

Online Learning for Prediction via Covariance Fitting: Computation, Performance and Robustness

Muhammad Osama

*Department of Information Technology,
Uppsala University*

muhammad.osama@it.uu.se

Dave Zachariah

*Department of Information Technology,
Uppsala University*

dave.zachariah@it.uu.se

Petre Stoica

*Department of Information Technology,
Uppsala University*

ps@it.uu.se

Thomas B. Schön

*Department of Information Technology,
Uppsala University*

thomas.schon@it.uu.se

Reviewed on OpenReview: <https://openreview.net/forum?id=nAr9PhyEbQ>

Abstract

We consider the problem of online prediction using linear smoothers that are functions of a nominal covariance model with unknown parameters. The model parameters are often learned using cross-validation or maximum-likelihood techniques. But when training data arrives in a streaming fashion, the implementation of such techniques can only be done in an approximate manner. Even if this limitation could be overcome, there appears to be no clear-cut results on the statistical properties of the resulting predictor.

Here we consider a covariance-fitting method to learn the model parameters, which was initially developed for spectral estimation. We first show that the use of this approach results in a computationally efficient online learning method in which the resulting predictor can be updated sequentially. We then prove that, with high probability, its out-of-sample error approaches the optimal level at a root- n rate, where n is the number of data samples. This is so even if the nominal covariance model is misspecified. Moreover, we show that the resulting predictor enjoys two robustness properties. First, it corresponds to a predictor that minimizes the out-of-sample error with respect to the least favourable distribution within a given Wasserstein distance from the empirical distribution. Second, it is robust against errors in the covariate training data. We illustrate the performance of the proposed method in a numerical experiment.

1 Introduction

We consider scenarios in which we observe a *stream* of randomly distributed data

$$\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\} \quad n = 1, 2, 3, \dots$$

Given covariate \mathbf{x}_{n+1} in a space \mathcal{X} , our goal is to predict the outcome y_{n+1} in \mathbb{R} . A large class of predictors (also known as linear smoothers) can be described as a weighted combination of observed outcomes:

$$\hat{y}(\mathbf{x}; \boldsymbol{\lambda}) = \sum_{i=1}^n w_i(\mathbf{x}; \boldsymbol{\lambda}) y_i, \quad (1)$$

where \mathbf{x} denotes any test point and the weights $\{w_i(\mathbf{x}; \boldsymbol{\lambda})\}$ are parameterized by $\boldsymbol{\lambda}$. The sensitivity of such a predictor function to noise in the training data is often characterized by how close the in-sample prediction $\hat{y}(\mathbf{x}_i; \boldsymbol{\lambda})$ is to y_i and quantified by the sum of in-sample weights,

$$0 < df_n \triangleq \sum_{i=1}^n w_i(\mathbf{x}_i; \boldsymbol{\lambda}), \quad (2)$$

also known as the ‘effective’ degrees of freedom (Ruppert et al., 2003; Wasserman, 2006; Hastie et al., 2009). These degrees of freedom are often tuned to avoid overfitting to the irreducible noise in the training data with the aim of achieving good out-of-sample performance. This includes learning the parameters $\boldsymbol{\lambda}$ from \mathcal{D}_n via distribution-free cross-validation or distribution-based maximum likelihood methods, which however can typically be implemented only approximately in the online scenario.

In this paper, we consider an alternative method using a covariance-based criterion first proposed in the context of spectral estimation (Stoica et al., 2010a;b). We show that this method

- enables sequential computation of a predictor with learned parameters,
- approaches an optimal out-of-sample performance at a root- n rate,
- enjoys two types of robustness properties.

For illustration of the online learning method, we include a numerical experiment.

Notation: $\|\mathbf{Z}\|_{\mathbf{W}} = \sqrt{\text{tr}\{\mathbf{Z}^{\top} \mathbf{W} \mathbf{Z}\}}$ is a weighted Frobenius norm of matrix \mathbf{Z} using a positive definite weight matrix \mathbf{W} . The element-wise Hadamard product between \mathbf{z} and \mathbf{z}' is denoted $\mathbf{z} \odot \mathbf{z}'$.

2 Problem formulation

The linear smoother predictor (1) can be written compactly as

$$\hat{y}(\mathbf{x}; \boldsymbol{\lambda}) = \mathbf{w}^{\top}(\mathbf{x}; \boldsymbol{\lambda}) \mathbf{y}, \quad \text{where } \mathbf{w}^{\top}(\mathbf{x}; \boldsymbol{\lambda}) = [w(\mathbf{x}_1; \boldsymbol{\lambda}), \dots, w(\mathbf{x}_n; \boldsymbol{\lambda})] \text{ and } \mathbf{y} = [y_1, \dots, y_n]^{\top}. \quad (3)$$

We will investigate a class of model-based weights for which (3) can be computed sequentially from the stream \mathcal{D}_n . Specifically, suppose y is modeled as a zero-mean stochastic process with a nominal covariance function parameterized as

$$\text{Cov}[y, y' | \mathbf{x}, \mathbf{x}'; \boldsymbol{\lambda}] = \lambda_0 \delta(\mathbf{x}, \mathbf{x}') + \sum_{k=1}^d \lambda_k \phi_k(\mathbf{x}) \phi_k(\mathbf{x}'), \quad (4)$$

where \mathbf{x} and \mathbf{x}' are two arbitrary covariates, $\delta(\mathbf{x}, \mathbf{x}')$ is the Kronecker delta function and $\{\phi_k(\mathbf{x})\}$ are real-valued features of \mathbf{x} . Here $\boldsymbol{\lambda}$ is the *unknown* set of $d+1$ nonnegative model covariance parameters. When \mathbf{x} belongs to a vector space, then periodic Fourier-type features provide a convenient choice due to their excellent covariance approximating properties (Ruppert et al., 2003; Rahimi & Recht, 2007; Stein, 2012; Hensman et al., 2017; Solin & Särkkä, 2020). Using the notation above, a set of optimal weights can be written as

$$\mathbf{w}(\mathbf{x}; \boldsymbol{\lambda}) = \mathbf{C}_{\boldsymbol{\lambda}}^{-1} \boldsymbol{\Phi} \boldsymbol{\Lambda} \phi(\mathbf{x}), \quad (5)$$

(Bishop, 2006; Rasmussen & Williams, 2006; Stein, 2012) where

$$\phi(\mathbf{x}) = \begin{bmatrix} \phi_1(\mathbf{x}) \\ \vdots \\ \phi_d(\mathbf{x}) \end{bmatrix}, \quad \boldsymbol{\Phi} = \begin{bmatrix} \phi^{\top}(\mathbf{x}_1) \\ \vdots \\ \phi^{\top}(\mathbf{x}_n) \end{bmatrix}$$

contain the features and

$$\mathbf{C}_\lambda = \Phi \Lambda \Phi^\top + \lambda_0 \mathbf{I}_n \succ \mathbf{0} \quad (6)$$

is a nominal covariance matrix where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$. The predictor function above includes a variety of penalized regression methods (see the references cited above). The degrees of freedom of (3) are controlled by λ and we have that

$$0 < df_n(\lambda) = \text{tr} \{ \Phi \Lambda \Phi^\top \mathbf{C}_\lambda^{-1} \} = \text{tr} \{ \Phi \Lambda \Phi^\top (\Phi \Lambda \Phi^\top + \lambda_0 \mathbf{I}_n)^{-1} \} \leq \min(n, d) \quad (7)$$

(Stoica & Stanasila, 1982; Ruppert et al., 2003).

