first needs to find an online learning rule for mean estimation with finite horizon as described by property F-TIME, which is then used as a subroutine in the optimistically universal learning rule for adversarial regression. We showed however that for totally-bounded value spaces, this learning rule can be *explicited* using ϵ -nets.
- If the value space does not satisfy F-TIME, we can only learn CS processes and there is an inherent gap between noiseless online learning and regression. We propose a learning rule in Section 7 which is optimistically universal—see Theorem 3.3. This rule is inspired by the proposed algorithm of [Han22] which is optimistically universal for metric losses $\alpha = 1$.

These two classes of learning rules use very different techniques. Specifically, under processes $\mathbb{X} \in \text{CS}$, [Han21a] showed that there exists a countable set \mathcal{F} of measurable functions $f : \mathcal{X} \rightarrow \mathcal{Y}$ which is “dense” within the space of all measurable functions along the realizations $f(X_t)$. We refer to Section 7 for a precise description of this density notion. Hence, under process \mathbb{X} , we can approximate f^* by functions in \mathcal{F} with arbitrary long-run average precision. However, such property is impossible to obtain for any process $\mathbb{X} \in \text{SMV} \setminus \text{CS}$: no process $\mathbb{X} \notin \text{CS}$ admits a “dense” countable sequence of measurable functions. Thus, to learn processes SMV for value spaces satisfying F-TIME, a fundamentally different learning rule than that proposed by [Han21a] or [Han22] was needed.

6 Adversarial universal learning for unbounded losses

We now turn to the case of unbounded losses, i.e., value spaces (\mathcal{Y}, ℓ) with $\bar{\ell} = \infty$. In this section, we consider universal learning without empirical integrability constraints, for which we introduced the notation SOLAR-U as the set of processes that admit universal learning (we recall that for bounded losses such distinction was

Input: Historical samples $(Y_t)_{t < T}$
Output: Predictions \hat{Y}_t for $t \leq T$
 $(y^i)_{i \geq 0}$ dense sequence in \mathcal{Y}
 $I_t := \{i \leq \ln t : \ell(y^0, y^i) \leq \ln t\}, \eta_t := \frac{1}{4\sqrt{t}}, t \geq 1; \quad t_i = \lceil \max(e^i, e^{\ell(y^0, y^i)}) \rceil, i \geq 0$
 $w_{0,0} := 1, \quad \hat{Y}_1 = y^0$ // Initialisation
for $t = 2, \dots, T$ **do**
 $L_{t-1,i} = \sum_{s=t_i}^{t-1} \ell(y^i, Y_s), \quad \hat{L}_{t-1,i} = \sum_{s=t_i}^{t-1} \hat{\ell}_s, \quad i \in I_t$
 $w_{t-1,i} := \exp(\eta_t (\hat{L}_{t-1,i} - L_{t-1,i})), \quad i \in I_t$
 $p_t(i) = \frac{w_{t-1,i}}{\sum_{j \in I_t} w_{t-1,j}}, \quad i \in I_t$
 $\hat{Y}_t \sim p_t(\cdot)$ // Prediction
 $\hat{\ell}_t := \frac{\sum_{j \in I_t} w_{t-1,j} \ell(y^j, Y_t)}{\sum_{j \in I_t} w_{t-1,j}}$
end

Algorithm 5: The mean estimation algorithm.

unnecessary). In this case, and for more general near-metrics, [BCH22] showed that $\text{SOUL} = \text{FS}$. In other terms, for unbounded losses, the learnable processes in the noiseless setting necessarily visit a finite number of distinct instance points of \mathcal{X} almost surely. Thus, universal learning on unbounded value spaces is very restrictive and in particular, $\text{SOLAR-U} \subset \text{FS}$. We will show that either $\text{SOLAR-U} = \text{FS}$ or $\text{SOLAR-U} = \emptyset$.

6.1 Adversarial regression for metric losses

In this section, we focus on metric losses ℓ , i.e., $\alpha = 1$. In this case, we show that we always have the equality $\text{SOLAR-U} = \text{FS}$ and that we can provide an optimistically universal learning rule. To do so, we first consider the fundamental estimation problem where one observes values \mathbb{Y} from a general separable metric value space and aims to sequentially predict a value \hat{Y}_t in order to minimize the long-run average loss. We refer to this problem as the mean estimation problem, which is equivalent to regression for the instance space $\mathcal{X} = \{0\}$. For instance, in the specific case of i.i.d. processes \mathbb{Y} , mean estimation is exactly the problem of Fréchet mean estimation for distributions on \mathcal{Y} . We show that even for adversarial processes \mathbb{Y} , we can achieve sublinear regret compared to the best single value prediction, even for unbounded value spaces (\mathcal{Y}, ℓ) .

If the space were finite, then we could use traditional Hedge algorithms [CL06]. Instead, given a separable value space, we have access to a dense countable sequence of values. We then select the best prediction among this dense sequence by introducing the values of the sequence one at a time, similarly to the argument we used in Lemma 4.2. The learning rule for mean estimation is described in Algorithm 5.

Theorem 3.6. *Let (\mathcal{Y}, ℓ) be a separable metric space. There exists an online learning rule f that is universally consistent for adversarial mean estimation, i.e., for any adversarial process \mathbb{Y} on \mathcal{Y} , almost surely, for all $y \in \mathcal{Y}$,*

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T (\ell(f_t(\mathbb{Y}_{\leq t-1}), Y_t) - \ell(y, Y_t)) \leq 0.$$

Remark 6.1. The above result guarantees that on the same event of probability one, the proposed learning rule achieves sublinear regret compared to any fixed value prediction. This was not the case for universal regression where, instead, for every fixed measurable function $f : \mathcal{X} \rightarrow \mathcal{Y}$, with probability one our learning rules achieved sublinear regret. This stems essentially from the fact that there exists a dense countable set of values \mathcal{Y} , but in general, there does not exist a countable set of measurable functions which are dense within all measurable functions in infinity norm.

We now return to the general regression problem on unbounded spaces. A simple learning rule would be to run in parallel the learning rule g_x for mean estimation on each distinct observed $x \in \mathcal{X}$, i.e., on the sub-process $\mathbb{Y}_{\{t: X_t=x\}}$. As a consequence of Theorem 3.6 we can show that this learning rule is universally consistent on FS processes.

Corollary 6.2. *Suppose that (\mathcal{Y}, ℓ) is an unbounded metric space. Then, $SOLAR-U = FS (= SOUL)$ and there exists an optimistically universal learning rule for adversarial regression, i.e., which achieves universal consistency with adversarial responses under any process $\mathbb{X} \in FS$.*

6.2 Negative result for real-valued adversarial regression with loss $\ell = |\cdot|^\alpha$ with $\alpha > 1$

Unfortunately, one cannot extend Corollary 6.2 to losses that are powers of metrics in general. Even in the classical setting of real-valued regression $\mathcal{Y} = \mathbb{R}$ with Euclidean norm, we show that adversarial regression with any loss $\ell = |\cdot|^\alpha$ for $\alpha > 1$ is not achievable, i.e., $SOLAR-U = \emptyset$.

Theorem 6.3. *Let $\alpha > 1$. For the Euclidean value space $(\mathbb{R}, |\cdot|)$ and loss $\ell = |\cdot|^\alpha$ we obtain $SOLAR-U = \emptyset$. In particular, there does not exist a consistent learning rule for mean estimation on \mathbb{R} with squared loss for adversarial responses.*

Proof sketch. The reason why mean estimation with adversarial responses is impossible for $\alpha > 1$ but possible for $\alpha = 1$ is that for $\alpha > 1$, predicting a value off by 1 unit of the best value in hindsight can yield unbounded excess loss for that specific prediction. In particular, we consider a sequence of values of the form $Y_t^{\mathbf{b}} = M_t b_t$ where $(M_t)_{t \geq 1}$ is a fixed sequence growing super-exponentially in t , and $\mathbf{b} = (b_t)$ is an i.i.d. Rademacher random variables in $\{\pm 1\}$. The sequence $(M_t)_{t \geq 1}$ is constructed so that if the prediction \hat{Y}_t and true value Y_t have different signs $\hat{Y}_t \cdot Y_t \leq 0$, the excess loss of the algorithm compared to the value $sign(Y_t^{\mathbf{b}}) = sign(b_t)$ is (super-)linear in t . Because the algorithm cannot know in advance the sign of b_t , there is a realization in which it makes an infinite number of mistakes and as a result has non-zero long-term excess loss compared to the value 1 or -1 .

The above of this result also shows that the same negative result holds more generally for unbounded metric value spaces which have some ‘‘symmetry’’. The main ingredients for this negative result were having a point from which there exist arbitrary far values from symmetric directions. In particular, this holds for a discretized value space $(\mathbb{N}, |\cdot|)$ with Euclidean metric, and any Euclidean space \mathbb{R}^d with $d \geq 1$.

6.3 An alternative for adversarial regression with unbounded losses

In the two previous sections, we gave examples of losses for which $SOLAR-U = \emptyset$ or $SOLAR-U = FS$. The following simple result is that this is the only alternative and that $SOLAR-U = FS$ is equivalent to achieving consistency for mean estimation with adversarial responses.

Proposition 6.4. *Let $(\mathcal{Y}, \rho_{\mathcal{Y}})$ be a separable metric value space. Suppose that there exists an online learning rule g which is consistent for mean estimation with adversarial responses for the loss $\ell = \rho_{\mathcal{Y}}^\alpha$, where $\alpha \geq 1$, i.e., for any adversarial process \mathbb{Y} on (\mathcal{Y}, ℓ) , we have for any $y^* \in \mathcal{Y}$,*

$$\limsup \frac{1}{T} \sum_{t=1}^T (\ell(f_t(\mathbb{Y}_{\leq t-1}), Y_t) - \ell(y^*, Y_t)) \leq 0, \quad (a.s),$$

then $SOLAR-U = FS$ and there exists an optimistically universal learning rule for adversarial regression. Otherwise, $SOLAR-U = \emptyset$.

Remark 6.5. There exists separable metric value spaces $(\mathcal{Y}, \rho_{\mathcal{Y}})$ for which powers of metrics losses still yield $SOLAR-U = FS$. For instance, consider $(\mathcal{Y}, \rho_{\mathcal{Y}}) = (\mathbb{R}, \sqrt{|\cdot|_2})$, where $|\cdot|_2$ denotes the Euclidean metric. One can check that this defines a metric on \mathcal{Y} and for any loss $\ell = \rho_{\mathcal{Y}}^\alpha$ with $\alpha \leq 2$, we have $SOLAR-U = FS$. However, for $\alpha > 2$, $SOLAR-U = \emptyset$.

7 Adversarial universal learning with moment constraint

In the previous section, we showed that learnable processes for adversarial regression are only in FS, i.e., visit a finite number of instance points. This shows that universal learning without restrictions on the adversarial responses \mathbb{Y} is extremely restrictive. For instance, it does not contain i.i.d. processes. A natural question is whether adding mild constraints on the process \mathbb{Y} would allow recovering the same results for unbounded losses as for bounded losses from Sections 4 and 5. This question also arises in noiseless regression since the set of learnable processes is reduced from $\text{SOUL} = \text{SMV}$ for bounded losses to $\text{SOUL} = \text{FS}$ for unbounded losses. Hence, [BCH22] posed as question whether having finite long-run empirical first-order moments would be sufficient to recover learnability in SMV. Precisely, they introduced the following constraint on noiseless processes $\mathbb{Y} = f^*(\mathbb{X})$: there exists $y_0 \in \mathcal{Y}$ with

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, f^*(X_t)) < \infty \quad (a.s.).$$

The question now becomes whether there exists an online learning rule which would be consistent under all $\mathbb{X} \in \text{SMV}$ processes for any noiseless responses $\mathbb{Y} = f^*(\mathbb{X})$ with f^* satisfying the above first-moment condition. We show that such an objective is not achievable whenever \mathcal{X} is infinite—if \mathcal{X} is finite, any process \mathbb{X} on \mathcal{X} is automatically FS and hence learnable in a noiseless or adversarial setting. In fact, under this first-order moment condition, we show the stronger statement that learning under all processes \mathbb{X} which admit pointwise convergent relative frequencies (CRF) is impossible even in this noiseless setting.

Condition CRF: For any measurable set $A \in \mathcal{B}$, $\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{1}_A(X_t)$ exists almost surely.

[Han21a] showed that $\text{CRF} \subset \text{CS}$. In particular, $\text{CRF} \subset \text{SMV}$. We show the following negative result on learning under CRF processes for noiseless regression under first-order moment constraint, which holds for unbounded near-metric spaces (\mathcal{Y}, ℓ) .

Theorem 7.1. *Suppose that \mathcal{X} is infinite and that (\mathcal{Y}, ℓ) is an unbounded separable near-metric space. There does not exist an online learning rule which would be consistent under all processes $\mathbb{X} \in \text{CRF}$ for all measurable target functions $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ such that there exists $y_0 \in \mathcal{Y}$ with*

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, f^*(X_t)) < \infty \quad (a.s.).$$

Proof sketch. We consider a sequence of values $(y_k)_{k \geq 0}$ such that $\ell(y_0, y_k)$ diverges as $k \rightarrow \infty$, then let $(t_k)_{k \geq 1}$ be a sequence of times such that $t_k \approx \sum_{k' \leq k} \ell(y_0, y_{k'})$. Next, let $(x_k)_{k \geq 0}$ be a sequence of distinct points. We construct a process \mathbb{X} such that $X_t = x_0$ except at sparse times $(t_k)_{k \geq 1}$ for which $X_{t_k} = x_k$. Because t_k has a super-linear growth, \mathbb{X} visits a sublinear number of distinct points and we can show that it satisfies the CRF property. Now for a random binary sequence $\mathbf{b} = (b_k)_{k \geq 1}$ we consider the function $f_{\mathbf{b}}^*$ which is equal to y_0 except at points x_k for $k \geq 1$ where $f_{\mathbf{b}}^*(x_k) = y_0 \mathbb{1}[b_k = 0] + y_k \mathbb{1}[b_k = 1]$. With these classes of functions, the algorithm cannot know in advance at time t_k whether to predict y_0 or y_k and incurs a loss $\mathcal{O}(\ell(y_0, y_k))$ in average as a result. Therefore, at time t_k , a total loss $\mathcal{O}(\sum_{k' \leq k} \ell(y_0, y_{k'})) = \mathcal{O}(t_k)$ is incurred compared to $f_{\mathbf{b}}^*$. On the other hand, by the construction of the sequence $(t_k)_{k \geq 1}$, $\frac{1}{T} \sum_{t=1}^T \ell(y_0, f_{\mathbf{b}}^*(X_t)) \leq \frac{1}{T} \sum_{t_k \leq T} \ell(y_0, y_{t_k})$ stays bounded. Thus the learning rule is not consistent under all target functions satisfying the specified moment constraint.

Theorem 7.1 answers negatively to the question posed in [BCH22]. A natural question is whether another meaningful constraint on responses can be applied to obtain positive results under large classes of processes on \mathcal{X} . To this means, we introduced the slightly stronger *empirical integrability* condition. We recall that an (adversarial) process \mathbb{Y} is *empirically integrable* if and only if there exists $y_0 \in \mathcal{Y}$ such that for any $\epsilon > 0$,

almost surely there exists $M \geq 0$ with

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M} \leq \epsilon.$$

Note that the threshold M may be *dependent* on the adversarial process \mathbb{Y} , but the guarantee should hold for any choice of predictions (in the case of adaptive adversaries). This is essentially the mildest condition on the sequence \mathbb{Y} for which we can still obtain results. For example, if the loss is bounded, this constraint is automatically satisfied using $M > \bar{\ell}$. More importantly, note that any process \mathbb{Y} which has bounded higher-than-first moments, i.e., such that there exists $p > 1$ and $y_0 \in \mathcal{Y}$ such that $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell^p(y_0, Y_t) < \infty$, (*a.s.*), is empirically integrable. Further, for stationary processes \mathbb{Y} , having bounded first moment $\mathbb{E}[\ell(y_0, Y_1)] < \infty$ is exactly being empirically integrable. Indeed, by the strong law of large numbers, almost surely $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M} = \mathbb{E}[\ell(y_0, Y_1) \mathbb{1}_{\ell(y_0, Y_1) \geq M}]$. Therefore, empirical integrability is a direct consequence of the dominated convergence theorem.

Lemma 7.2. *Let \mathbb{Y} an stationary process on \mathcal{Y} which has bounded first moment, i.e., there exists $y_0 \in \mathcal{Y}$ such that $\mathbb{E}[\ell(y_0, Y_1)] < \infty$. Then, \mathbb{Y} is empirically integrable.*

Proof. Let \mathbb{Y} an stationary process and $y_0 \in \mathcal{Y}$ with $\mathbb{E}[\ell(y_0, Y_1)] < \infty$. Then, by the dominated convergence theorem we have $\mathbb{E}[\ell(y_0, Y_1) \mathbb{1}_{\ell(y_0, Y_1) \geq M}] \rightarrow 0$ as $M \rightarrow \infty$. Hence, for $\epsilon > 0$, there exists M_ϵ such that $\mathbb{E}[\ell(y_0, Y_1) \mathbb{1}_{\ell(y_0, Y_1) \geq M_\epsilon}] \leq \epsilon$. Then, the sequence $(\ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M_\epsilon})_t$ is still stationary. hence, by the law of large numbers, almost surely,

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M_\epsilon} = \mathbb{E}[\ell(y_0, Y_1) \mathbb{1}_{\ell(y_0, Y_1) \geq M_\epsilon}] \leq \epsilon.$$

This ends the proof that \mathbb{Y} is empirically integrable. □

The goal of this section is to show that under this moment constraint, we can recover all results from [Bla22], [Han22] and this work in Sections 4 and 5, even for unbounded value spaces, leading up to Theorems 3.2 and 3.3. We will use the following simple equivalent formulation for empirical integrability.

Lemma 7.3. *A process \mathbb{Y} is empirically integrable if and only if there exists $y_0 \in \mathcal{Y}$ such that almost surely, for any $\epsilon > 0$ there exists $M > 0$ with*

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M} \leq \epsilon.$$

General strategy. First, the empirical integrability condition holds for some $y_0 \in \mathcal{Y}$ if and only if it holds for all $y_0 \in \mathcal{Y}$. Thus, we can fix $y_0 \in \mathcal{Y}$ independently of the instance or value process. Next, we define the restriction function $\phi_M : \mathcal{Y} \rightarrow \mathcal{Y}$ such that $\phi_M(y) = y$ if $\ell(y_0, y) < M$ and $\phi_M(y) = y_0$ otherwise. This function has values in the bounded set $B_\ell(y_0, M)$. Thus, we can apply our learning rules for the bounded loss case to learn the restricted values $\mathbb{Y}^M = (\phi_M(Y_t))_{t \geq 1}$. If we use these predictions to learn \mathbb{Y} , the excess loss compared to a fixed function mostly results from the restriction $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(Y_t, \phi_M(Y_t)) = \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M}$. This excess can then be bounded with the empirical integrability condition at y_0 . We then combine the resulting predictors for $M \geq 1$ using Lemma 4.2. While this general strategy allows to use learning rules for the bounded loss case as subroutine to solve the unbounded loss case with empirical integrability constraint, we can adapt it to each case to simplify the algorithms.

7.1 Noiseless universal learning with moment condition

We first apply this strategy to the noiseless case. The main result from [Bla22] showed that the 2C1NN learning rule achieves universal consistency on all SMV processes for bounded value spaces. Instead of using

the 2C1NN learning rule as subroutine as described in the strategy above, we show that we can readily use 2C1NN for empirically integrable noiseless responses in unbounded value spaces, as stated in Theorem 3.1.

To prove this result, we first observe that 2C1NN trained on the responses $\mathbb{Y} = (f^*(X_t))_{t \geq 1}$ or the restricted responses $(\phi_M \circ f^*(X_t))_{t \geq 1}$ gives the same prediction at time t provided that the representative $\phi(t)$ satisfied $\ell(y_0, Y_{\phi(t)}) < M$. By construction of the 2C1NN learning rule, points can be used as representatives at most twice. Hence, up to a factor 2, times when the predictions on unrestricted and restricted responses differ, can be associated with times when $\ell(y_0, Y_t) \geq M$. As a result, we show that the empirical integrability condition can be applied to bound the excess loss resulting from the difference between unrestricted and restricted responses.

7.2 Adversarial regression with moment condition under CS processes

We now turn to adversarial regression under CS processes. [Han22] showed that regression for arbitrary responses under all CS processes is achievable in bounded value spaces. We generalize this result to unbounded losses and to adversarial responses with empirical integrability constraint using the general strategy. In particular, our learning rule is also optimistically universal for adversarial regression for all bounded value spaces which do not satisfy F-TIME. Now consider the general case and suppose that there exists a ball $B_\ell(y, r)$ which does not satisfy F-TIME, Theorem 5.8 shows that universal learning for values falling in $B_\ell(y, r)$ cannot be achieved for processes $\mathbb{X} \notin \text{CS}$. Now because $B_\ell(y, r)$ is bounded, responses restricted to this set satisfy the empirical integrability constraint. In particular, this shows that the condition CS is also necessary for universal learning with adversarial responses with empirical integrability. Altogether, this proves Theorem 3.3.

This generalizes the main results from [Han22] to unbounded non-metric losses and from [CK22] to non-metric losses, arbitrary responses and CS instance processes \mathbb{X} . Indeed, they consider bounded first moment conditions on i.i.d. responses, which are empirically integrable by Lemma 7.2. Further, as a direct consequence of Theorem 3.3 and Lemma 7.2, we can significantly relax the conditions for universal consistency on stationary ergodic processes found in the literature. Precisely, [GO07] showed that for regression with squared loss, under the assumption $\mathbb{E}[Y_1^4] < \infty$, consistency on stationary ergodic processes is possible. We can relax this result to bounded second moments, matching the standard results for i.i.d. processes.

Corollary 7.4. *Let $(\mathcal{Y}, \ell) = (\mathbb{R}, |\cdot|^2)$. The learning rule of Theorem 3.3 is consistent on any stationary ergodic process $(X_t, Y_t)_{t \geq 1}$ with $\mathbb{E}[Y_1^2] < \infty$.*

7.3 Adversarial regression with moment condition under SMV processes

Last, we generalize our result Theorem 5.5 for value spaces satisfying F-TIME, to unbounded value spaces, with the same moment condition on responses using the general strategy. In order to apply Theorem 5.5 to bounded balls of the value space, we now ask that all balls $B_\ell(y, r)$ in the value space (\mathcal{Y}, ℓ) satisfy F-TIME. This proves Theorem 3.2.

Theorems 3.3 and 3.2 completely characterize learnability for adversarial regression with moment condition. Namely, if the value space (\mathcal{Y}, ℓ) is such that any bounded ball satisfies F-TIME (resp. there exists a ball $B_\ell(y, r)$ that disproves F-TIME), Theorem 3.2 (resp. 3.3) gives an optimistic learning rule which achieves consistency under all processes in SMV (resp. CS). This ends our analysis of adversarial regression for unbounded value spaces.

8 Open research directions

In this work, we provided a characterization of learnability for universal learning in the regression setting, for a class of losses satisfying specific relaxed triangle inequality identities, which contains powers of metrics $\ell = \rho_\alpha^\alpha$ for $\alpha \geq 1$. A natural question would be whether one can generalize these results to larger classes of losses, e.g. non-symmetric losses which may appear in classical machine learning problems.

The present work could also have some implications for adversarial contextual bandits. Specifically, one may consider the case of a learner who receives partial information on the rewards/losses as opposed to the traditional regression setting where the response is completely revealed at each iteration. In the latter case, the learner can for instance compute the loss of *all* values with respect to the response realization. On the other hand, in the contextual bandits framework, the reward/loss is revealed *only* for the pulled arm—or equivalently the prediction of the learner. In these partial information settings, exploration then becomes necessary. The authors are investigating whether the results presented in this work could have consequences in these related domains.

Acknowledgements. The authors are grateful to Prof. Steve Hanneke for enlightening discussions. This work is being partly funded by ONR grant N00014-18-1-2122.

References

- [BC12] Sébastien Bubeck and Nicolo Cesa-Bianchi. “Regret analysis of stochastic and nonstochastic multi-armed bandit problems”. In: *Machine Learning* 5.1 (2012), pp. 1–122.
- [BC22] Moïse Blanchard and Romain Cosson. “Universal online learning with bounded loss: reduction to binary classification”. In: *Conference on Learning Theory*. PMLR. 2022, pp. 479–495.
- [BCH22] Moïse Blanchard, Romain Cosson, and Steve Hanneke. “Universal online learning with unbounded losses: memory is all you need”. In: *International Conference on Algorithmic Learning Theory*. PMLR. 2022, pp. 107–127.
- [Bla22] Moïse Blanchard. “Universal online learning: an optimistically universal learning rule”. In: *Conference on Learning Theory*. PMLR. 2022, pp. 1077–1125.
- [BPS09] Shai Ben-David, Dávid Pál, and Shai Shalev-Shwartz. “Agnostic online learning.” In: *Conference on Learning Theory*. Vol. 3. 2009, p. 1.
- [Ces+97] Nicolo Cesa-Bianchi, Yoav Freund, David Haussler, David P Helmbold, Robert E Schapire, and Manfred K Warmuth. “How to use expert advice”. In: *Journal of the ACM (JACM)* 44.3 (1997), pp. 427–485.
- [Cha89] Ted Chang. “Spherical regression with errors in variables”. In: *The Annals of Statistics* (1989), pp. 293–306.
- [CK22] Dan Tsir Cohen and Aryeh Kontorovich. “Learning with metric losses”. In: *Conference on Learning Theory*. PMLR. 2022, pp. 662–700.
- [CL06] Nicolo Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge university press, 2006.
- [Dav+10] Brad C Davis, P Thomas Fletcher, Elizabeth Bullitt, and Sarang Joshi. “Population shape regression from random design data”. In: *International Journal of Computer Vision* 90.2 (2010), pp. 255–266.
- [Dev+94] Luc Devroye, Laszlo Györfi, Adam Krzyzak, and Gábor Lugosi. “On the strong universal consistency of nearest neighbor regression function estimates”. In: *The Annals of Statistics* (1994), pp. 1371–1385.
- [DGL13] Luc Devroye, László Györfi, and Gábor Lugosi. *A probabilistic theory of pattern recognition*. Vol. 31. Springer Science & Business Media, 2013.
- [EJ20] Steven N Evans and Adam Q Jaffe. “Strong laws of large numbers for Fréchet means”. In: *arXiv preprint arXiv:2012.12859* (2020).
- [Fle13] P. T. Fletcher. “Geodesic regression and the theory of least squares on Riemannian manifolds”. In: *International Journal of Computer Vision* 105.2 (2013), pp. 171–185.

