
Learning Theory and Algorithms for Forecasting Non-Stationary Time Series

Vitaly Kuznetsov
Courant Institute
New York, NY 10011
vitaly@cims.nyu.edu

Mehryar Mohri
Courant Institute and Google Research
New York, NY 10011
mohri@cims.nyu.edu

Abstract

We present data-dependent learning bounds for the general scenario of non-stationary non-mixing stochastic processes. Our learning guarantees are expressed in terms of a data-dependent measure of sequential complexity and a discrepancy measure that can be estimated from data under some mild assumptions. We use our learning bounds to devise new algorithms for non-stationary time series forecasting for which we report some preliminary experimental results.

1 Introduction

Time series forecasting plays a crucial role in a number of domains ranging from weather forecasting and earthquake prediction to applications in economics and finance. The classical statistical approaches to time series analysis are based on generative models such as the autoregressive moving average (ARMA) models, or their integrated versions (ARIMA) and several other extensions [Engle, 1982, Bollerslev, 1986, Brockwell and Davis, 1986, Box and Jenkins, 1990, Hamilton, 1994]. Most of these models rely on strong assumptions about the noise terms, often assumed to be i.i.d. random variables sampled from a Gaussian distribution, and the guarantees provided in their support are only asymptotic.

An alternative non-parametric approach to time series analysis consists of extending the standard i.i.d. statistical learning theory framework to that of stochastic processes. In much of this work, the process is assumed to be stationary and suitably mixing [Doukhan, 1994]. Early work along this approach consisted of the VC-dimension bounds for binary classification given by Yu [1994] under the assumption of stationarity and β -mixing. Under the same assumptions, Meir [2000] presented bounds in terms of covering numbers for regression losses and Mohri and Rostamizadeh [2009] proved general data-dependent Rademacher complexity learning bounds. Vidyasagar [1997] showed that PAC learning algorithms in the i.i.d. setting preserve their PAC learning property in the β -mixing stationary scenario. A similar result was proven by Shalizi and Kontorovitch [2013] for mixtures of β -mixing processes and by Berti and Rigo [1997] and Pestov [2010] for exchangeable random variables. Alquier and Wintenberger [2010] and Alquier et al. [2014] also established PAC-Bayesian learning guarantees under weak dependence and stationarity.

A number of algorithm-dependent bounds have also been derived for the stationary mixing setting. Lozano et al. [2006] studied the convergence of regularized boosting. Mohri and Rostamizadeh [2010] gave data-dependent generalization bounds for stable algorithms for φ -mixing and β -mixing stationary processes. Steinwart and Christmann [2009] proved fast learning rates for regularized algorithms with α -mixing stationary sequences and Modha and Masry [1998] gave guarantees for certain classes of models under the same assumptions.

However, stationarity and mixing are often not valid assumptions. For example, even for Markov chains, which are among the most widely used types of stochastic processes in applications, stationarity does not hold unless the Markov chain is started with an equilibrium distribution. Similarly,

long memory models such as ARFIMA, may not be mixing or mixing may be arbitrarily slow [Baillie, 1996]. In fact, it is possible to construct first order autoregressive processes that are not mixing [Andrews, 1983]. Additionally, the mixing assumption is defined only in terms of the distribution of the underlying stochastic process and ignores the loss function and the hypothesis set used. This suggests that mixing may not be the right property to characterize learning in the setting of stochastic processes.

A number of attempts have been made to relax the assumptions of stationarity and mixing. Adams and Nobel [2010] proved asymptotic guarantees for stationary ergodic sequences. Agarwal and Duchi [2013] gave generalization bounds for asymptotically stationary (mixing) processes in the case of stable on-line learning algorithms. Kuznetsov and Mohri [2014] established learning guarantees for fully non-stationary β - and φ -mixing processes.

In this paper, we consider the general case of non-stationary non-mixing processes. We are not aware of any prior work providing generalization bounds in this setting. In fact, our bounds appear to be novel even when the process is stationary (but not mixing). The learning guarantees that we present hold for both bounded and unbounded memory models. Deriving generalization bounds for unbounded memory models even in the stationary mixing case was an open question prior to our work [Meir, 2000]. Our guarantees cover the majority of approaches used in practice, including various autoregressive and state space models.

The key ingredients of our generalization bounds are a data-dependent measure of sequential complexity (*expected sequential covering number* or *sequential Rademacher complexity* [Rakhlin et al., 2010]) and a measure of *discrepancy* between the sample and target distributions. Kuznetsov and Mohri [2014] also give generalization bounds in terms of discrepancy. However, unlike the result of Kuznetsov and Mohri [2014], our analysis does not require any mixing assumptions which are hard to verify in practice. More importantly, under some additional mild assumption, the discrepancy measure that we propose can be estimated from data, which leads to data-dependent learning guarantees for non-stationary non-mixing case.

We devise new algorithms for non-stationary time series forecasting that benefit from our data-dependent guarantees. The parameters of generative models such as ARIMA are typically estimated via the maximum likelihood technique, which often leads to non-convex optimization problems. In contrast, our objective is convex and leads to an optimization problem with a unique global solution that can be found efficiently. Another issue with standard generative models is that they address non-stationarity in the data via a *differencing* transformation which does not always lead to a stationary process. In contrast, we address the problem of non-stationarity in a principled way using our learning guarantees.

The rest of this paper is organized as follows. The formal definition of the time series forecasting learning scenario as well as that of several key concepts is given in Section 2. In Section 3, we introduce and prove our new generalization bounds. In Section 4, we give data-dependent learning bounds based on the empirical discrepancy. These results, combined with a novel analysis of kernel-based hypotheses for time series forecasting (Appendix B), are used to devise new forecasting algorithms in Section 5. In Appendix C, we report the results of preliminary experiments using these algorithms.

2 Preliminaries

We consider the following general time series prediction setting where the learner receives a realization $(X_1, Y_1), \dots, (X_T, Y_T)$ of some stochastic process, with $(X_t, Y_t) \in \mathcal{Z} = \mathcal{X} \times \mathcal{Y}$. The objective of the learner is to select out of a specified family H a hypothesis $h: \mathcal{X} \rightarrow \mathcal{Y}$ that achieves a small generalization error $\mathbb{E}[L(h(X_{T+1}), Y_{T+1}) | Z_1, \dots, Z_T]$ conditioned on observed data, where $L: \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$ is a given loss function. The *path-dependent* generalization error that we consider in this work is a finer measure of the generalization ability than the *averaged* generalization error $\mathbb{E}[L(h(X_{T+1}), Y_{T+1})] = \mathbb{E}[\mathbb{E}[L(h(X_{T+1}), Y_{T+1}) | Z_1, \dots, Z_T]]$ since it only takes into consideration the realized history of the stochastic process and does not average over the set of all possible histories. The results that we present in this paper also apply to the setting where the time parameter t can take non-integer values and prediction lag is an arbitrary number $l \geq 0$. That is, the error is defined by $\mathbb{E}[L(h(X_{T+l}), Y_{T+l}) | Z_1, \dots, Z_T]$ but for notational simplicity we set $l = 1$.

Our setup covers a larger number of scenarios commonly used in practice. The case $\mathcal{X} = \mathcal{Y}^p$ corresponds to a large class of autoregressive models. Taking $\mathcal{X} = \cup_{p=1}^{\infty} \mathcal{Y}^p$ leads to growing memory models which, in particular, include state space models. More generally, \mathcal{X} may contain both the history of the process $\{Y_t\}$ and some additional side information.

To simplify the notation, in the rest of the paper, we will use the shorter notation $f(z) = L(h(x), y)$, for any $z = (x, y) \in \mathcal{Z}$ and introduce the family $\mathcal{F} = \{(x, y) \rightarrow L(h(x), y) : h \in H\}$ containing such functions f . We will assume a bounded loss function, that is $|f| \leq M$ for all $f \in \mathcal{F}$ for some $M \in \mathbb{R}_+$. Finally, we will use the shorthand \mathbf{Z}_a^b to denote a sequence of random variables Z_a, Z_{a+1}, \dots, Z_b .

The key quantity of interest in the analysis of generalization is the following supremum of the empirical process defined as follows:

$$\Phi(\mathbf{Z}_1^T) = \sup_{f \in \mathcal{F}} \left(\mathbb{E}[f(Z_{T+1}) | \mathbf{Z}_1^T] - \sum_{t=1}^T q_t f(Z_t) \right), \quad (1)$$

where q_1, \dots, q_T are real numbers, which in the standard learning scenarios are chosen to be uniform. In our general setting, different Z_t s may follow different distributions, thus distinct weights could be assigned to the errors made on different sample points depending on their relevance to forecasting the future Z_{T+1} . The generalization bounds that we present below are for an arbitrary sequence $\mathbf{q} = (q_1, \dots, q_T)$ which, in particular, covers the case of uniform weights. Remarkably, our bounds do not even require the non-negativity of \mathbf{q} .

