

Statistical bounds for Gaussian regression algorithms based on Karhunen-Loève expansions

Gianluigi Pillonetto Luca Schenato Damiano Varagnolo





Roadmap

Statistical bounds

for Gaussian regression algorithms

based on Karhunen-Loève expansions

Function estimation

$$f : \mathcal{X} \rightarrow \mathbb{R} \quad (1)$$

$$y_m = f(x_m) + \nu_m \quad m = 1, \dots, M \quad (2)$$

$$x_m \sim \mu(\mathcal{X}) \text{ i.i.d.} \quad \nu_m \sim \mathcal{N}(0, \sigma_\nu^2) \quad m = 1, \dots, M \quad (3)$$

$$\{x_m\}_{m=1}^M \quad \{\nu_m\}_{m=1}^M \quad \text{mutually independent} \quad (4)$$

Problem: estimate f starting from $\{x_m, y_m\}$

Function estimation – parametric approach

$$y_m = f(x_m; \theta) + \nu_m \quad (\text{known structure or set of alternative structures}) \quad (5)$$

Function estimation – parametric approach

$$y_m = f(x_m; \theta) + \nu_m \quad (\text{known structure or set of alternative structures}) \quad (5)$$

Least squares (classic approach):

$$\theta^* = \arg \min_{\tilde{\theta} \in \Theta} \sum_m \left(y_m - f(x_m; \tilde{\theta}) \right)^2 \quad (6)$$

Function estimation – parametric approach

$$y_m = f(x_m; \theta) + \nu_m \quad (\text{known structure or set of alternative structures}) \quad (5)$$

Least squares (classic approach):

$$\theta^* = \arg \min_{\tilde{\theta} \in \Theta} \sum_m \left(y_m - f(x_m; \tilde{\theta}) \right)^2 \quad (6)$$

Potential problems:

- non-convexity
- model order selection

Nonparametric approach *(cast the problem as a Gaussian regression)*

$$f \sim \mathcal{N}(0, K) \quad K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} \quad \text{so that} \quad \mathbb{E}[f(x)f(x')] = K(x, x') \quad (7)$$

Nonparametric approach *(cast the problem as a Gaussian regression)*

$$f \sim \mathcal{N}(0, K) \quad K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} \quad \text{so that} \quad \mathbb{E}[f(x)f(x')] = K(x, x') \quad (7)$$

Examples:

- Brownian motion: $K(x, x') = \min(x, x')$ $\mathcal{X} = [0, 1]$
- Radial basis: $K(x, x') = \exp(-\|x - x'\|^2)$ $\mathcal{X} \subseteq \mathbb{R}^m$

Nonparametric approach *(cast the problem as a Gaussian regression)*

$$f \sim \mathcal{N}(0, K) \quad K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R} \quad \text{so that} \quad \mathbb{E}[f(x)f(x')] = K(x, x') \quad (7)$$

$$x_m \sim \mu(\mathcal{X}) \text{ i.i.d.} \quad y_m = f(x_m) + \nu_m \quad \nu_m \sim \mathcal{N}(0, \sigma_\nu^2) \quad m = 1, \dots, M \quad (8)$$

$$\{x_m\}_{m=1}^M \quad \{\nu_m\}_{m=1}^M \quad f \quad \text{mutually independent} \quad (9)$$

Maximum a posteriori estimator

$$\hat{f}_{\text{MAP}}(x) = \mathbb{E} \left[f(x) \mid \{x_m, y_m\} \right] \quad (\text{also MV})$$

Maximum a posteriori estimator

$$\begin{aligned}\widehat{f}_{\text{MAP}}(x) &= \mathbb{E} \left[f(x) \mid \{x_m, y_m\} \right] && \text{(also MV)} \\ &= \sum_{m=1}^M K(x, x_m) c_m && \text{(a.k.a. regularization network)}\end{aligned} \tag{10}$$

