# On Universal Portfolios with Continuous Side Information

**Alankrita Bhatt**[*]
Caltech

**J. Jon Ryu**[*]
MIT

**Young-Han Kim**
UCSD / Gauss Labs Inc.

## Abstract

A new portfolio selection strategy that adapts to a continuous side-information sequence is presented, with a universal wealth guarantee against a class of state-constant rebalanced portfolios with respect to a state function that maps each side-information symbol to a finite set of states. In particular, given that a state function belongs to a collection of functions of finite Natarajan dimension, the proposed strategy is shown to achieve, asymptotically to first order in the exponent, the same wealth as the best state-constant rebalanced portfolio with respect to the best state function, chosen in hindsight from observed market. This result can be viewed as an extension of the seminal work of Cover and Ordentlich (1996) that assumes a single-state function.

## 1 INTRODUCTION

We study the classical problem of portfolio selection, formally defined as follows. Suppose that there exist $m \geq 2$ stocks in a stock market and let $\mathbf{x}_t = (x_{t1}, \ldots, x_{tm}) \in \mathbb{R}_{\geq 0}$ denote a market vector at time $t$, which encodes the *price relatives* of stocks on that day. That is, for each stock $i \in [m] := \{1, \ldots, m\}$, $x_{ti} \geq 0$ is the ratio of the end price to the start price on day $t$. Concretely, an investment strategy $a$, at each day $t$, outputs a nonnegative weight vector $a(\cdot|\mathbf{x}^{t-1}) \in \Delta^{m-1}$ over the stocks $[m]$, upon which the investor distributes her wealth accordingly; hereafter, we use $\Delta^{m-1} := \{(\theta_1, \ldots, \theta_m) \in \mathbb{R}_{\geq 0}^m \colon \sum_{i=1}^m \theta_i = 1\}$ to denote the standard $m$-simplex. That is, the multiplicative wealth gain on day $t$ (i.e., the ratio of wealth on day $t$ to the wealth on day $t-1$) is $\sum_{j \in [m]} a(j|\mathbf{x}^{t-1}) x_{tj}$. Thus, her cumulative wealth gain after $n$ days becomes

$$S_n(a, \mathbf{x}^n) := \prod_{t=1}^n \sum_{j \in [m]} a(j|\mathbf{x}^{t-1}) x_{tj} \qquad (1)$$

$$= \sum_{y^n \in [m]^n} \left( \prod_{t=1}^n a(y_t|\mathbf{x}^{t-1}) \right) \mathbf{x}(y^n), \qquad (2)$$

where $\mathbf{x}(y^n) := x_{1y_1} \cdots x_{ny_n}$ denotes the wealth gain of an extreme investment strategy that puts all money to the stock $y_t$ on day $t$, and the second equality follows from the distributive law.

An investor's goal is to design an investment strategy that maximizes her cumulative wealth $S_n(a, \mathbf{x}^n)$. For a stock market where $\mathbf{x}^n$ are i.i.d., it is known that the log-optimal portfolio $\theta^\star$ that maximizes $\mathsf{E}[\log \theta^T \mathbf{X}]$ is asymptotically and competitively optimal. A similar result is well-established for stationary ergodic markets, see, e.g., (Cover and Thomas, 2006, Chapter 16). The log-optimal portfolio theory with stochastic market assumptions, however, is unrealistic, as modeling a stock market could be harder than predicting the market.

As a more realistic alternative, Cover (1991) presented *universal portfolios* that asymptotically achieve the best wealth, to first order in the exponent, attained by a certain class of reference portfolios, with *no statistical assumptions* on the stock market. For the reference class, Cover considered a class of constant rebalanced portfolios (CRPs), where a CRP parameterized by a weight vector $\theta \in \Delta^{m-1}$ is defined to redistribute its wealth according to $\theta$ on every day. Note that CRPs are optimal in an i.i.d. stock market when the distribution is known.

Later, Cover and Ordentlich (1996) extended the theory to a setup where a discrete side information sequence is causally available to an investor; in practice, the side information sequence can be thought to encode an external information that may help predict the stock market. They proposed a variation of Cover (1991)'s universal portfolios that asymptotically achieves the best wealth attained by a class of *state-wise* CRPs that may play different weight vectors according to the side information.

Taking one step further, in this paper, we consider a more challenging scenario in which a side information sequence $z^n \in \mathcal{Z}^n$ is continuous-valued, which could even be the (truncated) market history itself—as a simple motivating example, note that whether or not the price relative of a certain stock was high yesterday may give a hint as to the price relative of that particular stock today. Thus, a ref-

erence portfolio we aim to compete with is parameterized by a state-wise CRP and a *state function* $g\colon \mathcal{Z} \to [S]$ for some $S \geq 2$ and plays the state-wise CRP according to the state sequence $g(z^n) := g(z_1)\ldots g(z_n)$, where we assume a class of state functions $\mathcal{G}$ from which $g$ is drawn; note that larger the $\mathcal{G}$, the richer the reference class.

As the main result, we propose a new investment strategy that asymptotically achieves the same wealth attained by the best state-constant rebalanced portfolios with a state function drawn from a class of functions of finite Natarajan dimension, under a mild regularity condition on the stochasticity of the side information sequence $Z^n$. The proposed strategy is based on a generalization of a universal probability assignment scheme recently proposed by Bhatt and Kim (2021). Note that we assume no transaction costs and that the investor's actions do not affect the market.

The rest of the paper is organized as follows. In Section 2, we review universal portfolios without and with discrete side information, highlighting the connection between universal compression (or probability assignment) and universal portfolios. Section 3 described the proposed algorithm and a crude approximation algorithm for its simulation, together with some concrete examples of side information sequence. We present the proof of the main theorem in Section 4. We conclude with discussing related work in Section 5. All deferred proofs and technical discussions can be found in Appendices.

## 2 A REVIEW OF UNIVERSAL PORTFOLIO THEORY

### 2.1 Universal Portfolios

In his seminal work, Cover (1991) set an ambitious goal that aims to design an investment strategy $b$ to compete with the best strategy in a class $\mathcal{A}$ of investment strategies for any stock market $\mathbf{x}^n$, in the sense that it minimizes the worst-case regret

$$\mathsf{Reg}_n^{\mathsf{port}}(b, \mathcal{A}) := \sup_{\mathbf{x}^n} \sup_{a \in \mathcal{A}} \log \frac{S_n(a, \mathbf{x}^n)}{S_n(b, \mathbf{x}^n)}.$$

We note in passing that by writing the regret as

$$\mathsf{Reg}_n^{\mathsf{port}}(b, \mathcal{A}) = \sup_{\mathbf{x}^n} \left[ \sum_{t=1}^n \log \frac{1}{\mathbf{b}_t^T \mathbf{x}_t} - \inf_{a \in \mathcal{A}} \sum_{t=1}^n \log \frac{1}{\mathbf{a}_t^T \mathbf{x}_t} \right],$$

we can view portfolio selection as online learning with the loss function $\ell_t(\mathbf{b}) = -\log \mathbf{b}^T \mathbf{x}$.

We call a portfolio $b$ *universal with respect to* $\mathcal{A}$ if $\mathsf{Reg}_n^{\mathsf{port}}(b, \mathcal{A}) = o(n)$, i.e., in words, $b$ achieves the same exponential wealth growth rate attained by the best strategy in $\mathcal{A}$ chosen in hindsight with observed market. Remarkably, Cover constructed a universal portfolio with respect to the class of CRPs and established its universality.

Cover's theory is based on the key observation that competing against CRPs in portfolio optimization is equivalent to competing against i.i.d. models in log-loss prediction problem. In what follows, we describe this relationship in a general form beyond between i.i.d. probabilities and CRPs.

For any sequential probability assignment scheme $q(\cdot|y^{t-1}) \in \Delta^{m-1}$ (where $y_i \in [m]$) the *probability induced portfolio* $a = \phi(p)$ is defined as

$$a(j|\mathbf{x}^{t-1}) := \frac{\sum_{y^{t-1} \in [m]^{t-1}} p(y^{t-1}j)\mathbf{x}(y^{t-1})}{\sum_{y^{t-1} \in [m]^{t-1}} p(y^{t-1})\mathbf{x}(y^{t-1})}. \quad (3)$$

Note that if $p$ is an i.i.d. probability, i.e., $p(\cdot|y^{t-1}) = \boldsymbol{\theta} \in \Delta^{m-1}$, it is easy to check from the expression (3) that the corresponding portfolio $\phi(p)$ is the CRP parameterized by $\boldsymbol{\theta}$; thus the class of CRPs $\mathcal{A}^{\mathsf{CRP}}$ is $\phi(\mathcal{P}^{\otimes})$, where we use $\mathcal{P}^{\otimes}$ to denote the class of i.i.d. probabilities.

A peculiar property of a probability induced portfolio $a = \phi(p)$ is that the daily gain can be written as

$$\sum_{y_t \in [m]} a(y_t|\mathbf{x}^{t-1})\mathbf{x}_t(y_t) = \frac{\sum_{y^t} p(y^t)\mathbf{x}(y^t)}{\sum_{y^{t-1}} p(y^{t-1})\mathbf{x}(y^{t-1})},$$

and thus by telescoping, the cumulative wealth gain (1) becomes

$$S_n(\phi(p), \mathbf{x}^n) = \sum_{y^n \in [m]^n} p(y^n)\mathbf{x}(y^n). \quad (4)$$

In view of this expression, a probability induced portfolio can be interpreted as a *fund-of-funds*, i.e., a mixture of the extremal portfolios $(\mathbf{x}(y^n)\colon y^n \in [m]^n)$ with weights $(p(y^n)\colon y^n \in [m]^n)$.

As alluded to earlier, there is an intimate connection between the portfolio optimization with respect to a class of probability induced portfolios and the corresponding log-loss prediction problem. In the log-loss prediction problem, given a class of probabilities $\mathcal{P}$, we define the worst-case regret of a probability $q$ with respect to $\mathcal{P}$ as

$$\mathsf{Reg}_n^{\mathsf{prob}}(q, \mathcal{P}) = \sup_{y^n} \sup_{p \in \mathcal{P}} \log \frac{p(y^n)}{q(y^n)} \quad (5)$$

and call a probability $q$ *universal* with respect to $\mathcal{P}$ if $\mathsf{Reg}_n^{\mathsf{prob}}(q, \mathcal{P}) = o(n)$.

Rather surprisingly, the portfolio selection is equivalent to the log-loss prediction problem for the class of probability induced portfolios.

