

Universal Regression with Adversarial Responses

Moïse Blanchard, Patrick Jaillet

March 2022

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 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 algorithm 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.

1 Introduction

We study the classical statistical problem of regression in general spaces. Given an instance metric space (\mathcal{X}, ρ) and a value space \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. These basic regression settings have then been extended to non-Euclidean value spaces, needed to analyze new types of data arising in numerous data analysis applications. Such examples include directional data on spherical and circular spaces [Cha89; MJM00], data lying on manifolds [Shi+09; Dav+10; Tho13], Banach spaces [Fer+11], Hilbert spaces [Bie+19] or Hadamard spaces [LM21]. This paper studies *metric-valued regression* where both instances and values lie in general metric spaces. This general setting adopted in the recent literature on universal learning [Han21b; TK22; Bla22] includes and extends the specific classification and regression settings mentioned above. In this context, we are interested in obtaining predictions with low average loss. It is well known, however, that obtaining vanishing average loss is impossible if the values are noisy. As a result, in the Bayesian version of this problem, where the samples $(\mathbb{X}, \mathbb{Y}) := (X_t, Y_t)_{t \geq 1}$ are drawn independent and identically distributed (i.i.d.) from an unknown distribution over $\mathcal{X} \times \mathcal{Y}$, a learning procedure long-term average loss is compared to the minimal loss (the term *risk* is more commonly used in the Bayesian literature) of a fixed predictor, where the minimum is taken over *all* predictor functions $f : \mathcal{X} \rightarrow \mathcal{Y}$. In the case of squared loss, and value space \mathbb{R} , the Bayes minimal risk is precisely $\text{Var}[Y|X]$. In this work, we study the considerably more general case of non-i.i.d. processes. Similarly to the Bayesian case, one aims to minimize the average *excess* loss of the predictions compared to some fixed measurable predictor function $f : \mathcal{X} \rightarrow \mathcal{Y}$, i.e., $\frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t)$. We are then interested in *consistent* learning rules which have vanishing long-run average excess loss almost surely.

There is a rich literature analyzing the general regression problem under i.i.d. sequences. In this setting, a classical result is that for the Euclidean space, k -nearest neighbor (kNN) with $k/\ln T \rightarrow \infty$ and $k/T \rightarrow 0$ is consistent under mild assumptions on the distribution of (X_1, Y_1) [Sto77; Dev+94; DGL13]. Other variants of kNN algorithms and simple learning rules have been proposed to obtain almost-sure consistency for larger classes of i.i.d. processes and spaces. For instance, in the case of binary classification and bounded regression, under mild conditions on the instance space \mathcal{X} , one can achieve consistency for any i.i.d. process, which we refer to as *universal Bayes consistency* [DGL13; Gyö+02]. More recently, this result was extended to any essentially separable metric space \mathcal{X} when the value space \mathcal{Y} is finite or countable and for 0-1 loss [Han+21; GW21], then generalized to any separable metric space (\mathcal{Y}, ℓ) [TK22] under the constraint that (X, Y) has a finite first order moment. A natural question then becomes whether such results can be obtained in the non-i.i.d. setting. Various assumptions on the sequence (\mathbb{X}, \mathbb{Y}) , which are natural relaxations of the i.i.d. condition, have been proposed. In particular, the universality result for binary classification was extended to *stationary ergodic* processes [MYG96; GLM99; Gyö+02] or processes satisfying the law of large numbers [MKN99; GG09; SHS09].

In this work, we are interested in the fundamental question of *learnability*. Namely, we aim to understand which are the minimal assumptions on the problem sequences for which consistency is still achievable. We follow the general *optimistic decision theory* introduced by [Han21a] which formalizes the general 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—this is the minimal assumption for any learner—and 2. if possible, provide learning rules that achieve this objective whenever it is achievable. These algorithms are called *optimistically universal* learning rules and enjoy the convenient property that if they failed the objective, any other algorithm would fail as well.

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 above universal consistency objective is equivalent to having vanishing long-run average error of the predictions for any measurable target function f^* . Indeed, in this specific case, the best fixed function to which we can compare the predictions of a learner is f^* , which always achieves zero loss. The two main questions of interest described above were very recently settled for noiseless responses, unveiling an important dichotomy between bounded and unbounded losses. For bounded losses, [BC21; Bla22] gave a characterization of the set SOUL (Strong Online Universal Learning) of stochastic processes \mathbb{X} for which universal learning is possible. The latter work also provides a learning rule 2C1NN, which is a simple variant of the 1-Nearest-Neighbor algorithm (1NN) and is optimistically universal for this noiseless setting with bounded losses. In particular, the set of learnable processes SOUL for noiseless online learning and bounded loss, is significantly larger than stationary processes or related assumptions. On the other hand, the case of unbounded losses is considerably more restrictive. Indeed, learnable processes necessarily visit a *finite* number of distinct instance points almost surely (Condition FS below) and simple memorization is optimistically universal [BCH22]. On the other hand, the more general non-realizable setting was not yet characterized. In this framework and for bounded losses, the very recent preprint [Han22] proposes an algorithm for metric losses which achieves consistency for *arbitrary* responses \mathbb{Y} under for a large family of instance processes \mathbb{Y} (condition CS below which intuitively asks that the sub-measure induced by empirical visits of the input sequence be continuous). Importantly, the response sequence \mathbb{Y} is completely unrestricted and can be arbitrarily correlated with the instance sequence \mathbb{X} . There is however a significant gap between the proposed CS condition and the learnable processes SOUL in the bounded noiseless setting. [Han22] then left as an open problem the question of identifying the precise provably-minimal conditions to achieve consistency, and whether there exists an optimistically universal learning rule.

In this paper, we solve this question and extend it to *adversarial* responses, which slightly generalize arbitrary responses. We will show that we can obtain the same results for adversarial processes as we would obtain if considering arbitrary responses, without any generalisability cost. Intuitively, adversarial responses can not only arbitrarily depend on the instance sequence \mathbb{X} , but may also depend on past predictions and (private) randomness used by the learner. Although adversarial processes coincide with arbitrary responses

if the learner is *deterministic*, this is a non-trivial generalization for randomized algorithms—note that randomization is necessary to obtain guarantees for general online learning frameworks [BC12; Sli19]. We now precise the distinction between arbitrary and adversarial responses. In the context of online learning, arbitrary responses correspond to an *oblivious opponent*. It is known that for some large classes of algorithms, having guarantees against any oblivious opponent yields the same guarantees against any adversarial opponents as well [CL06]. This statement holds for learners such that the immediate expected loss is only dependent on the past responses. The learning rule proposed in [Han22] falls into this category and, as a result, enjoys the same consistency guarantees for adversarial responses as shown for arbitrary responses in the original manuscript. However, for learning rules which may *explicitly* depend on the past prediction—which will be the case for some of our proposed optimistically universal learning rules—such result does not hold in general. Hence, at the level of generality considered in this work, adversarial responses seem to be a non-trivial generalization of arbitrary responses, for which we can obtain the same guarantees without any generalizability cost. We note that in this paper, we consider excess loss—or regret—compared to any fixed prediction function: in the case of adversarial responses, our regret guarantees hold against the losses obtained by any fixed prediction function along the *observed* response trajectory. As such, we do not analyze *counterfactual* excess loss—or regret—for which it is impossible to have sublinear rates against general unrestricted adversaries [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 consistency is achievable. Nevertheless, the algorithms we introduce for adversarial responses use as a 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].

Contributions. We first provide a complete characterization of the provably minimal assumptions for regression with adversarial responses, i.e., of the set of learnable processes SOLAR (Strong universal Online Learning with Adversarial Responses) in the *bounded loss* setting, which is the main interest of universal learning as shown in [BCH22]. An interesting discovery is that learnability for the general regression problem is fundamentally dependent on the value space (\mathcal{Y}, ℓ) . We show that the minimal condition for learnability SOLAR is either the CS condition which was considered in [Han22], or the significantly larger set of processes learnable in the noiseless setting SOUL for which a characterization is known (condition SMV below which intuitively asks that the process visit a sublinear number of sets from any measurable partition) [Bla22]. We further precisely identify this alternative by providing a property on the value space (\mathcal{Y}, ℓ) (Property F-TIME below) such that if satisfied, SOLAR = SMV and otherwise SOLAR = CS. This property intuitively asks that the mean-estimation problem on (\mathcal{Y}, ℓ) be achievable with a fixed error rate within a random time with fixed horizon. We show that this property is satisfied for all “reasonable” value spaces, e.g., totally-bounded spaces or countably-many-classes classification (\mathbb{N}, ℓ_{01}) . On the other hand, there exist “pathological” value spaces for which adversarial regression is inherently harder than noiseless regression.

For both cases, we provide optimistically universal learning rules. It is worth noting that the learning rules designed for each alternative are crucially different in their techniques and in nature. In the alternative when F-TIME is satisfied, the rule is *implicit* in general: it uses the existing algorithm for finite-time-mean-estimation as subroutine. However, given such algorithm, the rule is *explicit*. We show that this is the case for totally bounded value spaces and countable classification by providing explicit algorithms for finite-time-mean-estimation. On the other hand, the learning rule tailored for learning the more restrictive CS processes is always explicit and relies on specific properties of CS processes which are not satisfied by any other processes. This learning rule uses similar techniques to that introduced by [Han22] but generalizes this result to non-metric losses satisfying specific relaxed triangle-inequality properties. This allows encompassing any powers of a metric $|\cdot|^\alpha$ for $\alpha \geq 1$ and in particular the popular squared loss regression. In both alternatives,

an implication of these results is that for “reasonable” losses, learning in the regression framework can be achieved even in face of adversarially-chosen responses, and for a significantly larger class—CS or SMV—of instance processes family than *stationary* and *non-stationary* processes previously considered in the traditional statistical learning literature (e.g., [RB06]) outside of the optimistically universal learning theory.

For unbounded losses, we present a general result for mean-estimation when the loss is a metric, that holds for adversarial responses. Precisely, this is the problem of predicting values in \mathcal{Y} to minimize the loss on a sequence of samples \mathbb{Y} , which corresponds to regression without instances $\mathcal{X} = \{0\}$. For example in the case of i.i.d. sequences \mathbb{Y} , this is equivalent to the Fréchet means estimation problem for which different notions of consistency and generalizations have been recently examined [EJ20; Sch22; Jaf22]. Note that for $\mathcal{Y} = \mathbb{R}$ with Euclidean norm and i.i.d. sequences \mathbb{Y} with finite first moment, this is exactly the problem of estimating the median of Y . We show however, that mean estimation might not be achievable for adversarial responses, even in very simple cases. For instance we show that regression on Euclidean real-valued $\mathcal{Y} = \mathbb{R}$ with loss $|\cdot|^p$ adversarial responses is impossible, for any $p > 1$. As a simple consequence, this translates into an alternative for adversarial regression. Indeed, for unbounded losses, we show that we have either SOLAR = FS (= SOUL) when mean-estimation is achievable, and SOLAR = \emptyset otherwise. In particular, metric losses fall in the first alternative. Further, there always exists an optimistically universal learning rule which is inspired by the weighted average forecasters with expert advice algorithms [CL06]. The main difficulty is that the experts—fixed values in \mathcal{Y} —lie in a general metric space that may be infinite.

Last, we address the tremendous gap between learning for bounded losses compared to unbounded losses—where learnable processes necessarily visit a finite number of instance values almost surely. This issue was raised as an open problem in the noiseless setting in [BCH22]. Precisely, learning arbitrary functions even in this realizable case is too restrictive. Hence a natural question is whether imposing additional conditions on the responses would allow recovering the results from the bounded case. We propose a novel constraint asking that the sequence \mathbb{Y} be *empirically integrable*. Intuitively, this property asks that we can bound the tails of the empirical first moment of \mathbb{Y} . Under this additional constraint, we prove that learning under CS or SMV processes depending on the value space satisfying F-TIME, can be achieved with optimistically universal learning rules, even for unbounded losses. The empirical integrability property is essentially necessary in order to get such results. Indeed, we show that for the i.i.d. setting, this is exactly asking for the distribution to have finite first moment. As a result, our work significantly generalizes the main result from [TK22]—which consider i.i.d. processes (\mathbb{X}, \mathbb{Y}) —to adversarial responses, CS or SMV instance processes, and a larger class of losses which in particular encompass powers of a metric. Further, we also show that with only a finite first-order moment condition, even in the noiseless case, learning is not achievable with such a level of generality on the instance processes \mathbb{X} .

The two tables 1 and 2 summarize the known results in the literature and our contributions on learnability characterization and proposed learning rules in universal learning. For clarity, we state here the inclusions $\text{FS} \subset \text{CS} \subset \text{SMV}$ which are not equalities whenever \mathcal{X} is infinite [Han21a]. Further, CS contains in particular i.i.d., stationary ergodic or stationary processes.

Paper outline. After presenting the learning framework and definitions in Section 2, we describe in Section 3 our main results, where we give a complete characterization of learnable processes, and provide optimistically universal learning rules for all instance and value spaces. We first turn to the case of bounded losses which is less restrictive, and explicitly construct an optimistically universal learning rule for totally-bounded value spaces in Section 4. The alternative between non-totally-bounded value spaces is characterized in the following Section 5. We provide optimistically universal learning rules for each alternatives. We then turn to unbounded losses and show that for metric losses, the learnable processes are identical as for noiseless regression in Section 6, by proving an universality result for mean estimation in general spaces. We then give an alternative for adversarial regression in the general case. However, these learnable processes in unbounded value spaces are always very restrictive. Hence, in Section 7 we propose a mild moment constraint on the

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

Table 1: Characterization of the sets of learnable instance processes SOUL and SOLAR in universal consistency (ME = Mean-Estimation).

