

Contraction rates for conjugate gradient and Lanczos approximate posteriors in Gaussian process regression

Bernhard Stankewitz

Bernoulli-ims 11th World Congress in Probability and Statistics

Bochum, August 2024

*Department of Decision Sciences
Bocconi University*

Joint work with



Botond Szabo, Bocconi Milano

Gaussian process (GP) regression

Consider i.i.d. observations from the model

$$Y_i = F(X_i) + \varepsilon_i, \quad i = 1, \dots, n, \quad (1)$$

where

- ▶ $X_1, \dots, X_n \sim G$ i.i.d. on \mathbb{R}^d and $\varepsilon \sim N(0, \sigma^2 I_n)$;
- ▶ $F \sim \text{GP}(0, k)$ with p.s.d. kernel $k : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$ is a GP-prior on $L^2(G)$, i.e.

$$\mathbb{E}F(x) = 0, \quad \text{Cov}(F(x), F(x')) = k(x, x'), \quad x, x' \in \mathbb{R}^d. \quad (2)$$

Gaussian process (GP) regression

Consider i.i.d. observations from the model

$$Y_i = F(X_i) + \varepsilon_i, \quad i = 1, \dots, n, \quad (1)$$

where

- ▶ $X_1, \dots, X_n \sim G$ i.i.d. on \mathbb{R}^d and $\varepsilon \sim N(0, \sigma^2 I_n)$;
- ▶ $F \sim \text{GP}(0, k)$ with p.s.d. kernel $k : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$ is a GP-prior on $L^2(G)$, i.e.

$$\mathbb{E}F(x) = 0, \quad \text{Cov}(F(x), F(x')) = k(x, x'), \quad x, x' \in \mathbb{R}^d. \quad (2)$$

Setting $K := (k(X_i, X_j))_{i,j=1}^n \in \mathbb{R}^{n \times n}$ and $k(X, x) := (k(X_i, x))_{i=1}^n \in \mathbb{R}^n$, the posterior $\Pi(\cdot | X, Y)$ is given by the GP with mean and covariance function

$$\begin{aligned} x &\mapsto k(X, x)^\top (K + \sigma^2 I_n)^{-1} Y \\ (x, x') &\mapsto k(x, x') - k(X, x)^\top (K + \sigma^2 I_n)^{-1} k(X, x'). \end{aligned} \quad (3)$$

Gaussian process (GP) regression

Consider i.i.d. observations from the model

$$Y_i = F(X_i) + \varepsilon_i, \quad i = 1, \dots, n, \quad (1)$$

where

- ▶ $X_1, \dots, X_n \sim G$ i.i.d. on \mathbb{R}^d and $\varepsilon \sim N(0, \sigma^2 I_n)$;
- ▶ $F \sim \text{GP}(0, k)$ with p.s.d. kernel $k : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$ is a GP-prior on $L^2(G)$, i.e.

$$\mathbb{E}F(x) = 0, \quad \text{Cov}(F(x), F(x')) = k(x, x'), \quad x, x' \in \mathbb{R}^d. \quad (2)$$

Setting $K := (k(X_i, X_j))_{i,j=1}^n \in \mathbb{R}^{n \times n}$ and $k(X, x) := (k(X_i, x))_{i=1}^n \in \mathbb{R}^n$, the posterior $\Pi(\cdot | X, Y)$ is given by the GP with mean and covariance function

$$\begin{aligned} x &\mapsto k(X, x)^\top (K + \sigma^2 I_n)^{-1} Y \\ (x, x') &\mapsto k(x, x') - k(X, x)^\top (K + \sigma^2 I_n)^{-1} k(X, x'). \end{aligned} \quad (3)$$

Motivating problem.

The computation of $(K + \sigma^2 I)^{-1}$ has a computational complexity of $O(n^3)$, which becomes infeasible for large n .