Let λ_n denote the parameters that are *fitted* to the data stream \mathcal{D}_n in some way. The concern of this paper is two-fold: (i) *sequential* computation of $\hat{y}(\mathbf{x}; \lambda_n)$ and (ii) performance *guarantees* that hold even when the model class is not well-specified. That is, the nominal covariance model (4) using $\{\phi_k(\mathbf{x})\}$ may not match the unknown covariance structure of the data. Nevertheless, we aim to derive meaningful guarantees that remain valid also in this misspecified case.

Cross-validation or maximum likelihood methods are two popular approaches to fitting λ_n . However, these methods have challenges in the online learning setting. First, computing λ_n for the covariance model above is a non-convex problem that can be riddled with multiple local minima. Second, for each additional training data point $(\mathbf{x}_{n+1}, y_{n+1})$, the parameters λ_{n+1} will have to be refitted and therefore $\hat{y}(\mathbf{x}; \lambda_n)$ recomputed from scratch or approximated (see below). Third, to the best of our knowledge, finite out-of-sample prediction performance guarantees for $\hat{y}(\mathbf{x}; \lambda_n)$ neither exist for cross-validation nor maximum likelihood λ_n .

We will consider an alternative covariance-based fitting criterion for λ , used in another context, viz. spectral estimation (Stoica et al., 2010a;b). The main contribution is to derive performance guarantees of a sequentially computable predictor, thereby establishing predictive properties of the covariance-fitting methodology. Similar results do not appear to be available in the previous literature.

3 Other approaches to online model-based prediction

We identify two main lines of work that enable online or fast implementations of model-based linear smoothers, which include Gaussian process regression and kriging methods: The first line considers spectral model approximations (e.g., (Rahimi & Recht, 2007; Hensman et al., 2017; Solin & Särkkä, 2020)) which is covered by the class of covariance models in (4). These methods also enable efficient online computation of $\hat{y}(\mathbf{x}; \lambda)$, but for a fixed set of model parameters λ . The second line considers sparse variational approximations (e.g., (Titsias, 2009; Bui et al., 2017; Stanton et al., 2021)). These methods can recompute the predictor $\hat{y}(\mathbf{x}; \lambda)$ efficiently when new data arrives, but again for fixed λ .

Fitting λ_n and computing $\hat{y}(\mathbf{x}; \lambda_n)$ in a sequential manner requires recomputing past covariance quantities and is computationally prohibitive. For the covariance model above, it means recomputing (5) and the inverse of (6) at *each new sample*. Work on this problem appears to be scarce. Stanton et al. (2021) consider seeking maximum likelihood estimates λ_n by a type of gradient-based search that projects past quantities onto a lower-dimensional space. However, neither a convergence analysis towards the sought maximum likelihood estimate nor any resulting predictive properties are provided. By contrast, Cai & Yuan (2012) provide an asymptotic performance guarantee of a linear smoother with a learned regularization parameter, but their approach is restricted to the offline setting.

Our focus below is on deriving finite-sample performance guarantees for a sequentially computable predictor.

4 Learning via covariance fitting

The parameters λ can be learned by fitting the model covariance matrix \mathbf{C}_λ in (5) to the empirical covariance matrix $\mathbf{y}\mathbf{y}^\top$. Specifically, we will use the following fitting criterion, known as SPICE (Stoica et al., 2010a;b),

$$\lambda_n = \arg \min_{\lambda \geq \mathbf{0}} \|\mathbf{y}\mathbf{y}^\top - \mathbf{C}_\lambda\|_{\mathbf{C}_\lambda^{-1}}^2, \quad (8)$$

where λ_n is a function of the datastream \mathcal{D}_n . Since this criterion is convex in λ , we can be sure that a global minimizer can be determined.

4.1 Sequential computation

The covariance-fitting criterion in (8) can be connected to an alternative convex problem and this connection can be utilized to perform parameter estimation online, as shown in (Zachariah & Stoica, 2015). We first leverage this connection to compute $\hat{y}(\mathbf{x}; \lambda_n)$ in a sequential manner, and subsequently we derive a series of predictive performance guarantees of the proposed method.

Theorem 1 *The predictor function $\hat{y}(\mathbf{x}; \lambda_{n+1})$ can be updated from $\hat{y}(\mathbf{x}; \lambda_n)$ in a constant runtime $\mathcal{O}(d^2)$. The total memory requirement of the method is also on the order of $\mathcal{O}(d^2)$.*

Proof 1 *We first note that the predictor (3) has an equivalent form*

$$\hat{y}(\mathbf{x}; \lambda) = \phi^\top(\mathbf{x}) \underbrace{\Lambda \Phi^\top \mathbf{C}_\lambda^{-1} \mathbf{y}}_{\hat{\theta}(\lambda)} \quad (9)$$

For later use, we also note that (9) is invariant to a uniform rescaling of λ , since

$$\hat{\theta}(\lambda) \equiv \hat{\theta}(\alpha \lambda), \quad \forall \alpha > 0 \quad (10)$$

The covariance-fitting criterion in (8) can be expanded into

$$\|\mathbf{y}\mathbf{y}^\top - \mathbf{C}_\lambda\|_{\mathbf{C}_\lambda^{-1}}^2 = \mathbf{y}^\top \mathbf{C}_\lambda^{-1} \mathbf{y} \cdot \|\mathbf{y}\|^2 + \text{tr}\{\mathbf{C}_\lambda\} + K, \quad (11)$$

where K is a constant. Thus λ_n is also a minimizer of

$$V(\lambda) = \mathbf{y}^\top \mathbf{C}_\lambda^{-1} \mathbf{y} + \|\mathbf{y}\|^{-2} \cdot \text{tr}\{\mathbf{C}_\lambda\}$$

Next, we follow Zachariah & Stoica (2015, Appendix A) and consider the following augmented convex criterion,

$$V'(\theta, \lambda) = \frac{1}{\lambda_0} \|\mathbf{y} - \Phi \theta\|_2^2 + \theta^\top \Lambda^{-1} \theta + \text{tr}\{\mathbf{C}_\lambda\}, \quad (12)$$

and show that its minimizers produce the predictor $\hat{y}(\mathbf{x}; \lambda_n) = \phi^\top(\mathbf{x}) \hat{\theta}(\lambda_n)$.