Learning setting	Loss (and response/setting constraints)	Learning rule	Guarantees for which processes \mathbb{X} ?	Optimist. universal?	Reference
Bayesian (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	[TK22]
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-variant	FS	Yes	[This paper]
	Unbounded loss + not ME	N/A	\emptyset	N/A	[This paper]
	Unbounded + EI + local F-TiME	$(1 + \delta)$ C1NN-hedged 2	SMV	Yes	[This paper]
Unbounded + EI + not local F-TiME	Hedge-variant 3	CS	Yes	[This paper]	

Table 2: Proposed learning rules for universal consistency (ME = Mean-Estimation and EI = Empirical Integrability). Note: we refer to optimistical universality with respect to the set of processes \mathbb{X} for which each learning rule has guarantees. In this context, an algorithm is optimistically universal if it is universally consistent for all processes under which universal learning is possible in the considered setting. [Han+21] showed that OptiNet, Proto-NN and MedNet are optimistically universal in another sense: they show that their guarantees hold for all *essentially separable* metric instance spaces (\mathcal{X}, ρ) , which is exactly the optimistic set of metric spaces for which Bayesian universal learning is achievable. Our proposed learning rules also enjoy this property.

responses such that all known results on bounded losses can be recovered even in the unbounded loss case. Finally, we discuss open research directions in Section 8.

2 Formal setup and Preliminaries

Instance and value spaces. We recall that a metric space is *separable* if it contains a dense countable set. In the general *metric-valued* regression problem, we observe inputs from an *instance* separable metric space (\mathcal{X}, ρ) equipped with its Borel σ -algebra \mathcal{B} , and predict values from a *value* separable metric space $(\mathcal{Y}, |\cdot|)$ given with a loss ℓ . 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 $\ell = |\cdot|^\alpha$. All of the results in this work can be generalized to *essentially separable* metric instance spaces (\mathcal{X}, ρ) , but for the sake of exposition, we will consider separable metric spaces (\mathcal{X}, ρ) in the rest of this paper. This notion of essentially separable metric space was introduced by [Han+21] and asks that for every probability measure μ on the σ -algebra \mathcal{B} , the metric probability space (\mathcal{X}, ρ, μ) is separable, i.e., there exists $\mathcal{X}' \in \mathcal{B}$ with $\mu(\mathcal{X}') = 1$ such that (\mathcal{X}', ρ) is separable. [Han+21]

showed that this is the largest class of instance metric spaces for which learning is possible, even in the Bayesian i.i.d. setting. In the first Sections 4 and 5 of this work, we suppose that the loss ℓ is *bounded*, i.e., $\sup_{y_1, y_2 \in \mathcal{Y}} \ell(y_1, y_2) < \infty$. In the rest of this paper, we will use the notation $\bar{\ell} := \sup_{y_1, y_2 \in \mathcal{Y}} \ell(y_1, y_2)$. The case of *unbounded* losses is addressed in the next sections 6 and 7. As an example, binary classification corresponds to $\mathcal{Y} = \{0, 1\}$ together with the indicator loss $\ell_{01}(i, j) = \mathbb{1}_{i \neq j}$. Similarly, this setup covers finite multi-classification with $\mathcal{Y} = \{0, 1, \dots, k\}$ or with countable number of classes $\mathcal{Y} = \mathbb{N}$, and classical regression with $\mathcal{Y} = \mathbb{R}$ and any L^α loss $\ell(x, y) = |x - y|^\alpha$ for $\alpha \geq 1$. 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.

Online learning. 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 independent from the data generation process used to generate Y_t . Formally, an online learning rule $f := \{f_t\}_{t \geq 1}$ is defined as a sequence of measurable functions $f_t : \mathcal{R}_t \times \mathcal{X}^{t-1} \times \mathcal{Y}^{t-1} \times \mathcal{X} \rightarrow \mathcal{Y}$ together with a distribution R_t on \mathcal{R}_t , where \mathcal{R}_t denotes the space on which the randomness used by f_t lies. The prediction at time t 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 a sample according to the distribution of R_t , independent of the past history and the new value $(X_u, Y_u)_{u \leq t}$. For simplicity, we might omit the randomness r_t when possible and write directly $f_t : \mathcal{X}^{t-1} \times \mathcal{Y}^{t-1} \times \mathcal{X} \rightarrow \mathcal{Y}$.

Adversarial responses. 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 setting of arbitrarily dependent responses under instance processes \mathbb{X} , introduced in [Han22]. In this work, we consider *adversarial responses* \mathbb{Y} , which generalize arbitrarily dependent responses. Specifically, the difference is that 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 . Formally, the data generation is given by a stochastic process $\{(X_t, \mathbf{Y}_t)\}_{t \geq 1}$ where $\mathbf{Y}_t = \mathbf{Y}_t(\cdot | \cdot)$ is generated from a Markov kernel from $\mathcal{R}_1 \times \dots \times \mathcal{R}_{t-1}$ to \mathcal{Y} , using the realizations of the sampled randomness of the learning rule r_1, \dots, r_{t-1} . This data generation process for the values can be viewed as a randomized measurable function $\mathbf{Y}_t : \mathcal{R}_1 \times \dots \times \mathcal{R}_{t-1} \rightarrow \mathcal{Y}$ correlated with the instance process \mathbb{X} . Having observed the sampled randomness $r_1 \in \mathcal{R}_1, \dots, r_{t-1} \in \mathcal{R}_{t-1}$ used by the learning rule f , the target value at time t is given by $Y_t := \mathbf{Y}_t(r_1, \dots, r_{t-1})$. For simplicity, we will refer to this adversarial response process as \mathbb{Y} , which allows to view the data generating process as a usual stochastic process on $\mathcal{X} \times \mathcal{Y}$ where the responses can depend on the randomness of the learning rule. If the learning rule is *deterministic*, adversarial responses are equivalent to arbitrary dependent responses considered in [Han22], but this is not the case for general *randomized* algorithms. Note that only the responses can be adapted to the randomness used by the learning rule. The instances \mathbb{X} , however, are independent of this randomness.

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. Precisely, given a learning rule f and an adversarial process (\mathbb{X}, \mathbb{Y}) , for any measurable function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$, we define the long-run average excess loss as

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

In particular, we say that the learning rule f is strongly consistent under (\mathbb{X}, \mathbb{Y}) if for any measurable function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$ we have

$$\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f^*) \leq 0, \quad (a.s.).$$

In this paper, we will only study *strong* consistency, which we will then simply refer to as consistency. 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 Bayes-optimal risk, which is defined as

$$R^* := \inf_f R(f) = \inf_f \mathbb{E}_{(X,Y) \sim \mu} [\ell(f(X), Y)],$$

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 Bayes consistency to non-i.i.d. settings. The goal of universal learning under adversarial processes is to design learning rules which are consistent for any adversarial process \mathbb{Y} . Precisely, for any stochastic process \mathbb{X} on \mathcal{X} , we say that f is *universally consistent* for adversarial responses under \mathbb{X} if it is consistent for any adversarial process $(\tilde{\mathbb{X}}, \mathbb{Y})$ where $\tilde{\mathbb{X}} \sim \mathbb{X}$, i.e., follows the same distribution as \mathbb{X} . In the case of unbounded losses, we might need to impose additional moment restrictions on \mathbb{Y} . We refer to Sections 3 and 7 for details on the corresponding restrictions.

Optimistically 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* by some learning rule. Note that this learning rule is allowed to depend on the process \mathbb{X} . We are then interested in learning rules that would achieve universal consistency with adversarial responses under all processes $\mathbb{X} \in \text{SOLAR}$, i.e., that would achieve the regression objective whenever it is achievable. We refer to these algorithms as *optimistically universal* learning rules for adversarial regression. In particular, if such a learning rule fails to reach universal consistency, then any other algorithm would fail as well.

In this general setting under minimal assumptions, 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.

2.1 Preliminaries

We recall the following known identities, which we will use to analyze the loss $\ell = |\cdot|^\alpha$.

Lemma 2.1. *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 $c_\epsilon^\alpha = \left(1 + \frac{1}{(1+\epsilon)^{1/\alpha-1}}\right)^\alpha \leq \left(\frac{2}{\epsilon \cdot \frac{1}{2} 2^{1/\alpha-1}}\right)^\alpha \leq \left(\frac{4\alpha}{\epsilon}\right)^\alpha$. \square

In particular, we will use this identity to write for any $y_1, y_2, y_3 \in \mathcal{Y}$,

$$\ell(y_1, y_2) \leq 2^{\alpha-1}\ell(y_1, y_3) + 2^{\alpha-1}\ell(y_2, y_3) \quad \text{and} \quad \ell(y_1, y_2) \leq (1 + \epsilon)\ell(y_1, y_3) + c_\epsilon^\alpha \ell(y_2, y_3).$$

These will be the only used identities on the loss ℓ . Hence, except for Section 6.1 in which we assume that the loss is a metric $\alpha = 1$, our results can be generalized to any symmetric and discernible loss ℓ satisfying the following property: 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 . Note that this is a stronger assumption than having a near-metric ℓ , for which we also give some results in Section 4 and 7.

3 Main results

We introduce a first condition on stochastic processes on \mathcal{X} . For any process \mathbb{X} on \mathcal{X} , given any measurable set $A \in \mathcal{B}$ of \mathcal{X} , we define $\hat{\mu}_\mathbb{X}(A) = \limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \mathbb{1}_A(X_t)$. 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 $A_k \downarrow \emptyset$,

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

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. 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. Formally, the condition is defined as follows.

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}, ρ) . This dependence is omitted for the sake of 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 SMV processes are exactly those for which universal consistency is achievable. Precisely, defining SOUL (Strong Online Universal Learning) the optimistic set of processes \mathbb{X} for which universal consistency in the noiseless setting is achievable, we have $\text{SMV} = \text{SOUL}$ whenever the value space is bounded. [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. Indeed, for any process $\mathbb{X} \in \text{SMV}$ and any target function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$, we have

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(2\text{C1NN}_t(\mathbb{X}_{\leq t-1}, f^*(\mathbb{X}_{t-1}), X_t), f^*(X_t)) = 0, \quad (a.s.),$$

or equivalently $\mathcal{L}_{(\mathbb{X}, f^*(\mathbb{X}))}(2\text{C1NN}, f^*) = 0$ (a.s.). In other words, 2C1NN is universally consistent under all noiseless processes with $\mathbb{X} \in \text{SMV}$. Note that the above equation coincides with the notion of universal consistency introduced in Section 2. Indeed, using f^* as comparison to the learning rule 2C1NN is an optimal choice because the loss obtained with this fixed measurable function is null. Further, because $\text{SOUL} = \text{SMV}$, if 2C1NN fails to achieve universal consistency in the noiseless setting, then any other learning rule would fail as well. Because we consider more general (noisy) responses, this shows that in particular, universal consistency with adversarial responses cannot be achieved outside of SOUL, i.e., $\text{SOLAR} \subset \text{SOUL}$ in general. In particular, for bounded value spaces we obtain $\text{SOLAR} \subset \text{SMV}$. It was posed as open problem whether we could recover the complete set SMV for learning under adversarial—or arbitrary—processes [Han22].

Open problem [Han22]. For bounded losses, does there exist an online learning rule that is universally consistent for arbitrary responses under all processes $\mathbb{X} \in \text{SOUL}$?

We answer this question with an alternative. We show that depending on the bounded value space (\mathcal{Y}, ℓ) , we have either $\text{SOLAR} = \text{SOUL}$ or $\text{SOLAR} = \text{CS}$, but that 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 $g_{\leq T_\eta}$ and τ a random time with $1 \leq \tau \leq T_\eta$ 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

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

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, we can universally learn all the processes in SOUL even under adversarial responses. We also construct an optimistically universal learning rule for this case.

Theorem 3.1. *Suppose that ℓ is bounded and (\mathcal{Y}, ℓ) satisfies F-TiME. Then, $\text{SOLAR} = \text{SMV}(= \text{SOUL})$ and there exists an optimistically universal learning rule for adversarial regression, i.e., which achieves universal consistency with adversarial responses under any process $\mathbb{X} \in \text{SMV}$.*

In other words, value spaces satisfying F-TiME allow recovering the same learnable processes for adversarial responses as for noiseless responses. We show this includes a very large class of metric spaces. Specifically, we prove that any totally-bounded metric space satisfies F-TiME. We also show that $(\mathcal{Y}, \ell) = (\mathbb{N}, \ell_{01})$, which is not totally-bounded, still satisfies F-TiME. Hence, we can universally learn all SOUL processes with adversarial responses, for countable classification. However, this property defines a non-trivial alternative and we also explicitly construct a value space that disproves F-TiME. We now turn to this second case.

Theorem 3.2. *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 $\mathbb{X} \in \text{CS}$.*

In the case of metric losses $\alpha = 1$, it is already known [Han22] that universal learning under adversarial responses under all processes in CS is achievable by some learning rule. Hence, the above result shows that this learning rule is automatically optimistically universal for adversarial regression for all metric value spaces which do not satisfy F-TiME. The main result in [Han22] considered the case of regression under arbitrary responses, but the proof can easily be adapted to adversarial responses. We will give a stronger result that holds for any $\alpha \geq 1$ and unbounded value spaces, and hence, implies this statement. This completely closes the open problem of [Han22]: the answer is positive if and only if the value spaces satisfies F-TiME.

We then turn to the case of unbounded losses. 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 of processes \mathbb{X} which visit a finite number of different points of \mathcal{X} almost surely [BCH22]. This condition is referred to as the FS (Finite Support) condition.

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

Hence, because $\text{SOUL} = \text{FS}$ for unbounded losses, even the simple memorization learning rule is optimistically universal in the noiseless setting. We show that in the adversarial setting we still have $\text{SOLAR} = \text{SOUL} = \text{FS}$ when ℓ is a metric. We prove that 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. In the case of i.i.d. processes this is precisely the Fréchet means estimation problem. We now state our main result on mean estimation, which is of independent interest and holds for general separable metric value spaces and adversarial processes.

Theorem 3.3. *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 $\text{SOLAR} = \emptyset$. Specifically, for real-valued regression with Euclidean norm and loss $|\cdot|^\alpha$ and $\alpha > 1$, neither adversarial regression nor mean-estimation are achievable. We then show that we have an alternative. Either mean-estimation with adversarial responses is achievable and in this case $\text{SOLAR} = \text{FS}$ and we have an optimistically universal learning rule, or mean-estimation is not achievable and we obtain $\text{SOLAR} = \emptyset$.

Even in the best case scenario for unbounded losses, we obtain SOLAR = SOUL = FS, which is already extremely restrictive. Thus, [BCH22] asked whether imposing moment conditions on the responses—such as empirically bounded losses in the long-run average—would allow recovering the large set SMV as learnable processes instead of the restricted set SOUL = FS. Specifically, they formulated the following open problem.