Our generalization bounds are expressed in terms of data-dependent measures of sequential complexity such as expected sequential covering number or sequential Rademacher complexity [Rakhlin et al., 2010]. We give a brief overview of the notion of sequential covering number and refer the reader to the aforementioned reference for further details. We adopt the following definition of a complete binary tree: a \mathcal{Z} -valued complete binary tree \mathbf{z} is a sequence (z_1, \dots, z_T) of T mappings $z_t : \{\pm 1\}^{t-1} \rightarrow \mathcal{Z}, t \in [1, T]$. A path in the tree is $\sigma = (\sigma_1, \dots, \sigma_{T-1})$. To simplify the notation we will write $z_t(\sigma)$ instead of $z_t(\sigma_1, \dots, \sigma_{t-1})$, even though z_t depends only on the first $t-1$ elements of σ . The following definition generalizes the classical notion of covering numbers to sequential setting. A set V of \mathbb{R} -valued trees of depth T is a *sequential α -cover* (with respect to \mathbf{q} -weighted ℓ_p norm) of a function class \mathcal{G} on a tree \mathbf{z} of depth T if for all $g \in \mathcal{G}$ and all $\sigma \in \{\pm 1\}^T$, there is $\mathbf{v} \in V$ such that

$$\left(\sum_{t=1}^T |\mathbf{v}_t(\sigma) - g(\mathbf{z}_t(\sigma))|^p \right)^{\frac{1}{p}} \leq \|\mathbf{q}\|_q^{-1} \alpha,$$

where $\|\cdot\|_q$ is the dual norm. The *(sequential) covering number* $\mathcal{N}_p(\alpha, \mathcal{G}, \mathbf{z})$ of a function class \mathcal{G} on a given tree \mathbf{z} is defined to be the size of the minimal sequential cover. The *maximal covering number* is then taken to be $\mathcal{N}_p(\alpha, \mathcal{G}) = \sup_{\mathbf{z}} \mathcal{N}_p(\alpha, \mathcal{G}, \mathbf{z})$. One can check that in the case of uniform weights this definition coincides with the standard definition of sequential covering numbers. Note that this is a purely combinatorial notion of complexity which ignores the distribution of the process in the given learning problem.

Data-dependent sequential covering numbers can be defined as follows. Given a stochastic process distributed according to the distribution \mathbf{p} with $\mathbf{p}_t(\cdot | \mathbf{z}_1^{t-1})$ denoting the conditional distribution at time t , we sample a $\mathcal{Z} \times \mathcal{Z}$ -valued tree of depth T according to the following procedure. Draw two independent samples Z_1, Z'_1 from \mathbf{p}_1 : in the left child of the root draw Z_2, Z'_2 according to $\mathbf{p}_2(\cdot | Z_1)$ and in the right child according to $\mathbf{p}_2(\cdot | Z'_1)$. More generally, for a node that can be reached by a path $(\sigma_1, \dots, \sigma_t)$, we draw Z_t, Z'_t according to $\mathbf{p}_t(\cdot | S_1(\sigma_1), \dots, S_{t-1}(\sigma_{t-1}))$, where $S_t(1) = Z_t$ and $S_t(-1) = Z'_t$. Let \mathbf{z} denote the tree formed using Z_t s and define the *expected covering number* to be $\mathbb{E}_{\mathbf{z} \sim T(\mathbf{p})}[\mathcal{N}_p(\alpha, \mathcal{G}, \mathbf{z})]$, where $T(\mathbf{p})$ denotes the distribution of \mathbf{z} .

In a similar manner, one can define other measures of complexity such as sequential Rademacher complexity and the Littlestone dimension [Rakhlin et al., 2015] as well as their data-dependent counterparts [Rakhlin et al., 2011].

The final ingredient needed for expressing our learning guarantees is the notion of *discrepancy* between target distribution and the distribution of the sample:

$$\Delta = \sup_{f \in \mathcal{F}} \left(\mathbb{E}[f(Z_{T+1}) | \mathbf{Z}_1^T] - \sum_{t=1}^T q_t \mathbb{E}[f(Z_t) | \mathbf{Z}_1^{t-1}] \right). \quad (2)$$

The discrepancy Δ is a natural measure of the non-stationarity of the stochastic process \mathbf{Z} with respect to both the loss function L and the hypothesis set H . In particular, note that if the process \mathbf{Z} is i.i.d., then we simply have $\Delta = 0$ provided that q_t s form a probability distribution. It is also possible to give bounds on Δ in terms of other natural distances between distribution. For instance, Pinsker's inequality yields

$$\Delta \leq M \left\| \mathbf{P}_{T+1}(\cdot | \mathbf{Z}_1^T) - \sum_{t=1}^T q_t \mathbf{P}_t(\cdot | \mathbf{Z}_1^{t-1}) \right\|_{\text{TV}} \leq \sqrt{\frac{1}{2} D \left(\mathbf{P}_{T+1}(\cdot | \mathbf{Z}_1^T) \parallel \sum_{t=1}^T q_t \mathbf{P}_t(\cdot | \mathbf{Z}_1^{t-1}) \right)},$$

where $\|\cdot\|_{\text{TV}}$ is the total variation distance and $D(\cdot \parallel \cdot)$ the relative entropy, $\mathbf{P}_{t+1}(\cdot | \mathbf{Z}_1^t)$ the conditional distribution of Z_{t+1} , and $\sum_{t=1}^T q_t \mathbf{P}_t(\cdot | \mathbf{Z}_1^{t-1})$ the mixture of the sample marginals. Alternatively, if the target distribution at lag l , $\mathbf{P} = \mathbf{P}_{T+l}$ is a stationary distribution of an asymptotically stationary process \mathbf{Z} [Agarwal and Duchi, 2013, Kuznetsov and Mohri, 2014], then for $q_t = 1/T$ we have

$$\Delta \leq \frac{M}{T} \sum_{t=1}^T \|\mathbf{P} - \mathbf{P}_{t+l}(\cdot | \mathbf{Z}_{-\infty}^t)\|_{\text{TV}} \leq \phi(l),$$

where $\phi(l) = \sup_s \sup_{\mathbf{z}} [\|\mathbf{P} - \mathbf{P}_{l+s}(\cdot | \mathbf{z}_{-\infty}^s)\|_{\text{TV}}]$ is the coefficient of asymptotic stationarity. The process is asymptotically stationary if $\lim_{l \rightarrow \infty} \phi(l) = 0$. However, the most important property of the discrepancy Δ is that, as shown later in Section 4, it can be estimated from data under some additional mild assumptions. [Kuznetsov and Mohri, 2014] also give generalization bounds for non-stationary mixing processes in terms of a related notion of discrepancy. It is not known if the discrepancy measure used in [Kuznetsov and Mohri, 2014] can be estimated from data.

3 Generalization Bounds

In this section, we prove new generalization bounds for forecasting non-stationary time series. The first step consists of using *decoupled tangent* sequences to establish concentration results for the supremum of the empirical process $\Phi(\mathbf{Z}_1^T)$. Given a sequence of random variables \mathbf{Z}_1^T we say that \mathbf{Z}'_1^T is a decoupled tangent sequence if Z'_t is distributed according to $\mathbb{P}(\cdot | \mathbf{Z}_1^{t-1})$ and is independent of \mathbf{Z}_t^∞ . It is always possible to construct such a sequence of random variables [De la Peña and Giné, 1999]. The next theorem is the main result of this section.

Theorem 1. *Let \mathbf{Z}_1^T be a sequence of random variables distributed according to \mathbf{p} . Fix $\epsilon > 2\alpha > 0$. Then, the following holds:*

$$\mathbb{P}(\Phi(\mathbf{Z}_1^T) - \Delta \geq \epsilon) \leq \mathbb{E}_{\mathbf{v} \sim T(\mathbf{p})} [\mathcal{N}_1(\alpha, \mathcal{F}, \mathbf{v})] \exp \left(-\frac{(\epsilon - 2\alpha)^2}{2M^2 \|\mathbf{q}\|_2^2} \right).$$

Proof. The first step is to observe that, since the difference of the suprema is upper bounded by the supremum of the difference, it suffices to bound the probability of the following event

$$\left\{ \sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T q_t (\mathbb{E}[f(Z_t) | \mathbf{Z}_1^{t-1}] - f(Z_t)) \right) \geq \epsilon \right\}.$$

By Markov's inequality, for any $\lambda > 0$, the following inequality holds:

$$\begin{aligned} & \mathbb{P} \left(\sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T q_t (\mathbb{E}[f(Z_t) | \mathbf{Z}_1^{t-1}] - f(Z_t)) \right) \geq \epsilon \right) \\ & \leq \exp(-\lambda \epsilon) \mathbb{E} \left[\exp \left(\lambda \sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T q_t (\mathbb{E}[f(Z_t) | \mathbf{Z}_1^{t-1}] - f(Z_t)) \right) \right) \right]. \end{aligned}$$

Since \mathbf{Z}'_1^T is a tangent sequence the following equalities hold: $\mathbb{E}[f(Z_t)|\mathbf{Z}_1^{t-1}] = \mathbb{E}[f(Z'_t)|\mathbf{Z}_1^{t-1}] = \mathbb{E}[f(Z'_t)|\mathbf{Z}_1^T]$. Using these equalities and Jensen's inequality, we obtain the following:

$$\begin{aligned} & \mathbb{E} \left[\exp \left(\lambda \sup_{f \in \mathcal{F}} \sum_{t=1}^T q_t (\mathbb{E}[f(Z_t)|\mathbf{Z}_1^{t-1}] - f(Z_t)) \right) \right] \\ &= \mathbb{E} \left[\exp \left(\lambda \sup_{f \in \mathcal{F}} \mathbb{E} \left[\sum_{t=1}^T q_t (f(Z'_t) - f(Z_t)) | \mathbf{Z}_1^T \right] \right) \right] \\ &\leq \mathbb{E} \left[\exp \left(\lambda \sup_{f \in \mathcal{F}} \sum_{t=1}^T q_t (f(Z'_t) - f(Z_t)) \right) \right], \end{aligned}$$

where the last expectation is taken over the joint measure of \mathbf{Z}_1^T and \mathbf{Z}'_1^T . Applying Lemma 5 (Appendix A), we can further bound this expectation by