The first argument of (12) is minimized by

$$\hat{\theta}(\lambda) = (\lambda_0^{-1} \Phi^\top \Phi + \Lambda^{-1})^{-1} \Phi^\top \lambda_0^{-1} \mathbf{y} = \Lambda \Phi^\top \mathbf{C}_\lambda^{-1} \mathbf{y},$$

where the second equality follows from using the matrix inversion lemma. Inserting the minimizer it back into (12), we have a concentrated cost function:

$$V'(\hat{\theta}(\lambda), \lambda) = \mathbf{y}^\top \mathbf{C}_\lambda^{-1} \mathbf{y} + \text{tr}\{\mathbf{C}_\lambda\}$$

Let us now consider the minimizing λ . By rescaling the parameters by $\alpha = \|\mathbf{y}\|^{-1} > 0$, we have that

$$\begin{aligned} \alpha \cdot V'(\hat{\theta}(\alpha \lambda), \alpha \lambda) &= \alpha \cdot (\mathbf{y}^\top (\alpha \mathbf{C}_\lambda)^{-1} \mathbf{y} + \text{tr}\{\alpha \mathbf{C}_\lambda\}) \\ &= V(\lambda) \end{aligned}$$

where $\alpha = \|\mathbf{y}\|^{-1} > 0$. It follows that $\alpha \lambda_n$ is a minimizer of (12), since

$$V'(\hat{\theta}(\alpha \lambda_n), \alpha \lambda_n) = \frac{1}{\alpha} V(\lambda_n) \leq \frac{1}{\alpha} V(\lambda) = V'(\hat{\theta}(\alpha \lambda), \alpha \lambda) \quad \forall \lambda \geq 0.$$

This scaled minimizer can be related to the predictor with fitted parameters: $\hat{y}(\mathbf{x}; \boldsymbol{\lambda}_n) = \boldsymbol{\phi}^\top(\mathbf{x})\hat{\boldsymbol{\theta}}(\boldsymbol{\lambda}_n) \equiv \boldsymbol{\phi}^\top(\mathbf{x})\hat{\boldsymbol{\theta}}(\alpha\boldsymbol{\lambda}_n)$ using (10).

We now change the order of the minimization of $V'(\boldsymbol{\theta}, \boldsymbol{\lambda})$ to arrive at an alternative way of computation. The minimizing parameters are

$$\lambda_k(\boldsymbol{\theta}) = \begin{cases} \frac{\|\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\theta}\|_2}{\sqrt{n}}, & k = 0 \\ \frac{|\theta_k|}{\sqrt{n}\psi_k}, & k = 1, \dots, d \end{cases} \quad (13)$$

where $\psi_k = \sqrt{\frac{1}{n} \sum_{i=1}^n \phi_k^2(\mathbf{x}_i)}$. Inserting (13) into (12) yields the following equivalent convex cost function

$$\arg \min_{\boldsymbol{\theta}} \sqrt{\frac{1}{n} \|\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\theta}\|_2^2} + \frac{1}{\sqrt{n}} \|\boldsymbol{\psi} \odot \boldsymbol{\theta}\|_1, \quad (14)$$

where $\boldsymbol{\psi} = [\psi_1 \cdots \psi_d]^\top$. Let $\boldsymbol{\theta}_n$ denote the minimizer of (14), then $\hat{y}(\mathbf{x}; \boldsymbol{\lambda}_n) \equiv \boldsymbol{\phi}^\top(\mathbf{x})\boldsymbol{\theta}_n$.

Eq. (14) is a weighted square-root LASSO problem (Belloni et al., 2011) that can be solved in a runtime on the order of $\mathcal{O}(d^2)$ using variables

$$\mathbf{A}^n = \sum_{i=1}^n \boldsymbol{\phi}(\mathbf{x}_i)\boldsymbol{\phi}^\top(\mathbf{x}_i), \quad \mathbf{b}^n = \sum_{i=1}^n \boldsymbol{\phi}^\top(\mathbf{x}_i)y_i, \quad c^n = \sum_{i=1}^n y_i^2, \quad (15)$$

of fixed dimension that are updated recursively. Thus the memory requirement is dominated by the storing of the $d \times d$ -matrix \mathbf{A}^n .

The pseudocode for the method is provided in the appendix.

4.2 Out-of-sample performance

We now turn to evaluating the out-of-sample performance of the predictor learned from the data stream \mathcal{D}_n . Specifically, we consider the mean-squared error

$$\text{MSE} = \mathbb{E} \left[(y_{n+1} - \hat{y}(\mathbf{x}_{n+1}; \boldsymbol{\lambda}_n))^2 \right], \quad (16)$$

for the *subsequent* sample $(\mathbf{x}_{n+1}, y_{n+1})$ in the stream. The expectation is conditional on \mathcal{D}_n , thus the MSE will depend on the particular realization of the stream since the learned predictor is a function of all past samples.

To provide a performance reference, we note that all predictors of the form (3) with (5) belong to the following class of predictor functions

$$\mathcal{F} \triangleq \left\{ f(\mathbf{x}) = \sum_{k=1}^d \phi_k(\mathbf{x})\theta_k : \boldsymbol{\theta} \in \mathbb{R}^d \right\}, \quad (17)$$

i.e., a linear combination of all features $\{\phi_k(\mathbf{x})\}$ in the nominal model (4). We can now benchmark $\hat{y}(\mathbf{x}; \boldsymbol{\lambda}_n)$ against the minimal achievable error among all predictors in \mathcal{F} , even when the nominal covariance model (4) is misspecified. Specifically, we provide the following finite out-of-sample performance guarantee for the learned predictor.

Theorem 2 Assume the outcome y and features are $\boldsymbol{\phi}(\mathbf{x})$ bounded, and that the features are such that $\sum_{i=1}^n |\phi_k(\mathbf{x}_i)|^2 > 0$. If the data pairs (\mathbf{x}_i, y_i) in the stream are drawn i.i.d., then the out-of-sample error of $\hat{y}(\mathbf{x}; \boldsymbol{\lambda}_n)$ is given by

$$\mathbb{E} \left[(y_{n+1} - \hat{y}(\mathbf{x}_{n+1}; \boldsymbol{\lambda}_n))^2 \right] \leq \min_{\hat{y} \in \mathcal{F}} \mathbb{E} \left[(y_{n+1} - \hat{y}(\mathbf{x}_{n+1}))^2 \right] + K \sqrt{\frac{1}{n} \ln \frac{2(d+1)^2}{\varepsilon}} + b_n \quad (18)$$

with probability of at least $1 - \varepsilon$, where K is a constant and b_n is bounded as $\mathcal{O}(n^{-3/4})$. That is, with high probability, the out-of-sample error approaches the minimum achievable error at a root- n rate. Note that the number of features d increases only the second term at a logarithmic rate.