- [FS97] Yoav Freund and Robert E Schapire. “A decision-theoretic generalization of on-line learning and an application to boosting”. In: *Journal of Computer and System Sciences* 55.1 (1997), pp. 119–139.
- [GG09] Robert M Gray and RM Gray. *Probability, random processes, and ergodic properties*. Vol. 1. Springer, 2009.
- [GLM99] L Györfi, Gábor Lugosi, and Gusztáv Morvai. “A simple randomized algorithm for sequential prediction of ergodic time series”. In: *IEEE Transactions on Information Theory* 45.7 (1999), pp. 2642–2650.
- [GO07] László Györfi and György Ottucsák. “Sequential prediction of unbounded stationary time series”. In: *IEEE Transactions on Information Theory* 53.5 (2007), pp. 1866–1872.
- [GW21] László Györfi and Roi Weiss. “Universal consistency and rates of convergence of multiclass prototype algorithms in metric spaces”. In: *Journal of Machine Learning Research* 22.151 (2021), pp. 1–25.
- [Gyö+02] László Györfi, Michael Köhler, Adam Krzyżak, and Harro Walk. *A distribution-free theory of nonparametric regression*. Vol. 1. Springer, 2002.
- [Han+21] Steve Hanneke, Aryeh Kontorovich, Sivan Sabato, and Roi Weiss. “Universal Bayes consistency in metric spaces”. In: *The Annals of Statistics* 49.4 (2021), pp. 2129–2150.
- [Han21a] Steve Hanneke. “Learning whenever learning is possible: Universal learning under general stochastic processes”. In: *Journal of Machine Learning Research* 22.130 (2021), pp. 1–116.
- [Han21b] Steve Hanneke. “Open Problem: is there an online learning algorithm that learns whenever online learning is possible?” In: *Conference on Learning Theory*. PMLR. 2021, pp. 4642–4646.
- [Han22] Steve Hanneke. “Universally consistent online learning with arbitrarily dependent responses”. In: *International Conference on Algorithmic Learning Theory*. PMLR. 2022, pp. 488–497.
- [Jaf22] Adam Quinn Jaffe. “Strong consistency for a class of adaptive clustering procedures”. In: *arXiv preprint arXiv:2202.13423* (2022).
- [Lit88] Nick Littlestone. “Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm”. In: *Machine learning* 2.4 (1988), pp. 285–318.
- [LM21] Zhenhua Lin and Hans-Georg Müller. “Total variation regularized Fréchet regression for metric-space valued data”. In: *The Annals of Statistics* 49.6 (2021), pp. 3510–3533.
- [LW94] Nick Littlestone and Manfred K Warmuth. “The weighted majority algorithm”. In: *Information and Computation* 108.2 (1994), pp. 212–261.
- [MJM00] Kanti V Mardia, Peter E Jupp, and KV Mardia. *Directional statistics*. Vol. 2. Wiley Online Library, 2000.
- [MKN99] Gusztáv Morvai, Sanjeev R Kulkarni, and Andrew B Nobel. “Regression estimation from an individual stable sequence”. In: *Statistics: A Journal of Theoretical and Applied Statistics* 33.2 (1999), pp. 99–118.
- [MYG96] Gusztáv Morvai, Sidney Yakowitz, and László Györfi. “Nonparametric inference for ergodic, stationary time series”. In: *The Annals of Statistics* 24.1 (1996), pp. 370–379.
- [RST15] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. “Online learning via sequential complexities.” In: *Journal of Machine Learning Research* 16.1 (2015), pp. 155–186.
- [Sch22] Christof Schötz. “Strong laws of large numbers for generalizations of Fréchet mean sets”. In: *Statistics* (2022), pp. 1–19.
- [Shi+09] Xiaoyan Shi, Martin Styner, Jeffrey Lieberman, Joseph G Ibrahim, Weili Lin, and Hongtu Zhu. “Intrinsic regression models for manifold-valued data”. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer. 2009, pp. 192–199.

- [SHS09] Ingo Steinwart, Don Hush, and Clint Scovel. “Learning from dependent observations”. In: *Journal of Multivariate Analysis* 100.1 (2009), pp. 175–194.
- [Sli19] Aleksandrs Slivkins. “Introduction to multi-armed bandits”. In: *arXiv preprint arXiv:1904.07272* (2019).
- [Sto77] Charles J Stone. “Consistent nonparametric regression”. In: *The Annals of Statistics* (1977), pp. 595–620.
- [Zai+19] Hanan Zaichyk, Armin Biess, Aryeh Kontorovich, and Yury Makarychev. “Efficient Kirschbraun extension with applications to regression”. In: *arXiv preprint arXiv:1905.11930* (2019).

Contents

1	Introduction	1
1.1	Motivation and background	1
1.2	Optimistic universal learning	2
1.3	Related works in universal learning	2
1.4	Adversarial responses and related works in learning with experts	2
1.5	Contributions	3
1.5.1	Universal learning with empirically integrable responses	3
1.5.2	Universal learning with unrestricted responses	3
1.6	Organization of the paper	4
2	Formal setup	4
2.1	Instance and value spaces	4
2.2	Online learning on adversarial responses	4
2.3	Empirically integrable responses	5
2.4	Universal consistency	5
2.5	Optimistic universal learning	6
3	Main results	6
4	An optimistically universal learning rule for totally-bounded value spaces	9
5	Characterization of learnable processes for bounded losses	13
5.1	Negative result for non-totally-bounded spaces	13
5.2	Adversarial regression for classification with a countable number of classes	14
5.3	A complete characterization of universal regression on bounded spaces	15
6	Adversarial universal learning for unbounded losses	17
6.1	Adversarial regression for metric losses	18
6.2	Negative result for real-valued adversarial regression with loss $\ell = \cdot ^\alpha$ with $\alpha > 1$	19
6.3	An alternative for adversarial regression with unbounded losses	19
7	Adversarial universal learning with moment constraint	20
7.1	Noiseless universal learning with moment condition	21
7.2	Adversarial regression with moment condition under CS processes	22
7.3	Adversarial regression with moment condition under SMV processes	22
8	Open research directions	22
A	Identities on the loss function	27
B	Proofs of Section 4	27
B.1	Proof of Theorem 4.1	27
B.1.1	Step 1	28
B.1.2	Step 2	30
B.1.3	Step 3	31
B.1.4	Step 4	31
B.2	Proof of Theorem 4.3	32
B.3	Proof of Lemma 4.2	35

C	Proofs of Section 5	37
C.1	Proof of Theorem 5.1	37
C.2	Proof of Lemma 5.3	40
C.3	Proof of Theorem 5.4	41
C.4	Proof of Theorem 5.5	42
C.5	Proof of Lemma 5.7	45
C.6	Proof of Theorem 5.8	47
D	Proofs of Section 6	50
D.1	Proof of Theorem 3.6	50
D.2	Proof of Corollary 6.2	51
D.3	Proof of Theorem 6.3	52
D.4	Proof of Proposition 6.4	53
E	Proofs of Section 7	54
E.1	Proof of Theorem 7.1	54
E.2	Proof of Lemma 7.3	55
E.3	Proof of Theorem 3.1	55
E.4	Proof of Theorem 3.3	57
E.5	Proof of Theorem 3.2	61

A Identities on the loss function

We recall the following known identities, which we will use to analyze the loss $\ell = \rho_{\mathcal{Y}}^\alpha$.

Lemma A.1. *Let $\alpha \geq 1$. Then, $(a + b)^\alpha \leq 2^{\alpha-1}(a^\alpha + b^\alpha)$ for all $a, b \geq 0$. Let $0 < \epsilon \leq 1$ and $\alpha \geq 1$. There exists some constant $c_\epsilon^\alpha > 0$ such that $(a + b)^\alpha \leq (1 + \epsilon)a^\alpha + c_\epsilon^\alpha b^\alpha$ for all $a, b \geq 0$, and $c_\epsilon^\alpha \leq \left(\frac{4\alpha}{\epsilon}\right)^\alpha$.*

Proof. The first identity is classical. A proof of the second one can be found for example in [EJ20] (Lemma 2.3) where they obtain $c_\epsilon^\alpha = \left(1 + \frac{1}{(1+\epsilon)^{1/\alpha-1}}\right)^\alpha \leq \left(\frac{4\alpha}{\epsilon}\right)^\alpha$. \square

B Proofs of Section 4

B.1 Proof of Theorem 4.1

In this section, we prove that for any $\delta > 0$, the $(1 + \delta)$ C1NN learning rule is optimistically universal for the noiseless setting. The proof follows the same structure as the proof of the main result in [Bla22] which shows that 2C1NN is optimistically universal. We first focus on the binary classification setting and show that the learning rule $(1 + \delta)$ C1NN is consistent on functions representing open balls.

Proposition B.1. *Fix $0 < \delta \leq 1$. Let $(\mathcal{X}, \mathcal{B})$ be a separable Borel space constructed from the metric $\rho_{\mathcal{X}}$. We consider the binary classification setting $\mathcal{Y} = \{0, 1\}$ and the ℓ_{01} binary loss. For any input process $\mathbb{X} \in \text{SMV}$, for any $x \in \mathcal{X}$, and $r > 0$, the learning rule $(1 + \delta)$ C1NN is consistent for the target function $f^* = \mathbb{1}_{B_{\rho_{\mathcal{X}}}(x, r)}$.*

Proof. We fix $\bar{x} \in \mathcal{X}$, $r > 0$ and $f^* = \mathbb{1}_{B(\bar{x}, r)}$. We reason by the contrapositive and suppose that $(1 + \delta)$ C1NN is not consistent on f^* . Then, $\eta := \mathbb{P}(\mathcal{L}_{\mathbb{X}}((1 + \delta)\text{C1NN}, f^*) > 0) > 0$. Therefore, there exists $0 < \epsilon \leq 1$ such that $\mathbb{P}(\mathcal{L}_{\mathbb{X}}((1 + \delta)\text{C1NN}, f^*) > \epsilon) > \frac{\eta}{2}$. Denote by $\mathcal{A} := \{\mathcal{L}_{\mathbb{X}}((1 + \delta)\text{C1NN}, f^*) > \epsilon\}$. this event of probability at least $\frac{\eta}{2}$. Because \mathcal{X} is separable, let $(x^i)_{i \geq 1}$ a dense sequence of \mathcal{X} . We consider the same partition $(P_i)_{i \geq 1}$ of $B(\bar{x}, r)$ and the partition $(A_i)_{i \geq 0}$ of \mathcal{X} as in the original proof of [Bla22], but with the constant $c_\epsilon := \frac{1}{2 \cdot 2^{2^8/(\epsilon\delta)}}$ and changing the construction of the sequence $(n_l)_{l \geq 1}$ so that for all $l \geq 1$

$$\mathbb{P} \left[\forall n \geq n_l, |\{i, P_i(\tau_l) \cap \mathbb{X}_{<n} \neq \emptyset\}| \leq \frac{\epsilon\delta}{2^{10}n} \right] \geq 1 - \frac{\delta}{2 \cdot 2^{l+2}} \quad \text{and} \quad n_{l+1} \geq \frac{2^9}{\epsilon\delta} n_l.$$

Last, consider the product partition of $(P_i)_{i \geq 1}$ and $(A_i)_{i \geq 0}$ which we denote \mathcal{Q} . Similarly, we define the same events $\mathcal{E}_l, \mathcal{F}_l$ for $l \geq 1$. We aim to show that with nonzero probability, \mathbb{X} does not visit a sublinear number of sets of \mathcal{Q} .

We now denote by $(t_k)_{k \geq 1}$ the increasing sequence of all (random) times when $(1 + \delta)$ C1NN makes an error in the prediction of $f^*(X_t)$. Because the event \mathcal{A} is satisfied, $\mathcal{L}_{\mathbf{x}}((1 + \delta)\text{C1NN}, f^*) > \epsilon$, we can construct an increasing sequence of indices $(k_l)_{l \geq 1}$ such that $t_{k_l} < \frac{2k_l}{\epsilon}$. For any $t \geq 2$, we will denote by $\phi(t)$ the (random) index of the representative chosen by the $(1 + \delta)$ C1NN learning rule. Now let $l \geq 1$. Consider the tree \mathcal{G} where nodes are times $\mathcal{T} := \{t \leq t_{k_l}\}$ within horizon t_{k_l} , where the parent relations are given by $(t, \phi(t))$ for $t \in \mathcal{T} \setminus \{1\}$. In other words, we construct the tree in which the parent of each new input is its representative. Note that by construction of the $(1 + \delta)$ C1NN learning rule, each node has at most 2 children.

B.1.1 Step 1

In this step, we consider the case when the majority of input points on which $(1 + \delta)$ C1NN made a mistake belong to $B(\bar{x}, r)$, i.e., $|\{k \leq k_l, X_{t_k} \in B(\bar{x}, r)\}| \geq \frac{k_l}{2}$. We denote \mathcal{H}_1 this event. Let us now consider the subgraph $\tilde{\mathcal{G}}$ given by restricting \mathcal{G} only to nodes in the ball $B(\bar{x}, r)$ —which are mapped to the true value 1—i.e., on times $\mathcal{T} := \{t \leq t_{k_l}, X_t \in B(\bar{x}, r)\}$. In this subgraph, the only times with no parent are times t_k with $k \leq k_l$ and $X_{t_k} \in B(\bar{x}, r)$, and possibly time $t = 1$. Therefore, $\tilde{\mathcal{G}}$ is a collection of disjoint trees with roots times $\{t_k, k \leq k_l, X_{t_k} \in B(\bar{x}, r)\}$, and possibly $t = 1$ if $X_1 \in B(\bar{x}, r)$. For a given time t_k with $k \leq k_l$ and $X_{t_k} \in B(\bar{x}, r)$, we denote by \mathcal{T}_k the corresponding tree in $\tilde{\mathcal{G}}$ with root t_k . We now introduce the notion of *good* trees. We say that \mathcal{T}_k is a good tree if $\mathcal{T}_k \cap \mathcal{D}_{t_{k_l}+1} \neq \emptyset$, i.e., the tree survived until the last dataset. Conversely a tree is *bad* if all its nodes were deleted before time $t_{k_l} + 1$. We denote the set of good and bad trees by $G = \{k : \mathcal{T}_k \text{ good}\}$ and $B = \{k : \mathcal{T}_k \text{ bad}\}$. In particular, we have $|G| + |B| = |\{k \leq k_l, X_{t_k} \in B(\bar{x}, r)\}| \geq k_l/2$. We aim to upper bound the number of bad trees. We now focus on trees \mathcal{T}_k which induced a future first mistake, i.e., such that $\{l \in \mathcal{T}_k | \exists u \leq t_{k_l} : \phi(u) = l, \rho_{\mathcal{X}}(X_l, \bar{x}) \geq r \text{ and } \forall v < u, \phi(v) \neq l\} \neq \emptyset$. We denote the corresponding minimum time $l_k = \min\{l \in \mathcal{T}_k | \exists u \leq t_{k_l} : \phi(u) = l, \rho_{\mathcal{X}}(X_l, \bar{x}) \geq r, \forall v < u, \phi(v) \neq l\}$. The terminology first mistake refers to the fact that the first time which used l as representative corresponded to a mistake, as opposed to l already having a children $X_u \in B(\bar{x}, r)$ which continues descendents of l within the tree \mathcal{T}_k . Note that bad trees necessarily induce a future first mistake—otherwise, this tree would survive. For each of these times l_k two scenarios are possible.

1. The value U_{l_k} was never revealed within horizon t_{k_l} : as a result $l_k \in \mathcal{D}_{t_{k_l}+1}$.
2. The value U_{l_k} was revealed within horizon t_{k_l} . Then, U_{l_k} we revealed using a time t for which l_k was a potential representative. This scenario has two cases:
 - (a) $\rho_{\mathcal{X}}(X_t, \bar{x}) < r$. If used as representative $\phi(t) = l_k$, then l_k would not have induced a mistake in the prediction of Y_t .
 - (b) $\rho_{\mathcal{X}}(X_t, \bar{x}) \geq r$. If used as representative $\phi(t) = l_k$, then l_k would have induced a mistake in the prediction of Y_t .

In the case 2.a), if the point is used as representative $\phi(t) = l_k$ and if the corresponding tree \mathcal{T}_k was bad, at least another future mistake is induced by \mathcal{T}_k —otherwise this tree would survive. We consider times l_k for which the value was revealed, which corresponds to the only possible scenario for bad trees. We denote the corresponding set $K := \{k : U_{l_k} \text{ revealed within horizon } t_{k_l}\}$. We now consider the sequence k_1^a, \dots, k_α^a containing all indices of K for which scenario 2.a) was followed, ordered by chronological order for the reveal of $U_{l_{k_i^a}}$, i.e., $U_{l_{k_1^a}}$ was the first item of scenario 2.a) to be revealed, then $U_{l_{k_2^a}}$ etc. until $U_{l_{k_\alpha^a}}$. Similarly, we construct the sequence k_1^b, \dots, k_β^b of indices in K corresponding to scenario 2.b), ordered by order for the

reveal of $U_{l_{k_i^b}}$. We now consider the events

$$\mathcal{B} := \left\{ \alpha + \beta \leq \frac{k_l}{2} - \frac{k_l \delta}{32} \right\}, \quad \mathcal{C} := \left\{ \sum_{i=1}^{\min(\alpha, \lceil k_l/8 \rceil)} U_{l_{k_i^a}} \geq \frac{k_l \delta}{16} \right\},$$

$$\mathcal{D} := \left\{ \sum_{i=1}^{\min(\beta, \lceil k_l/8 \rceil)} U_{l_{k_i^b}} \geq \frac{k_l \delta}{16} \right\}.$$

We now show that for $l > 16$, under the event

$$\mathcal{M}_{k_l} := \mathcal{H}_1 \cap [\mathcal{B} \cup (\{\alpha \geq \lceil k_l/8 \rceil\} \cap \mathcal{C}) \cup (\{\alpha < \lceil k_l/8 \rceil\} \cap \mathcal{D})],$$

we have that $|G| \geq \frac{k_l \delta}{32}$. Suppose that \mathcal{M}_{k_l} is met. First note that because a bad tree can only fall into scenarios 2.a) or 2.b) we have $|B| \leq \alpha + \beta$. Hence $|G| \geq \frac{k_l}{2} - \alpha - \beta$ because of \mathcal{H}_1 . Thus, the result holds directly if \mathcal{B} is satisfied. We can now suppose that \mathcal{B}^c is satisfied, i.e., $\alpha + \beta > \frac{k_l}{2} - \frac{k_l \delta}{32}$. Now suppose that $\alpha \geq \lceil k_l/8 \rceil$ and \mathcal{C} are also satisfied. For all indices such that $U_{l_{k_i^a}} = 1$, i.e., we fall in case 2.a) and $l_{k_i^a}$ is used as representative, the corresponding tree $\mathcal{T}_{k_i^a}$ would need to induce at least an additional mistake to be bad. Recall that in total at most $k_l/2$ mistakes are induced by points of \mathcal{T} . Also, by definition of the set K , $\alpha + \beta$ mistakes are already induced by the times t_k for $k \in K$. These corresponded to the future first mistakes for all times $\{l_k : k \in K\}$. Hence, we obtain

$$|G| \geq \sum_{i=1}^{\alpha} U_{l_{k_i^a}} - \left(\frac{k_l}{2} - \alpha - \beta \right) \geq \frac{k_l \delta}{16} - \frac{k_l \delta}{32} = \frac{k_l \delta}{32}.$$

Now consider the case where \mathcal{H}_1 , \mathcal{B}^c , $\alpha < \lceil k_l/8 \rceil$ and \mathcal{D} are met. In particular, because $l > 16$ we have $k_l > 16$ hence $\frac{k_l}{2} - \frac{k_l \delta}{32} \geq 2\lceil k_l/8 \rceil$. Thus, because of \mathcal{B}^c we have $\beta > \frac{k_l}{2} - \frac{k_l \delta}{32} - \alpha \geq \lceil k_l/8 \rceil$. Now observe that for all indices such that $U_{l_{k_i^b}} = 1$, the time l_k induced two mistakes. Therefore, counting the total number of mistakes we obtain

$$\frac{k_l}{2} \geq \alpha + \beta + \sum_{i=1}^{\beta} U_{l_{k_i^b}} \geq \frac{k_l}{2} - \frac{k_l \delta}{32} + \frac{k_l \delta}{16}$$

which is impossible. This ends the proof that under \mathcal{M}_{k_l} we have $|G| \geq \frac{k_l \delta}{32}$.

We now aim to lower bound the probability of this event. To do so, we first upper bound the probability of the event $\{\alpha \geq \lceil k_l/8 \rceil\} \cap \mathcal{C}^c$. We introduce a process $(Z_i)_{i=1}^{\lceil k_l/8 \rceil}$ such that for all $i \leq \max(\alpha, \lceil k_l/8 \rceil)$, $Z_i = U_{l_{k_i^a}} - \delta$ and $Z_i = 0$ for $\alpha < i \leq \lceil k_l/8 \rceil$. Because of the specific ordering chosen k_1^a, \dots, k_α^a , this process is a sequence of martingale differences, with values bounded by 1 in absolute value. Therefore, for $l > 16$ the Azuma-Hoeffding inequality yields

$$\mathbb{P} \left[\sum_{i=1}^{\lceil k_l/8 \rceil} Z_i \leq -\frac{k_l \delta}{16} \right] \leq e^{-\frac{k_l^2 \delta^2}{2 \cdot 16^2 (\lceil k_l/8 \rceil + 1)}} \leq e^{-\frac{k_l \delta^2}{2^7}}.$$

But on the event $\{\alpha \geq \lceil k_l/8 \rceil\} \cap \mathcal{C}^c$ we have precisely

$$\sum_{i=1}^{\lceil k_l/8 \rceil} Z_i = \sum_{i=1}^{\min(\alpha, \lceil k_l/8 \rceil)} U_{l_{k_i^a}} - \lceil k_l/8 \rceil \delta \leq \frac{k_l \delta}{16} - \lceil k_l/8 \rceil \delta \leq -\frac{k_l \delta}{16}.$$

Therefore $\mathbb{P}[\mathcal{C}^c \cap \{\alpha \geq \lceil k_l/8 \rceil\}] \leq \mathbb{P} \left[\sum_{i=1}^{\lceil k_l/8 \rceil} Z_i \leq -\frac{k_l \delta}{16} \right] \leq e^{-k_l \delta^2 / 2^7}$. Similarly we obtain $\mathbb{P}[\mathcal{D}^c \cap \{\beta \geq$

$\lceil k_l/8 \rceil\} \leq e^{-k_l\delta^2/2^7}$. Finally we write for any $l > 16$,

$$\begin{aligned} \mathbb{P}[\mathcal{H}_1 \setminus \mathcal{M}_{k_l}] &= \mathbb{P}[\mathcal{H}_1 \cap \mathcal{B}^c \cap (\{\alpha < \lceil k_l/8 \rceil\} \cup \mathcal{C}^c) \cap (\{\alpha \geq \lceil k_l/8 \rceil\} \cup \mathcal{D}^c)] \\ &= \mathbb{P}[\mathcal{H}_1 \cap \mathcal{B}^c \cap [(\{\alpha < \lceil k_l/8 \rceil\} \cap \mathcal{D}^c) \cup (\{\alpha \geq \lceil k_l/8 \rceil\} \cap \mathcal{C}^c)]] \\ &\leq \mathbb{P}[\mathcal{C}^c \cap \{\alpha \geq \lceil k_l/8 \rceil\}] + \mathbb{P}[\mathcal{D}^c \cap \{\alpha < \lceil k_l/8 \rceil\} \cap \mathcal{B}^c] \\ &\leq \mathbb{P}[\mathcal{C}^c \cap \{\alpha \geq \lceil k_l/8 \rceil\}] + \mathbb{P}[\mathcal{D}^c \cap \{\beta \geq \lceil k_l/8 \rceil\}] \\ &\leq 2e^{-\frac{k_l\delta^2}{2^7}}. \end{aligned}$$

In particular, we obtain

$$\mathbb{P}\left[\left\{|G| \geq \frac{k_l\delta}{32}\right\} \cap \mathcal{H}_1\right] \geq \mathbb{P}[\mathcal{M}_{k_l}] \geq \mathbb{P}[\mathcal{H}_1] - 2e^{-\frac{k_l\delta^2}{2^7}}.$$

B.1.2 Step 2

We now consider the opposite case, when a majority of mistakes are made outside $B(\bar{x}, r)$, i.e., $|\{k \leq k_l, X_{t_k} \in B(\bar{x}, r)\}| < \frac{k_l}{2}$, which corresponds to the event \mathcal{H}_1^c . Similarly, we consider the subgraph $\tilde{\mathcal{G}}$ given by restricting \mathcal{G} only to nodes outside the ball $B(\bar{x}, r)$, i.e., on times $\mathcal{T} := \{t \leq t_{k_l}, \rho_{\mathcal{X}}(X_t, \bar{x}) \geq r\}$. Again, $\tilde{\mathcal{G}}$ is a collection of disjoint trees with roots times $\{t_k, k \leq k_l, \rho_{\mathcal{X}}(X_{t_k}, \bar{x}) \geq r\}$ —and possibly $t = 1$. For a given time t_k with $k \leq k_l$ and $\rho_{\mathcal{X}}(X_{t_k}, \bar{x}) \geq r$, we denote by \mathcal{T}_k the corresponding tree in $\tilde{\mathcal{G}}$ with root t_k . Similarly to the previous case, \mathcal{T}_k is a *good* tree if $\mathcal{T}_k \cap \mathcal{D}_{t_{k_l+1}} \neq \emptyset$ and *bad* otherwise. We denote the set of good and bad trees by $G = \{k : \mathcal{T}_k \text{ good}\}$. We can again focus on trees \mathcal{T}_k which induced a future first mistake, i.e., such that $\{l \in \mathcal{T}_k | \exists u \leq t_{k_l} : \phi(u) = l, \rho_{\mathcal{X}}(X_l, \bar{x}) < r \text{ and } \forall v < u, \phi(v) \neq l\} \neq \emptyset$ and more specifically their minimum time $l_k = \min\{l \in \mathcal{T}_k | \exists u \leq t_{k_l} : \phi(u) = l, \rho_{\mathcal{X}}(X_l, \bar{x}) < r, \forall v < u, \phi(v) \neq l\}$. The same analysis as above shows that

$$\mathbb{P}\left[\left\{|G| \geq \frac{k_l\delta}{32}\right\} \cap \mathcal{H}_1^c\right] \geq \mathbb{P}[\mathcal{H}_1^c] - 2e^{-\frac{k_l\delta^2}{2^7}}.$$

Therefore, if G denotes more generally the set of good trees (where we follow the corresponding case 1 or 2) we finally obtain that for any $l > 16$,

$$\mathbb{P}\left[|G| \geq \frac{k_l\delta}{32}\right] \geq 1 - 4e^{-\frac{k_l\delta^2}{2^7}}.$$

We denote by $\tilde{\mathcal{M}}_{k_l}$ this event. By Borel-Cantelli lemma, almost surely, there exists \hat{l} such that for any $l \geq \hat{l}$, the event $\tilde{\mathcal{M}}_{k_l}$ is satisfied. We denote $\mathcal{M} := \bigcup_{l \geq 1} \bigcap_{l' > l} \tilde{\mathcal{M}}_{k_{l'}}$ this event of probability one. The aim is to show that on the event $\mathcal{A} \cap \mathcal{M} \cap \bigcap_{l \geq 1} (\mathcal{E}_l \cap \mathcal{F}_l)$, which has probability at least $\frac{7}{4}$, \mathbb{X} disproves the SMV condition. In the following, we consider a specific realization \mathbf{x} of the process \mathbb{X} falling in the event $\mathcal{A} \cap \mathcal{M} \cap \bigcap_{l \geq 1} (\mathcal{E}_l \cap \mathcal{F}_l)$ — \mathbf{x} is not random anymore. Let \hat{l} be the index given by the event \mathcal{M} such that for any $l \geq \hat{l}$, \mathcal{M}_{k_l} holds. We consider $l \geq \hat{l}$ and successively consider different cases in which the realization \mathbf{x} may fall.

- In the first case, we suppose that a majority of mistakes were made in $B(\bar{x}, r)$, i.e., that we fell into event \mathcal{H}_1 similarly to Step 1. Because the event $\tilde{\mathcal{M}}_{k_l}$ is satisfied we have $|G| \geq \frac{k_l\delta}{2^5}$. Now note that trees are disjoint, therefore, $\sum_{k \in G} |\mathcal{T}_k| \leq t_{k_l} < \frac{2k_l}{\epsilon}$. Therefore,

$$\sum_{k \in G} \mathbb{1}_{|\mathcal{T}_k| \leq \frac{2^7}{\epsilon\delta}} = |G| - \sum_{k \in G} \mathbb{1}_{|\mathcal{T}_k| > \frac{2^7}{\epsilon\delta}} > |G| - \frac{\epsilon\delta}{2^7} \sum_{k \in G} |\mathcal{T}_k| \geq \frac{k_l\delta}{2^5} - \frac{k_l\delta}{2^6} = \frac{k_l\delta}{2^6}.$$

We will say that a tree $|\mathcal{T}_k|$ is *sparse* if it is good and has at most $\frac{2^7}{\epsilon\delta}$ nodes. With $S := \{k \in G, |\mathcal{T}_k| \leq \frac{2^7}{\epsilon\delta}\}$ the set of sparse trees, the above equation yields $|S| \geq \frac{k_l\delta}{2^6}$. The same arguments as in [Bla22] give

$$|\{i, A_i \cap \mathbf{x}_{\leq t_{k_l}} \neq \emptyset\}| \geq |S| \geq \frac{k_l\delta}{2^6} \geq \frac{\epsilon\delta}{2^7} t_{k_l}.$$

The only difference is that we chose c_ϵ so that $2^{2 \cdot \frac{27}{\epsilon\delta} - 1} \leq \frac{1}{4c_\epsilon}$ as needed in the original proof.