$$\begin{aligned} & \mathbb{E}_{(\mathbf{z}, \mathbf{z}') \sim T(\mathbf{p})} \mathbb{E}_{\sigma} \left[\exp \left(\lambda \sup_{f \in \mathcal{F}} \sum_{t=1}^T \sigma_t q_t (f(\mathbf{z}'_t(\sigma)) - f(\mathbf{z}_t(\sigma))) \right) \right] \\ &\leq \mathbb{E}_{(\mathbf{z}, \mathbf{z}') \sim T(\mathbf{p})} \mathbb{E}_{\sigma} \left[\exp \left(\lambda \sup_{f \in \mathcal{F}} \sum_{t=1}^T \sigma_t q_t f(\mathbf{z}'_t(\sigma)) + \lambda \sup_{f \in \mathcal{F}} \sum_{t=1}^T -\sigma_t q_t f(\mathbf{z}_t(\sigma)) \right) \right] \\ &\leq \frac{1}{2} \mathbb{E}_{(\mathbf{z}, \mathbf{z}') \sim T(\mathbf{p})} \mathbb{E}_{\sigma} \left[\exp \left(2\lambda \sup_{f \in \mathcal{F}} \sum_{t=1}^T \sigma_t q_t f(\mathbf{z}'_t(\sigma)) \right) \right] + \frac{1}{2} \mathbb{E}_{(\mathbf{z}, \mathbf{z}') \sim T(\mathbf{p})} \mathbb{E}_{\sigma} \left[\exp \left(2\lambda \sup_{f \in \mathcal{F}} \sum_{t=1}^T \sigma_t q_t f(\mathbf{z}_t(\sigma)) \right) \right] \\ &= \mathbb{E}_{\mathbf{z} \sim T(\mathbf{p})} \mathbb{E}_{\sigma} \left[\exp \left(2\lambda \sup_{f \in \mathcal{F}} \sum_{t=1}^T \sigma_t q_t f(\mathbf{z}_t(\sigma)) \right) \right], \end{aligned}$$

where for the second inequality we used Young's inequality and for the last equality we used symmetry. Given \mathbf{z} let C denote the minimal α -cover with respect to the \mathbf{q} -weighted ℓ_1 -norm of \mathcal{F} on \mathbf{z} . Then, the following bound holds

$$\sup_{f \in \mathcal{F}} \sum_{t=1}^T \sigma_t q_t f(\mathbf{z}_t(\sigma)) \leq \max_{\mathbf{c} \in C} \sum_{t=1}^T \sigma_t q_t \mathbf{c}_t(\sigma) + \alpha.$$

By the monotonicity of the exponential function,

$$\begin{aligned} \mathbb{E}_{\sigma} \left[\exp \left(2\lambda \sup_{f \in \mathcal{F}} \sum_{t=1}^T \sigma_t q_t f(\mathbf{z}_t(\sigma)) \right) \right] &\leq \exp(2\lambda\alpha) \mathbb{E}_{\sigma} \left[\exp \left(2\lambda \max_{\mathbf{c} \in C} \sum_{t=1}^T \sigma_t q_t \mathbf{c}_t(\sigma) \right) \right] \\ &\leq \exp(2\lambda\alpha) \sum_{\mathbf{c} \in C} \mathbb{E}_{\sigma} \left[\exp \left(2\lambda \sum_{t=1}^T \sigma_t q_t \mathbf{c}_t(\sigma) \right) \right]. \end{aligned}$$

Since $\mathbf{c}_t(\sigma)$ depends only on $\sigma_1, \dots, \sigma_{T-1}$, by Hoeffding's bound,

$$\begin{aligned} \mathbb{E}_{\sigma} \left[\exp \left(2\lambda \sum_{t=1}^T \sigma_t q_t \mathbf{c}_t(\sigma) \right) \right] &= \mathbb{E} \left[\exp \left(2\lambda \sum_{t=1}^{T-1} \sigma_t q_t \mathbf{c}_t(\sigma) \right) \mathbb{E}_{\sigma_T} \left[\exp \left(2\lambda \sigma_T q_T \mathbf{c}_T(\sigma) \right) \middle| \sigma_1^{T-1} \right] \right] \\ &\leq \mathbb{E} \left[\exp \left(2\lambda \sum_{t=1}^{T-1} \sigma_t q_t \mathbf{c}_t(\sigma) \right) \exp(2\lambda^2 q_T^2 M^2) \right] \end{aligned}$$

and iterating this inequality and using the union bound, we obtain the following:

$$\mathbb{P} \left(\sup_{f \in \mathcal{F}} \sum_{t=1}^T q_t (\mathbb{E}[f(Z_t)|\mathbf{Z}_1^{t-1}] - f(Z_t)) \geq \epsilon \right) \leq \mathbb{E}_{\mathbf{v} \sim T(\mathbf{p})} [\mathcal{N}_1(\alpha, \mathcal{G}, \mathbf{v})] \exp \left(-\lambda(\epsilon - 2\alpha) + 2\lambda^2 M^2 \|\mathbf{q}\|_2^2 \right).$$

Optimizing over λ completes the proof. \square

An immediate consequence of Theorem 1 is the following result.

Corollary 2. For any $\delta > 0$, with probability at least $1 - \delta$, for all $f \in \mathcal{F}$ and all $\alpha > 0$,

$$\mathbb{E}[f(Z_{T+1})|\mathbf{Z}_1^T] \leq \sum_{t=1}^T q_t f(Z_t) + \Delta + 2\alpha + M\|\mathbf{q}\|_2 \sqrt{2 \log \frac{\mathbb{E}_{\mathbf{v} \sim T(\mathbb{P})}[\mathcal{N}_1(\alpha, \mathcal{G}, \mathbf{v})]}{\delta}}.$$

We are not aware of other finite sample bounds in a non-stationary non-mixing case. In fact, our bounds appear to be novel even in the stationary non-mixing case. Using chaining techniques bounds, Theorem 1 and Corollary 2 can be further improved and we will present these results in the full version of this paper.

While Rakhlin et al. [2015] give high probability bounds for a different quantity than the quantity of interest in time series prediction,

$$\sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T q_t (\mathbb{E}[f(Z_t)|\mathbf{Z}_1^{t-1}] - f(Z_t)) \right), \quad (3)$$

their analysis of this quantity can also be used in our context to derive high probability bounds for $\Phi(\mathbf{Z}_1^T) - \Delta$. However, this approach results in bounds that are in terms of purely combinatorial notions such as maximal sequential covering numbers $\mathcal{N}_1(\alpha, \mathcal{F})$. While at first sight, this may seem as a minor technical detail, the distinction is crucial in the setting of time series prediction. Consider the following example. Let Z_1 be drawn from a uniform distribution on $\{0, 1\}$ and $Z_t \sim p(\cdot|Z_{t-1})$ with $p(\cdot|y)$ being a distribution over $\{0, 1\}$ such that $p(x|y) = 2/3$ if $x = y$ and $1/3$ otherwise. Let \mathcal{G} be defined by $\mathcal{G} = \{g(x) = \mathbf{1}_{x \geq \theta} : \theta \in [0, 1]\}$. Then, one can check that $\mathbb{E}_{\mathbf{v} \sim T(\mathbb{P})}[\mathcal{N}_1(\alpha, \mathcal{G}, \mathbf{v})] = 2$, while $\mathcal{N}_1(\alpha, \mathcal{G}) \geq 2^T$. The data-dependent bounds of Theorem 1 and Corollary 2 highlight the fact that the task of time series prediction lies in between the familiar i.i.d. scenario and adversarial on-line learning setting.

However, the key component of our learning guarantees is the discrepancy term Δ . Note that in the general non-stationary case, the bounds of Theorem 1 may not converge to zero due to the discrepancy between the target and sample distributions. This is also consistent with the lower bounds of Barve and Long [1996] that we discuss in more detail in Section 4. However, convergence can be established in some special cases. In the i.i.d. case our bounds reduce to the standard covering numbers learning guarantees. In the drifting scenario, with \mathbf{Z}_1^T being a sequence of independent random variables, our discrepancy measure coincides with the one used and studied in [Mohri and Muñoz Medina, 2012]. Convergence can also be established in asymptotically stationary and stationary mixing cases. However, as we show in Section 4, the most important advantage of our bounds is that the discrepancy measure we use can be estimated from data.

4 Estimating Discrepancy

In Section 3, we showed that the discrepancy Δ is crucial for forecasting non-stationary time series. In particular, if we could select a distribution \mathbf{q} over the sample \mathbf{Z}_1^T that would minimize the discrepancy Δ and use it to weight training points, then we would have a better learning guarantee for an algorithm trained on this weighted sample. In some special cases, the discrepancy Δ can be computed analytically. However, in general, we do not have access to the distribution of \mathbf{Z}_1^T and hence we need to estimate the discrepancy from the data. Furthermore, in practice, we never observe Z_{T+1} and it is not possible to estimate Δ without some further assumptions. One natural assumption is that the distribution \mathbf{P}_t of Z_t does not change drastically with t on average. Under this assumption the last s observations \mathbf{Z}_{T-s+1}^T are effectively drawn from the distribution close to \mathbf{P}_{T+1} . More precisely, we can write

$$\begin{aligned} \Delta \leq & \sup_{f \in \mathcal{F}} \left(\frac{1}{s} \sum_{t=T-s+1}^T \mathbb{E}[f(Z_t)|\mathbf{Z}_1^{t-1}] - \sum_{t=1}^T q_t \mathbb{E}[f(Z_t)|\mathbf{Z}_1^{t-1}] \right) \\ & + \sup_{f \in \mathcal{F}} \left(\mathbb{E}[f(Z_{T+1})|\mathbf{Z}_1^T] - \frac{1}{s} \sum_{t=T-s+1}^T \mathbb{E}[f(Z_t)|\mathbf{Z}_1^{t-1}] \right). \end{aligned}$$