Idea: Focussing on the posterior mean $k(X, x)^\top (K + \sigma^2 I_n)^{-1} Y$, iteratively solve $(K + \sigma^2 I_n)W = Y$ for the representer weights W .

- ▶ Consider a Bayesian updating scheme with initial beliefs $W = (K + \sigma^2 I_n)^{-1} Y \sim N(0, (K + \sigma^2 I_n)^{-1}) =: N(w_0, \Gamma_0)$.

¹J. Wenger et al. "Posterior and computational uncertainty in Gaussian processes." In: *Advances in Neural Information Processing Systems* (2022).

Idea: Focussing on the posterior mean $k(X, x)^\top (K + \sigma^2 I_n)^{-1} Y$, iteratively solve $(K + \sigma^2 I_n)W = Y$ for the representer weights W .

- ▶ Consider a Bayesian updating scheme updating scheme with initial beliefs $W = (K + \sigma^2 I_n)^{-1} Y \sim N(0, (K + \sigma^2 I_n)^{-1}) =: N(w_0, \Gamma_0)$.
- ▶ Consecutively update by conditioning on the information projection

$$\alpha_j := s_j^\top (Y - (K + \sigma^2 I_n)w_{j-1}), \quad j \leq m \quad (4)$$

where $s_j, j \leq m$ are search directions chosen by the user.

¹J. Wenger et al. "Posterior and computational uncertainty in Gaussian processes." In: *Advances in Neural Information Processing Systems* (2022).

Idea: Focussing on the posterior mean $k(X, x)^\top (K + \sigma^2 I_n)^{-1} Y$, iteratively solve $(K + \sigma^2 I_n)W = Y$ for the representer weights W .

- ▶ Consider a Bayesian updating scheme updating scheme with initial believes $W = (K + \sigma^2 I_n)^{-1} Y \sim N(0, (K + \sigma^2 I_n)^{-1}) =: N(w_0, \Gamma_0)$.
- ▶ Consecutively update by conditioning on the information projection

$$\alpha_j := s_j^\top (Y - (K + \sigma^2 I_n)w_{j-1}), \quad j \leq m \quad (4)$$

where $s_j, j \leq m$ are search directions chosen by the user.

- ▶ After m steps, believes are given by $N(w_m, \Gamma_m) = N(C_m Y, (K + \sigma^2)^{-1} - C_m)$. This yields the approximate Gaussian posterior $\Psi_m := \mathbb{P}^{F|W=w} N(w_m, \Gamma_m)(dw)$ with mean and covariance functions

$$x \mapsto k(X, x)^\top C_m Y \quad (x, x') \mapsto k(x, x') - k(X, x)^\top C_m k(X, x'), \quad (5)$$

where C_m is a rank m matrix approximating $(K + \sigma^2 I_n)^{-1}$.

¹J. Wenger et al. "Posterior and computational uncertainty in Gaussian processes." In: *Advances in Neural Information Processing Systems* (2022).

The Empirical eigenvector posterior

Consider the spectral decomposition of the empirical kernel matrix

$$K = \sum_{j=1}^n \hat{\mu}_j \hat{u}_j \hat{u}_j^\top \quad (6)$$

and choose the search directions $s_j := \hat{u}_j, j \leq m$.

²D. Nieman, B.Szabo and H. van Zanten. "Contraction rates for sparse variational approximations in Gaussian process regression". In: *Journal of Machine Learning Research* 23 (2022).

The Empirical eigenvector posterior

Consider the spectral decomposition of the empirical kernel matrix

$$K = \sum_{j=1}^n \hat{\mu}_j \hat{u}_j \hat{u}_j^\top \quad (6)$$

and choose the search directions $s_j := \hat{u}_j$, $j \leq m$. Then, the approximate posterior $\Psi_m = \Psi_m^{\text{EV}}$ is given by the mean and covariance function

$$x \mapsto k(X, x)^\top C_m Y \quad (x, x') \mapsto k(x, x') - k(X, x)^\top C_m k(X, x'), \quad (7)$$

where $(K + \sigma^2 I_n)^{-1}$ is approximated by

$$C_m = C_m^{\text{EV}} = \sum_{j=1}^m (\hat{\mu}_j + \sigma^2)^{-1} \hat{u}_j \hat{u}_j^\top. \quad (8)$$

²D. Nieman, B.Szabo and H. van Zanten. "Contraction rates for sparse variational approximations in Gaussian process regression". In: *Journal of Machine Learning Research* 23 (2022).

