
Beyond Least Squares: Uniform Approximation and the Hidden Cost of Misspecification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the problem of controlling worst-case errors in misspecified linear regres-
2 sion under the random design setting, where the regression function is estimated via
3 (penalized) least-squares. This setting arises naturally in value function approxima-
4 tion for bandit algorithms and reinforcement learning. Our first main contribution
5 is the observation that the amplification of the misspecification error when using
6 least-squares is governed by the *Lebesgue constant*, a classical quantity from ap-
7 proximation theory that depends on the choice of the feature subspace and the
8 covariate distribution. We also show that this dependence on the misspecification
9 error is tight for least-squares regression: in general, no method minimizing the
10 empirical squared loss can improve it substantially. As a second contribution, we
11 propose a method that augments the original feature set with auxiliary features
12 designed to reduce the error amplification. For this method we prove an oracle
13 inequality that shows that the method successfully competes with an “oracle” that
14 knows the best way of using the auxiliary features to reduce error amplification. As
15 an illustration, when the domain is a real interval and the features are monomials,
16 we prove that in the limit as $d \rightarrow \infty$, our method reduces the amplification factor
17 to $\mathcal{O}(1)$. Note that without our method, least-squares with the monomials (and in
18 fact polynomials) will suffer a worst-case error of order $\Omega(d)$ times the one of the
19 best uniform linear approximator.

20 1 Introduction

21 Value function approximation plays a central role in modern reinforcement learning (RL) and
22 contextual bandit algorithms [Sutton and Barto, 2018, Lattimore and Szepesvári, 2020]. In many such
23 settings, policies are evaluated or selected based on value estimates obtained by regressing observed
24 returns. In both cases, (penalized) linear regression—based on empirical squared loss—serves as a
25 core subroutine due to its simplicity and favorable computational properties [Ernst et al., 2005, Antos
26 et al., 2008]. A fundamental challenge arises, however, when the true value function or reward model
27 lies outside the span of the chosen features—a situation referred to as *model misspecification*. Recent
28 work by Du et al. [2020] highlighted that in this setting, the prediction error incurred by least-squares
29 regression can be *amplified* by a factor as large as \sqrt{d} , even when the misspecification error itself is
30 small and *the learner is allowed to choose the distribution of the features* (importance sampling). This
31 amplification phenomenon has since drawn significant attention in the RL and bandits communities,
32 due to its implications for the reliability of value estimation under function approximation [Lattimore
33 et al., 2020, Dong and Yang, 2023, Amortila et al., 2023, Maran et al., 2024, Amortila et al., 2024].

34 In this paper, we study the problem of controlling such worst-case errors in *misspecified linear regres-*
35 *sion* under the *random design setting*, where inputs are drawn from an unknown distribution. Our first
36 main result is a sharp characterization of how the *amplification of the misspecification error* depends

on the interaction between the sampling distribution and the feature subspace. Specifically, we show that this amplification is governed by the *Lebesgue constant*—a classical quantity in approximation theory that captures how well the 2-norm projection underlying least-squares regression projects arbitrary functions onto the span of the features. This result is a significant improvement on previous results. While previous works pointed out that with the *best* covariate distribution the misspecification amplification factor cannot be larger than \sqrt{d} regardless of the feature-map and at the same time, for some feature maps \sqrt{d} is actually the best factor, our approach reveals that the amplification factor can range from as low as 1 for favorable features. In such scenario, one can obtain significantly tighter finite-sample guarantees than previously known, which universally assumed a worst-case \sqrt{d} scaling. Moreover, we prove that this dependence is tight: no estimator based on least squares can substantially improve upon this bound in general.

Motivated by these insights, we propose a method for *reducing the misspecification error amplification* by augmenting the original feature set with extra features and then using a weighted ridge regression approach to regularize the corresponding projection operator. As an illustration of this idea, we show that when the domain is an interval and the base and extra features are monomials, our method reduces the amplification factor to 1 asymptotically as $d \rightarrow \infty$. In contrast, standard least squares remains susceptible to arbitrarily large worst-case errors for the same setting.

2 Problem Formulation

Our goal is to learn a linear predictor that enjoys *uniform* accuracy over the whole input space, even when the linear model is misspecified. We first detail the statistical setting, introduce the standing assumptions and define the performance criterion that will be used in the rest of the paper.

The learner receives a dataset $((x_t, y_t))_{t=1}^n$, where $x_t \in \mathcal{X}$ and $y_t \in \mathbb{R}$. Each x_t gives rise to a response $y_t = f(x_t) + \eta_t$ where $f : \mathcal{X} \rightarrow \mathbb{R}$ is an unknown function and η_t is a noise variable.

Assumption 1 (Sub-Gaussian Noise). *The noise variables $(\eta_t)_{t=1}^n$ are independent, centered and σ -sub-Gaussian, meaning that, for every $\lambda \in \mathbb{R}$, $\mathbb{E}[\exp(\lambda\eta_1)] \leq \exp(\sigma^2\lambda^2/2)$. The noise variables are independent of the inputs $(x_t)_{t=1}^n$.*

About $(x_t)_{t=1}^n$, we are going to make the following assumption.

Assumption 2 (Random design). *Samples $(x_t)_{t=1}^n$ are drawn i.i.d. from a probability distribution μ over \mathcal{X} (unknown to the learner).*

We are interested in the problem of *linear function approximation*, when the learner is given some feature map $\varphi_d : \mathcal{X} \rightarrow \mathbb{R}^d$ and aims to approximate f using $f_\theta(\cdot) = \varphi_d(\cdot)^\top \theta$ by selecting some $\theta \in \mathbb{R}^d$ based on the data available to it. In the rest of the paper, we use the short-hand $\varphi_i(x)$ for the i -th coordinate of $\varphi_d(x)$ and index data points by $t = 1, \dots, n$. Differently from most of the literature about this setting, motivated by the applications mentioned earlier, the performance is going to be evaluated via the uniform, or maximum-norm, which for a function $g : \mathcal{X} \rightarrow \mathbb{R}$ is defined via $\|g\|_\infty = \sup_{x \in \mathcal{X}} |g(x)|$. We let $L^\infty(\mathcal{X})$ denote the set of functions with finite maximum norm. In what follows, we assume that both f and our features φ_i belong to this set. For $f \in L^\infty(\mathcal{X})$ and $\theta \in \mathbb{R}^d$ we let

$$\mathcal{E}_\infty(\theta, f) := \|f_\theta - f\|_\infty, \quad \mathcal{E}_\infty(f) := \inf_{\theta \in \mathbb{R}^d} \mathcal{E}_\infty(\theta, f).$$

Thus, $\mathcal{E}_\infty(\theta, f)$ is the maximum error suffered when f is approximated using f_θ , while $\mathcal{E}_\infty(f)$ is the smallest possible value for this error; its value is *unknown* to the learner. When $\mathcal{E}_\infty(f) > 0$, we say that the problem of estimating f is *misspecified* and the error $\mathcal{E}_\infty(f)$ is known as the *misspecification error*. In the next section we will be interested in investigating how the error $\mathcal{E}_\infty(\hat{\theta}_n, f)$ behaves when $\hat{\theta}_n$ is given by *ordinary least-squares* (OLS) estimate:

$$\hat{\theta}_{n,\text{OLS}} = \arg \min_{\theta \in \mathbb{R}^d} \sum_{t=1}^n (y_t - f_\theta(x_t))^2.$$

71 3 Characterizing the behavior of OLS and Ridge Regression

72 Let $\mathcal{F}_d = \{f_\theta : \theta \in \mathbb{R}^d\}$ denote the subspace of $L^\infty(\mathcal{X})$ spanned by the basis functions underlying
 73 the feature-map φ_d . As $n \rightarrow \infty$, $\hat{f}_{\hat{\theta}_{n,\text{OLS}}}$ is known to converge to

$$\Pi_{d,\mu} f := \arg \min_{g \in \mathcal{F}_d} \|g - f\|_\mu^2 \quad (1)$$

74 with probability one (see Györfi et al. [2006]). Here, we define $\|\cdot\|_\mu^2$ to stand for the $L^2(\mu)$ -norm:
 75 For $g : \mathcal{X} \rightarrow \mathbb{R}$ measurable, $\|g\|_\mu^2 = \int_{\mathcal{X}} g^2(x) \mu(dx)$. Since μ is a probability measure, we have
 76 $\|g\|_\mu^2 \leq \|g\|_\infty^2$. The map in Eq. (1) is known to be *projection* onto \mathcal{F}_d : $\Pi_{d,\mu}$ is linear, idempotent and
 77 for all $f \in \mathcal{F}_d$, $\Pi_{d,\mu} f = f$ holds. Moreover, it is non-expansive in the $L^2(\mu)$ -norm.

78 By continuity, the previous comment on the convergence of the OLS estimate implies that
 79 $\lim_{n \rightarrow \infty} \mathcal{E}_\infty(\hat{\theta}_{n,\text{OLS}}, f) = \|\Pi_{d,\mu} f - f\|_\infty$. The first question then is thus how large $\|\Pi_{d,\mu} f - f\|_\infty$
 80 can be relative to $\mathcal{E}_\infty(f)$, or, in other words, by how much will the misspecification error $\mathcal{E}_\infty(f)$ be
 81 enlarged if one uses the linear projection of f to \mathcal{F}_d . The following definition will be useful:

82 **Definition 1** (Lebesgue constant). *Let $P : L^\infty(\mathcal{X}) \rightarrow L^\infty(\mathcal{X})$ be a linear operator. Then, the*
 83 *L^∞ -norm of P is called the Lebesgue constant associated with P :*

$$\Lambda(P) := (\|P\|_\infty =) \sup_{f \in L^\infty(\mathcal{X}) : \|f\|_\infty \leq 1} \|Pf\|_\infty.$$

84 The following result holds (see Proposition 4.1 from Chapter 2 of DeVore and Lorentz [1993]):

85 **Lemma 1** (Lebesgue's lemma). *Let $P : L^\infty(\mathcal{X}) \rightarrow \mathcal{F}_d$ be a linear map such that P is an identity*
 86 *on \mathcal{F}_d . In particular, assume that for any $f \in \mathcal{F}_d$, $Pf = f$. Then, for any $f \in L^\infty(\mathcal{X})$,*

$$\|f - Pf\|_\infty \leq (1 + \Lambda(P)) \inf_{g \in \mathcal{F}_d} \|f - g\|_\infty.$$

87 Since the Lebesgue constant of our projection operators will be frequently needed, to minimize
 88 clutter, we introduce the shorthand for them:

$$\Lambda_{d,\mu} := \Lambda(\Pi_{d,\mu}).$$

89 With the help of this notation, from Lemma 1 we thus have

$$\|f - \Pi_{d,\mu} f\|_\infty \leq \Lambda_{d,\mu} \mathcal{E}_\infty(f). \quad (2)$$

90 It is easy to see that $\Lambda_{d,\mu} \geq 1$: just take any $f \in \mathcal{F}_d$ such that $\|f\|_\infty = 1$ (such a function exist).
 91 Then $\Pi_{d,\mu} f = f$, which gives the result. Unfortunately, there is no upper limit on how large $\Lambda_{d,\mu}$
 92 can be in general. What is more, the bound in Lemma 1 is essentially tight:

93 **Theorem 3.** *For any $\varepsilon > 0$ there exist $f \in L^\infty(\mathcal{X})$ such that*

$$\|f - \Pi_{d,\mu} f\|_\infty \geq (\Lambda_{d,\mu} - 1 - \varepsilon) \mathcal{E}_\infty(f).$$

94 For the proof, see Appendix C.1. Because of the last result, we expect that any bound on $\mathcal{E}_\infty(\hat{\theta}_{n,\text{OLS}}, f)$
 95 where $\hat{\theta}_n$ is estimated from data will involve $\Lambda_{d,\mu} \mathcal{E}_\infty(f)$. Our main result of this section is indeed of
 96 this form. To state the result, let $(\bar{\varphi}_i)_{1 \leq i \leq d}$ be the orthonormal basis in $L^2(\mu)$ of \mathcal{F}_d given by the
 97 Gram-Schmidt procedure on the original features. We call $\bar{\varphi}_d(x) = (\bar{\varphi}_1(x), \dots, \bar{\varphi}_d(x))^\top$ and define
 98 $\bar{\varphi}_{d,2} = \sup_{x \in \mathcal{X}} \|\bar{\varphi}_d(x)\|_2$.

99 **Theorem 4.** *Let \mathcal{X} be finite. Let Assumptions 1 and 2 hold. Then, for any n positive integer and real*
 100 *$0 < \delta \leq 1/3$ such that $n \geq 20\bar{\varphi}_{d,2}^2 \log(d/\delta)$, letting $\hat{\theta}_{n,\text{OLS}}$ be the parameter vector returned by OLS,*
 101 *with probability at least $1 - 3\delta$,*

$$\begin{aligned} \mathcal{E}_\infty(\hat{\theta}_{n,\text{OLS}}, f) &\leq (1 + \Lambda_{d,\mu}) \mathcal{E}_\infty(f) + 3(\sigma + \Lambda_{d,\mu} \mathcal{E}_\infty(f)) \bar{\varphi}_{d,2} \sqrt{\frac{\log(|\mathcal{X}|/\delta)}{n}} \\ &\quad + \frac{\text{poly}(d, \bar{\varphi}_{d,2}, \Lambda_{d,\mu} \mathcal{E}_\infty(f))}{n}. \end{aligned}$$

The first quantity in the bound is the same as in Eq. (2) and this is the quantity that accounts for the gap between $\Pi_{d,\mu}f$ and f . The other terms bound the finite sample error. We gave the result for \mathcal{X} finite only for simplicity. Indeed, the result can be easily extended by means of a simple covering argument. For example, if $\mathcal{X} = [-1, 1]$ and the features are Lipschitz continuous with constant L_φ , we can achieve uniform error over \mathcal{X} by making a ε/L_φ -covering of \mathcal{X} . In this way, the bound increases just by a factor $\propto \log(L_\varphi/\varepsilon)$. Another possibility, which sometimes can lead to tighter results, is to cover $\{\varphi(x) : x \in \mathcal{X}\} \subset \mathbb{R}^d$.

In addition to $\Lambda_{d,\mu}$, another constant that depends on the feature map is $\bar{\varphi}_{d,2}$. As it turns out, this value “hides” the dimension d . In particular, we show that regardless of the feature map, $\bar{\varphi}_{d,2} \geq \sqrt{d}$ (see Proposition 19 in the appendix).

The terms and their scaling with the relevant quantities are as expected. We have already discussed the first term and argued that it cannot be significantly improved for OLS. The second (and the lower order third) term accounts for the sampling errors. In particular, the effect of the noise is shown through the term involving σ . The next term, which also involves the Lebesgue constant and $\mathcal{E}_\infty(f)$ accounts for the random design sampling error.

Below we show that when an *a priori* upper bound ε on $\mathcal{E}_\infty(f)$ is available (as can be the case in certain numerical applications when the target function belongs to some known class of functions, such as a smoothness class), we can obtain an empirical bound that has the potential to significantly reduce the terms of the bound shown in the last result.

A uniform, semi-empirical bound Define $\mu_n = \frac{1}{n} \sum_{t=1}^n \delta_{x_t}$ to be the empirical measure underlying the inputs (x_1, \dots, x_n) . We are interested in bounding the uniform error of the OLS estimate via empirical quantities. In particular, we will use the Lebesgue constant associated with the projection operator Π_{d,μ_n} corresponding to μ_n . Because of this, we can also remove the assumption that $(x_t)_t$ is sampled from μ ; in fact, we will not need any assumptions concerning how $(x_t)_t$ are selected.

The operator Π_{d,μ_n} takes the form

$$\Pi_{d,\mu_n}f(\cdot) := \varphi_d(\cdot)^\top (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{f} \quad \mathbf{f} := [f(x_1), \dots, f(x_t), \dots, f(x_n)]^\top,$$

where Φ is the $n \times d$ matrix storing, as rows, the features corresponding to every $\varphi_d(x_t)$. As before, we define $(\hat{\varphi}_i)_{1 \leq i \leq d}$ the orthonormal basis in $L^2(\mu_n)$ of \mathcal{F}_d . We let $\hat{\varphi}_d(x) = (\hat{\varphi}_1(x), \dots, \hat{\varphi}_d(x))^\top$ and define $\hat{\varphi}_{d,2} = \sup_{x \in \mathcal{X}} \|\hat{\varphi}_d(x)\|_2$.

Theorem 5. *Let Assumption 1 hold. Then, for any fixed $\delta > 0$, with probability at least $1 - \delta$,*

$$\mathcal{E}_\infty(\hat{\theta}_{n,\text{OLS}}) \leq (1 + \Lambda_{d,\mu_n})\mathcal{E}_\infty(f) + \frac{\sigma \hat{\varphi}_{d,2} \sqrt{2 \log(2\mathcal{X}/\delta)}}{\sqrt{n}}.$$

Compared to Theorem 4, we both removed Assumption 2 and the lower-order term $n^{-1} \text{poly}(d, \bar{\varphi}_{d,2}, \Lambda_{d,\mu} \mathcal{E}_\infty(f))$. At the same time, the Lebesgue constant $\Lambda_{d,\mu}$ is replaced with Λ_{d,μ_n} , which may be smaller or larger than $\Lambda_{d,\mu}$. When \mathcal{X} is finite then Λ_{d,μ_n} can be calculated in $O(n|\mathcal{X}|)$ time (it is just a matrix maximum norm).

If choosing the points (x_1, \dots, x_n) is an option, one may attempt to optimize the bound. Here, besides the term Λ_{d,μ_n} , $\hat{\varphi}_{d,2}$ also hides μ_n . In experimental optimal design, the G -optimal design is defined as the one that minimizes $\hat{\varphi}_{d,2}$ by choosing an appropriate distribution μ_n . Here, one of the main results is that for $n = \Omega(d)$ (or slightly larger), one can find μ_n such that $\hat{\varphi}_{d,2} = O(\sqrt{d})$ and this is the best possible value Kiefer and Wolfowitz [1960].

Under the assumption that μ_n is an optimal design, we can compare our result with Proposition 5.1 from Lattimore et al. [2020] (see equation (2) and the corresponding bound in high probability). Rephrasing their proposition in our notation, we get, roughly that, if μ is an optimal design for φ_d , then for $\sigma = 1$,

$$\mathcal{E}_\infty(\hat{\theta}_{n,\text{OLS}}, f) \leq \mathcal{O} \left(\sqrt{d} \mathcal{E}_\infty(f) + \sqrt{\frac{d \log(|\mathcal{X}|/\delta)}{n}} \right).$$

This result is a particular case of our Theorem 4: indeed $\hat{\varphi}_{d,2} = \sqrt{d}$ as we are using an optimal design, while it is not hard to see that $\Lambda_{d,\mu_n} \leq \hat{\varphi}_{d,2}$ holds without any assumptions. In this bound,

Basis functions	μ	$\Lambda_{d,\mu}$	Source	Note
Polynomial	uniform on regular d -grid	$\Omega(2^d)$	[Quarteroni et al., 2010]	
Polynomials	uniform	$\Theta(d)$	DeVore and Lorentz [1993]	$\bar{\varphi}_{d,2} \approx d$
Fourier	uniform	$O(\log(d))$	[Katznelson, 2004, p.59, Exercise 1]	
Continuous B-splines	uniform	$O(1)$	Huang [2003]	
Wavelets	uniform	$O(1)$	Chen and Christensen [2013]	

Table 1: Examples of Lebesgue constants. Domain is $\mathcal{X} = [-1, 1]$.

for large n , $\sqrt{d}\mathcal{E}_\infty(f)$ is the dominant term. Therefore, our tighter bound with the Lebesgue constant achieves a better result whenever the Lebesgue constant is significantly smaller than \sqrt{d} . As an example, of a feature-map with a small Lebesgue constant, consider any partitioning $(\mathcal{X}_i)_i$ of \mathcal{X} and set φ_i to be the indicator of part \mathcal{X}_i , $i = 1, \dots, d$. Then, the Lebesgue constant of φ_d is 1 regardless of the choice of μ . Note that with some extra work we can extract a more refined result from the proof of Proposition 5.1 when \sqrt{d} on the right-hand-side above is essentially replaced by $\hat{\varphi}_{d,2}$.

3.1 The Lebesgue constant: properties and particular cases

While as noted earlier $\Lambda_{d,\mu} \leq \bar{\varphi}_{d,2}$ always hold, since $\bar{\varphi}_{d,2}$ is never lower than \sqrt{d} , to get a better understanding of the Lebesgue constants associated with specific feature maps, we need to resort to feature-map dependent analysis. For many of the classical feature-maps, to under we have to resort to some feature map-dependent analysis. In the following, we enumerate few well-known/novel results in Table 1. As seen in the table, the range of values is quite large. Notably, polynomials with regular d -grids, show the worst-behavior (though this is somewhat unusual since here for every d one uses a different measure). Yet, polynomials with the uniform measure still exhibit quite big Lebesgue constants. Perhaps surprisingly, when switching to Fourier series, the Lebesgue constant decreases to $O(\log d)$. As such, if the convergence of a Fourier series to the target function is fast enough, there is little to no incentive to go beyond the L^2 -projections, or least-squares. It is interesting to note that some researchers have empirically found the Fourier series as a good “general” basis to be used in reinforcement learning [Konidaris et al., 2011]. This raises the hypothesis that this is primarily due to the reasonable error extrapolation properties of Fourier series, as attested by its relatively slow-growing Lebesgue constant. Finally, we mention that perhaps unsurprisingly localized basis functions such as wavelets and B-splines have Lebesgue constants that are constant independently of the number of basis functions used. As such, when uniform approximation is important, it seems that there are good reasons to prefer these systems. We speculate that tile-coding, which is itself localized and which is quite popular in reinforcement learning also shares the good extrapolation properties of these localized systems of basis functions.

