

Computation-Aware Gaussian Process Inference

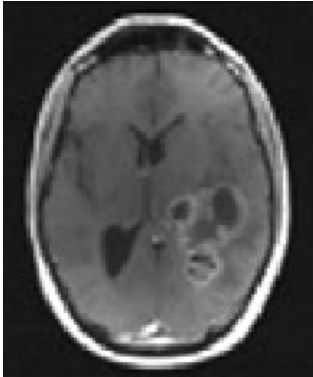
Jonathan Wenger

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



imprs-is

Accurate Reconstruction





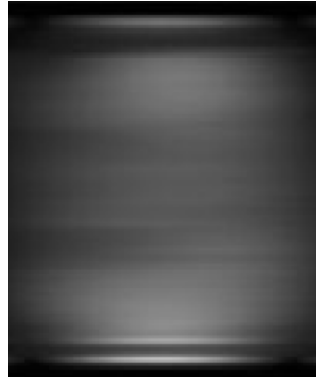
Accurate Reconstruction



Accurate Reconstruction



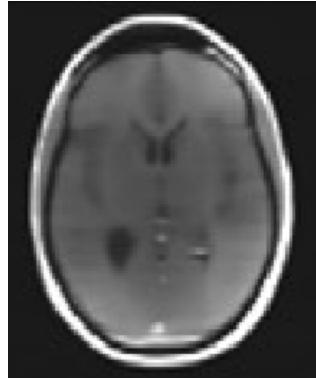
Subsampled Reconstruction (100x)



Accurate Reconstruction



Learned Reconstruction (100x)



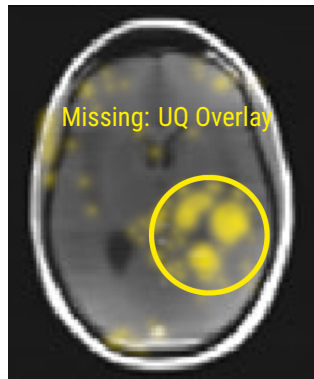
Importance of Uncertainty Quantification

Crucial information to benefit from the 100x acceleration is missing!

Accurate Reconstruction



Learned Reconstruction (100x)



Uncertainty quantification is essential to make critical decisions.

Gaussian Process Regression

Supervised learning of an unknown function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ with uncertainty quantification.

Gaussian Process Regression

Learning an unknown function from data.

Goal: Supervised learning from n data points (X, y)

Prior: Gaussian process $f \sim \mathcal{GP}(\mu, k)$

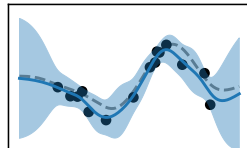
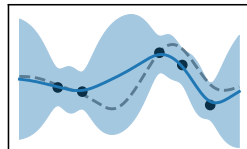
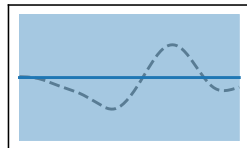
Likelihood: Observations $y = f(X) + \varepsilon \sim \mathcal{N}(f(X), \sigma^2 I)$

Posterior: $f | X, y \sim \mathcal{GP}(\mu_*, k_*)$ with

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, X) \hat{K}^{-1} (y - \mu(X))$$

$$k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, X) \hat{K}^{-1} k(X, \cdot)$$

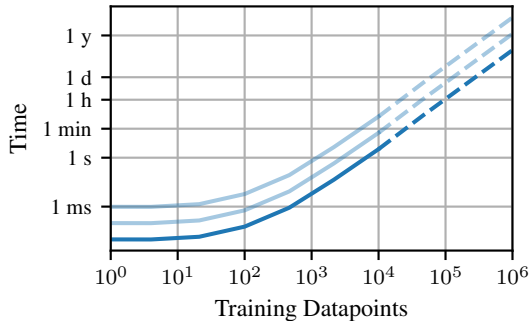
where $\hat{K} = K + \sigma^2 I \in \mathbb{R}^{n \times n}$.



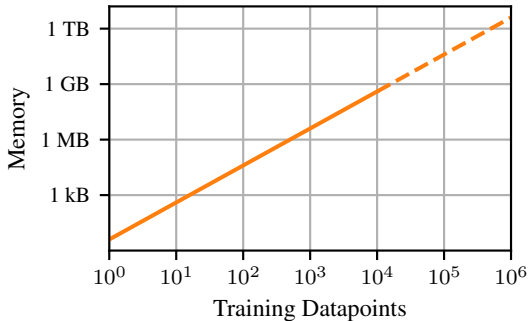
Computational Cost of Gaussian Processes

Uncertainty quantification can be expensive.

Time: $\mathcal{O}(n^3)$



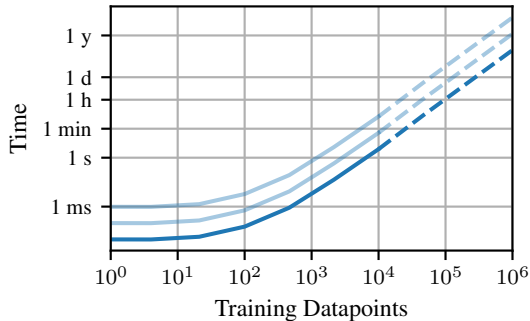
Space: $\mathcal{O}(n^2)$



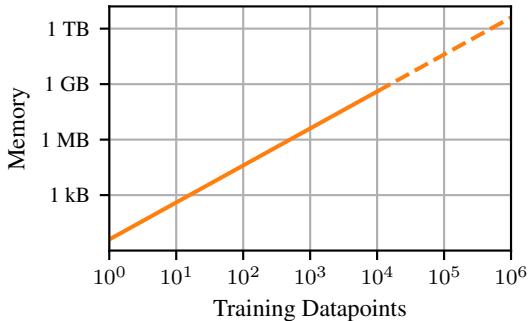
Computational Cost of Gaussian Processes

Uncertainty quantification can be expensive.

Time: $\mathcal{O}(n^3)$



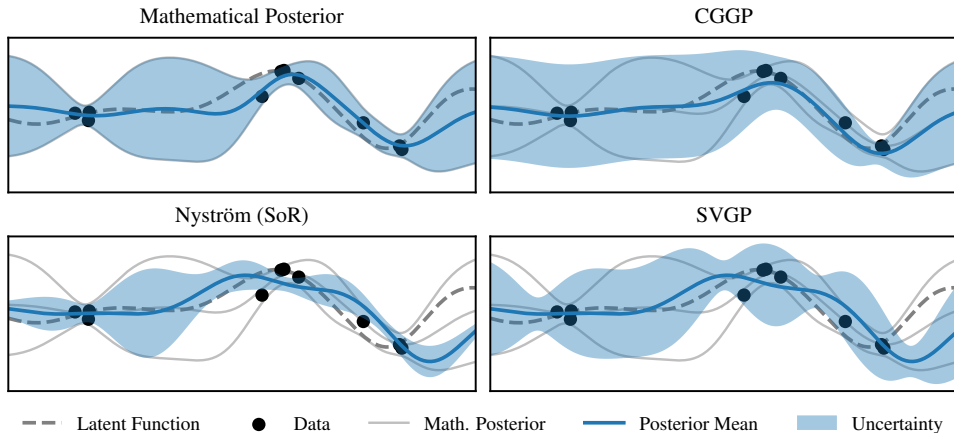
Space: $\mathcal{O}(n^2)$



We need to **approximate** the posterior.

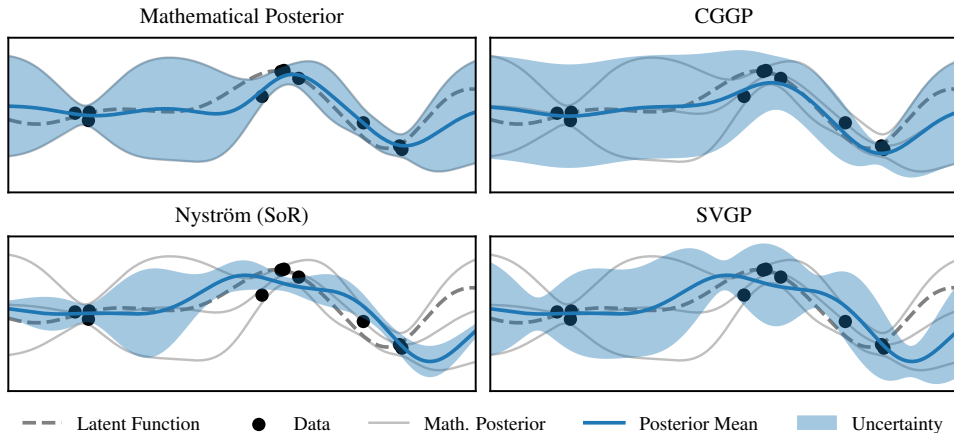
Approximate Gaussian Process Inference

Impact of approximations on uncertainty quantification and decision-making.



Approximate Gaussian Process Inference

Impact of approximations on uncertainty quantification and decision-making.



Approximations introduce **error**, which impacts downstream decisions.

Question 1:

How can we perform Gaussian process inference at scale?

Question 1:

How can we perform **Gaussian process inference at scale?**

Question 2:

How can we **quantify** the inevitable **approximation error**?

Q1: Gaussian Process Inference at Scale?

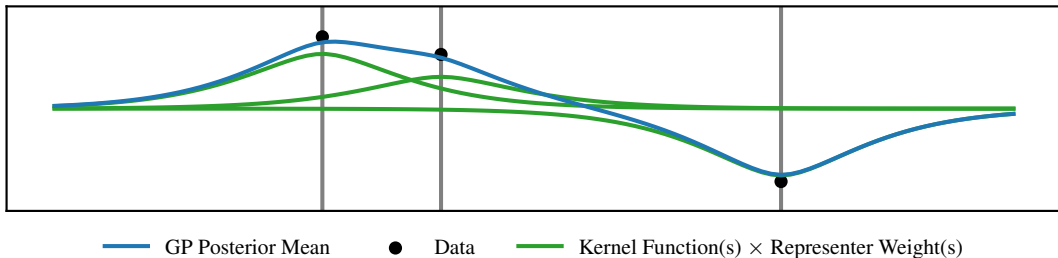
Efficiently approximating the posterior of a Gaussian process.

Representer Weights

The posterior mean is a linear combination of kernel functions centered at datapoints.

$$f \mid X, y \sim \mathcal{GP}(\mu_*, k_*)$$

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, X) \underbrace{\hat{K}^{-1}(y - \mu(X))}_{\text{representer weights } v_*} = \mu(\cdot) + \sum_{j=1}^n k(\cdot, x_j)(v_*)_j$$



Interlude: Method of Conjugate Gradients

Efficiently solving linear systems with positive definite system matrix via matrix-vector multiplies.

Goal: Approximately solve linear system $Ax = b$, where A symmetric positive definite.

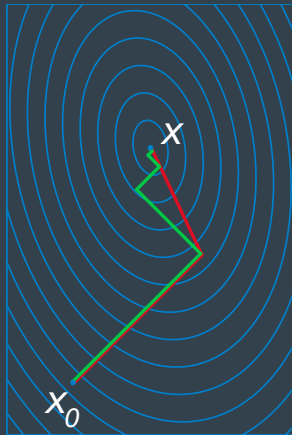
Idea: Rephrase as quadratic optimization problem and optimize. Let

$$f(x) = \frac{1}{2}x^T Ax - b^T x$$

then $\nabla f(x) = 0 \iff Ax = b \iff r(x) := b - Ax = 0$.

Question: How should we optimize?

- 1 **Gradient descent:** Follow $d_i = r(x_i) = -\nabla f(x_i)$ s.t. $\langle d_i, d_j \rangle = 0$.
- 2 **Conjugate direction method:** Follow d_i s. t. $\langle d_i^T d_j \rangle_A = d_i^T A d_j = 0$ for $i \neq j$.
 \Rightarrow convergence in at most n steps.
- 3 **Conjugate gradient method:** First step $d_0 = r(x_0)$.



Oleg Alexandrov, commons.wikimedia.org/w/index.php?curid=2267598

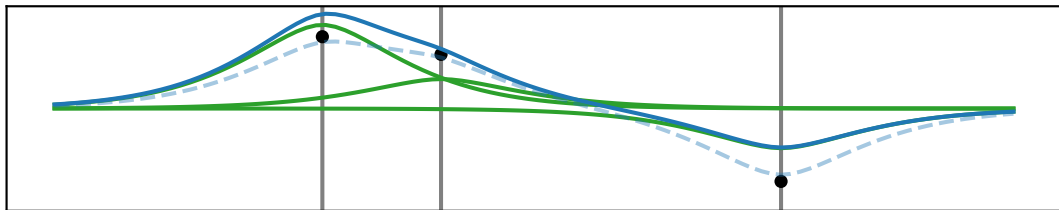
Approximating Representer Weights

Iterative linear solvers can approximate the representer weights.

(Gardner et al., 2018; Charlier et al., 2021)

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, X) \underbrace{\hat{K}^{-1}(y - \mu(X))}_{\text{representer weights } \mathbf{v}_*} \approx \mu(\cdot) + k(\cdot, X) \mathbf{v}_i$$

Observation: Can use iterative linear solvers (e.g. CG) to approximate the representer weights $\mathbf{v}_* \approx \mathbf{v}_i$.



— Approx. GP Posterior Mean ● Data — Kernel Function(s) \times Approx. Representer Weight(s)

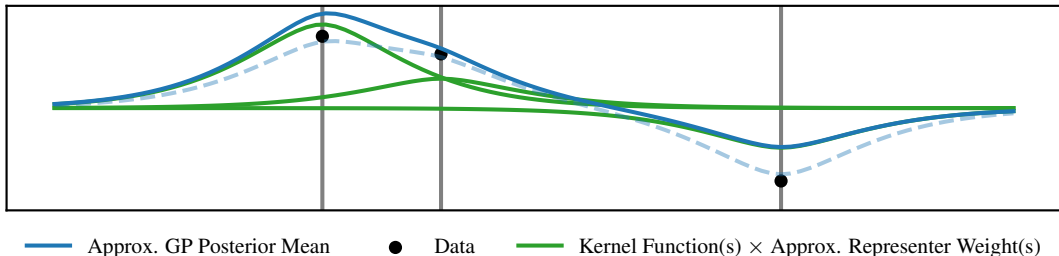
Approximating Representer Weights

Iterative linear solvers can approximate the representer weights.