Open Problem [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 in CS—is not achievable. We therefore introduce a novel condition on the responses, namely *empirical integrability*, under which, we can recover all positive results from the bounded losses case. Precisely, we ask that there exists $y_0 \in \mathcal{Y}$ such that for any $\epsilon > 0$, almost surely there exists $M \geq 0$ for which

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

Similarly as before with the definition of SOLAR, our objective is now to study the processes \mathbb{X} for which there exists an online learning rule which would be consistent for all adversarial processes $(\tilde{\mathbb{X}}, \mathbb{Y})$ with $\tilde{\mathbb{X}} \sim \mathbb{X}$ and \mathbb{Y} satisfying the above condition. We refer to this objective as universal consistency for adversarial responses with bounded moments. We then show that all results from the bounded loss case can be recovered with this additional bounded moments constraint. Interestingly, in the noiseless case, the same 2C1NN learning rule as introduced in [Bla22] for the bounded loss case, achieves this objective.

Theorem 3.4. *Let (\mathcal{Y}, ℓ) 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 there exists $y_0 \in \mathcal{Y}$ for which for all $\epsilon > 0$, there exists $M \geq 0$ with $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, f^*(X_t)) \mathbb{1}_{\ell(y_0, f^*(X_t)) \geq M} \leq \epsilon$ (a.s.), we have $\mathcal{L}_{\mathbb{X}}(2\text{C1NN}, f^*) = 0$ (a.s.).*

For adversarial responses in unbounded loss spaces, we can also recover similar results to the bounded loss case, but the learning rules have to be adapted. We obtain the following result for adversarial universal learning under CS processes.

Theorem 3.5. *There exists an online learning rule f that is universally consistent for adversarial empirically integrable responses under all processes in CS, 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.).*

As a result, this provides an optimistically universal learning rule for adversarial responses under moment condition, for all value spaces (\mathcal{Y}, ℓ) such that there exists a ball $B_\ell(y, r)$ with $r > 0$ which does not satisfy F-TIME. Otherwise, under the same moment condition, adversarial universal learning is achievable under all processes in SMV(= SOUL).

Theorem 3.6. *Suppose that any ball of (\mathcal{Y}, ℓ) , $B_\ell(y, r)$ satisfies F-TIME. Then, there exists an optimistically universal online learning rule f for adversarial empirically integrable responses with bounded moments, 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, 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.).*

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. 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 $\ell(y, y') < \epsilon$. Note in particular that 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 SOUL = SMV which correspond to learnable processes in the noiseless setting. Further, we explicitly construct an optimistically universal learning rule for adversarial responses. The main result of this section is stated below.

Theorem 4.1. *Suppose that (\mathcal{Y}, ℓ) is totally-bounded. Then, 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., such that 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.).*

We recall that in the noiseless setting, the 2C1NN learning rule achieves universal consistency for all SMV processes [Bla22]. Precisely, the 2C1NN learning rule performs, at each iteration t , the nearest neighbor rule over an updated 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 representative. This number is then capped at 2 by deleting of 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, [TK22] 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 need to generalize a result from [Bla22]. Instead of the 2C1NN learning rule, we will use $(1 + \delta)$ C1NN rules for $\delta > 0$ arbitrarily small. In the 2C1NN learning rule, to each new input point X_t is associated a representative $\phi(t)$ which is used for the prediction $\hat{Y}_t = Y_{\phi(t)}$. This rule was designed so that each point can be used at most twice as representative for future times. In the $(1 + \delta)$ C1NN rule, each point will be used as representative at most twice with probability δ and at most once with probability $1 - \delta$. In order to have this behaviour 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, as performed by the 2C1NN, and uses it as representative if allowed by the maximum number of children that this nearest neighbor 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 children.

- 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] up to minor memorization improvements.

Theorem 4.2. *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.*

Input: Historical samples $(X_t, Y_t)_{t < T}$ and new input point X_T
Output: Predictions $\hat{Y}_t = (1 + \delta)\text{C1NN}_t(\mathbf{X}_{<t}, \mathbf{Y}_{<t}, X_t)$ for $t \leq T$
 $\hat{Y}_1 := \bar{y}$ // Arbitrary prediction at $t = 1$
 $\mathcal{D}_2 \leftarrow \{1\}; n_1 \leftarrow 0; t \leftarrow 2;$ // Initialisation
while $t \leq T$ **do**
 $continue \leftarrow True$ // Begin search for available representative $\phi(t)$
 while $continue$ **do**
 $\phi(t) \leftarrow \min \{l \in \arg \min_{u \in \mathcal{D}_t} \rho(X_t, X_u)\}$
 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
 $\hat{Y}_t := Y_{\phi(t)}$
 $n_{\phi(t)} \leftarrow n_{\phi(t)} + 1$
 $n_t \leftarrow 0$
 $t \leftarrow t + 1$
end

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

The proof of this theorem is given in the following Section 4.1. 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 an ϵ -net of \mathcal{Y} for $\bar{\ell}$. Importantly, we can take \mathcal{Y}_ϵ finite because \mathcal{Y} is totally-bounded. Recall that we denote by $\bar{\ell}$ the supremum loss, i.e., $\bar{\ell} := \sup_{y_1, y_2 \in \mathcal{Y}} \ell(y_1, y_2)$. We pose

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

The quantity T_ϵ will be the horizon window used by our learning rule to make its prediction using the $(1 + \delta_\epsilon)\text{C1NN}$ learning rule. Precisely, let ϕ be the representative function from the $(1 + \delta_\epsilon)\text{C1NN}$ learning rule. Further, we denote by $d(t)$ the depth of time t within the graph constructed by $(1 + \delta_\epsilon)\text{C1NN}$. At time t , we define the horizon $L_t = d(t) \bmod T_\epsilon$. The learning rule performs its prediction based on the values $Y_{\phi^l(t)}$ for $l = 1, \dots, L_t$. We pose $\eta_\epsilon := \sqrt{\frac{8 \ln |\mathcal{Y}_\epsilon|}{\bar{\ell}^2 T_\epsilon}}$ and define the losses $L_y^t = \sum_{l=1}^{L_t} \ell(Y_{\phi^l(t)}, y)$. 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 . This ends the construction of the learning rules f^ϵ . We are now ready to define our final learning rule f_\cdot . Let $\epsilon_i = 2^{-i}$ for all $i \geq 0$. We define $I_t := \{i \leq \ln t\}$ for any $t \geq 1$. We also denote $t_i := \lceil e^i \rceil$ and pose $\eta_t = \sqrt{\frac{\ln t}{t}}$. For any $i \in I_t$ we define $L_{t-1, i} := \sum_{s=t_i}^{t-1} \ell(\hat{Y}_s(\epsilon_i), Y_s)$ and construct some weights $w_{t-1, i}$ 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 the weights $w_{s-1,i}$ have been constructed for all $i \in I_s$ and $1 \leq s \leq t-1$. We define

$$\hat{\ell}_s := \frac{\sum_{i \in I_s} w_{s-1,i} \ell(\hat{Y}_s(\epsilon_i), Y_s)}{\sum_{i \in I_s} w_{s-1,i}}.$$

Now 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 $w_{t-1,i} = e^{\eta_t(\hat{L}_{t-1,i} - L_{t-1,i})}$. We now define a random index \hat{i}_t , independent from the past history such that

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

The output of our learning rule is $\hat{Y}_t := \hat{Y}_t(\epsilon_{\hat{i}_t})$. Before analyzing this algorithm, we introduce the following helper lemma which can be found in [CL06].

Lemma 4.3 ([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 now compare the predictions of this learning rule f . to the predictions of the rules f^ϵ .

Lemma 4.4. *Almost surely, there exists $\hat{t} \geq 0$ such that*

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

Proof. 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

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

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 4.3. 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

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

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

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

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 $\ell(\hat{Y}_t, Y_t) - \hat{\ell}_t$ that form a sequence of martingale differences to obtain $\mathbb{P} \left[\sum_{s=t_i}^t \ell(\hat{Y}_s, Y_s) > \hat{L}_{t,i} + u \right] \leq e^{-\frac{2u^2}{i\bar{\ell}^2}}$. 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. \square

We are now ready to prove the main result that for totally-bounded value spaces, f_\cdot is a universal consistent for adversarial regression under all processes $\mathbb{X} \in \text{SMV}$.

Proof of Theorem 4.1. Let $0 < \epsilon \leq 1$. We first analyze the prediction of the learning rule f^ϵ . The learning rule was constructed so that we perform exactly the classical exponentially weighted average forecaster for the predictions $\hat{Y}_{\phi^{L_t(t)}}, \hat{Y}_{\phi^{L_t-1(t)}}, \dots, \hat{Y}_{\phi(t)}, \hat{Y}_t$. As a result, denoting $\bar{\ell}(\hat{Y}_t(\epsilon), Y_t) := \sum_{y \in \mathcal{Y}_\epsilon} \mathbb{P}(\hat{Y}_t(\epsilon) = y) \ell(y, Y_t)$, by online learning, we have that for any $t \geq 1$,

$$\frac{1}{\bar{\ell}} \sum_{u=0}^{L_t} \bar{\ell}(\hat{Y}_{\phi^u(t)}(\epsilon), Y_{\phi^u(t)}) \leq \frac{1}{\bar{\ell}} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u=0}^{L_t} \ell(y, Y_{\phi^u(t)}) + \sqrt{\frac{L_t \ln |\mathcal{Y}_\epsilon|}{2}}.$$

Now consider a horizon $T \geq 1$, 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 $i = 0, 1, 2$. Note that no time can have more than 2 children. Define $\mathcal{B}_T = \{t \leq T : L_t = T_\epsilon - 1 \text{ and } \forall u = 1, \dots, T_\epsilon - 1, \phi^u(t) \notin \mathcal{A}_1\}$, i.e., times congruent to $T_\epsilon - 1$ modulo T_ϵ and such that parents until the $(T_\epsilon - 1)$ -th generation have only childs. Note that by construction, for all $t \in \mathcal{B}_T$, the sets $\{\phi^u(t), u = 0, \dots, T_\epsilon - 1\}$ are disjoint which yields $|\mathcal{B}_T| T_\epsilon \leq T$. Last, we denote $\mathcal{E} = \bigcup_{t \in \mathcal{B}_T} \{\phi^u(t), u = 0, \dots, T_\epsilon - 1\}$. This set contains all times $t \leq T$ except for times close to leaves \mathcal{A}_0 or to times in \mathcal{A}_2 which had at least two children. Specifically, for time $t \in \mathcal{A}_2$, potentially, its children until generation $T_\epsilon - 1$ will not be in \mathcal{E} . Therefore, by summing the above equation for all times in \mathcal{B}_T , we obtain

$$\begin{aligned} \sum_{t=1}^T \bar{\ell}(\hat{Y}_t(\epsilon), Y_t) &\leq \sum_{t \in \mathcal{B}_T} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u=0}^{T_\epsilon-1} \ell(y, Y_{\phi^u(t)}) + |\mathcal{B}_T| \bar{\ell} \sqrt{\frac{T_\epsilon \ln |\mathcal{Y}_\epsilon|}{2}} + (T - |\mathcal{E}|) \bar{\ell} \\ &\leq \sum_{t \in \mathcal{B}_T} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u=0}^{T_\epsilon-1} \ell(y, Y_{\phi^u(t)}) + T \bar{\ell} \sqrt{\frac{\ln |\mathcal{Y}_\epsilon|}{2T_\epsilon}} + (|\mathcal{A}_2| 2^{T_\epsilon} + |\mathcal{A}_0| T_\epsilon) \bar{\ell} \\ &\leq \sum_{t \in \mathcal{B}_T} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u=0}^{T_\epsilon-1} \ell(y, Y_{\phi^u(t)}) + \epsilon T + (|\mathcal{A}_2| 2^{T_\epsilon} + |\mathcal{A}_0| T_\epsilon) \bar{\ell}. \end{aligned}$$

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

$$\begin{aligned} \sum_{t=1}^T \bar{\ell}(\hat{Y}_t(\epsilon), Y_t) &\leq \sum_{t \in \mathcal{B}_T} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u=0}^{T_\epsilon-1} \ell(y, Y_{\phi^u(t)}) + \epsilon T + (T_\epsilon + 2T\delta_\epsilon(2^{T_\epsilon} + T_\epsilon)) \bar{\ell}. \\ &\leq \sum_{t \in \mathcal{B}_T} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u=0}^{T_\epsilon-1} \ell(y, Y_{\phi^u(t)}) + T_\epsilon \bar{\ell} + 2\epsilon T. \end{aligned}$$

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

$$\sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) \leq \sum_{t \in \mathcal{B}_T} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u=0}^{T_\epsilon-1} \ell(y, Y_{\phi^u(t)}) + T_\epsilon \bar{\ell} + 2\epsilon T + 2\bar{\ell} \sqrt{T \ln T}.$$

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

$$\forall T \geq \hat{T}, \quad \sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) \leq \sum_{t \in \mathcal{B}_T} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u=0}^{T_\epsilon-1} \ell(y, Y_{\phi^u(t)}) + T_\epsilon \bar{\ell} + 2\bar{\ell} \sqrt{T \ln T} + 2\epsilon T.$$

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 4.2, the rule $(1 + \delta_\epsilon)C1NN$ is optimistically universal in the noiseless setting. Therefore, because $\mathbb{X} \in \text{SOUL}$ we have in particular

$$\frac{1}{T} \sum_{t=1}^T \ell((1 + \delta_\epsilon)C1NN_t(\mathbb{X}_{\leq t-1}, f(\mathbb{X}_{\leq t-1}), X_t), f(X_t)) \rightarrow 0 \quad (a.s.),$$

i.e., almost surely, $\frac{1}{T} \sum_{t=1}^T \ell(f(X_{\phi(t)}), f(X_t)) \rightarrow 0$. We denote by \mathcal{F}_ϵ this event of probability one. Using Lemma 2.1, we write for any $u = 1, \dots, T_\epsilon - 1$,

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

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 $u = 1, \dots, T_\epsilon - 1$