Maximum a posteriori estimator

$$\begin{aligned}\widehat{f}_{\text{MAP}}(x) &= \mathbb{E} \left[f(x) \mid \{x_m, y_m\} \right] \quad (\text{also MV}) \\ &= \sum_{m=1}^M K(x, x_m) c_m \quad (\text{a.k.a. regularization network})\end{aligned}\tag{10}$$

$$\begin{bmatrix} c_1 \\ \vdots \\ c_M \end{bmatrix} = H_{\text{MAP}} \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}\tag{11}$$

$$H_{\text{MAP}} := \left(\begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, x_M) \\ \vdots & & \vdots \\ K(x_M, x_1) & \cdots & K(x_M, x_M) \end{bmatrix} + \sigma_\nu^2 I \right)^{-1}\tag{12}$$

Gaussian regression – practical issues associated to the MAP

$$\widehat{f}_{\text{MAP}}(x) = [K(x, x_1) \ \dots \ K(x, x_M)] H_{\text{MAP}} \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} \quad (13)$$

$$H_{\text{MAP}} := \left(\begin{bmatrix} K(x_1, x_1) & \dots & K(x_1, x_M) \\ \vdots & & \vdots \\ K(x_M, x_1) & \dots & K(x_M, x_M) \end{bmatrix} + \sigma_\nu^2 I \right)^{-1} \quad (14)$$

computational cost $O(M^3)$

How may we tackle the $O(M^3)$ computational cost issue?

$$H_{\text{MAP}} := \left(\begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, x_M) \\ \vdots & & \vdots \\ K(x_M, x_1) & \cdots & K(x_M, x_M) \end{bmatrix} + \sigma_\nu^2 I \right)^{-1} \quad (15)$$

Typical approaches: low-rank / sparsification approximations



Smola & Schölkopf (2000)

Sparse greedy matrix approximations for machine learning



Quiñero-Candela & Rasmussen (2005)

A unifying view of sparse approximate Gaussian process regression



Bach & Jordan (2005)

Predictive low-rank decompositions for kernel methods



Snelson & Ghahramani (2006)

Sparse Gaussian processes using pseudo inputs



Culis et al. (2006)

Learning low-rank kernel matrices



Zhang & Kwok (2010)

Clustered Nyström method for large scale manifold learning and dimension reduction



Ambikasaran et al. (2016)

Fast direct methods for Gaussian processes

Our focus: Karhunen-Loève expansions

$$f(x) = \sum_{e=1}^{+\infty} \alpha_e \phi_e(x)$$

Our focus: Karhunen-Loève expansions

$$f(x) = \sum_{e=1}^{+\infty} \alpha_e \phi_e(x) = \underbrace{\sum_{e=1}^E a_e \phi_e(x)}_{=: \text{interesting}} + \underbrace{\sum_{e=1}^{+\infty} b_e \phi_{E+e}(x)}_{=: \text{remainder}} \quad (16)$$

Our focus: Karhunen-Loève expansions

$$f(x) = \sum_{e=1}^{+\infty} \alpha_e \phi_e(x) = \underbrace{\sum_{e=1}^E a_e \phi_e(x)}_{=: \text{interesting}} + \underbrace{\sum_{e=1}^{+\infty} b_e \phi_{E+e}(x)}_{=: \text{remainder}} \quad (16)$$

$$\lambda_e \phi_e(x) = \int_{\mathcal{X}} K(x, x') \phi_e(x') d\mu(x') \quad \lambda_1 \geq \lambda_2 \dots > 0 \quad (17)$$

$$K(x, x') = \sum_{e=1}^{+\infty} \lambda_e \phi_e(x) \phi_e(x') \quad \int_{\mathcal{X}} \phi_i(x) \phi_j(x) d\mu(x) = \delta_{ij} \quad (18)$$



Our focus: Karhunen-Loève expansions

$$f(x) = \sum_{e=1}^{+\infty} \alpha_e \phi_e(x) = \underbrace{\sum_{e=1}^E a_e \phi_e(x)}_{=:\text{interesting}} + \underbrace{\sum_{e=1}^{+\infty} b_e \phi_{E+e}(x)}_{=:\text{remainder}} \quad (16)$$

$$\lambda_e \phi_e(x) = \int_{\mathcal{X}} K(x, x') \phi_e(x') d\mu(x') \quad \lambda_1 \geq \lambda_2 \dots > 0 \quad (17)$$

$$K(x, x') = \sum_{e=1}^{+\infty} \lambda_e \phi_e(x) \phi_e(x') \quad \int_{\mathcal{X}} \phi_i(x) \phi_j(x) d\mu(x) = \delta_{ij} \quad (18)$$

$$a_e \sim \mathcal{N}(0, \lambda_e), \quad e = 1, \dots, E \quad b_e \sim \mathcal{N}(0, \lambda_{E+e}), \quad e = 1, 2, \dots \quad (19)$$