- We now turn to the case when the majority of input points on which $(1 + \delta)$ C1NN made a mistake are not in the ball $B(\bar{x}, r)$, similarly to Step 2. Using the same notion of sparse tree $S := \{k \in G, |\mathcal{T}_k| \leq \frac{2^7}{\epsilon\delta}\}$, we have again $|S| \geq \frac{k_l \delta}{2^6}$. We use the same arguments as in the original proof. Suppose $|\{k \in S, \rho_{\mathcal{X}}(x_{p_{d(k)}^k}, \bar{x}) > r\}| \geq \frac{|S|}{2}$, then we have

$$|\{i, A_i \cap \mathbf{x}_{\leq t_{k_l}} \neq \emptyset\}| \geq |\{k \in S, \rho_{\mathcal{X}}(x_{p_{d(k)}^k}, \bar{x}) > r\}| \geq \frac{|S|}{2} \geq \frac{k_l \delta}{2^7} \geq \frac{\epsilon\delta}{2^8} t_{k_l}.$$

B.1.3 Step 3

In this last step, we suppose again that the majority of input points on which $(1 + \delta)$ C1NN made a mistake are not in the ball $B(\bar{x}, r)$ but that $|\{k \in S, \rho_{\mathcal{X}}(x_{p_{d(k)}^k}, \bar{x}) > r\}| < \frac{|S|}{2}$. Therefore, we obtain

$$|\{k \in S, \rho_{\mathcal{X}}(x_{p_{d(k)}^k}, \bar{x}) = r\}| = |S| - |\{k \in S, \rho_{\mathcal{X}}(x_{p_{d(k)}^k}, \bar{x}) > r\}| \geq \frac{|S|}{2} \geq \frac{k_l \delta}{2^7} \geq \frac{\epsilon\delta}{2^8} t_{k_l}.$$

We will now make use of the partition $(P_i)_{i \geq 1}$. Because $(n_u)_{u \geq 1}$ is an increasing sequence, let $u \geq 1$ such that $n_{u+1} \leq t_{k_l} \leq n_{u+2}$ (we can suppose without loss of generality that $t_{k_0} > n_2$). Note that we have $n_u \leq \frac{\epsilon\delta}{2^9} n_{u+1} \leq \frac{\epsilon\delta}{2^9} t_{k_l}$. Let us now analyze the process between times n_u and t_{k_l} . In particular, we are interested in the indices $T = \{k \in S, \rho_{\mathcal{X}}(x_{p_{d(k)}^k}, \bar{x}) = r\}$ and times $\mathcal{U}_u = \{p_{d(k)}^k : n_u < p_{d(k)}^k \leq k_l, k \in T\}$. In particular, we have

$$|\mathcal{U}_u| \geq |\{k \in S, \rho_{\mathcal{X}}(x_{p_{d(k)}^k}, \bar{x}) = r\}| - n_u \geq \frac{\epsilon\delta}{2^8} t_{k_l} - \frac{\epsilon\delta}{2^9} t_{k_l} = \frac{\epsilon\delta}{2^9} t_{k_l}.$$

Defining $T' := \{k \in T, r - \frac{r}{2^{u+3}} \leq \rho_{\mathcal{X}}(x_{\phi(t_k)}, \bar{x}) < r\}$, the same arguments as in the original proof yield

$$|\{i, P_i \cap \mathbf{x}_{\leq t_{k_l}} \neq \emptyset\}| \geq |T'| \geq |\mathcal{U}_u| - |\{i, P_i(\tau_u) \cap \mathbf{x}_{\mathcal{U}_u} \neq \emptyset\}| \geq \frac{\epsilon\delta}{2^9} t_{k_l} - \frac{\epsilon\delta}{2^{10}} t_{k_l} = \frac{\epsilon\delta}{2^{10}} t_{k_l}.$$

B.1.4 Step 4

In conclusion, in all cases, we obtain

$$|\{Q \in \mathcal{Q}, Q \cap \mathbf{x}_{\leq t_{k_l}} \neq \emptyset\}| \geq \max(|\{i, A_i \cap \mathbf{x}_{\leq t_{k_l}} \neq \emptyset\}|, |\{i, P_i \cap \mathbf{x}_{\leq t_{k_l}} \neq \emptyset\}|) \geq \frac{\epsilon\delta}{2^{10}} t_{k_l}.$$

Because this is true for all $l \geq \hat{l}$ and t_{k_l} is an increasing sequence, we conclude that \mathbf{x} disproves the SMV condition for \mathcal{Q} . Recall that this holds whenever the event $\mathcal{A} \cap \mathcal{M} \cap \bigcap_{l \geq 1} (\mathcal{E}_l \cap \mathcal{F}_l)$ is met. Thus,

$$\mathbb{P}[\{Q \in \mathcal{Q}, Q \cap \mathbb{X}_{< T}\} = o(T)] \leq 1 - \mathbb{P} \left[\mathcal{A} \cap \mathcal{M} \cap \bigcap_{l \geq 1} (\mathcal{E}_l \cap \mathcal{F}_l) \right] \leq 1 - \frac{\eta}{4} < 1.$$

This shows that $\mathbb{X} \notin \text{SMV}$ which is absurd. Therefore $(1 + \delta)$ C1NN is consistent on f^* . This ends the proof of the proposition. \square

Using the fact that in the $(1 + \delta)$ C1NN learning rule, no time t can have more than 2 children, as the 2C1NN rule, we obtain with the same proof as in [Bla22] the following proposition.

Proposition B.2. *Fix $0 < \delta \leq 1$. Let $(\mathcal{X}, \mathcal{B})$ be a separable Borel space. For the binary classification setting, the learning rule $(1 + \delta)$ C1NN is universally consistent for all processes $\mathbb{X} \in \text{SMV}$.*

Finally, we use a result from [BC22] which gives a reduction from any near-metric bounded value space to binary classification.

Theorem B.3 ([BC22]). *If $(1 + \delta)$ C1NN is universally consistent under a process \mathbb{X} for binary classification, it is also universally consistent under \mathbb{X} for any separable near-metric setting (\mathcal{Y}, ℓ) with bounded loss.*

Together with Proposition B.2, Theorem B.3 ends the proof of Theorem 4.1.

B.2 Proof of Theorem 4.3

Let $0 < \epsilon \leq 1$. We first analyze the prediction of the learning rule f^ϵ . In the rest of the proof, we denote $\bar{\ell}(\hat{Y}_t(\epsilon), Y_t) := \sum_{y \in \mathcal{Y}_\epsilon} \mathbb{P}(\hat{Y}_t(\epsilon) = y) \ell(y, Y_t)$ the immediate expected loss at each iteration. The learning rule was constructed so that we perform exactly the classical Hedge / exponentially weighted average forecaster on each cluster of times $\mathcal{C}(t) = \{u \leq t : u \stackrel{\phi}{\sim} t\}$. As a result [CL06] (Theorem 2.2), we have that for any $t \geq 1$,

$$\begin{aligned} \frac{1}{\bar{\ell}} \sum_{u \in \mathcal{C}(t)} \bar{\ell}(\hat{Y}_u(\epsilon), Y_u) &\leq \frac{1}{\bar{\ell}} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u \in \mathcal{C}(t)} \ell(y, Y_u) + \frac{\ln |\mathcal{Y}_\epsilon|}{\bar{\ell} \eta_\epsilon} + \frac{|\mathcal{C}(t)| \bar{\ell} \eta_\epsilon}{8} \\ &\leq \frac{1}{\bar{\ell}} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u \in \mathcal{C}(t)} \ell(y, Y_u) + \sqrt{\frac{\ln |\mathcal{Y}_\epsilon|}{8 T_\epsilon}} (T_\epsilon + |\mathcal{C}(t)|) \\ &\leq \frac{1}{\bar{\ell}} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u \in \mathcal{C}(t)} \ell(y, Y_u) + \frac{\epsilon}{\bar{\ell}} \max(T_\epsilon, |\mathcal{C}(t)|) \end{aligned}$$

Now consider a horizon $T \geq 1$, and enumerate all the clusters $\mathcal{C}_1(T), \dots, \mathcal{C}_{p(T)}(T)$ at horizon T , i.e. the classes of equivalence of ϕ among the times $\{t \leq T\}$. Note that if a cluster $i \leq p$ has $|\mathcal{C}_i(T)| < T_\epsilon$, then either it must contain a time $t \in \mathcal{N}$ which is a leaf of the tree formed by ϕ until time T , or it is a cluster of duplicates of an instance X_u which has already had $\frac{T}{\epsilon}$ occurrences. As a result, the times falling into such clusters of duplicates with less than T_ϵ members form at most a proportion ϵ of the total T times. Denote by $\mathcal{A}_i := \{t \leq T : t \in \mathcal{N}, |\{u \leq T : \phi(u) = t\}| = i\}$ times which have exactly i children for $i \in \{0, 1, 2\}$. Note that no time can have more than 2 children. In particular \mathcal{A}_0 is the set of leaves. Then, by summing the above equations we obtain

$$\begin{aligned} \sum_{t=1}^T \bar{\ell}(\hat{Y}_t(\epsilon), Y_t) &\leq \sum_{i=1}^{p(T)} \left(\min_{y \in \mathcal{Y}_\epsilon} \sum_{u \in \mathcal{C}_i(T)} \ell(y, Y_u) + \epsilon \max(T_\epsilon, |\mathcal{C}_i(T)|) \right) \\ &\leq \sum_{i=1}^{p(T)} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u \in \mathcal{C}_i(T)} \ell(y, Y_u) + \epsilon T + T_\epsilon |\{1 \leq i \leq p : |\mathcal{C}_i(T)| < T_\epsilon\}| \\ &\leq \sum_{i=1}^{p(T)} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u \in \mathcal{C}_i(T)} \ell(y, Y_u) + \epsilon T + T_\epsilon |\mathcal{A}_0| + \epsilon T_\epsilon, \end{aligned}$$

where in the last inequality we used the fact that all clusters with $|\mathcal{C}_i(T)| < T_\epsilon$ contain a leaf from \mathcal{A}_0 , which is therefore distinct for each such cluster. Now note that by counting the number of edges of the tree structure we obtain $\frac{1}{2}(3|\mathcal{A}_2| + 2|\mathcal{A}_1| + |\mathcal{A}_0| - 1) = T - 1 = |\mathcal{A}_0| + |\mathcal{A}_1| + |\mathcal{A}_2| - 1$, where the -1 on the left-hand side accounts for the root of this tree which does not have a parent. Hence we obtain $|\mathcal{A}_0| = |\mathcal{A}_2| + 1$. Further, $|\mathcal{A}_2| \leq |\{t \leq T : U_t = 1\}|$ which follows a binomial distribution $\mathcal{B}(T, \delta_\epsilon)$. Therefore, using the Chernoff bound, with probability $1 - e^{-T\delta_\epsilon/3}$ we have

$$\begin{aligned} \sum_{t=1}^T \bar{\ell}(\hat{Y}_t(\epsilon), Y_t) &\leq \sum_{i=1}^{p(T)} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u \in \mathcal{C}_i(T)} \ell(y, Y_u) + 2\epsilon T + T_\epsilon(1 + 2T\delta_\epsilon) \\ &\leq \sum_{i=1}^{p(T)} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u \in \mathcal{C}_i(T)} \ell(y, Y_u) + T_\epsilon + 3\epsilon T. \end{aligned}$$

We now observe that the sequence $\{\ell(\hat{Y}_t(\epsilon), Y_t) - \bar{\ell}(\hat{Y}_t(\epsilon), Y_t)\}_{T \geq 1}$ is a sequence of martingale differences bounded by $\bar{\ell}$ in absolute value. Hence, the Hoeffding-Azuma inequality yields that for any $T \geq 1$, with

probability $1 - \frac{1}{T^2} - e^{-T\delta_\epsilon/3}$,

$$\sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) \leq \sum_{i=1}^{p(T)} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u \in \mathcal{C}_i(T)} \ell(y, Y_u) + T_\epsilon + 3\epsilon T + 2\bar{\ell}\sqrt{T \ln T}.$$

Because $\sum_{T \geq 1} \frac{1}{T^2} + e^{-T\delta_\epsilon/3} < \infty$ the Borel-Cantelli lemma implies that with probability one, there exists a time \hat{T} such that

$$\forall T \geq \hat{T}, \quad \sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) \leq \sum_{i=1}^{p(T)} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u \in \mathcal{C}_i(T)} \ell(y, Y_u) + T_\epsilon + 2\bar{\ell}\sqrt{T \ln T} + 3\epsilon T.$$

We denote by \mathcal{E}_ϵ this event. We are now ready to analyze the risk of the learning rule f^ϵ . Let $f : \mathcal{X} \rightarrow \mathcal{Y}$ a measurable function to which we compare the prediction of f^ϵ . By Theorem 4.1, the rule $(1 + \delta_\epsilon)\text{C1NN}$ is optimistically universal in the noiseless setting. Therefore, because $\mathbb{X} \in \text{SOUL}$ we have in particular

$$\frac{1}{T} \sum_{t=1}^T \ell((1 + \delta_\epsilon)\text{C1NN}_t(\mathbb{X}_{\leq t-1}, f(\mathbb{X}_{\leq t-1}), X_t), f(X_t)) \rightarrow 0 \quad (a.s.),$$

i.e., almost surely, $\frac{1}{T} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_{\phi(t)}), f(X_t)) \rightarrow 0$ — the times corresponding to duplicate instances incur a 0 loss by memorization. We denote by \mathcal{F}_ϵ this event of probability one. Using Lemma A.1, we write for any $u = 1, \dots, T_\epsilon - 1$,

$$\begin{aligned} & \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_{\phi^u(t)}), f(X_t)) \\ & \leq 2^{\alpha-1} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_{\phi^{u-1}(t)}), f(X_t)) + 2^{\alpha-1} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_{\phi^1(t)}), f(X_{\phi^{u-1}(t)})) \\ & \leq 2^{\alpha-1} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_{\phi^{u-1}(t)}), f(X_t)) \\ & \quad + 2^{\alpha-1} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_{\phi(t)}), f(X_t)) \cdot |\{l \leq T : \phi^{u-1}(l) = t\}| \\ & \leq 2^{\alpha-1} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_{\phi^{u-1}(t)}), f(X_t)) + 2^{\alpha+u-2} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_{\phi(t)}), f(X_t)) \end{aligned}$$

where we used the fact that times have at most 2 children. Therefore, iterating the above equations, we obtain that on \mathcal{F}_ϵ , for any $u = 1, \dots, T_\epsilon - 1$

$$\begin{aligned} \frac{1}{T} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_{\phi^u(t)}), f(X_t)) & \leq \left(\sum_{k=1}^u 2^{\alpha+k-2+(\alpha-1)(u-k)} \right) \frac{1}{T} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_{\phi(t)}), f(X_t)) \\ & \leq \frac{2^{u\alpha}}{T} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_{\phi(t)}), f(X_t)) \rightarrow 0. \end{aligned}$$

In the rest of the proof, for any $y \in \mathcal{Y}$, we will denote by y^ϵ a value in the ϵ -net \mathcal{Y}_ϵ such that $\ell(y, y^\epsilon) \leq \epsilon$. We now pose $\mu_\epsilon = \min\{0 < \mu \leq 1 : c_\mu^\alpha \leq \frac{1}{\sqrt{\epsilon}}\}$ if the corresponding set is non-empty and $\mu_\epsilon = 1$ otherwise. Note that because c_μ^α is non-increasing in μ , we have $\mu_\epsilon \rightarrow_{\epsilon \rightarrow 0} 0$. Now let $0 < \mu \leq 1$. $\mu := \epsilon^{\frac{1}{\alpha+1}}$. Finally, for any cluster $\mathcal{C}_i(T)$, let $t_i = \min\{u \in \mathcal{C}_i(T)\}$. Putting everything together, on the event $\mathcal{E}_\epsilon \cap \mathcal{F}_\epsilon$, for any

$T \geq \hat{T}$, we have

$$\begin{aligned}
\sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) &\leq \sum_{i=1}^{p(T)} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u \in \mathcal{C}_i(T)} \ell(y, Y_u) + T_\epsilon + 2\bar{\ell}\sqrt{T \ln T} + 3\epsilon T \\
&\leq \sum_{i=1}^{p(T)} \sum_{u \in \mathcal{C}_i(T)} \ell(f(X_{t_i})^\epsilon, Y_u) + T_\epsilon \bar{\ell} + 2\bar{\ell}\sqrt{T \ln T} + 3\epsilon T \\
&\leq \sum_{i=1}^{p(T)} \sum_{u \in \mathcal{C}_i(T)} [c_{\mu_\epsilon}^\alpha \ell(f(X_{t_i})^\epsilon, f(X_{t_i})) + (c_{\mu_\epsilon}^\alpha)^2 \ell(f(X_{t_i}), f(X_u)) \\
&\quad + (1 + \mu_\epsilon)^2 \ell(f(X_u), Y_u)] + T_\epsilon \bar{\ell} + 2\bar{\ell}\sqrt{T \ln T} + 3\epsilon T \\
&\leq (1 + \mu_\epsilon)^2 \sum_{t=1}^T \ell(f(X_t), Y_t) + (c_{\mu_\epsilon}^\alpha)^2 \frac{T_\epsilon}{\epsilon} \sum_{u=1}^{T_\epsilon-1} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_t), f(X_{\phi^u(t)})) \\
&\quad + T_\epsilon \bar{\ell} + 2\bar{\ell}\sqrt{T \ln T} + (3 + c_{\mu_\epsilon}^\alpha) \epsilon T \\
&\leq \sum_{t=1}^T \ell(f(X_t), Y_t) + \frac{(c_{\mu_\epsilon}^\alpha)^2 T_\epsilon}{\epsilon} \sum_{u=1}^{T_\epsilon-1} \sum_{t \leq T, t \in \mathcal{N}} \ell(f(X_t), f(X_{\phi^u(t)})) \\
&\quad + T_\epsilon \bar{\ell} + 2\bar{\ell}\sqrt{T \ln T} + (3\epsilon + \epsilon c_{\mu_\epsilon}^\alpha + 3\mu_\epsilon) T,
\end{aligned}$$

where in the third inequality we used Lemma A.1 twice, and in the fourth inequality we used the fact that clusters containing distinct instances have at most $\frac{T_\epsilon}{\epsilon}$ duplicates of each instance. Hence, for any $\epsilon < (c_1^\alpha)^{-2}$, on the event $\mathcal{E}_\epsilon \cap \mathcal{F}_\epsilon$, we obtain

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) - \ell(f(X_t), Y_t) \leq 3\epsilon + \epsilon c_{\mu_\epsilon}^\alpha + 3\mu_\epsilon \leq 3\epsilon + \sqrt{\epsilon} + 3\mu_\epsilon,$$

where $\mu_\epsilon \rightarrow_{\epsilon \rightarrow 0} 0$. We now denote $\delta_\epsilon := 2\epsilon + \sqrt{\epsilon} + 3\mu_\epsilon$ and $i_0 = \lceil \frac{2 \ln c_1^\alpha}{\ln 2} \rceil$. We now turn to the final learning rule and show that by using the predictions of the rules f^{ϵ_i} for $i \geq 0$, it achieves zero risk. First, by the union bound, on the event $\bigcap_{i \geq 0} \mathcal{E}_{\epsilon_i} \cap \mathcal{F}_{\epsilon_i}$ of probability one,

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t(\epsilon_i), Y_t) - \ell(f(X_t), Y_t) \leq \delta_{\epsilon_i}, \quad \forall i \geq i_0.$$

Now define \mathcal{H} the event probability one according to Lemma 4.2 such that there exists \hat{t} for which

$$\forall t \geq \hat{t}, \forall i \in I_t, \quad \sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) \leq \sum_{s=t_i}^t \ell(\hat{Y}_s(\epsilon_i), Y_s) + (2 + \bar{\ell} + \bar{\ell}^2) \sqrt{t \ln t}.$$

In the rest of the proof we will suppose that the event $\mathcal{H} \cap \bigcap_{i \geq 0} \mathcal{E}_{\epsilon_i} \cap \mathcal{F}_{\epsilon_i}$ is met. Let $i \geq i_0$. For any $T \geq \max(\hat{t}, t_i)$, we have

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) &\leq \frac{t_i}{T} \bar{\ell} + \frac{1}{T} \sum_{t=t_i}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \\
&\leq \frac{t_i}{T} \bar{\ell} + \frac{1}{T} \sum_{t=t_i}^T \ell(\hat{Y}_t(\epsilon_i), Y_t) - \ell(f(X_t), Y_t) + (2 + \bar{\ell} + \bar{\ell}^2) \sqrt{\frac{\ln T}{T}} \\
&\leq \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t(\epsilon_i), Y_t) - \ell(f(X_t), Y_t) + \frac{2t_i}{T} \bar{\ell} + (2 + \bar{\ell} + \bar{\ell}^2) \sqrt{\frac{\ln T}{T}}.
\end{aligned}$$

Therefore we obtain $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq \delta_{\epsilon_i}$. Because this holds for any $i \geq i_0$ on the event $\mathcal{H} \cap \bigcap_{i \geq 0} \mathcal{E}_{\epsilon_i} \cap \mathcal{F}_{\epsilon_i}$ of probability one, and $\delta_{\epsilon_i} \rightarrow 0$ for $i \rightarrow \infty$, we have

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq 0.$$

This ends the proof of the theorem.

B.3 Proof of Lemma 4.2

We first introduce the following helper lemma which can be found in [CL06].

Lemma B.4 ([CL06]). *For all $N \geq 2$, for all $\beta \geq \alpha \geq 0$ and for all $d_1, \dots, d_N \geq 0$ such that $\sum_{i=1}^N e^{-\alpha d_i} \geq 1$,*

$$\ln \frac{\sum_{i=1}^N e^{-\alpha d_i}}{\sum_{i=1}^N e^{-\beta d_i}} \leq \frac{\beta - \alpha}{\alpha} \ln N.$$

We are now ready to compare the predictions of the learning rule f to the predictions of the rules f^ϵ .

For any $t \geq 0$, we define the instantaneous regret $r_{t,i} = \hat{\ell}_t - \ell(\hat{Y}_t(\epsilon_i), Y_t)$. We first note that $|r_{t,i}| \leq \bar{\ell}$. We now define $w'_{t-1,i} := e^{\eta_{t-1}(\hat{L}_{t-1,i} - L_{t-1,i})}$. We also introduce $W_{t-1} = \sum_{i \in I_t} w_{t-1,i}$ and $W'_{t-1} = \sum_{i \in I_{t-1}} w'_{t-1,i}$. We denote the index $k_t \in I_t$ such that $\hat{L}_{t,k_t} - L_{t,k_t} = \max_{i \in I_t} \hat{L}_{t,i} - L_{t,i}$. Then we write

$$\begin{aligned} \frac{1}{\eta_t} \ln \frac{w_{t-1,k_{t-1}}}{W_{t-1}} - \frac{1}{\eta_{t+1}} \ln \frac{w_{t,k_t}}{W_t} &= \left(\frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \ln \frac{W_t}{w_{t,k_t}} + \frac{1}{\eta_t} \ln \frac{W_t/w_{t,k_t}}{W'_t/w'_{t,k_t}} \\ &\quad + \frac{1}{\eta_t} \ln \frac{w_{t-1,k_{t-1}}}{w'_{t,k_t}} + \frac{1}{\eta_t} \ln \frac{W'_t}{W_{t-1}}. \end{aligned}$$

By construction, we have $\ln \frac{W_t}{w_{t,k_t}} \leq \ln |I_t| \leq \ln(1 + \ln t)$. Further, we have that

$$\begin{aligned} \frac{1}{\eta_t} \ln \frac{W_t/w_{t,k_t}}{W'_t/w'_{t,k_t}} &= \frac{1}{\eta_t} \ln \frac{\sum_{i \in I_{t+1}} e^{\eta_{t+1}(\hat{L}_{t,i} - L_{t,i} - \hat{L}_{t,k_t} + L_{t,k_t})}}{\sum_{i \in I_t} e^{\eta_t(\hat{L}_{t,i} - L_{t,i} - \hat{L}_{t,k_t} + L_{t,k_t})}} \\ &= \frac{1}{\eta_t} \ln \frac{\sum_{i \in I_{t+1}} w_{t,i}}{\sum_{i \in I_t} w_{t,i}} + \frac{1}{\eta_t} \ln \frac{\sum_{i \in I_{t+1}} e^{\eta_{t+1}(\hat{L}_{t,i} - L_{t,i} - \hat{L}_{t,k_t} + L_{t,k_t})}}{\sum_{i \in I_{t+1}} e^{\eta_t(\hat{L}_{t,i} - L_{t,i} - \hat{L}_{t,k_t} + L_{t,k_t})}} \\ &\leq \frac{1}{\eta_t} \ln \frac{\sum_{i \in I_{t+1}} w_{t,i}}{\sum_{i \in I_t} w_{t,i}} + \frac{1}{\eta_t} \left(\frac{\eta_t - \eta_{t+1}}{\eta_{t+1}} \right) \ln |I_{t+1}| \\ &\leq \frac{|I_{t+1}| - |I_t|}{\eta_t \sum_{i \in I_t} w_{t,i}} + \left(\frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \ln(1 + \ln(t+1)), \end{aligned}$$

where in the first inequality we applied Lemma B.4. We also have

$$\frac{1}{\eta_t} \ln \frac{w_{t-1,k_{t-1}}}{w'_{t,k_t}} = (\hat{L}_{t-1,k_{t-1}} - L_{t-1,k_{t-1}}) - (\hat{L}_{t,k_t} - L_{t,k_t}).$$

Last, because $|r_{t,i}| \leq \bar{\ell}$ for all $i \in I_t$, we can use Hoeffding's lemma to obtain

$$\frac{1}{\eta_t} \ln \frac{W'_t}{W_{t-1}} = \frac{1}{\eta_t} \ln \sum_{i \in I_t} \frac{w_{t-1,i}}{W_{t-1}} e^{\eta_t r_{t,i}} \leq \frac{1}{\eta_t} \left(\eta_t \sum_{i \in I_t} r_{t,i} \frac{w_{t-1,i}}{W_{t-1}} + \frac{\eta_t^2 (2\bar{\ell})^2}{8} \right) = \frac{1}{2} \eta_t \bar{\ell}^2.$$

Putting everything together gives

$$\begin{aligned} \frac{1}{\eta_t} \ln \frac{w_{t-1, k_{t-1}}}{W_{t-1}} - \frac{1}{\eta_{t+1}} \ln \frac{w_{t, k_t}}{W_t} &\leq 2 \left(\frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \ln(1 + \ln(t+1)) + \frac{|I_{t+1}| - |I_t|}{\eta_t \sum_{i \in I_t} w_{t,i}} \\ &\quad + (\hat{L}_{t-1, k_{t-1}} - L_{t-1, k_{t-1}}) - (\hat{L}_{t, k_t} - L_{t, k_t}) + \frac{1}{2} \eta_t \bar{\ell}^2. \end{aligned} \quad (1)$$

First suppose that we have $\sum_{i \in I_t} w_{t,i} \leq 1$. Then either $k_t \in I_{t+1} \setminus I_t$ in which case $\hat{L}_{t, k_t} - L_{t, k_t} = 0$, or we have directly

$$\hat{L}_{t, k_t} - L_{t, k_t} \leq \frac{1}{\eta_{t+1}} \ln \left[\sum_{i \in I_t} w_{t,i} \right] \leq 0.$$