Proof 2 For notational simplicity, let (\mathbf{x}, y) denote the random $(n + 1)$ th sample. Let \hat{y} be any predictor in \mathcal{F} and express its out-of-sample mean-square error as:

$$R(\hat{y}) \equiv \mathbb{E} \left[(y - \phi^\top(\mathbf{x})\boldsymbol{\theta})^2 \right] = \mathbb{E} \left[\begin{array}{c} \left[\boldsymbol{\theta} \right]^\top \\ \left[-1 \right] \end{array} \underbrace{\begin{bmatrix} \phi(\mathbf{x}) \\ y \end{bmatrix}}_{\mathbf{z}} \underbrace{\begin{bmatrix} \phi(\mathbf{x}) \\ y \end{bmatrix}^\top}_{\tilde{\boldsymbol{\theta}}} \begin{bmatrix} \boldsymbol{\theta} \\ -1 \end{bmatrix} \right] = \tilde{\boldsymbol{\theta}}^\top \boldsymbol{\Sigma} \tilde{\boldsymbol{\theta}}, \quad (19)$$

where $\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{z}\mathbf{z}^\top]$. Similarly, the in-sample error can be expressed as $R_n(\hat{y}) = \tilde{\boldsymbol{\theta}}^\top \widehat{\boldsymbol{\Sigma}} \tilde{\boldsymbol{\theta}}$, where $\widehat{\boldsymbol{\Sigma}} = n^{-1}(\mathbf{z}_1\mathbf{z}_1^\top + \dots + \mathbf{z}_n\mathbf{z}_n^\top)$. The gap between in- and out-of-sample errors can be bounded as:

$$\begin{aligned} |R_n(\hat{y}) - R(\hat{y})| &= |\tilde{\boldsymbol{\theta}}^\top (\widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}) \tilde{\boldsymbol{\theta}}| \\ &\leq \sum_{i=1}^{d+1} \sum_{j=1}^{d+1} |\tilde{\theta}_i| |\tilde{\theta}_j| |\widehat{\boldsymbol{\Sigma}}_{ij} - \boldsymbol{\Sigma}_{ij}| \\ &\leq (\|\boldsymbol{\theta}\|_1 + 1)^2 \cdot \underbrace{\max_{i,j} |\widehat{\boldsymbol{\Sigma}}_{ij} - \boldsymbol{\Sigma}_{ij}|}_{\tilde{\sigma}} \end{aligned} \quad (20)$$

Next, we bound $\tilde{\sigma}$ (see also Greenshtein & Ritov (2004)). Since y and $\phi(\mathbf{x})$ are bounded random variables, we have that $|z_i z_j| \leq B$ for some B and using Hoeffding's inequality

$$\Pr \left\{ |\widehat{\boldsymbol{\Sigma}}_{ij} - \boldsymbol{\Sigma}_{ij}| \geq \sigma \right\} \leq 2 \exp \left(-\frac{n\sigma^2}{2B^2} \right) \quad (21)$$

Combining this result with the union bound over all $(d + 1)^2$ variables in $\tilde{\sigma}$, we have that

$$\Pr \{ \tilde{\sigma} \geq \sigma \} \leq (d + 1)^2 \cdot 2 \exp \left(-\frac{n\sigma^2}{2B^2} \right) \triangleq \varepsilon \quad (22)$$

Consequently, we can replace $\tilde{\sigma}$ by

$$\sigma = B \sqrt{\frac{2}{n}} \sqrt{\ln \frac{2(d+1)^2}{\varepsilon}} \quad (23)$$

in (20) so that

$$|R_n(\hat{y}) - R(\hat{y})| \leq (\|\boldsymbol{\theta}\|_1 + 1)^2 \sigma,$$

holds for any predictor in \mathcal{F} with a probability of at least $1 - \varepsilon$. Thus if $\boldsymbol{\theta}$ is bounded, $\|\boldsymbol{\theta}\|_1 \leq P$ for some P , then

$$R(\hat{y}) - (P + 1)^2 \sigma \leq R_n(\hat{y}) \leq R(\hat{y}) + (P + 1)^2 \sigma \quad (24)$$

holds with a probability of at least $1 - \varepsilon$.

Let us now study two specific predictor functions in \mathcal{F} : An optimal predictor that minimizes the out-of-sample error $y^*(\mathbf{x}) = \phi^\top(\mathbf{x})\boldsymbol{\theta}^*$, where $\boldsymbol{\theta}^* = \mathbb{E}[\phi(\mathbf{x})\phi^\top(\mathbf{x})]^\dagger \mathbb{E}[\phi(\mathbf{x})y]$ is a bounded vector. The learned predictor $y^*(\mathbf{x}; \boldsymbol{\lambda}_n) = \phi^\top(\mathbf{x})\boldsymbol{\theta}_n$, where $\boldsymbol{\theta}_n$ is a minimizer of (14). This vector is also bounded because the minimizer of (14) coincides with that of

$$\arg \min_{\boldsymbol{\theta}: \|\boldsymbol{\psi} \odot \boldsymbol{\theta}\|_1 \leq \gamma} \|\mathbf{y} - \boldsymbol{\Phi}\boldsymbol{\theta}\|_2$$

for some value of $0 \leq \gamma < \infty$. Thus both $\|\boldsymbol{\theta}^*\|_1$ and $\|\boldsymbol{\theta}_n\|_1$ are bounded by some P and (24) applies to the optimal and learned predictors, denoted by y^* and \hat{y}^n for brevity.

Since $\boldsymbol{\theta}_n$ minimizes the criterion in (14), it follows that

$$\sqrt{R_n(\hat{y}^n)} + n^{-1/2} \|\boldsymbol{\psi} \odot \boldsymbol{\theta}_n\|_1 \leq \sqrt{R_n(y^*)} + n^{-1/2} \|\boldsymbol{\psi} \odot \boldsymbol{\theta}^*\|_1, \quad \forall n$$

After rearranging, we have

$$\begin{aligned} \sqrt{R_n(\widehat{y}^n)} - \sqrt{R_n(y^*)} &\leq n^{-1/2}(\|\boldsymbol{\psi} \odot \boldsymbol{\theta}^*\|_1 - \|\boldsymbol{\psi} \odot \boldsymbol{\theta}_n\|_1) \\ &\leq n^{-1/2}\|\boldsymbol{\psi} \odot \boldsymbol{\theta}^*\|_1 \\ &\leq n^{-1/2}\beta P, \end{aligned} \quad (25)$$

where $\beta = \|\boldsymbol{\psi}\|_\infty$. Multiplying both sides of the equality by the positive quantity $(\sqrt{R_n(\widehat{y}^n)} + \sqrt{R_n(y^*)})$, we have

$$\begin{aligned} R_n(\widehat{y}^n) - R_n(y^*) &\leq (\sqrt{R_n(\widehat{y}^n)} + \sqrt{R_n(y^*)})n^{-1/2}\beta P \\ &\leq (2\sqrt{R_n(y^*)} + n^{-1/2}\beta P)n^{-1/2}\beta P, \end{aligned} \quad (26)$$

where the second inequality follows from using (25). Finally, by definition $R(y^*) \leq R(\widehat{y}^n)$ and we have that

$$\begin{aligned} R(\widehat{y}^n) &\leq R_n(\widehat{y}^n) + (P+1)^2\sigma \\ &\leq R_n(y^*) + (P+1)^2\sigma + (2\sqrt{R_n(y^*)} + n^{-1/2}\beta P)n^{-1/2}\beta P \\ &\leq R(y^*) + 2(P+1)^2\sigma + (2\sqrt{R(y^*)} + (P+1)^2\sigma + n^{-1/2}\beta P)n^{-1/2}\beta P \\ &= R(y^*) + 2(P+1)^2B\sqrt{2} \cdot \sqrt{\frac{1}{n} \ln \frac{2(d+1)^2}{\varepsilon}} + \mathcal{O}(n^{-3/4}), \end{aligned} \quad (27)$$

with a probability of at least $1 - \varepsilon$, where (24) was used in the first and third inequality and (26) was used in the second inequality.