Zhu et al. (1998)

Gaussian regression and optimal finite dimensional linear models

Our focus: Karhunen-Loève expansions

$$f(x) = \underbrace{\sum_{e=1}^E a_e \phi_e(x)}_{=: \text{interesting}} + \underbrace{\sum_{e=1}^{+\infty} b_e \phi_{E+e}(x)}_{=: \text{remainder}} \quad \lambda_e \phi_e(x) = \int_{\mathcal{X}} K(x, x') \phi_e(x') d\mu(x')$$

(20)

\implies first E ϕ_e 's = best a-priori E -dimensional approximation in a MSE sense

Our focus: Karhunen-Loève expansions

$$\mathbf{y} := [y_1, \dots, y_M]^T \quad \boldsymbol{\nu} := [\nu_1, \dots, \nu_M]^T \quad \mathbf{a} := [a_1, \dots, a_E]^T \quad \mathbf{b} := [b_1, b_2, \dots]^T \quad (21)$$

Our focus: Karhunen-Loève expansions

$$\mathbf{y} := [y_1, \dots, y_M]^T \quad \boldsymbol{\nu} := [\nu_1, \dots, \nu_M]^T \quad \mathbf{a} := [a_1, \dots, a_E]^T \quad \mathbf{b} := [b_1, b_2, \dots]^T \quad (21)$$

$$G := \begin{bmatrix} \phi_1(x_1) & \dots & \phi_E(x_1) \\ \vdots & & \vdots \\ \phi_1(x_M) & \dots & \phi_E(x_M) \end{bmatrix} \quad Z := \begin{bmatrix} \phi_{E+1}(x_1) & \phi_{E+2}(x_1) & \dots \\ \vdots & & \vdots \\ \phi_{E+1}(x_M) & \phi_{E+2}(x_M) & \dots \end{bmatrix} \quad (22)$$

Our focus: Karhunen-Loève expansions

$$\mathbf{y} := [y_1, \dots, y_M]^T \quad \boldsymbol{\nu} := [\nu_1, \dots, \nu_M]^T \quad \mathbf{a} := [a_1, \dots, a_E]^T \quad \mathbf{b} := [b_1, b_2, \dots]^T \quad (21)$$

$$G := \begin{bmatrix} \phi_1(x_1) & \dots & \phi_E(x_1) \\ \vdots & & \vdots \\ \phi_1(x_M) & \dots & \phi_E(x_M) \end{bmatrix} \quad Z := \begin{bmatrix} \phi_{E+1}(x_1) & \phi_{E+2}(x_1) & \dots \\ \vdots & & \vdots \\ \phi_{E+1}(x_M) & \phi_{E+2}(x_M) & \dots \end{bmatrix} \quad (22)$$

$$\mathbf{y} = G\mathbf{a} + Z\mathbf{b} + \boldsymbol{\nu} \quad (23)$$

Our focus: Karhunen-Loève expansions

$$\mathbf{y} := [y_1, \dots, y_M]^T \quad \boldsymbol{\nu} := [\nu_1, \dots, \nu_M]^T \quad \mathbf{a} := [a_1, \dots, a_E]^T \quad \mathbf{b} := [b_1, b_2, \dots]^T \quad (21)$$

$$G := \begin{bmatrix} \phi_1(x_1) & \dots & \phi_E(x_1) \\ \vdots & & \vdots \\ \phi_1(x_M) & \dots & \phi_E(x_M) \end{bmatrix} \quad Z := \begin{bmatrix} \phi_{E+1}(x_1) & \phi_{E+2}(x_1) & \dots \\ \vdots & & \vdots \\ \phi_{E+1}(x_M) & \phi_{E+2}(x_M) & \dots \end{bmatrix} \quad (22)$$