Otherwise, let $t' = \min\{1 \leq s \leq t : \forall s \leq s' \leq t, \sum_{i \in I_{s'}} w_{s',i} \geq 1\}$. We sum equation (1) for $s = t', \dots, t$ which gives

$$\begin{aligned} \frac{1}{\eta_1} \ln \frac{w_{t'-1, k_{t'-1}}}{W_{t'-1}} - \frac{1}{\eta_{t+1}} \ln \frac{w_{t, k_t}}{W_t} &\leq \frac{2}{\eta_{t+1}} \ln(1 + \ln(t+1)) + \frac{|I_{t+1}|}{\eta_t} \\ &\quad + (\hat{L}_{t'-1, k_{t'-1}} - L_{t'-1, k_{t'-1}}) - (\hat{L}_{t, k_t} - L_{t, k_t}) + \frac{\bar{\ell}^2}{2} \sum_{s=t'}^t \eta_s. \end{aligned}$$

Note that we have $\frac{w_{t, k_t}}{W_t} \leq 1$ and $\frac{w_{t'-1, k_{t'-1}}}{W_{t'-1}} \geq \frac{1}{|I_{t'-1}|} \geq \frac{1}{1 + \ln t}$. Also, assuming $t' \geq 2$, since $\sum_{i \in I_{t'-1}} w_{t'-1, i} < 1$, we have for any $i \in I_{t'-1}$ that $\hat{L}_{t'-1, i} - L_{t'-1, i} \leq 0$, hence $\hat{L}_{t'-1, k_{t'-1}} - L_{t'-1, k_{t'-1}} \leq 0$. If $t' = 1$ we have directly $\hat{L}_{0, k_0} - L_{0, k_0} = 0$. Finally, using the fact that $\sum_{s=1}^t \frac{1}{\sqrt{s}} \leq 2\sqrt{t}$, we obtain

$$\begin{aligned} \hat{L}_{t, k_t} - L_{t, k_t} &\leq \ln(1 + \ln(t+1)) \left(1 + 2\sqrt{\frac{t+1}{\ln(t+1)}} \right) + (1 + \ln(t+1)) \sqrt{\frac{t}{\ln t}} + \bar{\ell}^2 \sqrt{t \ln t} \\ &\leq (3/2 + \bar{\ell}^2) \sqrt{t \ln t}, \end{aligned}$$

for all $t \geq t_0$ where t_0 is a fixed constant. This in turn implies that for all $t \geq t_0$ and $i \in I_t$, we have $\hat{L}_{t, i} - L_{t, i} \leq (3/2 + \bar{\ell}^2) \sqrt{t \ln t}$. Now note that $|\ell(\hat{Y}_t, Y_t) - \hat{\ell}_t| \leq \bar{\ell}$. Hence, we can use Hoeffding-Azuma inequality for the variables $\ell(\hat{Y}_s, Y_s) - \hat{\ell}_s$ that form a sequence of martingale differences to obtain $\mathbb{P} \left[\sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) > \hat{L}_{t, i} + u \right] \leq e^{-\frac{2u^2}{t\bar{\ell}^2}}$. Hence, for $t \geq t_0$ and $i \in I_t$, with probability $1 - \delta$, we have

$$\sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) \leq \hat{L}_{t, i} + \bar{\ell} \sqrt{\frac{t}{2} \ln \frac{1}{\delta}} \leq L_{t, i} + (3/2 + \bar{\ell}^2) \sqrt{t \ln t} + \bar{\ell} \sqrt{\frac{t}{2} \ln \frac{1}{\delta}}.$$

Therefore, since $|I_t| \leq 1 + \ln t$, by union bound with probability $1 - \frac{1}{t^2}$ we obtain that for all $i \in I_t$,

$$\sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) \leq L_{t, i} + (3/2 + \bar{\ell}^2) \sqrt{t \ln t} + \bar{\ell} \sqrt{\frac{t}{2} \ln(1 + \ln t)} + \bar{\ell} \sqrt{t \ln t} \leq (2 + \bar{\ell} + \bar{\ell}^2) \sqrt{t \ln t},$$

for all $t \geq t_1$ where $t_1 \geq t_0$ is a fixed constant. The Borel-Cantelli lemma implies that almost surely, there exists $\hat{t} \geq 0$ such that

$$\forall t \geq \hat{t}, \forall i \in I_t, \quad \sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) \leq L_{t, i} + (2 + \bar{\ell} + \bar{\ell}^2) \sqrt{t \ln t}.$$

This ends the proof of the lemma.

C Proofs of Section 5

C.1 Proof of Theorem 5.1

We start by checking that with the defined loss (\mathbb{N}, ℓ) is indeed a metric space (\mathbb{N}, ℓ) . We only have to check that the triangular inequality is satisfied, the other properties of a metric being directly satisfied. By construction, the loss has values in $\{0, \frac{1}{2}, 1\}$. Now let $i, j, k \in \mathbb{N}$. The triangular inequality $\ell(i, j) \leq \ell(i, k) + \ell(k, j)$ is directly satisfied if two of these indices are equal. Therefore, we can suppose that they are all distinct and as a result $\ell(i, j), \ell(i, k), \ell(k, j) \in \{\frac{1}{2}, 1\}$. Therefore

$$\ell(i, j) \leq 1 \leq \ell(i, k) + \ell(k, j),$$

which ends the proof that ℓ is a metric.

Now let $(\mathcal{X}, \mathcal{B})$ be a separable metrizable Borel space. Let $\mathbb{X} \notin \text{CS}$. We aim to show that universal online learning under adversarial responses is not achievable under \mathbb{X} . Because $\mathbb{X} \notin \text{CS}$, there exists a sequence of decreasing measurable sets $\{A_i\}_{i \geq 1}$ with $A_i \downarrow \emptyset$ such that $\mathbb{E}[\hat{\mu}_{\mathbb{X}}(A_i)]$ does not converge to 0 for $i \rightarrow \infty$. In particular, there exist $\epsilon > 0$ and an increasing subsequence $(i_l)_{l \geq 1}$ such that $\mathbb{E}[\hat{\mu}_{\mathbb{X}}(A_{i_l})] \geq \epsilon$ for all $l \geq 1$. We now denote $B_l := A_{i_l} \setminus A_{i_{l+1}}$ for any $l \geq 1$. Then $\{B_l\}_{l \geq 1}$ forms a sequence of disjoint measurable sets such that

$$\mathbb{E} \left[\hat{\mu}_{\mathbb{X}} \left(\bigcup_{l' \geq l} B_{l'} \right) \right] \geq \epsilon, \quad l \geq 1.$$

Therefore, for any $l \geq 1$ because $\mathbb{E} \left[\hat{\mu}_{\mathbb{X}} \left(\bigcup_{l' \geq l} B_{l'} \right) \right] \leq \mathbb{P} \left[\hat{\mu}_{\mathbb{X}} \left(\bigcup_{l' \geq l} B_{l'} \right) \geq \frac{\epsilon}{2} \right] + \frac{\epsilon}{2}$ we obtain

$$\mathbb{P} \left[\hat{\mu}_{\mathbb{X}} \left(\bigcup_{l' \geq l} B_{l'} \right) \geq \frac{\epsilon}{2} \right] \geq \frac{\epsilon}{2}.$$

Now because $\hat{\mu}$ is increasing we obtain

$$\begin{aligned} \mathbb{P} \left[\hat{\mu}_{\mathbb{X}} \left(\bigcup_{l' \geq l} B_{l'} \right) \geq \frac{\epsilon}{2}, \forall l \geq 1 \right] &= \lim_{L \rightarrow \infty} \mathbb{P} \left[\hat{\mu}_{\mathbb{X}} \left(\bigcup_{l' \geq l} B_{l'} \right) \geq \frac{\epsilon}{2}, 1 \leq l \leq L \right] \\ &= \lim_{L \rightarrow \infty} \mathbb{P} \left[\hat{\mu}_{\mathbb{X}} \left(\bigcup_{l' \geq L} B_{l'} \right) \geq \frac{\epsilon}{2} \right] \geq \frac{\epsilon}{2}. \end{aligned}$$

We will denote by \mathcal{A} this event in which for all $l \geq 1$, we have $\hat{\mu}_{\mathbb{X}} \left(\bigcup_{l' \geq l} B_{l'} \right) \geq \frac{\epsilon}{2}$. Under the event \mathcal{A} , for any $l, t^0 \geq 1$, there always exists $t^1 > t^0$ such that $\frac{1}{t^1} \sum_{t=1}^{t^1} \mathbb{1}_{\bigcup_{l' \geq l} B_{l'}}(X_t) \geq \frac{3\epsilon}{8}$. We construct a sequence of times $(t_p)_{p \geq 1}$ and indices $(l_p)_{p \geq 1}, (u_p)_{p \geq 1}$ by induction as follows. We first pose $u_0 = t_0 = 0$. Now assume that for $p \geq 1$, the time t_{p-1} and index u_{p-1} are defined. We first construct an index $l_p > u_{p-1}$ such that

$$\mathbb{P} \left[\mathbb{X}_{\leq t_{p-1}} \cap \left(\bigcup_{l \geq l_p} B_l \right) \neq \emptyset \right] \leq \frac{\epsilon}{2^{p+3}}.$$

We will denote by \mathcal{E}_p the complementary of this event. Note that finding such index l_p is possible because the considered events $\{\mathbb{X}_{\leq t_{p-1}} \cap \left(\bigcup_{l' \geq l} B_{l'} \right) \neq \emptyset\}$ are decreasing as $l > u_{p-1}$ increases and we have $\bigcap_{l > u_{p-1}} \left\{ \mathbb{X}_{\leq t_{p-1}} \cap \left(\bigcup_{l' \geq l} B_{l'} \right) \neq \emptyset \right\} = \left\{ \mathbb{X}_{\leq t_{p-1}} \cap \left(\bigcap_{l > u_{p-1}} \bigcup_{l' \geq l} B_{l'} \right) \neq \emptyset \right\} = \emptyset$. We then construct $t_p > t_{p-1}$ such that

$$\mathbb{P} \left[\mathcal{A}^c \cup \bigcup_{t_{p-1} < t \leq t_p} \left\{ \frac{1}{t} \sum_{u=1}^t \mathbb{1}_{\bigcup_{l \geq l_p} B_l}(X_u) \geq \frac{3\epsilon}{8} \right\} \right] \geq 1 - \frac{\epsilon}{2^{p+4}}.$$

This is also possible because $\mathcal{A} \subset \bigcup_{t > \frac{8}{\epsilon} t_{p-1}} \left\{ \frac{1}{t} \sum_{u=1}^t \mathbb{1}_{\bigcup_{l \geq t_p} B_l}(X_u) \geq \frac{3\epsilon}{8} \right\}$. Last, we can now construct $u_p \geq l_p$ such that

$$\mathbb{P} \left[\mathcal{A}^c \cup \bigcup_{t_{p-1} < t \leq t_p} \left\{ \frac{1}{t} \sum_{u=1}^t \mathbb{1}_{\bigcup_{l_p \leq l \leq u_p} B_l}(X_u) \geq \frac{\epsilon}{4} \right\} \right] \geq 1 - \frac{\epsilon}{2^{p+3}},$$

which is possible using similar arguments as above. We denote \mathcal{F}_p this event. This ends the recursive construction of times t_p and indices l_p for all $p \geq 1$. Note that by construction, $\mathbb{P}[\mathcal{E}_p^c], \mathbb{P}[\mathcal{F}_p^c] \leq \frac{\epsilon}{2^{p+3}}$. Hence, by union bound, the event $\mathcal{A} \cap \bigcap_{p \geq 1} (\mathcal{E}_p \cap \mathcal{F}_p)$ has probability $\mathbb{P}[\mathcal{A} \cap \bigcap_{p \geq 1} (\mathcal{E}_p \cap \mathcal{F}_p)] \geq \mathbb{P}[\mathcal{A}] - \frac{\epsilon}{4} \geq \frac{\epsilon}{4}$. To simplify the rest of the proof, we denote $\tilde{B}_p = \bigcup_{l_p \leq l \leq u_p} B_l$ for any $p \geq 1$. Also, for any $t_1 \leq t_2$, we denote by

$$N_p(t_1, t_2) = \sum_{t=t_1}^{t_2} \mathbb{1}_{\tilde{B}_p}(X_t)$$

the number of times that set \tilde{B}_p has been visited between times t_1 and t_2 .

We now fix a learning rule f and construct a process \mathbb{Y} for which consistency will not be achieved on the event $\mathcal{A} \cap \bigcap_{p \geq 1} (\mathcal{E}_p \cap \mathcal{F}_p)$. Precisely, we first construct a family of processes \mathbb{Y}^b indexed by a sequence of binary digits $b = (b_i)_{i \geq 1}$. The process \mathbb{Y}^b is defined such that for any $p \geq 1$, and for all $t_{p-1} < t \leq t_p$,

$$Y_t^b := \begin{cases} n_{t_p} + 4u_p(t) + 2b_{i(p, u_p(t))} + b_{i(p, u_p(t))+1} & \text{if } X_t \in \tilde{B}_p, \\ n_{t_{p'}} + 4t_{p'} + \{b_{i(p', t_{p'}-1)} \cdots b_{i(p', 1)} b_{i(p', 0)}\} 2 & \text{if } X_t \in \tilde{B}_{p'}, p' < p, \\ 0 & \text{otherwise,} \end{cases}$$

where we denoted $u_p(t) = N_p(t_{p-1} + 1, t - 1)$ and posed for any $p \geq 1$ and $u \geq 1$:

$$i(p, u) = 2 \sum_{p' < p} t_{p'} + 2u.$$

Note in particular that conditionally on \mathbb{X} , \mathbb{Y}^b is deterministic: it does not depends on the random predictions of the learning rule. Because we always have $N_p(t_{p-1} + 1, t - 1) \leq t_p$ for any $t \leq t_p$, the process is designed so that we have $Y_t^b \in I_{t_p}$ if $X_t \in \tilde{B}_p$ and $t_{p-1} < t \leq t_p$. Further, for $t_{p-1} < t \leq t_p$, if $X_t \in \bigcup_{p' < p} \tilde{B}_{p'}$ then $Y_t^b \in J_{t_{p'}}$. We now consider an i.i.d. Bernoulli $\mathcal{B}(\frac{1}{2})$ sequence of random bits \mathbf{b} independent from the process \mathbb{X} —and any learning rule predictions. We analyze the responses of the learning rule for responses \mathbb{Y}^b . We first fix a realization \mathbf{x} of the process \mathbb{X} , which falls in the event $\mathcal{A} \cap \bigcap_{p \geq 1} (\mathcal{E}_p \cap \mathcal{F}_p)$. For any $p \geq 1$ we define $\mathcal{T}_p := \{t_{p-1} < t \leq t_p : x_t \in \tilde{B}_p\}$. For simplicity of notation, for any $t \in \mathcal{T}_p$ we denote $i(t) = i(p, u_p(t))$. We will also denote $\hat{Y}_t := f_t(\mathbf{x}_{<t}, \mathbb{Y}_{<t}^b, x_t)$. Last, denote by r_t the possible randomness used by the learning rule f_t at time t . For any $t \in \mathcal{T}_p$, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{b}, r} \ell(\hat{Y}_t, Y_t^b) &= \mathbb{E}_{\{b_{i(p', u')}, b_{i(p', u')+1}, p' \leq p, u' \leq t_{p'}\} \cup \{r_{t'}, t' \leq t\}} \ell(\hat{Y}_t, Y_t^b) \\ &= \mathbb{E} \left[\mathbb{E}_{b_{i(t)}, b_{i(t)+1}} \ell(\hat{Y}_t, Y_t^b) \Big| b_{i(t')}, b_{i(t')+1}, t' < t, t' \in \mathcal{T}_p; b_i, i < i(p, 0); r_{t'}, t' \leq t \right] \\ &= \mathbb{E} \left[\mathbb{E}_{b_{i(t)}, b_{i(t)+1}} \ell(\hat{Y}_t, Y_t^b) \Big| \hat{Y}_t \right] \\ &= \mathbb{E}_{\hat{Y}_t} \left[\frac{1}{4} \sum_{m=0}^3 \ell(\hat{Y}_t, n_{t_p} + 4u_p(t) + m) \right] \\ &= \mathbb{E}_{\hat{Y}_t} \left[\mathbb{1}_{\hat{Y}_t \notin \{n_{t_p} + 4u_p(t) + m, 0 \leq m \leq 3\}} \cup J_{t_p} + \frac{3}{4} \mathbb{1}_{\hat{Y}_t \in \{n_{t_p} + 4u_p(t) + m, 0 \leq m \leq 3\}} + \frac{3}{4} \mathbb{1}_{\hat{Y}_t \in J_{t_p}} \right] \\ &\geq \frac{3}{4}. \end{aligned}$$

where in the last equality, we used the fact that if $j \in J_{k(t)}$ then by construction $\ell(j, n_{t_p} + 4u_p(t)) = \ell(j, n_{t_p} + 4u_p(t) + 1)$, $\ell(j, n_{t_p} + 4u_p(t) + 2) = \ell(j, n_{t_p} + 4u_p(t) + 3)$, and $\{\ell(j, n_{t_p} + 4u_p(t)), \ell(j, n_{t_p} + 4u_p(t) + 2)\} = \{\frac{1}{2}, 1\}$. Summing all equations, we obtain for any $t_{p-1} < T \leq t_p$,

$$\mathbb{E}_{\mathbf{b}, \mathbf{r}} \left[\sum_{t=1}^T \ell(f_t(\mathbf{x}_{<t}, \mathbb{Y}_{<t}^{\mathbf{b}}, x_t), Y_t^{\mathbf{b}}) \right] \geq \frac{3}{4} \sum_{p' < p} |\mathcal{T}_{p'}| + \frac{3}{4} |\mathcal{T}_p \cap \{t \leq T\}|.$$

This holds for all $p \geq 1$. Let us now compare this loss to the best prediction of a fixed measurable function. Specifically, for any binary sequence b , we consider the following function $f^b : \mathcal{X} \rightarrow \mathbb{N}$:

$$f^b(x) = \begin{cases} n_{t_p} + 4t_p + \{b_{i(p, t_p-1)} \dots b_{i(p, 1)} b_{i(p, 0)}\}_2 & \text{if } x \in \tilde{B}_p \\ 0 & \text{if } x \notin \bigcup_{p \geq 1} \tilde{B}_p. \end{cases}$$

Now let $t_{p-1} < t \leq t_p$ and $p \geq 1$. If $x_t \in \bigcup_{p' < p} \tilde{B}_{p'}$ we have $f^b(x_t) = Y_t^{\mathbf{b}}$, hence $\ell(f^b(x_t), Y_t^{\mathbf{b}}) = 0$. Now if $X_t \in \tilde{B}_p$ by construction we have $\ell(f^b(x_t), Y_t^{\mathbf{b}}) = \frac{1}{2}$. Finally, observe that because the event \mathcal{E}_{p+1} is satisfied by \mathbf{x} there does not exist $t_{p-1} < t \leq t_p$ such that $t \in \bigcup_{p' > p} \tilde{B}_{p'} \subset \bigcup_{l \geq l_{p+1}} B_l$. As a result, we have $\ell(f^b(x_t), Y_t^{\mathbf{b}}) = \frac{1}{2} \mathbb{1}_{t \in \mathcal{T}_p}$ for any $t_{p-1} < t \leq t_p$. Thus, we obtain for any $t_{p-1} < T \leq t_p$,

$$\mathbb{E}_{\mathbf{b}, \mathbf{r}} \left[\sum_{t=1}^T \ell(\hat{Y}_t, Y_t^{\mathbf{b}}) - \ell(f^b(X_t), Y_t^{\mathbf{b}}) \right] \geq \frac{1}{4} \sum_{p' \leq p} |\mathcal{T}_{p'}| + \frac{1}{4} |\mathcal{T}_p \cap \{t \leq T\}| \geq \frac{1}{4} |\mathcal{T}_p \cap \{t \leq T\}|.$$

Recall that the event \mathcal{F}_p is satisfied by \mathbf{x} for any $p \geq 1$. Therefore, there exists a time $t_{p-1} < T_p \leq t_p$ such that $\sum_{t=1}^{T_p} \mathbb{1}_{\tilde{B}_p}(x_t) \geq \frac{\epsilon T_p}{4}$. Then, note that because the event \mathcal{E}_p is satisfied, we have $\sum_{t=1}^{t_{p-1}} \mathbb{1}_{\tilde{B}_p}(x_t) = 0$. Therefore, we obtain $|\mathcal{T}_p \cap \{t \leq T_p\}| \geq \frac{\epsilon T_p}{4}$, and as a result,

$$\mathbb{E}_{\mathbf{b}, \mathbf{r}} \left[\frac{1}{T_p} \sum_{t=1}^{T_p} \ell(\hat{Y}_t, Y_t^{\mathbf{b}}) - \ell(f^b(X_t), Y_t^{\mathbf{b}}) \right] \geq \frac{\epsilon}{16}.$$

Because this holds for any $p \geq 1$ and as $p \rightarrow \infty$ we have $T_p \rightarrow \infty$, we can now use Fatou lemma which yields

$$\mathbb{E}_{\mathbf{b}, \mathbf{r}} \left[\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t^{\mathbf{b}}) - \ell(f^b(X_t), Y_t^{\mathbf{b}}) \right] \geq \frac{\epsilon}{16}.$$

This holds for any realization in $\mathcal{A} \cap \bigcap_{p \geq 1} (\mathcal{E}_p \cap \mathcal{F}_p)$ which we recall has probability at least $\frac{\epsilon}{4}$. Therefore we finally obtain

$$\mathbb{E}_{\mathbf{b}, \mathbf{r}, \mathbb{X}} \left[\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t^{\mathbf{b}}) - \ell(f^b(X_t), Y_t^{\mathbf{b}}) \right] \geq \frac{\epsilon^2}{26}.$$

As a result, there exists a specific realization of \mathbf{b} which we denote b such that

$$\mathbb{E}_{\mathbf{r}, \mathbb{X}} \left[\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t^b) - \ell(f^b(X_t), Y_t^b) \right] \geq \frac{\epsilon^2}{26},$$

which shows that with nonzero probability $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t^b) - \ell(f^b(X_t), Y_t^b) > 0$. This ends the proof of the theorem. As a remark, one can note that the construction of our bad example \mathbb{Y}^b is a deterministic function of \mathbb{X} : it is independent from the realizations of the randomness used by the learning rule.

C.2 Proof of Lemma 5.3

We first construct our online learning algorithm, which is a simple variant of the classical exponential forecaster. We first define a step $\eta := \sqrt{2 \ln t_0 / t_0}$. At time $t = 1$ we always predict 0. For time step $t \geq 2$, we define the set $S_{t-1} := \{y \in \mathbb{N}, \sum_{u=1}^{t-1} \mathbb{1}_{y=y_u} > 0\}$ the set of values which have been visited. Then, we construct weights for all $y \in \mathbb{N}$ such that

$$w_{y,t-1} = \begin{cases} e^{\eta \sum_{u=1}^{t-1} \mathbb{1}_{y=y_u}}, & y \in S_{t-1} \\ 0 & \text{otherwise,} \end{cases}$$

and output a randomized prediction independent of the past history such that

$$\mathbb{P}(\hat{y}_t = y) = \frac{w_{y,t-1}}{\sum_{y' \in \mathbb{N}} w_{y',t-1}}.$$

This defines a proper online learning rule. Note that the denominator is well defined since $w_{y,t-1}$ is non-zero only for values in S_{t-1} , which contains at most $t-1$ elements. We now define the expected success at time $1 \leq t \leq T$ as $\hat{s}_t := \frac{w_{y_t,t-1}}{\sum_{y \in \mathbb{N}} w_{y,t-1}} \mathbb{1}_{y_t \in S_t}$. Note that $\hat{s}_t = \mathbb{E}[\mathbb{1}_{f_t(y_{\leq t-1})=y_t}]$. We first show that we have

$$\sum_{t=1}^T \hat{s}_t \geq \max_{y \in \mathbb{N}} \sum_{t=1}^T \mathbb{1}_{y=y_t} - \sqrt{T} \ln T.$$

To do so, we analyze the quantity $W_t := \frac{1}{\eta} \ln \left(\sum_{y \in S_t} e^{\eta \sum_{u=1}^t (\mathbb{1}_{y=y_u} - \hat{s}_u)} \right)$. Let $2 \leq t \leq T$. Supposing that $y_t \in S_{t-1}$, i.e., $S_t = S_{t-1}$, we define the operator $\Phi : \mathbf{x} \in \mathbb{R}^{|S_{t-1}|} \mapsto \frac{1}{\eta} \ln \left(\sum_{y \in S_{t-1}} e^{\eta x_y} \right)$ and use the Taylor expansion of Φ to obtain

$$\begin{aligned} W_t &= \frac{1}{\eta} \ln \left(\sum_{y \in S_{t-1}} e^{\eta \sum_{u=1}^{t-1} (\mathbb{1}_{y=y_u} - \hat{s}_u) + \eta (\mathbb{1}_{y=y_t} - \hat{s}_t)} \right) \\ &= W_{t-1} + \sum_{y \in S_{t-1}} (\mathbb{1}_{y=y_t} - \hat{s}_t) \frac{e^{\eta \sum_{u=1}^{t-1} \mathbb{1}_{y=y_u}}}{\sum_{y' \in S_{t-1}} e^{\eta \sum_{u=1}^{t-1} \mathbb{1}_{y'=y_u}}} \\ &\quad + \frac{1}{2} \sum_{y_1, y_2 \in S_{t-1}} \frac{\partial^2 \Phi}{\partial x_{y_1} \partial x_{y_2}} \Big|_{\xi} (\mathbb{1}_{y_1=y_u} - \hat{s}_u) (\mathbb{1}_{y_2=y_u} - \hat{s}_u) \\ &= W_{t-1} + \frac{1}{2} \sum_{y_1, y_2 \in S_{t-1}} \frac{\partial^2 \Phi}{\partial x_{y_1} \partial x_{y_2}} \Big|_{\xi} (\mathbb{1}_{y_1=y_t} - \hat{s}_u) (\mathbb{1}_{y_2=y_t} - \hat{s}_u) \\ &\leq W_{t-1} + \frac{1}{2} \sum_{y \in S_{t-1}} \frac{\eta e^{\eta \xi_y}}{\sum_{y' \in S_{t-1}} e^{\eta \xi_{y'}}} (\mathbb{1}_{y=y_t} - \hat{s}_u)^2 \\ &\leq W_{t-1} + \frac{\eta}{2}, \end{aligned}$$

for some vector $\xi \in \mathbb{R}^{|S_{t-1}|}$, where in the last inequality we used the fact $|\mathbb{1}_{y=y_t} - \hat{s}_u| \leq 1$. We now suppose that $y_t \notin S_{t-1}$ and $W_{t-1} \geq 1 + \frac{\ln 2 + 2 \ln \frac{1}{\eta}}{\eta}$. In that case, $e^{\eta W_t} = e^{\eta W_{t-1}} + e^{\eta(1 - \sum_{u=1}^{t-1} \hat{s}_u)}$. Hence, we obtain

$$W_t = W_{t-1} + \frac{\ln \left(1 + e^{\eta(1 - \sum_{u=1}^{t-1} \hat{s}_u)} \right)}{\eta} \leq W_{t-1} + \frac{e^{\eta(1 - \sum_{u=1}^{t-1} \hat{s}_u)}}{\eta} \leq W_{t-1} + \frac{\eta}{2}.$$

Now let $l = \max\{1\} \cup \left\{ 1 \leq t \leq T : W_t < 1 + \frac{\ln 2 + 2 \ln \frac{1}{\eta}}{\eta} \right\}$. Note that for any $l < t \leq T$ the above arguments yield $W_t \leq W_{t-1} + \frac{\eta}{2}$. As a result, noting that $W_1 \leq 1$, we finally obtain

$$W_T \leq W_l + \eta \frac{T-l}{2} \leq 1 + \frac{\ln 2 + 2 \ln \frac{1}{\eta}}{\eta} + \eta \frac{T}{2} \leq 1 + \ln 2 \sqrt{\frac{t_0}{2 \ln t_0}} + \sqrt{\frac{\ln t_0}{2 t_0}} (t_0 + T).$$

Therefore, for any $y \in S_T$, we have

$$\sum_{t=1}^T (\mathbb{1}_{y=y_t} - \hat{s}_t) \leq W_T \leq 1 + \ln 2 \sqrt{\frac{t_0}{2 \ln t_0}} + \sqrt{\frac{\ln t_0}{2t_0}} (t_0 + T).$$

In particular, this shows that

$$\sum_{t=1}^T \hat{s}_t \geq \max_{y \in S_T} \sum_{t=1}^T \mathbb{1}_{y=y_t} - 1 - \ln 2 \sqrt{\frac{t_0}{2 \ln t_0}} - \sqrt{\frac{\ln t_0}{2t_0}} (t_0 + T).$$

Now note that if $y \notin S_T$, then $\sum_{t=1}^T \mathbb{1}_{y=y_t} = 0$, which yields $\max_{y \in S_T} \sum_{t=1}^T \mathbb{1}_{y=y_t} = \max_{y \in \mathbb{N}} \sum_{t=1}^T \mathbb{1}_{y=y_t}$. For the sake of conciseness, we will now denote by \hat{y}_t the prediction of the online learning rule at time t . We observe that the variables $\mathbb{1}_{\hat{y}_t=y_t} - \hat{s}_t$ for $1 \leq t \leq T$ form a sequence of martingale differences. Further, $|\mathbb{1}_{\hat{y}_t=y_t} - \hat{s}_t| \leq 1$. Therefore, the Hoeffding-Azuma inequality shows that with probability $1 - \delta$,

$$\sum_{t=1}^T (\mathbb{1}_{\hat{y}_t=y_t} - \hat{s}_t) \geq -\sqrt{2T \ln \frac{1}{\delta}}.$$

Putting everything together yields that with probability $1 - \delta$,

$$\begin{aligned} \sum_{t=1}^T \mathbb{1}_{\hat{y}_t=y_t} &\geq \sum_{t=1}^T \hat{s}_t - \sqrt{2T \ln \frac{1}{\delta}} \\ &\geq \max_{y \in \mathbb{N}} \sum_{t=1}^T \mathbb{1}_{y=y_t} - 1 - \ln 2 \sqrt{\frac{t_0}{2 \ln t_0}} - \sqrt{\frac{\ln t_0}{2t_0}} (t_0 + T) - \sqrt{2T \ln \frac{1}{\delta}}. \end{aligned}$$

This ends the proof of the lemma.