4.3 Distributional robustness

In the previous section, we showed that the out-of-sample error of $\widehat{y}(\mathbf{x}; \boldsymbol{\lambda}_n)$ approaches the minimum achievable MSE at a root- n rate. We will now see that this predictor also provides robustness against distributional uncertainty for finite n .

The feature vector $\boldsymbol{\phi}(\mathbf{x}) : \mathcal{X} \rightarrow \mathbb{R}^d$ for any predictor in \mathcal{F} leads to a distribution of the random variables $(\boldsymbol{\phi}, y)$ which we denote $p(\boldsymbol{\phi}, y)$. The out-of-sample MSE can then be written

$$\mathbb{E} \left[(y_{n+1} - \widehat{y}(\mathbf{x}_{n+1}))^2 \right] \equiv \mathbb{E}_p \left[(y_{n+1} - \boldsymbol{\phi}_{n+1}^\top \boldsymbol{\theta})^2 \right], \quad (28)$$

Using n i.i.d. samples $(\boldsymbol{\phi}_i, y_i)$, we can define a predictor in \mathcal{F} that minimizes the MSE under the *least favourable distribution* among all plausible distributions that are consistent with the data. Such a predictor is called ‘distributionally robust’, see, e.g., Duchi & Namkoong (2018). To formalize a set of plausible distributions, we first define the empirical distribution

$$p_n(\boldsymbol{\phi}, y) = \frac{1}{n} \sum_{i=1}^n \delta(\boldsymbol{\phi} - \boldsymbol{\phi}_i, y - y_i) \quad (29)$$

Then we consider a set of distributions

$$\{p : D(p_n, p) \leq \epsilon_n\}, \quad (30)$$

where $D(p_n, p)$ is some divergence measure. A distributionally robust predictor minimizes the MSE under the least-favourable distribution in the set (30), viz.

$$\max_{p : D(p_n, p) \leq \epsilon_n} \mathbb{E}_p \left[(y_{n+1} - \boldsymbol{\phi}_{n+1}^\top \boldsymbol{\theta})^2 \right] \quad (31)$$

Several different divergence measures $D(p_n, p)$ have been considered in the literature, including Kullback-Leibler divergence, chi-square divergence, and so on. One popular divergence measure is the Wasserstein distance (Blanchet et al., 2019), which is defined as

$$D(p_n, p) = \inf_{\pi} \mathbb{E}_{\pi} [c(\boldsymbol{\phi}, y, \boldsymbol{\phi}', y')], \quad (32)$$

where $c(\phi, y, \phi', y')$ is a nonnegative cost function and π is a joint distribution over (ϕ, y, ϕ', y') whose marginals equal $p_n(\phi, y)$ and $p(\phi', y')$, respectively. Thus $D(p_n, p)$ can be interpreted as measuring the expected cost of moving probability mass from one distribution to the other.

Theorem 3 *Suppose the observed features are standardized so that $\frac{1}{n} \sum_{i=1}^n \phi_k^2(\mathbf{x}_i) = 1$. Then $\hat{y}(\mathbf{x}; \lambda_n)$ corresponds to a predictor that minimizes the out-of-sample MSE under the least favourable distribution in the set $\{p : D(p_n, p) \leq n^{-2}\}$, defined by a Wasserstein distance (32) with a cost function*

$$c(\phi, y, \phi', y') = \begin{cases} \|\phi - \phi'\|_\infty & y = y', \\ \infty & \text{otherwise.} \end{cases} \quad (33)$$

Thus the predictor is robust against distributional uncertainties in the features ϕ , which may be high-dimensional. Note that the size of the distribution set shrinks with n .

Proof 3 When the features are standardized, then $\psi = \mathbf{1}$ and (14) becomes

$$\arg \min_{\theta} \sqrt{\frac{1}{n} \|\mathbf{y} - \Phi \theta\|_2^2} + \frac{1}{\sqrt{n}} \|\theta\|_1. \quad (34)$$

Using Theorem 1 in (Blanchet et al., 2019), (34) corresponds to a predictor that minimizes (31) with divergence bound $\epsilon_n = n^{-2}$.

4.4 In-sample robustness

When learning $\hat{y}(\mathbf{x}; \lambda_n)$ it is possible that the observed covariates themselves are subject to errors so that the dataset is:

$$\tilde{\mathcal{D}}_n = \{(\tilde{\mathbf{x}}_1, y_1), \dots, (\tilde{\mathbf{x}}_n, y_n)\}$$

Then the true feature vector $\phi_i = \phi(\mathbf{x}_i)$ can be viewed as a perturbed version of the observed vector $\tilde{\phi}_i = \phi(\tilde{\mathbf{x}}_i)$, where the perturbation $\delta_i = \phi_i - \tilde{\phi}_i$ is unknown. This problem (aka. errors-in-variables) leads to yet another interpretation of the predictor $\hat{y}(\mathbf{x}; \lambda_n)$.

Theorem 4 *Consider the bounded set of possible in-sample perturbations:*

$$\mathcal{S}_n = \left\{ \delta_1, \dots, \delta_n : \mathbb{E}_{p_n} [\delta_k^2] \leq n^{-1} \mathbb{E}_{p_n} [\tilde{\phi}_k^2], \forall k = 1, \dots, d \right\}$$

Then $\hat{y}(\mathbf{x}; \lambda_n)$ corresponds to a predictor that minimizes the in-sample root-MSE under the least-favourable perturbations in \mathcal{S}_n :

$$\max_{\{\delta_i\} \in \mathcal{S}_n} \sqrt{\mathbb{E}_{p_n} [(y - (\tilde{\phi} + \delta)^\top \theta)^2]}, \quad (35)$$

where $\hat{y} = (\tilde{\phi} + \delta)^\top \theta \in \mathcal{F}$.