(Gardner et al., 2018; Charlier et al., 2021)

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, X) \underbrace{\hat{K}^{-1}(y - \mu(X))}_{\text{representer weights } \mathbf{v}_*} \approx \mu(\cdot) + k(\cdot, X) \mathbf{v}_i$$

Observation: Can use iterative linear solvers (e.g. CG) to approximate the representer weights $\mathbf{v}_* \approx \mathbf{v}_i$.



Benefit: Time complexity $\mathcal{O}(n^2)$ and space complexity $\mathcal{O}(nd)$.

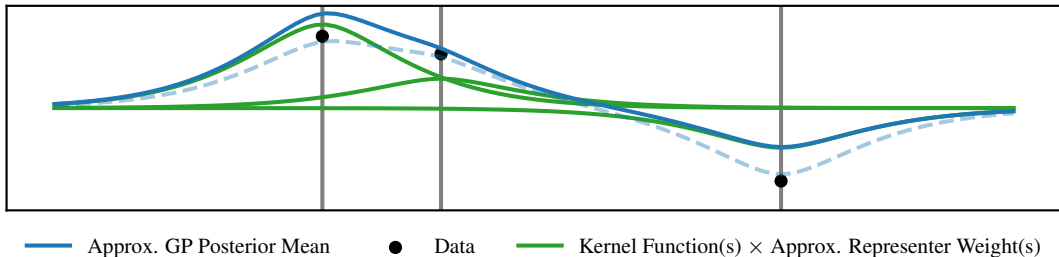
Approximating Representer Weights

Iterative linear solvers can approximate the representer weights.

(Gardner et al., 2018; Charlier et al., 2021)

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, X) \underbrace{\hat{K}^{-1}(y - \mu(X))}_{\text{representer weights } \mathbf{v}_*} \approx \mu(\cdot) + k(\cdot, X) \mathbf{v}_i$$

Observation: Can use iterative linear solvers (e.g. CG) to approximate the representer weights $\mathbf{v}_* \approx \mathbf{v}_i$.

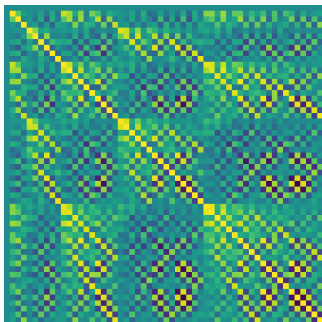


Question: Can we quantify the impact of this approximation on the posterior?

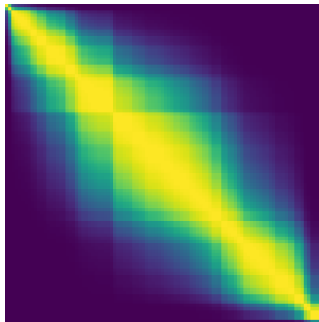
Q2: Can We Quantify Approximation Error?

Probabilistic error quantification at prediction time using probabilistic linear solvers.

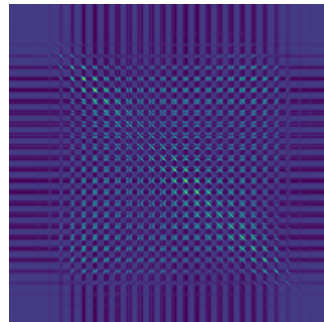
Problem: Solve linear system(s) $Ax_* = b$ for $x_* \in \mathbb{R}^n$.



(a) Gram matrix $X^T X$



(b) Kernel matrix $\hat{K} = k(X, X) + \sigma^2 I$

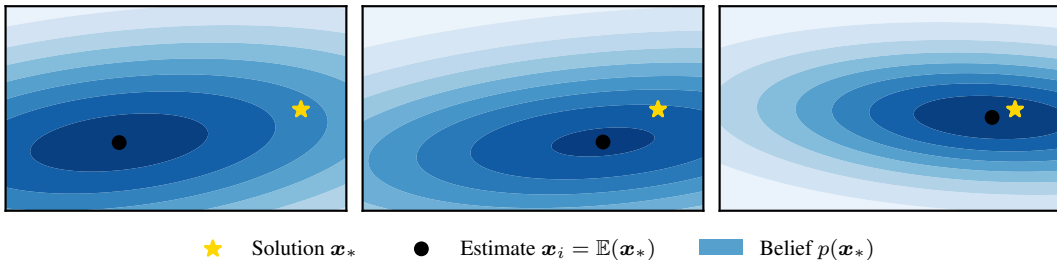


(c) Hessian matrix $\nabla^2 \ell(y, f(X))$

Linear systems in ML are **large-scale**, have **model-induced structure** and are often **solved repeatedly**.

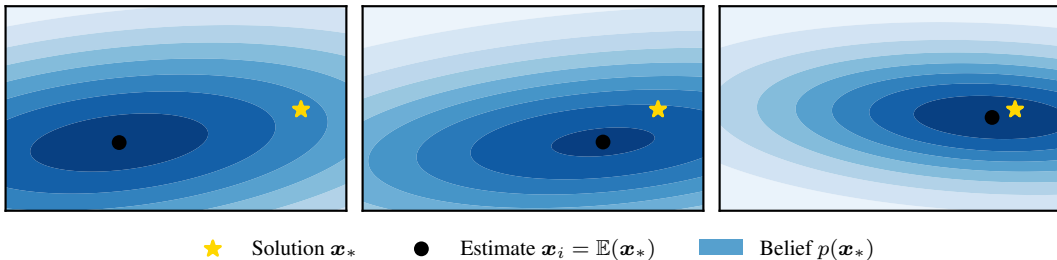
Core Insights of Probabilistic Numerics

- The solution to any numerical problem is fundamentally **uncertain**.



Core Insights of Probabilistic Numerics

- ▶ The solution to any numerical problem is fundamentally **uncertain**.
- ▶ Numerical algorithms are **learning agents**, which actively collect data and make predictions.





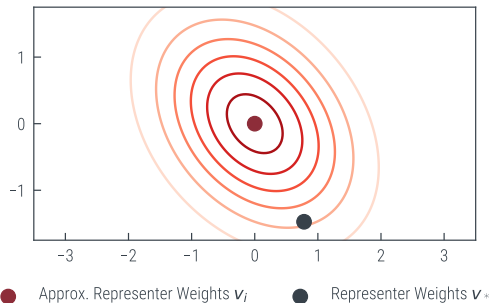
Learning Representer Weights

Estimating representer weights with a probabilistic linear solver.

(Wenger et al., 2022a)

Goal: Solve $\hat{K} \mathbf{v}_* = \mathbf{y} - \boldsymbol{\mu}$ approximately.

Prior: $\mathbf{v}_* \sim \mathcal{N}(\mathbf{v}_0, \boldsymbol{\Sigma}_0)$





Learning Representer Weights

Estimating representer weights with a probabilistic linear solver.

(Wenger et al., 2022a)

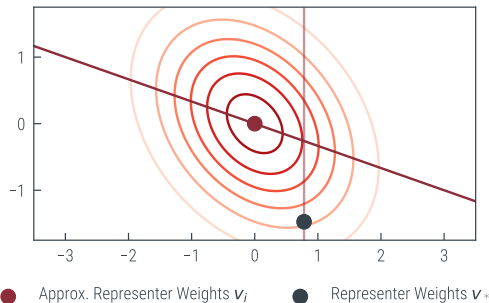
Goal: Solve $\hat{K} \mathbf{v}_* = \mathbf{y} - \boldsymbol{\mu}$ approximately.

Prior: $\mathbf{v}_* \sim \mathcal{N}(\mathbf{v}_0, \boldsymbol{\Sigma}_0)$

Likelihood: Observe representer weights via arbitrarily chosen actions $\mathbf{s}_i \in \mathbb{R}^n$:

$$\begin{aligned} \alpha_i &:= \mathbf{s}_i^\top \mathbf{r}_{i-1} = \mathbf{s}_i^\top ((\mathbf{y} - \boldsymbol{\mu}) - \hat{K} \mathbf{v}_{i-1}) \\ &= \mathbf{s}_i^\top \hat{K} (\mathbf{v}_* - \mathbf{v}_{i-1}) \end{aligned}$$

$$p(\alpha_i | \mathbf{v}_*) = \lim_{\varepsilon \rightarrow 0} \mathcal{N}(\alpha_i; 0, \varepsilon)$$



Learning Representer Weights

Estimating representer weights with a probabilistic linear solver.

(Wenger et al., 2022a)

Goal: Solve $\hat{K}\mathbf{v}_* = \mathbf{y} - \boldsymbol{\mu}$ approximately.

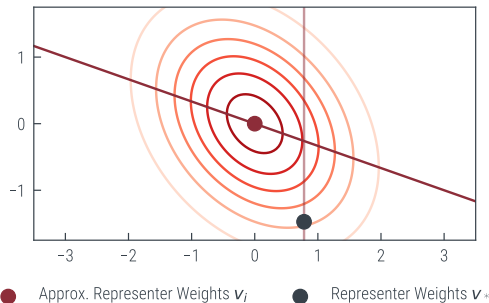
Prior: $\mathbf{v}_* \sim \mathcal{N}(\mathbf{v}_0, \boldsymbol{\Sigma}_0)$

Likelihood: Observe representer weights via arbitrarily chosen actions $\mathbf{s}_i \in \mathbb{R}^n$:

$$\begin{aligned}
 \alpha_i &:= \mathbf{s}_i^\top \mathbf{r}_{i-1} = \mathbf{s}_i^\top ((\mathbf{y} - \boldsymbol{\mu}) - \hat{K}\mathbf{v}_{i-1}) \\
 &= \mathbf{s}_i^\top \hat{K}(\mathbf{v}_* - \mathbf{v}_{i-1})
 \end{aligned}$$

$$p(\alpha_i \mid \mathbf{v}_*) = \lim_{\varepsilon \rightarrow 0} \mathcal{N}(\alpha_i; 0, \varepsilon)$$

Posterior: Affine Gaussian inference!





Learning Representer Weights

Estimating representer weights with a probabilistic linear solver.

(Wenger et al., 2022a)

Goal: Solve $\hat{K} \mathbf{v}_* = \mathbf{y} - \boldsymbol{\mu}$ approximately.

Prior: $\mathbf{v}_* \sim \mathcal{N}(\mathbf{v}_0, \boldsymbol{\Sigma}_0)$

Likelihood: Observe representer weights via arbitrarily chosen actions $\mathbf{s}_i \in \mathbb{R}^n$:

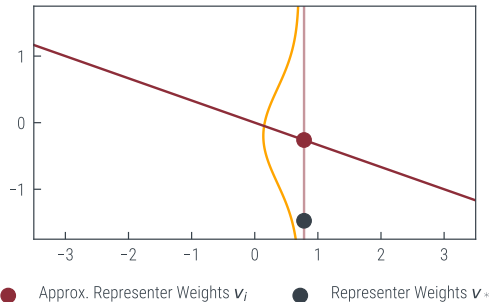
$$\begin{aligned} \alpha_i &:= \mathbf{s}_i^\top \mathbf{r}_{i-1} = \mathbf{s}_i^\top ((\mathbf{y} - \boldsymbol{\mu}) - \hat{K} \mathbf{v}_{i-1}) \\ &= \mathbf{s}_i^\top \hat{K} (\mathbf{v}_* - \mathbf{v}_{i-1}) \end{aligned}$$

$$p(\alpha_i | \mathbf{v}_*) = \lim_{\varepsilon \rightarrow 0} \mathcal{N}(\alpha_i; 0, \varepsilon)$$

Posterior: $\mathbf{v}_* | \alpha_i \sim \mathcal{N}(\mathbf{v}_i, \boldsymbol{\Sigma}_i)$, where

$$\mathbf{v}_i = \mathbf{v}_{i-1} + \boldsymbol{\Sigma}_{i-1} \hat{K} \mathbf{s}_i (\mathbf{s}_i^\top \hat{K} \boldsymbol{\Sigma}_{i-1} \hat{K} \mathbf{s}_i)^{-1} \mathbf{s}_i^\top \hat{K} (\mathbf{v}_* - \mathbf{v}_{i-1})$$

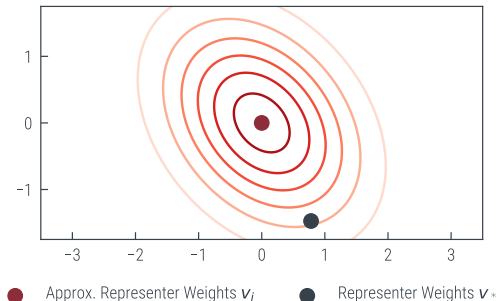
$$\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}_{i-1} - \boldsymbol{\Sigma}_{i-1} \hat{K} \mathbf{s}_i (\mathbf{s}_i^\top \hat{K} \boldsymbol{\Sigma}_{i-1} \hat{K} \mathbf{s}_i)^{-1} \mathbf{s}_i^\top \hat{K} \boldsymbol{\Sigma}_{i-1}$$



Choosing the Linear Solver Prior

The Gaussian process prior makes assumptions about the representer weights.

Question: How to choose the linear solver prior?



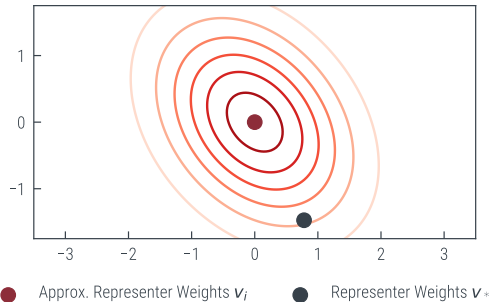
Choosing the Linear Solver Prior

The Gaussian process prior makes assumptions about the representer weights.

Question: How to choose the linear solver prior?

Remember: $y = f(X) + \varepsilon$, where $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$.

$$\Rightarrow y - \mu \sim \mathcal{N}(\mathbf{0}, k(X, X) + \sigma^2 I) = \mathcal{N}(\mathbf{0}, \hat{K})$$



Choosing the Linear Solver Prior

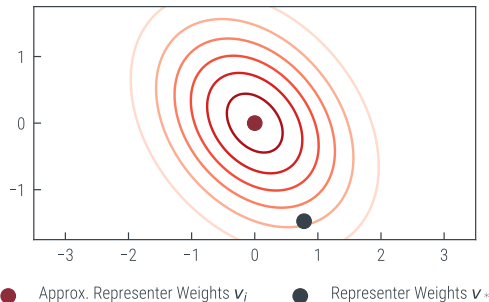
The Gaussian process prior makes assumptions about the representer weights.