The Empirical eigenvector posterior

Consider the spectral decomposition of the empirical kernel matrix

$$K = \sum_{j=1}^n \hat{\mu}_j \hat{u}_j \hat{u}_j^\top \quad (6)$$

and choose the search directions $s_j := \hat{u}_j$, $j \leq m$. Then, the approximate posterior $\Psi_m = \Psi_m^{\text{EV}}$ is given by the mean and covariance function

$$x \mapsto k(X, x)^\top C_m Y \quad (x, x') \mapsto k(x, x') - k(X, x)^\top C_m k(X, x'), \quad (7)$$

where $(K + \sigma^2 I_n)^{-1}$ is approximated by

$$C_m = C_m^{\text{EV}} = \sum_{j=1}^m (\hat{\mu}_j + \sigma^2)^{-1} \hat{u}_j \hat{u}_j^\top. \quad (8)$$

The Ψ_m^{EV} is equivalent to the Variational Bayes posterior based on spectral inducing variables [NSZ22]²

²D. Nieman, B.Szabo and H. van Zanten. "Contraction rates for sparse variational approximations in Gaussian process regression". In: *Journal of Machine Learning Research* 23 (2022).

The Lanczos posterior

Consider the spectral decomposition of the empirical kernel matrix

$$K = \sum_{j=1}^n \hat{\mu}_j \hat{u}_j \hat{u}_j^\top \quad (9)$$

and choose the search directions $s_j := \tilde{u}_j, j \leq m$, where $(\tilde{\mu}_j, \tilde{u}_j)_{j \leq m}$ is the Lanczos approximate eigensystem up to order m .

The Lanczos posterior

Consider the spectral decomposition of the empirical kernel matrix

$$K = \sum_{j=1}^n \hat{\mu}_j \hat{u}_j \hat{u}_j^\top \quad (9)$$

and choose the search directions $s_j := \tilde{u}_j, j \leq m$, where $(\tilde{\mu}_j, \tilde{u}_j)_{j \leq m}$ is the Lanczos approximate eigensystem up to order m . Then, the approximate posterior $\Psi_m = \Psi_m^L$ is given by the mean and covariance function

$$x \mapsto k(X, x)^\top C_m Y \quad (x, x') \mapsto k(x, x') - k(X, x)^\top C_m k(X, x'), \quad (10)$$

with

$$C_m = C_m^L = \sum_{j=1}^m (\tilde{\mu}_j + \sigma^2)^{-1} \tilde{u}_j \tilde{u}_j^\top. \quad (11)$$

The Lanczos posterior

Consider the spectral decomposition of the empirical kernel matrix

$$K = \sum_{j=1}^n \hat{\mu}_j \hat{u}_j \hat{u}_j^\top \quad (9)$$

and choose the search directions $s_j := \tilde{u}_j, j \leq m$, where $(\tilde{\mu}_j, \tilde{u}_j)_{j \leq m}$ is the Lanczos approximate eigensystem up to order m . Then, the approximate posterior $\Psi_m = \Psi_m^L$ is given by the mean and covariance function

$$x \mapsto k(X, x)^\top C_m Y \quad (x, x') \mapsto k(x, x') - k(X, x)^\top C_m k(X, x'), \quad (10)$$

with

$$C_m = C_m^L = \sum_{j=1}^m (\tilde{\mu}_j + \sigma^2)^{-1} \tilde{u}_j \tilde{u}_j^\top. \quad (11)$$

Randomness of the kernel matrix

Since $K = (k(X_i, X_j))_{i,j \leq n}$ is a random matrix, the spectral decomposition of K cannot be computed in advance.

