

# Probabilistic Linear Solvers

Jonathan Wenger

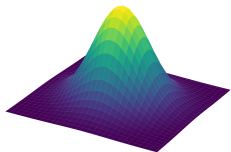
COLUMBIA | Zuckerman Institute

# Linear Systems are Everywhere in Scientific Computing

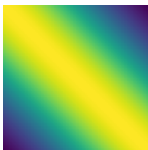


Arguably, the most fundamental numerical task in scientific computing and machine learning.

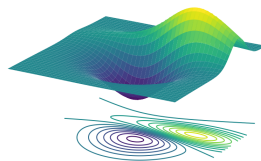
Basic Statistics



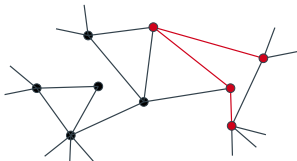
Probabilistic / Kernel Methods



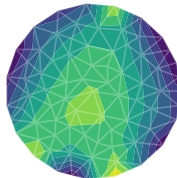
Optimization



Graphs and (Neural) Networks



Differential Equations



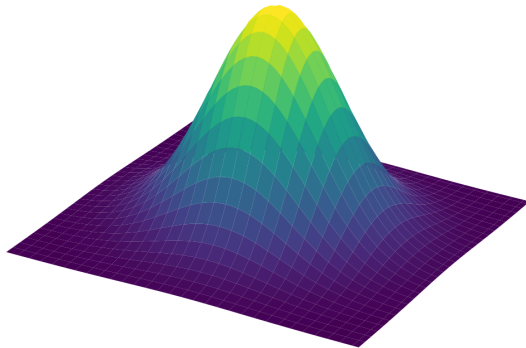
...and many more.



Example: Probability theory.

## Normal Distribution

$$x \sim \mathcal{N}(\mu, \Sigma)$$
$$p(x) = \frac{1}{\sqrt{(2\pi)^n \det(\Sigma)}} \exp\left(-\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu)\right)$$



# Linear Systems are Everywhere in Scientific Computing



Example: Probabilistic Models and Kernel Methods.

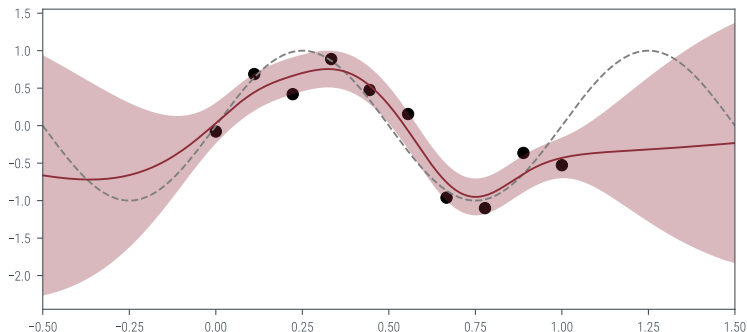
## Gaussian Processes

$$f \sim \mathcal{GP}(\mu, k)$$

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

$$\mu_{\text{post}}(x) = \mu(x) + k(x, X)(k(X, X) + \sigma^2 I)^{-1}(y - \mu(X))$$

$$k_{\text{post}}(x_0, x_1) = k(x_0, x_1) - k(x_0, X)(k(X, X) + \sigma^2 I)^{-1}k(X, x_1)$$





# Linear Systems are Everywhere in Scientific Computing



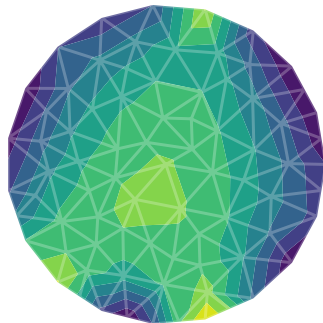
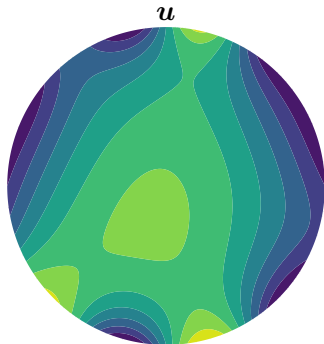
Example: Linear Differential Equations.