Question: How to choose the linear solver prior?

Remember: $y = f(X) + \varepsilon$, where $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$.

$$\Rightarrow y - \mu \sim \mathcal{N}(\mathbf{0}, k(X, X) + \sigma^2 I) = \mathcal{N}(\mathbf{0}, \hat{K})$$

$$\Rightarrow v_* = \hat{K}^{-1}(y - \mu) \sim \mathcal{N}\left(\underbrace{\mathbf{0}}_{=v_0}, \underbrace{\hat{K}^{-1}}_{=\Sigma_0}\right)$$



Choosing the Linear Solver Prior

The Gaussian process prior makes assumptions about the representer weights.

Question: How to choose the linear solver prior?

$$\Rightarrow y - \mu \sim \mathcal{N}(\mathbf{0}, k(X, X) + \sigma^2 I) = \mathcal{N}(\mathbf{0}, \hat{K})$$

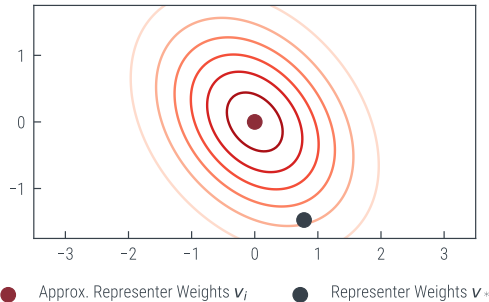
$$\Rightarrow \mathbf{v}_* = \hat{K}^{-1}(y - \mu) \sim \mathcal{N}\left(\underset{=\mathbf{v}_0}{\mathbf{0}}, \underset{=\Sigma_0}{\hat{K}^{-1}}\right)$$

Setting $\mathbf{v}_0 = \mathbf{0}$ and $\Sigma_0 = \hat{K}^{-1}$, we have

$$\mathbf{v}_i = \mathbf{S}_i(\mathbf{S}_i^\top \hat{K} \mathbf{S}_i)^{-1} \mathbf{S}_i^\top (y - \mu) = \mathbf{C}_i(y - \mu)$$

$$\Sigma_i = \Sigma_0 - \mathbf{S}_i(\mathbf{S}_i^\top \hat{K} \mathbf{S}_i)^{-1} \mathbf{S}_i^\top = \Sigma_0 - \mathbf{C}_i$$

where \mathbf{S}_i is the matrix of actions $\mathbf{s}_1, \dots, \mathbf{s}_i$.



Choosing the Linear Solver Prior

The Gaussian process prior makes assumptions about the representer weights.

Question: How to choose the linear solver prior?

$$\Rightarrow y - \mu \sim \mathcal{N}(\mathbf{0}, k(X, X) + \sigma^2 I) = \mathcal{N}(\mathbf{0}, \hat{K})$$

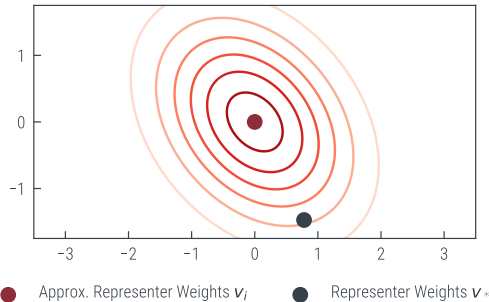
$$\Rightarrow \mathbf{v}_* = \hat{K}^{-1}(y - \mu) \sim \mathcal{N}\left(\underset{=\mathbf{v}_0}{\mathbf{0}}, \underset{=\Sigma_0}{\hat{K}^{-1}}\right)$$

Setting $\mathbf{v}_0 = \mathbf{0}$ and $\Sigma_0 = \hat{K}^{-1}$, we have

$$\mathbf{v}_i = \mathbf{S}_i(\mathbf{S}_i^\top \hat{K} \mathbf{S}_i)^{-1} \mathbf{S}_i^\top (y - \mu) = \mathbf{C}_i(y - \mu)$$

$$\Sigma_i = \Sigma_0 - \mathbf{S}_i(\mathbf{S}_i^\top \hat{K} \mathbf{S}_i)^{-1} \mathbf{S}_i^\top = \Sigma_0 - \mathbf{C}_i$$

where \mathbf{S}_i is the matrix of actions $\mathbf{s}_1, \dots, \mathbf{s}_i$.



Chicken & Egg Problem: How can we get a probabilistic error estimate for $\mathbf{v}_i \approx \mathbf{v}_*$, if we need \hat{K}^{-1} ?

IterGP: Computation-Aware Gaussian Process Inference

Quantifying uncertainty arising from observing finite data and performing a finite amount of computation.

(Wenger et al., 2022a)

Goal: Approximate the Gaussian process posterior $f \mid \mathbf{y} \sim \mathcal{GP}(\mu_*, k_*)$, where

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X})\hat{K}^{-1}(\mathbf{y} - \boldsymbol{\mu}), \quad \text{and} \quad k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X})\hat{K}^{-1}k(\mathbf{X}, \cdot)$$

IterGP: Computation-Aware Gaussian Process Inference

Quantifying uncertainty arising from observing finite data and performing a finite amount of computation.

(Wenger et al., 2022a)

Goal: Approximate the Gaussian process posterior $f \mid \mathbf{y} \sim \mathcal{GP}(\mu_*, k_*)$, where

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, X)\hat{K}^{-1}(\mathbf{y} - \boldsymbol{\mu}), \quad \text{and} \quad k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, X)\hat{K}^{-1}k(X, \cdot)$$

Obtained: Belief about representer weights $\mathbf{v}_* = \hat{K}^{-1}(\mathbf{y} - \boldsymbol{\mu}) \sim \mathcal{N}(\mathbf{v}_i, \boldsymbol{\Sigma}_i) = \mathcal{N}(\mathbf{v}_i, \hat{K}^{-1} - \mathbf{C}_i)$

IterGP: Computation-Aware Gaussian Process Inference

Quantifying uncertainty arising from observing finite data and performing a finite amount of computation.

(Wenger et al., 2022a)

Goal: Approximate the Gaussian process posterior $f \mid \mathbf{y} \sim \mathcal{GP}(\mu_*, k_*)$, where

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, X)\hat{K}^{-1}(\mathbf{y} - \boldsymbol{\mu}), \quad \text{and} \quad k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, X)\hat{K}^{-1}k(X, \cdot)$$

Obtained: Belief about representer weights $\mathbf{v}_* = \hat{K}^{-1}(\mathbf{y} - \boldsymbol{\mu}) \sim \mathcal{N}(\mathbf{v}_i, \boldsymbol{\Sigma}_i) = \mathcal{N}(\mathbf{v}_i, \hat{K}^{-1} - \mathbf{C}_i)$

Idea: Propagate uncertainty about representer weights to posterior.

Goal: Approximate the Gaussian process posterior $f \mid \mathbf{y} \sim \mathcal{GP}(\mu_*, k_*)$, where

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu}), \quad \text{and} \quad k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}k(\mathbf{X}, \cdot)$$

Obtained: Belief about representer weights $\mathbf{v}_* = \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu}) \sim \mathcal{N}(\mathbf{v}_i, \boldsymbol{\Sigma}_i) = \mathcal{N}(\mathbf{v}_i, \hat{\mathbf{K}}^{-1} - \mathbf{C}_i)$

Idea: Propagate uncertainty about representer weights to posterior.

1 Pathwise form of posterior: $(f \mid \mathbf{y})(\cdot) \stackrel{d}{=} f(\cdot) + k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$

Goal: Approximate the Gaussian process posterior $f | \mathbf{y} \sim \mathcal{GP}(\mu_*, k_*)$, where

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu}), \quad \text{and} \quad k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}k(\mathbf{X}, \cdot)$$

Obtained: Belief about representer weights $\mathbf{v}_* = \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu}) \sim \mathcal{N}(\mathbf{v}_i, \boldsymbol{\Sigma}_i) = \mathcal{N}(\mathbf{v}_i, \hat{\mathbf{K}}^{-1} - \mathbf{C}_i)$

Idea: Propagate uncertainty about representer weights to posterior.

1 Pathwise form of posterior: $(f | \mathbf{y})(\cdot) \stackrel{d}{=} f(\cdot) + k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$

2 Reparametrize posterior: $(f | \mathbf{v}_*)(\cdot) \stackrel{d}{=} f(\cdot) + k(\cdot, \mathbf{X})\mathbf{v}_*$

Goal: Approximate the Gaussian process posterior $f \mid \mathbf{y} \sim \mathcal{GP}(\mu_*, k_*)$, where

$$\mu_*(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu}), \quad \text{and} \quad k_*(\cdot, \cdot) = k(\cdot, \cdot) - k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}k(\mathbf{X}, \cdot)$$

Obtained: Belief about representer weights $\mathbf{v}_* = \hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu}) \sim \mathcal{N}(\mathbf{v}_i, \boldsymbol{\Sigma}_i) = \mathcal{N}(\mathbf{v}_i, \hat{\mathbf{K}}^{-1} - \mathbf{C}_i)$

Idea: Propagate uncertainty about representer weights to posterior.

- 1 Pathwise form of posterior: $(f \mid \mathbf{y})(\cdot) \stackrel{d}{=} f(\cdot) + k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})$
- 2 Reparametrize posterior: $(f \mid \mathbf{v}_*)(\cdot) \stackrel{d}{=} f(\cdot) + k(\cdot, \mathbf{X})\mathbf{v}_*$
- 3 Marginalize representer weights belief: $p(f(\cdot)) = \int p(f(\cdot) \mid \mathbf{v}_*)p(\mathbf{v}_*)d\mathbf{v}_* = \mathcal{GP}(f; \mu_i, k_i)$,

$$\mu_i(\cdot) = \mu(\cdot) + k(\cdot, \mathbf{X})\mathbf{v}_i$$

$$k_i(\cdot, \cdot) = \underbrace{k(\cdot, \cdot) - k(\cdot, \mathbf{X})\hat{\mathbf{K}}^{-1}k(\mathbf{X}, \cdot)}_{\text{mathematical uncertainty } \text{blue circle}} + \underbrace{k(\cdot, \mathbf{X})\boldsymbol{\Sigma}_i k(\mathbf{X}, \cdot)}_{\text{computational uncertainty } \text{green circle}} = \underbrace{k(\cdot, \cdot) - k(\cdot, \mathbf{X})\mathbf{C}_i k(\mathbf{X}, \cdot)}_{\text{combined uncertainty } \text{purple circle}}$$

Probabilistic Quantification of Approximation Error

The covariance can be interpreted as a squared error.

Combined Uncertainty

Belief about the latent function is captured by $f \sim \mathcal{GP}(\mu_i, k_i)$, s.t.

$$\begin{aligned}\mu_i(\cdot) &= \mu(\cdot) + k(\cdot, X)\mathbf{v}_i \\ k_i(\cdot, \cdot) &= \underbrace{k(\cdot, \cdot) - k(\cdot, X)\hat{K}^{-1}k(X, \cdot)}_{\text{mathematical uncertainty} \quad \text{blue circle}} + \underbrace{k(\cdot, X)\Sigma_i k(X, \cdot)}_{\text{computational uncertainty} \quad \text{green circle}} = \underbrace{k(\cdot, \cdot) - k(\cdot, X)\mathbf{C}_i k(X, \cdot)}_{\text{combined uncertainty} \quad \text{purple circle}}\end{aligned}$$

Remember: $k(\mathbf{x}, \mathbf{x}) = \text{Cov}(f(\mathbf{x}), f(\mathbf{x})) = \mathbb{E}((f(\mathbf{x}) - \mathbb{E}(f(\mathbf{x})))^2)$

Probabilistic Quantification of Approximation Error

The covariance can be interpreted as a squared error.

Combined Uncertainty

Belief about the latent function is captured by $f \sim \mathcal{GP}(\mu_i, k_i)$, s.t.

$$\begin{aligned}\mu_i(\cdot) &= \mu(\cdot) + k(\cdot, X)\mathbf{v}_i \\ k_i(\cdot, \cdot) &= \underbrace{k(\cdot, \cdot) - k(\cdot, X)\hat{K}^{-1}k(X, \cdot)}_{\text{mathematical uncertainty} \text{ (blue circle)}} + \underbrace{k(\cdot, X)\Sigma_i k(X, \cdot)}_{\text{computational uncertainty} \text{ (green circle)}} = \underbrace{k(\cdot, \cdot) - k(\cdot, X)\mathbf{C}_i k(X, \cdot)}_{\text{combined uncertainty} \text{ (purple circle)}}\end{aligned}$$

Remember: $k(\mathbf{x}, \mathbf{x}) = \text{Cov}(f(\mathbf{x}), f(\mathbf{x})) = \mathbb{E}((f(\mathbf{x}) - \mathbb{E}(f(\mathbf{x})))^2)$

$$k_*(\mathbf{x}, \mathbf{x}) = \underbrace{k(\mathbf{x}, \mathbf{x}) - k(\mathbf{x}, X)\hat{K}^{-1}k(X, \mathbf{x})}_{\text{mathematical uncertainty} \text{ (blue circle)}} = \mathbb{E}((f(\mathbf{x}) - \mu_*(\mathbf{x}))^2)$$

Probabilistic Quantification of Approximation Error

The covariance can be interpreted as a squared error.

Combined Uncertainty

Belief about the latent function is captured by $f \sim \mathcal{GP}(\mu_i, k_i)$, s.t.