$$\mathbf{y} = G\mathbf{a} + Z\mathbf{b} + \boldsymbol{\nu} \quad (23)$$

$$\widehat{f}_E(x) := [\phi_1(x) \quad \dots \quad \phi_E(x)] \widehat{\mathbf{a}} \quad \widehat{\mathbf{a}} = H\mathbf{y} \quad H := \left(\frac{G^T G}{M} + \frac{\sigma_\nu^2}{M} \Lambda_E^{-1} \right)^{-1} \frac{G^T}{M} \quad (24)$$

Our focus: Karhunen-Loève expansions

$$\mathbf{y} = G\mathbf{a} + Z\mathbf{b} + \boldsymbol{\nu} \quad (25)$$

$$\widehat{f}_E(x) := [\phi_1(x) \quad \cdots \quad \phi_E(x)] H \mathbf{y} \quad H := \left(\frac{G^T G}{M} + \frac{\sigma_{\nu}^2}{M} \Lambda_E^{-1} \right)^{-1} \frac{G^T}{M} \quad (26)$$

computational cost: $O(E^3)$

How well does \widehat{f}_E perform w.r.t. $\mathbb{E} \left[\|f(x) - \widehat{f}_E(x)\|^2 \right]$?

How well does \widehat{f}_E perform w.r.t. $\mathbb{E} \left[\|f(x) - \widehat{f}_E(x)\|^2 \right]$?

our aim: bound the statistical performance of \widehat{f}_E as a function of E

How well does \widehat{f}_E perform w.r.t. $\mathbb{E} \left[\|f(x) - \widehat{f}_E(x)\|^2 \right]$?

our aim: bound the statistical performance of \widehat{f}_E as a function of E

Key quantities:

How well does \widehat{f}_E perform w.r.t. $\mathbb{E} \left[\|f(x) - \widehat{f}_E(x)\|^2 \right]$?

our aim: bound the statistical performance of \widehat{f}_E as a function of E

Key quantities:

- $\mathbb{E} \left[\|f - \widehat{f}_E\|^2 \mid \boldsymbol{x} \right]$

How well does \widehat{f}_E perform w.r.t. $\mathbb{E} \left[\|f(x) - \widehat{f}_E(x)\|^2 \right]$?

our aim: bound the statistical performance of \widehat{f}_E as a function of E

Key quantities:

- $\mathbb{E} \left[\|f - \widehat{f}_E\|^2 \mid \mathbf{x} \right]$

- $\mathbb{E} \left[\mathbb{E} \left[\|f - \widehat{f}_E\|^2 \mid \mathbf{x} \right] \mid \mathcal{E} \right] \quad \mathbb{P}[\mathcal{E}] \geq 1 - \alpha$

How well does \widehat{f}_E perform w.r.t. $\mathbb{E} \left[\|f(x) - \widehat{f}_E(x)\|^2 \right]$?

our aim: bound the statistical performance of \widehat{f}_E as a function of E

Key quantities:

- $\mathbb{E} \left[\|f - \widehat{f}_E\|^2 \mid \mathbf{x} \right]$

- $\mathbb{E} \left[\mathbb{E} \left[\|f - \widehat{f}_E\|^2 \mid \mathbf{x} \right] \mid \mathcal{E} \right] \quad \mathbb{P}[\mathcal{E}] \geq 1 - \alpha$

$\alpha \in (0, 1) =$ *desired confidence level, e.g., $\alpha = 0.01$ or $\alpha = 0.05$*

How well does \widehat{f}_E perform w.r.t. $\mathbb{E} \left[\|f(x) - \widehat{f}_E(x)\|^2 \right]$?

our aim: bound the statistical performance of \widehat{f}_E as a function of E

Key quantities:

- $\mathbb{E} \left[\|f - \widehat{f}_E\|^2 \mid \mathbf{x} \right]$

- $\mathbb{E} \left[\mathbb{E} \left[\|f - \widehat{f}_E\|^2 \mid \mathbf{x} \right] \mid \mathcal{E} \right] \quad \mathbb{P}[\mathcal{E}] \geq 1 - \alpha$