## Galerkin Method

$$\underbrace{Du = f}_{\text{linear differential equation}}$$

$\Rightarrow$

$$\underbrace{\hat{D}\hat{u} = \hat{f}}_{\text{finite dimensional linear system}}$$



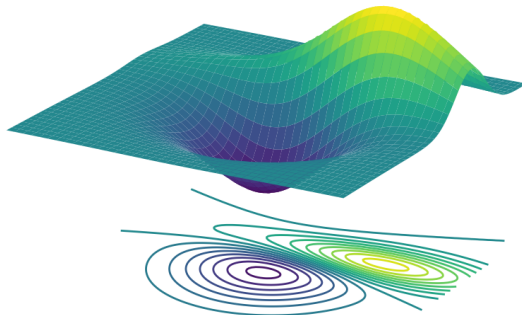
Example: Optimization.

## Iterative Optimization Methods

$$\boldsymbol{\theta}_i \approx \arg \min_{\boldsymbol{\theta} \in \Theta} \mathcal{L}(\boldsymbol{\theta})$$

$$\boldsymbol{\theta}_i = \boldsymbol{\theta}_{i-1} + \alpha_i \mathbf{M}_i \mathbf{d}_i$$

Examples: natural / conjugate / stochastic gradient descent, (Quasi-) Newton method, ...

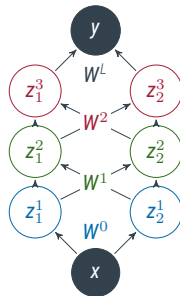


## Feedforward Neural Network

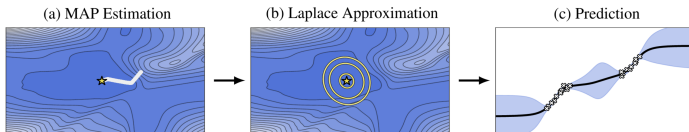
$$z^0(x, \theta) = x$$

$$z^{\ell+1}(x, \theta) = \sigma(W^\ell z^\ell + b^\ell)$$

$$y := f(x, \theta) = z^L(x, \theta)$$



Bayesian deep learning via Laplace approximation:  $p(\theta \mid \mathcal{D}) \approx \mathcal{N}(\theta; \theta_{\text{MAP}}, (\nabla_{\theta}^2 \mathcal{L}(\theta)|_{\theta_{\text{MAP}}})^{-1})$



Daxberger et al. [Dax+22]

# Probabilistic Linear Solvers

Learning the solution of a linear system.



## Goal

Solve **large-scale** linear system  $Ax_* = b$  for  $x_* \in \mathbb{R}^n$ .



## Goal

Solve **large-scale** linear system  $Ax_* = b$  for  $x_* \in \mathbb{R}^n$ .

## Core Insights of Probabilistic Numerics

[HOG15; Coc+19a; HOK22]

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

## Goal

Solve **large-scale** linear system  $Ax_* = b$  for  $x_* \in \mathbb{R}^n$ .

## Core Insights of Probabilistic Numerics

[HOG15; Coc+19a; HOK22]

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

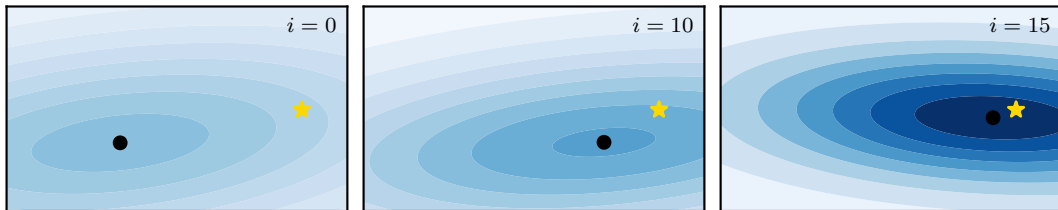
## Goal

Solve **large-scale** linear system  $Ax_* = b$  for  $x_* \in \mathbb{R}^n$ .