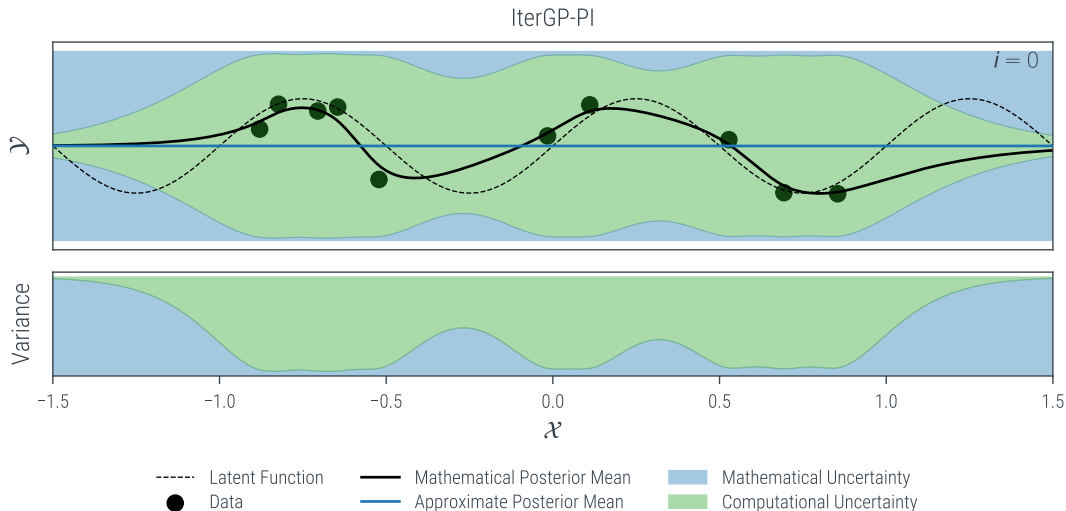
$$\begin{aligned}\mu_i(\cdot) &= \mu(\cdot) + k(\cdot, X)v_i \\ k_i(\cdot, \cdot) &= \underbrace{k(\cdot, \cdot) - k(\cdot, X)\hat{K}^{-1}k(X, \cdot)}_{\text{mathematical uncertainty} \text{ (blue circle)}} + \underbrace{k(\cdot, X)\Sigma_i k(X, \cdot)}_{\text{computational uncertainty} \text{ (green circle)}} = \underbrace{k(\cdot, \cdot) - k(\cdot, X)\mathcal{C}_i k(X, \cdot)}_{\text{combined uncertainty} \text{ (purple circle)}}\end{aligned}$$

Remember: $k(x, x) = \text{Cov}(f(x), f(x)) = \mathbb{E}((f(x) - \mathbb{E}(f(x)))^2)$

$$\begin{aligned}k_*(x, x) &= \underbrace{k(x, x) - k(x, X)\hat{K}^{-1}k(X, x)}_{\text{mathematical uncertainty} \text{ (blue circle)}} = \mathbb{E}((f(x) - \mu_*(x))^2) \\ k_i^{\text{comp}}(x, x) &= \underbrace{k(x, X)\Sigma_i k(X, x)}_{\text{computational uncertainty} \text{ (green circle)}} \stackrel{\text{}}{=} \mathbb{E}((\mu_*(x) - \mu_i(x))^2) \\ &\quad \Sigma_i = \text{Cov}(v_*) = \mathbb{E}((v_* - v_i)(v_* - v_i)^\top)\end{aligned}$$

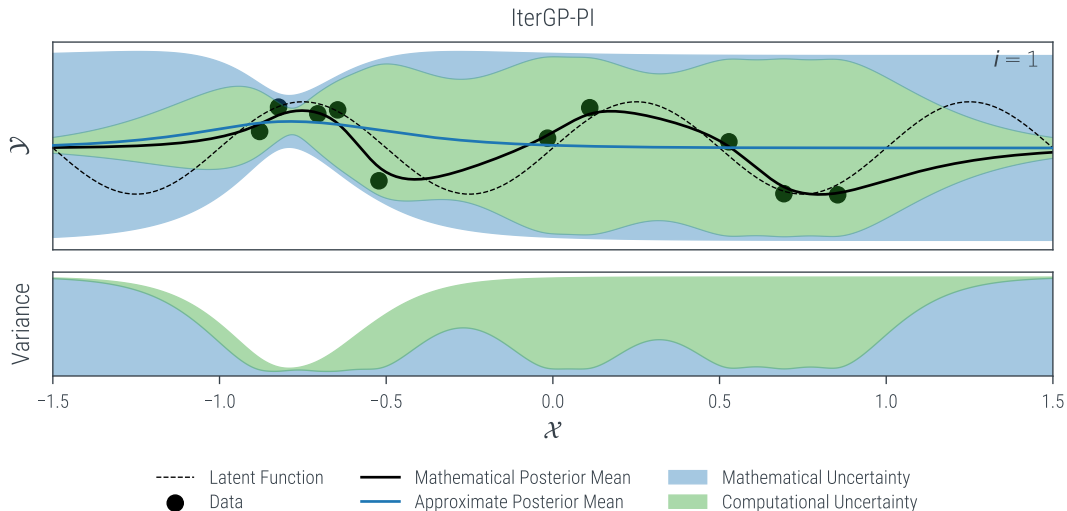
Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.



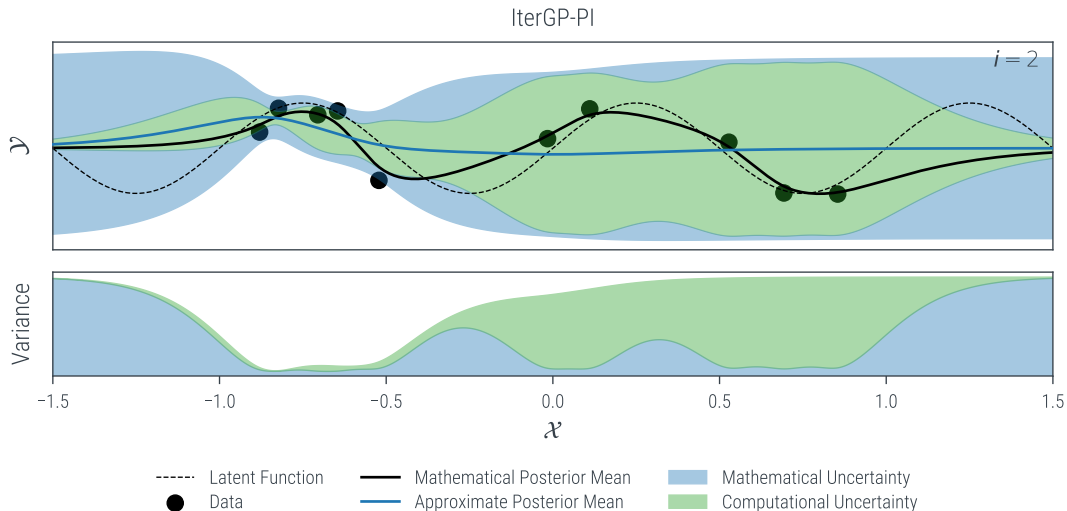
Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.



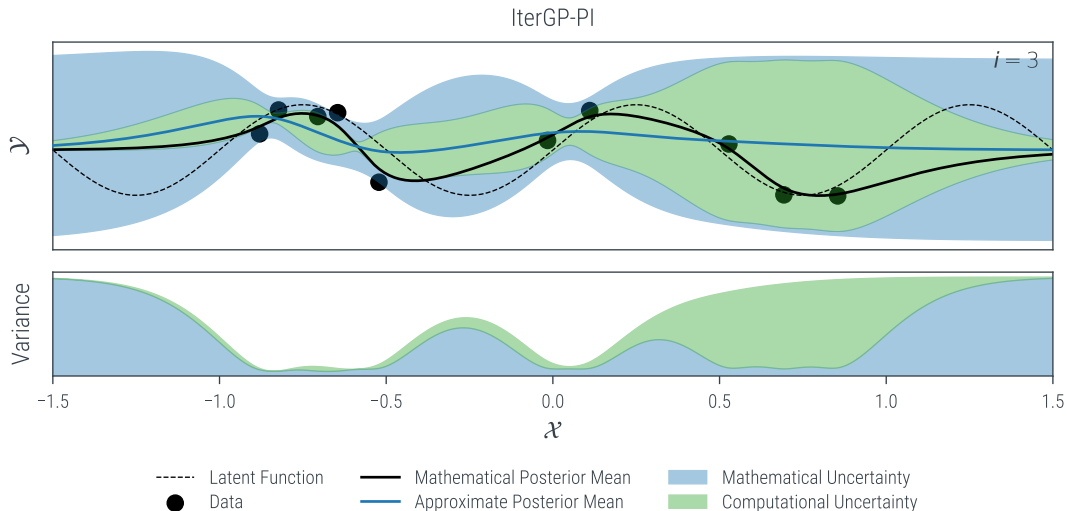
Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.



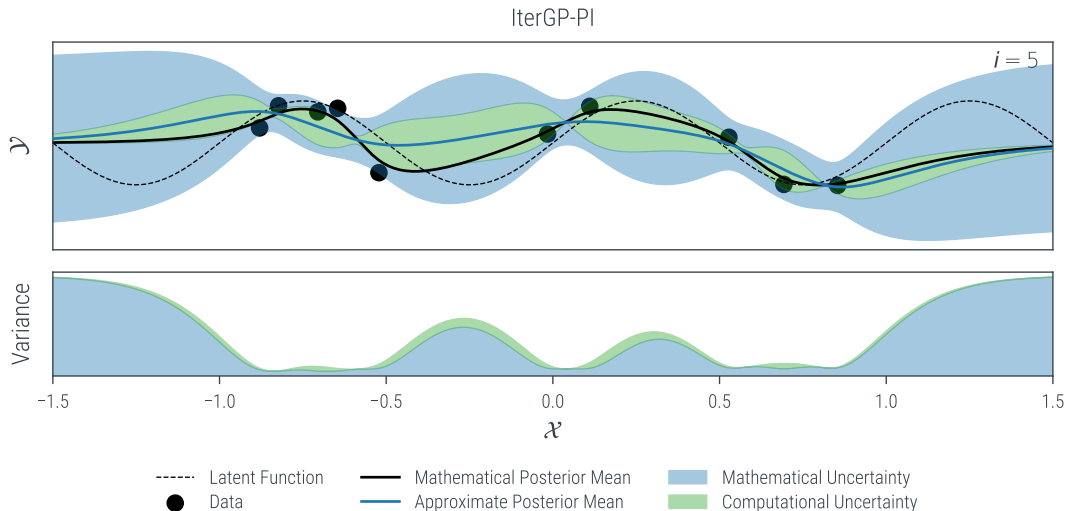
Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.



Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.



Theoretical Analysis

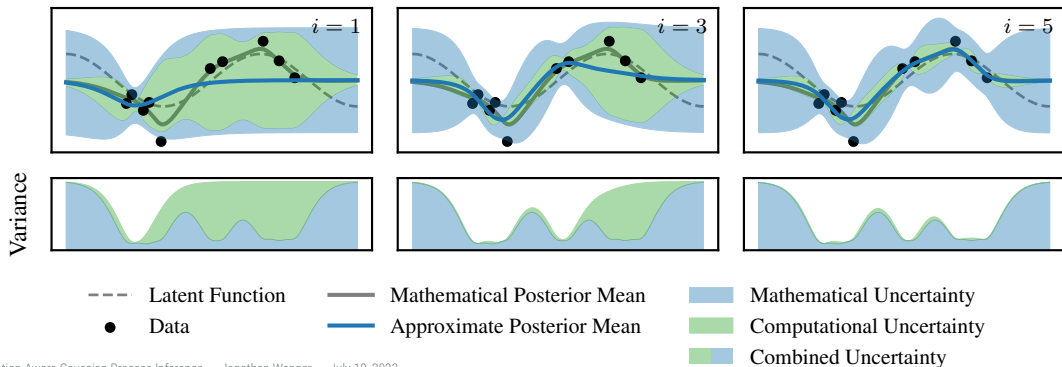
Uncertainty as a bound on the relative predictive error.

Theorem (Relative Error Bound)

$$\sup_{g \in \mathcal{H}_k \sigma : \|g\|_{\mathcal{H}_k \sigma} \leq 1} \underbrace{g(x) - \mu_*^g(x)}_{\text{error of math. post. mean}} = \sup_{g \in \mathcal{H}_k \sigma} \frac{|g(x) - \mu_*^g(x)|}{\|g\|_{\mathcal{H}_k \sigma}} = \sqrt{k_*(x, x) + \sigma^2} \quad (1)$$

Theorem (Relative Error Bound)

$$\sup_{g \in \mathcal{H}_{k\sigma} : \|g\|_{\mathcal{H}_{k\sigma}} \leq 1} \underbrace{g(x) - \mu_*^g(x)}_{\text{error of math. post. mean } \textcolor{blue}{\bullet}} + \underbrace{\mu_*^g(x) - \mu_i^g(x)}_{\text{computational error } \textcolor{green}{\bullet}} = \sqrt{k_i(x, x) + \sigma^2} \quad (1)$$



What Have We Learned?

So Far:

- ▶ Gaussian process inference is prohibitive for large datasets.

So Far:

- ▶ Gaussian process inference is prohibitive for large datasets.
- ▶ Iterative methods can reduce the necessary computations from $\mathcal{O}(n^3)$ to $\mathcal{O}(n^2)$.

So Far:

- ▶ Gaussian process inference is prohibitive for large datasets.
- ▶ Iterative methods can reduce the necessary computations from $\mathcal{O}(n^3)$ to $\mathcal{O}(n^2)$.
- ▶ Using probabilistic numerics we can quantify the error when approximating Gaussian processes.

So Far:

- ▶ Gaussian process inference is prohibitive for large datasets.
- ▶ Iterative methods can reduce the necessary computations from $\mathcal{O}(n^3)$ to $\mathcal{O}(n^2)$.
- ▶ Using probabilistic numerics we can quantify the error when approximating Gaussian processes.
- ▶ **Explicit trade-off** between computation and uncertainty.

What Have We Learned?

So Far:

- ▶ Gaussian process inference is prohibitive for large datasets.
- ▶ Iterative methods can reduce the necessary computations from $\mathcal{O}(n^3)$ to $\mathcal{O}(n^2)$.
- ▶ Using probabilistic numerics we can quantify the error when approximating Gaussian processes.
- ▶ **Explicit trade-off** between computation and uncertainty.

What About:

- ▶ How does IterGP relate to other numerical (approximation) methods, e.g. Cholesky, CGGP, SVGP?