Proof 4 The problem in (35) can be written as:

$$\max_{\{\delta_i\} \in \mathcal{S}} \frac{1}{\sqrt{n}} \|\mathbf{y} - (\tilde{\Phi} + \Delta) \theta\|_2, \quad \text{where } \Delta = \begin{bmatrix} \delta_1^\top \\ \vdots \\ \delta_n^\top \end{bmatrix} \quad (36)$$

Let $[\Delta]_k$ denote the k^{th} column of the matrix Δ . We can then upper bound the error as

$$\begin{aligned} \max_{\{\delta_i\} \in \mathcal{S}} \frac{1}{\sqrt{n}} \left\| \mathbf{y} - \tilde{\Phi} \theta - \sum_{k=1}^d [\Delta]_k \theta_k \right\|_2 &\leq \max_{\{\delta_i\} \in \mathcal{S}} \frac{1}{\sqrt{n}} \|\mathbf{y} - \tilde{\Phi} \theta\|_2 + \frac{1}{\sqrt{n}} \sum_{k=1}^d \|[\Delta]_k \theta_k\|_2, \\ &\leq \frac{1}{\sqrt{n}} \|\mathbf{y} - \tilde{\Phi} \theta\|_2 + \max_{\{\delta_i\} \in \mathcal{S}} \frac{1}{\sqrt{n}} \sum_{k=1}^d \|[\Delta]_k\|_2 |\theta_k|, \\ &\leq \frac{1}{\sqrt{n}} \|\mathbf{y} - \tilde{\Phi} \theta\|_2 + \frac{1}{\sqrt{n}} \sum_{k=1}^d \sqrt{\mathbb{E}_{p_n} [\tilde{\phi}_k^2]} |\theta_k|. \end{aligned} \quad (37)$$

where the bound is attainable when

$$[\Delta]_k = \sqrt{\mathbb{E}_{p_n}[\hat{\phi}_k^2]} \frac{\mathbf{y} - \tilde{\Phi}\boldsymbol{\theta}}{\|\mathbf{y} - \tilde{\Phi}\boldsymbol{\theta}\|_2}.$$

But the bound is of the same form as the cost function in (14). Thus solving problem (14) implies the minimization of (36). Theorem 1 in Xu et al. (2009) established a connection between a square-root maximization problem in the form of (36) and the penalized form on the right-hand side of (37).

5 Numerical Experiment

In the previous sections we have showed several computational and theoretical properties of the predictor function $\hat{y}(\mathbf{x}; \boldsymbol{\lambda}_n)$ which we shall call the SPICE-predictor. In this section we present a numerical experiment for sake of illustration.

5.1 Setup

We observe a stream of n samples generated by the following (unknown) process

$$\begin{aligned} \mathbf{x} &\sim \text{Uniform}([0, 10]^2), \\ y|\mathbf{x} &\sim \text{GP}(0, k(\mathbf{x}, \mathbf{x}') + \sigma^2\delta(\mathbf{x}, \mathbf{x}')), \end{aligned} \tag{38}$$

where

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \left(1 + \frac{\sqrt{3}}{l} \|\mathbf{x} - \mathbf{x}'\|_2 \right) \exp \left(-\frac{\sqrt{3}}{l} \|\mathbf{x} - \mathbf{x}'\|_2 \right).$$

with noise variance $\sigma = 2$ and scale $l = 7$. In other words, \mathbf{x} is a two-dimensional covariate drawn from a uniform distribution and y is drawn from a Gaussian process (GP) with zero mean and a Matérn covariance function. A realization of the process above and n training data points are shown in Figures 1a and 1e.

We consider a class \mathcal{F} with $d = 100$ periodic features $\{\phi_k(\mathbf{x})\}$ using the Laplacian eigenfunction basis (Solin & Särkkä, 2020). Note that this corresponds to a misspecified covariance model (4). We are interested in the online learning of predictors in \mathcal{F} , and use the least-squares (LS) and ridge regression methods as baseline references. Both methods can be implemented in an online fashion, but the latter requires fixing a regularization parameter. Here we simply set this parameter to 0.1 based on visual inspection.

5.2 Out-of-sample performance

For illustration, consider the predictions produced by the LS, ridge regression and SPICE methods, see Figure 1. As expected, the LS provides poor results at these sample sizes. Ridge regression with a fixed regularization parameter and SPICE with adaptively learned parameters appear to perform similarly here. To evaluate their out-of-sample errors, we compare the MSE against that of the oracle predictor based on the (oracle) Gaussian process regression (GPR) predictor in (38). Table 1 shows that the out-of-sample error of SPICE is lower than that of LS and ridge regression, and that the chosen class \mathcal{F} is capable of predicting the GP in (38) well.

Following the discussion of effective degrees of freedom df_n in Ruppert et al. (2003), we also provide a comparison between LS, SPICE and the oracle GPR predictors in Figure 2. While LS attains the maximum df_n at $n = 100$, SPICE moderates its growth rate in a data-adaptive and online manner. The degrees of freedom of the oracle predictor increases gracefully and remains below its maximum value, even when n increases beyond d .

5.3 Run-time

We report the runtimes of RIDGE (since LS is virtually identical) and SPICE, and consider a well-specified GPR predictor with covariance parameters learned using the maximum likelihood method as reference. As

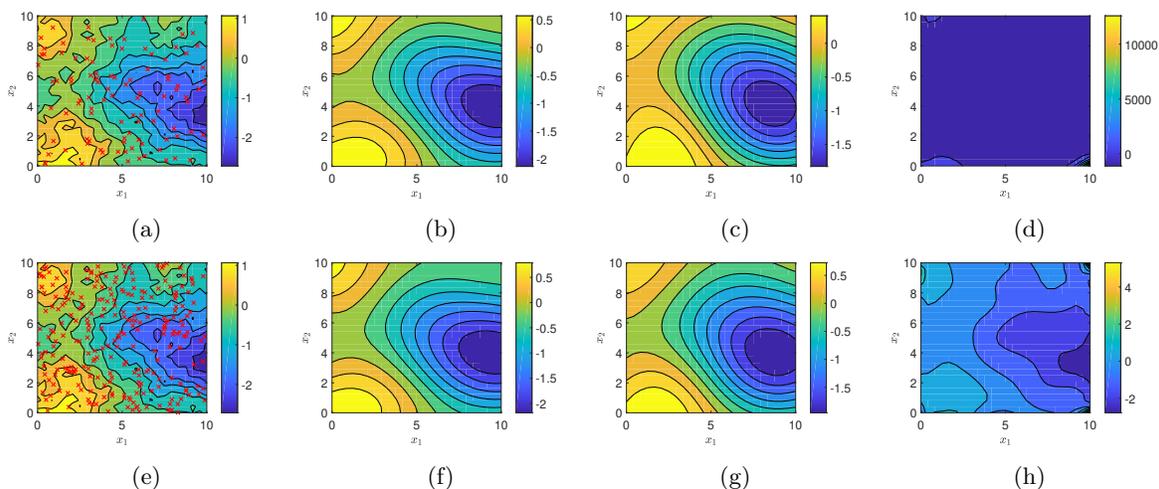


Figure 1: Contour plots. First column shows a realization of y in (38) along with sampling patterns $\{\mathbf{x}_i\}_{i=1}^n \in \mathcal{X}$ for $n = 100$ and $n = 250$ (top and bottom rows, respectively). Second, third and fourth columns show the contour plots of the SPICE-predictor, ridge regression and the LS-predictor, respectively. All three predictors belong to \mathcal{F} .

n	MSE/MSE*		
	LS	RIDGE	SPICE
50	4.38×10^4	1.71	1.11
100	21.12	1.47	1.09
250	1.47	1.19	1.06
500	1.11	1.06	1.02

Table 1: Mean-square error (MSE) for LS and SPICE methods, normalized by MSE* of an oracle predictor which is given the unknown covariance function in (38). For a given set of training data \mathcal{D}_n , we compute the averaged squared error over 250 test points. The mean of this error is the MSE and was approximated using 100 different realizations of \mathcal{D}_n .