## Core Insights of Probabilistic Numerics

[HOG15; Coc+19a; HOK22]

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



Solution  $x_*$



Estimate  $x_i = \mathbb{E}(x_*)$



Belief  $p(x_*)$



# Learning The Solution

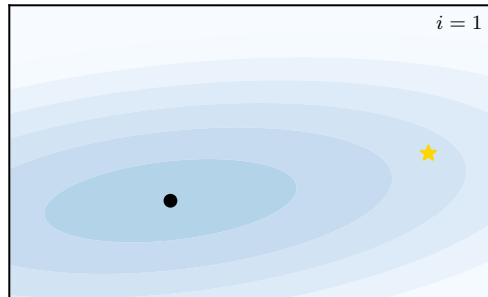


Estimating the solution of a linear system with a probabilistic linear solver.

[Hen15; Coc+19b]

Goal: Solve  $Ax_* = b$  for  $x_*$ .

Prior:  $x_* \sim \mathcal{N}(x_0, \Sigma_0)$



- ★ Solution  $x_*$
- Approximation  $x_{i-1}$
- Belief  $p(x_*) = \mathcal{N}(x_{i-1}, \Sigma_{i-1})$

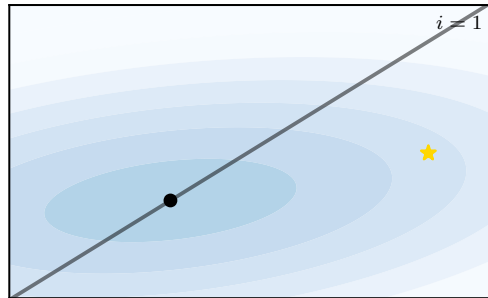
**Goal:** Solve  $Ax_* = b$  for  $x_*$ .

**Prior:**  $x_* \sim \mathcal{N}(x_0, \Sigma_0)$

**Likelihood:** Observe  $x_*$  via arbitrary actions  $s_j$ :

$$\alpha_j := s_j^\top A(x_* - x_{i-1}) = s_j^\top r_{i-1}$$

$$p(\alpha_j \mid x_*) = \lim_{\varepsilon \rightarrow 0} \mathcal{N}(\alpha_j; 0, \varepsilon)$$



- ★ Solution  $x_*$
- Approximation  $x_{i-1}$
- Belief  $p(x_*) = \mathcal{N}(x_{i-1}, \Sigma_{i-1})$
- Action  $s_i$

**Goal:** Solve  $A\mathbf{x}_* = \mathbf{b}$  for  $\mathbf{x}_*$ .

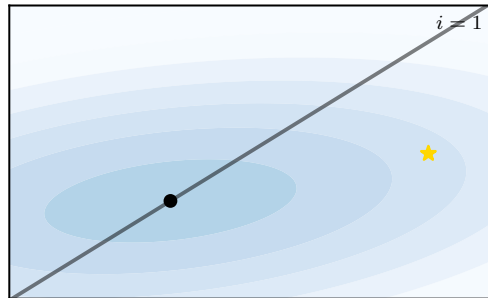
**Prior:**  $\mathbf{x}_* \sim \mathcal{N}(\mathbf{x}_0, \Sigma_0)$

**Likelihood:** Observe  $\mathbf{x}_*$  via arbitrary actions  $\mathbf{s}_i$ :

$$\alpha_i := \mathbf{s}_i^\top A(\mathbf{x}_* - \mathbf{x}_{i-1}) = \mathbf{s}_i^\top \mathbf{r}_{i-1}$$

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

**Posterior:** Bayes' rule gives a closed form update!



- ★ Solution  $\mathbf{x}_*$
- Approximation  $\mathbf{x}_{i-1}$
- Belief  $p(\mathbf{x}_*) = \mathcal{N}(\mathbf{x}_{i-1}, \Sigma_{i-1})$
- Action  $\mathbf{s}_i$