One weakness of the above results is that they are dependent on the choice of the sampling distribution μ . The following result shows that the Lebesgue constant changes gradually as one moves from one distribution to another, provided that the overlap between the two distributions is well-controlled:

Proposition 6. *Let μ, ν be two probability distributions on the discrete set \mathcal{X} such that for all $x \in \mathcal{X}$, $C \geq \frac{\mu(x)}{\nu(x)} \geq c > 0$. Then, $\Lambda_{d,\mu} \leq \frac{C}{c} \Lambda_{d,\nu}$.*

4 Regularized estimators

The previous Theorem 3 shows that, whatever the feature map, if we use the OLS estimator, the error is forced to scale with the misspecification multiplied by the Lebesgue constant. This is not a matter of overfitting, as the bound holds for infinite data; still, the problem is related to the LS solution becoming “too big” for some choices of f . Therefore, an idea would be to enforce a regularization on the LS loss, to limit the magnitude of the estimated function. In the next theorem, we show that the standard Ridge Regression approach is ineffective, even when knowing the orthonormal basis $\bar{\varphi}_d$.

179 **Theorem 7.** Let $\hat{\theta}_{n,\text{RIDGE}}$ the output of λ -ridge regression. For any feature map
 180 $\varphi_d(\cdot) : \mathcal{X} \rightarrow \mathbb{R}^d$ there is $f \in L^\infty(\mathcal{X})$ such that, for infinite data $\mathcal{E}_\infty(\hat{\theta}_{\infty,\text{RIDGE}}) =$
 181 $\Omega\left(\max\left\{(\Lambda_{d,\mu} - 2\lambda)\mathcal{E}_\infty(f), \frac{\lambda}{\lambda+1}\right\}\right).$

182 This result tells that the Ridge regularization is ineffective: if we take $\lambda \approx \Lambda_{d,\mu}/2$ the second term
 183 is close to one (we do not go to zero even for $\mathcal{E}_\infty(f) \rightarrow 0$), if we do not, we get the same lower
 184 bound of OLS. Crucially, this phenomenon persists even in the infinite data regime, indicating that
 185 it is not merely a sample size issue, but a geometric defect of the projection operator itself. Other
 186 regularization techniques like Cross-Validation and Early Stopping [Ghojogh and Crowley, 2019],
 187 that are designed for dealing with overfitting, cannot overcome this result, as they aim at minimizing
 188 the test error MSE, which is achieved by OLS for infinite data.

189 Let us dig deeper into the reason behind the failure of ridge regression. The proof of Theorem 7,
 190 builds on the fact that the corresponding operator $\Pi_{d,\mu}^{\text{Ridge}}$ can be written in the following form

$$\Pi_{d,\mu}^{\text{Ridge}} f(x) = \sum_{i=1}^d \alpha \bar{\varphi}_i(x) \int_{\mathcal{X}} \bar{\varphi}_i(z) f(z) d\mu(z) \quad \alpha = \frac{1}{1 + \lambda}. \quad (3)$$

191 In fact, this is *not* a projection operator. Indeed, this does to corresponds to the identity when applied
 192 over \mathcal{F}_d . F.e. taking $f(\cdot) = \bar{\varphi}_1(\cdot)$, we get $\bar{\varphi}_1(x) - \Pi_{d,\mu}^{\text{Ridge}} \bar{\varphi}_1(x) = \frac{\lambda}{1+\lambda} \bar{\varphi}_1(x)$. Indeed, in order to
 193 scale with $\mathcal{E}_\infty(f)$ we have to ensure that any function in \mathcal{F}_d is kept fixed by the operator associated
 194 to our estimator for θ . Still, keeping $\alpha = 1$, that means $\lambda = 0$ means going back to OLS.

195 **Extending the feature map** Let us assume to add one more feature to $\varphi_d(\cdot)$, which we now call
 196 $\varphi_{d+1}(\cdot)$. Indeed, regardless of the nature of the feature that we add, $\mathcal{E}_\infty(f)$ can only decrease, as we
 197 are taking the infimum over a larger set. Intuition, together with the results from Section 3.1, would
 198 suggest that the corresponding Lebesgue constant, that passes from $\Lambda_{d,\mu} \rightarrow \Lambda_{d+1,\mu}$ could only get
 199 bigger. Surprisingly, this is not the case: there are examples of feature maps such that adding one
 200 feature may correspond to the Lebesgue constant getting *smaller*.

201 To formalize this idea, let us assume to expand the original feature map with $D - d$ different functions,
 202 for some integer $D > d$, so that the full feature map can be written as $\varphi_D(\cdot) := [\varphi_d(\cdot), \varphi'_{D-d}(\cdot)]$.
 203 Even if the added features $\varphi'_{D-d}(\cdot)$ can be arbitrary, we argue that in many real problems there is a
 204 way to select them in a reasonable way. For example, in all the examples mentioned in Section 3.1
 205 (Fourier features, polynomials, splines...) the sequence considered can be easily extended up to
 206 infinity by enlarging the maximum degree. We call, as before, $\bar{\varphi}_D(\cdot)$ the feature map that we obtain
 207 from the Gram-Schmidt procedure on the basis $\varphi_D(\cdot)$ with measure $\mu(\cdot)$. We consider the operator
 208 associated to *weighted* ridge regression over this extended basis, which generalizes Eq. (3). This
 209 writes, for any sequence of D weights λ_i in $[0, \infty)$ as:

$$\Pi_{\alpha,\mu}^{\text{Ridge}} f(x) := \sum_{i=1}^D \alpha_i \bar{\varphi}_i(x) \int_{\mathcal{X}} \bar{\varphi}_i(z) f(z) d\mu(z). \quad \alpha_i = \frac{1}{1 + \lambda_i}. \quad (4)$$

210 In fact, it shall be proved that this operator is the minimizer of the weighted Ridge loss, when adding a
 211 penalization λ_i on each component $\bar{\varphi}_i(x)$. Not every value for the sequence α is meaningful: in fact,
 212 the original components $i \leq d$ must *not* be penalized. Moreover, as the ridge penalization of each
 213 component ranges in $[0, +\infty)$, the corresponding $\alpha_i \in (0, 1]$. These two constrains are formalized in
 214 the following set:

$$\mathcal{A}_d^D := \{\alpha \in [0, 1]^D : \forall i \leq d, \alpha_i = 1\}, \quad (5)$$

215 which we call the set of *attenuation parameters*.

216 4.1 Weighted ridge estimator and the Oracle Operator

217 Each of these operators $\Pi_{\alpha,\mu}^{\text{Ridge}}$, as long as $\alpha \in \mathcal{A}_d^D$, maintains every element of \mathcal{F}_d (the span of the
 218 original feature map $\varphi_d(\cdot)$) as fixed point. We call $\Lambda_{\alpha,\mu}$ the Lebesgue constant of the corresponding
 219 operator. The following result holds, generalizing Lebesgue's Lemma.

Proposition 8. Let $\alpha \in \mathcal{A}_d^D$ (equation (5)) and $\Pi_{\alpha,\mu}^{\text{Ridge}}$ be defined according to equation (4). Then,

$$\|f(\cdot) - \Pi_{\alpha,\mu}^{\text{Ridge}} f(\cdot)\|_\infty \leq (1 + \Lambda_{\alpha,\mu}) \mathcal{E}_\infty(f).$$

Crucial for the proof of Proposition 8 is the definition of \mathcal{A}_d^D , ensuring that the attenuation factor α_i is one for $i \leq d$, so that the projection on the original features is maintained exactly as it is. This ensures that the original features are *not* penalized by the Ridge regularization. Under our point of view, to minimize the expansion factor corresponds to using the value $\alpha \in \mathcal{A}_d^D$ which has the minimal Lebesgue constant. We call this value ORACLE:

$$\alpha_\mu^{\text{Oracle}} := \arg \min_{\alpha \in \mathcal{A}_d^D} \Lambda_{\alpha,\mu} \quad \Lambda_\mu^{\text{Oracle}} := \min_{\alpha \in \mathcal{A}_d^D} \Lambda_{\alpha,\mu}.$$

Unfortunately, $\alpha_\mu^{\text{Oracle}}$ is unknown to the learner, as it depends on the unknown distribution μ . Our questions for the rest of this section are the following:

- **Q1** Can we design a finite sample estimator that, for fixed $\alpha \in \mathcal{A}_d^D$, asymptotically, scales as $\Lambda_{\alpha,\mu}$?
- **Q2** Can we design a finite sample estimator that, asymptotically, scales as $\Lambda_\mu^{\text{Oracle}}$?

Q1 To answer both the previous questions, we start by generalizing Theorem 5 to the case of regularization with the desired parameter α . Indeed, even being μ unknown, we can define an empirical counterpart of the operator defined in Eq. (4) by means of the empirical measure $\mu_n(\cdot)$. In fact, recalling that $\widehat{\varphi}_D(\cdot)$ is the feature map obtained by orthogonalizing $\varphi_D(\cdot)$ w.r.t. $\mu_n(\cdot)$, we have

$$\Pi_{\alpha,\mu_n}^{\text{Ridge}} f(\cdot) := \sum_{i=1}^D \alpha_i \widehat{\varphi}_i(x) \frac{1}{n} \sum_{t=1}^n \widehat{\varphi}_i(x_t) f(x_t). \quad (6)$$

This operator 1) takes as argument only the evaluations of $f(\cdot)$ at x_t and 2) Outputs a linear combination of the features $\widehat{\varphi}_i(x)$. Thanks to the first point, we can estimate $\Pi_{\alpha,\mu_n}^{\text{Ridge}} f(\cdot)$ by our noisy samples by replacing $f(x_t)$ with y_t . Thanks to the second one, there exists an estimator $\widehat{\theta}$ such that the result of the operator is written as $\varphi_D(\cdot)^\top \widehat{\theta}$. Calling R_n the triangular matrix such that $\varphi_D(\cdot)^\top = \widehat{\varphi}_D(\cdot)^\top R_n$ (the matrix corresponding to Gram Schmidt procedure), we call

$$\widehat{\theta}_{n,\alpha} := R_n^{-1} I_\alpha \frac{1}{n} \sum_{t=1}^n \widehat{\varphi}_D(x_t) y_t, \quad (7)$$

where $I_\alpha = \text{diag}(\alpha)$ takes into account the regularization, we can prove the following result.

Theorem 9. Let assumption 1 hold. Then, for any $\delta > 0$, with probability $1 - \delta$,

$$\mathcal{E}_\infty(\widehat{\theta}_{n,\alpha}) \leq (1 + \Lambda_{\alpha,\mu_n}) \mathcal{E}_\infty(f) + \frac{\sigma \widehat{\varphi}_{2,D} \sqrt{2 \log(2\mathcal{X}/\delta)}}{\sqrt{n}}.$$

Therefore, our estimator $\widehat{\theta}_{n,\alpha}$ is able to scale with Λ_{α,μ_n} . To answer question **Q1** completely, we just need to show that, for large enough n , we can replace the Lebesgue constant of μ_n with the one of μ . This is done in the following proposition.

Proposition 10. Under assumption 2 we have, with probability $1 - \delta$ for every $\alpha \in \mathcal{A}_d^D$ at the same time, $|\widehat{\varphi}_{D,2} - \overline{\varphi}_{D,2}| \leq \tilde{O}(\overline{\varphi}_{D,2}^2 \sqrt{\log(1/\delta)/n})$, and $|\Lambda_{\alpha,\mu_n} - \Lambda_{\alpha,\mu}| \leq \tilde{O}\left(\frac{\sqrt{d} \overline{\varphi}_{D,2}^2 \sqrt{\log(1/\delta)}}{\sqrt{n}} + \frac{\sqrt{d} \overline{\varphi}_{D,2}^3 \log(1/\delta)}{n}\right)$.¹

¹The statement of this theorem is slightly different from the one in the main paper of the submission, as we have made the orders of magnitude more precise.

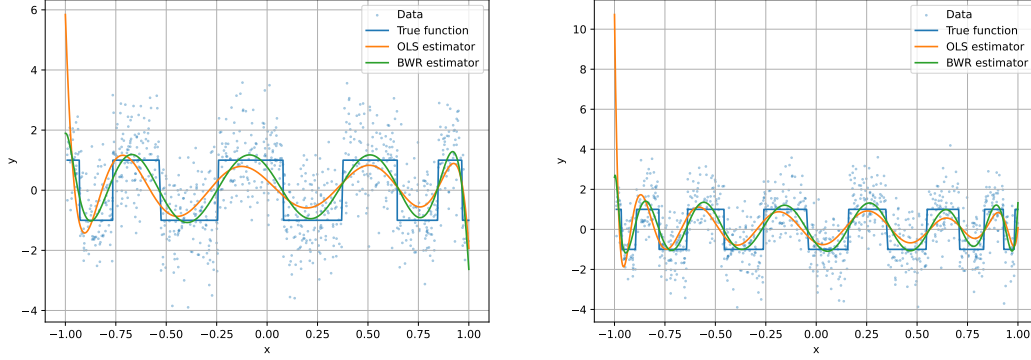


Figure 1: Comparison between the OLS estimator and the BWR estimator using polynomial features on $[-1, 1]$, with a basis of length $d = 10$ (left) and $d = 15$ (right). Even if the true function is bounded, OLS suffers from large oscillations near the boundaries due to the high Lebesgue constant. In contrast, BWR achieves a much more uniform approximation error across the domain by effectively controlling the amplification effect.

Q2 For this more challenging result, we have to optimize the value of α in order to converge to the one of the oracle, even not knowing the true distribution $\mu(\cdot)$. Our strategy, also for this point, is to work with we can compute: the operator $\Pi_{\alpha, \mu_n}^{\text{Ridge}}$ (equation (6)) and its Lebesgue constant Λ_{α, μ_n} . One observation is key for this goal: the Lebesgue constant is convex in α .

Proposition 11. The function $J : \mathcal{A}_d^D \rightarrow (0, +\infty)$ given by $J(\alpha) := \Lambda_{\alpha, \mu_n}$ is convex in α .

This result allows us to provably arrive to one minimizer of the Lebesgue constant in a finite number of iteration. The idea is what follows: we start from any $\alpha \in \mathcal{A}_d^D$ and update it iteratively with the sub-Gradient method until convergence. This algorithm is well-known [Boyd et al. \[2003\]](#), but, for completeness, we include it; see Algorithm 1 in appendix D.4.

Theorem 12. Fix $\epsilon > 0$. Algorithm 1, after a number of iterations $I = \tilde{\mathcal{O}}(\epsilon^{-2} \hat{\varphi}_{2,D}^2 (D - d))$ outputs $\alpha^{(I)} \in \mathcal{A}_d^D$ such that $J(\alpha^{(I)}) \leq \inf_{\alpha \in \mathcal{A}_d^D} J(\alpha) + \epsilon$.

By definition of $J(\cdot)$, the former result entails that $\alpha^{(I)}$ is close to be a minimizer of Λ_{α, μ_n} . To finally answer **Q2**, we define the following estimator, which corresponds to equation (7) for $\alpha^{(I)}$:

$$\hat{\theta}_{n, \text{BWR}} := R_n^{-1} I_{\alpha^{(I)}} \frac{1}{n} \sum_{t=1}^n \hat{\varphi}_D(x_t) y_t. \quad (8)$$

The estimator is called BWR, which stands for "Best Weighted Regularizer". We close this section with its performance guarantee.

Theorem 13. Let Assumptions 1 and 2 hold and fix $\delta > 0$. Then, with probability $1 - \delta$,

$$\mathcal{E}_{\infty}(\hat{\theta}_{n, \text{BWR}}) \leq (1 + \Lambda_{\mu}^{\text{Oracle}}) \mathcal{E}_{\infty}(f) + \tilde{\mathcal{O}} \left(\frac{\bar{\varphi}_{2,D} \sqrt{D \log(|\mathcal{X}|/\delta)}}{\sqrt{n}} + \frac{\bar{\varphi}_{2,D}^2 \log(|\mathcal{X}|/\delta)}{n} \right).$$

This oracle inequality answer also **Q2** in a positive way: our estimator is asymptotically able to compete with the Oracle Lebesgue constant. We close this paper with a case study of wide interest where this property allows $\hat{\theta}_{n, \text{BWR}}$ to get a much better result than the one of OLS.

5 Case study: polynomial basis

The method introduced in the previous section aims at reducing the amplification of the misspecification error by controlling the Lebesgue constant. While it applies broadly, its impact is best illustrated in settings where standard estimators suffer from poor uniform behavior. One such setting is the

classical case where the feature map φ_d consists of the first d monomials $\{1, x, x^2, \dots, x^{d-1}\}$ over a compact interval. We consider the scenario where $\mathcal{X} = [-1, 1]$ and the data-generating distribution μ is the uniform distribution on this interval. Even in this most favorable case, the Lebesgue constant associated with the polynomial basis grows linearly with the degree: $\Lambda_{d,\mu} \approx d$. This leads to a worst-case uniform error for the OLS estimator that scales as $\mathcal{O}(d \cdot \mathcal{E}_\infty(f))$, which can be arbitrarily large even when the misspecification bias $\mathcal{E}_\infty(f)$ is small.

In contrast, the BWR estimator described at the end in Section 4 augments the feature space with additional monomials and optimizes an attenuation vector to minimize the empirical Lebesgue constant. The result is a projection operator that preserves the behavior of the original features while drastically reducing the amplification of the uniform error. Theoretically, this allows reducing the amplification factor from $\mathcal{O}(d)$ down to $\mathcal{O}(1)$, as the following theorem shows.

Theorem 14. *Let $\mu(\cdot) = \mathcal{U}([-1, 1])$. There is a constant C independent of d such that, for $D = 2d$ and $\varphi_d(x) = [1, \dots, x^{d-1}]$, $\varphi_D(x) = [1, \dots, x^{2d-1}]$, we have $\Lambda_\mu^{Oracle} \leq C$.*

This improvement is evident in the numerical simulation shown in Fig. 1, where we compare the OLS estimator to the BWR estimator on synthetic data. While OLS exhibits large oscillations near the boundary of the interval—a manifestation of the classical Runge phenomenon—BWR remains stable across the domain and achieves significantly smaller uniform error. Despite both estimators using the same base features, the control of the Lebesgue constant yields a qualitative and quantitative advantage for BWR.

The above simulations visually demonstrate how the amplification factor is exacerbated by increasing d . We complement them with an asymptotic result that shows how heavy this factor is, even for a function such that $\mathcal{E}_\infty(f) \xrightarrow{d} 0$. In fact, there exist a bounded function that can be uniformly approximated with polynomial features but such that the OLS estimator diverges with uniform error roughly of order $\Omega(d)$.

Proposition 15. *Fix $\gamma > 0$. There is a function $f : [-1, 1] \rightarrow \mathbb{R}$ such that, $\mathcal{E}_\infty(f) \xrightarrow{d} 0$ and under assumptions 1 and 2 for $\mu = \mathcal{U}([-1, 1])$, with probability one,*

$$\lim_{d \rightarrow \infty} \lim_{n \rightarrow \infty} \|f(\cdot) - \varphi_d(\cdot)^\top \hat{\theta}_{n,BWR}\|_\infty = 0 \quad \lim_{n \rightarrow \infty} \|f(\cdot) - \varphi_d(\cdot)^\top \hat{\theta}_{n,OLS}\|_\infty \gtrsim d^{1-\gamma}.$$

6 Related works

The problem we deal in this paper, while motivated by Online Learning applications has roots in several different fields. Not just Mathematical Analysis and Fourier Series, but also Econometrics and Online Learning. Here, we give a short overview of the main papers, leaving an extended version for the appendix A.

In mathematical analysis the problem of projecting onto a linear subspace of $L^\infty(\mathcal{X})$ in a way that minimizes the uniform error have always been of great interest. Several results about orthogonal polynomials Szegő [1939] or Fourier Series Katznelson [2004] approximation have this goal. More recently, Kobos and Lewicki [2024] proposed an approach for general feature maps. Passing to the case when an unknown function is estimated via noisy samples, there is a line of research (Newey [1997], Belloni et al. [2015], and Li and Liao [2020]) that studies the properties of *pointwise* estimators based on LS. The latter can be naturally adapted to achieve a uniform convergence guarantee. A similar problem was recently studied, in a totally different context, by Online Learning papers (see Du et al. [2020], Lattimore et al. [2020], Maran et al. [2024], Dong and Yang [2023], Amortila et al. [2024]) under the name of *misspecified linear function approximation*.

The specific technique that we use in section 4 is inspired by an old method for regularizing Fourier series [de la Vallée Poussin, 1918, De La Vallée Poussin et al., 1919]. The technique he invented is still studied today in numerical mathematics [Németh, 2016, Themistoclakis and Van Barel, 2017, Occorsio and Themistoclakis, 2025].