**Proposition 1.** *For any probability $q$ and any class of probability assignments $\mathcal{P}$, we have*

$$\mathsf{Reg}_n^{\mathsf{port}}(\phi(q), \phi(\mathcal{P})) = \mathsf{Reg}_n^{\mathsf{prob}}(q, \mathcal{P}).$$

*Proof.* Recall (4) that the cumulative wealth of the probability induced portfolio $\phi(p)$ is written as $S_n(\phi(p), \mathbf{x}^n) =$

$\sum_{y^n} p(y^n)\mathbf{x}(y^n)$. First, note that

$$
\begin{aligned}
\mathsf{Reg}_n^{\mathsf{port}}(\phi(q), \phi(\mathcal{P})) &= \sup_{\mathbf{x}^n} \sup_{p \in \mathcal{P}} \log \frac{S_n(\phi(p), \mathbf{x}^n)}{S_n(\phi(q), \mathbf{x}^n)} \\
&\geq \max_{y^n \in [m]^n} \sup_{p \in \mathcal{P}} \log \frac{S_n(\phi(p), \mathbf{e}_{y_1} \ldots \mathbf{e}_{y_n})}{S_n(\phi(q), \mathbf{e}_{y_1} \ldots \mathbf{e}_{y_n})} \\
&\overset{(a)}{=} \max_{y^n \in [m]^n} \sup_{p \in \mathcal{P}} \log \frac{p(y^n)}{q(y^n)} \\
&= \mathsf{Reg}_n^{\mathsf{prob}}(q, \mathcal{P}),
\end{aligned}
$$

where the equality (a) follows since $S_n(\phi(p), \mathbf{e}_{y_1} \ldots \mathbf{e}_{y_n}) = p(y^n)$. Here, $\mathbf{e}_i$ denotes the $i$-th standard unit vector in $\mathbb{R}^m$. To see the converse, note that for any probability $q$, we can write

$$
\begin{aligned}
\mathsf{Reg}_n^{\mathsf{port}}(\phi(q), \phi(\mathcal{P})) &= \sup_{\mathbf{x}^n} \sup_{p \in \mathcal{P}} \log \frac{S_n(\phi(p), \mathbf{x}^n)}{S_n(\phi(q), \mathbf{x}^n)} \\
&= \sup_{\mathbf{x}^n} \sup_{p \in \mathcal{P}} \log \frac{\sum_{y^n} p(y^n)\mathbf{x}(y^n)}{\sum_{y^n} q(y^n)\mathbf{x}(y^n)} \\
&\overset{(b)}{\leq} \sup_{p \in \mathcal{P}} \max_{y^n} \log \frac{p(y^n)}{q(y^n)} \\
&= \mathsf{Reg}_n^{\mathsf{prob}}(q, \mathcal{P}),
\end{aligned}
$$

where (b) follows by Lemma 2 below. $\square$

**Lemma 2** (Cover and Thomas, 2006, Lemma 16.7.1). *Let $a_1, \ldots, a_n, b_1, \ldots, b_n$ be nonnegative real numbers. Then, defining $0/0 = 0$, we have $\frac{\sum_{i=1}^n a_i}{\sum_{i=1}^n b_i} \leq \max_{j \in [n]} \frac{a_j}{b_j}$.*

A direct implication of this statement is that if a probability assignment $q$ is universal with respect to $\mathcal{P}$ for the log-loss prediction problem, then the induced portfolio $\phi(q)$ is universal with respect to $\phi(\mathcal{P})$. If we consider the class of all i.i.d. probabilities $\mathcal{P}^\otimes$, it is well known that the Laplace probability assignment $q_{\mathsf{L}}(y^n) := \int_{\Delta^{m-1}} \mu(\boldsymbol{\theta}) p_{\boldsymbol{\theta}}(y^n) \, d\boldsymbol{\theta}$ is universal for $\mathcal{P}^\otimes$, where $\mu(\boldsymbol{\theta})$ is the uniform density over $\Delta^{m-1}$ and $p_{\boldsymbol{\theta}}(y^n)$ is the i.i.d. probability with parameter $\boldsymbol{\theta} = (\theta_1, \ldots, \theta_m) \in \Delta^{m-1}$, i.e., $p_{\boldsymbol{\theta}}(y^n) := \prod_{i=1}^n \theta_{y_n} = \prod_{j=1}^m \theta_j^{k_j}$ with $k_i = |\{t : y_t = i\}|$. We remark that while the Krichevsky–Trofimov (KT) probability assignment $q_{\mathsf{KT}}$ is universal with an optimal constant in the regret, we consider $q_{\mathsf{L}}$ for simplicity throughout this paper.

Indeed, we have:

**Lemma 3** (Cesa-Bianchi and Lugosi, 2006, Chapter 9).

$$
\sup_{\boldsymbol{\theta} \in \Delta^{m-1}} \sup_{y^n \in [m]^n} \log \frac{p_{\boldsymbol{\theta}}(y^n)}{q_{\mathsf{L}}(y^n)} \leq m \log n.
$$

Hence, $\phi(q_{\mathsf{L}})$ is a universal portfolio for $\mathcal{A}^{\mathsf{CRP}} = \phi(\mathcal{P}^\otimes)$—this is Cover (1991)'s universal portfolio. The universal portfolio $\phi(q_{\mathsf{L}})$ can be expressed as

$$
\phi(q_{\mathsf{L}})(\cdot | \mathbf{x}^{t-1}) = \frac{\int_{\Delta^{m-1}} \boldsymbol{\theta} S_{t-1}(\boldsymbol{\theta}, \mathbf{x}^{t-1}) \mu(\boldsymbol{\theta}) \, d\boldsymbol{\theta}}{\int_{\Delta^{m-1}} S_{t-1}(\boldsymbol{\theta}, \mathbf{x}^{t-1}) \mu(\boldsymbol{\theta}) \, d\boldsymbol{\theta}},
$$

and is thus also known as the $\mu$-weighted portfolio.

## 2.2 Universal Portfolios with Discrete Side Information

Let us now consider a scenario at each time $t$, the investor is additionally given a discrete side information $w_t \in [S]$ for some $S \geq 1$ and chooses a portfolio $a(\cdot | \mathbf{x}^{t-1}; w^t) \in \Delta^{m-1}$, as considered by Cover and Ordentlich (1996). Since the investor's multiplicative wealth gain is $\sum_{y \in [m]} a(y | \mathbf{x}^{t-1}; w^t) \mathbf{x}_t(y)$, similar to the no-side-information setting, the cumulative wealth factor is

$$
S_n(a, \mathbf{x}^n; w^n) := \prod_{t=1}^n \sum_{j \in [m]} a(j | \mathbf{x}^{t-1}; w^t) x_{tj} \quad (6)
$$

and we define the worst-case regret as

$$
\begin{aligned}
\mathsf{Reg}_n^{\mathsf{port}}(b, \mathcal{A}; w^n) &:= \sup_{\mathbf{x}^n} \mathsf{Reg}_n^{\mathsf{port}}(b, \mathcal{A}; \mathbf{x}^n, w^n) \\
&:= \sup_{\mathbf{x}^n} \sup_{a \in \mathcal{A}} \log \frac{S_n(a, \mathbf{x}^n; w^n)}{S_n(b, \mathbf{x}^n; w^n)}
\end{aligned}
$$

for a class $\mathcal{A}$ of portfolios that also adapt to $w^n$. Concretely, as a natural extension of CRPs, we consider a class of state-constant rebalanced portfolios (state-CRPs), denoted as $\mathcal{A}_S^{\mathsf{CRP}}$, where a state-CRP parameterized by a $S$-tuple $(\boldsymbol{\theta}_1, \ldots, \boldsymbol{\theta}_S) \in (\Delta^{m-1})^S$ plays a portfolio $\boldsymbol{\theta}_{w_t}$ at each time $t$.

Paralleling the connection between probability and portfolio in the no-side-information case, we can also define a probability induced portfolio in this setting. In the log-loss prediction with a causal side information sequence, a learner is asked to assign a probability $p(\cdot | y^{t-1}; w^t)$ over $[m]$ based on the causal information, i.e., past sequence $y^{t-1}$ and the side information sequence $w^t$. Here, we use $p(y^n \| w^n) := \prod_{t=1}^n p(y_t | y^{t-1}; w^t)$ to denote the joint probability over $y^n$ given $w^n$. The probability induced portfolio $a = \phi(p)$ is then defined as

$$
a(j | \mathbf{x}^{t-1}; w^t) := \frac{\sum_{y^{t-1}} p(y^{t-1} j \| w^t)\mathbf{x}(y^{t-1})}{\sum_{y^{t-1}} p(y^{t-1} \| w^{t-1})\mathbf{x}(y^{t-1})}, \quad (7)
$$

and as in the no-side information setting, we can write

$$
S_n(\phi(p), \mathbf{x}^n; w^n) = \sum_{y^n} p(y^n \| w^n)\mathbf{x}(y^n). \quad (8)
$$

For example, the class of $S$-state-CRPs $\mathcal{A}_S^{\mathsf{CRP}}$ is induced by the class of all $S$-state i.i.d. probabilities $\mathcal{P}_S^\otimes$, i.e., $\mathcal{A}_S^{\mathsf{CRP}} = \phi(\mathcal{P}_S^\otimes)$. To see this, note that every $S$-state-CRP parameterized by $\boldsymbol{\theta}_{1:S} = (\boldsymbol{\theta}_1, \ldots, \boldsymbol{\theta}_S)$ is the portfolio induced by the state-wise i.i.d. probability assignment $p_{\boldsymbol{\theta}_{1:S}}(y^n \| w^n) := \prod_{t=1}^n p_{\boldsymbol{\theta}_{w_t}}(y_t)$.

Similar to Proposition 1, portfolio optimization with side information with respect to a class of probability induced portfolios is equivalent to the corresponding log-loss prediction problem.

**Proposition 4.** *For any probability assignment $q$ and any class of probability assignment schemes $\mathcal{P}$ with side information sequence $w^n$, we have*

$$\mathsf{Reg}_n^{\mathsf{port}}(\phi(q), \phi(\mathcal{P}); w^n) = \mathsf{Reg}_n^{\mathsf{prob}}(q, \mathcal{P}; w^n),$$

*where we define*

$$\mathsf{Reg}_n^{\mathsf{prob}}(q, \mathcal{P}; w^n) := \sup_{p \in \mathcal{P}} \max_{y^n} \log \frac{p(y^n \| w^n)}{q(y^n \| w^n)}.$$

The proof can be found in Appendix B.1. Note that for the class of $S$-state-wise i.i.d. distributions $\mathcal{P}_S^{\otimes}$, the state-wise extension of the Laplace probability assignment $q_{\mathsf{L};S}$ that assigns

$$q_{\mathsf{L};S}(y^n \| w^n) := \prod_{s=1}^{S} q_{\mathsf{L}}(y^n(s; w^n)), \tag{9}$$

where $y^n(s; w^n) = (y_i : w_i = s, i \in [n])$, is universal, and so $\phi(q_{\mathsf{L};S})$ is universal for $\mathcal{A}_S^{\mathsf{CRP}} = \phi(\mathcal{P}_S^{\otimes})$—this is Cover and Ordentlich (1996)'s universal portfolio.

## 3 MAIN RESULTS

### 3.1 Universal Portfolios with Continuous Side Information

We now consider our main setting where a side information sequence $z^n \in \mathcal{Z}^n$ is continuous-valued. In this setup, for example, one may take $z_t$ as a suffix of the market history $\mathbf{x}_{t-k}^{t-1}$ for some $k \geq 1$. As described earlier in the introduction, we aim to design a universal portfolio that competes against a class of state-CRPs that adapts to the sequence $g(z_1) \dots g(z_n)$, where $g$ is a state function $g \colon \mathcal{Z} \to [S]$, assumed to belong to a class of functions $\mathcal{G}$. Note that a singleton $\mathcal{G} = \{g\}$ recovers the setting of Cover and Ordentlich (1996). Our goal is to design a portfolio that is universal for a largest possible $\mathcal{G}$ with a minimal assumption on the side information sequence. In this paper, we will assume that the *Natarajan dimension* (Shalev-Shwartz and Ben-David, 2014) of $\mathcal{G}$, denoted as $\mathrm{Ndim}(\mathcal{G})$, is finite. The Natarajan dimension can be seen as a generalization of the classic VC dimension, when the function class under consideration is not binary—a formal definition is provided in Appendix A for completeness.

Leveraging the established connection between probability and portfolio, we continue to view the class of state-wise CRPs $\mathcal{A}_S^{\mathsf{CRP}} = \phi(\mathcal{P}_S^{\otimes})$ as the class of portfolios induced by $\mathcal{P}_S^{\otimes}$ and describe the problem in an abstract setting. Since the equivalence between probability assignment and portfolio selection in Proposition 4 holds for *any* side information, we can describe our goal with continuous side information using the same language from Section 2.2 as follows. Given a state function class $\mathcal{G}$ and a reference

class $\mathcal{P}$ of probability assignment schemes with discrete side information, we define a class of probability assignment schemes with continuous side information $\mathcal{P}_{\mathcal{G}}$ as a collection of all probability assignment schemes induced by $\mathcal{P}$ and $\mathcal{G}$, where each probability assignment is parameterized by $p \in \mathcal{P}$ and $g \in \mathcal{G}$ and defined to be

$$p_g(y^n \| z^n) := p(y^n \| g(z^n)).$$

That is, continuous side information $z^n$ is quantized by a function $g$ and then a probability assignment scheme $p$ with discrete side information is deployed.

