

Modern Gaussian Processes: Scalable Inference and Novel Applications

(Part II-a) Model Approximations

Edwin V. Bonilla and **Maurizio Filippone**

CSIRO's Data61, Sydney, Australia and EURECOM, Sophia Antipolis, France

July 14th, 2019



① Random Feature Expansions

② Low-Rank Approximations

Inducing Variables

Structured Approximations

Random Feature Expansions

Bochner's theorem

- Continuous shift-invariant covariance function

$$k(\mathbf{x}_i - \mathbf{x}_j | \theta) = \sigma^2 \int p(\boldsymbol{\omega} | \theta) \exp\left(i(\mathbf{x}_i - \mathbf{x}_j)^\top \boldsymbol{\omega}\right) d\boldsymbol{\omega}$$

Bochner's theorem

- Continuous shift-invariant covariance function

$$k(\mathbf{x}_i - \mathbf{x}_j | \theta) = \sigma^2 \int p(\omega | \theta) \exp\left(\iota(\mathbf{x}_i - \mathbf{x}_j)^\top \omega\right) d\omega$$

- Monte Carlo estimate

$$k(\mathbf{x}_i - \mathbf{x}_j | \theta) \approx \frac{\sigma^2}{N_{\text{RF}}} \sum_{r=1}^{N_{\text{RF}}} \mathbf{z}(\mathbf{x}_i | \tilde{\omega}_r)^\top \mathbf{z}(\mathbf{x}_j | \tilde{\omega}_r)$$

with

$$\tilde{\omega}_r \sim p(\omega | \theta)$$

$$\mathbf{z}(\mathbf{x} | \omega) = [\cos(\mathbf{x}^\top \omega), \sin(\mathbf{x}^\top \omega)]^\top$$

GPs with Random Fourier Features

- Define

$$\Phi = \sqrt{\frac{\sigma^2}{N_{\text{RF}}}} [\cos(\mathbf{X}\Omega), \sin(\mathbf{X}\Omega)]$$

and

$$\mathbf{f} = \Phi \mathbf{w}$$

- GPs become Bayesian linear models with

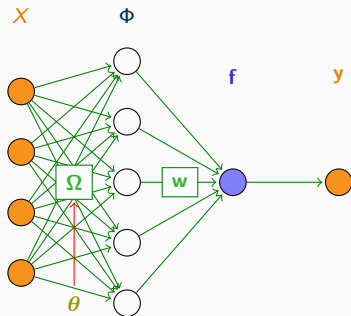
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- Low-rank approximation of \mathbf{K}

$$\text{cov}(\mathbf{f}) = \mathbb{E}[\Phi \mathbf{w} \mathbf{w}^\top \Phi^\top] = \Phi \Phi^\top \approx \mathbf{K}$$

GPs with Random Features become Bayesian Linear Models

- Neural Network-like diagram



Low-Rank Approximations

- Marginal likelihood GP regression:

$$-\frac{1}{2} \log |\mathbf{K}_y| - \frac{1}{2} \mathbf{y}^\top \mathbf{K}_y^{-1} \mathbf{y} + \text{const.}$$

- Most GP approximations aim to form a low-rank approximation to the covariance matrix

$$\mathbf{K}_y = \mathbf{K} + \sigma^2 \mathbf{I} \approx \mathbf{UCV} + \sigma^2 \mathbf{I}$$

Low-Rank Approximations

- Woodbury identity for the inverse

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1}$$

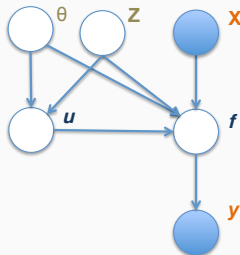
The diagram illustrates the Woodbury identity for matrix inversion using orange squares to represent matrix elements. The top part shows the inverse of the sum of a diagonal matrix \mathbf{A} and a low-rank matrix \mathbf{UCV} . The bottom part shows the expansion of this inverse as \mathbf{A}^{-1} minus a correction term involving \mathbf{A}^{-1} , \mathbf{U} , $(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}$, and \mathbf{VA}^{-1} .