7 Conclusion

We investigated the problem of uniform error control in misspecified linear regression under the random design setting. Our key insight is that the amplification of $\mathcal{E}_\infty(f)$ by least-squares methods

is governed by the Lebesgue constant, a concept from approximation theory. We showed that this amplification is tight and intrinsic to the projection geometry, thereby exposing a fundamental limitation of ordinary and ridge least-squares methods, even in the infinite data regime.

To overcome this limitation, we introduced a novel regularization framework based on weighted ridge regression over extended feature sets, which preserves the desirable properties of the base features while attenuating the contribution of auxiliary ones. We proved that this approach allows us to, asymptotically for $n \rightarrow \infty$, compete with the best possible (oracle) projection in terms of uniform error, and we proposed an efficient algorithm for learning such weights from data. In the polynomial basis case, we demonstrated a dramatic improvement: from $\Omega(d)$ amplification with OLS to the optimal $\mathcal{O}(1)$ with our method.

References

- Luigi Ambrosio, Nicola Fusco, and Diego Pallara. *Functions of bounded variation and free discontinuity problems*. Oxford university press, 2000.
- Philip Amortila, Nan Jiang, and Csaba Szepesvári. The optimal approximation factors in misspecified off-policy value function estimation. In *ICML*, pages 768–790, 2023.
- Philip Amortila, Tongyi Cao, and Akshay Krishnamurthy. Mitigating covariate shift in misspecified regression with applications to reinforcement learning. In *The Thirty Seventh Annual Conference on Learning Theory*, pages 130–160. PMLR, 2024.
- Andras Antos, Csaba Szepesvári, and Rémi Munos. Fitted Q-iteration in continuous action-space MDPs. In *Advances in neural information processing systems*, pages 9–16, 2008.
- Alexandre Belloni, Victor Chernozhukov, Denis Chetverikov, and Kengo Kato. Some new asymptotic theory for least squares series: Pointwise and uniform results. *Journal of Econometrics*, 186(2): 345–366, 2015.
- Stéphane Boucheron, Gábor Lugosi, and Olivier Bousquet. Concentration inequalities. In *Summer school on machine learning*, pages 208–240. Springer, 2003.
- Stephen Boyd, Lin Xiao, and Almir Mutapcic. Subgradient methods. *lecture notes of EE392o, Stanford University, Autumn Quarter*, 2004(01), 2003.
- Xiaohong Chen and Timothy Christensen. Optimal uniform convergence rates for sieve nonparametric instrumental variables regression. *arXiv preprint arXiv:1311.0412*, 2013.
- E.W. Cheney. *Introduction to approximation theory*. McGraw-Hill, 1966.
- Ch de la Vallée Poussin. Sur la meilleure approximation des fonctions d’une variable réelle par des expressions d’ordre donné. *CR Acad. Sci. Paris*, 166:799–802, 1918.
- Ch J De La Vallée Poussin et al. *Leçons sur l’approximation des fonctions d’une variable réelle*. Paris, 1919.
- Ronald A DeVore and George G Lorentz. *Constructive Approximation*, volume 303 of *Grundlehren der mathematischen Wissenschaften*. Springer Berlin Heidelberg, Berlin, Heidelberg, January 1993.
- Jialin Dong and Lin Yang. Does sparsity help in learning misspecified linear bandits? In *International Conference on Machine Learning*, pages 8317–8333. PMLR, 2023.
- Zlatko Drmač, Matjaž Omladič, and Krešimir Veselić. On the perturbation of the cholesky factorization. *SIAM Journal on Matrix Analysis and Applications*, 15(4):1319–1332, 1994.
- Simon S Du, Sham M Kakade, Ruosong Wang, and Lin F Yang. Is a good representation sufficient for sample efficient reinforcement learning? In *International Conference on Learning Representations*, 2020.
- Damien Ernst, Pierre Geurts, and Louis Wehenkel. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6, 2005.

Benyamin Ghogh and Mark Crowley. The theory behind overfitting, cross validation, regularization, bagging, and boosting: tutorial. *arXiv preprint arXiv:1905.12787*, 2019.

László Györfi, Michael Kohler, Adam Krzyzak, and Harro Walk. *A distribution-free theory of nonparametric regression*. Springer Science & Business Media, 2006.

Daniel Hsu, Sham M Kakade, and Tong Zhang. Random design analysis of ridge regression. *Foundations of Computational Mathematics*, 14:569–600, 2014.

Jianhua Z Huang. Local asymptotics for polynomial spline regression. *The Annals of Statistics*, 31(5):1600–1635, 2003.

Yitzhak Katznelson. *An introduction to harmonic analysis*. Cambridge University Press, 2004.

Jack Kiefer and Jacob Wolfowitz. The equivalence of two extremum problems. *Canadian Journal of Mathematics*, 12:363–366, 1960.

Tomasz Kobos and Grzegorz Lewicki. On the dimension of the set of minimal projections. *Journal of Mathematical Analysis and Applications*, 529(2):127250, 2024.

Andrei Kolmogoroff. Über die beste annäherung von funktionen einer gegebenen funktionenklasse. *Annals of Mathematics*, 37(1):107–110, 1936.

George Konidaris, Sarah Osentoski, and Philip Thomas. Value function approximation in reinforcement learning using the Fourier basis. In *Proceedings of the AAAI conference on artificial intelligence*, volume 25, pages 380–385, 2011.

Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.

Tor Lattimore, Csaba Szepesvári, and Gellért Weisz. Learning with good feature representations in bandits and in RL with a generative model. In *International conference on machine learning*, pages 5662–5670. PMLR, 2020.

Jia Li and Zhipeng Liao. Uniform nonparametric inference for time series. *Journal of Econometrics*, 219(1):38–51, 2020.

G.G. Lorentz. *Approximation of functions*. Athena. New York, Holt, Rinehart and Winston, 1966.

Davide Maran, Alberto Maria Metelli, Matteo Papini, and Marcello Restelli. Local linearity: the key for no-regret reinforcement learning in continuous MDPs. *arXiv preprint arXiv:2410.24071*, 2024.

Zsolt Németh. *De la Vallée Poussin Type Approximation Methods*. PhD thesis, Eötvös Loránd University, Hungary, 2016.

Whitney K Newey. Convergence rates and asymptotic normality for series estimators. *Journal of econometrics*, 79(1):147–168, 1997.

Donatella Occorsio and Woula Themistoclakis. De la Vallée Poussin filtered polynomial approximation on the half-line. *Applied Numerical Mathematics*, 207:569–584, 2025.

Allan Pinkus. *N-widths in Approximation Theory*, volume 7. Springer Science & Business Media, 2012.

Alfio Quarteroni, Riccardo Sacco, and Fausto Saleri. *Numerical mathematics*, volume 37. Springer Science & Business Media, 2010.

Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018.

Gábor Szegő. *Orthogonal polynomials*, volume 23. American Mathematical Society, 1939.

Woula Themistoclakis and Marc Van Barel. Generalized de la Vallée Poussin approximations on $[-1, 1]$. *Numerical Algorithms*, 75:1–31, 2017.

Joel A Tropp et al. An introduction to matrix concentration inequalities. *Foundations and Trends® in Machine Learning*, 8(1-2):1–230, 2015.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: —

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: We discussed the limits of the paper and the future research directions in order to address them.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: All the statements are provided with proofs in the appendix.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: We include the code in the supplementary material (very simple, just one very short Jupyter notebook)

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We include the code in the supplementary material (very simple, just one very short Jupyter notebook)

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We include the code in the supplementary material (very simple, just one very short Jupyter notebook)

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: -

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: The simulation is straightforward and its computational time is negligible.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

Answer: [Yes]

Justification: The paper is coherent with NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: —

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: —

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: —

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: —

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: —

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: —

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

710 Question: Does the paper describe the usage of LLMs if it is an important, original, or
711 non-standard component of the core methods in this research? Note that if the LLM is used
712 only for writing, editing, or formatting purposes and does not impact the core methodology,
713 scientific rigorousness, or originality of the research, declaration is not required.

714 Answer: [NA]

715 Justification: —

716 Guidelines:

- 717 • The answer NA means that the core method development in this research does not
718 involve LLMs as any important, original, or non-standard components.
- 719 • Please refer to our LLM policy ([https://neurips.cc/Conferences/2025/](https://neurips.cc/Conferences/2025/LLM)
720 [LLM](https://neurips.cc/Conferences/2025/LLM)) for what should or should not be described.

721 A Related Works

722 **Classical approximation theory** The idea of approximating a class of functions with a family of
 723 vector spaces in a uniform sense has always been an important topic in mathematical analysis. On
 724 the more general level, this theory takes the name of *Kolmogorov's n -width* (Kolmogoroff [1936];
 725 see Lorentz [1966] and Pinkus [2012] for a more modern formalization). The idea, central to this
 726 paper, of finding a *linear* operator that well approximates the *non-linear* L^∞ projection operator has
 727 also been the main topic of multiple line of research. In particular, many result about orthogonal
 728 polynomials Szegő [1939] or Fourier Series Katznelson [2004] approximation have this goal. More
 729 recently, Kobos and Lewicki [2024] studied the problem for general feature map, investigating the
 730 class of linear operators that achieve the lower bound.

731 **Asymptotic pointwise and uniform convergence of LS series in the econometric literature** In the
 732 econometric literature, the series least squares (LS) estimators have been analyzed primarily through
 733 an asymptotic lens: with the sample size $n \rightarrow +\infty$ and the basis dimension $d \rightarrow +\infty$, one studies
 734 asymptotic Gaussianity of the estimator of the function in each single point. Newey [1997] provided
 735 seminal results for this literature, which were then improved by Belloni et al. [2015], the first to use
 736 the Lebesgue constant in this field, and by Li and Liao [2020], who generalize the result to time
 737 series data. All these contributions, however, remain *asymptotic*: they provide limiting distributions
 738 or rates without explicit high-probability bounds, and—crucially—they do not propose algorithmic
 739 modifications capable of *reducing* the amplification factor induced by the Lebesgue constant.

740 **Uniform bounds for linear regression in the context of Online Learning** As anticipated in the
 741 introduction, the problem of getting L^∞ bounds for regression over a domain naturally arises in the
 742 context of Online Learning with linear function approximation; bandits and RL in particular. Du et al.
 743 [2020] established the first \sqrt{d} amplification lower bound in some specific cases, which was then
 744 refined by Lattimore et al. [2020], who also derives the corresponding an upper bound of \sqrt{d} , using
 745 an optimal design argument. In fact, it can be proved that the factor \sqrt{d} is precisely the maximal
 746 Lebesgue constant of any feature map for μ that is the optimal design. These lower bound hold
 747 for a worst-case feature map, but *allowing the learner to choose the data distribution*. Following
 748 these works, many papers tried to understand how this amplification factor could be reduced. Maran
 749 et al. [2024] shows how to remove it in case of a locally linear feature map; Dong and Yang [2023]
 750 improves the \sqrt{d} amplification in case of sparsity. Perhaps, the most similar paper to our one is
 751 Amortila et al. [2024], which proposes a method to mitigate the effect of misspecification w.r.t. the
 752 least-squares fitting. Still, the latter focuses on a different objective, i.e. the error under covariate
 753 shift (measuring the MSE under a distribution $\nu \neq \mu$), and scales with the density ratio $\nu(\cdot)/\mu(\cdot)$.
 754 Generalizing to the uniform error would mean to take $\nu(\cdot)$ as a Dirac's delta, which would make this
 755 bound vacuous.

756 **De la Vallée Poussin approach** The to reduce the Lebesgue constant by adding auxiliary features is
 757 rooted in a concept that dates back in the history of mathematics to Baron de la Vallée Poussin [de la
 758 Vallée Poussin, 1918, De La Vallée Poussin et al., 1919]. The technique he invented is still studied
 759 today in numerical mathematics [Németh, 2016, Themistoclakis and Van Barel, 2017, Occorsio and
 760 Themistoclakis, 2025].

761 **Finite-sample bounds for ridge regression** Hsu et al. [2014] gives finite-sample bounds for ridge
 762 regression under random design. The results, when translated into our setting, bound the error between
 763 $\hat{f}_{\hat{\theta}_n}$ and \bar{f} where $\bar{f} := g \circ \varphi$ and the bound is expressed in terms of $\bar{f} - \Pi_{\mu,d} f$. Here for $u \in \mathbb{R}^d$,
 764 $g(u) = \int f(x) \mu(dx|u)$ where $\mu(dx|u)$ is the disintegration of μ with respect to the push-forward
 765 of μ under φ . In particular, for $S \subset \mathcal{X}$, $u \in \mathbb{R}^d$, $\mu(S|u) = \int \mathbb{I}(x \in S, \varphi(x) = u) \mu(dx)$. In the
 766 special case when φ is injective, $\bar{f} = f$. Just like in the result that can be extracted from the work of
 767 Lattimore et al. [2020], the bounds in this work depend on $\bar{\varphi}_{d,2}$ (or $\hat{\varphi}_{d,2}$) and scale similarly.

768 B General-interest results

769 We start from the usual Bernstein's inequality Boucheron et al. [2003], here written for variables that
 770 are bounded in $[-B, B]$ and in the "high probability" form.

771 **Theorem 16.** *Let $\{X_t\}_{t=1}^n$ be a sequence of zero-mean random variable bounded in $[-B, B]$. Let*
 772 $\sigma^2 := \sum_{t=1}^n \text{Var}(X_t)$. *Then, with probability at least $1 - \delta$*

$$\left| \sum_{t=1}^n X_t \right| \leq \sqrt{2\sigma^2 \log(2/\delta)} + \frac{2B}{3} \log(2/\delta).$$

Lemma 2. Let $\bar{\varphi}_d$ be an orthonormal feature map w.r.t. ρ .

$$\mathbb{E}_{x \sim \rho} [\bar{\varphi}_d(x) \bar{\varphi}_d(x)^\top] = I_d,$$

where I_d is the d -dimensional identity matrix.

Proof. In this proof, let us denote with e_i , for $i = 1, \dots, d$, the standard basis of \mathbb{R}^d . By definition of outer product between two vectors we get what follows.

$$\begin{aligned} \mathbb{E}_{x \sim \rho} [\bar{\varphi}_d(x) \bar{\varphi}_d(x)^\top] &= \mathbb{E}_{x \sim \rho} \left[\sum_{i=1}^d \sum_{j=1}^d \bar{\varphi}_i(x) \bar{\varphi}_j(x) e_i e_j^\top \right] \\ &= \sum_{i=1}^d \sum_{j=1}^d \mathbb{E}_{x \sim \rho} [\bar{\varphi}_i(x) \bar{\varphi}_j(x)] e_i e_j^\top \\ &= \sum_{i=1}^d \sum_{j=1}^d \delta_{ij} e_i e_j^\top = I_d. \end{aligned}$$

This completes the proof. \square

Lemma 3. Let $\{v_t\}_{t=1}^k$ be a sequence of independent d -dimensional random vectors such that

$$\mathbb{E}[v_t v_t^\top] = \sigma I_d \quad \|v_t\|_2^2 \leq B.$$

Let $V := \sum_{t=1}^k v_t v_t^\top$. Then,

1. W.p. at least $1 - \delta$

$$\lambda_{\min}(V) \geq \left(1 - \sqrt{\frac{5B \log(d/\delta)}{k\sigma^2}}\right) k\sigma^2,$$

$$\text{if } \left(1 - \sqrt{\frac{5B \log(d/\delta)}{k\sigma^2}}\right) \leq 1/2.$$

2. W.p. at least $1 - \delta$

$$\lambda_{\max}(V) \leq \left(1 + \sqrt{\frac{2B \log(d/\delta)}{k\sigma^2}}\right) k\sigma^2,$$

$$\text{if } \left(1 + \sqrt{\frac{2B \log(d/\delta)}{k\sigma^2}}\right) \leq 1.$$

Proof. Note that, as $\lambda_{\max}(v_t v_t^\top) = \|v_t\|_2^2 \leq B$, we can then apply Theorem 5.1.1 from [Tropp et al. \[2015\]](#) taking

$$\mu_{\min} = \mu_{\max} = k\sigma^2 \quad L = B,$$

which ensures that

$$\forall \varepsilon \in (0, 1), \mathbb{P}(\lambda_{\min}(V) \leq (1 - \varepsilon)k\sigma^2) \leq d \left(\frac{e^{-\varepsilon}}{(1 - \varepsilon)^{1-\varepsilon}} \right)^{k\sigma^2/B},$$

while

$$\forall \varepsilon > 0, \mathbb{P}(\lambda_{\max}(V) \geq (1 + \varepsilon)k\sigma^2) \leq d \left(\frac{e^\varepsilon}{(1 + \varepsilon)^{1+\varepsilon}} \right)^{k\sigma^2/B}.$$

785 The thesis is going to follow by just simplifying the previous expressions. We recall from elementary
786 Taylor expansions that

$$\varepsilon < 0.5 \implies -\varepsilon - 4\varepsilon^2/5 \leq \log(1 - \varepsilon) \leq -\varepsilon - \frac{\varepsilon^2}{2}.$$

787 and

$$\varepsilon < 1 \implies \varepsilon - \frac{\varepsilon^2}{2} \leq \log(1 + \varepsilon) \leq \varepsilon - \frac{\varepsilon^2}{4}.$$

788 Therefore, we have, for $\varepsilon < 0.5$

$$\begin{aligned} \frac{e^{-\varepsilon}}{(1 - \varepsilon)^{1-\varepsilon}} &= \exp(-\varepsilon - (1 - \varepsilon)\log(1 - \varepsilon)) \\ &\leq \exp(-\varepsilon - (1 - \varepsilon)(-\varepsilon - 4\varepsilon^2/5)) \\ &= \exp(-\varepsilon + \varepsilon - \varepsilon^2/5 + \mathcal{O}(\varepsilon^3)) \approx e^{-\varepsilon^2/5}. \end{aligned}$$

789 On the other side, for $\varepsilon \leq 1$,

$$\begin{aligned} \frac{e^\varepsilon}{(1 + \varepsilon)^{1+\varepsilon}} &= \exp(\varepsilon - (1 + \varepsilon)\log(1 + \varepsilon)) \\ &\leq \exp(\varepsilon - (1 + \varepsilon)(\varepsilon - \varepsilon^2/2)) \\ &= \exp(-\varepsilon^2/2 + \mathcal{O}(\varepsilon^3)) \approx e^{-\varepsilon^2/2}. \end{aligned}$$

790 This tells us that

$$\forall \varepsilon \in (0, 1/2), \mathbb{P}(\lambda_{\min}(V) \leq (1 - \varepsilon)k\sigma^2) \leq de^{-k\sigma^2\varepsilon^2/(5B)},$$

791 and

$$\forall \varepsilon \in (0, 1), \mathbb{P}(\lambda_{\max}(V) \geq (1 + \varepsilon)k\sigma^2) \leq de^{-k\sigma^2\varepsilon^2/(2B)}.$$

792 We can reformulate the previous results in the high-probability notation. Indeed, taking $\delta =$
793 $de^{-k\sigma^2\varepsilon^2/(5B)}$, we get

$$\varepsilon = \sqrt{\frac{5B \log(d/\delta)}{k\sigma^2}},$$

794 which entails that

$$\sqrt{\frac{5B \log(d/\delta)}{k\sigma^2}} \leq 1/2 \implies \mathbb{P}\left(\lambda_{\min}(V) \leq \left(1 - \sqrt{\frac{5B \log(d/\delta)}{k\sigma^2}}\right)k\sigma^2\right) \leq \delta.$$

795 Doing the same for the other result, we get

$$\sqrt{\frac{2B \log(d/\delta)}{k\sigma^2}} \leq 1 \implies \mathbb{P}\left(\lambda_{\max}(V) \geq \left(1 + \sqrt{\frac{2B \log(d/\delta)}{k\sigma^2}}\right)k\sigma^2\right) \leq \delta,$$

796 which completes the proof. □

797

798 **Proposition 17.** *The Lebesgue constant satisfies $\Lambda_{d,\mu} = \sup_{x \in \mathcal{X}} \int_{\mathcal{X}} \left| \sum_{i=1}^d \bar{\varphi}_i(z) \bar{\varphi}_i(x) \right| d\mu(z)$.*

799 *Proof.* See [Cheney \[1966\]](#), chapter 4. □

800 C Proofs from section 3

801 C.1 Lower bound for LS

802 Recall that $\mathcal{F} = L^2(\mu) \cap L^\infty(\mathcal{X})$. Let $\Pi_\infty f = \arg \min_{g \in \mathcal{F}_d} \|f - g\|_\infty$ with ties broken arbitrarily.
 803 Note that Theorem 1.1 of Chapter 3 in the book of [DeVore and Lorentz \[1993\]](#) guarantees that at
 804 least one minimizer exists. (As discussed there, uniqueness may or may not hold.)