Precisely, we aim to design a strategy $b$ that achieves a sublinear expected worst-case regret the expected worst-case regret

$$\overline{\mathsf{Reg}}_n^{\mathsf{port}}(b, \phi(\mathcal{P}_{\mathcal{G}})) := \mathsf{E}\big[\mathsf{Reg}_n^{\mathsf{port}}(b, \phi(\mathcal{P}_{\mathcal{G}}); \mathbf{X}^n, Z^n)\big],$$

for a general state function class $\mathcal{G}$ of bounded Natarajan dimension and $\mathcal{P} = \mathcal{P}_S^{\otimes}$, the class of $S$-state i.i.d. probabilities. Here, we assume that the side information sequence $Z^n$ is stochastic with distribution $P_{Z^n}$ which may be arbitrarily correlated with the stock market $\mathbf{X}^n$. This stochastic assumption on $Z^n$ is necessitated to analyze our empirical-covering-based algorithm under the bounded Natarajan dimension condition in Theorem 5. In the individual-sequence assumption on $z^n$, it can be shown that there exists a class of functions with bounded Natarajan dimension, where the equivalent log-loss prediction problem with the function class suffers a linear regret lower bound, which in turn implies that there exists no universal portfolio under the conditions. We discuss the individual-sequence setting in Appendix D in more detail.

From Proposition 4, since the expected regret is always upper bounded by the worst-case regret, a universal portfolio can be readily derived as a plug-in strategy of a universal probability with respect to a continuous side information sequence with an unknown state function. In this paper, we deploy an extended version of the universal probability assignment $q_{\mathcal{G}}^*$ proposed by Bhatt and Kim (2021), which was designed for $m = 2$ and $S = 2$ with regret guarantee established when $y^n$ is random and the side information sequence $Z^n$ is i.i.d.. We will extend their scheme for arbitrary $m$ and $S$ with a guarantee for individual $y^n$ and non-i.i.d. $Z^n$.

**The Proposed Strategy.** Firstly, for any $\tilde{n} \in \mathbb{N}$ and any $\tilde{z}^{\tilde{n}} \in \mathcal{Z}^{\tilde{n}}$, let $\{\tilde{g}_1, \dots, \tilde{g}_\ell\} \subset \mathcal{G}$ be a *minimal empirical covering* of $\mathcal{G}$ with respect to $\tilde{z}^{\tilde{n}}$, i.e., a set of functions such that $\{\tilde{g}_i(\tilde{z}^{\tilde{n}}) : i \in [\ell]\} = \{g(\tilde{z}^{\tilde{n}}) : g \in \mathcal{G}\}$ with the minimum possible size $\ell = \ell(\tilde{z}^{\tilde{n}})$. Then, we define a mixture probability assignment

$$q_{\mathcal{G}; \tilde{z}^{\tilde{n}}}(y^i \| z^i) := \frac{1}{\ell} \sum_{j=1}^{\ell} q_{\mathsf{L};S}(y^i \| \tilde{g}_j(z^i)) \tag{10}$$

Alankrita Bhatt*, J. Jon Ryu*, Young-Han Kim

with respect to the empirical covering, and define the induced sequential probability assignment

$$q_{\mathcal{G};\tilde{z}^{\tilde{n}}}(y_i|y^{i-1};z^i) := \frac{q_{\mathcal{G};\tilde{z}^{\tilde{n}}}(y^i\|z^i)}{q_{\mathcal{G};\tilde{z}^{\tilde{n}}}(y^{i-1}\|z^{i-1})}.$$

The proposed probability assignment $q_{\mathcal{G}}^*$ is then defined as follows. First, we split the $n$ time steps into $\lceil \log_2 n \rceil$ epochs: starting from $j = 1$, define the $j$-th epoch to consist of the time steps $2^{j-1} + 1 \leq i \leq 2^j$. So, the first epoch consists of $z_2$, the second epoch consists of $z_3^4$, the third epoch consists of $z_5^8$ and so on. Then,

- For $i = 1$, $q_{\mathcal{G}}^*(\cdot|z_1) := 1/m$;
- For $i \geq 2$, if $2^{j-1} + 1 \leq i \leq 2^j$, i.e., if the time step $i$ falls within the $j$-th epoch, then $(q_{\mathcal{G};z^{2^{j-1}}}(\emptyset\|\emptyset) := 1)$

$$q_{\mathcal{G}}^*(y_i|y^{i-1};z^i) := \frac{q_{\mathcal{G};z^{2^{j-1}}}(y_{2^{j-1}+1}^i\|z_{2^{j-1}+1}^i)}{q_{\mathcal{G};z^{2^{j-1}}}(y_{2^{j-1}+1}^{i-1}\|z_{2^{j-1}+1}^{i-1})}.$$

Concretely, the probability assigned over $y^n$ given $z^n$ for some $n \in (2^{J-1}, 2^J]$ is

$$\begin{aligned}
q_{\mathcal{G}}^*(y^n\|z^n) &= \prod_{i=1}^n q_{\mathcal{G}}^*(y_i|y^{i-1};z^i) \\
&= q_{\mathcal{G};\emptyset}(y_1\|z_1) q_{\mathcal{G};z^1}(y_2\|z_2) q_{\mathcal{G};z^2}(y_3^4\|z_3^4) \\
&\quad \cdots q_{\mathcal{G};z^{2^{J-1}}}(y_{2^{J-1}+1}^n\|z_{2^{J-1}+1}^n). \quad (11)
\end{aligned}$$

Finally, a sequential portfolio $a = \phi(q_{\mathcal{G}}^*)$ follows from (7).

### 3.2 Universality and Examples

To provide a performance guarantee of the proposed algorithm for a class of $S$-state functions $\mathcal{G}$, we impose a structural condition on the sequence $Z^n \sim P_{Z^n}$ as a stochastic process. For any class of binary functions $\mathcal{H} \subset \{\mathcal{Z} \to \{0, 1\}\}$, define

$$\rho_{\mathcal{H}}(Z^n) := \sup_{h \in \mathcal{H}} \left| \sum_{i=1}^n \big(h(Z_i) - \mathsf{E}[h(Z_i)]\big) \right|, \quad (12)$$

which is a well-studied quantity in the empirical process theory. Specifically, we are interested in the binary function class $\mathbb{1}_{\{\mathcal{G} \times \mathcal{G}\}} := \{h_{g,g'}: \mathcal{Z} \to \{0, 1\}: g, g' \in \mathcal{G}\}$, where $h_{g,g'}(z) := \mathbb{1}_{\{g(z) \neq g'(z)\}}$. With a slight abuse of notation, we use $\rho_{\mathcal{G} \times \mathcal{G}}(Z^n)$ to denote $\rho_{\mathbb{1}_{\{\mathcal{G} \times \mathcal{G}\}}}(Z^n)$. We now state our main result.

**Theorem 5** (Asymptotic universality). *For any collection of functions $\mathcal{G}$ of finite Natarajan dimension and any stationary stochastic process $Z^n$ such that*

$$\mathsf{E}[\rho_{\mathcal{G} \times \mathcal{G}}(Z^n)] = o\Big(\frac{n}{\log^2 n}\Big), \quad (13)$$

*the induced portfolio $\phi(q_{\mathcal{G}}^*)$ satisfies*

$$\lim_{n \to \infty} \frac{1}{n} \overline{\mathsf{Reg}}^{\mathsf{port}}(\phi(q_{\mathcal{G}}^*), \phi((\mathcal{P}_S^{\otimes})_{\mathcal{G}})) = 0.$$

In Theorem 5, the condition $\mathsf{E}[\rho_{\mathcal{G} \times \mathcal{G}}(Z^n)] \ll \frac{n}{\log^2 n}$ on the marginal distribution $P_{Z^n}$ is crucial in ensuring consistency of the portfolio $\phi(q_{\mathcal{G}}^*)$. We now provide a few example cases of side information sequences $Z^n$ where this requirement (13) is satisfied. In fact, by controlling $\mathsf{E}[\rho_{\mathcal{G} \times \mathcal{G}}(Z^n)]$ we can also bound the *nonasymptotic regret* for these particularly interesting cases.

**Example 6** (i.i.d. processes). *When the joint distribution $P_{\mathbf{X}^n, Z^n}$ is such that $Z^n$ is i.i.d., it is well known that $\mathsf{E}[\rho_{\mathcal{H}}(Z^n)] \leq C\sqrt{VCdim(\mathcal{H})n}$ (for absolute constant $C$) for any binary function class $\mathcal{H}$ and any distribution $P_{Z^n}$; see Vershynin (2018, Theorem 8.3.23). Following the same proof, it can be shown that $\mathsf{E}[\rho_{\mathcal{G} \times \mathcal{G}}(Z^n)] \leq C\sqrt{(d \log S)n}$ and consequently $\overline{\mathsf{Reg}}^{\mathsf{port}} = \widetilde{O}(\sqrt{n})$; the only change to be made in the proof is in the growth function—rather than $\left(\frac{en}{d}\right)^d$, the growth function in this case is upper bounded by $(S^2 n)^{2d}$ by Natarajan's Lemma; see Section 4.1.*

**Example 7** ($\beta$-mixing processes). *The quantity $\mathsf{E}[\rho_{\mathcal{H}}(Z^n)]$ has also been studied for classes beyond i.i.d. sequences—in particular, Yu (1994) studied the case when $Z^n$ is a $\beta$-mixing process, which we now define. For the sigma-fields $\sigma_l := \sigma(Z_1, \ldots, Z_\ell)$ and $\sigma'_{l+k} := \sigma(Z_{\ell+k}, Z_{\ell+k+1}, \ldots,)$, we define $\beta_k := \frac{1}{2} \sup\{\mathsf{E} |P(B|\sigma_l) - P(B)|: B \in \sigma'_{\ell+k}, \ell \geq 1\}$. If $\beta_k = O(k^{-r_\beta})$ as $k \to \infty$, we call $r_\beta$ the $\beta$-mixing exponent and call the process $Z^n$ a $\beta$-mixing process with mixing exponent $r_\beta$. Note that a larger $r_\beta$ guarantees faster mixing. We can restate the main result of Yu (1994) for the case when $\mathcal{H}$ has a finite VC dimension. In what follows, $\xrightarrow{p}$ denotes convergence in probability.*

**Theorem 8** (Yu, 1994, Corollary 3.2 and Remark (i)). *Assume that a class of binary functions $\mathcal{H}$ is of finite VC dimension. Let $Z^n$ be a stationary $\beta$-mixing process with mixing exponent $r_\beta \in (0, 1]$. Then, for any given $s \in (0, r_\beta)$, we have*

$$n^{s/(1+s)} \frac{\rho_{\mathcal{H}}(Z^n)}{n} \xrightarrow{p} 0 \quad as \; n \to \infty. \quad (14)$$

*This theorem immediately implies that $\frac{1}{n}\mathsf{Reg}^{\mathsf{port}} \xrightarrow{p} 0$, i.e., $\phi(q_{\mathcal{G}}^*)$ is universal in probability. We can also establish its universality in expectation via Theorem 5, by showing (13) under the same assumption. The proof requires an additional technical argument and thus deferred to Appendix B.2.*

**Example 9** (Market history $z_t = \mathbf{x}_{t-k}^{t-1}$). *A canonical example of side information is the market history $z_t = \mathbf{x}^{t-1}$ or its truncated version with $k$ memory, i.e., $z_t = \mathbf{x}_{t-k}^{t-1}$. In this case, if the stock market $(\mathbf{x}_t)$ itself is $k$-th order Markov, then under an additional mild regularity condition, we can show a faster rate $\overline{\mathsf{Reg}}^{\mathsf{port}} \leq \widetilde{O}(\sqrt{n})$ than implied by the previous example; see Appendix C.*

## 3.3 Implementing Universal Portfolios

Note that computing a portfolio is equivalent to computing the wealth achieved by the portfolio. Recall from (4) that the cumulative wealth achieved by Cover's universal portfolio $\phi(q_{\mathsf{L}})$ can be written as

$$S_n(\phi(q_{\mathsf{L}}), \mathbf{x}^n) = \sum_{y^n \in [m]^n} q_{\mathsf{L}}(y^n) \mathbf{x}(y^n)$$
$$= \int_{\Delta^{m-1}} S_n(\boldsymbol{\theta}, \mathbf{x}^n) \mu(\boldsymbol{\theta}) \, \mathrm{d}\boldsymbol{\theta},$$

since the Laplace probability assignment $q_{\mathsf{L}}(y^n) = \int_{\Delta^{m-1}} \mu(\boldsymbol{\theta}) p_{\boldsymbol{\theta}}(y^n) \, \mathrm{d}\boldsymbol{\theta}$ is a mixture with respect to a uniform density $\mu(\boldsymbol{\theta})$ over the simplex $\Delta^{m-1}$.