C.3 Proof of Theorem 5.4

We use a similar learning rule to the one constructed in Section 4 for totally-bounded spaces. We only make a slight modification of the learning rules f^ϵ . Precisely, we pose for $0 < \epsilon \leq 1$,

$$T_\epsilon := \left\lceil \frac{2^4 \cdot 3^2 (1 + \ln \frac{1}{\epsilon})}{\epsilon^2} \right\rceil \quad \text{and} \quad \delta_\epsilon := \frac{\epsilon}{2T_\epsilon}.$$

Then, let ϕ be the representative function from the $(1 + \delta_\epsilon)\text{C1NN}$ learning rule. Similarly as for the ϵ -approximation learning rule from Section 4, we consider the same equivalence relation $\overset{\circ}{\sim}$ on times to define clusters. The learning rule then performs its prediction based on the values observed on the corresponding cluster using the learning rule from Lemma 5.3 using $t_0 = T_\epsilon$. Precisely, let $\eta_\epsilon := \sqrt{2 \ln T_\epsilon / T_\epsilon}$ and define the weights $w_{y,t} = e^{\eta_\epsilon \sum_{u < t: u \overset{\circ}{\sim} t} \mathbb{1}(Y_u = y)}$ for all $y \in \tilde{S} := \{y' \in \mathbb{N} : \sum_{u < t: u \overset{\circ}{\sim} t} \mathbb{1}(Y_u = y') > 0\}$ and $w_{y,t} = 0$ otherwise. The learning rule $f_t^\epsilon(\mathbb{X}_{\leq t-1}, \mathbb{Y}_{\leq t-1}, X_t)$ outputs a random value in \mathbb{N} independent of the past history such that

$$\mathbb{P}(\hat{Y}_t = y) = \frac{w_{y,t}}{\sum_{y' \in \mathbb{N}} w_{y',t}}, \quad y \in \mathbb{N}.$$

The final learning rule f is then defined similarly as before from the learning rules f^ϵ for $\epsilon > 0$. Therefore, Lemma 4.2 still holds. Also, using the same notations as in the proof of Theorem 4.3, Lemma 5.3 implies

that for any $t \geq 1$, we can write for any $t \geq 1$ on the cluster $\mathcal{C}(t) = \{u < t : u \overset{\phi}{\sim} t\}$,

$$\begin{aligned}
\sum_{u \in \mathcal{C}(t)} \bar{\ell}_{01}(\hat{Y}_u(\epsilon), Y_u) &\leq \min_{y \in \mathbb{N}} \sum_{u \in \mathcal{C}(t)} \ell_{01}(y, Y_u) + 1 + \ln 2 \sqrt{\frac{T_\epsilon}{2 \ln T_\epsilon}} + \sqrt{\frac{\ln T_\epsilon}{2 T_\epsilon}} (T_\epsilon + |\mathcal{C}(t)|) \\
&\leq \min_{y \in \mathbb{N}} \sum_{u \in \mathcal{C}(t)} \ell_{01}(y, Y_u) + \left(\frac{1}{T_\epsilon} + \frac{\ln 2}{\sqrt{2 T_\epsilon \ln T_\epsilon}} + \sqrt{\frac{2 \ln T_\epsilon}{T_\epsilon}} \right) \max(T_\epsilon, |\mathcal{C}(t)|) \\
&\leq \min_{y \in \mathbb{N}} \sum_{u \in \mathcal{C}(t)} \ell_{01}(y, Y_u) + \left(\frac{\epsilon}{3} + \frac{\epsilon}{3} + \frac{\epsilon}{3} \right) \max(T_\epsilon, |\mathcal{C}(t)|) \\
&= \min_{y \in \mathbb{N}} \sum_{u \in \mathcal{C}(t)} \ell_{01}(y, Y_u) + \epsilon \max(T_\epsilon, |\mathcal{C}(t)|)
\end{aligned}$$

Therefore, the same proof of Theorem 4.3 holds by replacing all ϵ -nets \mathcal{Y}_ϵ directly by \mathbb{N} . The martingale argument still holds since the learning rule used is indeed online. This ends the proof of this theorem.

C.4 Proof of Theorem 5.5

We first need the following simple result which intuitively shows that we can assume that the predictions of the learning rule for mean estimation $g_{\leq t_\epsilon}^\epsilon$ are unrelated for $t = 1, \dots, t_\epsilon$.

Lemma C.1. *Let (\mathcal{Y}, ℓ) satisfying F-TiME. For any $\eta > 0$, there exists a horizon time $T_\eta \geq 1$, an online learning rule $g_{\leq T_\eta}$ such that for any $\mathbf{y} := (y_t)_{t=1}^{T_\eta}$ of values in \mathcal{Y} and any value $y \in \mathcal{Y}$, we have*

$$\frac{1}{T_\eta} \mathbb{E} \left[\sum_{t=1}^{T_\eta} \ell(g_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t) \right] \leq \eta,$$

and such that the random variables $g_t(\mathbf{y}_{\leq t-1})$ are independent.

Proof. Fix $\eta > 0$, $T_\eta \geq 1$ and $g_{\leq T_\eta}$ such that this online learning rule satisfies the condition from F-TiME for $\eta > 0$. We consider the following learning rule \tilde{g} . For any $t \geq 1$ and $\mathbf{y} \in \mathcal{Y}^{t-1}$,

$$\tilde{g}_t(\mathbf{y}_{\leq t-1}) = g_t^t(\mathbf{y}_{\leq t-1}),$$

where (g^t) are i.i.d. samples of the learning rule g . By construction, we still have that for any sequence $\mathbf{y}_{T_\eta} \in \mathcal{Y}^{T_\eta}$,

$$\frac{1}{T_\eta} \mathbb{E} \left[\sum_{t=1}^{T_\eta} \ell(\tilde{g}_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t) \right] = \frac{1}{T_\eta} \mathbb{E} \left[\sum_{t=1}^{T_\eta} \ell(g_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t) \right] \leq \eta.$$

This ends the proof of the lemma. \square

From now on, by Lemma C.1, we will suppose without loss of generality that the learning rule g^ϵ has predictions that are independent at each step (conditionally on the observed values). For simplicity, we refer to the prediction of the defined learning rule f . (resp. f^ϵ) at time t as \hat{Y}_t (resp. $\hat{Y}_t(\epsilon)$). We now show that is optimistically universal for arbitrary responses. By construction of the learning rule f , Lemma 4.2 still holds. Therefore, we only have to focus on the learning rules f^ϵ and prove that we obtain similar results as before. Let $T \geq 1$ and denote by $\mathcal{A}_i := \{t \leq T : |\{u \leq T : \phi(u) = t\}| = i\}$ the set of times which have exactly i children within horizon T for $i = 0, 1, 2$. Then, we define

$$\mathcal{B}_T = \{t \leq T : L_t = 0 \text{ and } |\{t < u \leq T : u \overset{\phi}{\sim} t\}| \geq t_\epsilon\},$$

i.e., times that start a new learning block and such that there are at least t_ϵ future times falling in their cluster within horizon T . Note that the function ψ defines a parent-relation (similarly to ϕ , but defined for all times $t \geq 1$). To simplify notations, for any $t \in \mathcal{B}_T$, we denote t^u the ψ -children of t at generation $u - 1$ for $1 \leq u \leq t_\epsilon$, i.e., we have $\psi^{u-1}(t^u) = t$ for all $1 \leq u \leq t_\epsilon$. In particular $t = t^1$. By construction, blocks have length at most t_ϵ . More precisely, the block started at any $t \in \mathcal{B}_T$ has had time to finish completely, hence has length exactly t_ϵ . By construction of the indices L_t , the blocks $\{t^u, 1 \leq u \leq t_\epsilon\}$, for $t \in \mathcal{B}_T$, are all disjoint. This implies in particular $|\mathcal{B}_T|_{t_\epsilon} \leq T$. We first analyze the predictions along these blocks and for any $t \in \mathcal{B}_T$ and $y \in \mathcal{Y}$, we pose $\delta_t(y) := \frac{1}{t_\epsilon} \sum_{u=1}^{t_\epsilon} \left(\ell(\hat{Y}_{t^u}, Y_{t^u}) - \ell(y, Y_{t^u}) - \epsilon \right)$. Now by construction of the learning rule f^ϵ , we have

$$t_\epsilon \delta_t(y^t) = \sum_{u=1}^{t_\epsilon} \left(\ell(g_u^{\epsilon, t}(\{Y_{t^i}\}_{i=1}^{u-1}), Y_{t^u}) - \ell(y^t, Y_{t^u}) \right) - \epsilon t_\epsilon.$$

Next, for any $t \leq t_\epsilon$ and sequence $\mathbf{y}_{\leq t-1}$ and value $y \in \mathcal{Y}$, we write $\bar{\ell}(g_t^\epsilon(\mathbf{y}_{\leq t-1}), y) := \mathbb{E} [\ell(g_t^\epsilon(\mathbf{y}_{\leq t-1}), y)]$. Now by hypothesis on the learning rule $g_{\leq t_\epsilon}^\epsilon$,

$$\frac{1}{t_\epsilon} \sum_{u=1}^{t_\epsilon} \bar{\ell}(\hat{Y}_{t^u}, Y_{t^u}) - \ell(y^t, Y_{t^u}) \leq \epsilon. \quad (2)$$

Now consider the following sequence $(\ell(\hat{Y}_{t^u}, Y_{t^u}) - \bar{\ell}(\hat{Y}_{t^u}, Y_{t^u}))_{t \in \mathcal{B}_T, 1 \leq u \leq s(t)}$. Because of the definition of the learning rule, which uses i.i.d. copies of the learning rule g^ϵ , if we order the former sequence by increasing order of t^u , we obtain a sequence of martingale differences. We can continue this sequence by zeros to ensure that it has length exactly T . As a result, we obtain a sequence of T martingale differences, which are all bounded by $\bar{\ell}$ in absolute value. Now, the Azuma-Hoeffding inequality implies that for $\delta > 0$, with probability $1 - \delta$, we have

$$\sum_{t \in \mathcal{B}_T} \sum_{u=1}^{t_\epsilon} \ell(\hat{Y}_{t^u}, Y_{t^u}) \leq \sum_{t \in \mathcal{B}_T} \sum_{u=1}^{t_\epsilon} \bar{\ell}(\hat{Y}_{t^u}, Y_{t^u}) + \bar{\ell} \sqrt{2T \ln \frac{1}{\delta}}.$$

Thus, using Eq (2), with probability at least $1 - \delta$,

$$\sum_{t \in \mathcal{B}_T} t_\epsilon \delta_t(y^t) \leq \bar{\ell} \sqrt{2T \ln \frac{1}{\delta}}. \quad (3)$$

We also denote $\mathcal{T} = \bigcup_{t \in \mathcal{B}_T} \{t^u, 1 \leq u \leq t_\epsilon\}$ the union of all blocks within horizon T . This set contains all times $t \leq T$ except *bad* times close to the last times of their corresponding cluster $\{u \leq T : u \stackrel{\phi}{\sim} t\}$. Precisely, these are times t such that $|\{t < u \leq T : u \stackrel{\phi}{\sim} t\}| < t_\epsilon - L_t$. As a result, there are at most t_ϵ such times for each cluster. Using the same arguments as in the proof of Theorem 4.3, if we consider only clusters of duplicates (i.e., the cluster started for a specific instance which has high number of duplicates), the corresponding *bad* times contribute to a proportion $\leq \frac{t_\epsilon}{T_\epsilon/\epsilon} \leq \epsilon^2$ of times. Now consider clusters that have at least T_ϵ times. Their *bad* times contribute to a proportion $\leq \frac{t_\epsilon}{T_\epsilon} \leq \epsilon$ of times. Last, we need to account for clusters of size $< T_\epsilon$ which necessarily contain leaves of the tree ϕ : there are at most $|\mathcal{A}_0|$ such clusters. By the Chernoff bound, with probability at least $1 - e^{-T\delta_\epsilon/3}$ we have

$$T - |\mathcal{T}| \leq (\epsilon^2 + \epsilon)T + |\mathcal{A}_0|t_\epsilon \leq t_\epsilon + (\epsilon^2 + \epsilon + 2\delta_\epsilon t_\epsilon)T \leq t_\epsilon + 3\epsilon T.$$

By the Borel-Cantelli lemma, because $\sum_{T \geq 1} e^{-T\delta_\epsilon/3} < \infty$, almost surely there exists a time \hat{T} such that for $T \geq \hat{T}$ we have $T - |\mathcal{T}| \leq t_\epsilon + 3\epsilon T$. We denote by \mathcal{E}_ϵ this event. Then, on the event \mathcal{E}_ϵ , for any $T \geq \hat{T}$ and

for any sequence of values $(y^t)_{t \geq 1}$ we have

$$\begin{aligned}
\sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) &\leq \sum_{t \in \mathcal{B}_T} \sum_{u=1}^{t_\epsilon} \ell(\hat{Y}_{t^u}, Y_{t^u}) + (T - |\mathcal{T}|)\bar{\ell} \\
&\leq \sum_{t \in \mathcal{B}_T} \sum_{u=1}^{t_\epsilon} \ell(y^t, Y_{t^u}) + \sum_{t \in \mathcal{B}_T} t_\epsilon \delta_t(y^t) + \epsilon |\mathcal{B}_T| t_\epsilon + t_\epsilon \bar{\ell} + 3\epsilon T \\
&\leq \sum_{t \in \mathcal{B}_T} \sum_{u=1}^{t_\epsilon} \ell(y^t, Y_{t^u}) + \sum_{t \in \mathcal{B}_T} t_\epsilon \delta_t(y^t) + t_\epsilon \bar{\ell} + 4\epsilon T.
\end{aligned}$$

Now let $f : \mathcal{X} \rightarrow \mathcal{Y}$ be a measurable function to which we compare f^ϵ . By Theorem 4.1, because $(1+\delta_\epsilon)$ C1NN is optimistically universal without noise and $\mathbb{X} \in \text{SOUL}$, almost surely $\frac{1}{T} \sum_{t=1}^T \ell(f(X_{\phi(t)}), f(X_t)) \rightarrow 0$. We denote by \mathcal{F}_ϵ this event of probability one. The proof of Theorem 4.3 shows that on \mathcal{F}_ϵ , for any $0 \leq u \leq T_\epsilon - 1$ we have

$$\frac{1}{T} \sum_{t=1}^T \ell(f(X_{\phi^u(t)}), f(X_t)) \rightarrow 0.$$

We let $y^t = f(X_t)$ for all $t \geq 1$. Then, recalling that for any $t \in \mathcal{B}_T$, we have $t = \phi^{u-1}(t^u)$, on the event \mathcal{E}_ϵ , for any $T \geq \hat{T}$ we have

$$\begin{aligned}
&\sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) \\
&\leq \sum_{t \in \mathcal{B}_T} \sum_{u=1}^{t_\epsilon} ((1+\epsilon)\ell(f(X_{t^u}), Y_{t^u}) + c_\epsilon^\alpha \ell(f(X_t), f(X_{t^u}))) + \sum_{t \in \mathcal{B}_T} t_\epsilon \delta_t(y^t) + t_\epsilon \bar{\ell} + 4\epsilon T \\
&\leq \sum_{t=1}^T \ell(f(X_t), Y_t) + c_\epsilon^\alpha \frac{T_\epsilon}{\epsilon} \sum_{u=0}^{T_\epsilon-1} \sum_{t=1}^T \ell(f(X_{\phi^u(t)}), f(X_t)) + \sum_{t \in \mathcal{B}_T} t_\epsilon \delta_{\varphi(t)}(y^t) + t_\epsilon \bar{\ell} + 5\epsilon T,
\end{aligned}$$

where in the first inequality we used Lemma A.1, and in the second inequality we used the fact that cluster with distinct instance values have at most $\frac{T_\epsilon}{\epsilon}$ duplicates of each instance. Next, using Eq (3), with probability $1 - \frac{1}{T^2}$, we have

$$\sum_{t \in \mathcal{B}_T} t_\epsilon \delta_t(y^t) \leq 2\bar{\ell} \sqrt{T \ln T}.$$

Because $\sum_{T \geq 1} \frac{1}{T^2} < \infty$, the Borel-Cantelli lemma implies that on an event \mathcal{G}_ϵ of probability one, there exists \hat{T}_2 such that for all $T \geq \hat{T}_2$ the above inequality holds. As a result, on the event $\mathcal{E}_\epsilon \cap \mathcal{F}_\epsilon \cap \mathcal{G}_\epsilon$ we obtain for any $T \geq \max(\hat{T}, \hat{T}_2)$ that

$$\begin{aligned}
\sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) &\leq \sum_{t=1}^T \ell(f(X_t), Y_t) + \frac{c_\epsilon^\alpha T_\epsilon}{\epsilon} \sum_{u=0}^{T_\epsilon-1} \sum_{t=1}^T \ell(f(X_{\phi^u(t)}), f(X_t)) \\
&\quad + 2\bar{\ell} \sqrt{T \ln T} + t_\epsilon \bar{\ell} + 5\epsilon T.
\end{aligned}$$

where $\frac{1}{T} \sum_{u=0}^{T_\epsilon-1} \sum_{t=1}^T \ell(f(X_{\phi^u(t)}), f(X_t)) \rightarrow 0$ because the event \mathcal{F}_ϵ is met. Therefore, we obtain that on the event $\mathcal{E}_\epsilon \cap \mathcal{F}_\epsilon \cap \mathcal{G}_\epsilon$ of probability one,

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \left[\ell(\hat{Y}_t(\epsilon), Y_t) - \ell(f(X_t), Y_t) \right] \leq 5\epsilon,$$

i.e., almost surely, the learning rule f^ϵ achieves risk at most 5ϵ compared to the fixed function f . By union bound, on the event $\bigcap_{i \geq 0} (\mathcal{E}_{\epsilon_i} \cap \mathcal{F}_{\epsilon_i} \cap \mathcal{G}_{\epsilon_i})$ of probability one we have that

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \left[\ell(\hat{Y}_t(\epsilon_i), Y_t) - \ell(f(X_t), Y_t) \right] \leq 5\epsilon_i, \quad \forall i \geq 0.$$

The rest of the proof uses similar arguments as in the proof of Theorem 4.3. Precisely, let \mathcal{H} be the almost sure event of Lemma 4.2 such that there exists \hat{t} for which

$$\forall t \geq \hat{t}, \forall i \in I_t, \quad \sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) \leq \sum_{s=t_i}^t \ell(\hat{Y}_s(\epsilon_i), Y_s) + (2 + \bar{\ell} + \bar{\ell}^2) \sqrt{t \ln t}.$$

In the rest of the proof we will suppose that the event $\mathcal{H} \cap \bigcap_{i \geq 0} (\mathcal{E}_{\epsilon_i} \cap \mathcal{F}_{\epsilon_i} \cap \mathcal{G}_{\epsilon_i})$ of probability one is met. Let $i \geq 0$. For all $t \geq \max(\hat{t}, t_i)$ we have

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) &\leq \frac{t_i}{T} \bar{\ell} + \frac{1}{T} \sum_{t=t_i}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \\ &\leq \frac{t_i}{T} \bar{\ell} + \frac{1}{T} \sum_{t=t_i}^T \ell(\hat{Y}_t(\epsilon_i), Y_t) - \ell(f(X_t), Y_t) + (2 + \bar{\ell} + \bar{\ell}^2) \sqrt{\frac{\ln T}{T}} \\ &\leq \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t(\epsilon_i), Y_t) - \ell(f(X_t), Y_t) + \frac{2t_i}{T} \bar{\ell} + (2 + \bar{\ell} + \bar{\ell}^2) \sqrt{\frac{\ln T}{T}}. \end{aligned}$$

Therefore we obtain $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq 5\epsilon_i$. Because this holds for any $i \geq 0$ we finally obtain

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq 0.$$

As a result, f is universally consistent for adversarial responses under all SOUL processes. Hence, SOLAR = SOUL and f is in fact optimistically universal. This ends the proof of the theorem.

C.5 Proof of Lemma 5.7

We first note that with the same horizon time T_η , we have that F-TIME implies Property 2. We now show that Property 2 implies F-TIME. Let (\mathcal{Y}, ℓ) satisfying Property 2. We now fix $\eta > 0$ and let $T, g_{\leq \tau}$ such that for any $\mathbf{y} := (y_t)_{t=1}^T$ of values in \mathcal{Y} and any value $y \in \mathcal{Y}$, we have

$$\mathbb{E} \left[\frac{1}{\tau} \sum_{t=1}^{\tau} (\ell(g_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t)) \right] \leq \eta.$$

We now construct a random time $1 \leq \tilde{\tau} \leq T$ such that $\mathbb{P}[\tilde{\tau} = t] = \frac{\mathbb{P}[\tau=t]}{t\mathbb{E}[1/\tau]}$ for all $1 \leq t \leq T$. This indeed defines a proper random variable because $\sum_{t=1}^T \frac{\mathbb{P}[\tau=t]}{t\mathbb{E}[1/\tau]} = 1$. Let $\text{Supp}(\tau) := \{1 \leq t \leq T : \mathbb{P}[\tau = t] > 0\}$ be the support of τ . For any $t \in \text{Supp}(\tau)$, we denote by $g_{\leq t}^t$ the learning rule obtained by conditioning $g_{\leq \tau}$ on the event $\{\tau = t\}$, i.e., $g_{\leq t}^t = g_{\leq \tau} | \tau = t$. More precisely, recall that τ only uses the randomness of g_t . It is not an online random time. Hence, a practical way to simulate $g_{\leq t}^t$ for all $t \in \text{Supp}(\tau)$ is to first draw an i.i.d. sequence of learning rules $(g_{i, \leq \tau_i})_{i \geq 1}$. Then, for each $t \in \text{Supp}(\tau)$ we select the randomness which first satisfies $\tau = t$. Specifically, we define the time $i_t = \min\{i : \tau_i = t\}$ for all $t \in \text{Supp}(\tau)$. With probability one,

these times are finite for all $t \in \text{Supp}(\tau)$. Denote this event \mathcal{E} . Then, letting $\bar{y} \in \mathcal{Y}$ be an arbitrary fixed value, for all $1 \leq t \leq T$ we pose

$$g_{\leq t}^t = \begin{cases} g_{i_t, \leq t} & \text{if } \mathcal{E} \text{ is met,} \\ \bar{y}_{\leq t} & \text{otherwise,} \end{cases} \quad t \in \text{Supp}(\tau) \quad \text{and} \quad g_{\leq t}^t = \bar{y}_{\leq t}, \quad t \notin \text{Supp}(\tau).$$

where $\bar{y}_{\leq t}$ denotes the learning rules which always outputs value \bar{y} for all steps $u \leq t$. Intuitively, $g_{\leq t}^t$ has the same distribution as $g_{\leq \tau}$ conditioned on the event $\{\tau = t\}$. We are now ready to define a new learning rule $\tilde{g}_{\leq \tilde{\tau}}$, by $\tilde{g}_{\leq \tilde{\tau}} := g_{\leq \tilde{\tau}}^{\tilde{\tau}}$. Noting that for any $t \notin \text{Supp}(\tau)$ we have $\mathbb{P}[\tilde{\tau} = t] = 0$, we can write

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=1}^{\tau} (\ell(\tilde{g}_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t)) - \eta\tau \right] \\ &= \sum_{t=1}^T \mathbb{P}[\tilde{\tau} = t] \mathbb{E} \left[\sum_{u=1}^t (\ell(\tilde{g}_u(\mathbf{y}_{\leq u-1}), y_u) - \ell(y, y_u)) - \eta t \mid \tilde{\tau} = t \right] \\ &= \sum_{t \in \text{Supp}(\tau)} \mathbb{P}[\tilde{\tau} = t] \mathbb{E} \left[\sum_{u=1}^t (\ell(\tilde{g}_u(\mathbf{y}_{\leq u-1}), y_u) - \ell(y, y_u)) - \eta t \mid \tilde{\tau} = t, \mathcal{E} \right] \\ &= \frac{1}{\mathbb{E}[1/\tau]} \sum_{t \in \text{Supp}(\tau)} \mathbb{P}[\tau = t] \mathbb{E} \left[\frac{1}{t} \sum_{u=1}^t (\ell(g_{i_t, u}(\mathbf{y}_{\leq u-1}), y_u) - \ell(y, y_u)) - \eta \mid \tau = t, \mathcal{E} \right] \\ &= \frac{1}{\mathbb{E}[1/\tau]} \sum_{t \in \text{Supp}(\tau)} \mathbb{P}[\tau = t] \mathbb{E} \left[\frac{1}{t} \sum_{u=1}^t (\ell(g_{i_t, u}(\mathbf{y}_{\leq u-1}), y_u) - \ell(y, y_u)) - \eta \right] \\ &= \frac{1}{\mathbb{E}[1/\tau]} \sum_{t \in \text{Supp}(\tau)} \mathbb{P}[\tau = t] \mathbb{E} \left[\frac{1}{t} \sum_{u=1}^t (\ell(g_u(\mathbf{y}_{\leq u-1}), y_u) - \ell(y, y_u)) - \eta \mid \tau = t \right] \\ &= \frac{1}{\mathbb{E}[1/\tau]} \mathbb{E} \left[\frac{1}{\tau} \sum_{t=1}^{\tau} (\ell(g_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t)) - \eta \right] \leq 0. \end{aligned}$$

where in the second and fourth equality we used the fact that $\mathbb{P}[\mathcal{E}] = 1$. As a result, there exists a learning rule $\tilde{g}_{\leq \tilde{\tau}}$ such that $1 \leq \tilde{\tau} \leq T_\eta$, and for any $\mathbf{y}_{\leq T_\eta} \in \mathcal{Y}^{T_\eta}$ and $y \in \mathcal{Y}$ one has

$$\mathbb{E} \left[\sum_{t=1}^{\tilde{\tau}} (\ell(\tilde{g}_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t)) - \eta\tilde{\tau} \right] \leq 0.$$

We now pose $T'_\eta = \lceil T_\eta/\eta \rceil$ and draw an i.i.d. sequence of learning rules $(\tilde{g}_{\leq \tilde{\tau}_i}^i)_{i \geq 1}$. Denote $\theta_i = \sum_{j < i} \tilde{\tau}_j$ with the convention $\theta_1 = 0$. We are now ready to define a learning rule $h_{\leq T'_\eta}$ as follows. For any $1 \leq t \leq T'_\eta$ and $\mathbf{y}_{\leq t} \in \mathcal{Y}^t$,

$$h_t(\mathbf{y}_{\leq t-1}) = \tilde{g}_{\leq t-\theta_i}^i((y_{t'})_{\theta_i < t' \leq t-1}), \quad \theta_i < t \leq \theta_{i+1}, i \geq 1.$$

In other words, the learning rule performs independent learning rules $\tilde{g}_{\leq \tilde{\tau}}$ and when the time horizon $\tilde{\tau}$ is reached, we re-initialize the learning rule with a new randomness. Now let $\mathbf{y}_{\leq T'_\eta} \in \mathcal{Y}^{T'_\eta}$ and $y \in \mathcal{Y}$. We denote by $\hat{i} = \max\{i \geq 1, \theta_i \leq t\}$, the index of the last learning rule which had time to finish completely.