- Similar for the log-determinant
- This reduces complexity from $\mathcal{O}(N^3)$ to $\mathcal{O}(M^3) + \mathcal{O}(NM^2)$ with $M \ll N$

Sparse GPs with Nyström Approximation

- Introduce M pseudo-inputs collected in Z ...
- ... and corresponding inducing variables u
- Nyström approximation

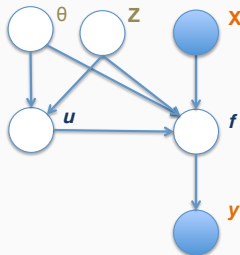
$$\mathbf{K} \approx \mathbf{K}_{\mathbf{x}Z} \mathbf{K}_{ZZ}^{-1} \mathbf{K}_{Z\mathbf{x}}$$



Sparse GPs with Nyström Approximation

- Introduce M pseudo-inputs collected in Z ...
- ... and corresponding inducing variables u
- Nyström approximations with diagonal correction

$$\mathbf{K} \approx \text{diag}(\mathbf{K} - \mathbf{K}_{\mathbf{xz}}\mathbf{K}_{\mathbf{zz}}^{-1}\mathbf{K}_{\mathbf{zx}}) + \mathbf{K}_{\mathbf{xz}}\mathbf{K}_{\mathbf{zz}}^{-1}\mathbf{K}_{\mathbf{zx}}$$



Structured Inputs

- When inputs lie on a regular 1D grid and $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \kappa(\mathbf{x}_i - \mathbf{x}_j)$
- \mathbf{K} is Toeplitz

$$\mathbf{K} = \begin{pmatrix} a & b & c & d \\ b & a & b & c \\ c & b & a & b \\ d & c & b & a \end{pmatrix}$$

- Solving \mathbf{K} exactly costs $\mathcal{O}(N \log N)$ time!

Structured Inputs

- When inputs lie on a regular grid and $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \kappa(\mathbf{x}_i - \mathbf{x}_j)$
- \mathbf{K} can be decomposed is Toeplitz

$$\mathbf{K} = \mathbf{K}_1 \otimes \dots \otimes \mathbf{K}_d$$

where \mathbf{K}_p has entries $\kappa(\mathbf{x}_i^{(p)}, \mathbf{x}_j^{(p)})$

- Algebraic operations for \mathbf{K} are based on faster ones for each factor \mathbf{K}_p in the Kronecker product

Structured Inducing Points

- Consider a sparse GP:

$$\mathbf{K} \approx \mathbf{K}_{\mathbf{xz}} \mathbf{K}_{\mathbf{zz}}^{-1} \mathbf{K}_{\mathbf{zx}}$$

- \mathbf{Z} on a grid makes the inverse fast (Toeplitz)!
- Can afford $M \gg N$
- Still expensive to deal with $\mathbf{K}_{\mathbf{zx}} \dots \mathcal{O}(NM^2)$

Structured Inducing Points

- Consider a sparse GP:

$$\mathbf{K} \approx \mathbf{K}_{\mathbf{xz}} \mathbf{K}_{\mathbf{zz}}^{-1} \mathbf{K}_{\mathbf{zx}}$$

- Kernel Interpolation (KISS-GP)

$$\mathbf{K}_{\mathbf{xz}} \approx \mathbf{W} \mathbf{K}_{\mathbf{zz}}$$

with \mathbf{W} a sparse “interpolation” matrix, so that

$$\mathbf{K} \approx \mathbf{K}_{\mathbf{xz}} \mathbf{K}_{\mathbf{zz}}^{-1} \mathbf{K}_{\mathbf{zx}} \approx \mathbf{W} \mathbf{K}_{\mathbf{zz}}^{-1} \mathbf{W}^{\top}$$

- All products/inverses are fast even if $M \gg N$!