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

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 $0 < \mu \leq 1$. $\mu := \epsilon^{\frac{1}{\alpha+1}}$. Putting everything together, on the event $\mathcal{E}_\epsilon \cap \mathcal{F}_\epsilon$, for any $T \geq \hat{T}$, we have

$$\begin{aligned} \sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) &\leq \sum_{t \in \mathcal{B}_T} \min_{y \in \mathcal{Y}_\epsilon} \sum_{u=0}^{T_\epsilon-1} \ell(y, Y_{\phi^u(t)}) + T_\epsilon \bar{\ell} + 2\bar{\ell} \sqrt{T \ln T} + 2\epsilon T \\ &\leq \sum_{t \in \mathcal{B}_T} \sum_{u=0}^{T_\epsilon-1} \ell(f(X_t)^\epsilon, Y_{\phi^u(t)}) + T_\epsilon \bar{\ell} + 2\bar{\ell} \sqrt{T \ln T} + 2\epsilon T \\ &\leq \sum_{t \in \mathcal{B}_T} \sum_{u=0}^{T_\epsilon-1} [c_{\mu_\epsilon}^\alpha \ell(f(X_t)^\epsilon, f(X_t)) + (c_{\mu_\epsilon}^\alpha)^2 \ell(f(X_t), f(X_{\phi^u(t)})) + (1 + \mu_\epsilon)^2 \ell(f(X_{\phi^u(t)}), Y_{\phi^u(t)})] \\ &\quad + T_\epsilon \bar{\ell} + 2\bar{\ell} \sqrt{T \ln T} + 2\epsilon T \\ &\leq (1 + \mu_\epsilon)^2 \sum_{t=1}^T \ell(f(X_t), Y_t) + (c_{\mu_\epsilon}^\alpha)^2 \sum_{u=1}^{T_\epsilon-1} \sum_{t=1}^T \ell(f(X_t), f(X_{\phi^u(t)})) + T_\epsilon \bar{\ell} + 2\bar{\ell} \sqrt{T \ln T} + (2 + c_{\mu_\epsilon}^\alpha) \epsilon T \\ &\leq \sum_{t=1}^T \ell(f(X_t), Y_t) + (c_{\mu_\epsilon}^\alpha)^2 \sum_{u=1}^{T_\epsilon-1} \sum_{t=1}^T \ell(f(X_t), f(X_{\phi^u(t)})) + T_\epsilon \bar{\ell} + 2\bar{\ell} \sqrt{T \ln T} + (2\epsilon + \epsilon c_{\mu_\epsilon}^\alpha + 3\mu_\epsilon) T, \end{aligned}$$

where in the thirds inequality we used Lemma 2.1 twice. Hence, for any $\epsilon < (c_1^\alpha)^{-2}$, on the event $\mathcal{E}_\epsilon \cap \mathcal{F}_\epsilon$, we obtain

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) - \ell(f(X_t), Y_t) \leq 2\epsilon + \epsilon c_{\mu_\epsilon}^\alpha + 3\mu_\epsilon \leq 2\epsilon + \sqrt{\epsilon} + 3\mu_\epsilon,$$

where $\mu_\epsilon \rightarrow_{\epsilon \rightarrow 0} 0$. We now denote $\delta_\epsilon := 2\epsilon + \sqrt{\epsilon} + 3\mu_\epsilon$ and $i_0 = \lceil \frac{2 \ln c^\alpha}{\ln 2} \rceil$. 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 4.4 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. \square

As a result, we obtain in particular $\text{SMV} \subset \text{SOLAR}$ 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, we obtain as corollary that the proposed learning rule from Theorem 4.1 is optimistically universal for adversarial regression. These results are summarized in the following result.

Corollary 4.5. *Suppose that (\mathcal{Y}, ℓ) is totally-bounded. Then, $\text{SOLAR} = \text{SMV}(= \text{SOUL})$ and there exists an optimistically universal learning rule for adversarial regression, i.e., which achieves universal consistency with adversarial responses under any process $\mathbb{X} \in \text{SOUL}$.*

This is a first step towards the more general Theorem 3.1. 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} . Intuitively, because this is a finite set, we can perform classical online learning—for instance with the exponentially weighted average forecaster [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 $\tau_\eta := T_\eta$.

4.1 Proof of Theorem 4.2

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

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.

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 $\{l \in \mathcal{T}_k \mid \exists u \leq t_{k_l} : \phi(u) = l, \rho(X_l, \bar{x}) \geq r \text{ and } \forall v < u, \phi(v) \neq l\} \neq \emptyset$. We denote the corresponding minimum time $l_k = \min\{l \in \mathcal{T}_k \mid \exists u \leq t_{k_l} : \phi(u) = l, \rho(X_l, \bar{x}) \geq r, \forall v < u, \phi(v) \neq l\}$. 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 the descendance 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) $\rho(X_t, \bar{x}) < r$. If used as representative $\phi(t) = l_k$, then l_k would not have induced a mistake in the prediction of Y_t .
- (b) $\rho(X_t, \bar{x}) \geq r$. 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

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

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

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

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 $\mathbb{P}[D^c \cap \{\beta \geq \lceil k_l/8 \rceil\}] \leq e^{-k_l \delta^2/2^7}$. Finally we write for any $l > 16$,

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

In particular, we obtain

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

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 $\mathcal{T} := \{t \leq t_{k_l}, \rho(X_t, \bar{x}) \geq r\}$. Again, $\tilde{\mathcal{G}}$ is a collection of disjoint trees with roots times $\{t_k, k \leq k_l, \rho(X_{t_k}, \bar{x}) \geq r\}$ —and possibly $t = 1$. For a given time t_k with $k \leq k_l$ and $\rho(X_{t_k}, \bar{x}) \geq r$, 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 $\{l \in \mathcal{T}_k | \exists u \leq t_{k_l} : \phi(u) = l, \rho(X_l, \bar{x}) < r \text{ and } \forall v < u, \phi(v) \neq l\} \neq \emptyset$ and more specifically their minimum time $l_k = \min\{l \in \mathcal{T}_k | \exists u \leq t_{k_l} : \phi(u) = l, \rho(X_l, \bar{x}) < r, \forall v < u, \phi(v) \neq l\}$. 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 $\mathcal{M} := \bigcup_{l \geq 1} \bigcap_{l' \geq l} \tilde{\mathcal{M}}_{k_{l'}}$ 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 $\frac{\eta}{4}$, \mathbb{X} disproves the $\text{SMV}_{(\mathcal{X}, \rho)}$ 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.

Case 1. In this 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 $|G| \geq \frac{k_l \delta}{32}$. 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 $2^{2 \cdot \frac{2^7}{\epsilon \delta} - 1} \leq \frac{1}{4c_\epsilon}$ as needed in the original proof.

Case 2. 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 $|\{k \in S, \rho(x_{p_{d(k)}^k}, \bar{x}) > r\}| \geq \frac{|S|}{2}$, then we have

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

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 $|\{k \in S, \rho(x_{p_{d(k)}^k}, \bar{x}) > r\}| < \frac{|S|}{2}$. Therefore, we obtain

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

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 $T = \{k \in S, \rho(x_{p_{d(k)}^k}, \bar{x}) = r\}$ and times $\mathcal{U}_u = \{p_{d(k)}^k : n_u < p_{d(k)}^k \leq t_{k_l}, k \in T\}$. In particular, we have

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

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

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

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 $\text{SMV}_{(\mathcal{X}, \rho)}$ 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 $\mathbb{X} \notin \text{SMV}_{(\mathcal{X}, \rho)}$ 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 4.7. *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 [BC21] which gives a reduction from any near-metric bounded value space to binary classification.

Theorem 4.8 ([BC21]). *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 4.7, Theorem 4.8 ends the proof of the Theorem 4.2.

5 Characterization of learnable processes for bounded losses

While the 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 for any *bounded* value space. We also provide optimistically universal algorithms for any bounded setting.

5.1 Negative result for non-totally-bounded spaces

In Corollary 4.5, we showed that for totally-bounded value spaces, $\text{SOLAR} = \text{SMV}$, i.e., adversarial regression is achievable for all processes $\mathbb{X} \in \text{SMV}$. Unfortunately, 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 4.5 to any bounded value space is impossible. Precisely, there exists value spaces for which the set of learnable processes for adversarial regression is 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 , 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 $\mathcal{L}_{(\mathbb{X}, \mathbb{Y})}(f, f^*) \geq \epsilon$.*

Proof. We start by constructing the space (\mathcal{Y}, ℓ) which we will use for the negative result. Specifically, we choose $\mathcal{Y} = \mathbb{N} := \{i \geq 0\}$ and construct the corresponding loss ℓ . For any $k \geq 1$, we pose $n_k := 2k(k-1) + 2^k - 1$ and define the sets

$$I_k := \{n_k, n_k + 1, \dots, n_k + 4k - 1\} \quad \text{and} \quad J_k := \{n_k + 4k, n_k + 4k + 1, \dots, n_k + (4k - 1) + 2^k = n_{k+1} - 1\}.$$

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 $\ell(0, i) = \mathbb{1}_{i=0}$ for all $i \in \mathbb{N}$. For any $1 \leq k < l$, and any $i \in I_k \cup J_k, j \in I_l \cup J_l$ we define $\ell(i, j) = 1$. Further, for any $k \geq 1$, for all $i, j \in I_k$ we pose $\ell(i, j) = \mathbb{1}_{i=j}$. Similarly, for all $i, j \in J_k$ we pose $\ell(i, j) = \mathbb{1}_{i=j}$. 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

$$\ell(n_k + 4u, j) = \ell(n_k + 4u + 1, j) = \frac{1 + b_j^u}{2} \quad \text{and} \quad \ell(n_k + 4u + 2, j) = \ell(n_k + 4u + 3, j) = \frac{2 - b_j^u}{2},$$

for all $u \in \{0, 1, \dots, k-1\}$ and $j \in J_k$. This ends the definition of ℓ . We first check that this indeed defines 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} for value space (\mathcal{Y}, ℓ) . 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

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

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 $\mathcal{A} \subset \bigcup_{t > \frac{8}{\epsilon} t_{p-1}} \left\{ \frac{1}{t} \sum_{u=1}^t \mathbb{1}_{\bigcup_{l \geq l_p} B_l}(X_u) \geq \frac{3\epsilon}{8} \right\}$. 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 $p \geq 1$,

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

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

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

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 $t_{p-1} < T \leq t_p$,

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

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 $f^b(x_t) = Y_t^b$, hence $\ell(f^b(x_t), Y_t^b) = 0$. Now if $X_t \in \tilde{B}_p$ by construction we have $\ell(f^b(x_t), Y_t^b) = \frac{1}{2}$. 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 $\ell(f^b(x_t), Y_t^b) = \frac{1}{2} \mathbb{1}_{t \in \mathcal{T}_p}$ for any $t_{p-1} < t \leq t_p$. Thus, we obtain for any $t_{p-1} < T \leq t_p$,

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

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,

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

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

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

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

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

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 we do not have almost surely $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t^b) - \ell(f^b(X_t), Y_t^b) \leq 0$. 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. \square

For the constructed value space, although we have $\text{SOLAR} = \text{CS} \subsetneq \text{SOUL}$, 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 $\mathbb{X} \in \text{SMV}(= \text{SOUL})$, 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 do not directly prove that this extension holds, because we will later prove a stronger result which also holds for any $\alpha \geq 1$ —even if ℓ is a power of metric—and for unbounded losses in Section 7. This shows that for any separable metric space (\mathcal{X}, ρ) , 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 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 exists non-totally-bounded value spaces for which we can recover SOLAR = SMV(= SOUL). 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 which allows to perform online learning with an infinite number of experts in the context of countable classification.

Lemma 5.3. *Let $T \geq 1$. 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*

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

or equivalently,

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

and with probability $1 - \delta$,

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

Proof. We first construct our online learning algorithm, which is a simple variant of the classical exponential forecaster. We first define a step $\eta := \sqrt{2 \ln T / T}$. 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 $\hat{s}_t = \mathbb{E}[\mathbb{1}_{f_t(\mathbf{y}_{\leq t-1})=y_t}]$. 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

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

for some vector $\xi \in \mathbb{R}^{|S_{t-1}|}$, where in the last inequality we used the fact $|\mathbb{1}_{y=y_t} - \hat{s}_t| \leq 1$. 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

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

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

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

In particular, this shows that

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

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 $1 - \delta$,

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

This ends the proof of the lemma. \square

We are now ready to prove our main result for countable classification.

Theorem 5.4. *Let $(\mathcal{X}, \mathcal{B})$ 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.).*

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

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

Then, let ϕ be the representative function from the $(1 + \delta_\epsilon)\text{C1NN}$ learning rule. Further, we denote by $d(t)$ the depth of time t within the graph constructed by $(1 + \delta_\epsilon)\text{C1NN}$. At time t , we define the horizon $L_t = d(t)$ mod T_ϵ . The learning rule performs its prediction based on the values $Y_{\phi^l(t)}$ for $l = 1, \dots, L_t$. We pose $\eta_\epsilon := \sqrt{2 \ln T_\epsilon / T_\epsilon}$ and define the weights $w_{y,t} = e^{\eta_\epsilon \sum_{l=1}^{L_t} \mathbb{1}(Y_{\phi^l(t)} = y)}$ for all $y \in \tilde{S} := \{y' \in \mathbb{N} : \sum_{l=1}^{L_t} \mathbb{1}(Y_{\phi^l(t)} = y') > 0\}$ 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 4.4 still holds. Also, using the same notations as in the proof of Theorem 4.1, Lemma 5.3 implies that for any $t \geq 1$, we can write

$$\sum_{u=0}^{L_t} \bar{\ell}_{01}(\hat{Y}_{\phi^u(t)}(\epsilon), Y_{\phi^u(t)}) \leq \min_{y \in \mathbb{N}} \sum_{u=0}^{L_t} \ell_{01}(y, Y_{\phi^u(t)}) + 1 + \ln 2 \sqrt{\frac{L_t + 1}{2 \ln(L_t + 1)}} + \sqrt{2(L_t + 1) \ln(L_t + 1)}.$$

Therefore, when summing these equations for all times in \mathcal{B}_T , the term corresponding to the margin regret of these predictions $\hat{Y}_t(\epsilon)$ becomes

$$\begin{aligned} |\mathcal{B}_T| \left(1 + \ln 2 \sqrt{\frac{T_\epsilon}{2 \ln T_\epsilon}} + \sqrt{2 T_\epsilon \ln T_\epsilon} \right) &\leq T \left(\frac{1}{T_\epsilon} + \frac{\ln 2}{\sqrt{2 T_\epsilon \ln T_\epsilon}} + \sqrt{\frac{2 \ln T_\epsilon}{T_\epsilon}} \right) \\ &\leq T \left(\frac{\epsilon}{3} + \frac{\epsilon}{3} + \frac{\epsilon}{3} \right) \\ &= \epsilon T. \end{aligned}$$

Thus the same estimates still hold starting from the inequality

$$\sum_{t=1}^T \bar{\ell}_{01}(\hat{Y}_t(\epsilon), Y_t) \leq \sum_{t \in \mathcal{B}_T} \min_{y \in \mathbb{N}} \sum_{u=0}^{T_\epsilon-1} \ell(y, Y_{\phi^u(t)}) + \epsilon T + (|\mathcal{A}_2| 2^{T_\epsilon-1} + |\mathcal{A}_0| T_\epsilon),$$

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. As a result, the rest of the proof of Theorem 4.1 holds, which ends the proof of this theorem. \square

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 $\text{SOLAR} = \text{CS}$ or $\text{SOLAR} = \text{SMV}$. In this section, we prove that there is an underlying alternative, defined by F-TIME, which enables to precisely characterize the set SOLAR of learnable processes for adversarial regression. For the sake of exposition, we recall the property F-TIME on separable value spaces (\mathcal{Y}, ℓ) .