We will assume that the second term, denoted by Δ_s , is sufficiently small and will show that the first term can be estimated from data. But, we first note that our assumption is necessary for learning in

this setting. Observe that

$$\begin{aligned} \sup_{f \in \mathcal{F}} \left(\mathbb{E}[Z_{T+1} | \mathbf{Z}_1^T] - \mathbb{E}[f(Z_r) | \mathbf{Z}_1^{r-1}] \right) &\leq \sum_{t=r}^T \sup_{f \in \mathcal{F}} \left(\mathbb{E}[f(Z_{t+1}) | \mathbf{Z}_1^t] - \mathbb{E}[f(Z_t) | \mathbf{Z}_1^{t-1}] \right) \\ &\leq M \sum_{t=r}^T \|\mathbf{P}_{t+1}(\cdot | \mathbf{Z}_1^t) - \mathbf{P}_t(\cdot | \mathbf{Z}_1^{t-1})\|_{\text{TV}}, \end{aligned}$$

for all $r = T - s + 1, \dots, T$. Therefore, we must have

$$\Delta_s \leq \frac{1}{s} \sum_{t=T-s+1}^T \sup_{f \in \mathcal{F}} \left(\mathbb{E}[Z_{T+1} | \mathbf{Z}_1^T] - \mathbb{E}[f(Z_t) | \mathbf{Z}_1^t] \right) \leq \frac{s+1}{2} M \gamma,$$

where $\gamma = \sup_t \|\mathbf{P}_{t+1}(\cdot | \mathbf{Z}_1^t) - \mathbf{P}_t(\cdot | \mathbf{Z}_1^{t-1})\|_{\text{TV}}$. Barve and Long [1996] showed that $[\text{VC-dim}(H)\gamma]^{\frac{1}{3}}$ is a lower bound on the generalization error in the setting of binary classification where \mathbf{Z}_1^T is a sequence of independent but not identically distributed random variables (drifting). This setting is a special case of the more general scenario that we are considering.

The following result shows that we can estimate the first term in the upper bound on Δ .

Theorem 3. *Let \mathbf{Z}_1^T be a sequence of random variables. Then, for any $\delta > 0$, with probability at least $1 - \delta$, the following holds for all $\alpha > 0$:*

$$\sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T (p_t - q_t) \mathbb{E}[f(Z_t) | \mathbf{Z}_1^{t-1}] \right) \leq \sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T (p_t - q_t) f(Z_t) \right) + B,$$

where $B = 2\alpha + M \|\mathbf{q} - \mathbf{p}\|_2 \sqrt{2 \log \frac{\mathbb{E}_{\mathbf{z} \sim T(\mathbf{p})} [\mathcal{N}_1(\alpha, \mathcal{G}, \mathbf{z})]}{\delta}}$ and where \mathbf{p} is the uniform distribution over the last s points.

The proof of this result is given in Appendix A. Theorem 1 and Theorem 3 combined with the union bound yield the following result.

Corollary 4. *Let \mathbf{Z}_1^T be a sequence of random variables. Then, for any $\delta > 0$, with probability at least $1 - \delta$, the following holds for all $f \in \mathcal{F}$ and all $\alpha > 0$:*

$$\begin{aligned} \mathbb{E}[f(Z_{T+1}) | \mathbf{Z}_1^T] &\leq \\ &\sum_{t=1}^T q_t f(Z_t) + \tilde{\Delta} + \Delta_s + 4\alpha + M [\|\mathbf{q}\|_2 + \|\mathbf{q} - \mathbf{p}\|_2] \sqrt{2 \log \frac{2 \mathbb{E}_{\mathbf{v} \sim T(\mathbf{p})} [\mathcal{N}_1(\alpha, \mathcal{G}, \mathbf{z})]}{\delta}}, \end{aligned}$$

where $\tilde{\Delta} = \sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T (p_t - q_t) f(Z_t) \right)$.

5 Algorithms

In this section, we use our learning guarantees to devise algorithms for forecasting non-stationary time series. We consider a broad family of kernel-based hypothesis classes with regression losses. We present the full analysis of this setting in Appendix B including novel bounds on the sequential Rademacher complexity. The learning bounds of Theorem 1 can be generalized to hold uniformly over \mathbf{q} at the price of an additional term in $O\left(\|\mathbf{q} - \mathbf{u}\|_1 \sqrt{\log_2 \log_2 \|\mathbf{q} - \mathbf{u}\|_1^{-1}}\right)$. We prove this result in Theorem 8 (Appendix B). Suppose L is the squared loss and $H = \{\mathbf{x} \rightarrow \mathbf{w} \cdot \Psi(\mathbf{x}) : \|\mathbf{w}\|_{\mathcal{H}} \leq \Lambda\}$, where $\Psi: \mathcal{X} \rightarrow \mathcal{H}$ is a feature mapping from \mathcal{X} to a Hilbert space \mathcal{H} . By Lemma 6 (Appendix B), we can bound the complexity term in our generalization bounds by

$$O\left((\log^3 T) \frac{\Lambda r}{\sqrt{T}} + (\log^3 T) \|\mathbf{q} - \mathbf{u}\|_1\right),$$

where K is a PDS kernel associated with \mathcal{H} such that $\sup_x K(x, x) \leq r$ and \mathbf{u} is the uniform distribution over the sample. Then, we can formulate a joint optimization problem over both \mathbf{q} and \mathbf{w} based on the learning guarantee of Theorem 8, which holds uniformly over all \mathbf{q} :

$$\min_{0 \leq \mathbf{q} \leq 1, \mathbf{w}} \left\{ \sum_{t=1}^T q_t (\mathbf{w} \cdot \Psi(x_t) - y_t)^2 + \lambda_1 \sum_{t=1}^T d_t q_t + \lambda_2 \|\mathbf{w}\|_{\mathcal{H}}^2 + \lambda_3 \|\mathbf{q} - \mathbf{u}\|_1 \right\}. \quad (4)$$

Here, we have upper bounded the empirical discrepancy term by $\sum_{t=1}^T d_t q_t$ with each d_t defined by $\sup_{\mathbf{w}' \leq \Lambda} |\sum_{s=1}^T p_s(\mathbf{w}' \cdot \Psi(x_s) - y_s)^2 - (\mathbf{w}' \cdot \Psi(x_t) - y_t)^2|$. Each d_t can be precomputed using DC-programming. For general loss functions, the DC-programming approach only guarantees convergence to a stationary point. However, for the squared loss, our problem can be cast as an instance of the trust region problem, which can be solved globally using the DCA algorithm of Tao and An [1998]. Note that problem (4) is not jointly convex in \mathbf{q} and \mathbf{w} . However, using the dual problem associated to \mathbf{w} yields the following equivalent problem, it can be rewritten as follows:

$$\min_{0 \leq \mathbf{q} \leq 1} \left\{ \max_{\alpha} \left\{ -\lambda_2 \sum_{t=1}^T \frac{\alpha_t^2}{q_t} - \alpha^T \mathbf{K} \alpha + 2\lambda_2 \alpha^T \mathbf{Y} \right\} + \lambda_1 (\mathbf{d} \cdot \mathbf{q}) + \lambda_3 \|\mathbf{q} - \mathbf{u}\|_1 \right\}, \quad (5)$$

where $\mathbf{d} = (d_1, \dots, d_T)^T$, \mathbf{K} is the kernel matrix and $\mathbf{Y} = (y_1, \dots, y_T)^T$. We use the change of variables $r_t = 1/q_t$ and further upper bound $\lambda_3 \|\mathbf{q} - \mathbf{u}\|_1$ by $\lambda'_3 \|\mathbf{r} - T^2 \mathbf{u}\|_2$, which follows from $|q_t - u_t| = |q_t u_t (r_t - T)|$ and Hölder's inequality. Then, this yields the following optimization problem:

$$\min_{\mathbf{r} \in \mathcal{D}} \left\{ \max_{\alpha} \left\{ -\lambda_2 \sum_{t=1}^T r_t \alpha_t^2 - \alpha^T \mathbf{K} \alpha + 2\lambda_2 \alpha^T \mathbf{Y} \right\} + \lambda_1 \sum_{t=1}^T \frac{d_t}{r_t} + \lambda_3 \|\mathbf{r} - T^2 \mathbf{u}\|_2^2 \right\}, \quad (6)$$

where $\mathcal{D} = \{\mathbf{r}: r_t \geq 1, t \in [1, T]\}$. The optimization problem (6) is convex since \mathcal{D} is a convex set, the first term in (6) is convex as a maximum of convex (linear) functions of \mathbf{r} . This problem can be solved using standard descent methods, where, at each iteration, we solve a standard QP in α , which admits a closed-form solution. Parameters λ_1 , λ_2 , and λ_3 are selected through cross-validation.

An alternative simpler algorithm based on the data-dependent bounds of Corollary 4 consists of first finding a distribution \mathbf{q} minimizing the (regularized) discrepancy and then using that to find a hypothesis minimizing the (regularized) weighted empirical risk. This leads to the following two-stage procedure. First, we find a solution \mathbf{q}^* of the following convex optimization problem:

$$\min_{\mathbf{q} \geq 0} \left\{ \sup_{\mathbf{w}' \leq \Lambda} \left(\sum_{t=1}^T (p_t - q_t) (\mathbf{w}' \cdot \Psi(x_t) - y_t)^2 \right) + \lambda_1 \|\mathbf{q} - \mathbf{u}\|_1 \right\}, \quad (7)$$

where λ_1 and Λ are parameters that can be selected via cross-validation. Our generalization bounds hold for arbitrary weights \mathbf{q} but we restrict them to being positive sequences. Note that other regularization terms such as $\|\mathbf{q}\|_2^2$ and $\|\mathbf{q} - \mathbf{p}\|_2^2$ from the bound of Corollary 4 can be incorporated in the optimization problem, but we discard them to minimize the number of parameters. This problem can be solved using standard descent optimization methods, where, at each step, we use DC-programming to evaluate the supremum over \mathbf{w}' . Alternatively, one can upper bound the supremum by $\sum_{t=1}^T q_t d_t$ and then solve the resulting optimization problem.