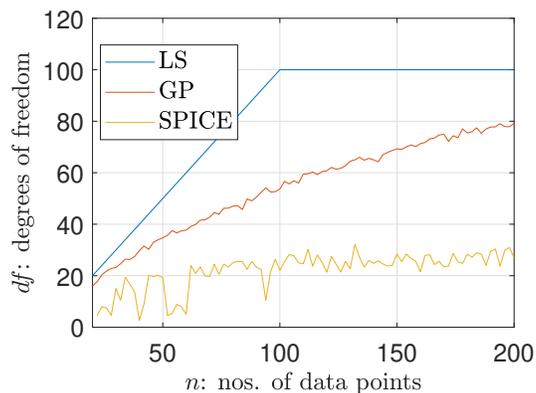


Figure 2: Plot of degrees of freedom df_n against number of data points n for LS, SPICE and oracle GP predictors.

can be seen in Table 2, the computational complexity of GPR is considerably higher than that of the online alternatives. RIDGE has a lower run-time than SPICE but both perform as $\mathcal{O}(n)$. A visualization of the trends for $0 < n \leq 500$ is given in Figure 3.

	Run times (ms)		
n	GPR	RIDGE	SPICE
50	59	0.15	6.0
100	113.8	0.16	11.6
250	231.3	0.20	29.2
500	1507.2	0.22	58.6

Table 2: Run times for oracle GPR, RIDGE and SPICE. Note that the run times are computed for a single realization.

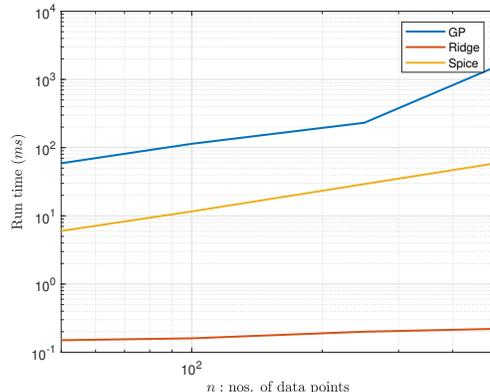


Figure 3: Plot of run time of a single run for learning parameters and prediction against versus of data points n for RIDGE, SPICE and oracle GP predictors.

5.4 Robustness

In Theorems 3 and 4 established two different types of robustness results with respect to the features $\phi(\mathbf{x})$. We therefore study the predictive performance when the training and test distributions over (ϕ, y) diverge.

Specifically, the covariates in the training data are drawn as $\tilde{\mathbf{x}}|\mathbf{x} \sim \mathcal{N}(\mathbf{x}, \sigma_x^2 \mathbf{I})$, where \mathbf{x} follows (38). This results in a distribution shift over the features (aka. errors-in-variables). Figure 4 evaluates the test MSE when data drawn from (38). We see that the out-of-sample error for RIDGE and SPICE increases consistently with σ_x , while LS produces very poor and inconsistent results. SPICE is notably more robust against this distributional shift and these results corroborated the derived robustness properties.

6 Conclusion

We considered the problem of learning model-based linear smoothers online. If the model parameters were fixed, the resulting predictor – which includes Gaussian process regression and kriging methods – can readily be computed sequentially. Since the model parameters unknown, however, they must be learned from data, using, e.g., maximum likelihood or cross-validation methods. But implementing them when the data arrives as a stream requires either recomputing the predictor which is computationally prohibitive or resorting to approximations. In either case, these approaches do not offer clear-cut results on the statistical properties of the resulting predictor.

We applied a covariance-fitting method to learn the model parameters, which was initially developed for spectral estimation. We first used its computational properties to show that the resulting predictor can be computed sequentially. We then derived finite out-of-sample performance guarantees of the resulting predictor and showed that its error approaches the minimum achievable level at root- n rate. Finally, we established connections to the distributional robustness literature by showing that the predictor is robust against distri-

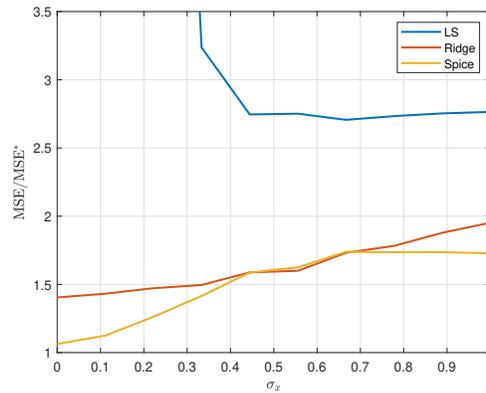


Figure 4: Out-of-sample error when the distribution of training data diverges from that of the test data. MSE (normalized by minimal MSE*) versus perturbation noise level σ_x on training data ($n = 100$). Evaluation based on 500 Monte Carlo runs.

butional uncertainties and errors in the covariate training data. The performance, computational complexity and robustness of the proposed method were illustrated in a numerical experiment.

Acknowledgments We thank the anonymous reviewers from TMLR for their constructive feedback. This research was partially supported by the Swedish Research Council under contracts 2018-05040 and 2021-05022, Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by Knut and Alice Wallenberg Foundation, and by Kjell och Märta Beijer Foundation.

References

- Alexandre Belloni, Victor Chernozhukov, and Lie Wang. Square-root lasso: pivotal recovery of sparse signals via conic programming. *Biometrika*, 98(4):791–806, 2011.
- Christopher Bishop. *Pattern recognition and machine learning*. Springer, 2006.
- Jose Blanchet, Yang Kang, and Karthyek Murthy. Robust wasserstein profile inference and applications to machine learning. *Journal of Applied Probability*, 56(3):830–857, 2019.
- Thang D Bui, Cuong Nguyen, and Richard E Turner. Streaming sparse gaussian process approximations. *Advances in Neural Information Processing Systems*, 30, 2017.
- T Tony Cai and Ming Yuan. Minimax and adaptive prediction for functional linear regression. *Journal of the American Statistical Association*, 107(499):1201–1216, 2012.
- John Duchi and Hongseok Namkoong. Learning models with uniform performance via distributionally robust optimization. *arXiv preprint arXiv:1810.08750*, 2018.
- Eitan Greenshtein and Ya’Acov Ritov. Persistence in high-dimensional linear predictor selection and the virtue of overparametrization. *Bernoulli*, 10(6):971–988, 2004.
- Trevor Hastie, Robert Tibshirani, and Jerome H Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- James Hensman, Nicolas Durrande, and Arno Solin. Variational fourier features for gaussian processes. *J. Mach. Learn. Res.*, 18(151–1), 2017.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *NIPS*, volume 3, pp. 5. Citeseer, 2007.