**Goal:** Solve  $A\mathbf{x}_* = \mathbf{b}$  for  $\mathbf{x}_*$ .

**Prior:**  $\mathbf{x}_* \sim \mathcal{N}(\mathbf{x}_0, \Sigma_0)$

**Likelihood:** Observe  $\mathbf{x}_*$  via arbitrary actions  $\mathbf{s}_i$ :

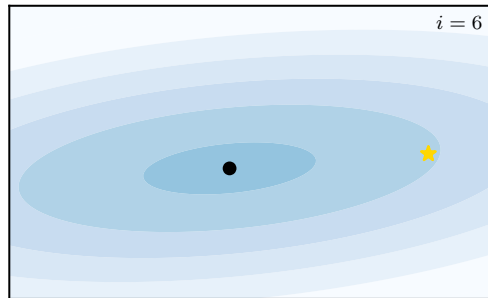
$$\alpha_i := \mathbf{s}_i^\top A(\mathbf{x}_* - \mathbf{x}_{i-1}) = \mathbf{s}_i^\top \mathbf{r}_{i-1}$$

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

**Posterior:**  $\mathbf{x}_* \mid \alpha_1, \dots, \alpha_i \sim \mathcal{N}(\mathbf{x}_i, \Sigma_i)$

$$\mathbf{x}_i = \mathbf{x}_0 + \Sigma_0 A \mathbf{S}_i (\mathbf{S}_i^\top A \Sigma_0 A \mathbf{S}_i)^{-1} \mathbf{S}_i^\top (\mathbf{b} - A \mathbf{x}_0)$$

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



- ★ Solution  $\mathbf{x}_*$
- Approximation  $\mathbf{x}_{i-1}$
- Belief  $p(\mathbf{x}_*) = \mathcal{N}(\mathbf{x}_{i-1}, \Sigma_{i-1})$

**Goal:** Solve  $Ax_* = b$  for  $x_*$ .

**Prior:**  $x_* \sim \mathcal{N}(x_0, \Sigma_0)$

**Likelihood:** Observe  $x_*$  via arbitrary actions  $s_i$ :

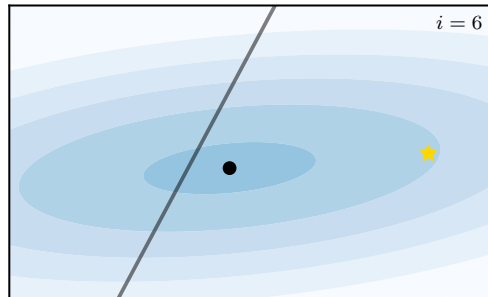
$$\alpha_i := s_i^\top A(x_* - x_{i-1}) = s_i^\top r_{i-1}$$

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

**Posterior:**  $x_* | \alpha_1, \dots, \alpha_i \sim \mathcal{N}(x_i, \Sigma_i)$

$$x_i = x_0 + \Sigma_0 A S_i (S_i^\top A \Sigma_0 A S_i)^{-1} S_i^\top (b - A x_0)$$

$$\Sigma_i = \Sigma_0 - \Sigma_0 A S_i (S_i^\top A \Sigma_0 A S_i)^{-1} S_i^\top A \Sigma_0$$



- ★ Solution  $x_*$
- Approximation  $x_{i-1}$
- Belief  $p(x_*) = \mathcal{N}(x_{i-1}, \Sigma_{i-1})$
- Action  $s_i$