Property F-TIME: For any $\eta > 0$, there exists a horizon time $T_\eta \geq 1$, an online learning rule $g_{\leq T_\eta}$ and τ a random time with $1 \leq \tau \leq T_\eta$ 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

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

As an important note, the random time τ may depend on the possible randomness of the learning rule g . However, τ is not online, 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. The random stopping time τ depends on this private randomness. This randomness is then never revealed to the adversary choosing the values \mathbf{y} , only through the realizations of the predictions. We now show that if F-TIME is satisfied by the value space, similarly to the case of countable classification, we can recover SOLAR = SMV(= SOUL) and there exists an optimistically universal rule for adversarial regression.

Theorem 3.1. 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 $\mathbb{X} \in \text{SMV}$.

Proof. Let us construct a learning rule which we will then prove is universally consistent for adversarial responses under all processes $\mathbb{X} \in \text{SMV}$. We first define the learning rules f^ϵ , which will be combined in a second step. For $\epsilon > 0$, 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) . Similarly as before, we then pose $\delta_\epsilon := \frac{\epsilon}{2\ell(2T_\epsilon + 2T_\epsilon)}$ and let ϕ be the representative function from the $(1 + \delta_\epsilon)\text{C1NN}$ learning rule. We define a sequence of i.i.d. copies $g^{\epsilon, t}$ of the learning rule g^ϵ for all $t \geq 1$. Precisely, 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. For any $t \geq 1$, we then construct an integer $0 \leq L_t < T_\epsilon$ which will correspond to the window span used for the prediction \hat{Y}_t . Precisely, we will make this prediction based on the values $Y_{\phi^l(t)}$ for $l = 1, \dots, L_t$. This window span is constructed recursively. Specifically, at time t , we draw a uniform $U_t \sim \mathcal{U}([0, 1])$ independent from the past, and pose

$$L_t = \begin{cases} L_{\phi(t)} + 1 & \text{if } U_t < \frac{\mathbb{P}[\tau_\epsilon \geq L_{\phi(t)} + 2]}{\mathbb{P}[\tau_\epsilon \geq L_{\phi(t)} + 1]}, \\ 0 & \text{otherwise.} \end{cases}$$

with the convention that if $\frac{0}{0} = 0$. In practice, this convention is not necessary, since if $\mathbb{P}[\tau_\epsilon \geq u] = 0$, the learning rule never creates indices $L_t \geq u - 1$ for any time $t \geq 1$. In particular, we always have $L_t \leq T_\epsilon - 1$. We now define the learning rule g^ϵ such that for any sequence \mathbf{x}, \mathbf{y} we have

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

For simplicity, we refer to this prediction as $\hat{Y}_t(\epsilon)$. The final learning rule f is defined similarly as before from the learning rules f^ϵ . This ends the construction of our learning rule.

We now show that it achieves Bayes optimistical universality for arbitrary responses. By construction of the learning rule f , Lemma 4.4 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

$$\mathcal{B}_T = \{t \leq T : L_t = 0 \text{ and } \nexists t' \in \mathcal{A}_0 \cup \mathcal{A}_2, 0 \leq u \leq T_\epsilon - 1, \phi^u(t') = t\},$$

i.e., times that start a new block and such that their descendance until generation $T_\epsilon - 1$ are neither leaves nor nodes with 2 children. In particular, any $t \in \mathcal{B}_T$ has a descendance of only children until generation $T_\epsilon - 1$. To simplify notations, for any $t \in \mathcal{B}_T$, we denote t^u its children at generation $u - 1$ for $1 \leq u \leq T_\epsilon$, i.e., we have $t^u = \phi^{T_\epsilon - u}(t^{T_\epsilon})$ for all $1 \leq u \leq T_\epsilon$, and $t = t^1$. By construction, because blocks cannot be

longer than T_ϵ , there exists $1 \leq u \leq T_\epsilon$ such that t^u ends the block started in t^1 . We have in particular $L_{t^u} = u - 1$ and the only child t' of t^u has $L_{t'} = 0$. To simplify notations, we denote $s(t) = u$ the size of the block started in t . By construction of the indices L_t , as well as the property that for any $t \in \mathcal{B}_T$, its children until generation $T_\epsilon - 1$ have exactly one child, the blocks $\{t^u, 1 \leq u \leq s(t)\}$, for $t \in \mathcal{B}_T$, are all disjoint. This implies in particular $\sum_{t \in \mathcal{B}_T} s(t) \leq T$. We first analyze the predictions along these blocks and for any $t \in \mathcal{B}_T$ and $y \in \mathcal{Y}$, we pose $\delta_t(y) := \frac{1}{s(t)} \sum_{u=1}^{s(t)} (\ell(\hat{Y}_{t^u}, Y_{t^u}) - \ell(y, Y_{t^u}) - \epsilon)$. We will now need the following lemma.

Lemma 5.5. *For any sequence $(y^t)_{t \geq 1}$ of values in \mathcal{Y} , with probability $1 - \delta$ we have*

$$\sum_{t \in \mathcal{B}_T} s(t) \delta_t(y^t) \leq [\bar{\ell}(T_\epsilon + 1) + T_\epsilon] \sqrt{2T \ln \frac{2}{\delta}}.$$

We also denote $\mathcal{T} = \bigcup_{t \in \mathcal{B}_T} \{t^u, 1 \leq u \leq s(t)\}$ the union of all blocks within horizon T . This set contains all times $t \leq T$ except times close to leaves \mathcal{A}_0 or times in \mathcal{A}_2 , within distance $T_\epsilon - 1 - L_t \leq T_\epsilon - 1$ distance. Therefore, using the same arguments as in the proof of Theorem 4.1, by the Chernoff bound, with probability at least $1 - e^{-T\delta_\epsilon/3}$ we have

$$T - |\mathcal{T}| \leq |\mathcal{A}_2| 2^{T_\epsilon} + (|\mathcal{A}_2| + |\mathcal{A}_0|) T_\epsilon = T_\epsilon + (2^{T_\epsilon} + 2T_\epsilon) |\mathcal{A}_2| \leq T_\epsilon + (2^{T_\epsilon} + 2T_\epsilon) (2T\delta_\epsilon) = T_\epsilon + \frac{\epsilon}{\ell} T.$$

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 $T - |\mathcal{T}| \leq T_\epsilon + \frac{\epsilon}{\ell} T$. 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

$$\begin{aligned} \sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) &\leq \sum_{t \in \mathcal{B}_T} \sum_{u=1}^{s(t)} \ell(\hat{Y}_{t^u}, Y_{t^u}) + (T - |\mathcal{T}|) \bar{\ell} \\ &\leq \sum_{t \in \mathcal{B}_T} \sum_{u=1}^{s(t)} \ell(y^t, Y_{t^u}) + \sum_{t \in \mathcal{B}_T} s(t) \delta_t(y^t) + \epsilon \sum_{t \in \mathcal{B}_T} s(t) + T_\epsilon \bar{\ell} + \epsilon T \\ &\leq \sum_{t \in \mathcal{B}_T} \sum_{u=1}^{s(t)} \ell(y^t, Y_{t^u}) + \sum_{t \in \mathcal{B}_T} s(t) \delta_t(y^t) + T_\epsilon \bar{\ell} + 2\epsilon T. \end{aligned}$$

Now let $f : \mathcal{X} \rightarrow \mathcal{Y}$ be a measurable function to which we compare f^ϵ . By Theorem 4.2, 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 4.1 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

$$\begin{aligned} \sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) &\leq \sum_{t \in \mathcal{B}_T} \sum_{u=1}^{s(t)} ((1 + \epsilon) \ell(f(X_{t^u}), Y_{t^u}) + c_\epsilon^\alpha \ell(f(X_t), f(X_{t^u}))) + \sum_{t \in \mathcal{B}_T} s(t) \delta_t(y^t) + T_\epsilon \bar{\ell} + 2\epsilon T \\ &\leq \sum_{t=1}^T \ell(f(X_t), Y_t) + c_\epsilon^\alpha \sum_{u=0}^{T_\epsilon-1} \sum_{t=1}^T \ell(f(X_{\phi^u(t)}), f(X_t)) + \sum_{t \in \mathcal{B}_T} s(t) \delta_{\phi(t)}(y^t) + T_\epsilon \bar{\ell} + 3\epsilon T, \end{aligned}$$

where in the first inequality we used Lemma 2.1. By Lemma 5.5, with probability $1 - \frac{1}{T^2}$, we have

$$\sum_{t \in \mathcal{B}_T} s(t) \delta_t(y^t) \leq 2 [\bar{\ell}(T_\epsilon + 1) + T_\epsilon] \sqrt{T(\ln T + \ln 2/2)}.$$

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

$$\sum_{t=1}^T \ell(\hat{Y}_t(\epsilon), Y_t) \leq \sum_{t=1}^T \ell(f(X_t), Y_t) + c_\epsilon^\alpha \sum_{u=0}^{T_\epsilon-1} \sum_{t=1}^T \ell(f(X_{\phi^u(t)}), f(X_t)) + 2 [\bar{\ell}(T_\epsilon + 1) + T_\epsilon] \sqrt{T(\ln T + \ln 2/2)} + T_\epsilon \bar{\ell} + 3\epsilon T.$$

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,

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

i.e., almost surely, the learning rule f^ϵ achieves risk at most 3ϵ 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

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T [\ell(\hat{Y}_t(\epsilon_i), Y_t) - \ell(f(X_t), Y_t)] \leq 3\epsilon_i, \quad \forall i \geq 0.$$

The rest of the proof uses similar arguments as in the proof of Theorem 4.1. Precisely, let \mathcal{H} be the almost sure event of Lemma 4.4 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 $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq 3\epsilon_i$. 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. \square

We now prove Lemma 5.5 which requires careful analysis.

Proof of Lemma 5.5. For any $t \in \mathcal{B}_T$ and $y \in \mathcal{Y}$, by construction of the learning rule f^ϵ , we have

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

Also, observe that the quantities L_t were constructed precisely so that for any $t \in \mathcal{B}_T$, $s(t)$ has the same distribution as τ_ϵ . The randomness in the construction of L_t for $t \geq 1$ will be used to view the above equation as a realization of the same sum where $s(t)$ is a random stopping time, which allows to use the hypothesis on $g_{\leq \tau_\epsilon}^\epsilon$. We will then use concentration inequalities on the sequence formed by $s(t)\delta_t(y^t)$ for all $t \in \mathcal{B}_T$. Unfortunately, neither of these steps can be performed directly, because the values Y_{t^u} for $t \in \mathcal{B}_T$ and $1 \leq u \leq s(t)$ are dependent on each other and do not form martingales.

We first fix the sequence of values $(y^t)_{t \geq 1}$. Also, for any $t \leq T_\epsilon$ and sequence $\mathbf{y}_{\leq t-1}$ and value $y \in \mathcal{Y}$, we pose

$$\bar{\ell}(g_t(\mathbf{y}_{\leq t-1}), y) := \mathbb{E} [\ell(g_t(\mathbf{y}_{\leq t-1}), y)],$$

where expectation is taken over the randomness of the learning rule. 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 $\phi^{L_t+1-u}(t)$, 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 with probability $1 - \delta/2$, we have

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

We denote by \mathcal{E} the event where the above inequality is met. We will now reason only on these averaged losses. We introduce for any $t \in \mathcal{B}_T$ the variable

$$\bar{\delta}_t(y) := \frac{1}{s(t)} \sum_{u=1}^{s(t)} (\bar{\ell}(\hat{Y}_{t^u}, Y_{t^u}) - \ell(y, Y_{t^u}) - \epsilon).$$

The hypothesis on the learning rule $g_{\leq \tau_\epsilon}^\epsilon$ now becomes: for any $y \in \mathcal{Y}$ and any sequence of values \mathbf{y} ,

$$\sum_{t=1}^{T_\epsilon} \mathbb{P}[\tau_\epsilon = t] \left[\sum_{u=1}^t (\bar{\ell}(g_t(\mathbf{y}_{\leq u-1}), y_u) - \ell(y, y_u)) - \epsilon t \right] \leq 0. \quad (2)$$

We now reason on the process \mathbb{Y} . This process is arbitrary and allowed to depend on the randomness of the learning rule when revealed. In particular, \mathbb{Y} may depend on the realizations of U_t for $t \geq 1$, which define the stopping times $s(t)$ for all $t \in \mathcal{B}_T$. We will only focus on the times in $\mathcal{T} = \bigcup_{t \in \mathcal{B}_T} \{t^u, 1 \leq u \leq s(t)\}$. We reason conditionally to a realization of the tree formed by ϕ . Formally, we consider a specific realization of the process \mathbb{X} and the corresponding rule $(1 + \delta_\epsilon)\text{C1NN}$ which is used in our learning rule to construct the tree with parent relations given by ϕ . Hence, we can now suppose that \mathbb{X} and ϕ are fixed. Now note that because times $t \in \mathcal{B}_T$ were chosen so that the starting point of their corresponding block did not have any descendance in $\mathcal{A}_0 \cup \mathcal{A}_2$ at distance $T_\epsilon - 1$, irrespective of the horizon value $s(t)$, which is bounded by T_ϵ , the block started in $t \in \mathcal{B}_T$ would always have time to end within time horizon T . Formally, on the realization of \mathbb{X} , \mathbb{Y} is an online process which depends on the past values of $U_{\leq t-1}$. As a result, we can view the ‘‘adversarial’’ process $\mathbb{Y}_{\cdot \in \mathcal{T}}$ as a decision process which takes into account the values given by the realizations U_t for $t \geq 1$ in an online fashion. Consider a realization \mathbf{Y} of this random decision process, and enumerate the values in $\mathcal{B}_T := \{t_1 \leq \dots \leq t_{|\mathcal{B}_T|}\}$. Recall that we have $|\mathcal{B}_T| \leq T$. Now, \mathbf{Y} can be written into a binary decision tree \mathcal{G} as follows.