Carl Edward Rasmussen and Chris Williams. *Gaussian processes for machine learning*, volume 1. MIT press Cambridge, 2006.

David Ruppert, Matt P Wand, and Raymond J Carroll. *Semiparametric regression*. Number 12. Cambridge university press, 2003.

Arno Solin and Simo Särkkä. Hilbert space methods for reduced-rank gaussian process regression. *Statistics and Computing*, 30(2):419–446, 2020.

Samuel Stanton, Wesley Maddox, Ian Delbridge, and Andrew Gordon Wilson. Kernel interpolation for scalable online gaussian processes. In *International Conference on Artificial Intelligence and Statistics*, pp. 3133–3141. PMLR, 2021.

Michael L Stein. *Interpolation of spatial data: some theory for kriging*. Springer Science & Business Media, 2012.

Petre Stoica and Octavian Stanasila. Some spectral properties of the matrix b/a. *Bul. Inst. Politehnic Bucuresti, ser. Electro.*, 44(3):3–8, 1982.

Petre Stoica, Prabhu Babu, and Jian Li. New method of sparse parameter estimation in separable models and its use for spectral analysis of irregularly sampled data. *IEEE Transactions on Signal Processing*, 59(1):35–47, 2010a.

Petre Stoica, Prabhu Babu, and Jian Li. Spice: A sparse covariance-based estimation method for array processing. *IEEE Transactions on Signal Processing*, 59(2):629–638, 2010b.

Michalis Titsias. Variational learning of inducing variables in sparse gaussian processes. In *Artificial intelligence and statistics*, pp. 567–574. PMLR, 2009.

Larry Wasserman. *All of nonparametric statistics*. Springer Science & Business Media, 2006.

Huan Xu, Constantine Caramanis, and Shie Mannor. Robust regression and lasso. In *Advances in Neural Information Processing Systems*, pp. 1801–1808, 2009.

Dave Zachariah and Petre Stoica. Online hyperparameter-free sparse estimation method. *IEEE Transactions on Signal Processing*, 63(13):3348–3359, 2015.

A Appendix: Sequential computation

Here we provide a pseudocode for computing $\hat{y}(\mathbf{x}; \boldsymbol{\lambda}_n) \equiv \boldsymbol{\phi}^\top(\mathbf{x})\boldsymbol{\theta}_n$ sequentially. This is accomplished via a cyclic minimization of the convex problem (14).

First, define the k th column as $\mathbf{c}_k = [\boldsymbol{\Phi}]_k$ and $\tilde{\mathbf{y}}_k = \mathbf{y} - \sum_{j \neq k} \mathbf{c}_j \theta_j$. Then the cost function in (14) can be equivalently expressed as

$$V(\theta_k) = (\|\tilde{\mathbf{y}}_k - \mathbf{c}_k \theta_k\|_2^2)^{1/2} + \psi_k |\theta_k| + C_k, \quad (39)$$

and minimized cyclically, one coordinate $k = 1, 2, \dots, d$ at a time. It was shown in (Zachariah & Stoica, 2015) that the minimizer of (39) is given by

$$\hat{\theta}_k = \begin{cases} s_k r_k, & \text{if } \sqrt{n-1} \gamma_k > \sqrt{\alpha_k \beta_k - \gamma_k^2} \\ 0, & \text{else} \end{cases} \quad (40)$$

where

$$\alpha_k = \|\tilde{\mathbf{y}}_k\|_2^2, \quad \beta_k = \|\mathbf{c}_k\|_2^2, \quad \gamma_k = |\mathbf{c}_k^\top \tilde{\mathbf{y}}_k| \quad (41)$$

and

$$s_k = \text{sign}(\mathbf{c}_k^\top \tilde{\mathbf{y}}_k), \quad r_k = \frac{\gamma_k}{\beta_k} - \frac{1}{\beta_k} \left(\frac{\alpha_k \beta_k - \gamma_k^2}{n-1} \right) \quad (42)$$

The key observation here is that the variables (41) are all expressible using the recursively computable quantities in (15) (along with a prior parameter iterate $\boldsymbol{\theta}'$ which is initialized at $\mathbf{0}$).

To arrive at this conclusion, define the following global variables

$$\mathbf{z} = \mathbf{y} - \Phi \boldsymbol{\theta}', \quad g_0 = \|\mathbf{z}\|_2^2, \quad \mathbf{g} = \Phi^\top \mathbf{z} \quad (43)$$

which can be expressed using (15). Then we have that $\tilde{\mathbf{y}}_k \equiv \mathbf{z} + \mathbf{c}_k \theta'_k$, so that we can express (41) in terms of the global variables:

$$\alpha_k = g_0 + A_{kk}^n \theta_k'^2 + 2g_k \theta'_k, \quad \beta_k = A_{kk}^n, \quad \gamma_k = |g_k + A_{kk}^n \theta'_k|, \quad s_k = \text{sign}(g_k + A_{kk}^n \theta'_k) \quad (44)$$

Finally we can express the cyclic minimization approach in recursive form as outlined in Algorithm 1, initializing $\boldsymbol{\theta}' = \mathbf{0}$ at $n = 1$. Additional computational considerations are developed in Sec. III of (Zachariah & Stoica, 2015). Link to implementation is provided here: <https://github.com/Muhammad-Osama>.

Algorithm 1 : SPICE-predictor

- 1: Input: (\mathbf{x}_n, y_n) and \mathbf{x} (and initial $\boldsymbol{\theta}'$)
 - 2: Update recursive variables $\mathbf{A}^n, \mathbf{b}^n$ and c^n in (15)
 - 3: Set global variables $g_0 = c^n + \boldsymbol{\theta}'^\top \mathbf{A}^n \boldsymbol{\theta}' - 2\boldsymbol{\theta}'^\top \mathbf{b}^n$ and $\mathbf{g} = \mathbf{b}^n - \mathbf{A}^n \boldsymbol{\theta}'$
 - 4: **repeat**
 - 5: $k = 1, \dots, d$
 - 6: Compute $\alpha_k, \beta_k, \gamma_k$ and s_k in (44)
 - 7: Compute $\hat{\theta}_k$ in (40)
 - 8: Update global variables $g_0 := g_0 + A_{kk}^n (\theta'_k - \hat{\theta}_k)^2 + 2(\theta'_k - \hat{\theta}_k) g_k$ and $\mathbf{g} := \mathbf{g} + [\mathbf{A}^n]_k (\theta'_k - \hat{\theta}_k)$
 - 9: Update $\theta'_k := \hat{\theta}_k$
 - 10: **until** number of iterations equals L
 - 11: Output: $\hat{y}(\mathbf{x}) = \boldsymbol{\phi}^\top(\mathbf{x}) \hat{\boldsymbol{\theta}}$
-