805 **Lemma 4.** *We have*

$$\sup_{f \in \mathcal{F}} \frac{\|\Pi_{d,\mu} f - f\|_\infty}{\mathcal{E}_\infty(f)} \geq \Lambda_{d,\mu} - 1.$$

Proof. By definition of Lebesgue constant, for every $\varepsilon > 0$ there is a function g such that

$$\|\Pi_{d,\mu} g\|_\infty \geq (\Lambda_{d,\mu} - \varepsilon) \|g\|_\infty.$$

806 Take $f = \Pi_\infty g - g$. We will use twice that for any $h \in \mathcal{F}$, $\|h\|_\infty = \|0 - h\|_\infty \geq \inf_{u \in \mathcal{F}_d} \|u -$
 807 $h\|_\infty = \|\Pi_\infty h - h\|_\infty$. Now,

$$\begin{aligned} \|\Pi_{d,\mu} f - f\|_\infty &= \|\Pi_{d,\mu}(\Pi_\infty g - g) - \Pi_\infty g + g\|_\infty \\ &= \|\Pi_\infty g - \Pi_{d,\mu} g - \Pi_\infty g + g\|_\infty \\ &= \|\Pi_{d,\mu} g - g\|_\infty \\ &\geq \|\Pi_{d,\mu} g\|_\infty - \|g\|_\infty \\ &\geq (\Lambda_{d,\mu} - 1 - \varepsilon) \|g\|_\infty \\ &\geq (\Lambda_{d,\mu} - 1 - \varepsilon) \|\Pi_\infty g - g\|_\infty \\ &= (\Lambda_{d,\mu} - 1 - \varepsilon) \|f\|_\infty \\ &\geq (\Lambda_{d,\mu} - 1 - \varepsilon) \|\Pi_\infty f - f\|_\infty. \end{aligned}$$

808 The result follows by letting $\varepsilon \rightarrow 0$. □

809 **Theorem 3.** *For any $\varepsilon > 0$ there exist $f \in L^\infty(\mathcal{X})$ such that*

$$\|f - \Pi_{d,\mu} f\|_\infty \geq (\Lambda_{d,\mu} - 1 - \varepsilon) \mathcal{E}_\infty(f).$$

810 *Proof.* The result is immediate from Lemma 4. □

811 C.2 Towards the proof of theorem 4

Lemma 5. *Fix $\delta > 0$, and $n \geq 20\bar{\varphi}_{d,2}^2 \log(d/\delta)$. Let*

$$V_n = \sum_{t=1}^n \bar{\varphi}_d(x_t) \bar{\varphi}_d(x_t)^\top.$$

812 *Then, $\lambda_{\min}(V_n) \geq n/2$.*

813 *Proof.* The matrices we are summing correspond to $\bar{\varphi}_d(x_t) \bar{\varphi}_d(x_t)^\top$ each one being semi-positive
 814 definite with the biggest eigenvalue bounded by $\bar{\varphi}_{d,2}^2$ almost surely (indeed, $v^\top \bar{\varphi}_d(x_t) \bar{\varphi}_d(x_t)^\top v$ is
 815 maximized for v parallel to $\bar{\varphi}_d(x_t)$ and produces $\|\bar{\varphi}_d(x_t)\|_2^2$). Moreover, as we have seen in lemma
 816 2,

$$\mathbb{E} \left[\sum_{t=1}^n \bar{\varphi}_d(x_t) \bar{\varphi}_d(x_t)^\top \right] = \sum_{t=1}^n \mathbb{E} [\bar{\varphi}_d(x_t) \bar{\varphi}_d(x_t)^\top] = nI_d.$$

These two ingredients allow us to apply lemma 3 part one, which ensures that with probability at least $1 - \delta$

$$\lambda_{\min}(V_n) \geq \left(1 - \sqrt{\frac{5\bar{\varphi}_{d,2}^2 \log(d/\delta)}{n}} \right) n,$$

817 if $\left(1 - \sqrt{\frac{5\bar{\varphi}_{d,2}^2 \log(d/\delta)}{n}} \right) \leq 1/2$. Therefore, taking $n \geq 20\bar{\varphi}_{d,2}^2 \log(d/\delta)$, we get $\lambda_{\min}(V_n) \geq n/2$,
 818 which completes the proof.

819

□

Lemma 6. Let $\zeta_{d,\mu}(\cdot) := f(\cdot) - \Pi_{d,\mu}f(\cdot)$. With probability at least $1 - \delta$,

$$\left| \bar{\varphi}_d(z)^\top V_n^{-1} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right| \leq \frac{2\Lambda_{d,\mu} \mathcal{E}_\infty(f) \bar{\varphi}_{d,2}}{\sqrt{n}} \sqrt{\log(1/\delta)},$$

820 plus a lower-order term depending on n^{-1} which takes the form of
 821 $\tilde{\mathcal{O}}(n^{-1}d^{1/2}\bar{\varphi}_{d,2}^2\Lambda_{d,\mu}\mathcal{E}_\infty(f) + n^{-3/2}d\bar{\varphi}_{d,2}^3\Lambda_{d,\mu}\mathcal{E}_\infty(f))$.

822 *Proof.* We start rearranging the equation as follows

$$\begin{aligned} \left| \bar{\varphi}_d(z)^\top V_n^{-1} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right| &= \left| \bar{\varphi}_d(z)^\top \left(\frac{1}{n} V_n \right)^{-1} \frac{1}{n} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right| \\ &= \left| \bar{\varphi}_d(z)^\top (I_d + \Delta_n) \frac{1}{n} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right| \\ &\leq \left| \bar{\varphi}_d(z)^\top \frac{1}{n} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right| \\ &\quad + \left| \bar{\varphi}_d(z)^\top \Delta_n \frac{1}{n} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right|. \end{aligned}$$

823 For $\Delta_n := (V_n/n)^{-1} - I_d$. To bound both parts, we start by giving a result for
 824 $\frac{1}{n} \sum_{t=1}^n v^\top \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t)$ that holds for one fixed $v \in \mathbb{R}^d$. Indeed,

825 1. Every random variable $v^\top \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t)$ is bounded by $\|v\|_2 \bar{\varphi}_{d,2} \Lambda_{d,\mu} \mathcal{E}_\infty(f)$ a.s.

826 2. The variance of the same random variable is

$$\begin{aligned} \mathbb{E}_{x \sim \rho} [(v^\top \bar{\varphi}_d(x) \zeta_{d,\mu}(x))^2] &= \mathbb{E}_{x \sim \rho} [\zeta_{d,\mu}(x)^2 v^\top \bar{\varphi}_d(x)^\top \bar{\varphi}_d(x) v] \\ &\leq (\Lambda_{d,\mu} \mathcal{E}_\infty(f))^2 v^\top \mathbb{E}_{x \sim \rho} [\bar{\varphi}_d(x)^\top \bar{\varphi}_d(x)] v \\ &= (\Lambda_{d,\mu} \mathcal{E}_\infty(f))^2 v^\top I_d v \\ &= (\Lambda_{d,\mu} \mathcal{E}_\infty(f))^2 \|v\|_2^2, \end{aligned}$$

827 the main step following from lemma 2.

828 So by Bernstein's inequality (theorem 16),

$$\frac{1}{n} \sum_{t=1}^n v^\top \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \leq \frac{2\Lambda_{d,\mu} \mathcal{E}_\infty(f) \|v\|_2}{\sqrt{n}} \sqrt{\log(1/\delta)} + \frac{2\|v\|_2 \bar{\varphi}_{d,2} \Lambda_{d,\mu} \mathcal{E}_\infty(f)}{3n} \log(1/\delta). \quad (9)$$

We can use the previous equation to bound both parts. For the first, we just take $v = \bar{\varphi}_d(z)$, which respects $\|v\|_2 \leq \bar{\varphi}_{d,2}$, in equation 9 and get

$$\left| \bar{\varphi}_d(z)^\top \frac{1}{n} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right| \leq \frac{2\Lambda_{d,\mu}\mathcal{E}_\infty(f)\bar{\varphi}_{d,2}}{\sqrt{n}} \sqrt{\log(1/\delta)} + \frac{2\bar{\varphi}_{d,2}\Lambda_{d,\mu}\mathcal{E}_\infty(f)}{3n} \log(1/\delta).$$

Let us now focus on the second part. Indeed,

$$\left| \bar{\varphi}_d(z)^\top \Delta_n \frac{1}{n} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right| \leq \bar{\varphi}_{d,2} \|\Delta_n\|_2 \left\| \frac{1}{n} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right\|_2$$

829 Now, using lemma 3 as done in the proof of lemma 5, we have

$$\|\Delta_n\|_2 \leq \bar{\varphi}_{d,2} \sqrt{\frac{5 \log(d/\delta)}{n}},$$

830 while for the last part we can write

$$\begin{aligned} \left\| \frac{1}{n} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right\|_2 &= \sup_{\|v\|_2=1} \frac{1}{n} \sum_{t=1}^n v^\top \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \\ &\leq \sup_{\|v\|_2 \in B_d^{1/n}} \frac{1}{n} \sum_{t=1}^n v^\top \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) + \frac{\bar{\varphi}_{d,2}\Lambda_{d,\mu}\mathcal{E}_\infty(f)}{n}, \end{aligned}$$

831 where $B_d^{1/n}$ is a $1/n$ covering of the set of vectors such that $\|v\|_2 = 1$. It is well-known that we can
832 choose $B_d^{1/n}$ so that $|B_d^{1/n}| \approx n^{-d}$, so that, making a union bound together with equation 9, we get

$$\left\| \frac{1}{n} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right\|_2 \leq \frac{2\Lambda_{d,\mu}\mathcal{E}_\infty(f)}{\sqrt{n}} \sqrt{d \log(1/\delta)} + \frac{2\bar{\varphi}_{d,2}\Lambda_{d,\mu}\mathcal{E}_\infty(f)}{3n} \log(1/\delta) + \frac{\Lambda_{d,\mu}\mathcal{E}_\infty(f)}{n}.$$

833 As a consequence,

$$\begin{aligned} &\left| \bar{\varphi}_d(z)^\top \Delta_n \frac{1}{n} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right| \\ &\leq \bar{\varphi}_{d,2}^2 \sqrt{\frac{5 \log(d/\delta)}{n}} \left(\frac{2\Lambda_{d,\mu}\mathcal{E}_\infty(f)}{\sqrt{n}} \sqrt{d \log(1/\delta)} + \frac{2\bar{\varphi}_{d,2}\Lambda_{d,\mu}\mathcal{E}_\infty(f)}{3n} \log(1/\delta) + \frac{\Lambda_{d,\mu}\mathcal{E}_\infty(f)}{n} \right) \\ &\leq \tilde{\mathcal{O}} \left(n^{-1} d^{1/2} \bar{\varphi}_{d,2}^2 \Lambda_{d,\mu} \mathcal{E}_\infty(f) + n^{-3/2} d \bar{\varphi}_{d,2}^3 \Lambda_{d,\mu} \mathcal{E}_\infty(f) \right). \end{aligned}$$

834 This completes the proof. \square

835 C.3 Proof of theorem 4

836 **Theorem 4.** Let \mathcal{X} be finite. Let Assumptions 1 and 2 hold. Then, for any n positive integer and real
837 $0 < \delta \leq 1/3$ such that $n \geq 20\bar{\varphi}_{d,2}^2 \log(d/\delta)$, letting $\hat{\theta}_{n,OLS}$ be the parameter vector returned by OLS,
838 with probability at least $1 - 3\delta$,

$$\begin{aligned} \mathcal{E}_\infty(\hat{\theta}_{n,OLS}, f) &\leq (1 + \Lambda_{d,\mu})\mathcal{E}_\infty(f) + 3(\sigma + \Lambda_{d,\mu}\mathcal{E}_\infty(f))\bar{\varphi}_{d,2} \sqrt{\frac{\log(|\mathcal{X}|/\delta)}{n}} \\ &\quad + \frac{\text{poly}(d, \bar{\varphi}_{d,2}, \Lambda_{d,\mu}\mathcal{E}_\infty(f))}{n}. \end{aligned}$$

839 *Proof.* In this proof, as before, we call Let $\zeta_{d,\mu}(\cdot) := f(\cdot) - \Pi_{d,\mu}f(\cdot)$.

840 Trough this proof we call $\hat{\theta}_n$ the OLS estimator and $\hat{f}_n(\cdot)$ the corresponding estimated function. We
841 start making the following decomposition:

$$\begin{aligned} |\bar{\varphi}_d(x)^\top \hat{\theta}_n - f(x)| &\leq |\bar{\varphi}_d(x)^\top \hat{\theta}_n - \Pi_{d,\mu}f(x)| + \|\Pi_{d,\mu}f - f\|_\infty \\ &\leq |\bar{\varphi}_d(x)^\top \hat{\theta}_n - \Pi_{d,\mu}f(x)| + (1 + \Lambda_{d,\mu})\mathcal{E}_\infty(f). \end{aligned}$$

842 To bound the first part, we let θ_* be such that $\Pi_{d,\mu}f(\cdot) = \bar{\varphi}_d(\cdot)^\top \theta_*$. By Assumption 1, the
843 samples take the form $y_t = \bar{\varphi}_d(x_t)^\top \theta_* + \zeta_{d,\mu}(x_t) + \eta_t$, where $\{\eta_t\}_{t=1}^n$ is a family of independent
844 σ -subgaussian random variables. By definition, letting $V_n = \sum_{t=1}^n \bar{\varphi}_d(x_t)\bar{\varphi}_d(x_t)^\top$, the LS solution
845 takes the form $\bar{\varphi}_d(x_t)^\top \hat{\theta}_n$, where

$$\begin{aligned} \hat{\theta}_n &= V_n^{-1} \sum_{t=1}^n \bar{\varphi}_d(x_t) y_t \\ &= V_n^{-1} \sum_{t=1}^n \bar{\varphi}_d(x_t) (\bar{\varphi}_d(x_t)^\top \theta_* + \eta_t + \zeta_{d,\mu}(x_t)) \\ &= \theta_* + V_n^{-1} \sum_{t=1}^n \bar{\varphi}_d(x_t) (\eta_t + \zeta_{d,\mu}(x_t)). \end{aligned}$$

Therefore, we have

$$|\bar{\varphi}_d(x)^\top \hat{\theta}_n - \Pi_{d,\mu}f(x)| \leq \underbrace{\left| \bar{\varphi}_d(x)^\top V_n^{-1} \sum_{t=1}^n \bar{\varphi}_d(x_t) \eta_t \right|}_{(I)} + \underbrace{\left| \bar{\varphi}_d(x)^\top V_n^{-1} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right|}_{(II)}.$$

846 We are going to bound the two terms separately. First, let $E := \{\lambda_{\min}(V_n) \geq n/2\}$. From lemma 5,
847 under the assumptions of this theorem, we have $\mathbb{P}(E) \geq 1 - \delta$.

848 (I) Since η_t are independent and σ -subgaussian, Lemma 5.4 and Theorem 5.3 from [Lattimore](#)
849 [and Szepesvári \[2020\]](#) ensure that, with probability at least $1 - 2\delta$

$$\begin{aligned} \left| \bar{\varphi}_d(x)^\top V_n^{-1} \sum_{t=1}^n \bar{\varphi}_d(x_t) \eta_t \right| &\leq \sqrt{2\sigma^2 \sum_{t=1}^n (\bar{\varphi}_d(x)^\top V_n^{-1} \bar{\varphi}_d(x_t))^2 \log(1/\delta)} \\ &= \sqrt{2\sigma^2 \|\bar{\varphi}_d(x)\|_{V_n^{-1}}^2 \log(1/\delta)} \\ &= \sqrt{2 \log(1/\delta)} \sigma \|\bar{\varphi}_d(x)\|_{V_n^{-1}}. \end{aligned}$$

Moreover, if event E holds,

$$\|\bar{\varphi}_d(x)\|_{V_n^{-1}} \leq \frac{2\|\bar{\varphi}_d(x)\|_2}{\sqrt{n}} \leq \frac{2\bar{\varphi}_{d,2}}{\sqrt{n}},$$

850 so that the full term is bounded by $\sqrt{8 \log(1/\delta)} \sigma \bar{\varphi}_{d,2} n^{-1/2}$.

851 (II) This term is bounded by lemma 6 which, with probability at least $1 - \delta$ gives

$$\left| \bar{\varphi}_d(z)^\top V_n^{-1} \sum_{t=1}^n \bar{\varphi}_d(x_t) \zeta_{d,\mu}(x_t) \right| \leq \frac{2\Lambda_{d,\mu}\mathcal{E}_\infty(f)\bar{\varphi}_{d,2}}{\sqrt{n}} \sqrt{\log(1/\delta)},$$

852 plus lower-order terms of the form $\frac{\text{poly}(d, \bar{\varphi}_{d,2}, \Lambda_{d,\mu}\mathcal{E}_\infty(f))}{n}$.

853 Note that, thanks to lemma 5, event E holds with probability $1 - \delta$ under the assumptions of this
 854 theorem. Moreover, imposing that both events in (I) and (II) verify, we get, with probability at least
 855 $1 - 3\delta$,

$$\begin{aligned} |\bar{\varphi}_d(x)^\top \hat{\theta}_n - f(x)|_\infty &\leq (1 + \Lambda_{d,\mu}) \mathcal{E}_\infty(f) + |\bar{\varphi}_d(x)^\top \hat{\theta}_n - \Pi_{d,\mu} f(x)| \\ &\leq (1 + \Lambda_{d,\mu}) \mathcal{E}_\infty(f) + \frac{3(\sigma + \Lambda_{d,\mu} \mathcal{E}_\infty(f)) \bar{\varphi}_{d,2}}{\sqrt{n}} \sqrt{\log(1/\delta)} \end{aligned}$$

856 plus lower-order terms of the form $\frac{\text{poly}(d, \bar{\varphi}_{d,2}, \Lambda_{d,\mu} \mathcal{E}_\infty(f))}{n}$. This completes the proof. \square

857 C.4 Bound scaling with the empirical lebesgue constant

Theorem 5. *Let Assumption 1 hold. Then, for any fixed $\delta > 0$, with probability at least $1 - \delta$,*

$$\mathcal{E}_\infty(\hat{\theta}_{n,\text{OLS}}) \leq (1 + \Lambda_{d,\mu_n}) \mathcal{E}_\infty(f) + \frac{\sigma \hat{\varphi}_{d,2} \sqrt{2 \log(2\mathcal{X}/\delta)}}{\sqrt{n}}.$$

Proof. Let $\hat{\theta}_n$ the estimator corresponding to $\hat{\theta}_{n,\text{OLS}}$ in the parameterization of $\hat{\varphi}_d(\cdot)$, so that

$$\hat{\varphi}_d(\cdot)^\top \hat{\theta}_n = \varphi_d(\cdot)^\top \hat{\theta}_{n,\text{OLS}} =: \hat{f}_n(\cdot).$$

858 The following decomposition holds:

$$\begin{aligned} \|f(\cdot) - \hat{f}_n(\cdot)\|_\infty &\leq \|f(\cdot) - \Pi_{d,\mu_n} f(\cdot)\|_\infty + \|\Pi_{d,\mu_n} f(\cdot) - \hat{f}_n(\cdot)\|_\infty \\ &\leq (1 + \hat{\Lambda}_{d,\mu}) \mathcal{E}_\infty(f) + \|\Pi_{d,\mu_n} f(\cdot) - \hat{f}_n(\cdot)\|_\infty. \end{aligned}$$

859 Now, we focus on the second term. As done in the previous proof of theorem 4, we let θ_\star be such
 860 that $\Pi_{d,\mu_n} f(\cdot) = \hat{\varphi}_d(\cdot)^\top \theta_\star$ and $\zeta_{d,\mu_n}(\cdot) := f(\cdot) - \hat{\varphi}_d(\cdot)^\top \theta_\star$. In this way, our samples take the form
 861 $y_t = \hat{\varphi}_d(x_t)^\top \theta_\star + \zeta_{d,\mu_n}(x_t) + \eta_t$.

862 For any fixed $x \in \mathcal{X}$ we have

$$\begin{aligned} \hat{f}_n(x) &= \hat{\varphi}_d(x)^\top \hat{\theta}_n \\ &= \hat{\varphi}_d(x)^\top \frac{1}{n} \sum_{t=1}^n \hat{\varphi}_d(x_t) y_t \\ &= \hat{\varphi}_d(x)^\top \frac{1}{n} \sum_{t=1}^n \hat{\varphi}_d(x_t) (\hat{\varphi}_d(x_t)^\top \theta_\star + \zeta_{d,\mu_n}(x_t) + \eta_t) \\ &= \underbrace{\hat{\varphi}_d(x)^\top \theta_\star + \hat{\varphi}_d(x)^\top \frac{1}{n} \sum_{t=1}^n \hat{\varphi}_d(x_t) \zeta_{d,\mu_n}(x_t)}_{(I)} + \underbrace{\hat{\varphi}_d(x)^\top \frac{1}{n} \sum_{t=1}^n \hat{\varphi}_d(x_t) \eta_t}_{(II)}. \end{aligned}$$

863 Here, the last passage is due to the fact that, being $\hat{\varphi}_d(\cdot)$ orthogonal w.r.t. $\mu_n(\cdot)$, it follows
 864 $\frac{1}{n} \sum_{t=1}^n \hat{\varphi}_d(x_t) \hat{\varphi}_d(x_t)^\top = I_d$. Now, we analyze the two terms (I) and (II) separately.

$$\begin{aligned} (I) &= \hat{\varphi}_d(x)^\top \frac{1}{n} \sum_{t=1}^n \hat{\varphi}_d(x_t) \zeta_{d,\mu_n}(x_t) \\ &= \hat{\varphi}_d(x)^\top \int_{\mathcal{X}} \hat{\varphi}_d(z) \zeta_{d,\mu_n}(z) d\mu_n(z) = \hat{\varphi}_d(x)^\top \mathbf{0} = 0. \end{aligned}$$

865 In fact, by definition of orthogonal projection, $\zeta_{d,\mu_n}(\cdot)$ is orthogonal in $L^2(\mu_n)$ to the span of $\hat{\varphi}_d(\cdot)$,
 866 so to each of its components in particular.