There are two types of nodes, splitting nodes S and leaves L so that $V = S \cup L$. Each splitting node $v \in S$ contains a value $z_v \in \mathcal{Y}$ and an index $i_v \in \{1, \dots, T\}$ and intuitively corresponds to the process trying to add value z_v to the block of index i_v , which started at t_{i_v} . Splitting nodes always have a right children and potentially a left children, while leaves do not have children. This tree allows to construct the process \mathbb{Y} taking into account the realizations of U_t for $t \geq 1$. For instance, the root r —which is always a splitting node—is such that historically, the first started block has index i_r , i.e., $i_r = \arg \min\{t_i, 1 \leq i \leq |\mathcal{B}_T|\} = 1$.

The process tries to add the value z_r to the current block i_r . Precisely, we have $Y_{t_{i_r}^1} = z_r$ if the corresponding sample $U_{t_{i_r}^1}$ allows to add a value to block i_r —in this case, this is almost surely the case because $\mathbb{P}[\tau_\epsilon \geq 1] = 1$ and the block i_r is still empty at this point. Depending on the result of this sample, the decision process can follow two scenarios. Either, the value was added to the block i_r successively, in which case we follow the next left children of the root; or the value was rejected, in which case we follow the right children of the root. More generally, arrived at a given splitting node $v \in S$, we try adding the value z_v to the current block i_v started at time t_{i_v} . If the block already contained $u - 1$ values, we draw $U_{t_{i_v}^u}$.

- If $U_{t_{i_v}^u} < \frac{\mathbb{P}[\tau_\epsilon \geq u]}{\mathbb{P}[\tau_\epsilon \geq u-1]}$, the value is successfully added $Y_{t_{i_v}^u} = z_v$. Provided that $s(t_{i_v}) \geq u - 1$, this case corresponds precisely to $s(t_{i_v}) \geq u$. We then move to the left children of v if it exists. If it does not exist, the process ends.
- Otherwise, the value is rejected. This corresponds to the case when $u - 1 = s(t_{i_v})$, where the block i_v was already complete. We then move to the right children of v .

In practice, the process never ends at a splitting node: we only allow splitting nodes to not have a left children if the value is always rejected, otherwise this does not form a complete decision rule for the process \mathbf{Y} . Thus, the process ends when it arrives at a leaf node. The binary decision tree is constructed so that from a node $v \in S$, we moved to the right node when the block i_v was completed, or in other words, we already added exactly $s(t_{i_v})$ values to this block.

As a remark, one can note that since \mathbb{X} and ϕ are already sampled, there is no choice from \mathbf{Y} for the indices on which to split at each node of the binary decision tree. Indeed, these are defined by the historical ordering of the variables t_i^u for all $i \leq |\mathcal{B}_T|$ and $1 \leq u \leq s(t)$ through the tree defined by ϕ —recall that by construction of \mathcal{B}_T , independently of the values of U_t for $t \geq 1$, started blocks from \mathcal{B}_T will be completed. The random process $\mathbb{Y}|\mathbb{X}, \phi$ can only induce randomness over the values z_v at each splitting node $v \in S$. Hence a realization \mathbf{Y} corresponds to a specific realization of the values z_v for $v \in S$. However, this additional constraint on \mathbf{Y} is not useful to derive our estimations. In other words, even if the process $\mathbb{Y}|\mathbb{X}, \phi$ was allowed to choose which variable to split upon at each node, the property that once a block is started it has enough time to be completed, would be sufficient for the following analysis to hold.

The corresponding decision tree has the property that from any node $v \in S$ the right sub-tree rooted in the right children of v (if it exists) does not have splitting nodes with the same index i_v as v . Indeed, we moved to this sub-tree if the block i_v was already complete. Hence, the future decision process never tries to add an additional value to block i_v . Further, because of the definition of \mathcal{B}_T which ensures that any block can end within time horizon T irrespective of the stopping times $s(t) \leq T_\epsilon$, for any leaf l of the decision tree and any node v in the path from the leaf to the root r , the block i_v was completed, i.e., we moved right on a node which had same index i_v . For a given node v , we will denote by $N_v(j)$ the number of values which have been previously added to block j in the path from the root to v . In particular, $N_v(i_v)$ counts the number of times a splitting on index i_v was performed, the vertex v not included. Last, because $\tau_\epsilon \leq T_\epsilon$, if a node v satisfies $N_v(i_v) = T_\epsilon$, then the value z_v will always be rejected and we can prune the left subtree rooted at the left child of v if it exists. Hence we can prune the decision tree and suppose without loss of generality that for any splitting node $v \in S$, we have $N_v(i_v) \leq T_\epsilon$. As a result of these remarks, one can note that leaves can only be right children.

Recall that by design of the learning rule f^ϵ , at any time $t \leq T$, the value Y_t is independent from the sampling U_t . As a result, the precise realization of \mathbb{Y} considering the realizations of U_t for $t \geq 1$, corresponds to a random walk on the decision tree. Precisely, arrived at a splitting node $v \in S$, the process follows its left children with probability $\frac{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)+1]}{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)]}$ —which is zero if $N_v(i_v) = T_\epsilon$ and more precisely, the process *never* follows its left children in this case—and to its right children with probability $1 - \frac{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)+1]}{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)]}$. Thanks to this decision tree, we will upper bound for this specific realization \mathbf{Y} and any $\lambda > 0$ the quantity

$$\Gamma_\lambda := \mathbb{E} \left[e^{\lambda \sum_{t \in \mathcal{B}_T} s(t) \bar{\delta}_t(y^t)} \right],$$

where the expectation is taken over U_t for $t \in \mathcal{T}$. At any node $v \in V$, we define the quantity

$$\Gamma_\lambda(v) := e^{-\lambda \sum_{j=1}^T \sum_{u=1}^{N_v(j)} (\bar{\ell}(\hat{Y}_{t_j^u}, z_{v,u}^j) - \ell(y^{t_j}, z_{v,u}^j) - \epsilon)} \mathbb{E} \left[e^{\lambda \sum_{t \in \mathcal{B}_T} s(t) \delta_\varphi(t)(y^t)} \middle| v \right]$$

where $z_{v,u}^j$ denotes the value $z_{v'}$ of the u -th splitting node on index j along the path from the root r to v , and the conditioning on v denotes the event that the decision process has passed through node v . In particular, on this event, for any $1 \leq j \leq T$ and $1 \leq u \leq N_v(j)$ we have $Y_{t_j^u} = z_{v,u}^j$, by construction of the decision tree process. The other values of Y are not revealed yet when the process arrives at node v . Hence, the quantity $\Gamma_\lambda(v)$ corresponds to the expected future contribution to the sum $\lambda \sum_{t \in \mathcal{B}_T} s(t) \delta_\varphi(t)(y)$ of the still unknown part at this step of the process. The first factor in $\Gamma_\lambda(s)$ precisely deletes the contribution of past decided values Y . An important remark is that $\Gamma_\lambda(r) = \Gamma_\lambda$.

We now define for all $0 \leq u \leq T_\epsilon$ and $y \in \mathcal{Y}$ and any sequence $\mathbf{z}_{\leq u}$,

$$\gamma_\lambda(y, \mathbf{z}_{\leq u}) = \sup_{z_{u+1}, \dots, z_{T_\epsilon} \in \mathcal{Y}} \frac{1}{\mathbb{P}[\tau_\epsilon \geq u]} \sum_{k=u}^{T_\epsilon} \mathbb{P}[\tau_\epsilon = k] e^{\lambda \sum_{i=u+1}^k (\bar{\ell}(g_i(\mathbf{z}_{\leq u}, \mathbf{z}_{u+1 \leq i \leq k-1}, z_i) - \ell(y, z_i) - \epsilon)},$$

with the convention $\frac{0}{0} = 1$. In other words, $\gamma_\lambda(y, \mathbf{z}_{\leq u})$ corresponds to the worst possible contribution of future values for the rule g , given the first values $\mathbf{z}_{\leq u}$. Note that if $u = T_\epsilon$, all values have been decided so the future worst case contribution is always null. For any $v \in V$, we denote A_v the set of indices $1 \leq i \leq T$ such that the block i was started in the past but potentially not completed, i.e., there exists a node $v' \neq v$ in the path from the root r to v with $i_{v'} = i$, but there does not exist a node v' from the root to v such that $i_{v'} = i$ and after v' we moved to his right child. We also introduce B_v the set of indices $1 \leq i \leq T$ which do not appear in a split from any node $v' \neq v$ in the path from r to v . This corresponds to blocks that have not been started but potentially could be in the future. We now prove by induction that for any node $v \in V$,

$$\Gamma_\lambda(v) \leq \prod_{j \in A_v} \gamma_\lambda \left(y^{t_j}, \mathbf{z}_{v, \leq N_v(j)}^j \right) \prod_{j \in B_v} \max(1, \gamma_\lambda(y^{t_j})). \quad (3)$$

We start the induction from leaves. For any leaf $v \in L$, because there are no future values to reveal we have $\Gamma_\lambda(v) = 1$. Also, because it is a leaf which may end the decision process, we have $A_v = \emptyset$. Hence Eq (3) holds. We now show the induction. Let $v \in S$ such that all its children satisfy Eq (3). We denote by v_l (resp. v_r) the left (resp. right) child of v , if it exists. Then, recalling that from v we move to v_l with probability $\frac{\mathbb{P}[\tau_\epsilon \geq N_v(i_v) + 1]}{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)]}$ and that this probability is null if v does not have a left children, we have

$$\begin{aligned} \Gamma_\lambda(v) &= \frac{\mathbb{P}[\tau_\epsilon \geq N_v(i_v) + 1]}{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)]} e^{-\lambda \sum_{j=1}^T \sum_{u=1}^{N_v(j)} (\bar{\ell}(\hat{Y}_{t_j^u}, z_{v,u}^j) - \ell(y^{t_j}, z_{v,u}^j) - \epsilon)} \mathbb{E} \left[e^{\lambda \sum_{t \in \mathcal{B}_T} s(t) \delta_t(y^t)} \middle| v_l \right] \\ &+ \frac{\mathbb{P}[\tau_\epsilon = N_v(i_v)]}{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)]} e^{-\lambda \sum_{j=1}^T \sum_{u=1}^{N_v(j)} (\bar{\ell}(\hat{Y}_{t_j^u}, z_{v,u}^j) - \ell(y^{t_j}, z_{v,u}^j) - \epsilon)} \mathbb{E} \left[e^{\lambda \sum_{t \in \mathcal{B}_T} s(t) \delta_t(y^t)} \middle| v_r \right] \\ &= \frac{\mathbb{P}[\tau_\epsilon \geq N_v(i_v) + 1]}{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)]} \exp \left[\lambda \left(\bar{\ell} \left(\hat{Y}_{t_{N_v(i_v)}}^{N_v(i_v)}, z_v \right) - \ell(y^{t_{N_v(i_v)}}, z_v) - \epsilon \right) \right] \Gamma_\lambda(v_l) + \frac{\mathbb{P}[\tau_\epsilon = N_v(i_v)]}{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)]} \Gamma_\lambda(v_r). \end{aligned}$$

Now note that if v_l exists, then we necessarily have $i_v \in A_{v_l}$. Further, we always have $i_v \notin A_{v_r}$ because if we moved right after node v , then the block i_v was completed. As a result, we have $A_{v_r} = A_{v_l} \setminus \{i_v\} = A_v \setminus \{i_v\}$ and also $B_{v_l} = B_{v_r} = B_v \setminus \{i_v\}$. Finally, for any $j \neq i_v$, we have $N_v(j) = N_{v_r}(j) = N_{v_l}(j)$. Therefore, we

obtain

$$\begin{aligned}
\Gamma_\lambda(v) &\leq \frac{\mathbb{P}[\tau_\epsilon \geq N_v(i_v) + 1]}{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)]} \exp \left[\lambda \left(\bar{\ell} \left(\hat{Y}_{t_{i_v}^{N_v(i_v)}}, z_v \right) - \ell(y^{t_{i_v}}, z_v) - \epsilon \right) \right] \gamma_\lambda \left(y^{t_{i_v}}, \mathbf{z}_{\mathbf{v}, \leq N_v(j)}^{i_v}, z_v \right) \\
&\quad \cdot \prod_{j \in A_{v_i} \setminus \{i_v\}} \gamma_\lambda \left(y^{t_j}, \mathbf{z}_{\mathbf{v}, \leq N_v(j)}^j \right) \prod_{j \in B_{v_i}} \max(1, \gamma_\lambda(y^{t_j})) \\
&\quad + \frac{\mathbb{P}[\tau_\epsilon = N_v(i_v)]}{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)]} \prod_{j \in A_{v_r}} \gamma_\lambda \left(y^{t_j}, \mathbf{z}_{\mathbf{v}, \leq N_v(j)}^j \right) \prod_{j \in B_{v_r}} \max(1, \gamma_\lambda(y^{t_j})) \\
&\leq \left[\frac{\mathbb{P}[\tau_\epsilon = N_v(i_v)]}{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)]} + \frac{\mathbb{P}[\tau_\epsilon \geq N_v(i_v) + 1]}{\mathbb{P}[\tau_\epsilon \geq N_v(i_v)]} \right] \exp \left[\lambda \left(\bar{\ell} \left(\hat{Y}_{t_{i_v}^{N_v(i_v)}}, z_v \right) - \ell(y^{t_{i_v}}, z_v) - \epsilon \right) \right] \gamma_\lambda \left(y^{t_{i_v}}, \mathbf{z}_{\mathbf{v}, \leq N_v(j)}^{i_v}, z_v \right) \\
&\quad \cdot \prod_{j \in A_v \setminus \{i_v\}} \gamma_\lambda \left(y^{t_j}, \mathbf{z}_{\mathbf{v}, \leq N_v(j)}^j \right) \prod_{j \in B_v \setminus \{i_v\}} \max(1, \gamma_\lambda(y^{t_j}))
\end{aligned}$$