Conjugate gradient descent. Iteratively solve $(K + \sigma^2 I_n)w = Y$ by setting $w_0 = 0$ and for $j \geq 1$,

$$\varrho(w_j) = \min_{t \in \mathbb{R}} \varrho(w_{j-1} + td_j^{\text{CG}}), \quad (12)$$

where $\varrho(w) := (w^\top (K + \sigma^2 I_n)w)/2 - Y^\top w$, and the $(d_j^{\text{CG}})_{j \geq 1}$ are conjugate search directions satisfying $(d_j^{\text{CG}})^\top (K + \sigma^2 I_n) d_k^{\text{CG}} = 0, j \neq k$.

Conjugate gradient descent. Iteratively solve $(K + \sigma^2 I_n)w = Y$ by setting $w_0 = 0$ and for $j \geq 1$,

$$\varrho(w_j) = \min_{t \in \mathbb{R}} \varrho(w_{j-1} + td_j^{\text{CG}}), \quad (12)$$

where $\varrho(w) := (w^\top (K + \sigma^2 I_n)w)/2 - Y^\top w$, and the $(d_j^{\text{CG}})_{j \geq 1}$ are conjugate search directions satisfying $(d_j^{\text{CG}})^\top (K + \sigma^2 I_n) d_k^{\text{CG}} = 0, j \neq k$.

For the policies $s_j := d_j^{\text{CG}}, j \leq m$, Bayesian updating is equivalent to the CG-iteration and we obtain the approximate posterior Ψ_m^{CG} given by

$$x \mapsto k(X, x)^\top C_m Y \quad (x, x') \mapsto k(x, x') - k(X, x)^\top C_m k(X, x'), \quad (13)$$

where $C_m = C_m^{\text{CG}}$ is given by the implicit approximation of the inverse provided by CG.

Conjugate gradient descent. Iteratively solve $(K + \sigma^2 I_n)w = Y$ by setting $w_0 = 0$ and for $j \geq 1$,

$$\varrho(w_j) = \min_{t \in \mathbb{R}} \varrho(w_{j-1} + td_j^{\text{CG}}), \quad (12)$$

where $\varrho(w) := (w^\top (K + \sigma^2 I_n)w)/2 - Y^\top w$, and the $(d_j^{\text{CG}})_{j \geq 1}$ are conjugate search directions satisfying $(d_j^{\text{CG}})^\top (K + \sigma^2 I_n) d_k^{\text{CG}} = 0, j \neq k$.

For the policies $s_j := d_j^{\text{CG}}, j \leq m$, Bayesian updating is equivalent to the CG-iteration and we obtain the approximate posterior Ψ_m^{CG} given by

$$x \mapsto k(X, x)^\top C_m Y \quad (x, x') \mapsto k(x, x') - k(X, x)^\top C_m k(X, x'), \quad (13)$$

where $C_m = C_m^{\text{CG}}$ is given by the implicit approximation of the inverse provided by CG.

GPU accelerated matrix vector multiplication

CG only relies on matrix vector multiplications, which can be GPU accelerated and makes CG particularly relevant for large scale applications, see Wang et al.

[Wan+19].

Conjugate gradient descent. Iteratively solve $(K + \sigma^2 I_n)w = Y$ by setting $w_0 = 0$ and for $j \geq 1$,

$$\varrho(w_j) = \min_{t \in \mathbb{R}} \varrho(w_{j-1} + td_j^{\text{CG}}), \quad (12)$$

where $\varrho(w) := (w^\top (K + \sigma^2 I_n)w)/2 - Y^\top w$, and the $(d_j^{\text{CG}})_{j \geq 1}$ are conjugate search directions satisfying $(d_j^{\text{CG}})^\top (K + \sigma^2 I_n) d_k^{\text{CG}} = 0, j \neq k$.