Let us look at the second term. Since η_t are independent and σ -subgaussian, Lemma 5.4 and Theorem 5.3 from [Lattimore and Szepesvári \[2020\]](#) ensure that, with probability at least $1 - 2\delta$

$$\begin{aligned} \left| \hat{\varphi}_d(x)^\top n^{-1} \sum_{t=1}^n \hat{\varphi}_d(x_t) \eta_t \right| &\leq \sqrt{2\sigma^2 n^{-1} \sum_{t=1}^n (\hat{\varphi}_d(x)^\top \hat{\varphi}_d(x_t))^2 \log(1/\delta)} \\ &= \sqrt{2\sigma^2 n^{-1} \|\hat{\varphi}_d(x)\|_2^2 \log(1/\delta)} \\ &= \sqrt{2 \log(1/\delta) \sigma n^{-1/2}} \|\hat{\varphi}_d(x)\|_2. \end{aligned}$$

Where the second passage comes once again from the fact that $\frac{1}{n} \sum_{t=1}^n \hat{\varphi}_d(x_t) \hat{\varphi}_d(x_t)^\top = I_d$. This proves that (II) is bounded by $\sqrt{2 \log(1/\delta) \sigma n^{-1/2}} \hat{\varphi}_{2,d}$. Making a union bound over $x \in \mathcal{X}$, this entails w.p. $1 - \delta$,

$$\sup_{x \in \mathcal{X}} \left| \hat{\varphi}_d(x)^\top n^{-1} \sum_{t=1}^n \hat{\varphi}_d(x_t) \eta_t \right| \leq \sqrt{2 \log(|\mathcal{X}|/\delta) \sigma n^{-1/2}} \hat{\varphi}_{2,d}.$$

We have proved that

$$\begin{aligned} \mathcal{E}_\infty(\hat{\theta}_{n,\text{OLS}}) &= \|f(\cdot) - \hat{f}_n(\cdot)\|_\infty \\ &\leq (1 + \hat{\Lambda}_{d,\mu}) \mathcal{E}_\infty(f) + \|\Pi_{d,\mu_n} f(\cdot) - \hat{f}_n(\cdot)\|_\infty \\ &\leq (1 + \hat{\Lambda}_{d,\mu}) \mathcal{E}_\infty(f) + \sqrt{2 \log(|\mathcal{X}|/\delta) \sigma n^{-1/2}} \hat{\varphi}_{2,d}. \end{aligned}$$

□

C.5 Proofs from section 3.1

Proposition 18. *The Lebesgue constant is bounded by $\Lambda_{d,\mu} \leq \bar{\varphi}_{d,2}$.*

Proof. Let $f \in L^\infty(\mathcal{X})$ with $\|f\|_\infty = 1$. We have, for any $x \in \mathcal{X}$,

$$\begin{aligned} |\Pi_{d,\mu} f(x)| &= \left| \sum_{i=1}^d \langle f, \bar{\varphi}_i \rangle \bar{\varphi}_i(x) \right| \\ &\leq \sqrt{\sum_{i=1}^d \langle f, \bar{\varphi}_i \rangle^2 \sum_{i=1}^d \bar{\varphi}_i(x)^2} \\ &\leq \sqrt{\|f\|_\mu^2 \|\bar{\varphi}_i(x)\|_2^2} \\ &\leq \|f\|_\infty \sqrt{\|\bar{\varphi}_i(x)\|_2^2} \leq \bar{\varphi}_{d,2}, \end{aligned}$$

the last passage coming from the fact that as ρ is a probability measure, $\|f\|_\mu \leq \|f\|_\infty$. The thesis follows taking the supremum on f, x . □

Proposition 19. *Let $\varphi_d : \mathcal{X} \rightarrow \mathbb{R}^d$ be any feature map, and ρ a probability measure. Then,*

$$\bar{\varphi}_2 \geq \sqrt{d}.$$

879 *Proof.* The key for this result is to note that, being ρ a probability measure, $\bar{\varphi}_{d,2}^2 \geq \mathbb{E}_{x \sim \rho} [\|\bar{\varphi}_d(x)\|_2^2]$
 880 (the supremum of a function upper bounds its integral on any probability measure). Then,

$$\begin{aligned} \bar{\varphi}_{d,2} &\geq \sqrt{\mathbb{E}_{x \sim \rho} [\|\bar{\varphi}_d(x)\|_2^2]} \\ &= \sqrt{\mathbb{E}_{x \sim \rho} [\bar{\varphi}_d(x)^\top \bar{\varphi}_d(x)]} \\ &= \sqrt{\mathbb{E}_{x \sim \rho} [\text{Tr}(\bar{\varphi}_d(x)^\top \bar{\varphi}_d(x))]} \\ &= \sqrt{\mathbb{E}_{x \sim \rho} [\text{Tr}(\bar{\varphi}_d(x) \bar{\varphi}_d(x)^\top)]} \\ &= \sqrt{\text{Tr}(\mathbb{E}_{x \sim \rho} [\bar{\varphi}_d(x) \bar{\varphi}_d(x)^\top])} \\ &\stackrel{*}{=} \sqrt{\text{Tr}(I_d)} = \sqrt{d}. \end{aligned}$$

881 Where the passage $(*)$ comes from lemma 2. □

Proposition 20. Let $\mathcal{X} = [k]$ and $\varphi_i(j) = X_{ij}$, with all the X_{ij} being independent bounded zero-mean unit variance random variables. Then, if $d = \mathcal{O}(\sqrt{k})$, the feature map φ_d , satisfies

$$\Lambda_{d,\mu} = \mathcal{O}(\sqrt{d \log(k/\delta)})$$

882 with probability at least $1 - \delta$. Moreover, $\mathbb{E}[\Lambda_{d,\mu}] \geq \Omega(\sqrt{d})$.

883 *Proof.* By convenience, we call $\Phi \in \mathbb{R}^{k \times d}$ the matrix having, as columns, the features of φ_d .
 884 Precisely, the i -th column of Φ corresponds to φ_i . It is well-known that, in a finite dimensional
 885 space the orthogonal projection operator writes as

$$\Pi_{d,\mu} := \Phi(\Phi^\top \Phi)^{-1} \Phi^\top.$$

886 We call Φ_m , the m -th row of Φ which, by assumption, is a random vector of independent entries
 887 bounded in $[-B, B]$ and with variance one. We have

$$\Phi^\top \Phi = \sum_{m=1}^k \Phi_m \Phi_m^\top, \quad \mathbb{E}[\Phi_m \Phi_m^\top] = \sigma^2 I_d, \quad \lambda_d(\Phi_m \Phi_m^\top) \leq dB^2.$$

888 At this point, we can apply lemma 3, that ensures with probability $1 - 2\delta$, for k sufficiently large,

$$\left(1 - \sqrt{\frac{5dB^2 \log(d/\delta)}{k\sigma^2}}\right) k\sigma^2 \leq \lambda_{\min}(\Phi^\top \Phi) \leq \lambda_{\max}(\Phi^\top \Phi) \leq \left(1 + \sqrt{\frac{2dB^2 \log(d/\delta)}{k\sigma^2}}\right) k\sigma^2.$$

889 Now, we can fix $\sigma = 1$ as in the assumption and rewrite the projection operator in the following form

$$\Pi_{d,\mu} := k^{-1} \Phi (k^{-1} \Phi^\top \Phi)^{-1} \Phi^\top = k^{-1} \Phi \Phi^\top + k^{-1} \Phi \Delta \Phi^\top,$$

890 where Δ has all the eigenvalues of magnitude less than $\sqrt{\frac{5dB^2 \log(d/\delta)}{k\sigma^2}}$, by the previous result.

891 We now bound the infinity norm of the two terms separately. First,

$$\begin{aligned} \|k^{-1} \Phi \Phi^\top\|_\infty &\stackrel{*}{=} \frac{1}{k} \max_{m=1, \dots, k} \|(\Phi \Phi^\top)_{m\cdot}\|_1 \\ &= \max_{m=1, \dots, k} \frac{1}{k} \sum_{n=1}^k \left| \sum_{i=1}^d \Phi_{mi} \Phi_{ni} \right|, \end{aligned}$$

where $*$ holds since the infinity norm of a matrix corresponds to the maximum 1–norm between its rows. Now, note that, as the rows are independent, each variable $\sum_{i=1}^d \Phi_{mi} \Phi_{ni}$, for $m \neq n$ is a sum of i.i.d. random variables such that

- $\Phi_{mi} \Phi_{ni}$ is bounded in $[-B^2, B^2]$ almost surely.
- The variance is

$$\mathbb{E}[(\Phi_{mi} \Phi_{ni})^2] = \mathbb{E}[\Phi_{mi}^2 \Phi_{ni}^2] = \mathbb{E}[\Phi_{mi}^2] \mathbb{E}[\Phi_{ni}^2] = 1.$$

Therefore, Bernstein's inequality 16 ensures that, w.p. $1 - \delta$

$$\left| \sum_{i=1}^d \Phi_{mi} \Phi_{ni} \right| \leq \sqrt{2d \log(2/\delta)} + \frac{2B^2}{3} \log(2/\delta).$$

Making a union bound over the $k^2 - k$ pairs $m \neq n$, we get, still with probability at least $1 - \delta$,

$$\forall n \neq m \quad \left| \sum_{i=1}^d \Phi_{mi} \Phi_{ni} \right| \leq \sqrt{4d \log(2k/\delta)} + \frac{4B^2}{3} \log(2k/\delta). \quad (10)$$

At this point, we simply have, with probability $1 - \delta$,

$$\begin{aligned} \|k^{-1} \Phi \Phi^\top\|_\infty &= \max_{m=1, \dots, k} \frac{1}{k} \sum_{n=1}^k \left| \sum_{i=1}^d \Phi_{mi} \Phi_{ni} \right| \\ &\leq \frac{dB^2}{k} + \max_{m=1, \dots, k} \frac{1}{k} \sum_{n=1, n \neq m}^k \left| \sum_{i=1}^d \Phi_{mi} \Phi_{ni} \right| \\ &\stackrel{10}{\leq} \frac{dB^2}{k} + \max_{m=1, \dots, k} \frac{1}{k} \sum_{n=1, n \neq m}^k \left(\sqrt{4d \log(2k/\delta)} + \frac{4B^2}{3} \log(2k/\delta) \right) \\ &= \sqrt{4d \log(2k/\delta)} + \frac{4B^2}{3} \log(2k/\delta) + \frac{dB^2}{k}. \end{aligned}$$

For the second term, we have

$$\begin{aligned} \|k^{-1} \Phi \Delta \Phi^\top\|_\infty &\leq k^{-1} \max_{m=1, \dots, k} \sum_{n=1}^k |\langle \Phi_{m\cdot}, (\Delta \Phi^\top)_{\cdot n} \rangle| \\ &\leq k^{-1} \max_{m=1, \dots, k} \sum_{n=1}^k \|\Phi_{m\cdot}\|_2 \|(\Delta \Phi^\top)_{\cdot n}\|_2 \\ &\stackrel{*}{\leq} k^{-1} \max_{m=1, \dots, k} \sum_{n=1}^k \frac{dB^2}{\sqrt{k}} \\ &\leq \frac{dB^2}{\sqrt{k}}, \end{aligned}$$

where $*$ comes from the bound on the eigenvalues of Δ . Putting everything together, we have proved that

$$\|\Pi_{d,\mu}\|_\infty \leq \sqrt{4d \log(2k/\delta)} + \frac{4B^2}{3} \log(2k/\delta) + \frac{dB^2}{k} + \frac{dB^2}{\sqrt{k}} = \sqrt{4d \log(2k/\delta)} + \mathcal{O}(d/\sqrt{k}).$$

902 To show that we cannot go much lower than this quantity, note that, even ignoring the contribution of
 903 Δ we have

$$\|\Pi_{d,\mu}\|_\infty \approx \|k^{-1}\Phi\Phi^\top\|_\infty = \max_{m=1,\dots,k} \frac{1}{k} \sum_{n=1}^k \left| \sum_{i=1}^d \Phi_{mi}\Phi_{ni} \right|.$$

904 Therefore,

$$\begin{aligned} \mathbb{E}[\|\Pi_{d,\mu}\|_\infty] &\approx \mathbb{E} \left[\max_{m=1,\dots,k} \frac{1}{k} \sum_{n=1}^k \left| \sum_{i=1}^d \Phi_{mi}\Phi_{ni} \right| \right] \\ &\geq \max_{m=1,\dots,k} \frac{1}{k} \sum_{n=1}^k \mathbb{E} \left[\left| \sum_{i=1}^d \Phi_{mi}\Phi_{ni} \right| \right] \\ &\geq \max_m \frac{1}{k} \sum_{n=1, n \neq m}^k \Omega(\sqrt{d}) = \Omega(\sqrt{d}). \end{aligned}$$

905 The last passage comes from the fact that, for $n \neq m$, we have the expected value of the modulus a
 906 sum of d independent random variables, which grows as \sqrt{d} . \square

907 **Proposition 21.** *Let μ, ν be two probability distributions on the discrete set \mathcal{X} such that for all*
 908 *$x \in \mathcal{X}, C \geq \frac{\mu(x)}{\nu(x)} \geq c > 0$. Then, $\Lambda_{d,\mu} \leq \frac{C}{c} \Lambda_{d,\nu}$.*

909 *Proof.* The following identity holds for the Lebesgue constant

$$\begin{aligned} \Lambda_{d,\mu} &= \sup_{x \in \mathcal{X}} \int_{\mathcal{X}} \bar{\varphi}_d(x)^\top \bar{\varphi}_d(z) d\mu(z) \\ &= \sup_{x \in \mathcal{X}} \int_{\mathcal{X}} \varphi_d(x)^\top R(\mu)^{-1} R(\mu)^{-\top} d\mu(z) \\ &= \sup_{x \in \mathcal{X}} \int_{\mathcal{X}} |\varphi_d(x)^\top G(\mu)^{-1} \varphi_d(z)| d\mu(z), \end{aligned}$$

where $G(\mu) = \int_{\mathcal{X}} \varphi_d(x) \varphi_d(x)^\top d\mu(x)$ and $R(\mu)$ is its Cholesky factor, such that $R(\mu)^\top R(\mu) = G(\mu)$; here, the second passage comes from the fact that the Cholesky factor of a matrix corresponds to the R factor in the QR factorization, which is the one giving Graham-Schmidt orthogonalization [Quarteroni et al. \[2010\]](#). In fact, letting $\bar{\varphi}_d(x)$ be the basis orthonormalized w.r.t. μ , we have

$$\bar{\varphi}_d(x)^\top \bar{\varphi}_d(z)^\top = \varphi_d(x)^\top G(\mu)^{-1} \varphi_d(z).$$

910 Note that, by absolute continuity, we have, for any $x \in \mathcal{X}$

$$\begin{aligned} \int_{\mathcal{X}} |\varphi_d(x)^\top G(\mu)^{-1} \varphi_d(z)| d\mu(z) &\leq C \int_{\mathcal{X}} |\varphi_d(x)^\top G(\mu)^{-1} \varphi_d(z)| d\nu(z) \\ &\leq C \int_{\mathcal{X}} \left| \varphi_d(x)^\top \left(\int_{\mathcal{X}} \varphi_d(z') \varphi_d(z')^\top d\mu(z') \right)^{-1} \varphi_d(z) \right| d\nu(z) \\ &\leq C \int_{\mathcal{X}} \left| \varphi_d(x)^\top c^{-1} \left(\int_{\mathcal{X}} \varphi_d(z') \varphi_d(z')^\top d\nu(z') \right)^{-1} \varphi_d(z) \right| d\nu(z) \\ &= \frac{C}{c} \int_{\mathcal{X}} \left| \varphi_d(x)^\top \left(\int_{\mathcal{X}} \varphi_d(z') \varphi_d(z')^\top d\nu(z') \right)^{-1} \varphi_d(z) \right| d\nu(z) \\ &= \frac{C}{c} \int_{\mathcal{X}} |\varphi_d(x)^\top G(\nu)^{-1} \varphi_d(z)| d\nu(z). \end{aligned}$$

911 Passing to the supremum, we get the thesis. \square

912 **D Proofs from section 4**

913 **D.1 Lower bound for standard ridge regression**

914 **Lemma 7.** Let $\Pi_{d,\mu}^\lambda$ be the operator defined in this way:

$$\Pi_{d,\mu}^\lambda f := \overline{\varphi}_d(\cdot)^\top \theta_\lambda \quad \theta_\lambda = \arg \min_{\theta} \|f(\cdot) - \overline{\varphi}_d(\cdot)^\top \theta\|_{L^2}^2 + \lambda \|\theta\|_2^2. \quad (11)$$

Then, we have

$$\Pi_{d,\mu}^\lambda f = \frac{\Pi_{d,\mu} f}{1 + \lambda}.$$

915 *Proof.* We start from the definition of θ_λ :

$$\begin{aligned} \theta_\lambda &= \arg \min_{\theta} \|f(\cdot) - \overline{\varphi}_d(\cdot)^\top \theta\|_{L^2}^2 + \lambda \|\theta\|_2^2 \\ &= \arg \min_{\theta} \|\Pi_{d,\mu} f(\cdot) + \zeta_{d,\mu}(\cdot) - \overline{\varphi}_d(\cdot)^\top \theta\|_{L^2}^2 + \lambda \|\theta\|_2^2 \\ &= \arg \min_{\theta} \|\zeta_{d,\mu}\|_{L^2}^2 + \|\Pi_{d,\mu} f(\cdot) - \overline{\varphi}_d(\cdot)^\top \theta\|_{L^2}^2 + \lambda \|\theta\|_2^2, \end{aligned}$$

916 where the last passage comes from Parseval's theorem, as $\zeta_{d,\mu}$ is orthogonal in L^2 to the span of φ_d ,
917 while $\Pi_{d,\mu} f(\cdot), \overline{\varphi}_d(\cdot)^\top \theta$ belongs to this vector space. We then write the operator $\Pi_{d,\mu} f$ explicitly:

$$\begin{aligned} \theta_\lambda &= \arg \min_{\theta} \|\Pi_{d,\mu} f(\cdot) - \overline{\varphi}_d(\cdot)^\top \theta\|_{L^2}^2 + \lambda \|\theta\|_2^2 \\ &= \arg \min_{\theta} \left\| \sum_{i=1}^d \langle f, \overline{\varphi}_i \rangle_{L^2} \overline{\varphi}_i(\cdot) - \overline{\varphi}_d(\cdot)^\top \theta \right\|_{L^2}^2 + \lambda \|\theta\|_2^2 \\ &= \arg \min_{\theta} \sum_{i=1}^d (\langle f, \overline{\varphi}_i \rangle_{L^2} - \theta_i)^2 + \lambda \theta_i^2. \end{aligned}$$

918 The last passage holds from Parseval's theorem since $\overline{\varphi}_i$ are orthonormal in L^2 . Note that, as the θ_i
919 in the last minimization problem are disentangled, we can find as explicit solution

$$\theta_{\lambda,i} = \frac{\langle f, \overline{\varphi}_i \rangle_{L^2}}{1 + \lambda}, \quad \Pi_{d,\mu}^\lambda f = \frac{\Pi_{d,\mu} f}{1 + \lambda}.$$

920 This completes the proof. \square

Lemma 8. Let $\Pi_{d,\mu}^\lambda$ be defined according to equation 11. For every feature map φ_d we have

$$\sup_{f \in L^\infty(\mathcal{X})} \frac{\|\Pi_{d,\mu}^\lambda f - f\|_\infty}{\|\Pi_\infty f - f\|_\infty} \geq \left(\frac{\Lambda_{d,\mu} - 1 - 2\lambda}{1 + \lambda} \right).$$

Proof. By definition of Lebesgue constant, for every $\varepsilon > 0$ there is a function g such that

$$\|\Pi_{d,\mu} g\|_\infty = (\Lambda_{d,\mu} - \varepsilon) \|g\|_\infty.$$

921 Take $f = \Pi_\infty g - g$. We have, by lemma 7,

$$\begin{aligned}
\|\Pi_{d,\mu}^\lambda f - f\|_\infty &= \left\| \frac{\Pi_{d,\mu} f}{1+\lambda} - f \right\|_\infty \\
&= \|(1+\lambda)^{-1} \Pi_{d,\mu} (P_\infty^d g - g) - \Pi_\infty g + g\|_\infty \\
&= \|(1+\lambda)^{-1} \Pi_\infty g - (1+\lambda)^{-1} \Pi_{d,\mu} g - P_\infty^d g + g\|_\infty \\
&= \left\| -(1+\lambda)^{-1} \Pi_{d,\mu} g - \frac{\lambda}{1+\lambda} P_\infty^d g + g \right\|_\infty.
\end{aligned}$$