So Far:

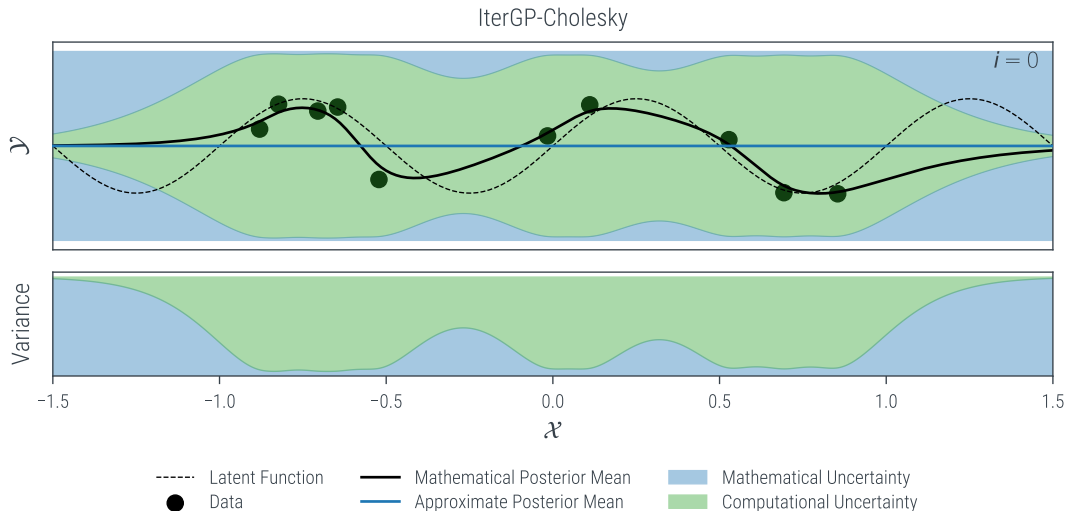
- ▶ Gaussian process inference is prohibitive for large datasets.
- ▶ Iterative methods can reduce the necessary computations from $\mathcal{O}(n^3)$ to $\mathcal{O}(n^2)$.
- ▶ Using probabilistic numerics we can quantify the error when approximating Gaussian processes.
- ▶ **Explicit trade-off** between computation and uncertainty.

What About:

- ▶ How does IterGP relate to other numerical (approximation) methods, e.g. Cholesky, CGGP, SVGP?
- ▶ Is quadratic time $\mathcal{O}(n^2)$ the limit? Can we approximate more cheaply?

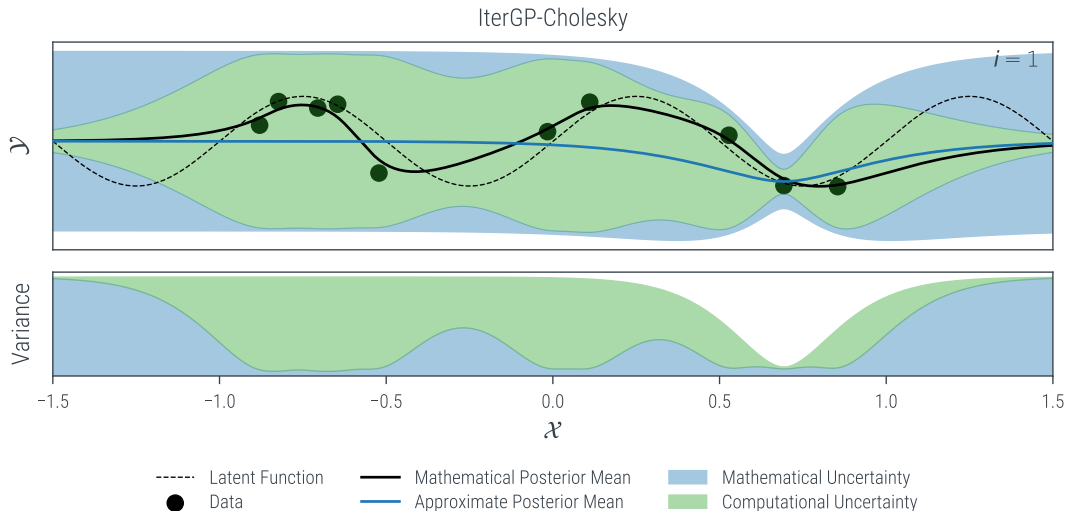
Subset of Data versus IterGP-Cholesky

IterGP with unit vector actions recovers vanilla GP inference.



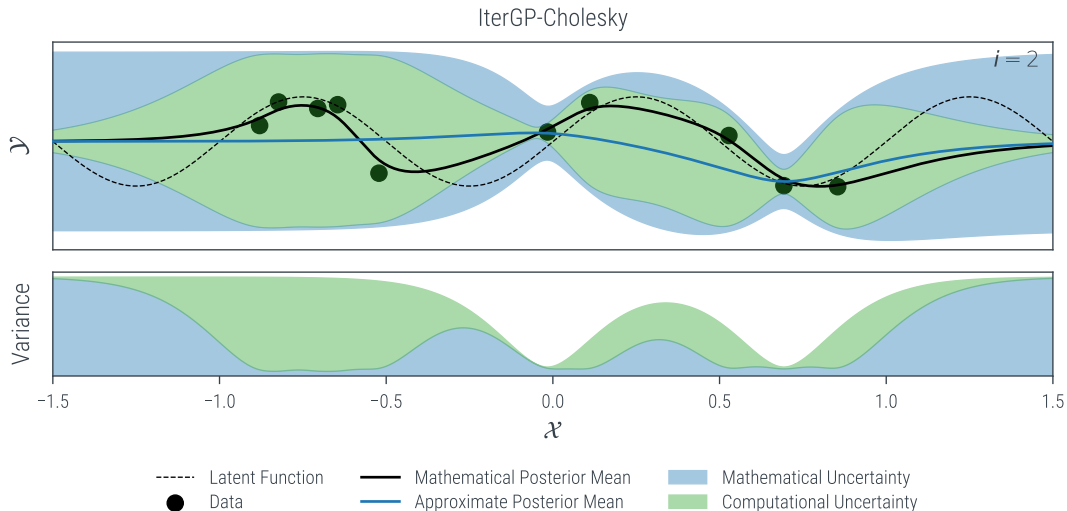
Subset of Data versus IterGP-Cholesky

IterGP with unit vector actions recovers vanilla GP inference.



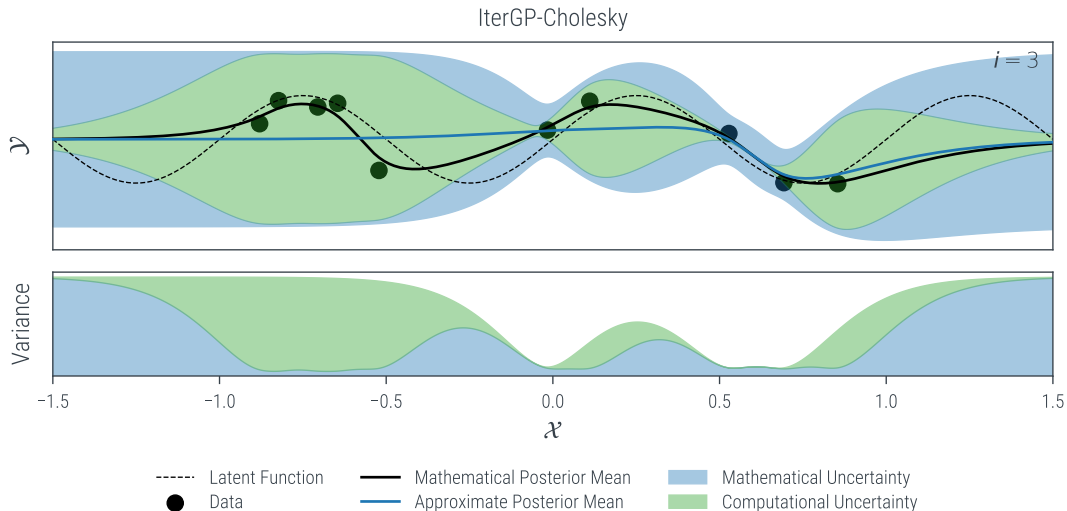
Subset of Data versus IterGP-Cholesky

IterGP with unit vector actions recovers vanilla GP inference.



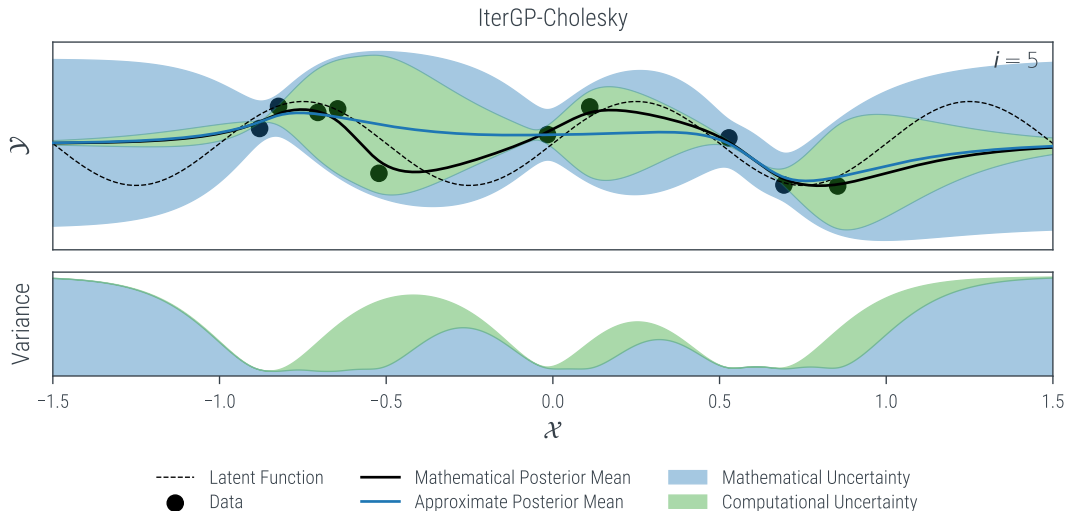
Subset of Data versus IterGP-Cholesky

IterGP with unit vector actions recovers vanilla GP inference.



Subset of Data versus IterGP-Cholesky

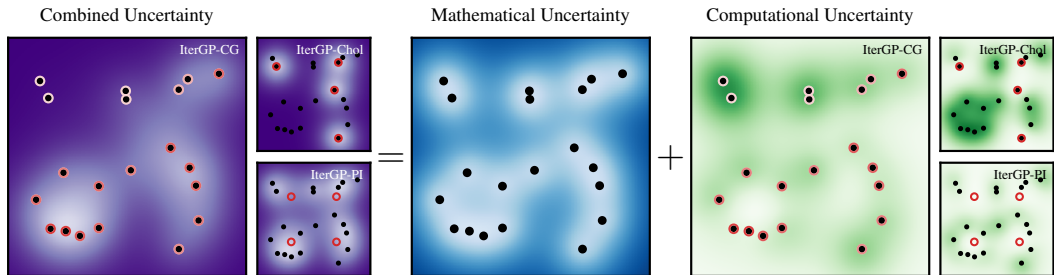
IterGP with unit vector actions recovers vanilla GP inference.



Policy Choice and Connection to Other Approximations

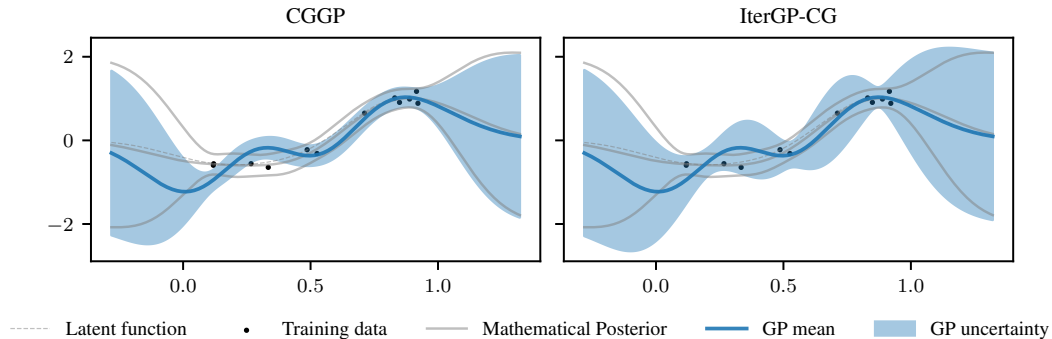
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

Method	Actions \mathbf{s}_i	Classic Analog
IterGP-Cholesky	\mathbf{e}_i	(Partial) Cholesky / subset of data
IterGP-EVD	$\text{ev}_i(\hat{\mathbf{K}})$	(Partial) Eigenvalue decomposition
IterGP-CG	$\mathbf{s}_i^{\text{PCG}}$ or $\hat{\mathbf{P}}^{-1}\mathbf{r}_i$	(Preconditioned) CG
IterGP-Pseudoinput	$k(\mathbf{X}, \mathbf{z}_i)$	\approx SVGP



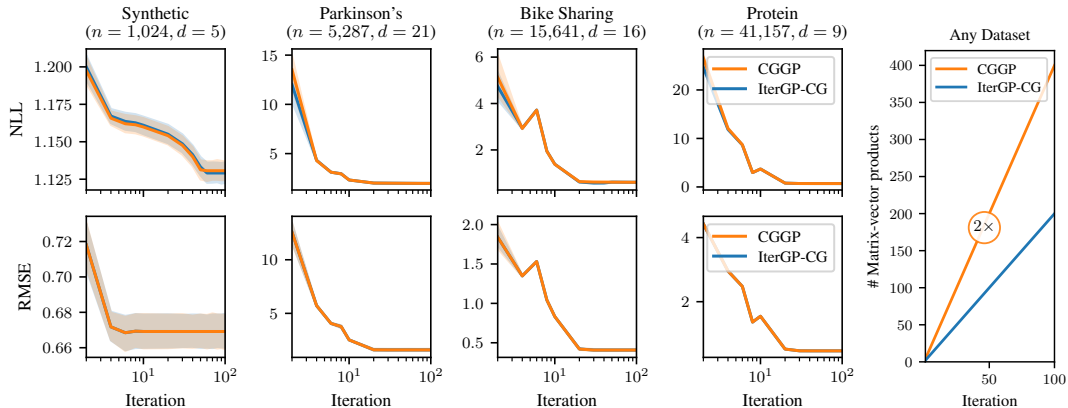
CGGP versus IterGP-CG

IterGP reduces the necessary computations for CG-based GP inference.