For the policies $s_j := d_j^{\text{CG}}, j \leq m$, Bayesian updating is equivalent to the CG-iteration and we obtain the approximate posterior Ψ_m^{CG} given by

$$x \mapsto k(X, x)^\top C_m Y \quad (x, x') \mapsto k(x, x') - k(X, x)^\top C_m k(X, x'), \quad (13)$$

where $C_m = C_m^{\text{CG}}$ is given by the implicit approximation of the inverse provided by CG.

Reduction in computational complexity

The approximate inversions C_m^L, C_m^{CG} have a computation cost of $O(mn^2)$, which is feasible when $m \ll n$.

A stylized approximate contraction result

Theorem (Approximate posterior contraction)

For $f_0 \in \overline{\mathbb{H}}$ with $\mathbb{H} = \text{ran } T_k^{1/2}$,

$$T_k : L^2(G) \rightarrow L^2(G), \quad f \mapsto \int f(y)k(\cdot, y) G(dy) = \sum_{j=1}^{\infty} \lambda_j \langle f, \phi_j \rangle_{L^2(G)} \phi_j, \quad (14)$$

let \mathbb{P}_{f_0} be the measure corresponding to the data generating process

$$Y_i = f_0(X_i) + \varepsilon_i, \quad i = 1, \dots, n \quad (15)$$

Then, the *true posterior* Π_n satisfies that for any sequence $M_n \rightarrow \infty$,

$$\Pi_n(\{f \in L^2(G) : d(f, f_0) \geq M_n \varepsilon_n\} | X, Y) \rightarrow 0 \quad (16)$$

in probability under $\mathbb{P}_{f_0}^{\otimes n}$ and $n \rightarrow \infty$, where ε_n is the optimal achievable rate implied by a concentration function inequality.

A stylized approximate contraction result

Theorem (Approximate posterior contraction)

For $f_0 \in \overline{\mathbb{H}}$ with $\mathbb{H} = \text{ran } T_k^{1/2}$,

$$T_k : L^2(G) \rightarrow L^2(G), \quad f \mapsto \int f(y)k(\cdot, y) G(dy) = \sum_{j=1}^{\infty} \lambda_j \langle f, \phi_j \rangle_{L^2(G)} \phi_j, \quad (14)$$

let \mathbb{P}_{f_0} be the measure corresponding to the data generating process

$$Y_i = f_0(X_i) + \varepsilon_i, \quad i = 1, \dots, n \quad (15)$$

Then, the *approximate posterior* Ψ_m satisfies that for any sequence $M_n \rightarrow \infty$,

$$\Psi_{m_n}(\{f \in L^2(G) : d(f, f_0) \geq M_n \varepsilon_n\} | X, Y) \rightarrow 0 \quad (16)$$

in probability under $\mathbb{P}_{f_0}^{\otimes n}$ and $n \rightarrow \infty$, where ε_n is the optimal achievable rate implied by a concentration function inequality and $m_n \rightarrow \infty$ is an *appropriate sequence*.

Example: Polynomially decaying eigenvalues

For an ONB $(\phi_j)_{j \geq 1}$ of $L^2(G)$ and $Z_j \sim N(0, 1)$ i.i.d., consider the random series prior

$$F(x) = \sum_{j=1}^{\infty} \tau j^{-1/2-\alpha/d} Z_j \phi_j(x), \quad x \in \mathbb{R}^d \quad (17)$$

where $\alpha > 0$ and τ are the regularity and scale hyperparameters of the process. Then, for any

$$f_0 \in S^\beta(L) := \{f \in L^2(G) : \|f\|_{S^\beta}^2 \leq L\} \quad \text{with} \quad \|f\|_{S^\beta}^2 := \sum_{j=1}^{\infty} j^{2\beta/d} \langle f, \phi_j \rangle^2, \quad (18)$$

with $d/2 < \beta \leq \alpha + d/2$ and an appropriate choice of τ , the approximate posterior satisfies that for any $M_n \rightarrow \infty$,

$$\Psi_{m_n} \{f : d_H(f, f_0) \geq M_n n^{-\beta/(d+2\beta)} | \mathcal{X}, Y\} \rightarrow 0, \quad (19)$$

in probability under $\mathbb{P}_{f_0}^{\otimes n}$ and $n \rightarrow \infty$ with $m_n \sim n^{d/(2\beta+d)} \log n$.