**Goal:** Solve  $A\mathbf{x}_* = \mathbf{b}$  for  $\mathbf{x}_*$ .

**Prior:**  $\mathbf{x}_* \sim \mathcal{N}(\mathbf{x}_0, \Sigma_0)$

**Likelihood:** Observe  $\mathbf{x}_*$  via arbitrary actions  $\mathbf{s}_i$ :

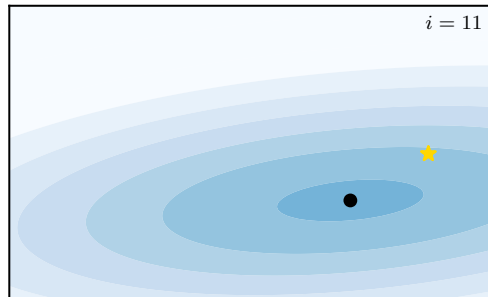
$$\alpha_i := \mathbf{s}_i^\top A(\mathbf{x}_* - \mathbf{x}_{i-1}) = \mathbf{s}_i^\top \mathbf{r}_{i-1}$$

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

**Posterior:**  $\mathbf{x}_* \mid \alpha_1, \dots, \alpha_i \sim \mathcal{N}(\mathbf{x}_i, \Sigma_i)$

$$\mathbf{x}_i = \mathbf{x}_0 + \Sigma_0 A \mathbf{S}_i (\mathbf{S}_i^\top A \Sigma_0 A \mathbf{S}_i)^{-1} \mathbf{S}_i^\top (\mathbf{b} - A \mathbf{x}_0)$$

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



- ★ Solution  $\mathbf{x}_*$
- Approximation  $\mathbf{x}_{i-1}$
- Belief  $p(\mathbf{x}_*) = \mathcal{N}(\mathbf{x}_{i-1}, \Sigma_{i-1})$

**Goal:** Solve  $Ax_* = b$  for  $x_*$ .

**Prior:**  $x_* \sim \mathcal{N}(x_0, \Sigma_0)$

**Likelihood:** Observe  $x_*$  via arbitrary actions  $s_i$ :

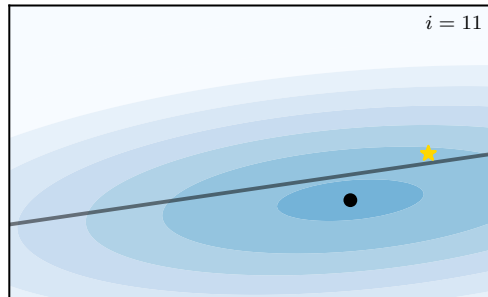
$$\alpha_i := s_i^\top A(x_* - x_{i-1}) = s_i^\top r_{i-1}$$

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

**Posterior:**  $x_* | \alpha_1, \dots, \alpha_i \sim \mathcal{N}(x_i, \Sigma_i)$

$$x_i = x_0 + \Sigma_0 A S_i (S_i^\top A \Sigma_0 A S_i)^{-1} S_i^\top (b - A x_0)$$

$$\Sigma_i = \Sigma_0 - \Sigma_0 A S_i (S_i^\top A \Sigma_0 A S_i)^{-1} S_i^\top A \Sigma_0$$



- ★ Solution  $x_*$
- Approximation  $x_{i-1}$
- Belief  $p(x_*) = \mathcal{N}(x_{i-1}, \Sigma_{i-1})$
- Action  $s_i$

**Goal:** Solve  $Ax_* = b$  for  $x_*$ .

**Prior:**  $x_* \sim \mathcal{N}(x_0, \Sigma_0)$

**Likelihood:** Observe  $x_*$  via arbitrary actions  $s_i$ :

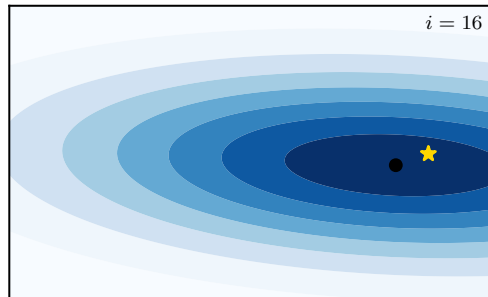
$$\alpha_i := s_i^\top A(x_* - x_{i-1}) = s_i^\top r_{i-1}$$