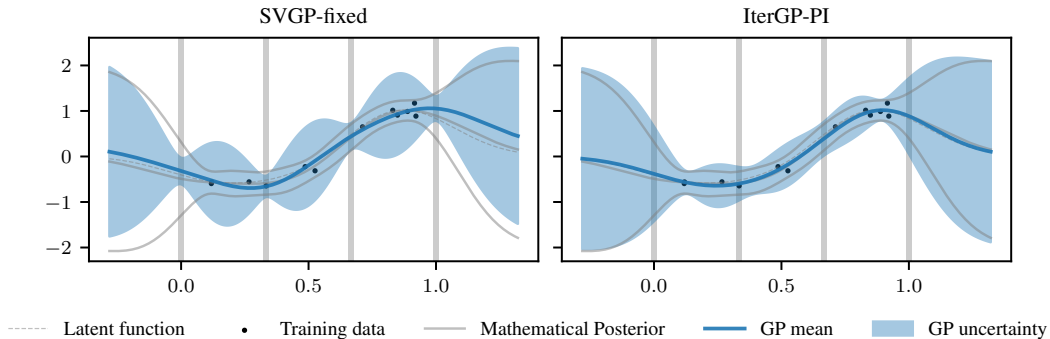
CGGP versus IterGP-CG

IterGP reduces the necessary computations for CG-based GP inference.



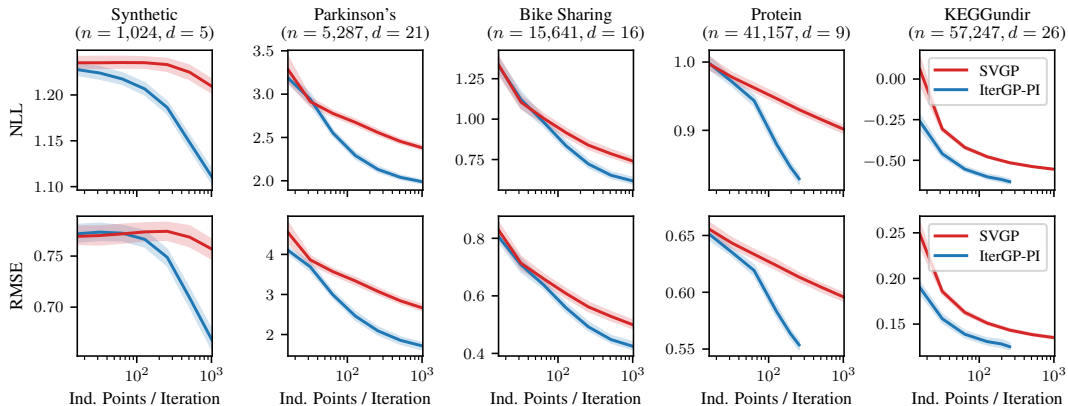
SVGP versus IterGP-PI

Quantifying computational uncertainty improves generalization of inducing point methods like SVGP (Titsias, 2009; Hensman et al., 2013).



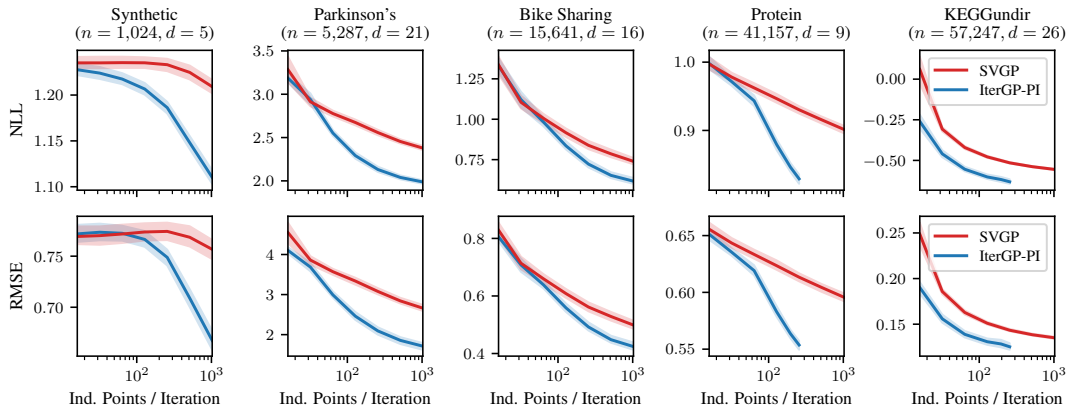
SVGP versus IterGP-PI

Quantifying computational uncertainty improves generalization of inducing point methods like SVGP (Titsias, 2009; Hensman et al., 2013).



SVGP versus IterGP-PI

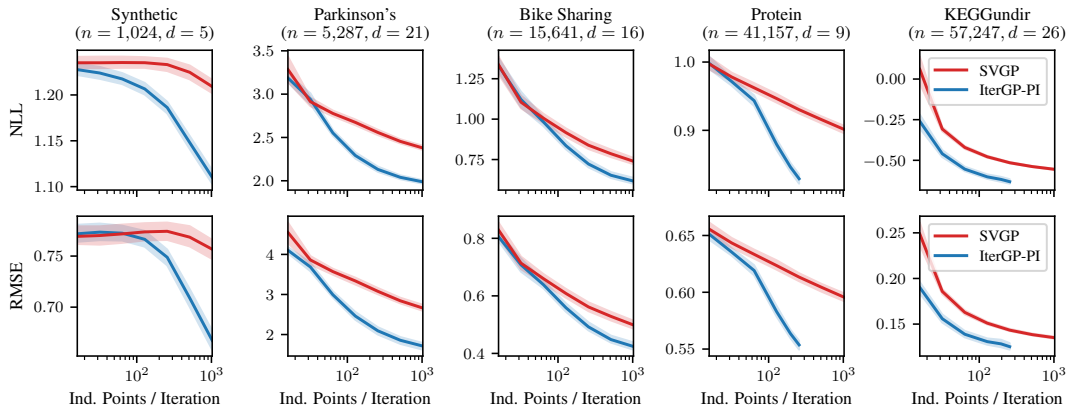
Quantifying computational uncertainty improves generalization of inducing point methods like SVGP (Titsias, 2009; Hensman et al., 2013).



What about optimizing inducing point locations?

SVGP versus IterGP-PI

Quantifying computational uncertainty improves generalization of inducing point methods like SVGP (Titsias, 2009; Hensman et al., 2013).



What about computational cost? SVGP: $\mathcal{O}(nm^2)$ versus IterGP-PI: $\mathcal{O}(n^2m)$.

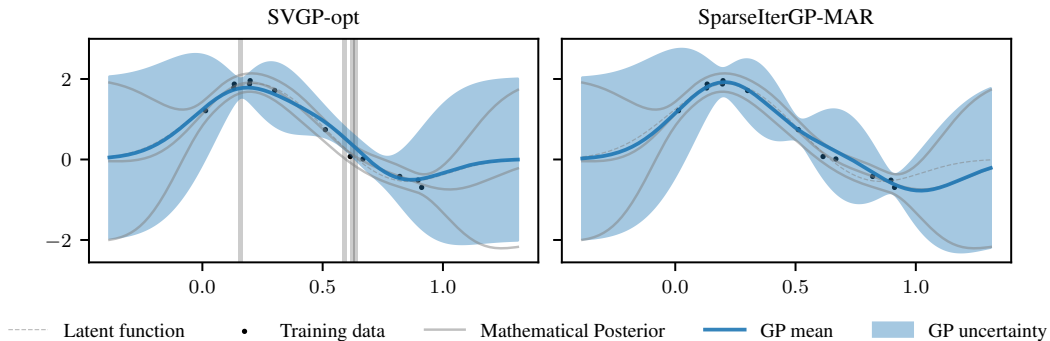
SVGP versus IterGP-MAR

Linear-time computation-aware GP inference with IterGP.



Unpublished work

Policy: Unit vector actions $\mathbf{s}_i = \mathbf{e}_j$ which select points greedily as $j = \arg \max r_{i-1} \Rightarrow \mathcal{O}(nm^2)$.



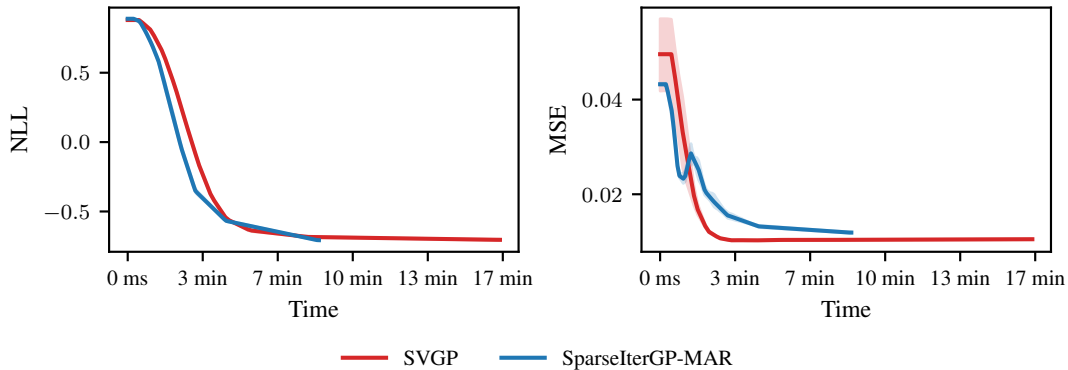


SVGP versus IterGP-MAR: Large-Scale Problem

Linear-time computation-aware GP inference with IterGP on a problem with $n \approx 10^5$ datapoints.

Unpublished work

Policy: Unit vector actions $\mathbf{s}_i = \mathbf{e}_j$ which select points greedily as $j = \arg \max r_{i-1} \Rightarrow \mathcal{O}(nm^2)$.



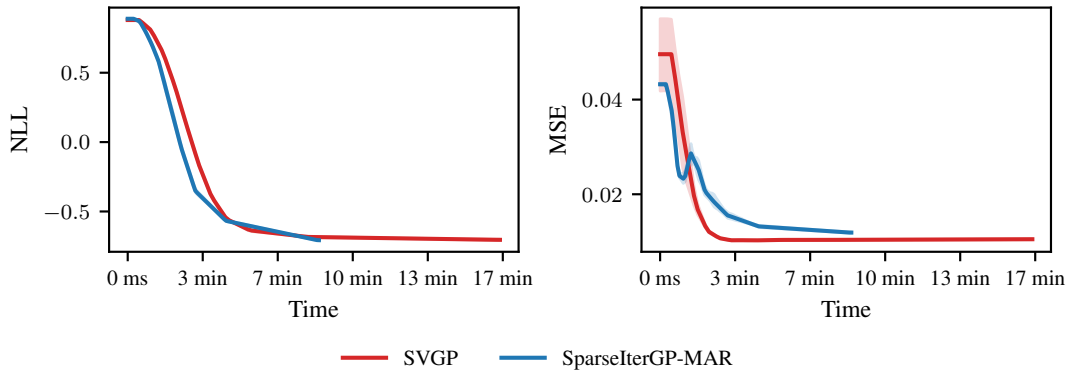


SVGP versus IterGP-MAR: Large-Scale Problem

Linear-time computation-aware GP inference with IterGP on a problem with $n \approx 10^5$ datapoints.

Unpublished work

Policy: Unit vector actions $\mathbf{s}_i = \mathbf{e}_j$ which select points greedily as $j = \arg \max r_{i-1} \Rightarrow \mathcal{O}(nm^2)$.



Scalable GP approximation without inadvertently compromising uncertainty quantification.

Bonus: Getting Philosophical

Blurring the lines between data and computation.

Working with Infinite Data

For IterGP it does not matter how large the dataset is, or whether we have it stored on our machine.

(Wenger et al., 2022a)

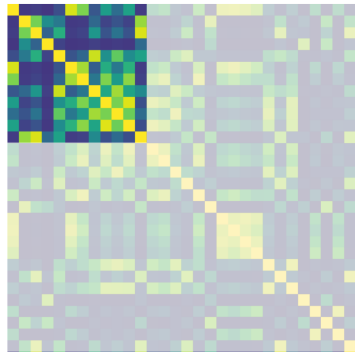
Theorem (Online GP Approximation with IterGP)

Let $n, n' \in \mathbb{N}$ and consider training data sets $\mathbf{X} \in \mathbb{R}^{n \times d}$, $\mathbf{y} \in \mathbb{R}^n$ and $\mathbf{X}' \in \mathbb{R}^{n' \times d}$, $\mathbf{y}' \in \mathbb{R}^{n'}$. Consider two sequences of actions $(\mathbf{s}_i)_{i=1}^n \in \mathbb{R}^n$ and $(\tilde{\mathbf{s}}_i)_{i=1}^{n+n'} \in \mathbb{R}^{n+n'}$ such that

$$\tilde{\mathbf{s}}_i = \begin{pmatrix} \mathbf{s}_i \\ \mathbf{0} \end{pmatrix} \quad (2)$$

Then the posterior returned by IterGP for the dataset (\mathbf{X}, \mathbf{y}) using actions \mathbf{s}_i is identical to the posterior returned by IterGP for the extended dataset using actions $\tilde{\mathbf{s}}_i$:

$$\text{ITERGP}(\mu, k, \mathbf{X}, \mathbf{y}, (\mathbf{s}_i)_i) = \text{ITERGP}\left(\mu, k, \begin{pmatrix} \mathbf{X} \\ \mathbf{X}' \end{pmatrix}, \begin{pmatrix} \mathbf{y} \\ \mathbf{y}' \end{pmatrix}, (\tilde{\mathbf{s}}_i)_i\right).$$



Data is as Data Does

An alternative view of IterGP as a better model for the way we do inference instead of an approximation.

Observation: Only once we perform computation on data, does it enter our prediction.



Data is as Data Does

An alternative view of IterGP as a better model for the way we do inference instead of an approximation.

Observation: Only once we perform computation on data, does it enter our prediction.



The distinction between data and computation vanishes from this perspective.

Data is as Data Does

An alternative view of IterGP as a better model for the way we do inference instead of an approximation.

Observation: Only once we perform computation on data, does it enter our prediction.



The distinction between data and computation vanishes from this perspective.

What if we modelled this situation with a Gaussian process?

$$\begin{aligned}
 f &\sim \mathcal{GP}(\mu, k) \\
 \tilde{y} \mid f(X) &\sim \mathcal{N}(\mathbf{S}_i^\top f(X), \sigma^2 \mathbf{S}_i^\top \mathbf{S}_i) \\
 f \mid X, \tilde{y} &\sim \mathcal{GP}(\mu_i, k_i)
 \end{aligned}$$

Data is as Data Does

An alternative view of IterGP as a better model for the way we do inference instead of an approximation.