The solution \mathbf{q}^* of (7) is then used to solve the following (weighted) kernel ridge regression problem:

$$\min_{\mathbf{w}} \left\{ \sum_{t=1}^T q_t^* (\mathbf{w} \cdot \Psi(x_t) - y_t)^2 + \lambda_2 \|\mathbf{w}\|_{\mathcal{H}}^2 \right\}. \quad (8)$$

Note that, in order to guarantee the convexity of this problem, we require $\mathbf{q}^* \geq 0$.

6 Conclusion

We presented a general theoretical analysis of learning in the broad scenario of non-stationary non-mixing processes, the realistic setting for a variety of applications. We discussed in detail several algorithms benefitting from the learning guarantees presented. Our theory can also provide a finer analysis of several existing algorithms and help devise alternative principled learning algorithms.

Acknowledgments

This work was partly funded by NSF IIS-1117591 and CCF-1535987, and the NSERC PGS D3.

References

- T. M. Adams and A. B. Nobel. Uniform convergence of Vapnik-Chervonenkis classes under ergodic sampling. *The Annals of Probability*, 38(4):1345–1367, 2010.
- A. Agarwal and J. Duchi. The generalization ability of online algorithms for dependent data. *Information Theory, IEEE Transactions on*, 59(1):573–587, 2013.
- P. Alquier and O. Wintenberger. Model selection for weakly dependent time series forecasting. Technical Report 2010-39, Centre de Recherche en Economie et Statistique, 2010.
- P. Alquier, X. Li, and O. Wintenberger. Prediction of time series by statistical learning: general losses and fast rates. *Dependence Modelling*, 1:65–93, 2014.
- D. Andrews. First order autoregressive processes and strong mixing. Cowles Foundation Discussion Papers 664, Cowles Foundation for Research in Economics, Yale University, 1983.
- R. Baillie. Long memory processes and fractional integration in econometrics. *Journal of Econometrics*, 73(1):5–59, 1996.
- R. D. Barve and P. M. Long. On the complexity of learning from drifting distributions. In *COLT*, 1996.
- P. Berti and P. Rigo. A Glivenko-Cantelli theorem for exchangeable random variables. *Statistics & Probability Letters*, 32(4):385 – 391, 1997.
- T. Bollerslev. Generalized autoregressive conditional heteroskedasticity. *J Econometrics*, 1986.
- G. E. P. Box and G. Jenkins. *Time Series Analysis, Forecasting and Control*. Holden-Day, Incorporated, 1990.
- P. J. Brockwell and R. A. Davis. *Time Series: Theory and Methods*. Springer-Verlag, New York, 1986.
- V. H. De la Peña and E. Giné. *Decoupling: from dependence to independence: randomly stopped processes, U-statistics and processes, martingales and beyond*. Probability and its applications. Springer, NY, 1999.
- P. Doukhan. *Mixing: properties and examples*. Lecture notes in statistics. Springer-Verlag, New York, 1994.
- R. Engle. Autoregressive conditional heteroscedasticity with estimates of the variance of United Kingdom inflation. *Econometrica*, 50(4):987–1007, 1982.
- J. D. Hamilton. *Time series analysis*. Princeton, 1994.
- V. Kuznetsov and M. Mohri. Generalization bounds for time series prediction with non-stationary processes. In *ALT*, 2014.
- A. C. Lozano, S. R. Kulkarni, and R. E. Schapire. Convergence and consistency of regularized boosting algorithms with stationary β -mixing observations. In *NIPS*, pages 819–826, 2006.
- R. Meir. Nonparametric time series prediction through adaptive model selection. *Machine Learning*, pages 5–34, 2000.
- D. Modha and E. Masry. Memory-universal prediction of stationary random processes. *Information Theory, IEEE Transactions on*, 44(1):117–133, Jan 1998.
- M. Mohri and A. Muñoz Medina. New analysis and algorithm for learning with drifting distributions. In *ALT*, 2012.
- M. Mohri and A. Rostamizadeh. Rademacher complexity bounds for non-i.i.d. processes. In *NIPS*, 2009.
- M. Mohri and A. Rostamizadeh. Stability bounds for stationary φ -mixing and β -mixing processes. *Journal of Machine Learning Research*, 11:789–814, 2010.
- V. Pestov. Predictive PAC learnability: A paradigm for learning from exchangeable input data. In *GRC*, 2010.
- A. Rakhlin, K. Sridharan, and A. Tewari. Online learning: Random averages, combinatorial parameters, and learnability. In *NIPS*, 2010.
- A. Rakhlin, K. Sridharan, and A. Tewari. Online learning: Stochastic, constrained, and smoothed adversaries. In *NIPS*, 2011.
- A. Rakhlin, K. Sridharan, and A. Tewari. Sequential complexities and uniform martingale laws of large numbers. *Probability Theory and Related Fields*, 2015.
- C. Shalizi and A. Kontorovitch. Predictive PAC learning and process decompositions. In *NIPS*, 2013.
- I. Steinwart and A. Christmann. Fast learning from non-i.i.d. observations. In *NIPS*, 2009.
- P. D. Tao and L. T. H. An. A D.C. optimization algorithm for solving the trust-region subproblem. *SIAM Journal on Optimization*, 8(2):476–505, 1998.
- M. Vidyasagar. *A Theory of Learning and Generalization: With Applications to Neural Networks and Control Systems*. Springer-Verlag New York, Inc., 1997.
- B. Yu. Rates of convergence for empirical processes of stationary mixing sequences. *The Annals of Probability*, 22(1):94–116, 1994.

A Proofs

Lemma 5. *Given a sequence of random variables \mathbf{Z}_1^T with joint distribution \mathbf{p} , let \mathbf{Z}'_1^T be a decoupled tangent sequence. Then, for any measurable function G , the following equality holds*

$$\mathbb{E} \left[G \left(\sup_{f \in \mathcal{F}} \sum_{t=1}^T q_t (f(Z'_t) - f(Z_t)) \right) \right] = \mathbb{E}_{\boldsymbol{\sigma}} \mathbb{E}_{\mathbf{z} \sim T(\mathbf{p})} \left[G \left(\sup_f \sum_{t=1}^T \sigma_t q_t (f(\mathbf{z}'_t(\boldsymbol{\sigma})) - f(\mathbf{z}_t(\boldsymbol{\sigma}))) \right) \right]. \quad (9)$$

The result also holds with the absolute value around the sums in (9).

Proof. The proof follows an argument in the proof of Theorem 3 of [Rakhlin et al., 2011]. We only need to check that every step holds for an arbitrary weight vector \mathbf{q} , in lieu of the uniform distribution vector \mathbf{u} , and for an arbitrary measurable function G , instead of the identity function. Observe that we can write the left-hand side of (9) as

$$\mathbb{E} \left[G \left(\sup_{f \in \mathcal{F}} \Sigma(\boldsymbol{\sigma}) \right) \right] = \mathbb{E}_{Z_1, Z'_1 \sim \mathbf{p}_1} \mathbb{E}_{Z_2, Z'_2 \sim \mathbf{p}_2(\cdot | S_1)} \cdots \mathbb{E}_{Z_T, Z'_T \sim \mathbf{p}_T(\cdot | \mathbf{Z}_1^{T-1})} \left[G \left(\sup_{f \in \mathcal{F}} \Sigma(\boldsymbol{\sigma}) \right) \right],$$

where $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_T) \in \{\pm 1\}^T$ and $\Sigma(\boldsymbol{\sigma}) = \sum_{t=1}^T \sigma_t q_t (f(Z'_t) - f(Z_t))$. Now, by definition of decoupled tangent sequences, the value of the last expression is unchanged if we swap the sign of any σ_{i-1} to -1 since that is equivalent to permuting Z_i and Z'_i . Thus, the last expression is in fact equal to

$$\mathbb{E}_{Z_1, Z'_1 \sim \mathbf{p}_1} \mathbb{E}_{Z_2, Z'_2 \sim \mathbf{p}_2(\cdot | S_1(\sigma_1))} \cdots \mathbb{E}_{Z_T, Z'_T \sim \mathbf{p}_T(\cdot | S_1(\sigma_1), \dots, S_{T-1}(\sigma_{T-1}))} \left[G \left(\sup_{f \in \mathcal{F}} \Sigma(\boldsymbol{\sigma}) \right) \right]$$

for any sequence $\boldsymbol{\sigma} \in \{\pm 1\}^T$, where $S_t(1) = Z_t$ and Z'_t otherwise. Since this equality holds for any $\boldsymbol{\sigma}$, it also holds for the mean with respect to uniformly distributed $\boldsymbol{\sigma}$. Therefore, the last expression is equal to

$$\mathbb{E}_{\boldsymbol{\sigma}} \mathbb{E}_{Z_1, Z'_1 \sim \mathbf{p}_1} \mathbb{E}_{Z_2, Z'_2 \sim \mathbf{p}_2(\cdot | S_1(\sigma_1))} \cdots \mathbb{E}_{Z_T, Z'_T \sim \mathbf{p}_T(\cdot | S_1(\sigma_1), \dots, S_{T-1}(\sigma_{T-1}))} \left[G \left(\sup_{f \in \mathcal{F}} \Sigma(\boldsymbol{\sigma}) \right) \right].$$

This last expectation coincides with the expectation with respect to drawing a random tree \mathbf{z} from $T(\mathbf{p})$ (and its tangent tree \mathbf{z}') and a random path $\boldsymbol{\sigma}$ to follow in that tree. That is, the last expectation is equal to

$$\mathbb{E}_{\boldsymbol{\sigma}} \mathbb{E}_{\mathbf{z} \sim T(\mathbf{p})} \left[G \left(\sup_f \sum_{t=1}^T \sigma_t q_t (f(\mathbf{z}'_t(\boldsymbol{\sigma})) - f(\mathbf{z}_t(\boldsymbol{\sigma}))) \right) \right],$$

which concludes the proof. \square

Theorem 3. *Let \mathbf{Z}_1^T be a sequence of random variables. Then, for any $\delta > 0$, with probability at least $1 - \delta$, the following holds for all $\alpha > 0$:*

$$\sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T (p_t - q_t) \mathbb{E}[f(Z_t) | \mathbf{Z}_1^{t-1}] \right) \leq \sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T (p_t - q_t) f(Z_t) \right) + \alpha + M \|\mathbf{q} - \mathbf{p}\|_2 \sqrt{\log \frac{\mathbb{E}_{\mathbf{z} \sim T(\mathbf{p})} [\mathcal{N}_1(\alpha, \mathcal{G}, \mathbf{z})]}{\delta}},$$

where \mathbf{p} is the distribution the uniform on the last s points.