It is not hard to see that the *exact* computation of Cover's universal portfolio requires, on the $t$-th day of investment over $m$ stocks, $O(t^{m-1})$ time complexity; see (Cover and Ordentlich, 1996) for a detailed argument. A naive Monte Carlo approximation can be used to approximately estimate the wealth of universal portfolios: if we draw $N$ CRPs $\boldsymbol{\theta}_1, \ldots, \boldsymbol{\theta}_N$ from $\mu$ and *buy-and-hold* uniformly over the CRPs, we will attain approximately similar wealth to $\mu$-universal portfolio, where a larger $N$ leads to better approximation. Note, however, that this naive approximation requires $N = \Omega(\frac{1}{\epsilon^m})$ to achieve an approximation error $\epsilon$. Both exact and approximate computation quickly become infeasible, especially for a large stock market with $m \gg 1$ and/or for a long investment period.

### 3.3.1 Implementing Universal Portfolios with Discrete Side Information

The universal portfolio $\phi(q_{\mathsf{L};S})$ with discrete side information $w^n$ achieves the wealth can be written as

$$S_n(\phi(q_{\mathsf{L};S}), \mathbf{x}^n; w^n) = \sum_{y^n \in [m]^n} q_{\mathsf{L};S}(y^n \| w^n) \mathbf{x}(y^n)$$
$$= \prod_{s=1}^{S} S_{|\mathbf{x}^n(s;w^n)|}(\phi(q_{\mathsf{L}}), \mathbf{x}^n(s; w^n)),$$

where $\mathbf{x}^n(s; w^n) = (\mathbf{x}_i \colon w_i = s, i \in [n])$, since $q_{\mathsf{L};S}(y^n \| w^n) = \prod_{s=1}^{S} q_{\mathsf{L}}(y^n(s; w^n))$. That is, we run Cover's UP for each state separately.

### 3.3.2 Implementing Universal Portfolios with Continuous Side Information

We now describe how to implement the proposed strategy $\phi(q_{\mathcal{G}}^*)$. By the epoch-wise construction of $q_{\mathcal{G}}^*$ as explicitly shown in (11), the cumulative wealth can be factorized as

$$S_n(\phi(q_{\mathcal{G}}^*), \mathbf{x}^n; z^n)$$
$$= \prod_{j=1}^{J} \sum_{y_{2^{j-1}+1}^{2^j}} q_{\mathcal{G}; z^{2^{j-1}}}(y_{2^{j-1}+1}^{2^j} \| z_{2^{j-1}+1}^{2^j}) \mathbf{x}_{2^{j-1}+1}^{2^j}(y_{2^{j-1}+1}^{2^j})$$

$$= \prod_{j=1}^{J} S_{2^{j-1}}(\phi(q_{\mathcal{G}; z^{2^{j-1}}}), \mathbf{x}_{2^{j-1}+1}^{2^j}; z_{2^{j-1}+1}^{2^j}).$$

where we assume $n = 2^J$ for simplicity. Here, for each $j \in [J]$, if $\{\tilde{g}_1, \ldots, \tilde{g}_{\ell_j}\}$ is a minimal empirical covering of $\mathcal{G}$ with respect to $z^{2^{j-1}}$, we can write

$$S_{2^{j-1}}(q_{\mathcal{G}; z^{2^{j-1}}}, \mathbf{x}_{2^{j-1}+1}^{2^j}; z_{2^{j-1}+1}^{2^j})$$
$$= \frac{1}{\ell_j} \sum_{k=1}^{\ell_j} S_{2^{j-1}}(\phi(q_{\mathsf{L};S}), \mathbf{x}_{2^{j-1}+1}^{2^j}; \tilde{g}_k(z_{2^{j-1}+1}^{2^j})).$$

For each state function $\tilde{g}_k$, the summand is the cumulative wealth of the UP with the side information $\tilde{g}_k(z_{2^{j-1}+1}^{2^j})$. We can summarize the algorithm as follows:

---

For each epoch $j = 1, 2, \ldots$:

1. Find an empirical covering $\{\tilde{g}_1, \ldots, \tilde{g}_{\ell_j}\} \subseteq \mathcal{G}$ with respect to $z^{2^{j-1}}$.

2. For each $k \in [\ell_j]$, run UP with the discrete side information $\tilde{g}_k(z_{2^{j-1}+1}^{2^j})$.

3. During the $j$-th investment epoch, i.e., $t \in (2^{j-1}, 2^j]$, run the buy-and-hold strategy uniformly over all UPs with side information $\tilde{g}_k(z_{2^{j-1}+1}^{2^j})$ for each $k \in [\ell_j]$.

4. At the end of the epoch, sell all stocks.

---

### 3.3.3 An Example with Real Data

In the following, we study a simple example for concreteness, which admits an easy construction of minimal empirical coverings. Note that, for a richer class of state functions, finding a minimal empirical covering may be another computational bottleneck.

**Example 10.** *As a simple case of the canonical side information considered in Example 9, we choose the price relative of the stock 1 on the previous day as the continuous side information, i.e., $z_t = x_{t-1,1}$, and a class of 1D threshold functions $\mathcal{G} = \{x \mapsto g_a(x) = \mathbb{1}\{x \geq a\} \colon a > 0\}$ of $\mathrm{Ndim}(\mathcal{G}) = 1$. Note that we consider a binary state space ($S = 2$). In this case, it is easy to show that $\{g_{x_{0,1}}, \ldots, g_{x_{t-1,1}}\}$ is a minimal empirical covering given $z_t = x_{t-1,1}$. More generally, we can consider $z_t = \mathbf{x}_{t-1}$ with a class of product of 1D threshold functions $\mathcal{G} = \{x \mapsto g_{\mathbf{a}}(\mathbf{x}) = (\mathbb{1}\{x_1 \geq a_1\}, \ldots, \mathbb{1}\{x_m \geq a_m\}) \colon \mathbf{a} = (a_1, \ldots, a_m) \in \mathbb{R}_{++}^m\}$ of $\mathrm{Ndim}(\mathcal{G}) \leq m \log m$ (Shalev-Shwartz and Ben-David, 2014, Lemma 29.6) and $S = 2^m$. Given $z_t = \mathbf{x}_{t-1}$, $\{g_{\mathbf{x}_0}, \ldots, g_{\mathbf{x}_{t-1}}\}$ is a minimal empirical covering.*

We briefly demonstrate how the proposed portfolio performs on two real stocks. We collected the 6-year period

from Jan-01-2012 to Dec-31-2017 (total 1508 trading days) of two stocks Ford (F) and Macy's (M). Over the period, Ford went up by a factor of 1.11, while Macy's went down by a factor of 0.77. The best CRP in hindsight, which turns out to be the buy-and-hold of Ford, achieves a growth factor of 1.11. The uniform CRP achieves a growth factor of 0.99. While the universal portfolio without side information achieves a growth factor of only 0.98, the proposed algorithm with the yesterday's prices and the class of thresholding functions achieves a growth factor of 1.15.

## 4 PROOFS

In this section, we prove Theorem 5. Before that, we review the main methodological and technical challenges associated with the problem. Methodologically, the main challenge with having continuous side information is that one does not know the state from just observing the side information and hence cannot construct the Laplace probability assignment directly. As an example, let's say there were two states (namely 1 and 2) that were visible to the investor as $z_t$: if the state is 1 on day $t$, the investor can look at the past days when the side information was 1, construct the corresponding Laplace probability assignment, and subsequently the portfolio strategy. If there are an infinite number of possible mappings of side information into states (represented by the function class $\mathcal{G}$), however, it is not clear how to construct a portfolio strategy. In terms of theoretical challenges, it is clear that some notion of discretization of the infinite class of functions $\mathcal{G}$ is required in order to use a mixture probability assignment: if the distribution of the side information of $z_t$ were known, the problem would be considerably simplified; the main challenge is creating such a cover in a data-dependent way (i.e., using $z_1, \ldots, z_{t-1}$ to construct a cover).

We now proceed to the proof. We proceed by noting that the probability assignment $q_{\mathcal{G}}^*$ used to derive the proposed portfolio guarantees the following regret bound.

**Theorem 11.** *For the probability assignment $q_{\mathcal{G}}^*$, if the Natarajan dimension $\mathrm{Ndim}(\mathcal{G}) = d$ of $\mathcal{G}$ is finite and $Z^n \sim P_{Z^n}$ is stationary, we have*

$$\mathsf{E}\left[\sup_{g \in \mathcal{G}} \sup_{p \in \mathcal{P}_S^{\otimes}} \sup_{y^n \in [m]^n} \log \frac{p(y^n \| g(Z^n))}{q_{\mathcal{G}}^*(y^n \| Z^n)}\right] \qquad (15)$$

$$\leq S(d+m)(\log^2 n) + 2.5Sm \sum_{j=0}^{\log n - 1} j \, \mathsf{E}[\rho_{\mathcal{G} \times \mathcal{G}}(Z^{2^j})].$$

*Here, $\log n$ is assumed to be an integer for simplicity, which can be easily rectified at the cost of an absolute constant factor in the regret; see Section 4.1.*

We will first prove Theorem 11; Theorem 5 then follows as a corollary of Theorem 11 via the established connection between a probability and the induced portfolio in Proposition 4.

### 4.1 Proof of Theorem 11

Note that the key building block of the proposed probability assignment scheme $q_{\mathcal{G}}^*$ is $q_{\tilde{z}^n}(y^i \| z^i)$ defined in (10), the uniform mixture based on a minimal empirical covering of $\mathcal{G}$ with respect to $\tilde{z}^n$. The proof consists of three steps. In Step 1, we first consider the simplest case where the whole side information sequence $z^n$ is provided *noncausally* by an oracle, where we can use $z^n$ as $\tilde{z}^n$ to build the empirical covering. We then analyze the performance of $q_{\tilde{z}^n}(y^i \| z^i)$ for an arbitrary auxiliary sequence $\tilde{z}^n$ in Step 2. Finally, in Step 3, we analyze $q_{\mathcal{G}}^*$ based on the analysis of $q_{\tilde{z}^n}(y^i \| z^i)$.

**Step 1. Side Information Given Noncausally**

Suppose that $z^n$ is available noncausally so that it can be used to construct a minimal empirical covering in $q_{z^n}(y^i \| z^i)$ for $i \in [n]$. First, note that since $|\{(g(z^n) \colon g \in \mathcal{G}\}| \leq S^n$, we can construct an empirical covering $\{g_1, \ldots, g_\ell\}$ of $\mathcal{G}$ with respect to $z^n$ with $\ell \leq S^n$. Assuming $\mathrm{Ndim}(\mathcal{G}) = d < \infty$, however, we can even do so with $\ell \leq (S^2 n)^d$ by Natarajan's Lemma ([Shalev-Shwartz and Ben-David](), 2014, Lemma 29.4). Hence, for the mixture probability assignment $q_{\tilde{z}^n}(y^i \| z^i)$ defined in (10) with $\tilde{z}^n \leftarrow z^n$, i.e., $q_{z^n}(y^i \| z^i) = \frac{1}{\ell} \sum_{j=1}^{\ell} q_{\mathsf{L};S}(y^i \| g_j(z^i))$, it readily follows that for any $g \in \mathcal{G}$,

$$\sup_{p \in P_S^{\otimes}} \sup_{y^n \in [m]^n} \log \frac{p(y^n \| g(z^n))}{q_{z^n}(y^n \| z^n)} \leq d \log(S^2 n) + Sm \log n \qquad (16)$$

by invoking that $\ell \leq (S^2 n)^d$ and applying the regret bound for the $m$-ary Laplace probability assignment in Lemma 3 for each state.