At this point, note that

$$\|\Pi_\infty g\|_\infty \leq 2\|g\|_\infty,$$

922 as follows from

$$\begin{aligned}
\|\Pi_\infty g\|_\infty &\leq \|g - \Pi_\infty g\|_\infty + \|g\|_\infty \\
&\leq \|g - 0\|_\infty + \|g\|_\infty = 2\|g\|_\infty.
\end{aligned}$$

923 Using this property, we have

$$\begin{aligned}
\|\Pi_{d,\mu}^\lambda f - f\|_\infty &\geq \left\| -(1+\lambda)^{-1} \Pi_{d,\mu} g - \frac{\lambda}{1+\lambda} \Pi_\infty g + g \right\|_\infty \\
&\geq \| -(1+\lambda)^{-1} \Pi_{d,\mu} g \|_\infty - \frac{1+2\lambda}{1+\lambda} \|g\|_\infty.
\end{aligned}$$

924 At this point, using the definition of g ,

$$\begin{aligned}
\| -(1+\lambda)^{-1} \Pi_{d,\mu} g \|_\infty - \frac{1+2\lambda}{1+\lambda} \|g\|_\infty &\geq \left(\frac{\Lambda_{d,\mu}}{1+\lambda} - \varepsilon - \frac{1+2\lambda}{1+\lambda} \right) \|\Pi_\infty g - g\|_\infty \\
&= \left(\frac{\Lambda_{d,\mu}}{1+\lambda} - \varepsilon - \frac{1+2\lambda}{1+\lambda} \right) \|f\|_\infty \\
&\geq \left(\frac{\Lambda_{d,\mu}}{1+\lambda} - \varepsilon - \frac{1+2\lambda}{1+\lambda} \right) \|\Pi_\infty f - f\|_\infty.
\end{aligned}$$

925 The thesis follows letting $\varepsilon \rightarrow 0$. □

926 **Theorem 7.** Let $\hat{\theta}_{n,\text{RIDGE}}$ the output of λ -ridge regression. For any feature map
927 $\varphi_d(\cdot) : \mathcal{X} \rightarrow \mathbb{R}^d$ there is $f \in L^\infty(\mathcal{X})$ such that, for infinite data $\mathcal{E}_\infty(\hat{\theta}_{\infty,\text{RIDGE}}) =$
928 $\Omega\left(\max\left\{(\Lambda_{d,\mu} - 2\lambda)\mathcal{E}_\infty(f), \frac{\lambda}{\lambda+1}\right\}\right).$

929 *Proof.* Let \hat{f}_n be the output of λ -ridge regression, that is the function $\bar{\varphi}_d(\cdot)^\top \hat{\theta}_n$, where

$$\hat{\theta}_n := \arg \min_{\theta \in \mathbb{R}^d} \sum_{t=1}^n (\bar{\varphi}_d(x_t)^\top \theta - y_t)^2 + \lambda n \|\theta\|_2^2 \quad x_t \stackrel{i.i.d.}{\sim} \mu.$$

930 By the uniform law of large numbers, in the limit, the minimizer \hat{f}_n converges to $\Pi_{d,\mu}^\lambda f$, the
931 regularized projection operator is defined as follows

$$\Pi_{d,\mu}^\lambda f(\cdot) := \bar{\varphi}_d(\cdot)^\top \theta_\lambda \quad \theta_\lambda = \arg \min_{\theta} \|f(\cdot) - \bar{\varphi}_d(\cdot)^\top \theta\|_{L_2}^2 + \lambda \|\theta\|_2^2.$$

932 We start showing the $\frac{\lambda}{\lambda+1}$ lower bound. Taking any function in the span of $\varphi_d(\cdot)$ with $\|f\|_\infty = 1$ we
933 have, by lemma 7,

$$\|f - \Pi_{d,\mu}^\lambda f\|_\infty = \|f - (1+\lambda)^{-1} f\|_\infty = \frac{\lambda}{\lambda+1}.$$

934 To show the other part, use lemma 8 to define a function f such that

$$\|\Pi_{d,\mu}^\lambda f - f\|_\infty \geq \left(\frac{\Lambda_{d,\mu} - 2 - 2\lambda}{1 + \lambda} \right) \|\Pi_\infty f - f\|_\infty.$$

935 Replacing $\|\Pi_\infty f - f\|_\infty = \mathcal{E}_\infty(f)$ completes the proof. \square

936 D.2 Proofs from section 4.1

Proposition 22. Let $\alpha \in \mathcal{A}_d^D$ (equation (5)) and $\Pi_{\alpha,\mu}^{Ridge}$ be defined according to equation (4). Then,

$$\|f(\cdot) - \Pi_{\alpha,\mu}^{Ridge} f(\cdot)\|_\infty \leq (1 + \Lambda_{\alpha,\mu}) \mathcal{E}_\infty(f).$$

937 *Proof.* We have

$$\begin{aligned} \|f - P_{\mathbf{h}}^2 f\|_\infty &= \|\Pi_{d,\infty} f + \xi_d - P_{\mathbf{h}}^2[\Pi_{d,\infty} f + \xi_d]\|_\infty \\ &\stackrel{*}{=} \|\xi_d - P_{\mathbf{h}}^2[\xi_d]\|_\infty \\ &\leq \|\xi_d\|_\infty + \|P_{\mathbf{h}}^2[\xi]\|_\infty \\ &= \|\xi_d\|_\infty + \Lambda_{\mathbf{h}} \|\xi_d\|_\infty. \end{aligned}$$

938 Here, the key passage (*) holds since $P_{\mathbf{h}}^2$ by definition ?? is the identity over everything in the span
939 of the first d features, so $\Pi_{d,\infty} f$ in particular. \square

Theorem 9. Let assumption 1 hold. Then, for any $\delta > 0$, with probability $1 - \delta$,

$$\mathcal{E}_\infty(\hat{\theta}_{n,\alpha}) \leq (1 + \Lambda_{\alpha,\mu_n}) \mathcal{E}_\infty(f) + \frac{\sigma \hat{\varphi}_{2,D} \sqrt{2 \log(2\mathcal{X}/\delta)}}{\sqrt{n}}.$$

Proof. Let $\hat{\theta}_n$ the estimator corresponding to P_α^2 in the parameterization of $\hat{\varphi}_D(\cdot)$, so that

$$\hat{\varphi}_d(\cdot)^\top \hat{\theta}_n = \hat{P}_\alpha^2 \mathbf{f} =: \hat{f}_n(\cdot).$$

940 The following decomposition holds:

$$\begin{aligned} \|f(\cdot) - \hat{f}_n(\cdot)\|_\infty &\leq \|f(\cdot) - \hat{P}_\alpha^2 f(\cdot)\|_\infty + \|\hat{P}_\alpha^2 f(\cdot) - \hat{f}_n(\cdot)\|_\infty \\ &\leq (1 + \hat{\Lambda}_\alpha) \mathcal{E}_\infty(f) + \|\hat{P}_\alpha^2 f(\cdot) - \hat{f}_n(\cdot)\|_\infty. \end{aligned}$$

941 where we have applied proposition 8 for $\mu_n(\cdot)$. Let us focus on the second term. As in the proof of
942 the previous theorems, we call θ_* the vector corresponding to the orthogonal projection over $\hat{\varphi}_D(\cdot)$
943 so that we have, for every $x \in \mathcal{X}$

$$\begin{aligned} \hat{f}_n(x) &= \hat{\varphi}_D(x)^\top I_\alpha \frac{1}{n} \sum_{t=1}^n y_t \hat{\varphi}_D(x_t) \\ &= \hat{\varphi}_D(x)^\top I_\alpha \frac{1}{n} \sum_{t=1}^n (\hat{\varphi}_D(x_t)^\top \theta_* + \zeta_D(x_t) + \eta_t) \hat{\varphi}_D(x_t) \\ &= \hat{\varphi}_D(x)^\top I_\alpha \frac{1}{n} \sum_{t=1}^n \hat{\varphi}_D(x_t) \hat{\varphi}_D(x_t)^\top \theta_* \\ &\quad + \hat{\varphi}_D(x)^\top I_\alpha \frac{1}{n} \sum_{t=1}^n \zeta_D(x_t) \hat{\varphi}_D(x_t) \\ &\quad + \hat{\varphi}_D(x)^\top I_\alpha \frac{1}{n} \sum_{t=1}^n \eta_t \hat{\varphi}_D(x_t). \end{aligned}$$

944 By orthogonality, the first term corresponds to

$$\widehat{\varphi}_D(x)^\top I_\alpha \underbrace{\frac{1}{n} \sum_{t=1}^n \widehat{\varphi}_D(x_t) \widehat{\varphi}_D(x_t)^\top}_{I_D} \theta_\star = \widehat{\varphi}_D(x)^\top I_\alpha \theta_\star = P_\alpha^2(x).$$

945 The second term is

$$\widehat{\varphi}_D(x)^\top I_\alpha \frac{1}{n} \sum_{t=1}^n \zeta_D(x_t) \widehat{\varphi}_D(x_t) = \widehat{\varphi}_D(x)^\top I_\alpha \underbrace{\int_{\mathcal{X}} \zeta_D(z) \widehat{\varphi}_D(z) d\mu_n(z)}_{\mathbf{0} \text{ vector}} = 0,$$

946 by definition of orthogonal projection. The third term is

$$\widehat{\varphi}_D(x)^\top I_\alpha \frac{1}{n} \sum_{t=1}^n \eta_t \widehat{\varphi}_D(x_t),$$

947 which can be bounded as the corresponding terms in theorems 4 and 5: as η_t are independent and
 948 σ -subgaussian, Lemma 5.4 and Theorem 5.3 from [Lattimore and Szepesvári \[2020\]](#) ensure that,
 949 with probability at least $1 - 2\delta$

$$\begin{aligned} \left| \widehat{\varphi}_d(x)^\top I_\alpha n^{-1} \sum_{t=1}^n \widehat{\varphi}_d(x_t) \eta_t \right| &\leq \sqrt{2\sigma^2 n^{-1} \sum_{t=1}^n (\widehat{\varphi}_d(x)^\top I_\alpha \widehat{\varphi}_d(x_t))^2 \log(1/\delta)} \\ &= \sqrt{2\sigma^2 n^{-1} \|\widehat{\varphi}_d(x)\|_2^2 \log(1/\delta)} \\ &= \sqrt{2 \log(1/\delta) \sigma n^{-1/2}} \|\widehat{\varphi}_d(x)\|_2. \end{aligned}$$

950 Where the only difference w.r.t. the other proofs is the presence of I_α , which is erased after the
 951 first step since, being $\alpha \in \mathcal{A}_d^D$, its norm is ≤ 1 . This proves that the last term is bounded by
 952 $\sqrt{2 \log(1/\delta) \sigma n^{-1/2}} \widehat{\varphi}_{2,D}$. Making a union bound over \mathcal{X} gives, w.p. $1 - \delta$,

$$\sup_{x \in \mathcal{X}} |\widehat{P}_\alpha^2 f(x) - \widehat{f}_n(x)| \leq \sqrt{2 \log(1/\delta) \sigma n^{-1/2}} \widehat{\varphi}_{2,D}.$$

953 Putting everything together, we have proved that

$$\mathcal{E}_\infty(\widehat{\theta}_{n,\alpha}) \leq \|f(\cdot) - \widehat{f}_n(\cdot)\|_\infty \leq (1 + \widehat{\Lambda}_\alpha) \mathcal{E}_\infty(f) + \frac{\sigma \widehat{\varphi}_{2,D} \sqrt{2 \log(2\mathcal{X}/\delta)}}{\sqrt{n}}.$$

954

□

955 **Proposition 23.** *Under assumption 2 we have, with probability $1 - \delta$ for every $\alpha \in$
 956 \mathcal{A}_d^D at the same time, $|\widehat{\varphi}_{D,2} - \overline{\varphi}_{D,2}| \leq \widetilde{O}(\overline{\varphi}_{D,2}^2 \sqrt{\log(1/\delta)/n})$, and $|\Lambda_{\alpha,\mu_n} - \Lambda_{\alpha,\mu}| \leq$
 957 $\widetilde{O}\left(\frac{\sqrt{d} \overline{\varphi}_{D,2}^2 \sqrt{\log(1/\delta)}}{\sqrt{n}} + \frac{\sqrt{d} \overline{\varphi}_{D,2}^3 \log(1/\delta)}{n}\right)$.²*

958 We prove this theorem for a generic $d \in \mathbb{N}$. The result follows for $d = D$.

959 We define $V_n := \frac{1}{n} \sum_{t=1}^n \overline{\varphi}_d(x_t) \overline{\varphi}_d(x_t)^\top$. Let $\widehat{\varphi}_d(\cdot)$ the basis obtained from φ_d by Gram-Schmidt
 960 orthogonalization w.r.t. μ_n , the empirical distribution of the $\{x_t\}_t$. As in the main paper, we let
 961 $R_n = \text{Chol}(V_n)$ and, since the Cholesky factor corresponds to the matrix given by Graham Schmidt
 962 orthogonalization (proposition 3.4 in [Quarteroni et al. \[2010\]](#)),

$$\overline{\varphi}_d(x_t) = R_n^\top \widehat{\varphi}_d(x_t) \quad \widehat{\varphi}_d(x_t) = R_n^{-\top} \overline{\varphi}_d(x_t). \quad (12)$$

²The statement of this theorem is slightly different from the one in the main paper of the submission, as we have made the orders of magnitude more precise.

so that, under this convenient normalization, we can pass from $\bar{\varphi}_d(x_t)$ to $\hat{\varphi}_d(x_t)$ through a matrix that is exactly the Cholesky factor of V_n . In this setting, Theorem 2.1. in [Drmač et al. \[1994\]](#), which provides a stability result for the Cholesky decomposition which, combined with our theorem gives

$$1 - \mathcal{O}\left(\bar{\varphi}_{d,2}\sqrt{\log(1/\delta)/n\log(d)}\right) \leq \lambda_{\min}(R_n) \leq \lambda_{\max}(R_n) \leq 1 + \mathcal{O}\left(\bar{\varphi}_{d,2}\sqrt{\log(1/\delta)/n\log(d)}\right) \quad (13)$$

We can now proceed with the proof.

Proof. Bounding norm difference

We have to measure

$$\sup_{x \in \mathcal{X}} \|\hat{\varphi}_d(x) - \bar{\varphi}_d(x)\|_2.$$

As we said, the relation between the two is $\bar{\varphi}_d(x) = R_n^\top \hat{\varphi}_d(x)$ which we can also write as $R_n^{-\top} \bar{\varphi}_d(x) = \hat{\varphi}_d(x)$, so that

$$\sup_{x \in \mathcal{X}} \|\hat{\varphi}_d(x) - \bar{\varphi}_d(x)\|_2 = \sup_{x \in \mathcal{X}} \|(I_d - R_n^{-\top})\bar{\varphi}_d(x)\|_2.$$

At this point, equation (13) ensures that $\|I_d - R_n^{-\top}\|_{2 \rightarrow 2} = \mathcal{O}\left(\bar{\varphi}_{d,2}\sqrt{\log(1/\delta)/n\log(d)}\right)$, so we get

$$\sup_{x \in \mathcal{X}} \|\hat{\varphi}_d(x) - \bar{\varphi}_d(x)\|_2 \leq \mathcal{O}\left(\bar{\varphi}_{d,2}^2\sqrt{\log(1/\delta)/n\log(d)}\right). \quad (14)$$

A simple yet useful consequence of this result is

$$|\bar{\varphi}_{d,2} - \hat{\varphi}_{d,2}| = \sup_{x \in \mathcal{X}} \|\hat{\varphi}_d(x)\|_2 - \sup_{x \in \mathcal{X}} \|\bar{\varphi}_d(x)\|_2 \quad (15)$$

$$\leq \sup_{x \in \mathcal{X}} \|\hat{\varphi}_d(x)\|_2 - \|\bar{\varphi}_d(x)\|_2 \quad (16)$$

$$\leq \sup_{x \in \mathcal{X}} \|\hat{\varphi}_d(x) - \bar{\varphi}_d(x)\|_2 \quad (17)$$

$$= \mathcal{O}\left(\bar{\varphi}_{d,2}^2\sqrt{\log(1/\delta)/n\log(d)}\right) \quad (18)$$

Lebesgue constants difference

Let us bound the distance between the estimated and the true Lebesgue constant, for any $\alpha \in \mathcal{A}_d^D$,

$$\begin{aligned}
|\Lambda_{\alpha, \mu_n} - \Lambda_{\alpha, \mu}| &= \left| \sup_{x \in \mathcal{X}} \frac{1}{n} \sum_{t=1}^n \left| \sum_{i=1}^d \alpha_i \widehat{\varphi}_i(x) \widehat{\varphi}_i(x_t) \right| - \sup_{x \in \mathcal{X}} \int_{\mathcal{X}} \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(z) \right| d\mu(z) \right| \\
&\leq \sup_{x \in \mathcal{X}} \left| \frac{1}{n} \sum_{t=1}^n \left| \sum_{i=1}^d \alpha_i \widehat{\varphi}_i(x) \widehat{\varphi}_i(x_t) \right| - \int_{\mathcal{X}} \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(z) \right| d\mu(z) \right| \\
&\leq \sup_{x \in \mathcal{X}} \left| \frac{1}{n} \sum_{t=1}^n \left| \sum_{i=1}^d \alpha_i \widehat{\varphi}_i(x) \widehat{\varphi}_i(x_t) \right| - \frac{1}{n} \sum_{t=1}^n \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(x_t) \right| \right| \\
&\quad + \sup_{x \in \mathcal{X}} \left| \frac{1}{n} \sum_{t=1}^n \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(x_t) \right| - \int_{\mathcal{X}} \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(z) \right| d\mu(z) \right| \\
&= \sup_{x \in \mathcal{X}} \left| \frac{1}{n} \sum_{t=1}^n \left| \sum_{i=1}^d \alpha_i \widehat{\varphi}_i(x) \widehat{\varphi}_i(x_t) \right| - \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(x_t) \right| \right| \\
&\quad + \sup_{x \in \mathcal{X}} \left| \frac{1}{n} \sum_{t=1}^n \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(x_t) \right| - \int_{\mathcal{X}} \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(z) \right| d\mu(z) \right|.
\end{aligned}$$

976 In the following, we call

$$\text{First term} := \sup_{\alpha \in \mathcal{A}_d^D, x \in \mathcal{X}} \left| \frac{1}{n} \sum_{t=1}^n \left| \sum_{i=1}^d \alpha_i \widehat{\varphi}_i(x) \widehat{\varphi}_i(x_t) \right| - \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(x_t) \right| \right|$$

977 and

$$\text{Second term} := \sup_{\alpha \in \mathcal{A}_d^D, x \in \mathcal{X}} \left| \frac{1}{n} \sum_{t=1}^n \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(x_t) \right| - \int_{\mathcal{X}} \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(z) \right| d\mu(z) \right|.$$

978 **Bound the first term.**

979 Fix $\alpha \in \mathcal{A}_d^D$,

$$\begin{aligned}
\text{First part} &= \frac{1}{n} \sum_{t=1}^n \left| \sum_{i=1}^d \alpha_i \widehat{\varphi}_i(x) \widehat{\varphi}_i(x_t) \right| - \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(x_t) \right| \\
&\leq \frac{1}{n} \sum_{t=1}^n \left| \sum_{i=1}^d \alpha_i \widehat{\varphi}_i(x) \widehat{\varphi}_i(x_t) - \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(x_t) \right| \\
&= \frac{1}{n} \sum_{t=1}^n \left| \widehat{\varphi}_d(x)^\top I_\alpha \widehat{\varphi}_d(x_t) - \overline{\varphi}_d(x)^\top I_\alpha \overline{\varphi}_d(x_t) \right|.
\end{aligned}$$

980 Where, $I_\alpha = \text{diag}(\alpha)$. At this point, we can replace the result of equation 12: getting

$$\begin{aligned}
\text{First part} &\leq \frac{1}{n} \sum_{t=1}^n |\widehat{\varphi}_d(x)^\top I_\alpha \widehat{\varphi}_d(x_t) - \overline{\varphi}_d(x)^\top I_\alpha \overline{\varphi}_d(x_t)| \\
&= \frac{1}{n} \sum_{t=1}^n |\widehat{\varphi}_d(x)^\top I_\alpha \widehat{\varphi}_d(x_t) - \widehat{\varphi}_d(x)^\top R_n I_\alpha R_n^\top \widehat{\varphi}_d(x_t)| \\
&\leq \frac{1}{n} \sum_{t=1}^n |\widehat{\varphi}_d(x)^\top (I_\alpha - R_n I_\alpha R_n^\top) \widehat{\varphi}_d(x_t)| \\
&\leq \frac{1}{n} \sum_{t=1}^n \|\widehat{\varphi}_d(x)\|_2 \|I_\alpha - R_n I_\alpha R_n^\top\|_2 \|\widehat{\varphi}_d(x_t)\|_2.
\end{aligned}$$

981 This formulation allows us to apply equation 13: As I_α is diagonal matrix with elements in $[0, 1]$, we
982 have

$$\|I_\alpha - R_n I_\alpha R_n^\top\|_2 = \mathcal{O}\left(\overline{\varphi}_{d,2} \sqrt{\log(1/\delta)/n} \log(d)\right).$$

983 This gives the following

$$\begin{aligned}
\text{First part} &\leq \mathcal{O}\left(\frac{1}{n} \sum_{t=1}^n \|\widehat{\varphi}_d(x)\|_2 \|H - R_n H R_n^\top\|_2 \|\widehat{\varphi}_d(x_t)\|_2\right) \\
&\leq \mathcal{O}\left(\frac{\overline{\varphi}_{d,2} \sqrt{\log(1/\delta)} \log(d)}{\sqrt{n}} \frac{1}{n} \sum_{t=1}^n \|\widehat{\varphi}_d(x)\|_2 \|\widehat{\varphi}_d(x_t)\|_2\right) \\
&\leq \mathcal{O}\left(\frac{\overline{\varphi}_{d,2} \widehat{\varphi}_{d,2} \sqrt{\log(1/\delta)} \log(d)}{\sqrt{n}} \frac{\sum_{t=1}^n \|\widehat{\varphi}_d(x_t)\|_2}{n}\right) \\
&\leq \mathcal{O}\left(\frac{\overline{\varphi}_{d,2} \widehat{\varphi}_{d,2} \sqrt{\log(1/\delta)} \log(d)}{\sqrt{n}} \frac{\sqrt{n \sum_{t=1}^n \|\widehat{\varphi}_d(x_t)\|_2^2}}{n}\right) \\
&= \mathcal{O}\left(\frac{\overline{\varphi}_{d,2} \widehat{\varphi}_{d,2} \sqrt{\log(1/\delta)} \log(d)}{\sqrt{n}} \frac{\sqrt{n^2 d}}{n}\right) \\
&= \mathcal{O}\left(\frac{\sqrt{d} \overline{\varphi}_{d,2} \widehat{\varphi}_{d,2} \sqrt{\log(1/\delta)} \log(d)}{\sqrt{n}}\right).
\end{aligned}$$

984 Here, the first equality is due to the fact that, being $\widehat{\varphi}_d$ orthonormal w.r.t. μ_n , we have
985 $\sum_{t=1}^n \|\widehat{\varphi}_d(x_t)\|_2^2 = nd$. This holds uniformly for every α , as we have only used the fact that
986 $\|I_\alpha\|_2 \leq 1$.