Proof. First, observe that

$$\begin{aligned} \sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T (p_t - q_t) \mathbb{E}[f(Z_t) | \mathbf{Z}_1^{t-1}] \right) - \sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T (p_t - q_t) f(Z_t) \right) \\ \leq \sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T (p_t - q_t) (\mathbb{E}[f(Z_t) | \mathbf{Z}_1^{t-1}] - f(Z_t)) \right). \end{aligned}$$

The result then follows using similar arguments to those used in the proof of Theorem 1. \square

B Generalization Bounds for Kernel-Based Hypotheses with Regression Losses

In this section, we present generalization bounds for kernel-based hypothesis with regression losses. One of the main technical tools used in our analysis is the notion of *sequential Rademacher complexity*. Let \mathcal{G} be a set of functions from \mathcal{Z} to \mathbb{R} . The sequential Rademacher complexity of a function class \mathcal{Z} is defined as the following:

$$\mathfrak{R}_T^{\text{seq}}(\mathcal{G}) = \sup_{\mathbf{z}} \mathbb{E} \left[\sup_{g \in \mathcal{G}} \sum_{t=1}^T \sigma_t q_t g(z_t(\boldsymbol{\sigma})) \right], \quad (10)$$

where the supremum is taken over all complete binary trees of depth T with values in \mathcal{Z} and where $\boldsymbol{\sigma}$ is a sequence of Rademacher random variables. This a combinatorial measure of complexity which makes bounds based on this notion coarser than those of Theorem 1, which are stated in terms of expected covering numbers. However, it turns out that this coarser analysis is sufficient for the derivation of our algorithms in Section 5. We also remark that most of the results in this section can be tightened using the notion of *distribution-dependent* Rademacher complexity, but we defer these results to the full version of the paper.

Our first result is a bound on the sequential Rademacher complexity of the kernel-based hypothesis with regression losses.

Lemma 6. *Let $p \geq 1$ and $\mathcal{F} = \{(\mathbf{x}, y) \rightarrow (\mathbf{w} \cdot \Psi(\mathbf{x}) - y)^p : \|\mathbf{w}\|_{\mathcal{H}} \leq \Lambda\}$ where \mathcal{H} is a Hilbert space and $\Psi: \mathcal{X} \rightarrow \mathcal{H}$ a feature map. Assume that the condition $|\mathbf{w} \cdot \mathbf{x} - y| \leq M$ holds for all $(\mathbf{x}, y) \in \mathcal{Z}$ and all \mathbf{w} such that $\|\mathbf{w}\|_{\mathcal{H}} \leq \Lambda$. Then, the following inequalities hold:*

$$\mathfrak{R}_T^{\text{seq}}(\mathcal{F}) \leq pM^{p-1} C_T \mathfrak{R}_T^{\text{seq}}(H) \leq C_T \left(pM^{p-1} \frac{\Lambda r}{\sqrt{T}} + pM^p \|\mathbf{q} - \mathbf{u}\|_1 \right), \quad (11)$$

where K is a PDS kernel associated to \mathcal{H} , $H = \{\mathbf{x} \rightarrow \mathbf{w} \cdot \Psi(\mathbf{x}) : \|\mathbf{w}\|_{\mathcal{H}} \leq \Lambda\}$, $r = \sup_x K(x, x)$, and $C_T = 8(1 + 4\sqrt{2} \log^{3/2}(eT^2))$.

Proof. We begin the proof by setting $q_t f(\mathbf{z}_t(\boldsymbol{\sigma})) = q_t (\mathbf{w} \cdot \Psi(\mathbf{x}_t(\boldsymbol{\sigma})) - \mathbf{y}_t(\boldsymbol{\sigma}))^2 = \frac{1}{T} (\mathbf{w} \cdot \mathbf{x}'_t(\boldsymbol{\sigma}) - \mathbf{y}'_t(\boldsymbol{\sigma}))^2$, where $\mathbf{x}'_t(\boldsymbol{\sigma}) = \sqrt{T q_t} \Psi(\mathbf{x}_t(\boldsymbol{\sigma}))$ and $\mathbf{y}'_t(\boldsymbol{\sigma}) = \sqrt{T q_t} \mathbf{y}_t(\boldsymbol{\sigma})$. We let $\mathbf{z}'_t = (\mathbf{x}'_t, \mathbf{y}'_t)$. Then we observe that

$$\begin{aligned} \mathfrak{R}_T^{\text{seq}}(\mathcal{F}) &= \sup_{\mathbf{z}' = (\mathbf{x}', \mathbf{y}')} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\mathbf{w}} \frac{1}{T} \sum_{t=1}^T \sigma_t (\mathbf{w} \cdot \mathbf{x}'_t(\boldsymbol{\sigma}) - \mathbf{y}'_t(\boldsymbol{\sigma}))^p \right] \\ &= \sup_{\mathbf{z} = (\mathbf{x}, \mathbf{y})} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\mathbf{w}} \sum_{t=1}^T q_t \sigma_t (\mathbf{w} \cdot \mathbf{x}_t(\boldsymbol{\sigma}) - \mathbf{y}_t(\boldsymbol{\sigma}))^p \right]. \end{aligned}$$

Since $x \rightarrow |x|^p$ is pM^{p-1} -Lipschitz over $[-M, M]$, by Lemma 13 in [Rakhlin et al., 2015], the following bound holds:

$$\mathfrak{R}_T^{\text{seq}}(\mathcal{F}) \leq pM^{p-1} C_T \mathfrak{R}_T^{\text{seq}}(H'),$$

where $H' = \{(\mathbf{x}, y) \rightarrow \mathbf{w} \cdot \Psi(\mathbf{x}) - y : \|\mathbf{w}\|_{\mathcal{H}} \leq \Lambda\}$. Note that Lemma 13 requires that $\mathfrak{R}_T^{\text{seq}}(H') > 1/T$ which is guaranteed by Khintchine's inequality. By definition of the sequential Rademacher complexity

$$\begin{aligned} \mathfrak{R}_T^{\text{seq}}(H') &= \sup_{(\mathbf{x}, y)} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\mathbf{w}} \sum_{t=1}^T \sigma_t q_t (\mathbf{w} \cdot \Psi(\mathbf{x}_t(\boldsymbol{\sigma})) - y(\boldsymbol{\sigma})) \right] \\ &= \sup_{\mathbf{x}} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\mathbf{w}} \sum_{t=1}^T \sigma_t q_t \mathbf{w} \cdot \Psi(\mathbf{x}_t(\boldsymbol{\sigma})) \right] + \sup_y \mathbb{E}_{\boldsymbol{\sigma}} \left[\sum_{t=1}^T \sigma_t q_t y(\boldsymbol{\sigma}) \right] = \mathfrak{R}_T^{\text{seq}}(H), \end{aligned}$$

where for the last equality we used the fact that σ_t s are mean zero random variables and σ_t is independent of $y(\boldsymbol{\sigma}) = y(\sigma_1, \sigma_2, \dots, \sigma_{t-1})$. This proves the first result. To prove the second bound we first observe that

$$\mathfrak{R}_T^{\text{seq}}(H) \leq \frac{1}{T} \sup_{\mathbf{x}} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\mathbf{w}} \sum_{t=1}^T \sigma_t \mathbf{w} \cdot \Psi(\mathbf{x}_t(\boldsymbol{\sigma})) \right] + M \|\mathbf{q} - \mathbf{u}\|_1.$$

Next, the first term on the right-hand side can be bounded as follows:

$$\begin{aligned}
\frac{1}{T} \sup_{\mathbf{x}} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\mathbf{w}} \sum_{t=1}^T \sigma_t \mathbf{w} \cdot \Psi(\mathbf{x}_t(\boldsymbol{\sigma})) \right] &\leq \frac{\Lambda}{T} \sup_{\mathbf{x}} \mathbb{E}_{\boldsymbol{\sigma}} \left\| \sum_{t=1}^T \sigma_t \Psi(\mathbf{x}_t(\boldsymbol{\sigma})) \right\|_{\mathcal{H}} \\
&\leq \frac{\Lambda}{T} \sup_{\mathbf{x}} \sqrt{\mathbb{E}_{\boldsymbol{\sigma}} \left\| \sum_{t=1}^T \sigma_t \Psi(\mathbf{x}_t(\boldsymbol{\sigma})) \right\|_{\mathcal{H}}^2} \\
&= \frac{\Lambda}{T} \sup_{\mathbf{x}} \sqrt{\mathbb{E}_{\boldsymbol{\sigma}} \left[\sum_{t,s=1}^T \sigma_t \sigma_s \Psi(\mathbf{x}_t(\boldsymbol{\sigma})) \cdot \Psi(\mathbf{x}_s(\boldsymbol{\sigma})) \right]} \\
&\leq \frac{\Lambda}{T} \sup_{\mathbf{x}} \sqrt{\sum_{t=1}^T \mathbb{E}_{\boldsymbol{\sigma}} [K(x_t(\boldsymbol{\sigma}), x_t(\boldsymbol{\sigma}))]} \\
&\leq \frac{\Lambda r}{\sqrt{T}},
\end{aligned}$$

where again we are using the fact that if $s < t$ then

$$\mathbb{E}_{\boldsymbol{\sigma}} [\sigma_t \sigma_s K(x_t(\boldsymbol{\sigma}), x_s(\boldsymbol{\sigma}))] = \mathbb{E}_{\boldsymbol{\sigma}} [\sigma_t] \mathbb{E}_{\boldsymbol{\sigma}} [\sigma_s K(x_t(\boldsymbol{\sigma}), x_s(\boldsymbol{\sigma}))] = 0$$

by the independence of σ_t from σ_s , $x_t(\boldsymbol{\sigma}) = x_t(\sigma_1, \dots, \sigma_{t-1})$ and $x_s(\boldsymbol{\sigma}) = x_s(\sigma_1, \dots, \sigma_s)$. \square

Our next result establishes a high-probability learning guarantee for kernel-based hypothesis.