**Step 2. Auxiliary Side Information Given Noncausally**

We now analyze the mixture probability $q_{\tilde{z}^n}(y^n \| z^n)$ for an arbitrary auxiliary sequence $\tilde{z}^n$, possibly being different from $z^n$. Intuitively, the sequence $\tilde{z}^n$ will also reduce the class $\mathcal{G}$ to at most $(S^2 n)^d$ functions, and if $z^n$ and $\tilde{z}^n$ are "not too far apart", the two reductions each obtained by $z^n$ and $\tilde{z}^n$ may be also close. The following lemma provides the performance of the mixture probability $q_{\tilde{z}^n}(y^n \| z^n)$ with respect to the auxiliary sequence $\tilde{z}^n$, capturing the expected gap from the intuition by the Hamming distance (denoted by $\mathrm{d_H}$) between $g(z^n)$ and $\tilde{g}(z^n)$.

**Lemma 12.** *For any $\tilde{z}^n$, $z^n$, and $g \in \mathcal{G}$ with $\mathrm{Ndim}(\mathcal{G}) = d < \infty$, we have*

$$\sup_{p \in \mathcal{P}_S^{\otimes}} \sup_{y^n \in [m]} \log \frac{p(y^n \| g(z^n))}{q_{\tilde{z}^n}(y^n \| z^n)}$$

$$\leq d \log(S^2 n) + Sm(\log n)(1 + 2.5 d_H(g(z^n), \tilde{g}(z^n)))$$

$$\leq S(\log n)(d + m + 2.5m \, \mathrm{d_H}(g(z^n), \tilde{g}(z^n))). \qquad (17)$$

Note that setting $\mathrm{d_H}(g(z^n), \tilde{g}(z^n)) = 0$ recovers (16) as expected.

*Proof.* Let $p_{\boldsymbol{\theta}_{1:S}}$ be a state-wise i.i.d. probability assignment characterized by $\boldsymbol{\theta}_{1:S} = (\boldsymbol{\theta}_1, \ldots, \boldsymbol{\theta}_S) \in (\Delta^{m-1})^S$, where $\boldsymbol{\theta}_i = (\theta_{i1}, \theta_{i2}, \ldots, \theta_{im}) \in \Delta^{m-1}$ for each $i \in [S]$. For any state function $g \in \mathcal{G}$, by definition of the empirical covering, there exists a function $\tilde{g} \in \{\tilde{g}_1, \ldots, \tilde{g}_\ell\}$ such that $\tilde{g}(\tilde{z}^n) = g(\tilde{z}^n)$. Hence, we first have

$$\log \frac{p_{\boldsymbol{\theta}_{1:S}}(y^n \| g(z^n))}{q_{\tilde{z}^n}(y^n \| z^n)}$$
$$\leq d \log(S^2 n) + \log \frac{p_{\boldsymbol{\theta}_{1:S}}(y^n \| g(z^n))}{q_{\mathsf{L};S}(y^n \| \tilde{g}(z^n))}. \quad (18)$$

It only remains to analyze $q_{\mathsf{L};S}(y^n \| \tilde{g}(z^n))$. For each $i \in [S]$ and $j \in [m]$, we define $n_i := |t: g(Z_t) = i|$ and $k_{ij} := |t: g(Z_t) = i, y_t = j|$. Moreover let $\tilde{n}_i, \tilde{k}_{ij}$ be defined in a similar way as $\tilde{n}_i := |t: \tilde{g}(Z_t) = i|$ and $\tilde{k}_{ij} := |t: \tilde{g}(Z_t) = i, y_t = j|$). We can then write $p_{\boldsymbol{\theta}_{1:S}}(y^n \| g(z^n)) = \prod_{s=1}^S \theta_{s1}^{k_{s1}} \cdots \theta_{sm}^{k_{sm}}$.

Further, we can explicitly write the expression for the Laplace probability assignment as $q_{\mathsf{L}}(y^n) = (\binom{n+m-1}{m-1}\binom{n}{k_1, \ldots, k_m})^{-1}$, where $k_i = |\{t: y_t = i\}|$, and thus its state-wise extension as

$$q_{\mathsf{L};S}(y^n \| \tilde{g}(z^n))$$
$$= \Big( \prod_{s=1}^S \binom{\tilde{n}_s + m - 1}{m - 1} \binom{\tilde{n}_s}{\tilde{k}_{s1}, \ldots, \tilde{k}_{s,m-1}} \Big)^{-1}.$$

Now, consider

$$\log \frac{p_{\boldsymbol{\theta}_{1:S}}(y^n \| g(z^n))}{q_{\mathsf{L};S}(y^n \| \tilde{g}(\tilde{z}^n))}$$
$$= \sum_{i=1}^S \log \binom{\tilde{n}_i + m - 1}{m - 1} \binom{\tilde{n}_i}{\tilde{k}_{i1}, \ldots, \tilde{k}_{i,m-1}} \theta_{i1}^{k_{i1}} \cdots \theta_{im}^{k_{im}}$$
$$\leq Sm \log n + \sum_{i=1}^S \log \frac{\binom{\tilde{n}_i}{\tilde{k}_{i1}, \ldots, \tilde{k}_{i,m-1}}}{\binom{n_i}{k_{i1}, \ldots, k_{i,m-1}}} \quad (19)$$
$$= Sm \log n + \sum_{i=1}^S \log \frac{\tilde{n}_i!}{n_i!} + \sum_{i=1}^S \sum_{j=1}^m \log \frac{k_{ij}!}{\tilde{k}_{ij}}, \quad (20)$$

where (19) follows since $\binom{n_i}{k_{i1}, \ldots, k_{i,m-1}} \theta_{i1}^{k_{i1}} \cdots \theta_{im}^{k_{im}} \leq 1$.

Now, since for all $i \in [S]$ and $j \in [m]$, we have $|n_i - \tilde{n}_i| \leq d_{\mathrm{H}}(g(z^n), \tilde{g}(z^n))$ and $|k_{ij} - \tilde{k}_{ij}| \leq d_{\mathrm{H}}(g(z^n), \tilde{g}(z^n))$, we have that $\tilde{n}_i \leq n_i + d_{\mathrm{H}}(g(z^n), \tilde{g}(z^n))$ and consequently $\frac{\tilde{n}_i!}{n_i!} \leq \frac{(n_i + d_{\mathrm{H}}(g(z^n), \tilde{g}(z^n)))!}{n_i!}$. Thus, we can invoke the exact same calculations as in (Bhatt and Kim, 2021, Propositions 5 and 6) to bound the second and third terms in (20) as

$$\log \frac{p_{\boldsymbol{\theta}_{1:S}}(y^n \| g(z^n))}{q_{\mathsf{L};S}(y^n \| \tilde{g}(\tilde{z}^n))}$$
$$\leq Sm \log n + S(m+3) \, d_{\mathrm{H}}(g(z^n), \tilde{g}(z^n)) \log n$$
$$\leq Sm(\log n)(1 + 2.5 d_H(g(z^n), \tilde{g}(z^n))), \quad (21)$$

since $m \geq 2$. Plugging this into (18) establishes the first bound. The second bound follows by observing $\log(S^2 n) \leq S \log n$. $\square$

When $Z^n$ is stationary as a stochastic process and if $\tilde{Z}^n$ is a statistical copy of $Z^n$, the following lemma shows that the Hamming distance can be bounded by $\rho_{\mathcal{G} \times \mathcal{G}}(Z^n)$, which can be controlled in expectation as $o(n / \log^2 n)$ under mild regularity conditions on $P_{Z^n}$ and $\mathcal{G}$. The proof is deferred to Appendix B.3.

**Lemma 13.** *If $Z^n$ is stationary, $\tilde{Z}^n \overset{(d)}{=} Z^n$, and $\tilde{g}(\tilde{Z}^n) = \tilde{g}(\tilde{Z}^n)$, then*

$$d_{\mathrm{H}}(g(Z^n), \tilde{g}(Z^n)) \leq \rho_{\mathcal{G} \times \mathcal{G}}(Z^n) + \rho_{\mathcal{G} \times \mathcal{G}}(\tilde{Z}^n).$$

**Step 3. Side Information Given Causally**

In view of Lemma 13, provided that $Z^n$ is stationary, we can *bootstrap* the history sequence to construct such an auxiliary sequence, which motivates the epoch-based construction of $q_{\mathcal{G}}^*$. That is, we split the $n$ time steps into $\log n$ epochs. For simplicity, we assume that $\log n$ is an integer; if not, we may "extend" the horizon of the game from $n$ to $2^{\lceil \log n \rceil} < 2n$, and follow the same analysis incurring at most a constant factor extra in the regret bound. Defining the $j$-the epoch to consist of the time steps $2^{j-1} + 1 \leq i \leq 2^j$ starting from $j = 1$, while we define $q_{\mathcal{G}}^*(\cdot | Z_1) = 1/m$ for the 0-th epoch. For $i \geq 2$, if the time step $i$ falls within the $j$-th epoch, i.e., $2^{j-1} + 1 \leq i \leq 2^j$, then

$$q_{\mathcal{G}}^*(y_i | y^{i-1}; Z^i) = \frac{q_{Z^{2^{j-1}}}(y_{2^{j-1}+1}^i \| Z_{2^{j-1}+1}^i)}{q_{Z^{2^{j-1}}}(y_{2^{j-1}+1}^{i-1} \| Z_{2^{j-1}+1}^{i-1})} \quad (22)$$

where we can recall the definition of $q_{Z^{2^{j-1}}}$ from (10). For any $p \in \mathcal{P}_S^{\otimes}$, we then have

$$\sum_{i=1}^n \log \frac{p(y_i | g(Z_i))}{q_{\mathcal{G}}^*(y_i | y^{i-1}; Z^i)}$$
$$\leq \sum_{i=2}^n \log \frac{p(y_i | g(Z_i))}{q_{\mathcal{G}}^*(y_i | y^{i-1}; Z^i)} + \log m$$
$$= \sum_{j=1}^{\log n} \sum_{i=2^{j-1}+1}^{2^j} \log \frac{p(y_i | g(Z_i))}{q_{\mathcal{G}}^*(y_i | y^{i-1}; Z^i)}$$
$$= \sum_{j=1}^{\log n} \log \frac{p(y_{2^{j-1}+1}^{2^j} \| g(Z_{2^{j-1}+1}^{2^j}))}{q_{Z^{2^{j-1}}}(y_{2^{j-1}+1}^{2^j} \| Z_{2^{j-1}+1}^{2^j})} \quad (23)$$
$$\leq S(d + m)(\log^2 n)$$
$$+ 2.5 Sm \sum_{j=0}^{\log n - 1} j d_H(g(Z_1^{2^j}), g(Z_{2^j+1}^{2^{j+1}})), \quad (24)$$

where (23) follows by (22) and (24) follows from Lemma 12. Finally, taking supremum over $y^n, p$ and $g$ and expectation over $Z^n$ leads to the desired inequality by Lemma 13. $\square$

Alankrita Bhatt*, J. Jon Ryu*, Young-Han Kim

## 4.2 Proof of Theorem 5

By Proposition 4 and Theorem 11, we have

$$\overline{\mathsf{Reg}}^{\mathsf{port}}(\phi(q_{\mathcal{G}}^*), \mathcal{A}_S^{\mathsf{CRP}}, \mathcal{G})$$

$$= \mathsf{E}[\mathsf{Reg}_n^{\mathsf{port}}(\phi(q_{\mathcal{G}}^*), \mathcal{A}_S^{\mathsf{CRP}}, \mathcal{G}; \mathbf{X}^n, Z^n)]$$

$$\leq \mathsf{E}[\mathsf{Reg}_n^{\mathsf{prob}}(q; \mathcal{P}, \mathcal{G}; Z^n)]$$

$$= \mathsf{E}\left[\sup_{g \in \mathcal{G}} \sup_{p \in \mathcal{P}} \max_{y^n} \log \frac{p(y^n \| g(Z^n))}{q(y^n \| Z^n)}\right]$$

$$\leq S(d+m)(\log^2 n) + 2.5Sm \sum_{j=0}^{\log n - 1} j \, \mathsf{E}[\rho(Z^{2^j})],$$

where we omit the subscript in $\rho_{\mathcal{G} \times \mathcal{G}}(\cdot)$ for brevity. Since the first term in the bound is sublinear in $n$ when $d$ and $S$ are fixed, it then suffices to show that $\sum_{j=0}^{\log n - 1} j \, \mathsf{E}[\rho(Z^{2^j})] = o(n)$. Using the change of variables $n' = \log n$, observe

$$\sum_{j=0}^{\log n - 1} j \, \mathsf{E}[\rho(Z^{2^j})] = \frac{1}{n'} \sum_{j=0}^{n'-1} j \, \mathsf{E}[\rho(Z^{2^j})] \frac{n'}{2^{n'}}$$

$$\leq \frac{1}{n'} \sum_{j=0}^{n'-1} \frac{j^2}{2^j} \, \mathsf{E}[\rho(Z^{2^j})],$$

where the inequality follows since $\frac{n'}{2^{n'}} \leq \frac{j}{2^j}$ for all $j \leq n'$. Now, since $\frac{(\log n)^2}{n} \mathsf{E}[\rho(Z^n)] = \frac{n'^2}{2^{n'}} \mathsf{E}[\rho(Z^{2^{n'}})] \to 0$ as $n \to \infty$ is assumed, we also have $\frac{1}{n'} \sum_{j=0}^{n'-1} \frac{j^2}{2^j} \mathsf{E}[\rho(Z^{2^j})] \to 0$ as $n' \to \infty$, by the Cesàro mean Theorem. A final change of variables concludes. □