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

**Posterior:**  $x_* | \alpha_1, \dots, \alpha_i \sim \mathcal{N}(x_i, \Sigma_i)$

$$x_i = x_0 + \Sigma_0 A S_i (S_i^\top A \Sigma_0 A S_i)^{-1} S_i^\top (b - A x_0)$$

$$\Sigma_i = \Sigma_0 - \Sigma_0 A S_i (S_i^\top A \Sigma_0 A S_i)^{-1} S_i^\top A \Sigma_0$$



- ★ Solution  $x_*$
- Approximation  $x_{i-1}$
- Belief  $p(x_*) = \mathcal{N}(x_{i-1}, \Sigma_{i-1})$



**Goal:** Solve  $A\mathbf{x}_* = \mathbf{b}$  for  $\mathbf{x}_*$ .

**Prior:**  $\mathbf{x}_* \sim \mathcal{N}(\mathbf{x}_0, \Sigma_0)$

**Likelihood:** Observe  $\mathbf{x}_*$  via arbitrary actions  $\mathbf{s}_i$ :

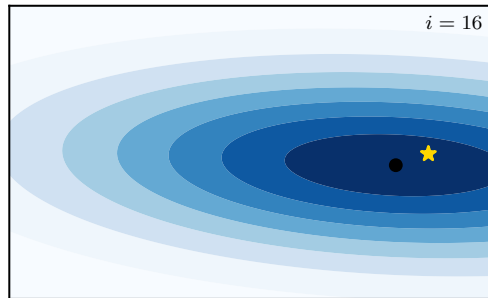
$$\alpha_i := \mathbf{s}_i^\top A(\mathbf{x}_* - \mathbf{x}_{i-1}) = \mathbf{s}_i^\top \mathbf{r}_{i-1}$$

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

**Posterior:**  $\mathbf{x}_* | \alpha_1, \dots, \alpha_i \sim \mathcal{N}(\mathbf{x}_i, \Sigma_i)$

$$\mathbf{x}_i = \mathbf{x}_0 + \Sigma_0 A \mathbf{S}_i (\mathbf{S}_i^\top A \Sigma_0 A \mathbf{S}_i)^{-1} \mathbf{S}_i^\top (\mathbf{b} - A \mathbf{x}_0)$$

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



- ★ Solution  $\mathbf{x}_*$
- Approximation  $\mathbf{x}_{i-1}$
- Belief  $p(\mathbf{x}_*) = \mathcal{N}(\mathbf{x}_{i-1}, \Sigma_{i-1})$

How do we choose the linear solver actions  $\mathbf{S}$  and the prior  $\mathcal{N}(\mathbf{x}_0, \Sigma_0)$ ?



**Observation:** Actions “weigh” entries in the residual:  $\alpha_i := \mathbf{s}_i^\top \mathbf{r}_{i-1} = \mathbf{s}_i^\top \mathbf{A}(\mathbf{x}_* - \mathbf{x}_{i-1})$



**Observation:** Actions “weigh” entries in the residual:  $\alpha_i := \mathbf{s}_i^\top \mathbf{r}_{i-1} = \mathbf{s}_i^\top \mathbf{A}(\mathbf{x}_* - \mathbf{x}_{i-1})$

**Idea:** Focus computation where residual is large:  $\mathbf{s}_i = \mathbf{r}_{i-1} \Rightarrow \alpha_i = \|\mathbf{r}_{i-1}\|_2^2$

$\Rightarrow$  **BayesCG** [Coc+19b]

# Interlude: Method of Conjugate Gradients



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

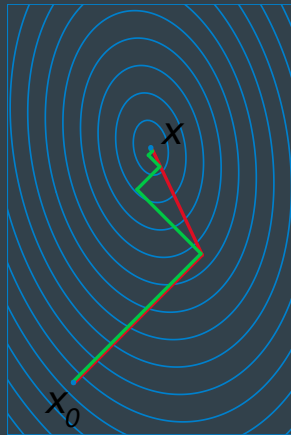
**Goal:** Approximately solve linear system  $Ax_* = b$ .