Theorem 7. *Let $p \geq 1$ and $\mathcal{F} = \{(\mathbf{x}, y) \rightarrow (\mathbf{w} \cdot \Psi(\mathbf{x}) - y)^p : \|\mathbf{w}\|_{\mathcal{H}} \leq \Lambda\}$ where \mathcal{H} is a Hilbert space and $\Psi: \mathcal{X} \rightarrow \mathcal{H}$ a feature map. Assume that the condition $|\mathbf{w} \cdot \mathbf{x} - y| \leq M$ holds for all $(\mathbf{x}, y) \in \mathcal{Z}$ and all \mathbf{w} such that $\|\mathbf{w}\|_{\mathcal{H}} \leq \Lambda$. If $\mathbf{Z}_1^T = (\mathbf{X}_1^T, \mathbf{Y}_1^T)$ is a sequence of random variables then, for any $\delta > 0$, with probability at least $1 - \delta$ the following holds for all $h \in \{\mathbf{x} \rightarrow \mathbf{w} \cdot \Psi(\mathbf{x}) : \|\mathbf{w}\|_{\mathcal{H}} \leq \Lambda\}$:*

$$\begin{aligned}
\mathbb{E}[(h(X_{T+1}) - Y_{T+1})^p | \mathbf{Z}_1^T] &\leq \sum_{t=1}^T (h(X_t) - Y_t)^p + \Delta \\
&+ M \tilde{C}_T \sqrt{\log \frac{8L}{\delta}} \left(p M^{p-1} \frac{\Lambda r}{\sqrt{T}} + p M^p \|\mathbf{q} - \mathbf{u}\|_1 \right),
\end{aligned}$$

where $L = \max \{e^4, \sum_{j=1}^{\infty} \mathcal{N}_{\infty}(2^{-j}, \mathcal{F})\}$, $\tilde{C}_T = c \log^{3/2} T (1 + 4\sqrt{2} \log^{3/2}(eT^2))$ and where c is an absolute constant. Thus, for $p = 2$,

$$\mathbb{E}[(h(X_{T+1}) - Y_{T+1})^2 | \mathbf{Z}_1^T] \leq \sum_{t=1}^T (h(X_t) - Y_t)^2 + \Delta + O\left((\log^3 T) \frac{\Lambda r}{\sqrt{T}} + (\log^3 T) \|\mathbf{q} - \mathbf{u}\|_1 \right).$$

Note that for this result to be non-trivial we need $\sum_{j=1}^{\infty} \mathcal{N}_{\infty}(2^{-j}, \mathcal{F}) < \infty$. This condition is easy to verify in our case. First, observe that for any set of linear functions the inequality $\mathcal{N}_{\infty}(\alpha, H) > \Lambda r / \alpha$ holds and it follows that $\sum_{j=1}^{\infty} \mathcal{N}_{\infty}(2^{-j}, \mathcal{F}) < 2\Lambda r$. The case of composition of H with ℓ_p loss can be handled by realizing that this composition leads to a linear function in a higher dimensional space corresponding to a polynomial kernel of degree p .

Proof. The beginning of the proof closely follows that of Theorem 1. The first step is to observe that since the difference of the suprema is bounded by the supremum of the difference, it suffices to bound the probability of the following event

$$\left\{ \sup_{f \in \mathcal{F}} \left(\sum_{t=1}^T q_t (\mathbb{E}[f(Z_t) | \mathbf{Z}_1^{t-1}] - f(Z_t)) \right) \geq \epsilon \right\}.$$

Next, we note that $q_t f(Z_t) = q_t(\mathbf{w} \cdot \Psi(X_t) - Y_t)^2 = \frac{1}{T}(\mathbf{w} \cdot X'_t - Y'_t)^2$, where $X'_t = \sqrt{T}q_t\Psi(X_t)$ and $Y'_t = \sqrt{T}q_t Y_t$. We let $Z'_t = (X'_t, Y'_t)$. Applying Lemma 15 from Rakhlin et al. [2015], we obtain that for any $\delta > 0$ with probability at least $1 - \delta$, the following holds for all $f \in \mathcal{F}$:

$$\Phi(\mathbf{Z}_1^T) - \Delta \leq M\mathfrak{R}_T^{\text{seq}}(\mathcal{F})(\log^{3/2} T)\sqrt{c \log \frac{8L}{\delta}},$$

where

$$\begin{aligned} \mathfrak{R}_T^{\text{seq}}(\mathcal{F}) &= \sup_{\mathbf{z}'=(\mathbf{x}',\mathbf{y}')} \mathbb{E} \left[\sup_{\mathbf{w}} \frac{1}{T} \sum_{t=1}^T \sigma_t(\mathbf{w} \cdot \mathbf{x}'_t(\sigma) - \mathbf{y}'_t(\sigma))^p \right] \\ &= \sup_{\mathbf{z}=(\mathbf{x},\mathbf{y})} \mathbb{E} \left[\sup_{\mathbf{w}} \sum_{t=1}^T q_t \sigma_t(\mathbf{w} \cdot \mathbf{x}_t(\sigma) - \mathbf{y}_t(\sigma))^p \right] \end{aligned}$$

is the sequential Rademacher complexity of \mathcal{F} . Note that $\mathfrak{R}_T^{\text{seq}}(\mathcal{F}) > 1/T$ as in the proof of Lemma 6. The desired result follows from Lemma 6. \square

The final result of this section extends Theorem 7 to hold uniformly over \mathbf{q} s.

Theorem 8. *Let $p \geq 1$ and $\mathcal{F} = \{(\mathbf{x}, y) \rightarrow (\mathbf{w} \cdot \Psi(\mathbf{x}) - y)^p : \|\mathbf{w}\|_{\mathcal{H}} \leq \Lambda\}$ where \mathcal{H} is a Hilbert space and $\Psi: \mathcal{X} \rightarrow \mathcal{H}$ a feature map. Assume that the condition $|\mathbf{w} \cdot \mathbf{x} - y| \leq M$ holds for all $(\mathbf{x}, y) \in \mathcal{Z}$ and all \mathbf{w} such that $\|\mathbf{w}\|_{\mathcal{H}} \leq \Lambda$. Then, if $\mathbf{Z}_1^T = (\mathbf{X}_1^T, \mathbf{Y}_1^T)$ is a sequence of random variables, for any $\delta > 0$, with probability at least $1 - \delta$, the following bound holds for all $h \in H = \{\mathbf{x} \rightarrow \mathbf{w} \cdot \Psi(\mathbf{x}) : \|\mathbf{w}\|_{\mathcal{H}} \leq \Lambda\}$ and all \mathbf{q} such that $0 < \|\mathbf{q} - \mathbf{u}\|_1 \leq 1$:*

$$\begin{aligned} \mathbb{E}[(h(X_{T+1}) - Y_{T+1})^p | \mathbf{Z}_1^T] &\leq \sum_{t=1}^T (h(X_t) - Y_t)^p + \Delta + 4M\|\mathbf{q} - \mathbf{u}\|_1 \\ &+ M\tilde{C}_T \left(\sqrt{\log \frac{16L}{\delta}} + \sqrt{\log \log_2 2\|\mathbf{q} - \mathbf{u}\|^{-1}} \right) \left(pM^{p-1} \frac{\Lambda r}{\sqrt{T}} + pM^p \|\mathbf{q} - \mathbf{u}\|_1 \right), \end{aligned}$$

where $\tilde{C}_T = c \log^{3/2} T (1 + 4\sqrt{2} \log^{3/2}(eT^2))$, c is an absolute constant and $L = \max\{e^4, \sum_{j=1}^{\infty} \mathcal{N}_{\infty}(2^{-j}, \mathcal{F})\}$. Thus, for $p = 2$,

$$\begin{aligned} \mathbb{E}[(h(X_{T+1}) - Y_{T+1})^2 | \mathbf{Z}_1^T] &\leq \sum_{t=1}^T (h(X_t) - Y_t)^2 + \Delta \\ &+ O\left((\log^3 T) \sqrt{\log \log_2 2\|\mathbf{q} - \mathbf{u}\|^{-1}} \left(\frac{\Lambda r}{\sqrt{T}} + \|\mathbf{q} - \mathbf{u}\|_1 \right) \right). \end{aligned}$$

This result suggests that we should try to minimize $\sum_{t=1}^T q_t f(Z_t) + \Delta$ over \mathbf{q} and \mathbf{w} making sure that \mathbf{q} does not deviate from \mathbf{u} by more than $O(T^{-1/2})$. Theorem 1 can be extended in a similar way to hold uniformly over \mathbf{q} s and we will provide this result in the full version of the paper.