## 5 RELATED WORK AND DISCUSSION

Portfolio selection has been extensively studied in information theory since the seminal work of Cover (1991) and Cover and Ordentlich (1996), both of which established close connections between portfolio selection and the classically studied information theoretic problem of universal compression (Rissanen, 1996; Ziv and Lempel, 1978; Merhav and Feder, 1998; Xie and Barron, 2000). A number of variations have been considered since Helmbold et al. (1998). For example, incorporating transaction costs (Blum and Kalai, 1999; Uziel and El-Yaniv, 2020) using other probability assignments than i.i.d. (Kozat et al., 2008; Tavory and Feder, 2010), and considering space complexity issues (Tavory and Feder, 2008). Cross and Barron (2003) and Györfi et al. (2006) proposed portfolio selection techniques incorporating continuous side information; however, the competitor classes considered in both are disparate from ours.

As demonstrated, portfolio selection with side information is closely related to sequential prediction with side information and log-loss. This problem has attracted recent interest (Rakhlin and Sridharan, 2015; Bilodeau et al., 2020;

Fogel and Feder, 2017; Bhatt and Kim, 2021), with the first two focused on obtaining fundamental limits via the sequential complexities approach of Rakhlin et al. (2015a). More recently, the preprint of Bilodeau et al. (2021) proposed a mixture-based conditional density estimator, which specifically achieves $\mathsf{E}[\mathsf{Reg}^{\mathsf{prob}}] = O(\log^2 n)$ for the binary probability assignment problem with i.i.d. side information with a VC class, which tightens the regret $\tilde{O}(\sqrt{n})$ established in (Bhatt and Kim, 2021). Therefore, it is natural to consider applying the probability assignment of Bilodeau et al. (2021) in hoping to relax the technical condition (13) and establish Theorem 5 for *all* stationary ergodic $Z^n$—it is known that $\mathsf{E}[\rho_{\mathcal{G} \times \mathcal{G}}(Z^n)] = o(n)$ for any stationary ergodic process; see e.g., (Adams and Nobel, 2010). We note, however, that analyzing their method in our setting of non-i.i.d. side information sequences seems to involve a significant amount of additional work. More precisely, their analysis needs to be extended to (1) individual-sequence $y^n$ and (2) stationary ergodic side information with a dependence of the regret on $\rho_{\mathcal{H}}(Z^n)$ similar to that of the method of Bhatt and Kim (2021). In their words, we would need to relax the assumption of the data being *well-specified*. At a high level, they use a similar covering approach (with respect to the Hellinger metric over distributions) as well as a smoothing of probabilities in order to avoid unbounded likelihood ratios (we, in contrast, have used the Laplace/KT probability assignment). Using a similar epoch-based analysis they establish regret bounds in (Bilodeau et al., 2021, Appendix D) by first upper bounding the KL divergence in terms of the Hellinger divergence and then leveraging local Rademacher complexities in conjunction with an inequality of Bousquet (2002). In order to extend their method to individual-sequence $y^n$ and stationary ergodic $Z^n$, one would need to either extend the aforementioned inequality to these cases, or to bypass the step of upper-bounding the KL divergence in terms of the Hellinger divergence altogether. We leave these directions for future work.

The stochasticity assumption on $Z^n$ cannot be completely removed for a state function class with bounded Natarajan dimension. We show in Appendix D that there exists no universal portfolio for a class of 1-dimensional (1D) threshold functions. This is a consequence of the fact that the 1D threshold has infinite *Littlestone dimension* (Ldim), which is a combinatorial dimension that characterizes whether classification is possible in (sequential and adversarial) online learning. Therefore, it seems that we either have to consider function classes with finite Ldim in order to achieve vanishing regret with fully adversarial side information; or consider stationary side information along with function classes with finite VC dimension (VCdim). Since Ldim $\geq$ VCdim in general with several examples where the gap is infinite, both settings are complementary with neither one subsuming the other. We thus leave an in-depth study of the setting with adversarial side information to future work.

## References

Adams, T. M. and Nobel, A. B. (2010). Uniform convergence of Vapnik–Chervonenkis classes under ergodic sampling. *Ann. Probab.*, pages 1345–1367.

Bertail, P. and Portier, F. (2019). Rademacher complexity for Markov chains: Applications to kernel smoothing and Metropolis–Hastings. *Bernoulli*, 25(4B):3912–3938.

Bhatt, A. and Kim, Y.-H. (2021). Sequential prediction under log-loss with side information. In *Algo. Learn. Theory*, pages 340–344. PMLR.

Bilodeau, B., Foster, D., and Roy, D. (2020). Tight bounds on minimax regret under logarithmic loss via self-concordance. In *Proc. Int. Conf. Mach. Learn.*, pages 919–929. PMLR.

Bilodeau, B., Foster, D. J., and Roy, D. M. (2021). Minimax rates for conditional density estimation via empirical entropy. *arXiv preprint arXiv:2109.10461*.

Blum, A. and Kalai, A. (1999). Universal portfolios with and without transaction costs. *Mach. Learn.*, 35(3):193–205.

Bousquet, O. (2002). *Concentration inequalities and empirical processes theory applied to the analysis of learning algorithms*. PhD thesis, École Polytechnique: Department of Applied Mathematics Paris, France.

Cesa-Bianchi, N. and Lugosi, G. (2006). *Prediction, learning, and games*. Cambridge University Press.

Cover, T. M. (1991). Universal portfolios. *Math. Financ.*, 1(1):1–29.

Cover, T. M. and Ordentlich, E. (1996). Universal portfolios with side information. *IEEE Trans. Inf. Theory*, 42(2):348–363.

Cover, T. M. and Thomas, J. A. (2006). *Elements of information theory*. John Wiley & Sons.

Cross, J. E. and Barron, A. R. (2003). Efficient universal portfolios for past-dependent target classes. *Math. Financ.*, 13(2):245–276.

Fogel, Y. and Feder, M. (2017). On the problem of online learning with log-loss. In *Proc. IEEE Int. Symp. Inf. Theory*, pages 2995–2999. IEEE.

Györfi, L., Lugosi, G., and Udina, F. (2006). Nonparametric kernel-based sequential investment strategies. *Math. Financ.*, 16(2):337–357.

Haghtalab, N., Han, Y., Shetty, A., and Yang, K. (2022a). Oracle-efficient online learning for beyond worst-case adversaries. *arXiv preprint arXiv:2202.08549*.

Haghtalab, N., Roughgarden, T., and Shetty, A. (2022b). Smoothed analysis with adaptive adversaries. In *Proc. IEEE Ann. Symp. Found. Comput. Sci.*, pages 942–953. IEEE.

Hanneke, S. and Yang, L. (2019). Statistical learning under nonstationary mixing processes. In *Int. Conf. Artif. Int. Statist.*, pages 1678–1686. PMLR.

Helmbold, D. P., Schapire, R. E., Singer, Y., and Warmuth, M. K. (1998). On-line portfolio selection using multiplicative updates. *Math. Financ.*, 8(4):325–347.

Karandikar, R. L. and Vidyasagar, M. (2002). Rates of uniform convergence of empirical means with mixing processes. *Stat. Probab. Lett.*, 58(3):297–307.

Kozat, S. S., Singer, A. C., and Bean, A. J. (2008). Universal portfolios via context trees. In *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, pages 2093–2096. IEEE.

Merhav, N. and Feder, M. (1998). Universal prediction. *IEEE Trans. Inf. Theory*, 44(6):2124–2147.

Rakhlin, A. and Sridharan, K. (2014). Statistical learning theory and sequential prediction. Lecture Notes at University of Pennsyvania.

Rakhlin, A. and Sridharan, K. (2015). Sequential probability assignment with binary alphabets and large classes of experts. *arXiv preprint arXiv:1501.07340*.

Rakhlin, A., Sridharan, K., and Tewari, A. (2015a). Online learning via sequential complexities. *J. Mach. Learn. Res.*, 16(1):155–186.

Rakhlin, A., Sridharan, K., and Tewari, A. (2015b). Sequential complexities and uniform martingale laws of large numbers. *Probab. Theory Relat. Fields*, 161(1-2):111–153.

Rissanen, J. J. (1996). Fisher information and stochastic complexity. *IEEE Trans. Inf. Theory*, 42(1):40–47.

Shalev-Shwartz, S. and Ben-David, S. (2014). *Understanding machine learning: From theory to algorithms*. Cambridge university press.

Tavory, A. and Feder, M. (2008). Finite memory universal portfolios. In *Proc. IEEE Int. Symp. Inf. Theory*, pages 1408–1412. IEEE.

Tavory, A. and Feder, M. (2010). Universal portfolio algorithms in realistic-outcome markets. In *Proc. IEEE Inf. Theory Workshop*, pages 1–5. IEEE.

Uziel, G. and El-Yaniv, R. (2020). Long-and short-term forecasting for portfolio selection with transaction costs. In *Int. Conf. Artif. Int. Statist.*, pages 100–110. PMLR.

Vershynin, R. (2018). *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press.

Xie, Q. and Barron, A. R. (2000). Asymptotic minimax regret for data compression, gambling, and prediction. *IEEE Trans. Inf. Theory*, 46(2):431–445.

Yu, B. (1994). Rates of convergence for empirical processes of stationary mixing sequences. *Ann. Probab.*, pages 94–116.

Ziv, J. and Lempel, A. (1978). Compression of individual sequences via variable-rate coding. *IEEE Trans. Inf. Theory*, 24(5):530–536.

# A DEFINITION OF NATARAJAN DIMENSION

We use the definitions from Shalev-Shwartz and Ben-David (2014, Definitions 29.1, 29.2).

**Definition 14** (Shattering). *Let $\mathcal{G} \subset \{\mathcal{Z} \to [S]\}$. Then, a set $C \subset \mathcal{Z}$ is said to be* shattered *by the function class $\mathcal{G}$ if there exist two functions $g_0, g_1 \in \mathcal{G}$ such that*

- *For each $z \in C, g_0(z) \neq g_1(z)$, and*

- *For each $B \subset C$ there exists a function $g \in \mathcal{G}$ such that*

$$\forall z \in B, g(x) = g_0(x) \text{ and } \forall z \in C \setminus B, g(x) = g_1(x).$$

We can now define the Natarajan dimension.