Then, because $\tilde{\tau}_i \leq T_\eta$,

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=1}^{T'_\eta} (\ell(h_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t)) - 2\eta T'_\eta \right] \\ & \leq \mathbb{E} \left[\sum_{i \leq \hat{i}} \sum_{t=1}^{\tilde{\tau}_i} (\ell(\tilde{g}_{t-\theta_i}^i(\mathbf{y}_{\theta_i < \cdot \leq t-1}), y_t) - \ell(y, y_t)) - \eta T'_\eta \right] - \eta T'_\eta + T_\eta \\ & \leq \mathbb{E} \left[\sum_{i \leq \hat{i}} \left(\sum_{t=1}^{\tilde{\tau}_i} (\ell(\tilde{g}_{t-\theta_i}^i(\mathbf{y}_{\theta_i < \cdot \leq t-1}), y_t) - \ell(y, y_t)) - \eta \tilde{\tau}_i \right) \right]. \end{aligned}$$

We now analyze the last term. First, note that by construction, the sequence

$$\left\{ S_j := \sum_{i \leq j} \left(\sum_{t=1}^{\tilde{\tau}_j} (\ell(\tilde{g}_{t-\theta_j}^j(\mathbf{y}_{\theta_j < \cdot \leq t-1}), y_t) - \ell(y, y_t)) - \eta \tilde{\tau}_j \right) \right\}_{j \geq 1}$$

is a super-martingale. Now, note that $\hat{i} \leq 1 + T'_\eta$ since for all i , $\theta_i = \sum_{j < i} \tau_j \geq i - 1$. As a result, \hat{i} is bounded, is a stopping time for the considered filtration (after finishing period \hat{i} we stop if and only we exceed time T'_η) and we can apply Doob's optimal sampling theorem to obtain $\mathbb{E}[S_{\hat{i}}] \leq 0$. Thus, combining the above equations gives

$$\frac{1}{T'_\eta} \mathbb{E} \left[\sum_{t=1}^{T'_\eta} (\ell(h_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t)) \right] \leq 2\eta.$$

Because this holds for all $\eta > 0$, F-TIME is satisfied. This ends the proof of the lemma.

C.6 Proof of Theorem 5.8

We first prove that adversarial regression for processes outside of CS is not achievable. Precisely, we show that for any $\mathbb{X} \notin \text{CS}$, for any online learning rule f , there exists a process \mathbb{Y} on \mathcal{Y} , a measurable function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ and $\delta > 0$ such that with non-zero probability $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f^*) > \delta$.

Because F-TIME is not satisfied by (\mathcal{Y}, ℓ) , by Lemma 5.7, Property 2 is not satisfied either. Hence, we can fix $\eta > 0$ such that for any horizon $T \geq 1$ and any online learning rule $g_{\leq \tau}$ with $1 \leq \tau \leq T$, there exist a sequence $\mathbf{y} := (y_t)_{t=1}^T$ of values in \mathcal{Y} and a value y such that

$$\mathbb{E} \left[\frac{1}{\tau} \sum_{t=1}^{\tau} (\ell(g_t(\mathbf{y}_{\leq t-1}), y_t) - \ell(y, y_t)) \right] > \eta,$$

as in the assumption of the space (\mathcal{Y}, ℓ) . Let $\mathbb{X} \notin \text{CS}$. The proof of Theorem 5.1 shows that there exist $0 < \epsilon < 1$, a sequence of disjoint measurable sets $\{B_p\}_{p \geq 1}$ and a sequence of times $(t_p)_{p \geq 0}$ with $t_0 = 0$ and such that with $\mu := \max(1, \frac{8\bar{\ell}}{\epsilon\eta})$, for any $p \geq 1$, $t_p > \mu t_{p-1}$, and defining the events

$$\mathcal{E}_p = \left\{ \mathbb{X}_{\leq t_{p-1}} \cap \left(\bigcup_{p' \geq p} B_{p'} \right) = \emptyset \right\} \text{ and } \mathcal{F}_p := \bigcup_{\mu t_{p-1} < t \leq t_p} \left\{ \frac{1}{t} \sum_{u=1}^t \mathbb{1}_{B_p}(X_u) \geq \frac{\epsilon}{4} \right\},$$

we have $\mathbb{P}[\bigcap_{p \geq 1} (\mathcal{E}_p \cap \mathcal{F}_p)] \geq \frac{\epsilon}{4}$. We now fix a learning rule f and construct a ‘‘bad’’ process \mathbb{Y} recursively. Fix $\bar{y} \in \mathcal{Y}$ an arbitrary value. We start by defining the random variables $N_p(t) = \sum_{u=t_{p-1}+1}^t \mathbb{1}_{B_p}(X_u)$ for any $p \geq 1$. We now construct (deterministic) values y_p and sequences $(y_p^u)_{u=1}^{t_p}$ for all $p \geq 1$, of values in \mathcal{Y} .

Suppose we have already constructed the values y_q as well as the sequences $(y_q^u)_{u=1}^{t_q}$ for all $q < p$. We will now construct y_p and $(y_p^u)_{u=1}^{t_p}$. Assuming that the event $\mathcal{E}_p \cap \mathcal{F}_p$ is met, there exists $\mu t_{p-1} < t \leq t_p$ such that

$$N_p(t) = \sum_{u=t_{p-1}+1}^t \mathbb{1}_{B_p}(X_u) = \sum_{u=1}^t \mathbb{1}_{B_p}(X_u) \geq \frac{\epsilon}{4}t,$$

where in the first equality we used the fact that on \mathcal{E}_p , the process $\mathbb{X}_{\leq t_{p-1}}$ does not visit B_p . In the rest of the construction, we will denote

$$T_p = \begin{cases} \min\{\mu t_{p-1} < t \leq t_p : N_p(t) \geq \frac{\epsilon}{4}t\} & \text{if } \mathcal{E}_p \cap \mathcal{F}_p \text{ is met.} \\ t_p & \text{otherwise.} \end{cases}$$

Now consider the process $\mathbb{Y}_{t \leq t_{p-1}}(\mathbb{X})$ defined as follows. For any $1 \leq q < p$ we pose

$$Y_t(\mathbb{X}) = \begin{cases} y_q^{N_q(t)} & \text{if } t \leq T_q \text{ and } X_t \in B_q, \\ y_q & \text{if } t > T_q \text{ and } X_t \in B_q, \\ y_{q'} & \text{if } X_t \in B_{q'}, q' < q, \\ \bar{y} & \text{otherwise,} \end{cases} \quad t_{q-1} < t \leq t_q.$$

Similarly, for $M \geq 1$ and given any sequence $\{\tilde{y}_i\}_{i=1}^M$, we define the following process $\mathbb{Y}_{t_{p-1} < u \leq t_p}(\mathbb{X}, \{\tilde{y}_i\}_{i=1}^M)$ by

$$Y_u(\mathbb{X}, \{\tilde{y}_i\}_{i=1}^M) = \begin{cases} \tilde{y}_{\min(N_p(u), M)} & \text{if } X_t \in B_p, \\ y_q & \text{if } X_t \in B_q, q < p, \\ \bar{y} & \text{otherwise.} \end{cases}$$

We now construct a learning rule g^p . First, we define the event $\mathcal{B} := \bigcap_{p \geq 1} (\mathcal{E}_p \cap \mathcal{F}_p)$. We will denote by $\tilde{\mathbb{X}} = \mathbb{X}|_{\mathcal{B}}$ a sampling of the process \mathbb{X} on the event \mathcal{B} which has probability at least $\frac{\epsilon}{4}$. For instance we draw i.i.d. samplings following the same distribution as \mathbb{X} then select the process which first falls into \mathcal{B} . We are now ready to define a learning rule $(g_u^p)_{u \leq \tau}$ where τ is a random time. To do so, we first draw a sample $\tilde{\mathbb{X}}$ which is now fixed for the learning rule g^p . We define the stopping time as $\tau = N_p(T_p)$. Finally, for all $1 \leq u \leq \tau$, and any sequence of values $\tilde{\mathbf{y}}_{\leq u-1}$, we pose

$$g_u^p(\tilde{\mathbf{y}}_{\leq u-1}) = f_{T_p(u)}\left(\tilde{\mathbb{X}}_{\leq T_p(u)-1}, \left\{ \mathbb{Y}_{\leq t_{p-1}}(\tilde{\mathbb{X}}), \mathbb{Y}_{t_{p-1} < u \leq T_p(u)-1}\left(\tilde{\mathbb{X}}, \{\tilde{y}_i\}_{i=1}^{u-1}\right) \right\}, \tilde{X}_{T_p(u)}\right),$$

where we used the notation $T_p(u) := \min\{t_{p-1} < t' \leq t_p : N_p(t') = u\}$ for the time of the u -th visit of B_p , which exists because $u \leq \tau = N_p(T_p) \leq N_p(t_p)$ since the event \mathcal{B} is satisfied by $\tilde{\mathbb{X}}$. Note that the prediction of the rule g^p is random because of the dependence on $\tilde{\mathbb{X}}$. Also, observe that the random time τ is bounded by $1 \leq \tau \leq T_p \leq t_p$. Therefore, by hypothesis on the value space (\mathcal{Y}, ℓ) , there exists a sequence $\{y_p^u\}_{u=1}^{t_p}$ and a value $y_p \in \mathcal{Y}$ such that

$$\mathbb{E} \left[\frac{1}{\tau} \sum_{u=1}^{\tau} (\ell(g_u^p(\mathbf{y}_p^{\leq u-1}), y_p^u) - \ell(y_p, y_p^u)) \right] \geq \eta.$$

This ends the recursive construction of the values y_p and the sequences $(y_p^u)_{u=1}^{t_p}$ for all $p \geq 1$. We are now ready to define the process $\mathbb{Y}(\mathbb{X})$, using a similar construction as before. For any $p \geq 1$ we define

$$Y_t(\mathbb{X}) = \begin{cases} y_p^{N_p(t)} & \text{if } t \leq T_p \text{ and } X_t \in B_p, \\ y_p & \text{if } t > T_p \text{ and } X_t \in B_p, \\ y_q & \text{if } X_t \in B_q, q < p, \\ \bar{y} & \text{otherwise,} \end{cases} \quad t_{p-1} < t \leq t_p.$$

We also define a function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ by

$$f^*(x) = \begin{cases} y_p & \text{if } x \in B_p, \\ \bar{y} & \text{otherwise.} \end{cases}$$

This function is simple hence measurable. From now, we will suppose that the event \mathcal{B} is met. For simplicity, we will denote by $\hat{Y}_t := f_t(\mathbb{X}_{\leq t-1}, \mathbb{Y}_{\leq t-1}, X_t)$ the prediction of the learning rule at time t . For any $p \geq 1$, because $\mathcal{E}_p \cap \mathcal{F}_p$ is met, for all $1 \leq u \leq N_p(T_p)$, we have $t_{p-1} < T_p(u) \leq T_p$, and $X_{T_p(u)} \in B_p$. Hence, by construction, we have $\hat{Y}_{T_p(u)} = y_p^u$ and we can write

$$\begin{aligned} \sum_{t=1}^{T_p} \ell(\hat{Y}_t, Y_t) &\geq \sum_{t=t_{p-1}+1}^{T_p} \ell(\hat{Y}_t, Y_t) \\ &\geq \sum_{u=1}^{N_p(T_p)} \ell(\hat{Y}_{T_p(u)}, Y_{T_p(u)}) \\ &= \sum_{u=1}^{\tau} \ell(f_{T_p(u)}(\mathbb{X}_{\leq T_p(u)-1}, \mathbb{Y}_{\leq T_p(u)-1}, X_{T_p(u)}), y_p^u). \end{aligned}$$

Now note that because the construction was similar to the construction of g^p , we have $\mathbb{Y}_{\leq T_p(u)-1} = \{\mathbb{Y}_{\leq t_{p-1}}(\mathbb{X}), \mathbb{Y}_{t_{p-1} < t \leq T_p(u)-1}(\mathbb{X}, \{y_p^i\}_{i=1}^{u-1})\}$, i.e., $\hat{Y}_{T_p(u)}$ coincides with the prediction $g_u^p(\{y_p^i\}_{i=1}^{u-1})$ provided that g_u^p precisely used the realization \mathbb{X} . Hence, conditioned on \mathcal{B} for all $u \leq \tau$, $\hat{Y}_{T_p(u)}$ has the same distribution as $g_u^p(\mathbf{y}_p^{\leq u-1})$. Therefore we obtain

$$\begin{aligned} \mathbb{E} \left[\frac{1}{\tau} \sum_{t=1}^{T_p} \ell(\hat{Y}_t, Y_t) - \frac{1}{\tau} \sum_{u=1}^{\tau} \ell(y_p, y_p^u) \middle| \mathcal{B} \right] &\geq \mathbb{E} \left[\frac{1}{\tau} \sum_{u=1}^{\tau} \left(\ell(g_u^p(\hat{Y}_{T_p(u)}, y_p^u) - \ell(y_p, y_p^u)) \right) \middle| \mathcal{B} \right] \\ &= \mathbb{E} \left[\frac{1}{\tau} \sum_{u=1}^{\tau} \left(\ell(g_u^p(\mathbf{y}_p^{\leq u-1}), y_p^u) - \ell(y_p, y_p^u) \right) \right] \\ &\geq \eta. \end{aligned}$$

We now turn to the loss obtained by the simple function f^* . By construction, assuming that the event \mathcal{B} is met, we have

$$\sum_{t=1}^{T_p} \ell(f^*(X_t), Y_t) \leq \bar{\ell} t_{p-1} + \sum_{u=1}^{N_p(T_p)} \ell(f^*(X_{T_p(u)}), y_p^u) = \bar{\ell} t_{p-1} + \sum_{u=1}^{\tau} \ell(y_p, y_p^u).$$

Recalling that $T_p > \mu t_{p-1} \geq \frac{8\bar{\ell}}{\epsilon\eta} t_{p-1}$ and noting that $\tau = N_p(T_p) \geq \frac{\epsilon}{4} T_p$, we obtain

$$\begin{aligned} &\mathbb{E} \left[\sup_{t_{p-1} < T \leq t_p} \frac{1}{T} \sum_{t=1}^T (\ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t)) \middle| \mathcal{B} \right] \\ &\geq \mathbb{E} \left[\frac{\tau}{T_p} \frac{1}{\tau} \left(\sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \sum_{u=1}^{\tau} \ell(y_p, y_p^u) \right) - \bar{\ell} \frac{t_{p-1}}{T_p} \middle| \mathcal{B} \right] \\ &\geq \frac{\epsilon}{4} \mathbb{E} \left[\frac{1}{\tau} \sum_{t=1}^{T_p} \ell(\hat{Y}_t, Y_t) - \frac{1}{\tau} \sum_{u=1}^{\tau} \ell(y_p, y_p^u) \middle| \mathcal{B} \right] - \frac{\epsilon\eta}{8} \\ &\geq \frac{\epsilon\eta}{8}. \end{aligned}$$

Because this holds for any $p \geq 1$, Fatou lemma yields

$$\begin{aligned} & \mathbb{E} \left[\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \right] \\ & \geq \mathbb{E} \left[\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T (\ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t)) \middle| \mathcal{B} \right] \mathbb{P}[\mathcal{B}] \\ & \geq \frac{\epsilon^2 \eta}{32}. \end{aligned}$$

Hence, we do not have almost surely $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq 0$. This shows that $\mathbb{X} \notin \text{SOLAR}$, which in turn implies $\text{SOLAR} \subset \text{CS}$. This ends the proof that $\text{SOLAR} \subset \text{CS}$. The proof that $\text{CS} \subset \text{SOLAR}$ and the construction of an optimistically universal learning rule for adversarial regression is deferred to Section 7 where we give a stronger result which also holds for unbounded losses. Note that generalizing Theorem 5.2 to adversarial responses already shows that $\text{CS} \subset \text{SOLAR}$ and provides an optimistically universal learning rule when the loss ℓ is a metric $\alpha = 1$.

D Proofs of Section 6

D.1 Proof of Theorem 3.6

We first show that there exists $t_1 \geq 1$ such that for any $t \geq t_1$, with high probability, for all $i \in I_t$,

$$\sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) \leq L_{t,i} + 3 \ln^2 t \sqrt{t}.$$

For any $t \geq 0$, note that we have $\hat{\ell}_t = \mathbb{E}[\ell(\hat{Y}_t, Y_t) \mid \mathbb{Y}_{\leq t}]$. We define the instantaneous regret $r_{t,i} = \hat{\ell}_t - \ell(y^i, Y_t)$. We now define $w'_{t-1,i} := e^{\eta_{t-1}(\hat{L}_{t-1,i} - L_{t-1,i})}$ and pose $W_{t-1} = \sum_{i \in I_t} w_{t-1,i}$ and $W'_{t-1} = \sum_{i \in I_{t-1}} w'_{t-1,i}$, i.e., which induces the most regret. We also denote the index $k_t \in I_t$ such that $\hat{L}_{t,k_t} - L_{t,k_t} = \max_{i \in I_t} \hat{L}_{t,i} - L_{t,i}$. We first note that for any $i, j \in I_t$, we have $\ell(y^i, Y_t) - \ell(y^j, Y_t) \leq \ell(y^i, y^0) + \ell(y^0, y^j) \leq 2 \ln t$. Therefore, we also have $|r_{t,i}| \leq 2 \ln t$. Hence, we can apply Hoeffding's lemma to obtain

$$\frac{1}{\eta_t} \ln \frac{W'_t}{W_{t-1}} = \frac{1}{\eta_t} \ln \sum_{i \in I_t} \frac{w_{t-1,i}}{W_{t-1}} e^{\eta_t r_{t,i}} \leq \frac{1}{\eta_t} \left(\eta_t \sum_{i \in I_t} r_{t,i} \frac{w_{t-1,i}}{W_{t-1}} + \frac{\eta_t^2 (4 \ln t)^2}{8} \right) = 2 \eta_t \ln^2 t.$$

The same computations as in the proof of Lemma 4.2 then show that

$$\begin{aligned} \frac{1}{\eta_t} \ln \frac{w_{t-1,k_{t-1}}}{W_{t-1}} - \frac{1}{\eta_{t+1}} \ln \frac{w_{t,k_t}}{W_t} & \leq 2 \left(\frac{1}{\eta_{t+1}} - \frac{1}{\eta_t} \right) \ln(1 + \ln(t+1)) + \frac{|I_{t+1}| - |I_t|}{\eta_t \sum_{i \in I_t} w_{t,i}} \\ & \quad + (\hat{L}_{t-1,k_{t-1}} - L_{t-1,k_{t-1}}) - (\hat{L}_{t,k_t} - L_{t,k_t}) + 2 \eta_t \ln^2 t. \end{aligned} \quad (4)$$

First suppose that we have $\sum_{i \in I_t} w_{t,i} \leq 1$. Similarly to Lemma 4.2, we obtain $\hat{L}_{t,k_t} - L_{t,k_t} \leq 0$. Otherwise, let $t' = \min\{1 \leq s \leq t : \forall s \leq s' \leq t, \sum_{i \in I_{s'}} w_{s',i} \geq 1\}$. We sum equation (4) for $s = t', \dots, t$ which gives

$$\begin{aligned} \frac{1}{\eta_1} \ln \frac{w_{t'-1,k_{t'-1}}}{W_{t'-1}} - \frac{1}{\eta_{t+1}} \ln \frac{w_{t,k_t}}{W_t} & \leq \frac{2}{\eta_{t+1}} \ln(1 + \ln(t+1)) + \frac{|I_{t+1}|}{\eta_t} \\ & \quad + (\hat{L}_{t'-1,k_{t'-1}} - L_{t'-1,k_{t'-1}}) - (\hat{L}_{t,k_t} - L_{t,k_t}) + 2 \sum_{s=t'}^t \eta_s \ln^2 s. \end{aligned}$$

Similarly as in Lemma 4.2, we have $\frac{w_{t,k_t}}{W_t} \leq 1$, $\frac{w_{t'-1,k_{t'-1}}}{W_{t'-1}} \geq \frac{1}{1+\ln t}$ and $\hat{L}_{t'-1,k_{t'-1}} - L_{t'-1,k_{t'-1}} \leq 0$. Finally, using the fact that $\sum_{s=1}^t \frac{1}{\sqrt{s}} \leq 2\sqrt{t}$, we obtain

$$\hat{L}_{t,k_t} - L_{t,k_t} \leq \ln(1 + \ln(t+1))(4 + 8\sqrt{t+1}) + 4(1 + \ln(t+1))\sqrt{t} + \ln^2 t\sqrt{t} \leq 2\ln^2 t\sqrt{t},$$

for all $t \geq t_0$ where t_0 is a fixed constant, and as a result, for all $t \geq t_0$ and $i \in I_t$, we have $\hat{L}_{t,i} - L_{t,i} \leq 2\ln^2 t\sqrt{t}$.

Now note that $|\ell(\hat{Y}_t, Y_t) - \mathbb{E}[\ell(\hat{Y}_t, Y_t) | \mathbb{Y}_{\leq t}]| \leq 2\ln t$ because for all $i \in I_t$, we have $\ell(y^i, y^0) \leq \ln t$. Hence, we can apply Hoeffding-Azuma inequality to the variables $\ell(\hat{Y}_t, Y_t) - \hat{\ell}_t$ that form a sequence of differences of a martingale, which yields

$$\mathbb{P}\left[\sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) > \hat{L}_{t,i} + u\right] \leq e^{-\frac{u^2}{8t\ln^2 t}}.$$

Hence, for $t \geq t_0$ and $i \in I_t$, with probability $1 - \delta$, we have

$$\sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) \leq \hat{L}_{t,i} + \ln t \sqrt{2t \ln \frac{1}{\delta}} \leq L_{t,i} + 2\ln^2 t\sqrt{t} + \ln t \sqrt{2t \ln \frac{1}{\delta}}.$$

Therefore, since $|I_t| \leq 1 + \ln t$, by union bound with probability $1 - \frac{1}{t^2}$ we obtain that for all $i \in I_t$,

$$\sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) \leq L_{t,i} + 2\ln^2 t\sqrt{t} + \ln t \sqrt{2t \ln(1 + \ln t)} + \ln t \sqrt{4t \ln t} \leq 3\ln^2 t\sqrt{t}$$

for all $t \geq t_1$ where $t_1 \geq t_0$ is a fixed constant. Now because $\sum_{t \geq 1} \frac{1}{t^2} < \infty$, the Borel-Cantelli lemma implies that almost surely, there exists $\hat{t} \geq 0$ such that

$$\forall t \geq \hat{t}, \forall i \in I_t, \quad \sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) \leq L_{t,i} + 3\ln^2 t\sqrt{t}.$$

We denote by \mathcal{A} this event. Now let $y \in \mathcal{Y}$, $\epsilon > 0$ and consider $i \geq 0$ such that $\ell(y^i, y) < \epsilon$. On the event \mathcal{A} , we have for all $t \geq \max(\hat{t}, t_i)$,

$$\sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) \leq \sum_{s=t_i}^t \ell(y^i, Y_s) + 3\ln^2 t\sqrt{t} \leq \sum_{s=t_i}^t \ell(y, Y_s) + \epsilon t + 3\ln^2 t\sqrt{t}.$$

Therefore, $\limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{s=1}^t \left(\ell(\hat{Y}_s, Y_s) - \ell(y, Y_s) \right) \leq \epsilon$ on \mathcal{A} . Because this holds for any $\epsilon > 0$ we finally obtain $\limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{s=1}^t \left(\ell(\hat{Y}_s, Y_s) - \ell(y, Y_s) \right) \leq 0$ on the event \mathcal{A} of probability one, which holds for all $y \in \mathcal{Y}$ simultaneously. This ends the proof of the theorem.

D.2 Proof of Corollary 6.2

We denote by g the learning rule on values \mathcal{Y} for mean estimation described in Theorem 3.6. Because processes in $\mathbb{X} \in \text{FS}$ visit only finite number of different instance points in \mathcal{X} almost surely, we can simply perform the learning rule g on each sub-process $\mathbb{Y}_{\{t: X_t=x\}}$ separately for any $x \in \mathcal{X}$. Note that the learning rule g does not explicitly re-use past randomness for its prediction. Hence, we will not specify that the randomness used for all learning rules—for each x visited by \mathbb{X} —should be independent. Let us formally describe our learning rule. Consider a sequence $\mathbf{x}_{\leq t-1}$ of instances in \mathcal{X} and $\mathbf{y}_{\leq t-1}$ of values in \mathcal{Y} . We denote by $S_{t-1} = \{x : \mathbf{x}_{\leq t-1} \cap \{x\} \neq \emptyset\}$ the support of $\mathbf{x}_{\leq t-1}$. Further, for any $x \in S_{t-1}$, we denote $N_{t-1}(x) = \sum_{u \leq t-1} \mathbb{1}_{x_u=x}$ the number of times that the specific instance x was visited by the sequence $\mathbf{x}_{\leq t-1}$. Last, for any $x \in S_{t-1}$, we denote $\mathbf{y}_{\leq N(x)}^x$ the values $\mathbf{y}_{\{u \leq t: X_u=x\}}$ obtained when the instance was

precisely x in the sequence $\mathbf{x}_{\leq t-1}$, ordered by increasing time u . We are now ready to define our learning rule f_t at time t . Given a new instance point x_t , we pose

$$f_t(\mathbf{x}_{\leq t-1}, \mathbf{y}_{\leq t-1}, x_t) = \begin{cases} g_{N_{t-1}(x)+1}(\mathbf{y}_{\leq N_{t-1}(x)}^x) & \text{if } x_t \in S_{t-1}, \\ g_1(\emptyset) & \text{otherwise.} \end{cases}$$

Recall that for any $u \geq 1$, g_u uses some randomness. The only subtlety is that at each iteration $t \geq 1$ of the learning rule f , the randomness used by the subroutine call to g should be independent from the past history. We now show that f is universally consistent for adversarial regression under all processes $\mathbb{X} \in \text{FS}$.