Observation: Only once we perform computation on data, does it enter our prediction.



The distinction between data and computation vanishes from this perspective.

What if we modelled this situation with a Gaussian process?

$$\begin{aligned}
 f &\sim \mathcal{GP}(\mu, k) \\
 \tilde{y} \mid f(X) &\sim \mathcal{N}(\mathbf{S}_i^\top f(X), \sigma^2 \mathbf{S}_i^\top \mathbf{S}_i) \\
 f \mid X, \tilde{y} &\sim \mathcal{GP}(\mu_i, k_i)
 \end{aligned}$$

IterGP's combined posterior is equivalent to exact GP regression for linearly projected data.

Takeaways

- ▶ Large-scale models are often as much about the approximation as they are about the data.

Takeaways

- ▶ Large-scale models are often as much about the approximation as they are about the data.
- ▶ **Uncertainty** arises from **limited data** and from **limited computation**.

Takeaways

- ▶ Large-scale models are often as much about the approximation as they are about the data.
- ▶ **Uncertainty** arises from **limited data** and from **limited computation**.
- ▶ For GPs, we can exactly quantify the approximation error given arbitrary resources \Rightarrow IterGP.

Takeaways

- ▶ Large-scale models are often as much about the approximation as they are about the data.
- ▶ **Uncertainty** arises from **limited data** and from **limited computation**.
- ▶ For GPs, we can exactly quantify the approximation error given arbitrary resources \Rightarrow IterGP.
- ▶ **Explicit trade-off** between computation and uncertainty.

Takeaways

- ▶ Large-scale models are often as much about the approximation as they are about the data.
- ▶ **Uncertainty** arises from **limited data** and from **limited computation**.
- ▶ For GPs, we can exactly quantify the approximation error given arbitrary resources \Rightarrow IterGP.
- ▶ **Explicit trade-off** between computation and uncertainty.

Open Questions

- ▶ Model selection / hyperparameter optimization?

Takeaways

- ▶ Large-scale models are often as much about the approximation as they are about the data.
- ▶ **Uncertainty** arises from **limited data** and from **limited computation**.
- ▶ For GPs, we can exactly quantify the approximation error given arbitrary resources \Rightarrow IterGP.
- ▶ **Explicit trade-off** between computation and uncertainty.

Open Questions

- ▶ Model selection / hyperparameter optimization?
- ▶ Policy design for downstream tasks and decision making problems.
 - ▶ Active learning
 - ▶ Bayesian optimization
 - ▶ ...

Takeaways

- ▶ Large-scale models are often as much about the approximation as they are about the data.
- ▶ **Uncertainty** arises from **limited data** and from **limited computation**.
- ▶ For GPs, we can exactly quantify the approximation error given arbitrary resources \Rightarrow IterGP.
- ▶ **Explicit trade-off** between computation and uncertainty.

Open Questions

- ▶ Model selection / hyperparameter optimization?
- ▶ Policy design for downstream tasks and decision making problems.
 - ▶ Active learning
 - ▶ Bayesian optimization
 - ▶ ...
- ▶ Extension to non-Gaussian likelihoods.

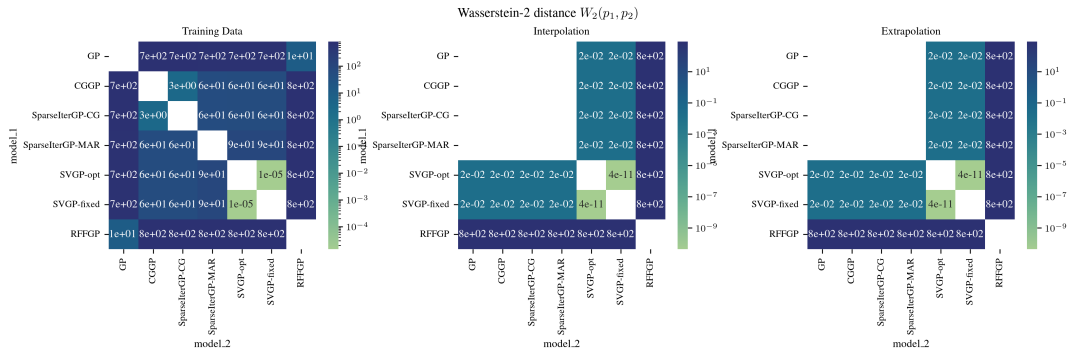
Comparison of GP Approximations

Gaussian Process Classification

Large-scale Model Selection

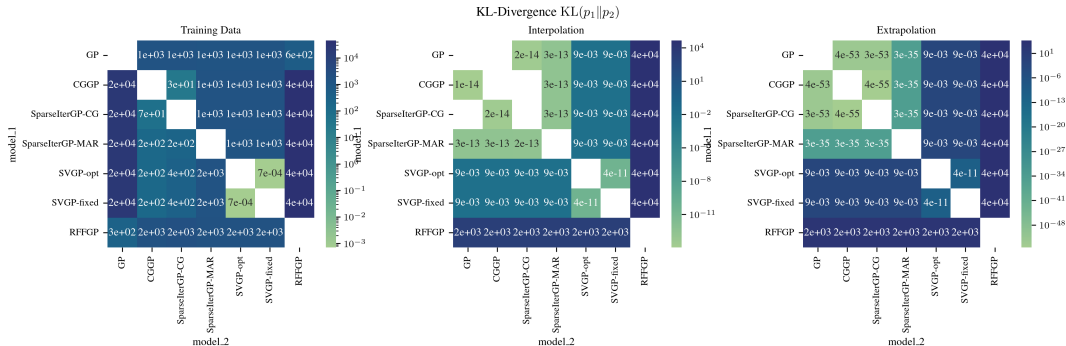
Comparison of GP Approximations: Wasserstein-2 Distance

Comparison of different GP approximations at the training data, for interpolation and extrapolation.



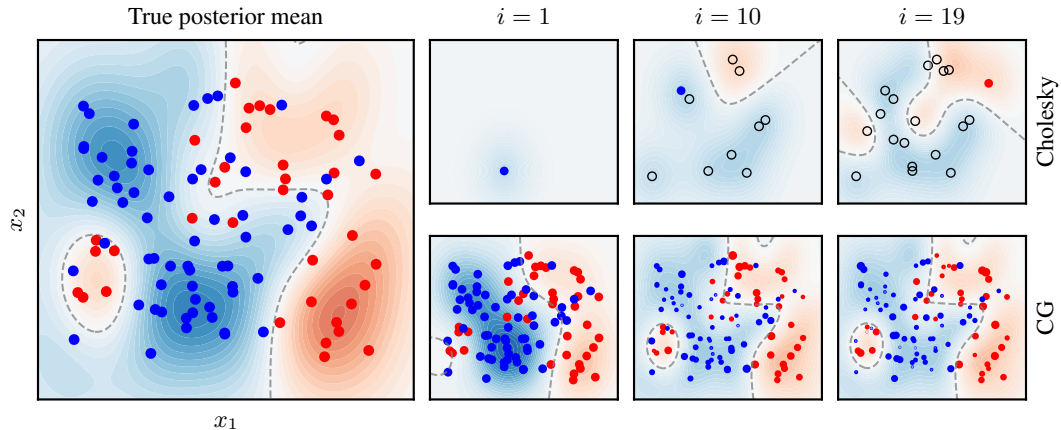
Comparison of GP Approximations: KL-Divergence

Comparison of different GP approximations at the training data, for interpolation and extrapolation.



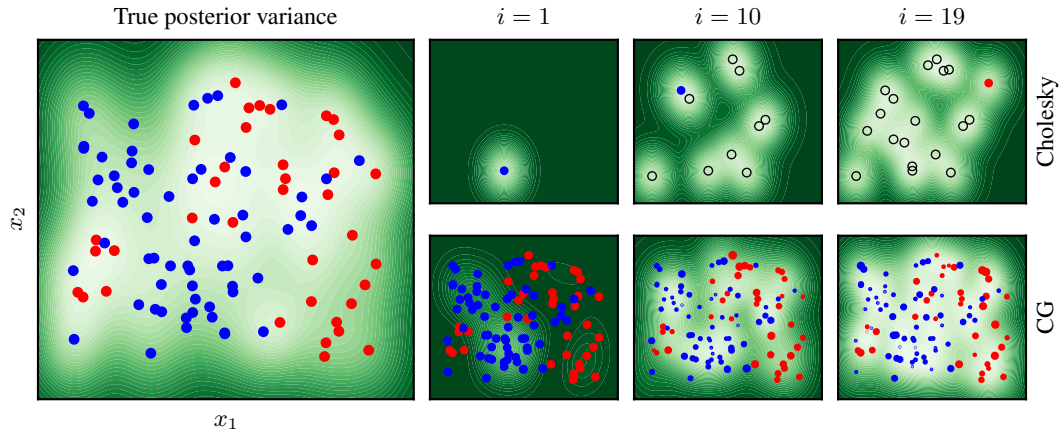
Gaussian Process Classification

Extension to non-Gaussian likelihoods via Laplace Approximation.



Gaussian Process Classification

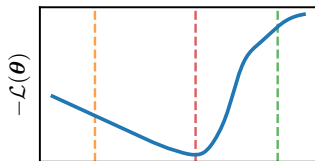
Extension to non-Gaussian likelihoods via Laplace Approximation.



Model Selection for Gaussian Processes

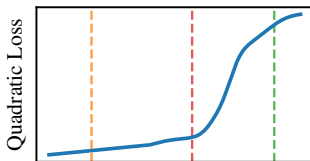
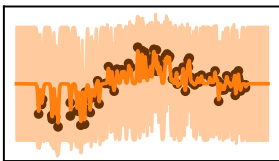
We can identify kernel hyperparameters by optimizing the log-marginal likelihood.

$$\theta_* = \arg \max_{\theta} \mathcal{L}(\theta) = \arg \max_{\theta} \log p(\mathbf{y} \mid \theta) = \arg \min_{\theta} \left(\underbrace{(\mathbf{y} - \boldsymbol{\mu})^\top \hat{\mathbf{K}}^{-1} (\mathbf{y} - \boldsymbol{\mu})}_{\text{quadratic loss}} + \underbrace{\log \det(\hat{\mathbf{K}})}_{\text{model complexity}} \right)$$



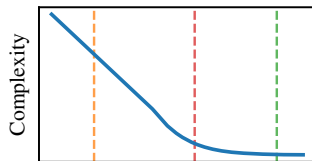
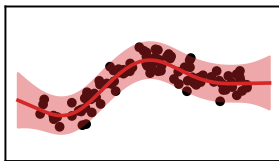
Lengthscale θ

Overfitting



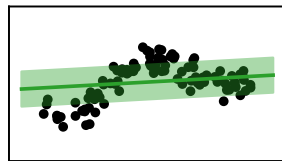
Lengthscale θ

Appropriate Fit



Lengthscale θ

Underfitting



Large-scale GP Hyperparameter Optimization

A numerical linear algebra bottleneck.

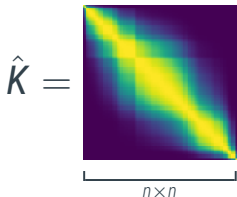
(Ubaru et al., 2017; Gardner et al., 2018)

Need to: Evaluate log-marginal likelihood and its derivative repeatedly.

- ▶ log-marginal likelihood $\mathcal{L}(\theta) = -\frac{1}{2} (y^\top \hat{K}^{-1} y + \log \det(\hat{K}) + n \log(2\pi))$
- ▶ derivative $\frac{\partial}{\partial \theta} \mathcal{L}(\theta) = \frac{1}{2} y^\top \hat{K}^{-1} \frac{\partial \hat{K}}{\partial \theta} \hat{K}^{-1} y - \frac{1}{2} \text{tr}(\hat{K}^{-1} \frac{\partial \hat{K}}{\partial \theta})$

Challenge: Computationally costly operations with the kernel matrix.

- ▶ linear solves $v \mapsto \hat{K}^{-1} v$
- ▶ matrix traces $\log \det(\hat{K}) = \text{tr}(\log(\hat{K}))$ and $\text{tr}(\hat{K}^{-1} \frac{\partial \hat{K}}{\partial \theta_i})$



Linear solves and matrix traces can be computed solely via *matrix-vector multiplication*!

This is great because ...

- ▶ matrix-vector multiplies have complexity $\mathcal{O}(n^2)$.
- ▶ structured or sparse matrices are efficient to multiply with.
- ▶ the kernel matrix does not need to be stored in memory explicitly (Charlier et al., 2021).
- ▶ we can exploit parallelization and modern hardware (GPUs).

lower time and space complexity

Preconditioning

How to encode and leverage structural prior knowledge about matrices.

Preconditioner

$$\hat{P} \approx \hat{K}$$

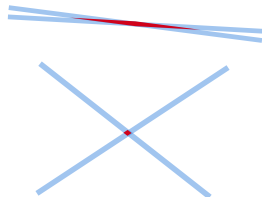
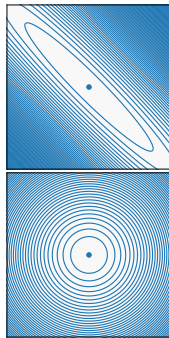
such that $\kappa(\hat{P}^{-1}\hat{K}) \ll \kappa(\hat{K})$ and \hat{P} is computationally tractable.

- Computing and storing \hat{P} is cheap.
- Linear solves $\mathbf{v} \mapsto \hat{P}^{-1}\mathbf{v}$ are efficient.
- Derived properties, such as the determinant or spectrum are known.

Asymptotic approx. error $g(\ell) \rightarrow 0$ of sequence of preconditioners $\hat{P}_\ell \rightarrow \hat{K}$:

$$\kappa(\hat{P}_\ell^{-1}\hat{K}) \leq (1 + \mathcal{O}(g(\ell))\|\hat{K}\|_F)^2$$

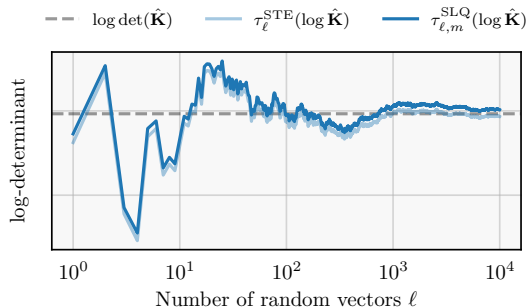
Known Use: Accelerate and stabilize linear solves via CG \Rightarrow **bias reduction**



Stochastic Trace Estimation

Computing matrix traces $\text{tr}(f(\hat{\mathbf{K}}))$ via matrix-vector multiplication.