987 **Bounding the second term.**

988 The second term corresponds to

$$\text{Second term} = \sup_{x \in \mathcal{X}} \left| \frac{1}{n} \sum_{t=1}^n \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(x_t) \right| - \int_{\mathcal{X}} \left| \sum_{i=1}^d \alpha_i \overline{\varphi}_i(x) \overline{\varphi}_i(z) \right| d\mu(z) \Big|.$$

989 First, we fix $x \in \mathcal{X}$ and $\alpha \in \mathcal{A}_d^D$ and use the scalar product to write it as

$$\left| \frac{1}{n} \sum_{t=1}^n |\overline{\varphi}_d(x)^\top I_\alpha \overline{\varphi}_d(x_t)| - \int_{\mathcal{X}} |\overline{\varphi}_d(x)^\top I_\alpha \overline{\varphi}_d(z)| d\mu(z) \right|. \quad (19)$$

990 Note that by definition

$$\mathbb{E}[|\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(x_t)|] = \int_{\mathcal{X}} |\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(z)| d\mu(z).$$

991 Moreover,

$$\begin{aligned} \text{Var}(|\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(x_t)|) &\leq \mathbb{E}[|\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(x_t)|^2] \\ &= \mathbb{E}[\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(x_t) \bar{\varphi}_d(x_t)^\top I_\alpha \bar{\varphi}_d(x)] \\ &= \bar{\varphi}_d(x)^\top I_\alpha \underbrace{\mathbb{E}[\bar{\varphi}_d(x_t) \bar{\varphi}_d(x_t)^\top]}_{=I_d} I_\alpha \bar{\varphi}_d(x) \\ &= \bar{\varphi}_d(x)^\top I_\alpha^2 \bar{\varphi}_d(x) \\ &\leq \bar{\varphi}_2^2, \end{aligned}$$

992 where the last step comes from the fact that $I_\alpha^2 \preceq I_d$. For the same reason, we also have
 993 $|\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(x_t)| \leq \bar{\varphi}_2^2$ almost surely. These three results allow us to apply Bernstein's inequality
 994 16 for

$$\begin{aligned} 995 \quad &\bullet X_t = |\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(x_t)| - \mathbb{E}[|\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(x_t)|]. \\ 996 \quad &\bullet \sigma^2 = \sum_{t=1}^n \text{Var}(|\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(x_t)|) \leq n\bar{\varphi}_2^2. \\ 997 \quad &\bullet B = \bar{\varphi}_2^2. \end{aligned}$$

998 This gives, with probability at least $1 - \delta$,

$$\left| \sum_{t=1}^n X_t \right| \leq \sqrt{2n\bar{\varphi}_2^2 \log(2/\delta)} + \frac{2\bar{\varphi}_2^2}{3} \log(2/\delta).$$

999 So, we can bound equation 19, which corresponds to $\frac{1}{n} |\sum_{t=1}^n X_t|$, as follows.

$$\left| \frac{1}{n} \sum_{t=1}^n |\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(x_t)| - \int_{\mathcal{X}} |\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(z)| d\mu(z) \right| \leq \sqrt{\frac{2\bar{\varphi}_2^2 \log(2/\delta)}{n}} + \frac{2\bar{\varphi}_2^2}{3n} \log(2/\delta).$$

1000 The former holds for any fixed $\alpha \in \mathcal{A}_d^D$. To have a uniform bound, let

$$\mathcal{A}' = \varepsilon - \text{Cover of } \mathcal{A}_d^D \quad \varepsilon = (n\bar{\varphi}_{d,2})^{-1},$$

1001 so that $\log |\mathcal{A}'| \leq d \log(n\bar{\varphi}_{d,2})$. Making a union bound gives, $\forall \alpha \in \mathcal{A}'$

$$\begin{aligned} &\left| \frac{1}{n} \sum_{t=1}^n |\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(x_t)| - \int_{\mathcal{X}} |\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(z)| d\mu(z) \right| \\ &\leq \sqrt{\frac{2d\bar{\varphi}_{d,2}^2 \log(2n\bar{\varphi}_{d,2}/\delta)}{n}} + \frac{2d\bar{\varphi}_2^2}{3n} \log(2n\bar{\varphi}_{d,2}/\delta). \end{aligned}$$

1002 To pass to the general case, note that for every $\alpha \in \mathcal{A}_d^D$ there is $\alpha' \in \mathcal{A}'$ such that
 1003 $|\frac{1}{n} \sum_{t=1}^n |\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(x_t)| - \int_{\mathcal{X}} |\bar{\varphi}_d(x)^\top I_\alpha \bar{\varphi}_d(z)| d\mu(z)|$ changes no more than $2\bar{\varphi}_{d,2}$ between
 1004 the two, by definition of ε -cover. Therefore, we have, with probability at least $1 - \delta$ over all $\alpha \in \mathcal{A}_d^D$
 1005 at the same time

$$\begin{aligned} & \left| \frac{1}{n} \sum_{t=1}^n |\bar{\varphi}_d(x)^\top I_{\alpha} \bar{\varphi}_d(x_t)| - \int_{\mathcal{X}} |\bar{\varphi}_d(x)^\top I_{\alpha} \bar{\varphi}_d(z)| d\mu(z) \right| \\ & \leq \sqrt{\frac{2d\bar{\varphi}_{d,2}^2 \log(2n\bar{\varphi}_{d,2}/\delta)}{n}} + \frac{2d\bar{\varphi}_2^2}{3n} \log(2n\bar{\varphi}_{d,2}/\delta) + 2\bar{\varphi}_{d,2}. \end{aligned}$$

1006 This means,

$$\text{Second term} \leq \tilde{\mathcal{O}} \left(\sqrt{\frac{d\bar{\varphi}_{d,2}^2 \log(1/\delta)}{n}} + \frac{d\bar{\varphi}_2^2}{n} \log(1/\delta) \right).$$

1007 **Putting the two results together.** By the two bounds that we got for the two terms, it follows with
1008 probability at least $1 - \delta$

$$\sup_{\alpha \in \mathcal{A}_d^D} |\Lambda_{\alpha, \mu_n} - \Lambda_{\alpha, \mu}| \leq \tilde{\mathcal{O}} \left(\frac{\sqrt{d\bar{\varphi}_{d,2}\hat{\varphi}_{d,2}\sqrt{\log(1/\delta)}}}{\sqrt{n}} + \sqrt{\frac{d\bar{\varphi}_{d,2}^2 \log(1/\delta)}{n}} + \frac{d\bar{\varphi}_2^2}{n} \log(1/\delta) \right).$$

1009 To end the proof, note that, using equation 18, the difference between $\bar{\varphi}_{d,2}$ and $\hat{\varphi}_{d,2}$ is of order
1010 $\bar{\varphi}_{d,2}^2 \sqrt{\log(1/\delta)/n}$, so that

$$\begin{aligned} \frac{\sqrt{d\bar{\varphi}_{d,2}\hat{\varphi}_{d,2}\sqrt{\log(1/\delta)}}}{\sqrt{n}} & \leq \frac{\sqrt{d\bar{\varphi}_{d,2}(\bar{\varphi}_{d,2} + \bar{\varphi}_{d,2}^2 \sqrt{\log(1/\delta)/n})\sqrt{\log(1/\delta)}}}{\sqrt{n}} \\ & = \frac{\sqrt{d\bar{\varphi}_{d,2}^2 \sqrt{\log(1/\delta)}}}{\sqrt{n}} + \frac{\sqrt{d\bar{\varphi}_{d,2}^3 \log(1/\delta)}}{n}. \end{aligned}$$

1011 Finally, note that, as $\sqrt{d} \leq \bar{\varphi}_{d,2}$, the term $\frac{\sqrt{d\bar{\varphi}_{d,2}^3 \log(1/\delta)}}{n}$ dominates over $\frac{d\bar{\varphi}_2^2}{n} \log(1/\delta)$ that we had
1012 before. \square

1013 D.3 Proofs about gradient method

1014 **Proposition 24.** The function $J : \mathcal{A}_d^D \rightarrow (0, +\infty)$ given by $J(\alpha) := \Lambda_{\alpha, \mu_n}$ is convex in α .

Proof. By definition,

$$J(\alpha) = \|M(\alpha)\|_{\infty},$$

1015 where $M(\alpha) = \frac{1}{n} \sum_{i=1}^d \alpha_i \hat{\varphi}_i(x) \hat{\varphi}_i(x_t)$. Therefore, in particular

$$J(\alpha) = \sup_{x \in \mathcal{X}, \mathbf{f} \in \{-1, 1\}^n} \left| \frac{1}{n} \sum_{i=1}^d \alpha_i \hat{\varphi}_i(x) \hat{\varphi}_i(x_t) \mathbf{f} \right|.$$

1016 This function is convex, being the supremum of a family of linear functions $\frac{1}{n} \sum_{i=1}^d \alpha_i \hat{\varphi}_i(x) \hat{\varphi}_i(x_t)$
1017 in α . \square

1018 **Theorem 12.** Fix $\epsilon > 0$. Algorithm 1, after a number of iterations $I = \tilde{\mathcal{O}}(\epsilon^{-2} \hat{\varphi}_{2,D}^2 (D - d))$ outputs
1019 $\alpha^{(I)} \in \mathcal{A}_d^D$ such that $J(\alpha^{(I)}) \leq \inf_{\alpha \in \mathcal{A}_d^D} J(\alpha) + \epsilon$.

1020 *Proof.* The first step of this proof consists in finding an upper bound for any sub-gradient of α . As
1021 we said,

$$J(\alpha) = \sup_{x \in \mathcal{X}, \mathbf{f} \in \{-1, 1\}^n} \left| \frac{1}{n} \sum_{i=1}^d \alpha_i \hat{\varphi}_i(x) \hat{\varphi}_i(x_t) \mathbf{f} \right| = \sup_{x \in \mathcal{X}, \mathbf{f} \in \{-1, 1\}^n} \left| \frac{1}{n} \hat{\varphi}_D(x)^\top I_\alpha \hat{\Phi}^\top \mathbf{f} \right|,$$

1022 where $I_\alpha = \text{diag}(\alpha)$ is a $D \times D$ diagonal matrix and $\hat{\Phi}$ is the $n \times d$ matrix having, as rows, $\hat{\varphi}_D(x_t)$
 1023 for each $t = 1, \dots, n$. At this point note that, by duality

$$J(\alpha) = \sup_{x \in \mathcal{X}, \mathbf{f} \in \{-1, 1\}^n} \left| \frac{1}{n} \hat{\varphi}_D(x)^\top I_\alpha \hat{\Phi}^\top \mathbf{f} \right| = \sup_{x \in \mathcal{X}} \frac{1}{n} \sum_{t=1}^n |\{\hat{\varphi}_D(x)^\top I_\alpha \hat{\Phi}^\top\}_t|,$$

1024 where $\{\}_t$ denotes the t -th component of $\hat{\varphi}_D(x)^\top I_\alpha \hat{\Phi}^\top$, which is a $1 \times n$ row vector. Now,
 1025 assuming³ that the supremum is obtained by just one value $x_* \in \mathcal{X}$, we can compute the gradient as

$$\begin{aligned} \nabla J(\alpha) &= \nabla \frac{1}{n} \sum_{t=1}^n |\{\hat{\varphi}_D(x_*)^\top I_\alpha \hat{\Phi}^\top\}_t| \\ &= \frac{1}{n} \sum_{t=1}^n \text{sign}(\{\hat{\varphi}_D(x_*)^\top I_\alpha \hat{\Phi}^\top\}_t) \nabla \{\hat{\varphi}_D(x_*)^\top I_\alpha \hat{\Phi}^\top\}_t \\ &= \frac{1}{n} \sum_{t=1}^n \text{sign}(\{\hat{\varphi}_D(x_*)^\top I_\alpha \hat{\Phi}^\top\}_t) \hat{\varphi}_D(x_*)^\top \odot \{\hat{\Phi}\}_t^\top. \end{aligned}$$

1026 In the last line, we have used the Hadamard product \odot , that is defined, for two vectors of length D
 1027 like $\hat{\varphi}_D(x_*)^\top$ and $\{\hat{\Phi}\}_t^\top$, as the component-wise product, generating another vector of length D .

1028 Now, we are going to bound the two-norm of this gradient:

$$\begin{aligned} \|\nabla J(\alpha)\|_2^2 &= \sum_{i=1}^D \left\{ \frac{1}{n} \sum_{t=1}^n \text{sign}(\{\hat{\varphi}_D(x_*)^\top I_\alpha \hat{\Phi}^\top\}_t) \hat{\varphi}_D(x_*)^\top \odot \{\hat{\Phi}\}_t^\top \right\}_i^2 \\ &\leq \sum_{i=1}^D \frac{1}{n} \sum_{t=1}^n \left\{ \text{sign}(\{\hat{\varphi}_D(x_*)^\top I_\alpha \hat{\Phi}^\top\}_t) \hat{\varphi}_D(x_*)^\top \odot \{\hat{\Phi}\}_t^\top \right\}_i^2 \\ &\leq \sum_{i=1}^D \frac{1}{n} \sum_{t=1}^n \left\{ \hat{\varphi}_D(x_*)^\top \odot \{\hat{\Phi}\}_t^\top \right\}_i^2 \\ &\leq \sum_{i=1}^D \frac{1}{n} \sum_{t=1}^n \hat{\varphi}_i(x_*)^2 \hat{\varphi}_i(x_t)^2 \\ &\leq \sum_{i=1}^D \hat{\varphi}_i(x_*)^2 \underbrace{\frac{1}{n} \sum_{t=1}^n \hat{\varphi}_i(x_t)^2}_{=1} = \hat{\varphi}_{D,2}^2, \end{aligned}$$

1029 where the last passage holds since the features $\hat{\varphi}_i(\cdot)$ are orthonormal w.r.t. $\mu_n(\cdot)$. Under these
 1030 assumption, namely

- 1031 1. J is convex
- 1032 2. Each sub-gradient has norm bounded by $G := \hat{\varphi}_{2,D}$
- 1033 3. The diameter of the optimization space \mathcal{H}_d^D is $R := \sqrt{D-d}$

³if there are ties, the argument applied to each of them still holds bounding the norm of the sub-gradient

equation (3) on [Boyd et al. \[2003\]](#) guarantees that running the subgradient method for I iterations with step size

$$\gamma_\ell = \frac{R}{G\sqrt{\ell+1}}$$

1034 (corresponding to line 7), achieves suboptimality ϵ_I bounded by

$$\epsilon_I \leq \frac{R^2 + G^2 \sum_{\ell=1}^I \gamma_\ell^2}{2 \sum_{\ell=1}^I \gamma_\ell} \leq \frac{R^2 + R^2(\log(I) + 1)}{(R/G)\sqrt{I}} \leq \frac{2RG \log(I)}{\sqrt{I}} = \frac{2\hat{\varphi}_{2,D}\sqrt{D-d}\log(I)}{\sqrt{I}}.$$

Therefore, a number of iterations $I = 4\epsilon^{-2}\hat{\varphi}_{2,D}^2(D-d)\log^3(4\hat{\varphi}_{2,D}^2(D-d))$ allows to ensure $\epsilon_I \leq \epsilon$. In this way, we have

$$\hat{\Lambda}_{\alpha^*} - \inf_{\alpha \in \mathcal{A}_d^D} \hat{\Lambda}_{\alpha} = J(\alpha^{(I)}) - \inf_{\alpha \in \mathcal{A}_d^D} J(\alpha) \leq \epsilon_I \leq \epsilon,$$

1035 which completes the proof. □

Theorem 13. *Let Assumptions 1 and 2 hold and fix $\delta > 0$. Then, with probability $1 - \delta$,*

$$\mathcal{E}_{\infty}(\hat{\theta}_{n,BWR}) \leq (1 + \Lambda_{\mu}^{Oracle})\mathcal{E}_{\infty}(f) + \tilde{\mathcal{O}}\left(\frac{\bar{\varphi}_{2,D}\sqrt{D\log(|\mathcal{X}|/\delta)}}{\sqrt{n}} + \frac{\bar{\varphi}_{2,D}^2\log(|\mathcal{X}|/\delta)}{n}\right).$$

1036 *Proof.* By theorem 9 and the definition of $\hat{\theta}_{n,BWR}$,

$$\mathcal{E}_{\infty}(\hat{\theta}_{n,BWR}) \leq (1 + \Lambda_{\alpha^{(I)},\mu_n})\mathcal{E}_{\infty}(f) + \frac{\sigma\hat{\varphi}_{2,D}\sqrt{2\log(2\mathcal{X}/\delta)}}{\sqrt{n}}. \quad (20)$$

1037 By theorem 12, for fixed ϵ , we have $\Lambda_{\alpha^{(I)},\mu_n} \leq \min_{\alpha \in \mathcal{A}_d^D} \Lambda_{\alpha,\mu_n} + \epsilon$. Moreover, note that

$$\begin{aligned} \Lambda_{\mu}^{Oracle} &= \Lambda_{\alpha_{\mu}^{Oracle},\mu} \\ &\geq \Lambda_{\alpha_{\mu}^{Oracle},\mu_n} - \tilde{\mathcal{O}}\left(\frac{\bar{\varphi}_{2,D}\sqrt{D\log(|\mathcal{X}|/\delta)}}{\sqrt{n}} + \frac{\bar{\varphi}_{2,D}^2\log(|\mathcal{X}|/\delta)}{n}\right) \\ &\geq \min_{\alpha \in \mathcal{A}_d^D} \Lambda_{\alpha,\mu_n} - \tilde{\mathcal{O}}\left(\frac{\bar{\varphi}_{2,D}\sqrt{D\log(|\mathcal{X}|/\delta)}}{\sqrt{n}} + \frac{\bar{\varphi}_{2,D}^2\log(|\mathcal{X}|/\delta)}{n}\right) \\ &\geq \Lambda_{\alpha^{(I)},\mu_n} - \epsilon - \tilde{\mathcal{O}}\left(\frac{\bar{\varphi}_{2,D}\sqrt{D\log(|\mathcal{X}|/\delta)}}{\sqrt{n}} + \frac{\bar{\varphi}_{2,D}^2\log(|\mathcal{X}|/\delta)}{n}\right). \end{aligned}$$

1038 Replacing this relation in equation 20 we get the thesis. □

1039

1040 D.4 Gradient method

1041 E Proofs of section 5

1042 **Theorem 14.** *Let $\mu(\cdot) = \mathcal{U}([-1, 1])$. There is a constant C independent of d such that, for $D = 2d$*
 1043 *and $\varphi_d(x) = [1, \dots, x^{d-1}]$, $\varphi_D(x) = [1, \dots, x^{2d-1}]$, we have $\Lambda_{\mu}^{Oracle} \leq C$.*