Let $\mathbb{X} \in \text{FS}$. For simplicity, we will denote by \hat{Y}_t the prediction of the learning rule f at time t . We denote $S = \{x : \{x\} \cap \mathbb{X} \neq \emptyset\}$ the random support of \mathbb{X} . By hypothesis, we have $|S| < \infty$ with probability one. Denote by \mathcal{E} this event. We now consider a specific realization \mathbf{x} of \mathbb{X} falling in the event \mathcal{E} . Then, S is a fixed set. We also denote $\tilde{S} := \{x \in S : \lim_{t \rightarrow \infty} N_t(x) = \infty\}$ the instances which are visited an infinite number of times by the sequence \mathbf{x} . Now, we can write for any function $f : \mathcal{X} \rightarrow \mathcal{Y}$,

$$\begin{aligned} \sum_{t=1}^T \left(\ell(\hat{Y}_t, Y_t) - \ell(f(x_t), Y_t) \right) &= \sum_{x \in S} \sum_{u=1}^{N_t(x)} \left(\ell(g_u(\mathbb{Y}_{\leq u-1}^x), Y_u^x) - \ell(f(x), Y_u) \right) \\ &\leq \sum_{s \in S \setminus \tilde{S}} \bar{\ell} |\{t \geq 1 : x_t = x\}| + \sum_{s \in \tilde{S}} \sum_{u=1}^{N_t(x)} \left(\ell(g_u(\mathbb{Y}_{\leq u-1}^x), Y_u^x) - \ell(f(x), Y_u) \right). \end{aligned}$$

Now, because the randomness in g was taken independently from the past at each iteration, we can apply directly Theorem 3.6. For all $x \in \tilde{S}$, with probability one, for all $y^x \in \mathcal{Y}$,

$$\limsup_{t' \rightarrow \infty} \frac{1}{t'} \sum_{u=1}^{t'} \left(\ell(g_u(\mathbb{Y}_{\leq u-1}^x), Y_u^x) - \ell(y^x, Y_u) \right) \leq 0.$$

We denote by \mathcal{E}_x this event. Then, on the event $\bigcap_{x \in \tilde{S}} \mathcal{E}_x$ of probability one, we have for any measurable function $f : \mathcal{X} \rightarrow \mathcal{Y}$,

$$\begin{aligned} \limsup_{T \rightarrow \infty} \frac{1}{T} \left(\ell(\hat{Y}_T, Y_T) - \ell(f(x_T), Y_T) \right) &\leq \sum_{s \in \tilde{S}} \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{u=1}^{N_t(x)} \left(\ell(g_u(\mathbb{Y}_{\leq u-1}^x), Y_u^x) - \ell(f(x), Y_u) \right) \\ &\leq \sum_{s \in \tilde{S}} \limsup_{T \rightarrow \infty} \frac{1}{N_t(x)} \sum_{u=1}^{N_t(x)} \left(\ell(g_u(\mathbb{Y}_{\leq u-1}^x), Y_u^x) - \ell(f(x), Y_u) \right) \leq 0. \end{aligned}$$

As a result, averaging on realisations of \mathbb{X} , we obtain that with probability one, we have that $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f) \leq 0$ for all measurable functions $f : \mathcal{X} \rightarrow \mathcal{Y}$. Note that this is stronger than the notion of universal consistency which we defined in Section 2, where we ask that for all measurable function $f : \mathcal{X} \rightarrow \mathcal{Y}$, we have almost surely $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f) \leq 0$. In particular, this shows that $\text{FS} \subset \text{SOLAR-U}$. As result $\text{SOLAR-U} = \text{FS}$ and f is optimistically universal. This ends the proof of the result.

D.3 Proof of Theorem 6.3

We first show that mean-estimation is not achievable. To do so, let f be a learning rule. For simplicity, we will denote by \hat{Y}_t its prediction at step t . We aim to construct a process \mathbb{Y} on \mathbb{R} and a value $y^* \in \mathbb{R}$ such that with non-zero probability we have

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(f_t(\mathbb{Y}_{\leq t-1}), Y_t) - \ell(y^*, Y_t) > 0.$$

We now pose $\beta := \frac{2\alpha}{\alpha-1} > 2$. For any sequence $\mathbf{b} := (b_t)_{t \geq 1}$ in $\{-1, 1\}$, we consider the following process $\mathbb{Y}^{\mathbf{b}}$ such that for any $t \geq 1$ we have $Y_t^{\mathbf{b}} = 2^{\beta t} b_t$. Let $\mathbf{B} := (B_t)_{t \geq 1}$ be an i.i.d. sequence of Rademacher random variables, i.e., such that $B_1 = 1$ (resp. $B_1 = -1$) with probability $\frac{1}{2}$. We consider the random variables $e_t := \mathbb{1}_{\hat{Y}_t \cdot Y_t \leq 0}$ which intuitively correspond to flags for large mistakes of the learning rule f . at time t . Because f is an online learning rule, we have

$$\mathbb{E}[e_t \mid \mathbb{Y}_{\leq t-1}] = \mathbb{E}_{\hat{Y}_t} \left[\mathbb{E}_{B_t} [\mathbb{1}_{\hat{Y}_t \cdot Y_t \leq 0} \mid \hat{Y}_t] \right] = \mathbb{E}_{\hat{Y}_t} \left[\mathbb{1}_{\hat{Y}_t=0} + \frac{1}{2} \mathbb{1}_{\hat{Y}_t \neq 0} \right] \geq \frac{1}{2}.$$

where the expectation $\mathbb{E}_{\hat{Y}_t}$ refers to the expectation on the randomness of the rule f_t . As a result, the random variables $e_t - \frac{1}{2}$ form a sequence of differences of a sub-martingale bounded by $\frac{1}{2}$ in absolute value. By the Azuma-Hoeffding inequality, we obtain $\mathbb{P} \left[\sum_{t=1}^T e_t \leq \frac{T}{4} \right] \leq e^{-T/8}$. Because $\sum_{t \geq 1} e^{-t/8} < \infty$, the Borel-Cantelli lemma implies that on an event \mathcal{E} of probability one, we have $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T e_t \geq \frac{1}{4}$. As a result, there exists a specific realization \mathbf{b} of \mathbf{B} such that on an event $\tilde{\mathcal{E}}$ of probability one, we have $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T e_t \geq \frac{1}{4}$. Note that the sequence $\mathbb{Y}^{\mathbf{b}}$ is now deterministic. Then, writing $e_t = e_t \mathbb{1}_{Y_t > 0} + e_t \mathbb{1}_{Y_t < 0}$, we obtain

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T e_t \mathbb{1}_{Y_t > 0} + \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T e_t \mathbb{1}_{Y_t < 0} \geq \frac{1}{4}.$$

Without loss of generality, we can suppose that $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{1}_{\hat{Y}_t \cdot Y_t \leq 0} \mathbb{1}_{Y_t > 0} \geq \frac{1}{8}$. We now pose $y^* = 1$. In the other case, we pose $y^* = -1$. We now compute for any $T \geq 1$ such that $\hat{Y}_t \cdot Y_t \leq 0$ and $Y_t > 0$,

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T (\ell(f_t(\mathbb{Y}_{\leq t-1}), Y_t) - \ell(y^*, Y_t)) &\geq \frac{\ell(0, 2^{\beta T}) - \ell(1, 2^{\beta T})}{T} - \frac{1}{T} \sum_{t=1}^{T-1} \ell(1, -2^{\beta t}). \\ &= \frac{\alpha}{T} 2^{(\alpha-1)\beta T} + O\left(\frac{1}{T} 2^{(\alpha-2)\beta T}\right) - 2^{\alpha(1+\beta T-1)} \\ &= \frac{\alpha}{T} 2^{2\alpha\beta T-1} (1 + o(1)). \end{aligned}$$

Because, by construction $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{1}_{\hat{Y}_t \cdot Y_t \leq 0} \mathbb{1}_{Y_t > 0} \geq \frac{1}{8}$, we obtain

$$\limsup \frac{1}{T} \sum_{t=1}^T (\ell(f_t(\mathbb{Y}_{\leq t-1}), Y_t) - \ell(y^*, Y_t)) = \infty,$$

on the event $\tilde{\mathcal{E}}$ of probability one. This ends the proof that mean-estimation is not achievable. Because mean-estimation is the easiest regression setting, this directly implies $\text{SOLAR-U} = \emptyset$. Formally, let \mathbb{X} a process on \mathcal{X} . and f a learning rule for regression. We consider the same processes $\mathbb{Y}^{\mathbf{B}}$ where \mathbf{B} is i.i.d. Rademacher and independent from \mathbb{X} . The same proof shows that there exists a realization \mathbf{b} for which we have almost surely $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f^* := y^*) = \infty$, where $f^* = y^*$ denotes the constant function equal to y^* where $y^* \in \mathbb{R}$ is the value constructed as above. Hence, $\mathbb{X} \notin \text{SOLAR-U}$, and as a result, $\text{SOLAR-U} = \emptyset$.

D.4 Proof of Proposition 6.4

Suppose that there exists an online learning rule g for mean-estimation. In the proof of Corollary 6.2, instead of using the learning rule for mean-estimation for metric losses introduced in Theorem 3.6, we can use the learning rule g to construct the learning rule f for adversarial regression on FS instance processes, which simply performs f separately on each subprocess $\mathbb{Y}_{t: X_t=x}$ with the same instance $x \in \mathcal{X}$ for all visited $x \in \mathcal{X}$ in the process \mathbb{X} . The same proof shows that because almost surely \mathbb{X} visits a finite number of

different instances, f is universally consistent under any process $\mathbb{X} \in \text{FS}$. Hence, $\text{FS} \subset \text{SOLAR-U}$. Because $\text{SOLAR-U} \subset \text{SOUL} = \text{FS}$, we obtain directly $\text{SOLAR-U} = \text{FS}$ and f is optimistically universal.

On the other hand, if mean-estimation with adversarial responses is not achievable, we can use similar arguments as for the proof of Theorem 6.3. Let f a learning rule for regression, and consider the following learning rule g for mean-estimation. We first draw a process $\tilde{\mathbb{X}}$ with same distribution as \mathbb{X} . Then, we pose

$$g_t(\mathbf{y}_{\leq t-1}) := f_t(\tilde{\mathbb{X}}_{\leq t-1}, \mathbf{y}_{\leq t-1}, \tilde{X}_t).$$

Then, because mean-estimation is not achievable, there exists an adversarial process \mathbb{Y} on (\mathcal{Y}, ℓ) such that with non-zero probability,

$$\limsup \frac{1}{T} \sum_{t=1}^T (\ell(g_t(\mathbb{Y}_{\leq t-1}), Y_t) - \ell(y^*, Y_t)) > 0.$$

Then, we obtain that with non-zero probability, $\mathcal{L}_{(\tilde{\mathbb{X}}, \mathbb{Y})} > 0$. Hence, f is not universally consistent. Note that the “bad” process \mathbb{Y} is not correlated with $\tilde{\mathbb{X}}$ in this construction.

E Proofs of Section 7

E.1 Proof of Theorem 7.1

Let $(x^k)_{k \geq 0}$ a sequence of distinct points of \mathcal{X} . Now fix a value $y_0 \in \mathcal{Y}$ and construct a sequence of values y_k^1, y_k^2 for $k \geq 1$ such that $\ell(y_k^1, y_k^2) \geq c_\ell 2^{k+1}$. Because $\ell(y_k^1, y_k^2) \leq c_\ell \ell(y_0, y_k^1) + c_\ell \ell(y_0, y_k^2)$, there exists $i_k \in \{1, 2\}$ such that $\ell(y_0, y_k^{i_k}) \geq 2^k$. For simplicity, we will now write $y_k := y_k^{i_k}$ for all $k \geq 1$. We define

$$t_k = \left\lceil \sum_{l=1}^k \ell(y_0, y_l) \right\rceil.$$

This forms an increasing sequence of times because $t_{k+1} - t_k \geq \ell(y_0, y_{k+1}) \geq 1$. Consider the deterministic process \mathbb{X} that visits x^k at time t_k and x^0 otherwise, i.e., such that

$$X_t = \begin{cases} x^k & \text{if } t = t_k, \\ x^0 & \text{otherwise.} \end{cases}$$

The process \mathbb{X} visits $\mathcal{X} \setminus \{x^0\}$ a sublinear number of times. Hence we have for any measurable set A :

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{1}_A(X_t) = \begin{cases} 1 & \text{if } x^0 \in A \\ 0 & \text{otherwise.} \end{cases}$$

As a result, $\mathbb{X} \in \text{CRF}$. We will now show that universal learning under \mathbb{X} with the first moment condition on the responses is not achievable. For any sequence $b := (b_k)_{k \geq 1}$ of binary variables $b_k \in \{0, 1\}$, we define the function $f_b^* : \mathcal{X} \rightarrow \mathcal{Y}$ such that

$$f_b^*(x^k) = \begin{cases} y_0 & \text{if } b_k = 0, \\ y_k & \text{otherwise,} \end{cases} \quad k \geq 0 \quad \text{and} \quad f_b^*(x) = y_0 \text{ if } x \notin \{x_k, k \geq 0\}.$$

These functions are simple, hence measurable. We will first show that for any binary sequence b , the function f_b^* satisfies the moment condition on the target functions. Indeed, we note that for any $T \geq t_1$, with $k := \max\{l \geq 1 : t_l \leq T\}$, we have

$$\frac{1}{T} \sum_{t=1}^T \ell(y_0, f_b^*(X_t)) \leq \frac{1}{T} \sum_{l=1}^k \ell(y_0, y_l) \leq \frac{t_k + 1}{T} \leq \frac{T + 1}{T}.$$

Therefore, $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, f_b^*(X_t)) \leq 1$. We now consider any online learning rule f . Let $B = (B_k)_{k \geq 1}$ be an i.i.d. sequence of Bernoulli variables independent from the learning rule randomness. For any $k \geq 1$, denoting by $\hat{Y}_{t_k} := f_{t_k}(\mathbb{X}_{\leq t_k-1}, f_B^*(\mathbb{X}_{\leq t_k-1}), X_{t_k})$ we have

$$\mathbb{E}_{B_k} \ell(\hat{Y}_{t_k}, f_B^*(X_{t_k})) = \frac{\ell(\hat{Y}_{t_k}, y_0) + \ell(\hat{Y}_{t_k}, y_k)}{2} \geq \frac{1}{2c_\ell} \ell(y_0, y_k).$$

In particular, taking the expectation over both B and the learning rule, we obtain

$$\mathbb{E} \left[\frac{1}{t_k} \sum_{t=1}^{t_k} \ell(f_t(\mathbb{X}_{\leq t-1}, f_B^*(\mathbb{X}_{\leq t-1}), X_t), f_B^*(X_t)) \right] \geq \frac{1}{2c_\ell t_k} \sum_{l=1}^k \ell(y_0, y_k) \geq \frac{1}{2c_\ell}.$$

As a result, using Fatou's lemma we obtain

$$\begin{aligned} & \mathbb{E} \left[\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(f_t(\mathbb{X}_{\leq t-1}, f_B^*(\mathbb{X}_{\leq t-1}), X_t), f_B^*(X_t)) \right] \\ & \geq \limsup_{T \rightarrow \infty} \mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \ell(f_t(\mathbb{X}_{\leq t-1}, f_B^*(\mathbb{X}_{\leq t-1}), X_t), f_B^*(X_t)) \right] \\ & \geq \frac{1}{2c_\ell}. \end{aligned}$$

Therefore, the learning rule f is not consistent under \mathbb{X} for all target functions of the form f_b^* for some sequence of binary variables b . Indeed, otherwise for all binary sequence $b = (b_k)_{k \geq 1}$, we would have $\mathbb{E}_{\mathbb{X}} \left[\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(f_t(\mathbb{X}_{\leq t-1}, f_b^*(\mathbb{X}_{\leq t-1}), X_t), f_b^*(X_t)) \right] = 0$ and as a result

$$\mathbb{E}_B \mathbb{E}_{\mathbb{X}} \left[\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(f_t(\mathbb{X}_{\leq t-1}, f_B^*(\mathbb{X}_{\leq t-1}), X_t), f_B^*(X_t)) \right] = 0.$$

This ends the proof of the theorem.

E.2 Proof of Lemma 7.3

It suffices to prove that empirical integrability implies the latter property. We pose $\epsilon_i = 2^{-i}$ for any $i \geq 0$. By definition, there exists an event \mathcal{E}_i of probability one such that on \mathcal{E}_i we have

$$\exists M_i \geq 0, \quad \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M_i} \leq \epsilon_i.$$

As a result, on $\bigcap_{i \geq 0} \mathcal{E}_i$ of probability one, we obtain

$$\forall \epsilon > 0, \exists M := M_{\lceil \log_2 \frac{1}{\epsilon} \rceil} \geq 0, \quad \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M} \leq \epsilon.$$

This ends the proof of the lemma.

E.3 Proof of Theorem 3.1

Let $\mathbb{X} \in \text{SOUL}$ and $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ such that $f^*(\mathbb{X})$ is empirically integrable. By Lemma 7.3, there exists some value $y_0 \in \mathcal{Y}$ such that on an event \mathcal{A} of probability one, for all $\epsilon > 0$ there exists $M_\epsilon \geq 0$ such that

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, f^*(X_t)) \mathbb{1}_{\ell(y_0, f^*(X_t)) \geq M_\epsilon} \leq \epsilon.$$

For any $M \geq 1$ we define the function f_M^* by

$$f_M^*(x) = \begin{cases} f^*(x) & \text{if } \ell(y_0, f^*(x)) \leq M, \\ y_0 & \text{otherwise.} \end{cases}$$

We know that 2C1NN is optimistically universal in the noiseless setting for bounded losses. Therefore, restricting the study to the output space $(B_\ell(y_0, M), \ell)$ we obtain that 2C1NN is consistent for f_M^* under \mathbb{X} , i.e.

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(2C1NN_t(\mathbb{X}_{t-1}, f_M^*(\mathbb{X}_{\leq t-1}), X_t), f_M^*(X_t)) = 0 \quad (a.s.).$$

For any $t \geq 1$, we denote $\phi(t)$ the representative used by the 2C1NN learning rule. We denote \mathcal{E}_M the above event such that $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(f_M^*(X_{\phi(t)}), f_M^*(X_t)) = 0$. We now write for any $T \geq 1$ and $M \geq 1$,

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \ell(f^*(X_{\phi(t)}), f^*(X_t)) &\leq \frac{c_\ell^2}{T} \sum_{t=1}^T \ell(f_M^*(X_{\phi(t)}), f_M^*(X_t)) + \frac{c_\ell^2}{T} \sum_{t=1}^T \ell(f^*(X_t), f_M^*(X_t)) \\ &\quad + \frac{c_\ell}{T} \sum_{t=1}^T \ell(f^*(X_{\phi(t)}), f_M^*(X_{\phi(t)})). \end{aligned}$$

We now note that by construction of the 2C1NN learning rule,

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \ell(f^*(X_{\phi(t)}), f_M^*(X_{\phi(t)})) &= \frac{1}{T} \sum_{u=1}^T \ell(f^*(X_u), f_M^*(X_u)) |\{u < t \leq T : \phi(t) = u\}| \\ &\leq \frac{2}{T} \sum_{t=1}^T \ell(f^*(X_t), f_M^*(X_t)). \end{aligned}$$

Hence, we obtain

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \ell(f^*(X_{\phi(t)}), f^*(X_t)) &\leq \frac{c_\ell^2}{T} \sum_{t=1}^T \ell(f_M^*(X_{\phi(t)}), f_M^*(X_t)) \\ &\quad + \frac{c_\ell(2 + c_\ell)}{T} \sum_{t=1}^T \ell(y_0, f^*(X_t)) \mathbb{1}_{\ell(y_0, f^*(X_t)) > M}. \end{aligned}$$

As a result, on the event $\mathcal{A} \cap \bigcap_{M \geq 1} \mathcal{E}_M$ of probability one, for any $M \geq 1$, we obtain

$$\begin{aligned} \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(f^*(X_{\phi(t)}), f^*(X_t)) \\ \leq c_\ell(2 + c_\ell) \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, f^*(X_t)) \mathbb{1}_{\ell(y_0, f^*(X_t)) \geq M}. \end{aligned}$$

In particular, if $\epsilon > 0$ we can apply this result with $M := \lceil M_\epsilon \rceil$, which in turn shows that we have $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(f^*(X_{\phi(t)}), f^*(X_t)) \leq c_\ell(2 + c_\ell)\epsilon$. Because this holds for any $\epsilon > 0$ we finally obtain that on the event $\mathcal{A} \cap \bigcap_{M \geq 1} \mathcal{E}_M$ we have

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(f^*(X_{\phi(t)}), f^*(X_t)) = 0.$$

This ends the proof of the theorem.

E.4 Proof of Theorem 3.3

We first define the learning rule. Using Lemma 23 of [Han21a], let $\mathcal{T} \subset \mathcal{B}$ a countable set such that for all $\mathbb{X} \in \text{CS}, A \subset \mathcal{B}$ we have

$$\inf_{G \in \mathcal{T}} \mathbb{E}[\hat{\mu}_{\mathbb{X}}(G \triangle A)] = 0.$$

Now let $(y^i)_{i \geq 0}$ be a dense sequence in \mathcal{Y} . For any $k \geq 0$, any indices $l_1, \dots, l_k \in \mathbb{N}$ and any sets $A_1, \dots, A_k \in \mathcal{T}$, we define the function $f_{\{l_1, \dots, l_k\}, \{A_1, \dots, A_k\}} : \mathcal{X} \rightarrow \mathcal{Y}$ as

$$f_{\{l_1, \dots, l_k\}, \{A_1, \dots, A_k\}}(x) = y^{\max\{0 \leq j \leq k: x \in A_j\}}$$

where $A_0 = \mathcal{X}$. These functions are simple hence measurable. Because the set of such functions is countable, we enumerate these functions as f^0, f^1, \dots . Without loss of generality, we suppose that $f^0 = y^0$. For any $i \geq 0$, we denote $k^i \geq 0$, $\{l_1^i, \dots, l_{k^i}^i\}$ and $\{A_1^i, \dots, A_{k^i}^i\}$ such that f^i was defined as $f^i := f_{\{l_1^i, \dots, l_{k^i}^i\}, \{A_1^i, \dots, A_{k^i}^i\}}$. We now define a sequence of sets $(I_t)_{t \geq 1}$ of indices and a sequence of sets $(\mathcal{F}_t)_{t \geq 1}$ of measurable functions by

$$I_t := \{i \leq \ln t : \ell(y^{l_p^i}, y^0) \leq 2^{-\alpha+1} \ln t, \forall 1 \leq p \leq k^i\} \quad \text{and} \quad \mathcal{F}_t := \{f^i : i \in I_t\}.$$

Then, clearly I_t is finite and $\bigcup_{t \geq 1} I_t = \mathbb{N}$. For any $i \geq 0$, we define $t_i = \min\{t : i \in I_t\}$. We are now ready to construct our learning rule. Let $\eta_t = \frac{1}{\ln t \sqrt{t}}$. Fix any sequences $(x_t)_{t \geq 1}$ in \mathcal{X} and $(y_t)_{t \geq 1}$ in \mathcal{Y} . At step $t \geq 1$, after observing the values x_i for $1 \leq i \leq t$ and y_i for $1 \leq i \leq t-1$, we define for any $i \in I_t$ the loss $L_{t-1,i} := \sum_{s=t_i}^{t-1} \ell(f^i(x_s), y_s)$. For any $M \geq 1$ we define the function $\phi_M : \mathcal{Y} \rightarrow \mathcal{Y}$ such that

$$\phi_M(y) = \begin{cases} y & \text{if } \ell(y, y^0) < M, \\ y^0 & \text{otherwise.} \end{cases}$$

We now construct some weights $w_{t,i}$ for $t \geq 1$ and $i \in I_t$ recursively in the following way. Note that $I_1 = \{0\}$. Therefore, we pose $w_{0,0} = 1$. Now let $t \geq 2$ and suppose that $w_{s-1,i}$ have been constructed for all $1 \leq s \leq t-1$. We define

$$\hat{\ell}_s := \frac{\sum_{j \in I_s} w_{s-1,j} \ell(f^j(x_s), \phi_{2^{-\alpha+1} \ln s}(y_s))}{\sum_{j \in I_s} w_{s-1,j}}$$

and for any $i \in I_t$ we note $\hat{L}_{t-1,i} := \sum_{s=t_i}^{t-1} \hat{\ell}_s$. In particular, if $t_i = t$ we have $\hat{L}_{t-1,i} = L_{t-1,i} = 0$. The weights at time t are constructed as $w_{t-1,i} := e^{\eta_t(\hat{L}_{t-1,i} - L_{t-1,i})}$ for any $i \in I_t$. Last, let $\{\hat{i}_t\}_{t \geq 1}$ a sequence of independent random \mathbb{N} -valued variables such that

$$\mathbb{P}(\hat{i}_t = i) = \frac{w_{t-1,i}}{\sum_{j \in I_t} w_{t-1,j}}, \quad i \in I_t.$$

Finally, the prediction is defined as $\hat{y}_t := f^{\hat{i}_t}(x_t)$. The learning rule is summarized in Algorithm 6.

For simplicity, we will refer to the predictions of the learning rule as $(\hat{Y}_t)_{t \geq 1}$. Now consider a process (\mathbb{X}, \mathbb{Y}) with $\mathbb{X} \in \text{CS}$ and such that \mathbb{Y} is empirically integrable. By Lemma 7.3, there is $y_0 \in \mathcal{Y}$ such that on an event \mathcal{A} of probability one, for any $\epsilon > 0$, there exists $M_\epsilon \geq 0$ with $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M_\epsilon} \leq \epsilon$. We will now denote $\tilde{\mathbb{Y}}$ the process defined by $\tilde{Y}_t = \phi_{2^{-\alpha+1} \ln t}(Y_t)$ for all $t \geq 1$. Then, for any $i \in I_t$, note that using Lemma A.1 we have

$$0 \leq \ell(f^i(x_t), \tilde{Y}_t) \leq 2^{\alpha-1} \left(\ell(f^i(x_t), y^0) + \ell(y^0, \tilde{Y}_t) \right) \leq 2 \ln t,$$

by construction of the set I_t . As a result, for any $i, j \in I_t$, we obtain $|\ell(f^i(x_t), \tilde{Y}_t^M) - \ell(f^j(x_t), \tilde{Y}_t^M)| \leq 2 \ln t$. Hence, we can use the same proof as for Theorem 3.6 and show that almost surely, there exists $\hat{t} \geq 1$ such that

$$\forall t \geq \hat{t}, \forall i \in I_t, \quad \sum_{s=t_i}^t \ell(\hat{Y}_s, \tilde{Y}_s^M) \leq L_{t,i} + 3 \ln^2 t \sqrt{t}.$$

Input: Historical samples $(X_t, Y_t)_{t < T}$ and new input point X_T
Output: Predictions \hat{Y}_t for $t \leq T$
Construct the sequence of measurable functions $\{f^i, i \geq 0\}$ with $f^i = f_{\{l_1^i, \dots, l_k^i\}, \{A_1^i, \dots, A_k^i\}}$
 $I_t := \{i \leq \ln t, \ell(y^{l_p^i}, y^0) \leq 2^{-\alpha+1} \ln t, \forall 1 \leq p \leq k^i\}, \mathcal{F}_t := \{f^i, i \in I_t\}, \eta_t := \frac{1}{\ln t \sqrt{t}}, t \geq 1$
 $t_i = \min\{t : i \in I_t\}, i \geq 0$
 $w_{0,0} := 1, \hat{Y}_1 = y^0 (= f^0(X_0))$ // Initialisation
for $t = 2, \dots, T$ **do**
 $L_{t-1,i} = \sum_{s=t_i}^{t-1} \ell(f^i(X_s), \phi_{2^{-\alpha+1} \ln t}(Y_s)), \hat{L}_{t-1,i} = \sum_{s=t_i}^{t-1} \hat{\ell}_s, i \in I_t$
 $w_{t-1,i} := \exp(\eta_t (\hat{L}_{t-1,i} - L_{t-1,i})), i \in I_t$
 $p_t(i) = \frac{w_{t-1,i}}{\sum_{j \in I_t} w_{t-1,j}}, i \in I_t$
 $\hat{i}_t \sim p_t(\cdot)$ // Function selection
 $\hat{Y}_t = f^{\hat{i}_t}(X_t)$
 $\hat{\ell}_t := \frac{\sum_{j \in I_t} w_{t-1,j} \ell(f^j(X_s), \phi_{2^{-\alpha+1} \ln t}(Y_t))}{\sum_{j \in I_t} w_{t-1,j}}$
end

Algorithm 6: A learning rule for adversarial empirically integrable responses under CS processes.