(Hutchinson, 1989; Golub et al., 2009; Ubaru et al., 2017)



$$\begin{aligned}\text{tr}(f(\hat{\mathbf{K}})) &= n\mathbb{E}[\mathbf{z}_i^\top f(\hat{\mathbf{K}})\mathbf{z}_i] \\ &\approx \tau_\ell^{\text{STE}}(f(\hat{\mathbf{K}})) = \frac{n}{\ell} \sum_{i=1}^{\ell} \mathbf{z}_i^\top f(\hat{\mathbf{K}})\mathbf{z}_i \\ &\approx \tau_{\ell,m}^{\text{SLQ}}(f(\hat{\mathbf{K}}))\end{aligned}$$

Problems:

- ▶ Worst-case convergence in the number of random vectors is $\mathcal{O}(\ell^{-\frac{1}{2}})$ \Rightarrow slows down training
- ▶ Introduces stochasticity into hyperparameter optimization

Preconditioned Log-Determinant Estimation

Variance-reduced stochastic trace estimation via preconditioning.

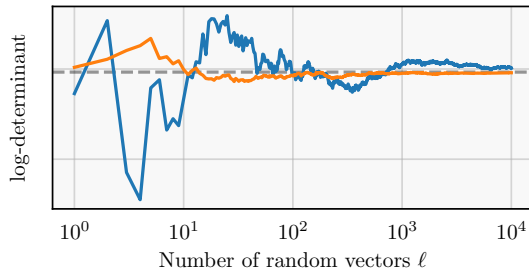
(Wenger et al., 2022b)

Idea: Decompose log-determinant into deterministic and stochastic approximation.

$$\log \det(\hat{K}) = \log \det(\hat{P}_\ell \hat{P}_\ell^{-1} \hat{K}) = \underbrace{\log \det(\hat{P}_\ell)}_{\text{known}} + \underbrace{\text{tr}(\log(\hat{K}) - \log(\hat{P}_\ell))}_{\approx \text{stochastic trace estimate}}$$

The better the preconditioner, the smaller the stochastic approximation \Rightarrow **variance reduction**

— $\log \det(\hat{K})$ $\text{— } \tau_{\ell,m}^{\text{SLQ}}(\log \hat{K})$ $\text{— } \log \det(\hat{P}) + \tau_{\ell,m}^{\text{SLQ}}(\log \hat{P}^{-1} \hat{K})$



- Backward pass analogously via automatic differentiation.
- If we compute a preconditioner for CG, we can simply reuse it at negligible overhead.
- If $\hat{P}_\ell \rightarrow \hat{K}$ at rate $g(\ell)$, then the STE only requires $\mathcal{O}(\ell^{-\frac{1}{2}} g(\ell))$ random vectors.

Convergence Rates for Kernel – Preconditioner Combinations

The faster the preconditioner converges to the kernel matrix (i.e. $g(\ell) \rightarrow 0$) the fewer random vectors are needed.

If $\hat{P}_\ell \rightarrow \hat{K}$ at rate $g(\ell)$, then the STE only requires $\mathcal{O}(\ell^{-\frac{1}{2}} g(\ell))$ random vectors.

Kernel	d	Preconditioner	$g(\ell)$	Condition
any	\mathbb{N}	none	1	
any	\mathbb{N}	truncated SVD	$\ell^{-\frac{1}{2}}$	
any	\mathbb{N}	random. SVD	$\ell^{-\frac{1}{2}} + \mathcal{O}(\ell^{\frac{1}{4}} s^{-\frac{1}{4}})$	w/ high prob. for s samples
any	\mathbb{N}	random. Nyström	$\ell^{-\frac{1}{2}} + \mathcal{O}(\ell^{\frac{1}{4}} s^{-\frac{1}{4}})$	w/ high prob. for s samples
any	\mathbb{N}	RFF	$\ell^{-\frac{1}{2}}$	w/ high prob.
RBF	1	partial Cholesky	$\exp(-c\ell)$	for some $c > 0$
RBF	\mathbb{N}	QFF	$\exp(-b\ell^{\frac{1}{d}})$	for some $b > 0$ if $\ell^{\frac{1}{d}} > 2\gamma^{-2}$
Matérn(ν)	\mathbb{N}	partial Cholesky	$\ell^{-(\frac{2\nu}{d}+1)}$	$2\nu \in \mathbb{N}$, maximin ordering Schaefer2021a
Matérn(ν)	1	QFF	$\ell^{-(s(\nu)+1)}$	where $s(\nu) \in \mathbb{N}$
mod. Matérn(ν)	\mathbb{N}	QFF	$\ell^{-\frac{s(\nu)+1}{d}}$	where $s(\nu) \in \mathbb{N}$
additive	\mathbb{N}	any	$dg(\ell)$	all summands have rate $g(\ell)$
any	\mathbb{N}	any kernel approx.	$g(\ell)$	\exists uniform convergence bound

Theoretical Guarantees

Probabilistic error bounds for the estimates of the log-marginal likelihood and its derivative.

Theorem (Log-marginal likelihood)

[...] Then with probability $1 - \delta$, the error in the estimate η of the log-marginal likelihood \mathcal{L} satisfies

$$|\eta - \mathcal{L}| \leq \varepsilon_{\text{CG}} + \frac{1}{2}(\varepsilon_{\text{Lanczos}} + \varepsilon_{\text{STE}}) \|\log(\hat{K})\|_F,$$

where the individual errors are bounded by

$$\varepsilon_{\text{CG}}(\kappa, i) \leq K_3 \left(\frac{\sqrt{\kappa-1}}{\sqrt{\kappa+1}} \right)^i \quad (3)$$

$$\varepsilon_{\text{Lanczos}}(\kappa, i) \leq K_1 \left(\frac{\sqrt{2\kappa+1}-1}{\sqrt{2\kappa+1}+1} \right)^{2i} \quad (4)$$

$$\varepsilon_{\text{STE}}(\delta, \ell) \leq C_1 \sqrt{\log(\delta^{-1})} \ell^{-\frac{1}{2}} g(\ell) \quad (5)$$

Theorem (Derivative)

[...] Then with probability $1 - \delta$, the error in the estimate ϕ of the derivative of the log-marginal likelihood $\frac{\partial}{\partial \theta} \mathcal{L}$ satisfies

$$|\phi - \frac{\partial}{\partial \theta} \mathcal{L}| \leq \varepsilon_{\text{CG}} + \frac{1}{2}(\varepsilon_{\text{CG}'} + \varepsilon_{\text{STE}}) \|\hat{K}^{-1} \frac{\partial \hat{K}}{\partial \theta}\|_F$$

where the individual errors are bounded by

$$\varepsilon_{\text{CG}}(\kappa, i) \leq K_4 \left(\frac{\sqrt{\kappa-1}}{\sqrt{\kappa+1}} \right)^i \quad (6)$$

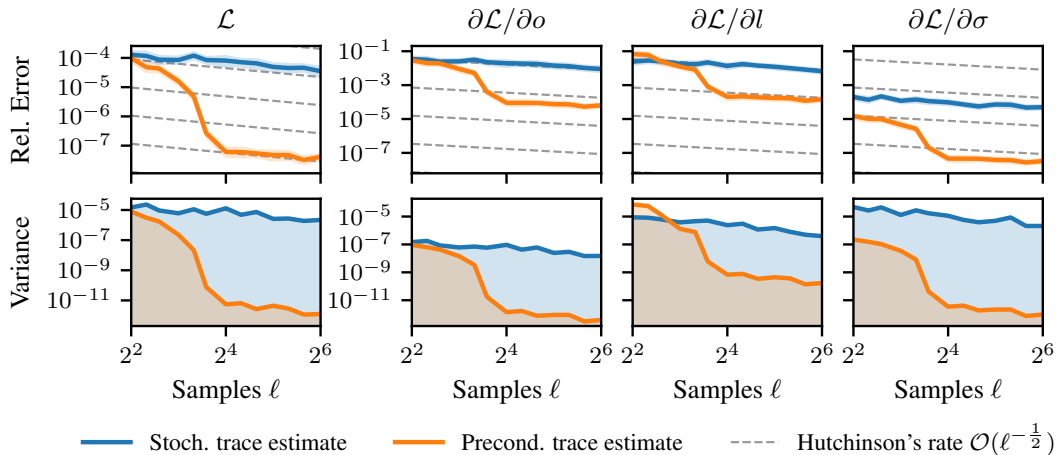
$$\varepsilon_{\text{CG}'}(\kappa, i) \leq K_2 \left(\frac{\sqrt{\kappa-1}}{\sqrt{\kappa+1}} \right)^i \quad (7)$$

$$\varepsilon_{\text{STE}}(\delta, \ell) \leq C_1 \sqrt{\log(\delta^{-1})} \ell^{-\frac{1}{2}} g(\ell) \quad (8)$$

We leverage preconditioning not only to **reduce bias**, but crucially also to **reduce variance**.

Preconditioning Reduces Bias and Variance

Estimating the log-marginal likelihood and its derivatives on synthetic data.



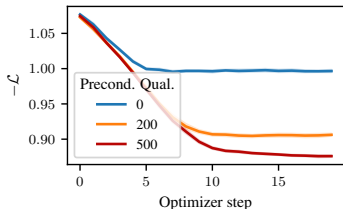
Experiment Details:

► Randomly sampled synthetic data ($n = 10,000, d = 1$)

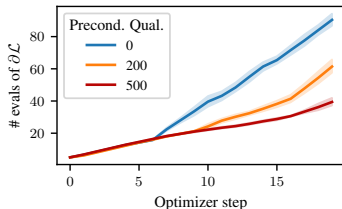
► RBF kernel with noise scale $\sigma^2 = 10^{-2}$

Preconditioning Accelerates Hyperparameter Optimization

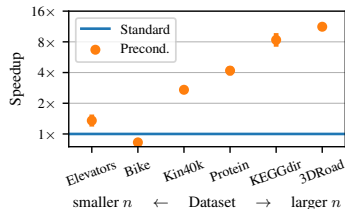
Gaussian process hyperparameter optimization on UCI data.



(a) Training loss (Protein).



(b) Line search computations (Protein).



(c) Speedup on UCI datasets.

Experiment Details:

- ▶ UCI datasets ($n = 12,449$ to $n = 326,155$)
- ▶ Matérn($\frac{3}{2}$) kernel with noise scale $\sigma^2 = 10^{-2}$
- ▶ Partial Cholesky preconditioner of size 500
- ▶ $\ell = 50$ random vectors

- ▶ Alireza Radmanesh, Matthew J. Muckley, Tullie Murrell, Emma Lindsey, Anuroop Sriram, Florian Knoll, Daniel K. Sodickson, and Yvonne W. Lui. “Exploring the Acceleration Limits of Deep Learning Variational Network-based Two-dimensional Brain MRI”. In: *Radiology: Artificial Intelligence* 4.6 (2022). DOI: [10.1148/ryai.210313](https://doi.org/10.1148/ryai.210313) (cit. on pp. 2–5).
- ▶ Jacob R Gardner, Geoff Pleiss, David Bindel, Kilian Q Weinberger, and Andrew Gordon Wilson. “GPYtorch: Blackbox matrix-matrix Gaussian process inference with GPU acceleration”. In: *Advances in Neural Information Processing Systems (NeurIPS)* (2018) (cit. on pp. 18–20, 90).
- ▶ Benjamin Charlier, Jean Feydy, Joan Alexis Glaunès, François-David Collin, and Ghislain Durif. “Kernel Operations on the GPU, with Autodiff, without Memory Overflows”. In: *Journal of Machine Learning Research* 22.74 (2021), pp. 1–6 (cit. on pp. 18–20, 90).
- ▶ Philipp Hennig. “Probabilistic Interpretation of Linear Solvers”. In: *SIAM Journal on Optimization* 25.1 (2015), pp. 234–260 (cit. on pp. 22–24).
- ▶ Jon Cockayne, Chris J. Oates, Ilse C.F. Ipsen, and Mark Girolami. “A Bayesian Conjugate Gradient Method (with Discussion)”. In: *Bayesian Analysis* 14.3 (2019), pp. 937–1012. DOI: [10.1214/19-BA1145](https://doi.org/10.1214/19-BA1145) (cit. on pp. 22–24).

- ▶ Jonathan Wenger and Philipp Hennig. “Probabilistic Linear Solvers for Machine Learning”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2020 (cit. on pp. 22–24).
- ▶ Jonathan Wenger, Geoff Pleiss, Marvin Pförtner, Philipp Hennig, and John P. Cunningham. “Posterior and Computational Uncertainty in Gaussian Processes”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2022 (cit. on pp. 25–28, 34–39, 49, 72).
- ▶ Michalis Titsias. “Variational learning of inducing variables in sparse Gaussian processes”. In: *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2009 (cit. on pp. 64–67).
- ▶ James Hensman, Nicolò Fusi, and Neil D Lawrence. “Gaussian processes for big data”. In: *Conference on Uncertainty in Artificial Intelligence (UAI)*. 2013 (cit. on pp. 64–67).
- ▶ Shashanka Ubaru, Jie Chen, and Yousef Saad. “Fast estimation of $\text{tr}(f(A))$ via stochastic Lanczos quadrature”. In: *SIAM Journal on Matrix Analysis and Applications* 38.4 (2017), pp. 1075–1099 (cit. on pp. 90, 92).
- ▶ Michael F Hutchinson. “A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines”. In: *Communications in Statistics-Simulation and Computation* 18.3 (1989), pp. 1059–1076 (cit. on p. 92).

- ▶ Gene H Golub and Gérard Meurant. *Matrices, moments and quadrature with applications*. Vol. 30. Princeton University Press, 2009 (cit. on p. 92).
- ▶ Jonathan Wenger, Geoff Pleiss, Philipp Hennig, John P. Cunningham, and Jacob R. Gardner. “Preconditioning for Scalable Gaussian Process Hyperparameter Optimization”. In: *International Conference on Machine Learning (ICML)*. 2022 (cit. on p. 93).