Proof. Let $(\epsilon_k)_{k=0}^{\infty}$ and $(\mathbf{q}(k))_{k=0}^{\infty}$ be infinite sequences specified below. By Theorem 7, the following holds for each k

$$\mathbb{P}\left(\mathbb{E}[f(Z_{T+1}) | \mathbf{Z}_1^T] > \sum_{t=1}^T q_t(k) f(Z_t) + \Delta(\mathbf{q}(k)) + C\epsilon_k \right) \leq 8L \exp(-\epsilon_k^2),$$

where $\Delta(\mathbf{q}(k))$ denotes the discrepancy computed with respect to the weights $\mathbf{q}(k)$ and C is equal to

$$M\tilde{C}_T \left(pM^{p-1} \frac{\Lambda r}{\sqrt{T}} + pM^p \|\mathbf{q}(k) - \mathbf{u}\|_1 \right).$$

Table 1: Average squared error (standard deviation)

	ads1	ads2	ads3
DBF	0.0001 (0.0001)	0.0002 (0.0001)	0.0047 (0.0001)
WAR	0.0099 (0.0155)	0.0997 (0.1449)	0.1026 (0.1509)
ARIMA	0.1432 (0.2091)	0.4797 (0.6942)	0.2598 (0.3696)

Let $\epsilon_k = \epsilon + \sqrt{2 \log k}$. Then, by the union bound we can write

$$\begin{aligned} \mathbb{P}\left(\exists k: \mathbb{E}[f(Z_{T+1})|\mathbf{Z}_1^T] > \sum_{t=1}^T q_t(k)f(Z_t) + \Delta(\mathbf{q}(k)) + C\epsilon_k\right) &\leq \sum_{k=1}^{\infty} 8L \exp(-\epsilon_k^2) \\ &\leq \sum_{k=1}^{\infty} 8L \exp(-\epsilon^2 - \log k^2) \\ &\leq 16L \exp(-\epsilon^2). \end{aligned}$$

We choose the sequence $\mathbf{q}(k)$ to satisfy $\|\mathbf{q}(k) - \mathbf{u}\|_1 = 2^{-k}$. Then, for any \mathbf{q} such that $0 < \|\mathbf{q} - \mathbf{u}\|_1 \leq 1$, there exists k such that

$$\|\mathbf{q}(k) - \mathbf{u}\|_1 < \|\mathbf{q} - \mathbf{u}\|_1 \leq \|\mathbf{q}(k-1) - \mathbf{u}\|_1 = 2\|\mathbf{q}(k) - \mathbf{u}\|_1.$$

Thus, the following inequality holds:

$$\sqrt{2 \log k} \leq \sqrt{2 \log \log_2 2 \|\mathbf{q} - \mathbf{u}\|_1^{-1}}.$$

Combining this with the observation that the following two inequalities hold:

$$\begin{aligned} \sum_{t=1}^T q_t(k)f(Z_t) &\leq \sum_{t=1}^T q_t f(Z_t) + 2M\|\mathbf{q} - \mathbf{u}\|_1 \\ \Delta(\mathbf{q}(k)) &\leq \Delta(\mathbf{q}) + 2M\|\mathbf{q} - \mathbf{u}\|_1, \end{aligned}$$

shows that the event

$$\left\{ \mathbb{E}[f(Z_{T+1})|\mathbf{Z}_1^T] > \sum_{t=1}^T q_t f(Z_t) + \Delta + C\left(\epsilon + \sqrt{2 \log \log_2 2 \|\mathbf{q} - \mathbf{u}\|_1^{-1}}\right) + 4M\|\mathbf{q} - \mathbf{u}\|_1 \right\}$$

implies the following one

$$\left\{ \mathbb{E}[f(Z_{T+1})|\mathbf{Z}_1^T] > \sum_{t=1}^T q_t(k)f(Z_t) + \Delta(\mathbf{q}(k)) + C\epsilon_k \right\},$$

which completes the proof. \square

C Experiments

In Section 5, we described an algorithm benefitting from our learning guarantees based on solving the convex optimization problem (6). Due to submission time constraints, our experiments were carried out instead by solving directly problem (4) using an alternating optimization method. This is based on the observation that for a fixed \mathbf{q} , problem (4) is a simple QP over \mathbf{w} and, for a fixed \mathbf{w} , the problem reduces to an LP in \mathbf{q} . This suggests an iterative scheme where we alternate between each of these two problems.

We have compared our algorithm against a standard ARIMA model that is commonly used in practice for forecasting non-stationary time series, as well as a weighted autoregression algorithm (WAR) that solves optimization problem in (8) with \mathbf{q} tuned manually.

In our experiments, we have used three artificial datasets: ads1, ads2, ads3. For each dataset, we have generated time series with 2,000 sample points, trained on the first 1,999 points and tested on

the last point. To gain statistical significance, we repeat this procedure 1,000 times. To generate these time series the following autoregressive processes have been used:

$$\begin{aligned} \text{ads1: } & Y_t = \alpha_t Y_{t-1} + \epsilon_t, \quad \alpha_t = 1 \text{ if } t < 1800 \text{ and } -1 \text{ otherwise,} \\ \text{ads2: } & Y_t = \alpha_t Y_{t-1} + \epsilon_t, \quad \alpha_t = 0.9 - 1.8(t/2000), \\ \text{ads3: } & Y_t = \alpha_t Y_{t-1} + (1 - \alpha_t)Y_{t-2} + \epsilon_t, \quad \alpha_t = 0.9t/2000, \end{aligned}$$

where ϵ_t are independent standard Gaussian random variables.

The results of our experiments are summarized in Table 1. Observe that the results of each experiment are statistically significant using paired t -test and in each case our discrepancy-based forecaster (DBF) significantly outperforms other algorithms. Moreover, DBF has a better performance on at least 90% of individual runs in each experiment.

D Optimization Problem

In this section, we provide a detailed derivation of the optimization problem in (6) starting with optimization problem in (4). The first step is to appeal to the following chain of equalities:

$$\begin{aligned} & \min_{\mathbf{w}} \left\{ \sum_{t=1}^T q_t (\mathbf{w} \cdot \Psi(x_t) - y_t)^2 + \lambda_2 \|\mathbf{w}\|_{\mathcal{H}}^2 \right\} \\ &= \min_{\mathbf{w}} \left\{ \sum_{t=1}^T (\mathbf{w} \cdot x'_t - y'_t)^2 + \lambda_2 \|\mathbf{w}\|_{\mathcal{H}}^2 \right\} \\ &= \max_{\beta} \left\{ -\lambda_2 \sum_{t=1}^T \beta_t^2 - \sum_{s,t=1}^T \beta_s \beta_t x'_s x'_t + 2\lambda_2 \sum_{t=1}^T \beta_t y'_t \right\} \\ &= \max_{\beta} \left\{ -\lambda_2 \sum_{t=1}^T \beta_t^2 - \sum_{s,t=1}^T \beta_s \beta_t \sqrt{q_s} \sqrt{q_t} K_{s,t} + 2\lambda_2 \sum_{t=1}^T \beta_t \sqrt{q_t} y_t \right\} \\ &= \max_{\alpha} \left\{ -\lambda_2 \sum_{t=1}^T \frac{\alpha_t^2}{q_t} - \alpha^T \mathbf{K} \alpha + 2\lambda_2 \alpha^T \mathbf{Y} \right\}, \end{aligned} \tag{12}$$

where the first equality follows by substituting $x'_t = \sqrt{q_t} \Psi(x_t)$ and $y'_t = \sqrt{q_t} y_t$ the second equality uses the dual formulation of the kernel ridge regression problem and the last equality follows from the following change of variables: $\alpha_t = \sqrt{q_t} \beta_t$.

By (12), optimization problem in (4) is equivalent to the following optimization problem

$$\min_{0 \leq \mathbf{q} \leq \mathbf{1}} \left\{ \max_{\alpha} \left\{ -\lambda_2 \sum_{t=1}^T \frac{\alpha_t^2}{q_t} - \alpha^T \mathbf{K} \alpha + 2\lambda_2 \alpha^T \mathbf{Y} \right\} + \lambda_1 (\mathbf{d} \cdot \mathbf{q}) + \lambda_3 \|\mathbf{q} - \mathbf{u}\|_1 \right\}.$$

Next, we apply the change of variables $r_t = 1/q_t$, and upper bound the last term in the objective $\lambda_3 \|\mathbf{q} - \mathbf{u}\|_1$ by $\lambda_3 \|\mathbf{r} - T^2 \mathbf{u}\|_2$, where we use the fact that $|q_t - u_t| = |q_t u_t (r_t - T)|$ and Hölder's inequality. This leads to the following convex optimization problem:

$$\min_{\mathbf{r} \in \mathcal{D}} \left\{ \max_{\alpha} \left\{ -\lambda_2 \sum_{t=1}^T r_t \alpha_t^2 - \alpha^T \mathbf{K} \alpha + 2\lambda_2 \alpha^T \mathbf{Y} \right\} + \lambda_1 \sum_{t=1}^T \frac{d_t}{r_t} + \lambda_3 \|\mathbf{r} - T^2 \mathbf{u}\|_2 \right\}.$$

This optimization problem is convex, since the domain $\mathcal{D} = \{\mathbf{r}: r_t \geq 1, \forall t \in [1, T]\}$ and the first term in the objective is a maximum of convex (linear) functions of \mathbf{r} and hence is a convex function of \mathbf{r} . The last term in the objective is equivalent to a constraint $\|\mathbf{r} - T^2 \mathbf{u}\|_2 \leq \Lambda$ or $\|\mathbf{r} - T^2 \mathbf{u}\|_2^2 \leq \Lambda^2$, for some Λ . This allows us to write the optimization equivalently as

$$\min_{\mathbf{r} \in \mathcal{D}} \left\{ \max_{\alpha} \left\{ -\lambda_2 \sum_{t=1}^T r_t \alpha_t^2 - \alpha^T \mathbf{K} \alpha + 2\lambda_2 \alpha^T \mathbf{Y} \right\} + \lambda_1 \sum_{t=1}^T \frac{d_t}{r_t} + \lambda_3 \|\mathbf{r} - T^2 \mathbf{u}\|_2^2 \right\},$$

which is exactly the problem in (6).