1044 *Proof.* See Theorem 3.1 by [Themistoclakis and Van Barel \[2017\]](#) □

Algorithm 1 Subgradient Method

Require: Feature map φ_D, d , Number I of iterations

Ensure: Sequence $\alpha^* \in \mathcal{A}_d^D$

- 1: Compute $\widehat{\varphi}_D$ from φ_D via Gram-Schmidt orthogonalization
- 2: Define loss as in equation (??):

$$J(\alpha) = \|M(\alpha)\|_\infty$$

- 3: Initialize $\alpha^{(0)} \leftarrow [\text{ones}(d), \text{zeros}(D-d)]^\top$
 - 4: **for** $\ell = 1$ to I **do**
 - 5: Compute step size $\gamma_\ell = \frac{\sqrt{D-d}}{\widehat{\varphi}_{2,d}\sqrt{\ell+1}}$
 - 6: Compute a subgradient $g_\ell \in \partial J(\alpha^{(\ell-1)})$
 - 7: Update: $\alpha^{(\ell)} = \alpha^{(\ell-1)} - \gamma_\ell g_\ell$
 - 8: **if** $\alpha^{(\ell)} \notin \mathcal{A}_d^D$ **then**
 - 9: Project: $\mathbf{h}^{(\ell)} = \Pi_{\mathcal{H}_d^D} \alpha^{(\ell)}$
 - 10: **end if**
 - 11: **end for**
 - 12: **return** $\alpha^* = \alpha^{(I)}$
-

Proposition 25. Fix $\gamma > 0$. There is a function $f : [-1, 1] \rightarrow \mathbb{R}$ such that, $\mathcal{E}_\infty(f) \xrightarrow{d} 0$ and under assumptions 1 and 2 for $\mu = \mathcal{U}([-1, 1])$, with probability one,

$$\lim_{d \rightarrow \infty} \lim_{n \rightarrow \infty} \|f(\cdot) - \varphi_d(\cdot)^\top \widehat{\theta}_{n,BWR}\|_\infty = 0 \quad \lim_{n \rightarrow \infty} \|f(\cdot) - \varphi_d(\cdot)^\top \widehat{\theta}_{n,OLS}\|_\infty \gtrsim d^{1-\gamma}.$$

1045 Most of the proof of this proposition is about in building the function, that we are calling $f(\cdot)$.

1046 The construction of the function in this proof is going to be quite involved. The function is going to
1047 be a sum over n of terms of the form $\widetilde{f}_n(\cdot)$. The following notation will be used

- 1048 1. Let d_n dimension of the basis function used at step n
- 1049 2. Let $a_n = d_n^{-\gamma}$, for a parameter $\gamma > 0$ to be defined
- 1050 3. Let h_n width of the mollifier
- 1051 4. Let $M_n(\cdot) = M(\cdot/h_n)$, where $M(\cdot)$ is the standard mollifier, that is, a nonnegative function
1052 $M(\cdot) \in C^\infty((-1, 1))$ with integral one and compact support.
5. $f_n(\cdot) := \text{sgn}(\overline{\varphi}_{d_n}(\cdot)^\top \overline{\varphi}_{d_n}(x_n))$, where x_n is such that

$$\|\overline{\varphi}_{d_n}(\cdot)^\top \overline{\varphi}_{d_n}(x_n)\|_{L^1} \geq \sup_{x \in (-1, 1)} \|\overline{\varphi}_{d_n}(\cdot)^\top \overline{\varphi}_{d_n}(x)\|_{L^1} - 1.$$

- 1053 6. $\widetilde{f}_n := f_n * M_n$

1054 We are able to prove the following lemmas:

Lemma 9. For every n ,

$$\|f_n - \widetilde{f}_n\|_{L^2} = \|f_n - f_n * M_n\|_{L^2} \leq 4\sqrt{h_n}d_n$$

1055 *Proof.* In order to perform this proof, we need one result from mathematical analysis. In fact, call
1056 *bounded variation* a function $\mathcal{X} = (-1, 1) \rightarrow \mathbb{R}$ such that the following norm is bounded

$$\|f\|_{BV} := \sup_{\{x_n\}_n \subset \mathcal{X}} \sum_n |f(x_{n+1}) - f(x_n)|.$$

1057 A well-known characterization of this space Ambrosio et al. [2000] ensures that the former norm is
1058 equivalent to

$$\|f_n\|_{BV} \propto \|f\|_{L^1} + \|f'\|_{\mathcal{M}} \quad \|f'\|_{\mathcal{M}} := \sup_{g \in \mathcal{C}^0(\mathcal{X}), \|g\|_{\infty}=1} \int_{\mathcal{X}} g(x) f'(x) dx.$$

1059 Now, we can proceed to the proof. First, note that by definition f_n is in the $BV((-1, 1))$ class with
 1060 $\|f_n\|_{BV} = \mathcal{O}(d_n)$. Indeed, $f_n(\cdot)$ takes only values in $\{-1, +1\}$, and can only jump between the two
 1061 values when $\overline{\varphi}_{d_n}(\cdot)^\top \overline{\varphi}_{d_n}(x_n) = 0$, which happens at most d_n times, as the previous is a polynomial
 1062 of degree d_n . At this point, by the properties of convolution,

$$\begin{aligned} f_n(y) - f_n * M_n(y) &= f_n * (M_n(y) - \delta(y)) \\ &= f'_n * \left(\int_{-1}^y M_n(t) - \delta(t) dt \right), \end{aligned}$$

1063 Where we have moved the derivative in the first term. At this point, the properties of convolution
 1064 allow us to say that for any pair of functions g_1, g_2 , $\|g_1 * g_2\|_{L^2} \leq \|g_1\|_{\mathcal{M}} \|g_2\|_{L^2}$. Therefore, we
 1065 have

$$\begin{aligned} \|f_n(\cdot) - f_n * M_n(\cdot)\|_{L^2} &\leq \|f'_n(\cdot)\|_{\mathcal{M}} \left\| \int_{-1}^y M_n(t) - \delta(t) dt \right\|_{L^2} \\ &\leq \underbrace{\|f_n(\cdot)\|_{BV}}_{\leq d_n} \left\| \int_{-1}^y M_n(t) - \delta(t) dt \right\|_{L^2}. \end{aligned}$$

1066 At this point, note that by definition $M_n(t) \geq 0$, its integral is one and its support is contained in
 1067 $(-h_n, h_n)$. Therefore,

$$\left| \int_{-1}^y M_n(t) - \delta(t) dt \right| \leq \begin{cases} 0 & y \geq h_n \\ 2 & -h_n < y < h_n \\ 0 & y \leq -h_n \end{cases},$$

1068 so that its L^2 norm is bounded by $4\sqrt{h_n}$. This completes the proof.

1069

□

Lemma 10. For every $m \leq n$, and $s > 0$

$$\|\tilde{f}_m - \Pi_{d_{n+1}, \infty}^\infty \tilde{f}_m\|_{\infty} \leq \mathcal{O}(d_{n+1}^{-s} h_m^{-s}).$$

1070 *Proof.* First, let us examine the smoothness of \tilde{f}_m . Indeed, we have, for any $s > 0$

$$\begin{aligned} \|\tilde{f}_m\|_{\mathcal{C}^s} &= \|f_m * M_m\|_{\mathcal{C}^s} \\ &\leq \|f_m\|_{\infty} \|M_m\|_{\mathcal{C}^s} \\ &= \|M_m\|_{\mathcal{C}^s} = \mathcal{O}(h_m^{-s}). \end{aligned}$$

1071 Therefore, by Jackson's theorem, we have for any s ,

$$\|\tilde{f}_m - \Pi_{d_{n+1}, \infty}^\infty \tilde{f}_m\| \leq \mathcal{O}(d_{n+1}^{-s} \|\tilde{f}_m\|_{\mathcal{C}^s}) = \mathcal{O}(d_{n+1}^{-s} h_m^{-s}).$$

1072

□

1073 **Theorem 26.** For any $\gamma < 1/4$ there is f^* such that

1074 • $\lim_d \|f^* - \Pi_{d, \infty} f^*\|_{\infty} = 0$

1075 $\bullet \limsup_d \frac{\|f^* - \Pi_{d,\mu} f^*\|_\infty}{d^{1-\gamma}} > 0.$

Proof. Let

$$f^*(\cdot) = \sum_{n=1}^{\infty} a_n \tilde{f}_n(\cdot).$$

1076 **First part**

1077 Fix $\varepsilon > 0$. As a_n goes to zero faster than exponentially and $\|\tilde{f}_n(\cdot)\|_\infty \leq 1$, we can find n_0 such that

$$\left\| f^*(\cdot) - \sum_{n=1}^{n_0} a_n \tilde{f}_n(\cdot) \right\|_\infty \leq \varepsilon/2.$$

1078 Now, $\sum_{n=1}^{n_0} a_n \tilde{f}_n(\cdot)$ is a finite sum of $C^\infty([-1, 1])$ functions, so it is uniformly continuous, in
 1079 particular. Therefore, by Stone-Weierstrass theorem, for sufficiently large d ,

$$\left\| \sum_{n=1}^{n_0} a_n \tilde{f}_n(\cdot) - \Pi_{d,\infty} \sum_{n=1}^{n_0} a_n \tilde{f}_n(\cdot) \right\|_\infty \leq \varepsilon/2.$$

1080 Putting the two results together, we have proved that, for sufficiently large d ,

$$\begin{aligned} \|f^* - \Pi_{d,\infty} f^*\|_\infty &\leq \left\| f^*(\cdot) - \Pi_{d,\infty} \sum_{n=1}^{n_0} a_n \tilde{f}_n(\cdot) \right\|_\infty \\ &\leq \varepsilon/2 + \left\| \sum_{n=1}^{n_0} a_n \tilde{f}_n(\cdot) - \Pi_{d,\infty} \sum_{n=1}^{n_0} a_n \tilde{f}_n(\cdot) \right\|_\infty \\ &\leq \varepsilon. \end{aligned}$$

1081 **Second part** Let us fix $n = \ell$ and consider

$$\begin{aligned} \|\Pi_{d_\ell,\mu} f^*\|_\infty &= \left\| \Pi_{d_\ell,\mu} \sum_{n=1}^{\infty} a_n \tilde{f}_n(\cdot) \right\|_\infty \\ &= \left\| \Pi_{d_\ell,\mu} \sum_{n=1}^{\ell-1} a_n \tilde{f}_n(\cdot) + \Pi_{d_\ell,\mu} a_\ell \tilde{f}_\ell(\cdot) + \Pi_{d_\ell,\mu} \sum_{n=\ell+1}^{\infty} a_n \tilde{f}_n(\cdot) \right\|_\infty \\ &\geq \underbrace{\left\| \Pi_{d_\ell,\mu} a_\ell \tilde{f}_\ell(\cdot) \right\|_\infty}_A - \underbrace{\left\| \Pi_{d_\ell,\mu} \sum_{n=1}^{\ell-1} a_n \tilde{f}_n(\cdot) \right\|_\infty}_B - \underbrace{\left\| \Pi_{d_\ell,\mu} \sum_{n=\ell+1}^{\infty} a_n \tilde{f}_n(\cdot) \right\|_\infty}_C. \end{aligned}$$

1082 We are going to analyze the three terms separately.

1083 (A) We start bounding the first term from below,

$$\begin{aligned} A &= a_\ell \left\| \Pi_{d_\ell,\mu} \tilde{f}_\ell(\cdot) \right\|_\infty \\ &\geq a_\ell \left\| \Pi_{d_\ell,\mu} f_\ell(\cdot) \right\|_\infty - a_\ell \left\| \Pi_{d_\ell,\mu} (\tilde{f}_\ell(\cdot) - f_\ell(\cdot)) \right\|_\infty \\ &\geq \alpha_\ell \Lambda_{d_\ell,\mu} - \alpha_\ell \bar{\varphi}_{2,d_\ell} \|\mathbf{c}_{d_\ell}(\tilde{f}_\ell(\cdot) - f_\ell(\cdot))\|_2 \\ &= \alpha_\ell \Lambda_{d_\ell,\mu} - \alpha_\ell \bar{\varphi}_{2,d_\ell} \|\tilde{f}_\ell(\cdot) - f_\ell(\cdot)\|_{L^2} \\ &\geq \alpha_\ell \Lambda_{d_\ell,\mu} - 4\alpha_\ell \bar{\varphi}_{2,d_\ell} d_\ell \sqrt{h_\ell}. \end{aligned}$$

1084 Here, the second inequality comes from Cauchy-Schwartz, the sequent equality from
 1085 Parseval's theorem and the last comes from lemma 9. Note that, for the polynomial basis,
 1086 $\bar{\varphi}_{2,d_\ell} \approx \Lambda_{d_\ell,\mu} \approx d_\ell$, so we get

$$A \geq \Omega \left(\alpha_\ell d_\ell (1 - d_\ell \sqrt{h_\ell}) \right)$$

1087 (B) This term is

$$\begin{aligned} B &= \left\| \Pi_{d_\ell,\mu} \sum_{n=1}^{\ell-1} a_n \tilde{f}_n(\cdot) \right\|_\infty \\ &\leq \sum_{n=1}^{\ell-1} a_n \left\| \Pi_{d_\ell,\mu} \tilde{f}_n(\cdot) \right\|_\infty \\ &\leq \sum_{n=1}^{\ell-1} a_n \left\| \Pi_{d_\ell,\mu} (\tilde{f}_n(\cdot) - \Pi_{d_\ell,\infty} \tilde{f}_n(\cdot)) \right\|_\infty + a_n \left\| \Pi_{d_\ell,\infty} \tilde{f}_n(\cdot) \right\|_\infty. \end{aligned}$$

1088 The last passage holds as $\Pi_{d_\ell,\mu} \Pi_{d_\ell,\infty} \tilde{f}_n(\cdot) = \Pi_{d_\ell,\infty} \tilde{f}_n(\cdot)$. Now, we can apply lemma 10,
 1089 as $n < \ell$, which ensures

$$\begin{aligned} B &\leq \sum_{n=1}^{\ell-1} a_n \left\| \Pi_{d_\ell,\mu} (\tilde{f}_n(\cdot) - \Pi_{d_\ell,\infty} \tilde{f}_n(\cdot)) \right\|_\infty + a_n \left\| \Pi_{d_\ell,\infty} \tilde{f}_n(\cdot) \right\|_\infty \\ &\leq \sum_{n=1}^{\ell-1} a_n \left\| \Pi_{d_\ell,\mu} (\tilde{f}_n(\cdot) - \Pi_{d_\ell,\infty} \tilde{f}_n(\cdot)) \right\|_\infty + a_n \left\| \Pi_{d_\ell,\infty} \tilde{f}_n(\cdot) \right\|_\infty \\ &\leq \sum_{n=1}^{\ell-1} a_n \Lambda_\ell \left\| \tilde{f}_n(\cdot) - \Pi_{d_\ell,\infty} \tilde{f}_n(\cdot) \right\|_\infty + a_n \left\| \Pi_{d_\ell,\infty} \tilde{f}_n(\cdot) \right\|_\infty \\ &\leq \sum_{n=1}^{\ell-1} a_n d_\ell^{-s+1} h_n^{-s} + a_n \left\| \Pi_{d_\ell,\infty} \tilde{f}_n(\cdot) \right\|_\infty. \end{aligned}$$

1090 (C) The last term can be simply bounded due to the fact that $\|\tilde{f}_n\|_\infty \leq 1$:

$$C \leq d_\ell \sum_{n=\ell+1}^{\infty} a_n.$$

1091 Now, fix any $\gamma < 1/4$ and take

$$s = 2; \quad d_n = \exp(1/\gamma^n); \quad h_n = \exp(-1/(2\gamma^{n+1})); \quad a_n = \exp(-1/\gamma^{n-1}).$$

1092 We get

$$\begin{aligned} A &\geq \Omega \left(\alpha_\ell d_\ell (1 - d_\ell \sqrt{h_\ell}) \right) \\ &\geq \Omega \left(\exp((1 - \gamma)/\gamma^\ell) (1 - \exp(1/\gamma^n - 1/(4\gamma^{n+1}))) \right) \\ &\geq \Omega \left(\exp((1 - \gamma)/\gamma^\ell) (1 - \underbrace{\exp(1/\gamma^n (1 - 1/(4\gamma)))}_{\leq 0}) \right) \\ &\geq \Omega \left(\exp((1 - \gamma)/\gamma^\ell) \right) = \Omega(d_\ell^{1-\gamma}). \end{aligned}$$

1093 For term B , we have

$$\begin{aligned}
B &\leq \mathcal{O} \left(\sum_{n=1}^{\ell-1} a_n d_\ell^{-s+1} h_n^{-s} + a_n \left\| \Pi_{d_\ell, \infty} \tilde{f}_n(\cdot) \right\|_\infty \right) \\
&\leq \mathcal{O} \left(\sum_{n=1}^{\ell-1} a_n d_\ell^{-s+1} h_n^{-s} + a_n \right) \\
&\leq \mathcal{O} \left(\sum_{n=1}^{\ell-1} a_n \exp((-s+1)/\gamma^\ell) \exp(s/(2\gamma^{n+1})) + a_n \right) \\
&\leq \mathcal{O} \left(\sum_{n=1}^{\ell-1} a_n \exp((-s+1)/\gamma^\ell) \exp(s/(2\gamma^\ell)) + a_n \right) \\
&\leq \mathcal{O} \left(\sum_{n=1}^{\ell-1} a_n \underbrace{\exp((-s/2+1)/\gamma^\ell)}_{\leq 1} + a_n \right).
\end{aligned}$$

1094 Last term:

$$C \leq \mathcal{O} \left(d_\ell \sum_{n=\ell+1}^{\infty} a_n \right) \leq \mathcal{O} \left(\exp(1/\gamma^\ell) \sum_{n=\ell+1}^{\infty} \exp(-1/\gamma^{n-1}) \right) = \mathcal{O} \left(\sum_{m=0}^{\infty} \exp(-1/\gamma^m) \right).$$

1095 Again, this term satisfies $C = \mathcal{O}(1)$, as the term $\exp(-1/\gamma^m)$ in the last sum decays faster than
1096 2^{-m} . \square

1097 where the last passage holds as $s = 2$. Therefore we get $B \leq \mathcal{O}(\sum_{n=1}^{\ell-1} a_n) = \mathcal{O}(1)$, since a_n decays
1098 faster than 2^{-n} which already generates a convergent sequence.

1099 All together, these passages prove

$$\|\Pi_{d_n, \mu} f^*\|_\infty \geq \Omega(d_n^{1-\gamma}).$$

1100 Therefore, taking this d_n sequence entails $\limsup_{d \rightarrow \infty} \frac{\|f^* - \Pi_{d, \mu} f^*\|_\infty}{d^{1-\gamma}} > 0$.

1101 *Proof. (of proposition 15).* Let $f = f^*$ defined before, for the specific value of $\gamma > 0$. Thanks to
1102 part one of theorem 26⁴, assumption $\mathcal{E}_\infty(f) \xrightarrow{d} 0$ is satisfied:

$$\mathcal{E}_\infty(f) = \|f^* - \Pi_{d, \infty} f^*\|_\infty \xrightarrow{d} 0.$$

1103 Then, we prove the two theses point by point. Point one: for fixed d , theorem 13 gives

$$\|f(\cdot) - \varphi_d(\cdot)^\top \hat{\theta}_{n, \text{BWR}}\|_\infty \leq (1 + \Lambda_\mu^{\text{Oracle}}) \mathcal{E}_\infty(f) + \tilde{\mathcal{O}} \left(\frac{\bar{\varphi}_{2,D} \sqrt{D \log(|\mathcal{X}|/\delta)}}{\sqrt{n}} + \frac{\bar{\varphi}_{2,D}^2 \log(|\mathcal{X}|/\delta)}{n} \right).$$

1104 As \mathcal{X} is $[-1, 1]$ and the feature map is Lipschitz continuous, we can get rid of the $|\mathcal{X}|$ by a covering
1105 argument. As $n \rightarrow \infty$, the former gives

$$\lim_n \|f(\cdot) - \varphi_d(\cdot)^\top \hat{\theta}_{n, \text{BWR}}\|_\infty \leq (1 + \Lambda_\mu^{\text{Oracle}}) \mathcal{E}_\infty(f).$$

⁴formally, the result holds for $\gamma > 1/4$ but, for what we are trying to prove, the validity of the statement for γ implies its validity for every $\gamma' > \gamma$, therefore we can proceed w.l.o.g.

1106 For $\mu = \mathcal{U}([-1, 1])$, theorem 14 ensured that $\Lambda_\mu^{\text{Oracle}} < C$, a universal constant independent on d .
 1107 Therefore,

$$\lim_d \lim_n \|f(\cdot) - \varphi_d(\cdot)^\top \hat{\theta}_{n, \text{BWR}}\|_\infty \leq \lim_n (1 + C) \mathcal{E}_\infty(f) = 0.$$

1108 Let us pass to the second thesis:

$$\lim_{n \rightarrow \infty} \|f(\cdot) - \varphi_d(\cdot)^\top \hat{\theta}_{n, \text{OLS}}\|_\infty \gtrsim d^{1-\gamma}.$$

1109 This follows from the fact that, for $n \rightarrow \infty$, $\varphi_d(\cdot)^\top \hat{\theta}_{n, \text{OLS}} \rightarrow \Pi_{d, \mu} f(\cdot)$ and that theorem 26 ensures

1110 $\limsup_d \frac{\|f^* - \Pi_{d, \mu} f^*\|_\infty}{d^{1-\gamma}} > 0.$

1111

□