**Definition 15** (Natarajan dimension). *For any function class $\mathcal{G} \subset \{\mathcal{Z} \to [S]\}$ the* Natarajan dimension *of $\mathcal{G}$ is the maximal size of a shattered set $C \subset \mathcal{Z}$.*

# B DEFERRED PROOFS

## B.1 Proof of Proposition 4

In this proof, we assume that the side information sequence $z^n$ may take values from an arbitrary alphabet $\mathcal{Z}$. Recall that for a probability assignment $q(y_i | y^{i-1}; z^i)$, we can write the wealth of the portfolio induced by $q$ as

$$S_n(\phi(q), \mathbf{x}^n; z^n) = \sum_{y^n \in [m]^n} q(y^n \| z^n) \mathbf{x}(y^n). \tag{8}$$

To see this, recall that the probability induced portfolio $a = \phi(q)$ is defined as

$$a(j | \mathbf{x}^{t-1}; z^t) := \frac{\sum_{y^{t-1}} q(y^{t-1} j \| z^t) \mathbf{x}(y^{t-1})}{\sum_{y^{t-1}} q(y^{t-1} \| z^{t-1}) \mathbf{x}(y^{t-1})},$$

where recall for $t \in [n], q(y^t \| z^t) = \prod_{i=1}^t q(y_i | y^{i-1}; z^i)$. We then have

$$\sum_{y_t \in [m]} a(y_t | \mathbf{x}^{t-1}; z^t) \mathbf{x}_t(y_t) = \frac{\sum_{y^{t-1}} q(y^t \| z^t) \mathbf{x}(y^t)}{\sum_{y^{t-1}} q(y^{t-1} \| z^{t-1}) \mathbf{x}(y^{t-1})},$$

and (8) consequently follows from telescoping.

Now, we note that

$$
\begin{aligned}
\mathsf{Reg}_n^{\mathsf{port}}(\phi(q); \phi(\mathcal{P}); z^n) &= \sup_{p \in \mathcal{P}} \sup_{\mathbf{x}^n} \log \frac{S_n(\phi(p), \mathbf{x}^n; z^n)}{S_n(\phi(q), \mathbf{x}^n; z^n)} \\
&\geq \sup_{p \in \mathcal{P}} \sup_{y^n \in [m]^n} \log \frac{S_n(\phi(p), \mathbf{e}_{y_1} \ldots \mathbf{e}_{y_n}; z^n)}{S_n(\phi(q), \mathbf{e}_{y_1} \ldots \mathbf{e}_{y_n}; z^n)} \\
&= \sup_{p \in \mathcal{P}} \sup_{y^n \in [m]^n} \log \frac{p(y^n \| z^n)}{q(y^n \| z^n)} \\
&= \mathsf{Reg}_n^{\mathsf{prob}}(q; \mathcal{P}; z^n). \tag{25}
\end{aligned}
$$

Conversely, we can also see that

$$
\begin{aligned}
\mathsf{Reg}_n^{\mathsf{port}}(\phi(q); \phi(\mathcal{P}); \mathbf{x}^n, z^n) &= \sup_{p \in \mathcal{P}} \log \frac{S_n(\phi(p), \mathbf{x}^n; z^n)}{S_n(\phi(q), \mathbf{x}^n; z^n)} \\
&= \sup_{p \in \mathcal{P}} \log \frac{\sum_{y^n \in [m]^n} p(y^n \| z^n) \mathbf{x}(y^n)}{\sum_{y^n \in [m]^n} q(y^n \| z^n) \mathbf{x}(y^n)} \\
&\leq \sup_{p \in \mathcal{P}} \max_{y^n \in [m]^n} \log \frac{p(y^n \| z^n)}{q(y^n \| z^n)} \\
&= \mathsf{Reg}_n^{\mathsf{prob}}(q; \mathcal{P}; z^n),
\end{aligned}
\tag{26}
$$

where (26) follows from Lemma 2. The desired inequality follows since this inequality holds for any $\mathbf{x}^n$. $\qquad\square$

### B.2 Proof of Universality in Expectation in Example 7

Recall that by Theorem 5, it suffices to show that

$$
\mathsf{E}[\rho_{\mathcal{G} \times \mathcal{G}}(Z^n)] = o\Big(\frac{n}{\log^2 n}\Big)
\tag{13}
$$

to establish that the induced portfolio $\phi(q_{\mathcal{G}}^*)$ is universal in expectation. Indeed, for a $\beta$-mixing process $Z^n$ with $\beta$-mixing coefficient $\beta_k$ and $\beta$-mixing exponent $r > 0$, i.e., $\beta_k = O(k^{-r})$ as $k \to \infty$, we can prove a stronger statement:

$$
\mathsf{E}[\rho_{\mathcal{G} \times \mathcal{G}}(Z^n)] = O(n^{(3+r)/(3+2r)}).
\tag{27}
$$

The argument below to show (27) is based on the techniques of Karandikar and Vidyasagar (2002) and Hanneke and Yang (2019).

Pick $k \geq 1$ which divides $n$ for simplicity; the divisibility can be easily lifted by elongating the game from $n$ steps to the next number divisible by $k$. We will choose $k$ as a function of $n$ at the end of proof. We define the nonoverlapping $k$ subsequences $Z^{(1)}, \ldots, Z^{(k)}$ of length $n/k$ as

$$
\begin{aligned}
Z^{(1)}{}_1^{n/k} &= Z_1, Z_{k+1}, Z_{2k+1} \ldots, Z_{(n/k-1)+1}, \\
Z^{(2)}{}_1^{n/k} &= Z_2, Z_{k+2}, Z_{2k+2} \ldots, Z_{(n/k-1)k+2}, \\
&\;\;\vdots \\
Z^{(k)}{}_1^{n/k} &= Z_k, Z_{2k}, Z_{3k} \ldots, Z_{(n/k)k}.
\end{aligned}
$$

We will invoke the classical result on $\beta$-mixing processes that states that

$$
d_{\mathsf{TV}}\Big( P_{Z^{(j)}{}_1^{n/k}}, \prod_{i=1}^{n/k} P_{Z_i^{(j)}} \Big) \leq \Big(\frac{n}{k} - 1\Big) \beta_k
\tag{28}
$$

for each $j \in [k]$, where $d_{\mathsf{TV}}(\cdot, \cdot)$ denotes the total variation distance; see, for example, (Hanneke and Yang, 2019, Lemma 1) and the references therein.

Now, we consider

$$
\begin{aligned}
\mathsf{E}[\rho_{\mathcal{H}}(Z^n)] &= \mathsf{E}\Big[ \sup_{h \in \mathcal{H}} \Big| \sum_{i=1}^{n} (h(Z_i) - E[h(Z_i)]) \Big| \Big] \\
&\leq \mathsf{E}\Big[ \sum_{j=1}^{k} \sup_{h \in \mathcal{H}} \Big| \sum_{i=1}^{n/k} (h(Z_i^{(j)}) - \mathsf{E}[h(Z_i^{(j)})]) \Big| \Big] \\
&= \sum_{j=1}^{k} \mathsf{E}\Big[ \sup_{h \in \mathcal{H}} \Big| \sum_{i=1}^{n/k} (h(Z_i^{(j)}) - \mathsf{E}[h(Z_i^{(j)})]) \Big| \Big].
\end{aligned}
\tag{29}
$$

Let $Z_1', \ldots, Z_{n/k}'$ be an i.i.d. process with the same marginal distribution of the stationary process $Z^n$, i.e., $P_{Z_1'} = P_{Z_1}$. Continuing from the summand in (29), we then have

$$\mathsf{E}\left[\sup_{h\in\mathcal{H}}\left|\sum_{i=1}^{n/k}(h(Z_i^{(1)}) - \mathsf{E}[h(Z_i^{(1)})])\right|\right] = \mathsf{E}\left[\sup_{h\in\mathcal{H}}\left|\sum_{i=1}^{n/k}(h(Z_i^{(1)}) - h(Z_i') + h(Z_i') - \mathsf{E}[h(Z_i^{(1)})])\right|\right]$$

$$= \mathsf{E}\left[\sup_{h\in\mathcal{H}}\left|\sum_{i=1}^{n/k}(h(Z_i^{(1)}) - h(Z_i') + h(Z_i') - \mathsf{E}[h(Z_i')])\right|\right] \tag{30}$$

$$\leq \mathsf{E}\left[\sup_{h\in\mathcal{H}}\left|\sum_{i=1}^{n/k}(h(Z_i^{(1)}) - h(Z_i'))\right|\right] + \mathsf{E}\left[\sup_{h\in\mathcal{H}}\left|\sum_{i=1}^{n/k}(h(Z_i') - E[h(Z_i')])\right|\right] \tag{31}$$

$$\leq \mathsf{E}\left[\sup_{h\in\mathcal{H}}\left|\sum_{i=1}^{n/k}(h(Z_i^{(1)}) - h(Z_i'))\right|\right] + C\sqrt{\frac{dn}{k}} \tag{32}$$

$$\leq \frac{n}{k}\sup_{h\in\mathcal{H}}\mathsf{E}\left|\frac{k}{n}\sum_{i=1}^{n/k}h(Z_i^{(1)}) - \frac{k}{n}\sum_{i=1}^{n/k}h(Z_i')\right| + C\sqrt{\frac{dn}{k}}$$

$$\leq \frac{n}{k}d_{\mathsf{TV}}\left(P_{{Z^{(j)}}_1^{n/k}}, \prod_{i=1}^{n/k}P_{Z_i^{(j)}}\right) + C\sqrt{\frac{dn}{k}} \tag{33}$$

$$\leq \frac{n^2\beta_k}{k^2} + C\sqrt{\frac{dn}{k}}. \tag{34}$$

Here, (30) follows since the marginal distribution $Z_i' \stackrel{(d)}{=} Z_i^{(1)}$, (32) follows since the distribution $Z'^n$ is i.i.d. and from (Vershynin, 2018, Theorem 8.3.23), (33) follows from the following variational form of the total variation distance $d_{\mathsf{TV}}(P, P')$ between two measures $P$ and $P'$ defined over the same measure space, i.e.,

$$d_{\mathsf{TV}}(P, P') = \sup_{f:|f|\leq 1} |\mathsf{E}_{X\sim P}[f(X)] - \mathsf{E}_{X\sim P'}[f(X)]|,$$

and lastly (34) follows from (28). Substituting (34) into (29) yields that

$$\mathsf{E}[\rho_{\mathcal{H}}(Z^n)] \leq \frac{n^2\beta_k}{k} + C\sqrt{dnk} \leq \frac{C'n^2k^{-r}}{k} + C\sqrt{dnk}$$

for $k$ sufficiently large with some $C' > 0$, where we use the definition of the $\beta$-mixing exponent $r$ in the second inequality. Finally, choosing $k = O(n^{\frac{3}{3+2r}})$ yields the claimed rate $\mathsf{E}[\rho_{\mathcal{H}}(Z^n)] = O(n^{\frac{3+r}{3+2r}})$. $\square$

## B.3 Proof of Lemma 13

Note that for any $Z^n$ and $\tilde{Z}^n$, we can write

$$d_{\mathrm{H}}(g(Z^n), \tilde{g}(Z^n)) = d_{\mathrm{H}}(g(Z^n), \tilde{g}(Z^n)) - d_{\mathrm{H}}(g(\tilde{Z}^n), \tilde{g}(\tilde{Z}^n)) \tag{35}$$

$$\leq \sup_{g_1,g_2}\left|d_{\mathrm{H}}(g_1(Z^n), g_2(Z^n)) - d_{\mathrm{H}}(g_1(\tilde{Z}^n), g_2(\tilde{Z}^n))\right|$$

$$\leq \sup_{g_1,g_2}|d_{\mathrm{H}}(g_1(Z^n), g_2(Z^n)) - n\mathsf{P}(g_1(Z_1) \neq g_2(Z_2))|$$

$$+ \sup_{g_1,g_2}\left|d_{\mathrm{H}}(g_1(\tilde{Z}^n), g_2(\tilde{Z}^n)) - n\mathsf{P}(g_1(\tilde{Z}_1) \neq g_2(\tilde{Z}_1))\right|$$

$$= \rho_{\mathcal{G}\times\mathcal{G}}(z^n) + \rho_{\mathcal{G}\times\mathcal{G}}(\tilde{z}^n) \tag{36}$$

where (35) follows since $d_{\mathrm{H}}(g(\tilde{Z}^n), \tilde{g}(\tilde{Z}^n)) = 0$ by design and (36) follows since by stationarity of $Z^n \stackrel{(d)}{=} \tilde{Z}^n$, we have $n\mathsf{P}(g_1(Z_1) \neq g_2(Z_1)) = n\mathsf{P}(g_1(\tilde{Z}_1) \neq g_2(\tilde{Z}_1)) = \sum_{i=1}^n \mathsf{P}(g_1(\tilde{Z}_i) \neq g_2(\tilde{Z}_i)) = \sum_{i=1}^n \mathsf{E}[\mathbb{1}_{\{g_1(\tilde{Z}_i)\neq g_2(\tilde{Z}_i)\}}]$. Finally, substituting (36) into (21) yields the lemma. $\square$

## C   A DETAILED DISCUSSION ON EXAMPLE 9

For the side information $z_t = \mathbf{x}_{t-k}^{t-1}$ in Example 9, if the market $(\mathbf{X}_t)$ itself is $k$-th order Markov, then we can establish the following guarantee.