We denote by \mathcal{B} this event. Now let $f : \mathcal{X} \rightarrow \mathcal{Y}$ to which we compare the predictions of our learning rule. For any $M \geq 1$, the function $\phi_M \circ f$ is measurable and has values in the ball $B_\ell(y_0, M)$ where the loss is bounded by $2^\alpha M$. Hence, by Lemma 24 from [Han21a] because $\mathbb{X} \in \mathcal{C}_1$ we have

$$\inf_{i \geq 0} \mathbb{E} [\hat{\mu}_{\mathbb{X}}(\ell(\phi_M \circ f(\cdot), f^i(\cdot)))] = 0.$$

Now for any $k \geq 0$, let $i_k \geq 0$ such that $\mathbb{E} [\hat{\mu}_{\mathbb{X}}(\ell(\phi_M \circ f(\cdot), f^{i_k}(\cdot)))] < 2^{-2k}$. By Markov inequality, we have

$$\mathbb{P} [\hat{\mu}_{\mathbb{X}}(\ell(\phi_M \circ f(\cdot), f^{i_k}(\cdot))) < 2^{-k}] \geq 1 - 2^{-k}.$$

Because $\sum_k 2^{-k} < \infty$, the Borel-Cantelli lemma implies that almost surely there exists \hat{k} such that for any $k \geq \hat{k}$, the above inequality is met. We denote \mathcal{E}_M this event. On the event $\mathcal{B} \cap \mathcal{E}_M$ of probability one, for $k \geq \hat{k}$ and any $T \geq \max(t_{i_k}, \hat{t})$ we have for any $\epsilon > 0$,

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \left(\ell(\hat{Y}_t, \tilde{Y}_t) - \ell(\phi_M \circ f(X_t), \tilde{Y}_t) \right) \\ &= \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, \tilde{Y}_t) - \ell(f^{i_k}(X_t), \tilde{Y}_t) + \frac{1}{T} \sum_{t=1}^T \ell(f^{i_k}(X_t), \tilde{Y}_t) - \ell(\phi_M \circ f(X_t), \tilde{Y}_t) \\ &\leq \frac{1}{T} \sum_{t=1}^{t_{i_k}-1} \ell(\hat{Y}_t, \tilde{Y}_t) + \frac{1}{T} \left(\sum_{t=t_{i_k}}^T \ell(\hat{Y}_t, \tilde{Y}_t) - L_{T, i_k} \right) + \frac{\epsilon}{T} \sum_{t=1}^T \ell(\phi_M \circ f(X_t), \tilde{Y}_t) \\ &\quad + \frac{c_\epsilon^\alpha}{T} \sum_{t=1}^T \ell(f^{i_k}(X_t), \phi_M \circ f(X_t)) \\ &\leq \frac{2 \ln t_{i_k}}{T} + \frac{3 \ln^2 T}{\sqrt{T}} + \epsilon 2^{\alpha-1} M + \epsilon 2^{\alpha-1} \frac{1}{T} \sum_{t=1}^T \ell(y^0, \tilde{Y}_t) + \frac{c_\epsilon^\alpha}{T} \sum_{t=1}^T \ell(f^{i_k}(X_t), \phi_M \circ f(X_t)) \\ &\leq \frac{2 \ln t_{i_k}}{T} + \frac{3 \ln^2 T}{\sqrt{T}} + \epsilon 2^{\alpha-1} M + \epsilon 2^{\alpha-1} \frac{1}{T} \sum_{t=1}^T \ell(y^0, Y_t) + \frac{c_\epsilon^\alpha}{T} \sum_{t=1}^T \ell(f^{i_k}(X_t), \phi_M \circ f(X_t)), \end{aligned}$$

where in the last inequality we used the inequality $\ell(y^0, \tilde{Y}_t) \leq \ell(y^0, Y_t)$ by construction of $\tilde{Y}_t = \phi_{2^{-\alpha+1} \ln t}(Y_t)$. Now on the event \mathcal{A} , we have

$$\begin{aligned} Z_1 &:= \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y^0, Y_t) \leq 2^{\alpha-1} \ell(y_0, y^0) + 2^{\alpha-1} \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \\ &\leq 2^{\alpha-1} \ell(y_0, y^0) + 2^{\alpha-1} \left(M_1 + \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M_1} \right) \\ &\leq 2^{\alpha-1} \ell(y_0, y^0) + 2^{\alpha-1} (M_1 + 1) < \infty. \end{aligned}$$

Thus, on the event $\mathcal{A} \cap \mathcal{B} \cap \mathcal{E}_M$, for any $k \geq \hat{k}$ we have for any $\epsilon > 0$,

$$\limsup_T \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, \tilde{Y}_t) - \ell(\phi_M \circ f(X_t), \tilde{Y}_t) \leq \epsilon 2^{\alpha-1} M + \epsilon 2^{\alpha-1} Z_1 + \frac{\epsilon^\alpha}{2^k}.$$

Let $\delta > 0$. Now taking $\epsilon = \frac{1}{2^{\alpha(M+Z_1)}}$, we obtain that on the event $\mathcal{A} \cap \mathcal{B} \cap \mathcal{E}_M$, for any $k \geq \hat{k}$, we have $\limsup_T \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, \tilde{Y}_t) - \ell(\phi_M \circ f(X_t), \tilde{Y}_t) \leq \delta + \frac{\epsilon^\alpha}{2^k}$. This yields $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, \tilde{Y}_t) - \ell(\phi_M \circ f(X_t), \tilde{Y}_t) \leq \delta$. Because this holds for any $\delta > 0$ we obtain $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, \tilde{Y}_t) - \ell(\phi_M \circ f(X_t), \tilde{Y}_t) \leq 0$. Finally, on the event $\mathcal{A} \cap \mathcal{B} \cap \bigcap_{M=1}^{\infty} \mathcal{E}_M$ of probability one, we have

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \left(\ell(\hat{Y}_t, \tilde{Y}_t) - \ell(\phi_M \circ f(X_t), \tilde{Y}_t) \right) \leq 0, \quad \forall M \geq 1,$$

where M is an integer. We now observe that on the event \mathcal{A} , the same guarantee for y_0 also holds for y^0 . Indeed, let ϵ . For $\tilde{M}_\epsilon := 2^{\alpha-1}(M_{2^{-\alpha\epsilon}} + \ell(y^0, y_0)) + \ell(y_0, y^0)$ we have

$$\begin{aligned} &\frac{1}{T} \sum_{t=1}^T \ell(y^0, Y_t) \mathbb{1}_{\ell(y^0, Y_t) \geq \tilde{M}_\epsilon} \\ &\leq 2^{\alpha-1} \ell(y^0, y_0) \frac{1}{T} \sum_{t=1}^T \mathbb{1}_{\ell(y^0, Y_t) \geq \tilde{M}_\epsilon} + 2^{\alpha-1} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y^0, Y_t) \geq \tilde{M}_\epsilon} \\ &\leq 2^{\alpha-1} \ell(y^0, y_0) \frac{1}{T} \sum_{t=1}^T \mathbb{1}_{\ell(y_0, Y_t) \geq 2^{-\alpha+1} M - \ell(y_0, y^0)} \\ &\quad + 2^{\alpha-1} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq 2^{-\alpha+1} M - \ell(y_0, y^0)} \\ &\leq 2^\alpha \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M_{2^{-\alpha\epsilon}}} \end{aligned}$$

Hence, we obtain $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y^0, Y_t) \mathbb{1}_{\ell(y^0, Y_t) \geq \tilde{M}_\epsilon} \leq \epsilon$. We now write

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T \ell(\phi_M \circ f(X_t), \tilde{Y}_t) - \ell(f(X_t), Y_t) \\
& \leq \frac{1}{T} \sum_{t=1}^T (\ell(y^0, Y_t) - \ell(f(X_t), Y_t)) \mathbb{1}_{\ell(f(X_t), y^0) \geq M} \mathbb{1}_{\ell(Y_t, y^0) \leq \ln t} \\
& \quad + \frac{1}{T} \sum_{t=1}^T (\ell(f(X_t), y^0) - \ell(f(X_t), Y_t)) \mathbb{1}_{\ell(f(X_t), y^0) \leq M} \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha+1} \ln t} \\
& \leq \frac{1}{T} \sum_{t=1}^T (2\ell(y^0, Y_t) - 2^{-\alpha+1} \ell(f(X_t), y^0)) \mathbb{1}_{\ell(f(X_t), y^0) \geq M} \\
& \quad + \frac{1}{T} \sum_{t=1}^T (2\ell(f(X_t), y^0) - 2^{-\alpha+1} \ell(y^0, Y_t)) \mathbb{1}_{\ell(f(X_t), y^0) \leq M} \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha+1} \ln t} \\
& \leq \frac{2}{T} \sum_{t=1}^T \ell(y^0, Y_t) \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha} M} + \frac{2Me^{2^{2\alpha-1}M}}{T}.
\end{aligned}$$

As a result, on the event $\mathcal{A} \cap \mathcal{B} \cap \bigcap_{M=1}^{\infty} \mathcal{E}_M$, for any $M \geq 1$,

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\phi_M \circ f(X_t), \tilde{Y}_t) - \ell(f(X_t), Y_t) \leq 2 \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y^0, Y_t) \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha} M}.$$

Last, we compute

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(\hat{Y}_t, \tilde{Y}_t) &= \frac{1}{T} \sum_{t=1}^T (\ell(\hat{Y}_t, Y_t) - \ell(\hat{Y}_t, y^0)) \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha+1} \ln t} \\
&\leq \frac{1}{T} \sum_{t=1}^T (2^{\alpha-1} \ell(\hat{Y}_t, y^0) + 2^{\alpha-1} \ell(Y_t, y^0)) \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha+1} \ln t} \\
&\leq \frac{1}{T} \sum_{t=1}^T (\ln t + 2^{\alpha-1} \ell(Y_t, y^0)) \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha+1} \ln t} \\
&\leq \frac{2^\alpha}{T} \sum_{t=1}^T \ell(Y_t, y^0) \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha+1} \ln t}.
\end{aligned}$$

Note that for any $\epsilon > 0$, we have on the event \mathcal{A} that for any $M \geq 1$,

$$\begin{aligned}
\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(Y_t, y^0) \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha+1} \ln t} &\leq \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t \geq e^{2^{\alpha-1}M}}^T \ell(Y_t, y^0) \mathbb{1}_{\ell(Y_t, y^0) \geq M} \\
&= \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(Y_t, y^0) \mathbb{1}_{\ell(Y_t, y^0) \geq M}.
\end{aligned}$$

Hence, because this holds for any $M \geq 1$, if $\epsilon > 0$ we can apply this to the integer $M := \lceil \tilde{M}_\epsilon \rceil$ which yields $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(Y_t, y^0) \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha+1} \ln t} \leq \epsilon$. This holds for any $\epsilon > 0$. Hence we obtain on the event \mathcal{A} that $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(Y_t, y^0) \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha+1} \ln t} \leq 0$, which implies that $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) -$

$\ell(\hat{Y}_t, \tilde{Y}_t) \leq 0$. Putting everything together, we obtain on $\mathcal{A} \cap \mathcal{B} \cap \bigcap_{M=1}^{\infty} \mathcal{E}_M$ that for any $M \geq 1$,

$$\begin{aligned} \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) &\leq \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(\hat{Y}_t, \tilde{Y}_t) \\ &+ \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, \tilde{Y}_t) - \ell(\phi_M \circ f(X_t), \tilde{Y}_t) \\ &+ \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\phi_M \circ f(X_t), \tilde{Y}_t) - \ell(f(X_t), Y_t) \\ &\leq 2 \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y^0, Y_t) \mathbb{1}_{\ell(Y_t, y^0) \geq 2^{-\alpha} M}. \end{aligned}$$

Because this holds for all $M \geq 1$, we can again apply this result to $M := \lceil \tilde{M}_\epsilon \rceil$ which yields the result $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq \epsilon$. This holds for any $\epsilon > 0$. Therefore, we finally obtain on the event $\mathcal{A} \cap \mathcal{B} \cap \bigcap_{M=1}^{\infty} \mathcal{E}_M$ of probability one, one has $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq 0$. This ends the proof that Algorithm 6 is universally consistent under CS processes for adversarial empirically integrable responses. Now because there exists a ball $B_\ell(y, r)$ of (\mathcal{Y}, ℓ) that does not satisfy F-TIME, from Theorem 5.8, universal learning with responses restricted on this ball cannot be achieved for processes $\mathbb{X} \notin \text{CS}$. However, these responses are empirically integrable because they are bounded. Hence, CS is still necessary for universal learning with adversarial empirically integrable responses. Thus SOLAR = CS and the provided learning rule is optimistically universal. This ends the proof of the theorem.

E.5 Proof of Theorem 3.2

Fix $(\mathcal{X}, \rho_{\mathcal{X}})$ and a value space (\mathcal{Y}, ℓ) such that any ball satisfies F-TIME. We now construct our learning rule. Let $\bar{y} \in \mathcal{Y}$ be an arbitrary value. For any $M \geq 1$, because $B_\ell(\bar{y}, M)$ is bounded and satisfies F-TIME, there exists an optimistically universal learning rule f^M for value space $(B_\ell(y_0, M), \ell)$. For any $M \geq 1$, we define the function $\phi_M : \mathcal{Y} \rightarrow \mathcal{Y}$ defined by restricting the space to the ball $B_\ell(\bar{y}, M)$ as follows

$$\phi_M(y) := \begin{cases} y & \text{if } \ell(y, \bar{y}) < M \\ \bar{y} & \text{otherwise.} \end{cases}$$

For simplicity, we will denote by $\hat{Y}_t^M := f_t^M(\mathbb{X}_{\leq t-1}, \phi_M(\mathbb{Y}_{\leq t-1}), X_t)$ the prediction of f^M at time t for the responses which are restricted to the ball $B_\ell(\bar{y}, M)$. We now combine these predictors using online learning into a final learning rule f . Specifically, we define $I_t := \{0 \leq M \leq 2^{-\alpha+1} \ln t\}$ for all $t \geq 1$. We also denote $t_M = \lceil e^{2^{\alpha-1} M} \rceil$ for $M \geq 0$ and pose $\eta_t = \frac{1}{4\sqrt{t}}$. For any $M \in I_t$, we define

$$L_{t-1, M} := \sum_{s=t_M}^{t-1} \ell(\hat{Y}_s^M, \phi_{2^{-\alpha+1} \ln s}(Y_s)).$$

For simplicity, we will denote by $\tilde{\mathbb{Y}}$ the process defined by $\tilde{Y}_t = \phi_{2^{-\alpha+1} \ln t}(Y_t)$ for all $t \geq 1$. We now construct recursive weights as $w_{0,0} = 1$ and for $t \geq 2$ we pose for all $1 \leq s \leq t-1$

$$\hat{l}_s := \frac{\sum_{M \in I_s} w_{s-1, M} \ell(\hat{Y}_s^M, \tilde{Y}_s)}{\sum_{M \in I_s} w_{s-1, M}}.$$

Now for any $M \in I_t$ we note $\hat{L}_{t-1, M} := \sum_{s=t_M}^{t-1} \hat{l}_s$, and pose $w_{t-1, M} := e^{\eta_t (\hat{L}_{t-1, M} - L_{t-1, M})}$. We then choose a random index \hat{M}_t independent from the past history such that

$$\mathbb{P}(\hat{M}_t = M) := \frac{w_{t-1, M}}{\sum_{M' \in I_t} w_{t-1, M'}}, \quad M \in I_t.$$

Input: Historical samples $(X_t, Y_t)_{t < T}$ and new input point X_T
 Optimistically universal learning rule f_t^M for value space $B_\ell(y_0, M), \ell$, where $y_0 \in \mathcal{Y}$ fixed.
Output: Predictions \hat{Y}_t for $t \leq T$
 $I_t := \{0 \leq M \leq 2^{-\alpha+1} \ln t\}, \eta_t := \frac{1}{4\sqrt{t}}, t \geq 1$
 $t_M = \lceil e^{2^{\alpha-1} M} \rceil, M \geq 0$
 $w_{0,0} := 1, \hat{Y}_1 = y^0 (= f^0(X_0))$ // Initialisation
for $t = 2, \dots, T$ **do**
 $L_{t-1, M} = \sum_{s=t_M}^{t-1} \ell(f_s^M(\mathbb{X}_{\leq s-1}, \phi_M(\mathbb{Y})_{\leq s-1}, X_s), \phi_{2^{-\alpha+1} \ln s}(Y_s)), \hat{L}_{t-1, M} = \sum_{s=t_M}^{t-1} \hat{\ell}_s, M \in I_t$
 $w_{t-1, M} := \exp(\eta_t(\hat{L}_{t-1, M} - L_{t-1, M})), M \in I_t$
 $p_t(M) = \frac{w_{t-1, M}}{\sum_{M' \in I_t} w_{t-1, M'}}, M \in I_t$
 $\hat{M}_t \sim p_t(\cdot)$ // Model selection
 $\hat{Y}_t = f_t^{\hat{M}_t}(\mathbb{X}_{\leq t-1}, \phi_M(\mathbb{Y})_{\leq t-1}, X_t)$
 $\hat{\ell}_t := \frac{\sum_{j \in I_t} w_{t-1, j} \ell(f_t^M(\mathbb{X}_{\leq t-1}, \phi_M(\mathbb{Y})_{\leq t-1}, X_t), \phi_{2^{-\alpha+1} \ln t}(Y_t))}{\sum_{j \in I_t} w_{t-1, j}}$
end

Algorithm 7: A learning rule for adversarial empirically integrable responses under SMV processes for value spaces (\mathcal{Y}, ℓ) such that any ball satisfies F-TIME.

The output the learning rule is $f_t(\mathbb{X}_{\leq t-1}, \mathbb{Y}_{\leq t-1}, X_t) := \hat{Y}_t^{\hat{M}_t}$. For simplicity, we will denote by $\hat{Y}_t := f_t(\mathbb{X}_{\leq t-1}, \mathbb{Y}_{\leq t-1}, X_t)$ the prediction of f at time t . This ends the construction of our learning rule which is summarized in Algorithm 7.

Now let (\mathbb{X}, \mathbb{Y}) be such that $\mathbb{X} \in \text{SOUL}$ and \mathbb{Y} empirically integrable. By Lemma 7.3, there exists some value $y_0 \in \mathcal{Y}$ such that on an event \mathcal{A} of probability one, we have for any ϵ , a threshold $M_\epsilon \geq 0$ with $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M_\epsilon} \leq \epsilon$. We fix a measurable function $f : \mathcal{X} \rightarrow \mathcal{Y}$. Also, for any $t \geq 1$ and $M \in I_t$ we have $0 \leq \ell(\hat{Y}_t^M, \tilde{Y}_t) \leq 2^{\alpha-1} \ell(\hat{Y}_t^M, \bar{y}) + 2^{\alpha-1} \ell(\tilde{Y}_t, \bar{y}) \leq 2 \ln t$. As a result, for any $M, M' \in I_t$ we have $|\ell(\hat{Y}_t^M, \tilde{Y}_t) - \ell(\hat{Y}_t^{M'}, \tilde{Y}_t)| \leq 2 \ln t$. Because $|I_t| \leq 1 + \ln t$ for all $t \geq 1$, the same proof as Theorem 3.6 shows that on an event \mathcal{B} of probability one, there exists $\hat{t} \geq 0$ such that

$$\forall t \geq \hat{t}, \forall M \in I_t, \quad \sum_{s=t_M}^t \ell(\hat{Y}_s, \tilde{Y}_s) \leq \sum_{s=t_M}^t \ell(\hat{Y}_s^M, \tilde{Y}_s) + 3 \ln^2 t \sqrt{t}.$$

Further, we know that f_t^M is Bayes optimistically universal for value space $(B_\ell(\bar{y}, M), \ell)$. In particular, because $\mathbb{X} \in \text{SOUL}$ and $\phi_M \circ f : \mathcal{X} \rightarrow B_\ell(\bar{y}, M)$, we have

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t^M, \phi_M(Y_t)) - \ell(\phi_M \circ f(X_t), \phi_M(Y_t)) \leq 0 \quad (a.s.).$$

For simplicity, we introduce $\delta_T^M := \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t^M, \phi_M(Y_t)) - \ell(\phi_M \circ f(X_t), \phi_M(Y_t))$ and define \mathcal{E}_M as the event of probability one where the above inequality is satisfied, i.e., $\limsup_{T \rightarrow \infty} \delta_T^M \leq 0$. Because we always have $\ell(\hat{Y}_t, \bar{y}) \leq 2^{-\alpha+1} \ln t$, we can write

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(\hat{Y}_t, \bar{y}) &= \frac{1}{T} \sum_{t=1}^T \left(\ell(\hat{Y}_t, Y_t) - \ell(\hat{Y}_t, \bar{y}) \right) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq 2^{-\alpha+1} \ln t} \\ &\leq \frac{1}{T} \sum_{t=1}^T \left(2^{\alpha-1} \ell(\hat{Y}_t, \bar{y}) + 2^{\alpha-1} \ell(Y_t, \bar{y}) \right) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq 2^{-\alpha+1} \ln t} \\ &\leq \frac{2^\alpha}{T} \sum_{t=1}^T \ell(Y_t, \bar{y}) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq 2^{-\alpha+1} \ln t}. \end{aligned}$$

The proof of Theorem 3.3 shows that on the event \mathcal{A} ,

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(Y_t, \bar{y}) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq 2^{-\alpha+1} \ln t} \leq 0,$$

which implies $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(\hat{Y}_t, \tilde{Y}_t) \leq 0$. Now let $M \geq 1$. We write

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t^M, \tilde{Y}_t) - \ell(\hat{Y}_t^M, \phi_M(Y_t)) \\ & \leq \frac{1}{T} \sum_{t=1}^{t_M-1} \ell(\hat{Y}_t^M, \tilde{Y}_t) + \frac{1}{T} \sum_{t=t_M}^T \left(\ell(\hat{Y}_t^M, Y_t) - \ell(\hat{Y}_t^M, \bar{y}) \right) \mathbb{1}_{M \leq \ell(Y_t, \bar{y}) < 2^{-\alpha+1} \ln t} \\ & \leq \frac{e^{2^{\alpha-1}M} 2^\alpha M}{T} + \frac{1}{T} \sum_{t=1}^T \left(2^{\alpha-1} \ell(\hat{Y}_t^M, \bar{y}) + 2^{\alpha-1} \ell(Y_t, \bar{y}) \right) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq M} \\ & \leq \frac{e^{2^{\alpha-1}M} 2^\alpha M}{T} + \frac{2^\alpha}{T} \sum_{t=1}^T \ell(Y_t, \bar{y}) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq M}. \end{aligned}$$

Hence, on the event \mathcal{A} , we obtain

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t^M, \tilde{Y}_t) - \ell(\hat{Y}_t^M, \phi_M(Y_t)) \leq 2^\alpha \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(Y_t, \bar{y}) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq M}.$$

Finally, we compute

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \ell(\phi_M \circ f(X_t), \phi_M(Y_t)) - \ell(f(X_t), Y_t) \\ & \leq \frac{1}{T} \sum_{t=1}^T (\ell(\bar{y}, Y_t) - \ell(f(X_t), Y_t)) \mathbb{1}_{\ell(f(X_t), \bar{y}) \geq M} \mathbb{1}_{\ell(Y_t, \bar{y}) \leq M} \\ & \quad + \frac{1}{T} \sum_{t=1}^T (\ell(f(X_t), \bar{y}) - \ell(f(X_t), Y_t)) \mathbb{1}_{\ell(f(X_t), \bar{y}) \leq M} \mathbb{1}_{\ell(Y_t, \bar{y}) \geq M} \\ & \leq \frac{1}{T} \sum_{t=1}^T \ell(\bar{y}, Y_t) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq 2^{-\alpha} M} + \frac{M}{T} \sum_{t=1}^T \mathbb{1}_{\ell(Y_t, \bar{y}) \geq M} \\ & \quad + \frac{1}{T} \sum_{t=1}^T (\ell(\bar{y}, Y_t) - \ell(f(X_t), Y_t)) \mathbb{1}_{\ell(f(X_t), \bar{y}) \geq M} \mathbb{1}_{\ell(Y_t, \bar{y}) \leq 2^{-\alpha} M} \\ & \leq \frac{1}{T} \sum_{t=1}^T \ell(\bar{y}, Y_t) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq 2^{-\alpha} M} + \frac{1}{T} \sum_{t=1}^T \ell(Y_t, \bar{y}) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq M} \\ & \quad + \frac{1}{T} \sum_{t=1}^T (2\ell(\bar{y}, Y_t) - 2^{-\alpha+1} \ell(f(X_t), \bar{y})) \mathbb{1}_{\ell(f(X_t), \bar{y}) \geq M} \mathbb{1}_{\ell(Y_t, \bar{y}) \leq 2^{-\alpha} M} \\ & \leq \frac{1}{T} \sum_{t=1}^T \ell(\bar{y}, Y_t) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq 2^{-\alpha} M} + \frac{1}{T} \sum_{t=1}^T \ell(Y_t, \bar{y}) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq M}. \end{aligned}$$

We now put all these estimates together. On the event $\mathcal{A} \cap \mathcal{B} \cap \bigcap_{M=1}^{\infty} \mathcal{E}_M$, for any $M \geq 1$ and $t \geq \max(\hat{t}, t_M)$

we can write

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq \frac{1}{T} \sum_{t=1}^T \left(\ell(\hat{Y}_t, Y_t) - \ell(\hat{Y}_t, \tilde{Y}_t) \right) \\
& + \frac{1}{T} \sum_{t=1}^T \left(\ell(\hat{Y}_t, \tilde{Y}_t) - \ell(\hat{Y}_t^M, \tilde{Y}_t) \right) + \frac{1}{T} \sum_{t=1}^T \left(\ell(\hat{Y}_t^M, \tilde{Y}_t) - \ell(\hat{Y}_t^M, \phi_M(Y_t)) \right) + \delta_T^M \\
& + \frac{1}{T} \sum_{t=1}^T \left(\ell(\phi_M \circ f(X_t), \phi_M(Y_t)) - \ell(f(X_t), Y_t) \right) \\
& \leq \frac{1}{T} \sum_{t=1}^T \left(\ell(\hat{Y}_t, Y_t) - \ell(\hat{Y}_t, \tilde{Y}_t) \right) + \frac{3 \ln^2 T}{\sqrt{T}} + \frac{1}{T} \sum_{t=1}^T \left(\ell(\hat{Y}_t^M, \tilde{Y}_t) - \ell(\hat{Y}_t^M, \phi_M(Y_t)) \right) \\
& \quad + \delta_T^M + \frac{1}{T} \sum_{t=1}^T \left(\ell(\phi_M \circ f(X_t), \phi_M(Y_t)) - \ell(f(X_t), Y_t) \right).
\end{aligned}$$

Thus, we obtain on the event $\mathcal{A} \cap \mathcal{B} \cap \bigcap_{M=1}^{\infty} \mathcal{E}_M$, for any $M \geq 1$,

$$\begin{aligned}
\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) & \leq \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\bar{y}, Y_t) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq 2^{-\alpha} M} \\
& \quad + (1 + 2^\alpha) \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(Y_t, \bar{y}) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq M}
\end{aligned}$$

On the event \mathcal{A} , the same arguments as in the proof of Theorem 3.3 show that we have same guarantees for y_0 as for \bar{y} , i.e., for any $\epsilon > 0$, there exists \tilde{M}_ϵ such that $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(Y_t, \bar{y}) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq \tilde{M}_\epsilon} \leq \epsilon$. Therefore, for any $\epsilon > 0$, we can apply the above equation to $M := \lceil 2^\alpha \tilde{M}_\epsilon + M_{2^{-\alpha-1}\epsilon} \rceil$ to obtain

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq \epsilon + \frac{1 + 2^\alpha}{2^{\alpha+1}} \leq 2\epsilon.$$

Because this holds for all $\epsilon > 0$, we can in finally get

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \left(\ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \right) \leq 0,$$

on the event $\mathcal{A} \cap \mathcal{E} \cap \bigcap_{M \geq 1} \mathcal{F}_M$ of probability one. This ends the proof of the theorem.