Now recall that γ_λ was constructed as the worst future contribution possible. Essentially, y_v could be optimized to get a worse contribution which yields the following inequality.

$$\begin{aligned}
\Gamma_\lambda(v) &\leq \gamma_\lambda \left(y^{t_{i_v}}, \mathbf{z}_{\mathbf{v}, \leq N_v(j)}^{i_v} \right) \prod_{j \in A_v \setminus \{i_v\}} \gamma_\lambda \left(y^{t_j}, \mathbf{z}_{\mathbf{v}, \leq N_v(j)}^j \right) \prod_{j \in B_v \setminus \{i_v\}} \max(1, \gamma_\lambda(y^{t_j})) \\
&= \prod_{j \in A_v} \gamma_\lambda \left(y^{t_j}, \mathbf{z}_{\mathbf{v}, \leq N_v(j)}^j \right) \prod_{j \in B_v \setminus \{i_v\}} \max(1, \gamma_\lambda(y^{t_j})) \cdot [\mathbb{1}(i_v \notin B_v) + \mathbb{1}(i_v \in B_v) \gamma_\lambda(y^{t_{i_v}})] \\
&\leq \prod_{j \in A_v} \gamma_\lambda \left(y^{t_j}, \mathbf{z}_{\mathbf{v}, \leq N_v(j)}^j \right) \prod_{j \in B_v} \max(1, \gamma_\lambda(y^{t_j})).
\end{aligned}$$

This ends the recursion and shows that Eq (3) holds for all $v \in V$, in particular for the root r . At the root, because no values have yet been revealed, we have $A_r = \emptyset$ and $B_r = \{1 \leq i \leq T\}$. As a result we obtain the desired result

$$\Gamma_\lambda = \Gamma_\lambda(r) \leq \prod_{1 \leq j \leq T} \max(1, \gamma_\lambda(y^{t_j})) \leq \max(1, \gamma_\lambda)^T,$$

where $\gamma_\lambda := \sup_{y \in \mathcal{Y}} \gamma_\lambda(y)$. We now give an upper bound on γ_λ . To do so, we note that for any $y, z_1, \dots, z_{T_\epsilon} \in \mathcal{Y}$, we have $-T\bar{\ell} - T\epsilon \leq \sum_{u=1}^t (\bar{\ell}(g_u(\mathbf{z}_{\leq u-1}), z_u) - \ell(y, z_u) - \epsilon) \leq T\bar{\ell}$ for any $1 \leq t \leq T_\epsilon$, so we can apply Hoeffding's lemma.

$$\begin{aligned}
\gamma_\lambda &= \sup_{y, y_1, \dots, y_{T_\epsilon}} \mathbb{E}_{\tau_\epsilon} \left[e^{\lambda \sum_{u=1}^{\tau_\epsilon} (\bar{\ell}(g_u(\mathbf{y}_{\leq u-1}), y_u) - \ell(y, y_u) - \epsilon)} \right] \\
&\leq \sup_{y, y_1, \dots, y_{T_\epsilon}} \exp \left(\lambda \mathbb{E}_{\tau_\epsilon} \left[\sum_{u=1}^{\tau_\epsilon} (\bar{\ell}(g_u(\mathbf{y}_{\leq u-1}), y_u) - \ell(y, y_u) - \epsilon) \right] + \lambda^2 \frac{T_\epsilon^2 (2\bar{\ell} + \epsilon)^2}{8} \right) \\
&\leq e^{\lambda^2 T_\epsilon^2 (\bar{\ell} + 1)^2 / 2},
\end{aligned}$$

where in the second inequality we used the hypothesis Eq (2) on the learning rule $g_{\leq \tau_\epsilon}$. As a result, we obtain $\Gamma_\lambda \leq e^{\lambda^2 \bar{\ell}^2 T / 2}$. We can now apply standard Markov inequalities for concentration. Let $\alpha > 0$. Then,

$$\mathbb{P} \left[\sum_{t \in \mathcal{B}_T} s(t) \bar{\delta}_t(y^t) \geq \alpha \right] \leq \min_{\lambda > 0} \exp \left(-\lambda \alpha + \lambda^2 \frac{T_\epsilon^2 (\bar{\ell} + 1)^2 T}{2} \right) = \exp \left(-\frac{\alpha^2}{2T_\epsilon^2 (\bar{\ell} + 1)^2 T} \right),$$

and as a result, with probability at least $1 - \delta/2$ we have

$$\sum_{t \in \mathcal{B}_T} s(t) \bar{\delta}_t(y^t) < T_\epsilon (\bar{\ell} + 1) \sqrt{2T \ln \frac{2}{\delta}}.$$

Because this holds for any realization of \mathbb{X} and \mathbf{Y} , we finally obtain that with probability $1 - \delta/2$ the same inequality holds where the probability is taken over \mathbb{X} , \mathbb{Y} and the learning rule together. We denote \mathcal{F} the event where the above inequality holds. Then, on $\mathcal{E} \cap \mathcal{F}$ which has probability at least $1 - \delta$ by the union bound, we have

$$\sum_{t \in \mathcal{B}_T} s(t) \delta_t(y^t) \leq \sum_{t \in \mathcal{B}_T} s(t) \bar{\delta}_t(y^t) + \bar{\ell} \sqrt{2T \ln \frac{2}{\delta}} \leq [\bar{\ell}(T_\epsilon + 1) + T_\epsilon] \sqrt{2T \ln \frac{2}{\delta}}.$$

This ends the proof of the lemma. \square

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 $\mathbf{y} := (y_t)_{t=1}^T$ of values in \mathcal{Y} and any value $y \in \mathcal{Y}$, we have*

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

This equation can be conveniently rewritten as

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

which hints to the similarity with F-TIME. We prove their equivalence in the next lemma.

Lemma 5.6. *Property F-TIME is equivalent to Property 2.*

Proof. We first start by showing that Condition 1 implies Condition 2. Let (\mathcal{Y}, ℓ) satisfying Condition 1. 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

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

We now construct a random time $1 \leq \tilde{\tau} \leq T$ such that $\mathbb{P}[\tilde{\tau} = t] = \frac{t \mathbb{P}[\tau = t]}{\mathbb{E}[\tau]}$ for all $1 \leq t \leq T$. This indeed defines a proper random variable because $\sum_{t=1}^T \frac{t \mathbb{P}[\tau = t]}{\mathbb{E}[\tau]} = 1$. Let $Supp(\tau) := \{1 \leq t \leq T : \mathbb{P}[\tau = t] > 0\}$ be the support of τ . For any $t \in Supp(\tau)$, 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 $t \in Supp(\tau)$ is to first draw an i.i.d. sequence of learning rules $(g_{i, \leq \tau_i})_{i \geq 1}$. Then, for each $t \in Supp(\tau)$ and select the randomness which first satisfies $\tau = t$. Specifically, we define the time $i_t = \min\{i : \tau_i = t\}$ for all $t \in Supp(\tau)$. With probability one, these times are finite for all $t \in Supp(\tau)$. 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

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

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

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

where in the third and fifth equality we used the fact that $\mathbb{P}[\mathcal{E}] = 1$. This ends the proof that Condition 1 implies Condition 2.

The other implication can be proved using the same technique. Suppose that (\mathcal{Y}, ℓ) satisfies Condition 2. Let $\eta > 0$ and $T \geq 1, g_{\leq \tau}$, where $1 \leq \tau \leq T$ such that

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

We first construct a random time $1 \leq \tau \leq T$ such that $\mathbb{P}[\tilde{\tau}] = \frac{\mathbb{P}[\tau=t]}{t\mathbb{E}[\frac{1}{\tau}]}$ for all $1 \leq t \leq T$. We then construct a learning rule $\tilde{g}_{\leq \tilde{\tau}}$ similarly as before. Using the same arguments, we obtain

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

This ends the proof of the lemma. \square

We are now ready to prove our main result for the second alternative when F-TIME is not satisfied. Specifically, we show that in this case, universal learning for adversarial responses under processes outside CS under arbitrary responses is not achievable.

Theorem 3.2. *Suppose that ℓ is bounded and (\mathcal{Y}, ℓ) does not satisfy F-TIME. Then, SOLAR = 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. 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 5.6, 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

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

Because 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 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

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

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

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

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 $u \leq \tau_p$, $\hat{Y}_{T_p(u)}$ has the same distribution as $g_u^p(\mathbf{y}_p^{\leq u-1})$. Therefore we obtain

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

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

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

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 $\alpha = 1$. \square

This closes our study of universal learning with adversarial responses for bounded value spaces. Interestingly, we showed that there always exists an optimistically universal learning rule for adversarial regression, however this rule highly depends on the value space. Namely, if (\mathcal{Y}, ℓ) satisfies F-TIME, we can learn all SMV = SOUL processes. The proposed learning rule of Theorem 3.1 is *implicit* in general. Indeed, to construct it one first needs to find an online learning rule for mean estimation with finite horizon as explicated by property F-TIME, which is then used as subroutine in the optimistically universal learning rule for adversarial regression. We showed however that for totally-bounded value spaces, this learning rule can be *explicitated* using ϵ -nets for decreasing $\epsilon > 0$.

If the value space does not satisfy F-TIME, we showed that we can only learn $\text{CS} \subsetneq \text{SOUL}$ processes, and we propose a learning rule in Section 7 which is optimistically universal—see Theorem 3.5. This learning rule is inspired by the proposed algorithm of [Han22] which is optimistically universal for metric losses $\alpha = 1$. It is worth noting that this learning rule uses very different techniques than our first proposed algorithm for value spaces satisfying F-TIME. 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. Intuitively it asks that for any measurable function $f^* : \mathcal{X} \rightarrow \mathcal{Y}$, the long-run empirical average losses of f^* compared to a function $f \in \mathcal{F}$ along the instances X_1, X_2, \dots are arbitrarily small. Hence, under process \mathbb{X} , we can approximate f^* by functions in \mathcal{F} with arbitrary long-run average precision. Such property is impossible to obtain for any process $\mathbb{X} \in \text{SMV} \setminus \text{CS}$. Indeed, [Han21a] showed that having such a “dense” countable set of measurable functions under a given process \mathbb{X} is a sufficient condition for universal inductive or self-adaptive learning in the noiseless setting. However, in the same work, it is shown that the set of processes learnable for inductive or self-adaptive noiseless learning is precisely CS. As a result, no process $\mathbb{X} \notin \text{CS}$ admits a “dense” countable sequence of measurable functions. This implies that in order to learn

processes SMV for value spaces satisfying F-TIME, a *fundamentally* different learning rule than that proposed by [Han21a] or [Han22] was needed. Further, the alternative SOLAR = CS or SOLAR = SMV(= SOUL) shows that there is an inherent gap between noiseless and noisy regression for some non-totally-bounded value spaces.

6 Adversarial universal learning for unbounded losses

We now turn to the case of unbounded losses. We say that the value space (\mathcal{Y}, ℓ) is unbounded when $\sup_{y_1, y_2 \in \mathcal{Y}} \ell(y_1, y_2) = \infty$. In this case, and for more general near-metrics, [BCH22] showed that SOUL = 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. This shows that universal learning on unbounded value spaces is very restrictive and in particular we have SOLAR \subset FS. The main question is whether we can recover the complete set of processes FS for adversarial regression. We will show that there is again an alternative. We will obtain either SOLAR = FS or SOLAR = \emptyset .

6.1 Adversarial regression for metric losses

In this section, we focus on metric losses ℓ , i.e., $\alpha = 1$. We will show in this section that we have in fact an equality SOLAR = FS and that we can provide an optimistically universal learning rule. To do so, we first prove a general result on mean estimation. We refer as mean estimation 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. This is equivalent to regression where the input space is $\mathcal{X} = \{0\}$. Note that 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) .

Theorem 3.3. *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 $y \in \mathcal{Y}$,*

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

Proof. Consider the following algorithm. Fix any $(y^i)_{i \geq 0}$ a sequence of values dense in \mathcal{Y} . Define $I_t = \{i : i \leq \ln t, \ell(y^0, y^i) \leq \ln t\}$, for any $i \geq 0$, denote $t_i = \lceil \max(e^i, e^{\ell(y^0, y^i)}) \rceil$ and pose $\eta_t = \frac{1}{4\sqrt{t}}$. For any $i \in I_t$ we pose $L_{t-1, i} := \sum_{s=t_i}^{t-1} \ell(y^i, Y_s)$ and 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 $\hat{\ell}_s := \frac{\sum_{j \in I_s} w_{s-1, j} \ell(y^j, Y_s)}{\sum_{j \in I_s} w_{s-1, j}}$ 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 $w_{t-1, i} := e^{\eta_t(\hat{L}_{t-1, i} - L_{t-1, i})}$. Finally, we name indices $I_t = \{i_1, \dots, i_{|I_t|}\}$ and construct the prediction \hat{Y}_t such that this prediction is independent of the past history and for any $i \in I_t$, we have

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

Note that the random prediction \hat{Y}_t only uses the values Y_1, \dots, Y_{t-1} hence defines a proper online learning rule. We will now 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 4.4 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 4.4, 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

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

Similarly as in Lemma 4.4, 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 $|\ell(\hat{Y}_t, Y_t) - \mathbb{E}[\ell(\hat{Y}_t, Y_t) \mid \mathbb{Y}_{\leq t}]| \leq 2 \ln t$ 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 $1 - \frac{\delta}{t^2}$ 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. \square

Note that, unlike all the results that we showed in the previous sections for universal regression, on the same event of probability one, the propose 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 achieve 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. As a consequence of Theorem 3.3 we obtain the following result for universal regression with adversarial responses.