$\alpha \in (0, 1) =$ *desired confidence level, e.g., $\alpha = 0.01$ or $\alpha = 0.05$*

- $k := \sup_{e \in \mathbb{N}, x \in \mathcal{X}} |\phi_e(x)|^2$

- $\varepsilon \in (0, 1] =$ opportune distance index between $\frac{G^T G}{M}$ and I

How well does \widehat{f}_E perform w.r.t. $\mathbb{E} \left[\|f(x) - \widehat{f}_E(x)\|^2 \right]$?

in words

if M is sufficiently¹ big w.r.t. E

then with at least probability $1 - \alpha$ the expected performance of \widehat{f}_E is upper bounded by the following Bnd, that is computable a-priori

$$\text{Bnd} := \frac{kM}{1 - \alpha} \left(\sum_{e=1}^E \frac{\lambda_e^2}{(\varepsilon M \lambda_e + \sigma_\nu^2)^2} \right) \left(\sum_{e=E+1}^{+\infty} \lambda_e \right) + \frac{\sigma_\nu^2}{1 - \alpha} \left(\sum_{e=1}^E \frac{\lambda_e}{\varepsilon M \lambda_e + \sigma_\nu^2} \right) + \left(\sum_{e=E+1}^{+\infty} \lambda_e \right) \quad (27)$$

¹With different k, α and ε defining different “sufficiently”

How well does \widehat{f}_E perform w.r.t. $\mathbb{E} \left[\|f(x) - \widehat{f}_E(x)\|^2 \right]$?

formally

How well does \widehat{f}_E perform w.r.t. $\mathbb{E} \left[\|f(x) - \widehat{f}_E(x)\|^2 \right]$?

formally

if $E, M, k, \alpha, \varepsilon$ satisfy $1 - \varepsilon + \varepsilon \log(\varepsilon) \geq \frac{Ek}{M} \log\left(\frac{E}{\alpha}\right)$ then

$$\mathbb{P} \left[\mathbb{E} \left[\mathbb{E} \left[\|f - \widehat{f}_E\|^2 \mid \mathbf{x} \right] \mid \mathcal{E} \right] \leq \text{Bnd} \right] \geq 1 - \alpha \quad (28)$$

with

$$\text{Bnd} := \frac{kM}{1 - \alpha} \left(\sum_{e=1}^E \frac{\lambda_e^2}{(\varepsilon M \lambda_e + \sigma_\nu^2)^2} \right) \left(\sum_{e=E+1}^{+\infty} \lambda_e \right) + \frac{\sigma_\nu^2}{1 - \alpha} \left(\sum_{e=1}^E \frac{\lambda_e}{\varepsilon M \lambda_e + \sigma_\nu^2} \right) + \left(\sum_{e=E+1}^{+\infty} \lambda_e \right) \quad (29)$$

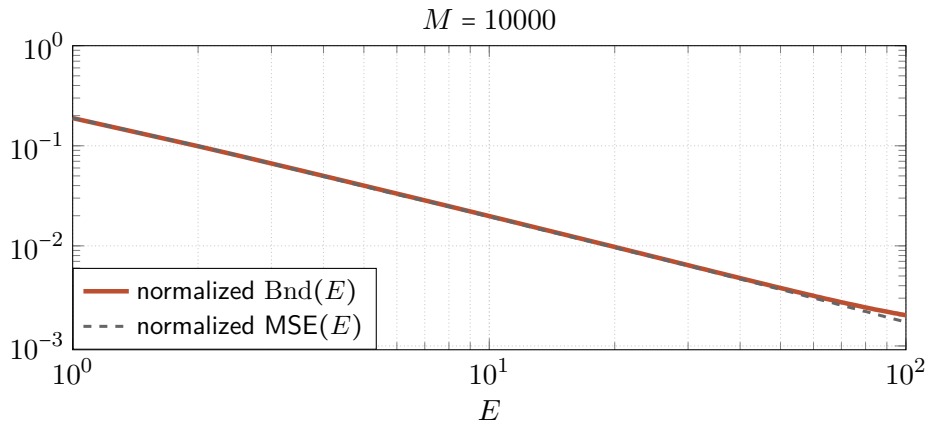
What do we enable?