Simulation example

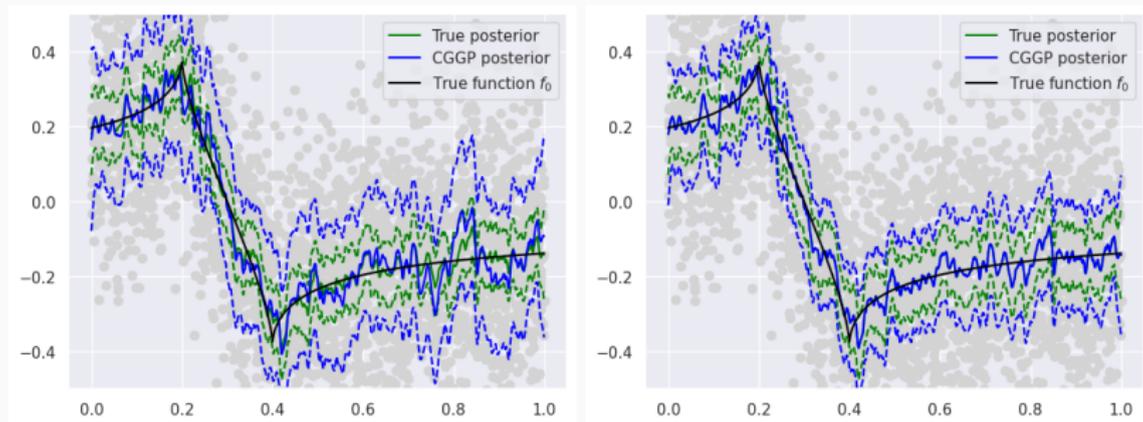


Figure 1: Simulation results for $n = 3000$, $m = 20, 40$.

Simulation example

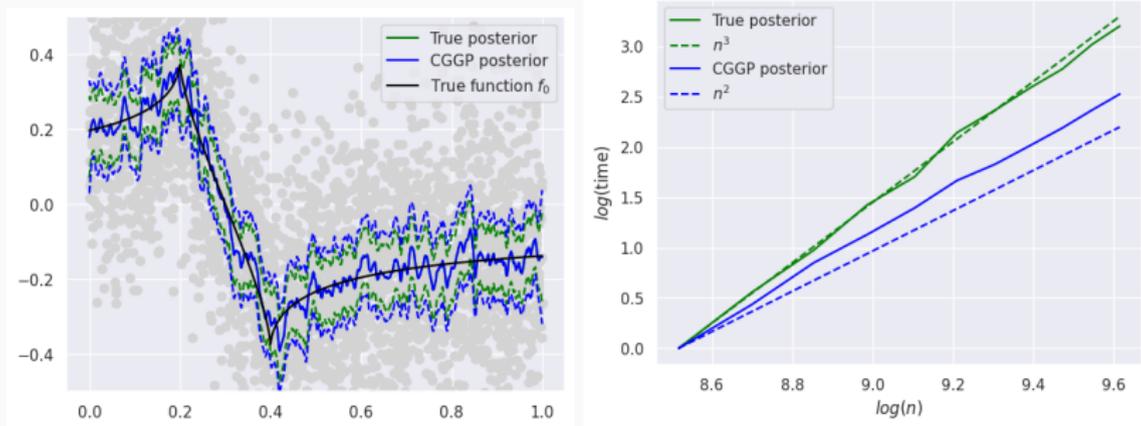


Figure 2: Simulation results for $n = 3000$, $m = 80$ and scaling of computation times.