Corollary 6.1. *Suppose that (\mathcal{Y}, ℓ) is an unbounded metric space. Then, $\text{SOLAR} = \text{FS} (= \text{SOUL})$ and there exists an optimistically universal learning rule for adversarial regression, i.e., which achieves universal consistency with adversarial responses under any process $\mathbb{X} \in \text{FS}$.*

Proof. We denote by g the learning rule on values \mathcal{Y} for mean estimation described in Theorem 3.3. 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 3.3. For all $x \in \tilde{S}$, with probability one, for all $y^x \in \mathcal{Y}$, we have

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

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 $f : \mathcal{X} \rightarrow \mathcal{Y}$,

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

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 consistence 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 FS \subset SOLAR. As result SOLAR = FS and f is optimistically universal. This ends the proof of the result. \square

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

In the previous section, we observed that in general metric spaces, we can recover SOLAR = FS for metric spaces in adversarial regression. Unfortunately, even though FS is an already extremely restrictive set of processes, we will now show that in the classical setting of real-valued regression $\mathcal{Y} = \mathbb{R}$ with Euclidean norm, adversarial regression with any loss $\ell = |\cdot|^\alpha$ for $\alpha > 1$ is not achievable. Specifically, we will show that in that case SOLAR = \emptyset . As a consequence, adversarial regression and even mean-estimation for the classical setting of real-valued regression and squared loss is never achievable.

Theorem 6.2. *Let $\alpha > 1$. For the Euclidean value space $(\mathbb{R}, |\cdot|)$ and loss $\ell = |\cdot|^\alpha$ we obtain SOLAR = \emptyset . In particular, there does not exist a consistent learning rule for mean-estimation on \mathbb{R} with squared loss for adversarial responses.*

Proof. 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

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

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 $\text{SOLAR} = \emptyset$. 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, $\mathbb{X} \notin \text{SOLAR}$, and as a result, $\text{SOLAR} = \emptyset$. \square

The above proof 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 exists 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

The previous two sections show that there exist losses for which we obtain $\text{SOLAR} = \emptyset$ or $\text{SOLAR} = \text{FS}$. In this section, we show the simple result that this is the only alternative, and that having non-empty $\text{SOLAR} = \text{FS}$ is equivalent to achieving consistency for mean-estimation with adversarial responses.

Proposition 6.3. *Suppose that there exists an online learning rule g which is consistent for mean-estimation with adversarial responses, i.e., for any adversarial process \mathbb{Y} on (\mathcal{Y}, ℓ) , we have for any $y \in \mathcal{Y}$,*

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

then $\text{SOLAR} = \text{FS}$ and there exists an optimistically universal learning rule for adversarial regression. On the other hand, if mean-estimation is not achievable for adversarial responses, $\text{SOLAR} = \emptyset$.

Proof. Suppose that there exists an online learning rule g for mean-estimation. In the proof of Corollary 6.1, instead of using the learning rule for mean-estimation for metric losses introduced in Theorem 3.3, 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, $\text{FS} \subset \text{SOLAR}$. Because $\text{SOLAR} \subset \text{SOUL} = \text{FS}$, we obtain directly $\text{SOLAR} = \text{FS}$ 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 6.2. 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. \square

As a note, although in the Euclidean real-valued case we obtained a simple alternative $\text{SOLAR} = \text{FS}$ if $\alpha = 1$ and $\text{SOLAR} = \emptyset$ for $\alpha > 1$, one cannot hope to simplify the characterization of Corollary 6.3 so that in general for any power of a metric with $\alpha > 1$ we would obtain $\text{SOLAR} = \emptyset$. Indeed, consider the case $\mathcal{Y} = \mathbb{R}$ but with the metric $|\cdot| = \sqrt{|\cdot|_2}$ square root of the classical Euclidean metric. One can check that this does define a proper metric because it satisfies the triangular inequality. In that case, we obtain that for $\alpha \leq 2$, we have $\text{SOLAR} = \text{FS}$ and $\text{SOLAR} = \emptyset$ if $\alpha > 2$.

7 Adversarial universal learning with moment constraint

In the previous section, we showed that the only learnable processes for adversarial regression are processes in FS, i.e., which visit 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 account for i.i.d. processes. A natural question is whether adding mild constraints on the process \mathbb{Y} would allow to recover the same results for unbounded losses as for bounded losses from Section 4 and 5. In fact, this question 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 open problem 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 open question now becomes whether there exist 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 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 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. Formally, set CRF is defined as follows.

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 we have $\text{CRF} \subset \text{CS}$. In particular, we have $\text{CRF} \subset \text{SMV}$. 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 $\mathbb{X} \in \text{CRF}(\subset \text{SMV})$ 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. 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

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

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

Theorem 7.1 completely answers (negatively to) the open problem 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 introduce a novel constraint, similar to that introduced by [BCH22], but slightly stronger, which we will refer to as the *empirical integrability* constraint. An (adversarial) process \mathbb{Y} is *empirically integrable* i.f 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} . 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 $M > \bar{\ell}$. Hence, for bounded value spaces, the moment constraint is not restrictive. 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

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

is empirically integrable. Note that for i.i.d. 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, in this case, we obtain almost surely $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, Y_t) \mathbb{1}_{\ell(y_0, Y_t) \geq M} = \mathbb{E}[\ell(y_0, Y_1) \mathbb{1}_{\ell(y_0, Y_1) \geq M}]$. Then, using the dominated convergence theorem allows to find a suitable M for any fixed tolerance $\epsilon > 0$. We formally prove this in the next lemma.

Lemma 7.2. *Let \mathbb{Y} an i.i.d. 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 i.i.d. 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 $\mathbb{E}[\ell(y_0, Y_1) \mathbb{1}_{\ell(y_0, Y_1) \geq M}] \leq \epsilon$. Now by the law of large numbers, almost surely, we have

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

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

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. We first prove a simple equivalent formulation for empirical integrability.

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

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

Proof. 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. \square

7.1 Noiseless universal learning with moment condition

The main result from [Bla22] showed that for bounded value spaces, the 2C1NN learning rule is optimistically universal, i.e., achieves universal consistency on all SMV processes. We now show that the same learning rule is consistent under all SMV processes for noiseless responses $\mathbb{Y} = f^*(\mathbb{X})$ which are empirically integrable, even for unbounded value spaces.

Theorem 3.4. *Let (\mathcal{Y}, ℓ) 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 there exists $y_0 \in \mathcal{Y}$ for which for all $\epsilon > 0$, there exists $M \geq 0$ with $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(y_0, f^*(X_t)) \mathbb{1}_{\ell(y_0, f^*(X_t)) \geq M} \leq \epsilon$ (a.s.), we have $\mathcal{L}_{\mathbb{X}}(2\text{C1NN}, f^*) = 0$ (a.s.).*

Proof. 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.

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(2\text{C1NN}_t(\mathbb{X}_{t-1}, f_M^*(\mathbb{X}_{\leq t-1}), X_t), f_M^*(X_t)) = 0 \quad (\text{a.s.}).$$

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 $M \geq 1$,

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

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

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

Hence, we obtain

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

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

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

In particular, if $\epsilon > 0$ we can apply this result to $M := \lceil M_\epsilon \rceil$, which yields $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(f^*(X_{\phi(t)}), f^*(X_t)) \leq c_\ell(2+c_\ell)\epsilon$. 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. \square

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 a similar online learning rule. In particular, our proposed learning rule will also optimistically universal for adversarial regression for all bounded value spaces which do not satisfy F-TIME.

Theorem 3.5. *There exists an online learning rule f that is universally consistent for adversarial empirically integrable responses under all processes in CS, 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.).*

Proof. Using Lemma 23 of [Han21a], let $\mathcal{T} \subset \mathcal{B}$ a countable set such that for all $\mathbb{X} \in \mathcal{C}_1, A \in \mathcal{B}$ 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 $\hat{y}_t := f^{\hat{i}_t}(x_t)$. Note that the random prediction \hat{y}_t only uses the values $x_1, \dots, x_{t-1}, y_1, \dots, y_{t-1}, x_t$ hence can be used to create an online learning rule which we denote by simplicity $(\hat{Y}_t)_{t \geq 1}$. Now consider a process (\mathbb{X}, \mathbb{Y}) with $\mathbb{X} \in \mathcal{C}_1$ and such that \mathbb{Y} is empirically integrable. By Lemma 7.3, there exists $y_0 \in \mathcal{Y}$ 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 2.1 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 3.3 and show that almost surely, there exists $t \geq 1$ such that

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

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 $T \geq \max(t_{i_k}, \hat{\ell})$ we have for any $\epsilon > 0$,

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

where in the last inequality we used $\ell(y^0, \tilde{Y}_t) \leq \ell(y^0, Y_t)$ 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 $\epsilon > 0$,

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

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

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

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

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

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 $M \geq 1$,

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

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 finally 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 $M \geq 1$,

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

Because this holds for all $M \geq 1$, we can again apply this result to $M := \lceil \tilde{M}_\epsilon \rceil$ which yields $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq \epsilon$. Because this holds for any $\epsilon > 0$, we finally obtain on the event $\mathcal{A} \cap \mathcal{B} \cap \bigcap_{M=1}^{\infty} \mathcal{E}_M$ of probability one, that $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t, Y_t) - \ell(f(X_t), Y_t) \leq 0$. This ends the proof of the theorem. \square

This generalizes the main results from [Han22] to unbounded non-metric losses and from [TK22] 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.

7.3 Adversarial regression with moment condition under SMV processes

Last, we generalize our result Theorem 3.1 for value spaces satisfying F-TIME, to unbounded value spaces, with the same moment condition on responses. In order to apply Theorem 3.1 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. For such value spaces, we will be able to recover learning for adversarial responses under all SMV processes. Note that if this property is not satisfied, then the set of learnable processes for adversarial responses under this moment condition is automatically reduced to CS. Indeed, given a ball $B_\ell(y, r)$ which does not satisfy F-TIME, one can focus on responses which have value in this ball and directly apply the negative result from Theorem 3.2 to show that adversarial regression under processes $\mathbb{X} \notin \text{CS}$ is not achievable.

Theorem 3.6. *Suppose that any ball of (\mathcal{Y}, ℓ) , $B_\ell(y, r)$ satisfies F-TIME. Then, there exists an optimistically universal online learning rule f for adversarial empirically integrable responses with bounded moments, 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, 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. Fix (\mathcal{X}, ρ) and a value space (\mathcal{Y}, ℓ) such that any ball satisfies Condition 1. 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 Condition 1, there exists a Bayes optimistically universal learning rule f^M for value space $(B_\ell(\bar{y}, M), \ell)$. 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 $w_{0,1} = 0$ 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

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

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.

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 3.3 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^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 the quantity $\delta_T^M := \frac{1}{T} \sum_{t=1}^T \ell(\hat{Y}_t^M, \phi_M(Y_t)) - \ell(\phi_M \circ f(X_t), \phi_M(Y_t))$ 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

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

The proof of Theorem 3.5 shows that on the event \mathcal{A} we have $\limsup_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \ell(Y_t, \bar{y}) \mathbb{1}_{\ell(Y_t, \bar{y}) \geq 2^{-\alpha+1} \ln t} \leq 0$, 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

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

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

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

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 3.5 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 $M := \lceil 2^\alpha M_\epsilon + M_{2^{-\alpha-1}\epsilon} \rceil$ 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. \square

Theorems 3.5 and 3.6 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, then Theorem 3.6 provides an optimistic learning rule which achieves consistency under all processes in SMV. On the other hand, if there exists a ball $B_\ell(y, r)$ which disproves F-TIME, then Theorem 3.5 provides an optimistic learning rule which achieves consistency under all processes in 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 $\ell = |\cdot|^\alpha$ for $\alpha \geq 1$. A natural question would be whether one can generalize these results to larger classes of losses, for instance 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 learner which 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.
- [BC21] Moïse Blanchard and Romain Cosson. “Universal Online Learning with Bounded Loss: Reduction to Binary Classification”. In: *arXiv preprint arXiv:2112.14638* (2021).
- [BCH22] Moïse Blanchard, Romain Cosson, and Steve Hanneke. “Universal Online Learning with Unbounded Losses: Memory Is All You Need”. In: *arXiv preprint arXiv:2201.08903* (2022).
- [Bie+19] Armin Biess, Aryeh Kontorovich, Yury Makarychev, and Hanan Zaichyk. “Regression via Kirschbraun Extension”. In: *arXiv preprint arXiv:1905.11930* (2019).
- [Bla22] Moïse Blanchard. “Universal Online Learning: an Optimistically Universal Learning Rule”. In: *arXiv preprint arXiv:2201.05947* (2022).
- [BPS09] Shai Ben-David, Dávid Pál, and Shai Shalev-Shwartz. “Agnostic Online Learning.” In: *COLT*. 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.
- [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 Gyorfi, 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).
- [Fer+11] Frédéric Ferraty, Ali Laksaci, Amel Tadj, and Philippe Vieu. “Kernel regression with functional response”. In: *Electronic Journal of Statistics* 5 (2011), pp. 159–171.
- [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 Gyorfi, 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.
- [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] L. Györfi, M. Kohler, A. Krzyzak, and H. Walk. *A Distribution-Free Theory of Nonparametric Regression*. Springer-Verlag New York, 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: *ALT 2022* (2022).
- [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.
- [RB06] Daniil Ryabko and Peter Bartlett. “Pattern Recognition for Conditionally Independent Data.” In: *Journal of Machine Learning Research* 7.4 (2006).
- [RST15] Alexander Rakhlin, Karthik Sridharan, and Ambuj Tewari. “Online learning via sequential complexities.” In: *J. Mach. Learn. Res.* 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.
- [Tho13] P Thomas Fletcher. “Geodesic regression and the theory of least squares on Riemannian manifolds”. In: *International journal of computer vision* 105.2 (2013), pp. 171–185.
- [TK22] Dan Tsir Cohen and Aryeh Kontorovich. “Metric-valued regression”. In: *arXiv e-prints* (2022), arXiv–2202.