given M and K , how big should E be
so to have a certain expected statistical performance?

$$\text{Bnd} := \frac{kM}{1-\alpha} \left(\sum_{e=1}^E \frac{\lambda_e^2}{(\varepsilon M \lambda_e + \sigma_\nu^2)^2} \right) \left(\sum_{e=E+1}^{+\infty} \lambda_e \right) + \frac{\sigma_\nu^2}{1-\alpha} \left(\sum_{e=1}^E \frac{\lambda_e}{\varepsilon M \lambda_e + \sigma_\nu^2} \right) + \left(\sum_{e=E+1}^{+\infty} \lambda_e \right) \quad (30)$$

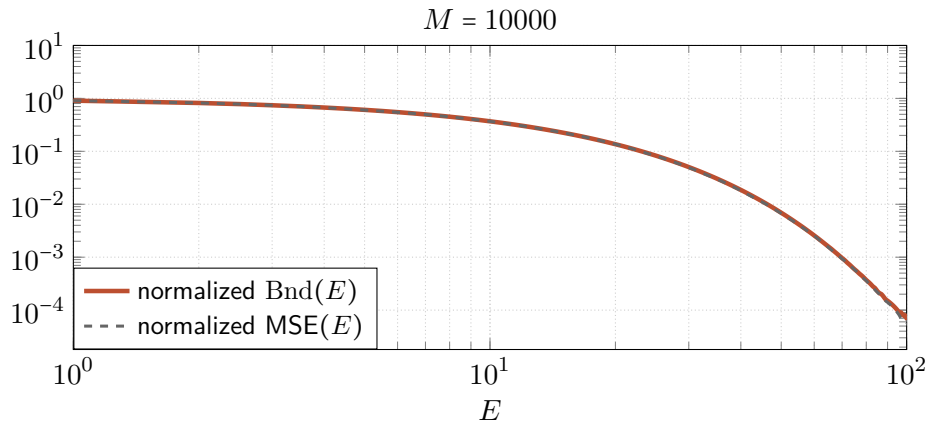
How significant is Bnd?

Case splines, i.e., $K(x, x') = \min(x, x')$ $\mathcal{X} = [0, 1]$ $\lambda_e = \frac{1}{(e\pi - \pi/2)^2}$

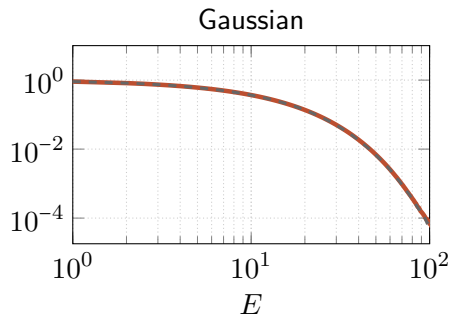
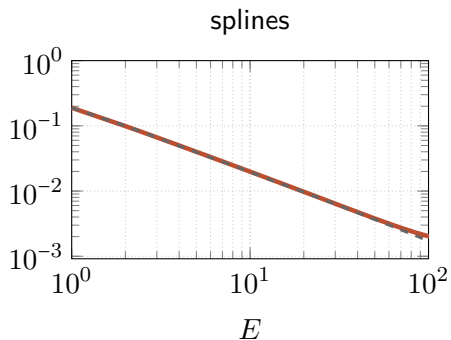


How significant is Bnd?

Case Gaussian, i.e., $K(x, x') = \exp(-\|x - x'\|^2)$ $\mathcal{X} = [0, 1]$ $\lambda_e = \exp(-0.1e)$



What can we do with this result?



understand how big E should be a priori so to obtain certain statistical performance

What do we enable?

PAMI extension

*strategies for distributedly tuning
the hyperparameters of distributed estimators
through the minimization of the bound
(both a priori and a posteriori)*

Statistical bounds for Gaussian regression algorithms based on Karhunen-Loève expansions

Gianluigi Pillonetto Luca Schenato Damiano Varagnolo



`damiano.varagnolo@ltu.se`