**Lemma 16.** *Let $\mathbf{X}^n$ be a stationary $k$-th order Markov process and let $Z_t = \mathbf{X}_{t-k}^{t-1} \in (\mathbb{R}_+^m)^k$. Suppose that (1) the density of $Z_0 = \mathbf{X}_{-k}^{-1}$ exists and is bounded and supported over a bounded, convex set $E \subset (\mathbb{R}_+^m)^k$ with nonempty interior and (2) there exist $b > 0$ and $\epsilon > 0$ such that the time-invariant conditional density satisfies*

$$p_{\mathbf{X}_{t-k+1}^t | \mathbf{X}_{t-k}^{t-1}}(z'|z) \geq b \mathbb{1}_{B(z,\epsilon)}(z')$$

*for any $z \in (\mathbb{R}_+^m)^k$, where $B(z,\epsilon)$ denotes the open ball of radius $\epsilon$ centered at $z \in (\mathbb{R}_+^m)^k$ with respect to Euclidean distance. Then, we have $\mathsf{E}[\rho_{\mathcal{H}}(Z^n)] = \widetilde{O}(\sqrt{n})$.*

*Proof.* This is a direct consequence of (Bertail and Portier, 2019, Proposition 11), which establishes an upper bound on $\mathsf{E}[\rho_{\mathcal{H}}(Z'^n)]$ for a *Metropolis–Hastings* (MH) walk $Z'^n$. First, note that $Z^n$ forms a Markov chain due to the $k$-th order Markovity of $\mathbf{X}^n$. To apply the proposition over the Markov chain $Z^n$, we set the proposal distribution $q$ in the MH algorithm to be the actual transition kernel of the Markov chain $Z^n$, so that the MH walk becomes the process $Z^n$ of our interest. Then, under the assumptions above, we can apply the result of Bertail and Portier (2019) and conclude that $\mathsf{E}[\rho_{\mathcal{H}}(Z^n)] = \tilde{O}(\sqrt{n})$ for a VC-class $\mathcal{H}$. $\qquad\square$

## D   UNIVERSAL PORTFOLIO DOES NOT ALWAYS EXIST FOR ADVERSARIAL CONTINUOUS SIDE INFORMATION

In this paper, we studied a universal portfolio optimization with continuous side information, under a certain stochastic assumption on the market sequence $\mathbf{X}^n$ and the side information sequence $Z^n$. A natural question is whether such stochasticity assumptions on the continuous side information sequence $Z^n$ can be further relaxed or completely removed, as was assumed by Cover and Ordentlich (1996) for the discrete side information case. In this section, we provide a simple counterexample with a binary function class having bounded VC dimension, where no universal portfolio exists for an adversarial side information sequence $z^n$.

Under the adversarial continuous side information assumption, we can equivalently consider the corresponding log-loss prediction problem as shown in Proposition 4. We will consider the simple binary-state and binary-alphabet case, i.e., $S = 2$ and $m = 2$. As a reference class $\mathcal{P}$ of probability assignment schemes with side information, we consider state-wise i.i.d. probabilities, where each element $p_{\boldsymbol{\theta}}$ can be parameterized by a pair $\boldsymbol{\theta} = (\theta_0, \theta_1) \in [0,1]^2$ and

$$p_{\boldsymbol{\theta}}(1|w) := \theta_w$$

for $w \in \{0,1\}$. Assume that side information $z_t$ takes a value from $\mathcal{Z} = [0,1]$. As a state function class $\mathcal{G}$, we consider the 1-dimensional (1D) threshold function class

$$\mathcal{G} = \{g_a \mapsto \mathbb{1}_{(a,\infty)}(z) \colon a \in [0,1]\}.$$

Hence, in this setting, we wish to compete with any reference probability assignment with continuous side information that has the form

$$p_{\boldsymbol{\theta},a}(1|w) = \theta_{g_a(w)},$$

for some $a \in [0,1]$ and $\boldsymbol{\theta} \in [0,1]^2$.

Now, in the corresponding sequential prediction with side information, the adversary first picks $\boldsymbol{\theta} = (\theta_0, \theta_1)$ and $a \in [0,1]$. At each time step, the adversary chooses $z_t \in [0,1]$ and present to the player as side information. The player then must assign probability assignment $q(\cdot|z^t, y^{t-1})$ based on all the observations so far, i.e., $z^t$ and $y^{t-1}$. Afterwards, the adversary reveals $y_t \in \{0,1\}$ and suffers a loss of $-\log p(y_t|g(z_t))$, while the player suffers a loss of $-\log q(y_t|z^t, y^{t-1})$. The worst-case regret of the strategy $q$ employed by the player is thus written as

$$\mathsf{Reg}^{\mathsf{prob}}(q, \mathcal{P}_{\mathcal{G}}) := \sup_{y^n \in [m]^n} \sup_{z^n} \sup_{a \in [0,1]} \sup_{\boldsymbol{\theta} \in [0,1]^2} \left\{ \sum_{t=1}^n \log \frac{1}{q(y_t|z^t, y^{t-1})} - \sum_{t=1}^n \log \frac{1}{p_{\boldsymbol{\theta},a}(y_t|z_t)} \right\}.$$
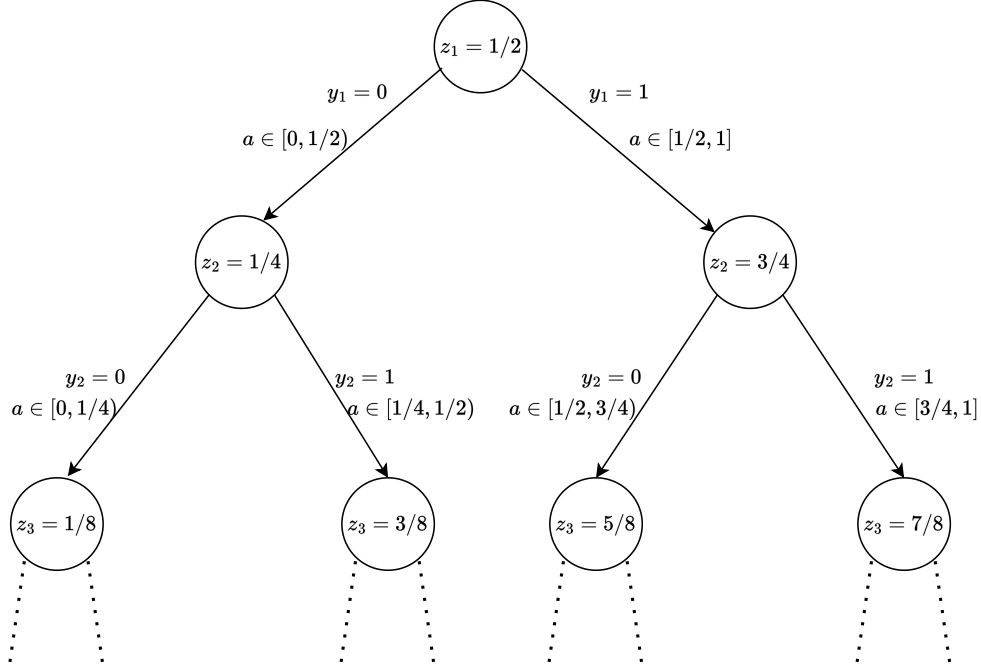
Figure 1: A tree representing the adversary's strategy. At time step $t$, the adversary presents $y_t = 0$ if $q(1|z^t, y^{t-1}) > 1/2$ and $y_t = 1$ otherwise, and then presents $z_{t+1}$ following the path in the tree. The values of $a$ at each branch of the tree show the values satisfying $g_a(z_i) = y_i$ for $i \in \{1, \ldots, t\}$.

We now construct a sequence $z^n$ and $y^n$ as well as choices of $a$ and $\boldsymbol{\theta}$ that ensure the player to suffer a cumulative regret of $n$ regardless of the player's strategy $q$, by mimicking the adversarial construction used to show that the 1D threshold class is not learnable under the 0-1 loss; see (Shalev-Shwartz and Ben-David, 2014, Chapter 21). First, set $z_1 = 1/2$. Since the adversary knows the player's strategy, the adversary picks $y_1 = 0$ if if $q(1|z_1) > 1/2$ and $y_1 = 1$ otherwise. With this choice, the player will suffer loss of at least $\log(1/(1/2)) = \log 2 = 1$ at time step 1. In the next round ($t = 2$), if $y_1 = 0$, choose $z_2 = 1/4$ and $z_2 = 3/4$ otherwise. Similar to the previous time step, the adversary sets $y_2 = 0$ if $q(1|z_1^2, y_1) > 1/2$ and $y_2 = 1$ otherwise. This again ensures the player to incur loss at least 1 at $t = 2$. We can visualize the strategy employed by the adversary as a binary tree as shown in Figure D. Once we proceed in the following manner and reach the final time step, the loss of the player has added up to $n$.

We now only need to specify $\boldsymbol{\theta}$ and $a$. Note that for the sequence $\{(z_t, y_t)\}_{t=1}^n$ constructed above, there exists a threshold $a^* \in [0, 1]$ such that $g_{a^*}(z_t) = y_t$ for all $t$. By setting $\boldsymbol{\theta}^* = (0, 1)$, the player assigns $p_{\boldsymbol{\theta}^*, a^*}(y_t|z_t) = 1$ at each time step, and thus incurs zero cumulative loss. Therefore, the cumulative regret of the player, $\mathsf{Reg}^{\mathsf{prob}}(q, \mathcal{P}_G) \geq n$, regardless of $q$ for the above outlined strategy of the adversary.

The problem of studying fundamental limits on the regret in sequential prediction with the log-loss and side information, with the state function classes possibly being even more general than the VC class we have considered in this work, has previously been studying using the notion of *sequential complexity measures* proposed by Rakhlin et al. (2015b). The current state of the art result in this direction is due to Bilodeau et al. (2020), who establish a regret bound of $\widetilde{O}(\mathfrak{R}_n^{2/3}(\mathcal{G})n^{1/3})$ where $\mathfrak{R}_n(\mathcal{G})$ is the *sequential Rademacher complexity* of the function class $\mathcal{G}$. However, for a class $\mathcal{G}$ with finite VC dimension, it is not necessarily true that $\mathfrak{R}_n(\mathcal{G}) = o(n)$. For example, for the 1D threshold class discussed in the previous example, we can show that the sequential fat-shattering dimension $\mathsf{fat}_\alpha(\mathcal{G}) = \infty$; see (Rakhlin et al., 2015b, Definition 7). This is because for $0 < \alpha \leq 1$ $\mathsf{fat}_\alpha(\mathcal{G})$ coincides with the *Littlestone dimension* of $\mathcal{G}$, which is well known to be infinite (Shalev-Shwartz and Ben-David, 2014, Chapter 21). Using (Rakhlin et al., 2015b, Lemma 7), we see that $\mathfrak{R}_n(\mathcal{G}) \geq \Omega(n)$ thereby implying that the regret bound of Bilodeau et al. (2020) is vacuous in this particular case as one may expect; see also (Rakhlin and Sridharan, 2014, Chapter 8) for a more detailed discussion on learning thresholds.

The above discussion establishes that there are two settings one can consider for this problem: that of completely adversarial side information and a function class with bounded Littlestone dimension; or that of stochastic side information and a function class with bounded VC dimension. Both settings are complementary and neither subsumes the other, since the

Littlestone dimension of a function class is always greater than the VC dimension, with the gap possibly being infinite as in the 1D threshold class. In this paper we pursue the latter direction, in the spirit of recent investigations Haghtalab et al. (2022b,a) into the *beyond worst case analysis* in online learning.