Some conclusions

- ▶ Our theory is the first providing statistical guarantees for fully numerical algorithms.
- ▶ Particular relevance in the CG posterior. Default method in the GPyTorch library, see Gardner et al. [Gar+18].

- ▶ For $\text{KL}(\Psi_{m_n}, \Pi_n(\cdot|X, Y)) \leq n\varepsilon_n^2$, the approximate posterior Ψ_{m_n} inherits the contraction rate ε_n , see Ray and Szabó [RS19].
- ▶ For the empirical eigenvector posterior with $s_j = \hat{u}_j$, $j \leq m$ this bound is available via elementary tools.
- ▶ Analyze the Lanczos posterior as an approximation.

Theorem (Lanczos: Eigenvalue bound, [Saa80])

Under Assumption (LWdf), for any fixed integer $i \leq \tilde{m} < n$ with $\tilde{\lambda}_{i-1} > \hat{\lambda}_i$ if $i > 1$ and any integer $\tilde{p} \leq \tilde{m} - i$, the eigenvalue approximation satisfies

$$0 \leq \hat{\lambda}_i - \tilde{\lambda}_i \leq (\hat{\lambda}_i - \hat{\lambda}_n) \left(\frac{\tilde{\kappa}_i \kappa_{i,\tilde{p}} \tan(\hat{u}_i, v_0)}{T_{\tilde{m}-i-\tilde{p}}(\gamma_i)} \right)^2, \quad (20)$$

where $\gamma_i := 1 + 2(\hat{\lambda}_i - \hat{\lambda}_{i+\tilde{p}+1})/(\hat{\lambda}_{i+\tilde{p}+1} - \hat{\lambda}_n)$,

$$\tilde{\kappa}_i := \prod_{j=1}^{i-1} \frac{\tilde{\lambda}_j - \hat{\lambda}_n}{\tilde{\lambda}_j - \hat{\lambda}_i}, \quad \kappa_{i,\tilde{p}} := \prod_{j=i+1}^{i+\tilde{p}} \frac{\hat{\lambda}_j - \hat{\lambda}_n}{\hat{\lambda}_i - \hat{\lambda}_j}, \quad (21)$$

and T_l denotes the l -th Tschebychev polynomial.

Theorem (Eigenvalue concentration, Shawe-Taylor and Williams [STW02])

The empirical eigenvalues $(\hat{\lambda}_j)_{j \leq n}$ of the normalized kernel matrix K/n satisfy

(i) For any $t > 0$ and any fixed $m \geq 1$, both

$$\mathbb{P}\{|\hat{\lambda}_m - \mathbb{E}\hat{\lambda}_m| \geq t\} \leq 2 \exp\left(\frac{-2nt^2}{\max_x k(x, x)^4}\right) \quad (22)$$

and

$$\mathbb{P}\left\{\left|\sum_{j=m+1}^n \hat{\lambda}_j - \mathbb{E} \sum_{j=m+1}^n \hat{\lambda}_j\right| \geq t\right\} \leq 2 \exp\left(\frac{-2nt^2}{\max_x k(x, x)^4}\right). \quad (23)$$

(ii) For any fixed $m \geq 1$,

$$\mathbb{E} \sum_{j=1}^m \hat{\lambda}_j \geq \sum_{j=1}^m \lambda_j \quad \text{and} \quad \mathbb{E} \sum_{j=m+1}^n \hat{\lambda}_j \leq \sum_{j=m+1}^{\infty} \lambda_j. \quad (24)$$

Proposition (Relative perturbation bounds, Jirak and Wahl [JW23])

Under appropriate assumptions, fix $m \in \mathbb{N}$ and further assume that

$$r_j(T_k) := \sum_{k \neq j} \frac{\lambda_k}{|\lambda_j - \lambda_k|} + \frac{\lambda_j}{(\lambda_{j-1} - \lambda_j) \wedge (\lambda_j - \lambda_{j+1})} \leq C \sqrt{\frac{n}{\log n}}, \quad (22)$$

for all $j \leq m$. Then, the eigenvalues of K/n satisfy the relative perturbation bound

$$\left| \frac{\hat{\lambda}_j - \lambda_j}{\lambda_j} \right| \leq C \sqrt{\frac{\log n}{n}} \quad \text{for all } j \leq m \quad (23)$$

with high probability.

Challenges from spectral concentration

Proposition (Relative perturbation bounds, Jirak and Wahl [JW23])

Under appropriate assumptions, fix $m \in \mathbb{N}$ and further assume that

$$r_j(T_k) := \sum_{k \neq j} \frac{\lambda_k}{|\lambda_j - \lambda_k|} + \frac{\lambda_j}{(\lambda_{j-1} - \lambda_j) \wedge (\lambda_j - \lambda_{j+1})} \leq C \sqrt{\frac{n}{\log n}}, \quad (22)$$

for all $j \leq m$. Then, the eigenvalues of K/n satisfy the relative perturbation bound

$$\left| \frac{\hat{\lambda}_j - \lambda_j}{\lambda_j} \right| \leq C \sqrt{\frac{\log n}{n}} \quad \text{for all } j \leq m \quad (23)$$

with high probability.



Martin Wahl

Ongoing joint work on perturbation series for empirical eigenvalues and eigenprojectors.

- ▶ For $\text{KL}(\Psi_{m_n}, \Pi_n(\cdot|X, Y)) \leq n\varepsilon_n^2$, the approximate posterior Ψ_{m_n} inherits the contraction rate ε_n , see Ray and Szabó [RS19].
- ▶ For the empirical eigenvector posterior with $s_j = \hat{u}_j$, $j \leq m$ this bound is available via elementary tools.
- ▶ Analyze the Lanczos posterior as an approximation.
- ▶ Establish the equivalence of the CG and the Lanczos posterior.

- ▶ Our theory is the first providing statistical guarantees for fully numerical algorithms.
- ▶ Particular relevance lies in the CG posterior. Default method in the GPyTorch library, see Gardner et al. [Gar+18].
- ▶ New interpretation of the CG posterior as a numerical approximation of a variational Bayes method.



Preprint



Author page

Thank you!

References

- [Gar+18] J. Gardner et al. **“GPYtorch: Blackbox Matrix-Matrix Gaussian Process Inference with GPU Acceleration”**. In: *Advances in Neural Information Processing Systems*. Vol. 31. Curran Associates, Inc., 2018.
- [JW23] M. Jirak and M. Wahl. **“Relative perturbation bounds with applications to empirical covariance operators”**. In: *Advances in Mathematics* 412 (2023), p. 108808.
- [NSZ22] D. Nieman, B. Szabo, and H. van Zanten. **“Contraction rates for sparse variational approximations in Gaussian process regression”**. In: *Journal of Machine Learning Research* 23 (2022), pp. 1–26.
- [RS19] K. Ray and B. Szabó. **“Variational Bayes for High-Dimensional Linear Regression With Sparse Priors”**. In: *Journal of the American Statistical Association* (2019).
- [STW02] J. Shawe-Taylor and C. K. I. Williams. **“The stability of kernel Principal component analysis and its relation to the process eigenspectrum”**. In: *Advances in Neural Information Processing Systems*. 2002.

- [Saa80] Y. Saad. **“On the rates of convergence of the Lanczos and the Block-Lanczos methods”**. In: *SIAM Journal on Numerical Analysis* 17.5 (1980), pp. 687–706.
- [Wan+19] K. Wang et al. **“Exact Gaussian Processes on a Million Data Points”**. In: *Advances in Neural Information Processing Systems*. Vol. 32. Curran Associates, Inc., 2019.
- [Wen+22] J. Wenger et al. **“Posterior and computational uncertainty in Gaussian processes”**. In: *Advances in Neural Information Processing Systems*. 2022.