**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?



Oleg Alexandrov, [commons.wikimedia.org/w/index.php?curid=2267598](https://commons.wikimedia.org/w/index.php?curid=2267598)

# Interlude: Method of Conjugate Gradients



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

**Goal:** Approximately solve linear system  $Ax_* = b$ .

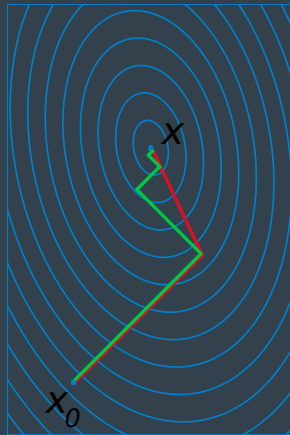
**Idea:** Rephrase as quadratic optimization problem and optimize. Let

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

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.  $\langle d_i, d_j \rangle = 0$ .



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

# Interlude: Method of Conjugate Gradients



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

**Goal:** Approximately solve linear system  $Ax_* = b$ .

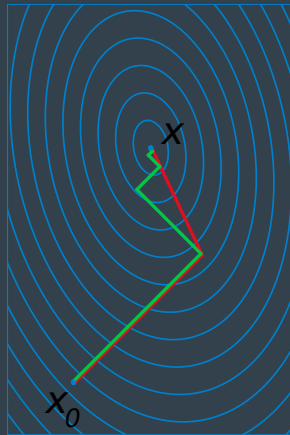
**Idea:** Rephrase as quadratic optimization problem and optimize. Let

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

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.



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

# Interlude: Method of Conjugate Gradients



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

**Goal:** Approximately solve linear system  $Ax_* = b$ .

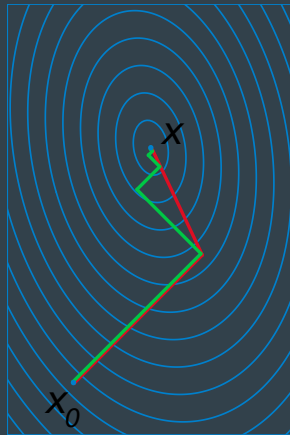
**Idea:** Rephrase as quadratic optimization problem and optimize. Let

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

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



**Observation:** Actions  $\mathbf{s}_i$  “weigh” entries in the residual:  $\alpha_i := \mathbf{s}_i^\top \mathbf{r}_{i-1} = \mathbf{s}_i^\top \mathbf{A}(\mathbf{x}_* - \mathbf{x}_{i-1})$

**Idea:** Focus computation where residual is large:  $\mathbf{s}_i = \mathbf{r}_{i-1} \Rightarrow \alpha_i = \|\mathbf{r}_{i-1}\|_2^2$

$\Rightarrow$  **BayesCG** [Coc+19b]

**Theorem (Equivalence to Conjugate Gradient Method [Coc+19b; Wen+22])**

*If  $\mathbf{x}_0 = \mathbf{0}$ ,  $\Sigma_0 = \mathbf{A}^{-1}$  and the actions are either conjugate gradients  $\mathbf{s}_i = \mathbf{d}_i^{\text{CG}}$  or gradients  $\mathbf{s}_i = \mathbf{r}_{i-1}$ , then the posterior mean  $\mathbf{x}_i = \mathbf{x}_i^{\text{CG}}$  of BayesCG is equivalent to the approximation returned by CG.*

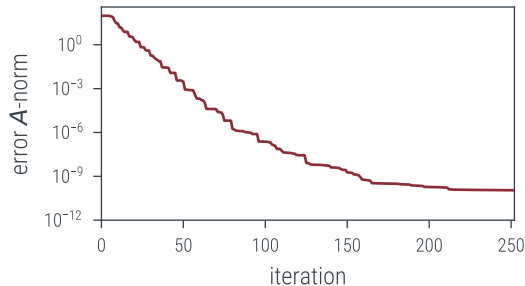


# Convergence Behavior of the Conjugate Gradient Method



The spectrum of the matrix determines the convergence speed.

$$n = 10^3 \quad \kappa(A) \approx 7 \cdot 10^5$$



Theorem (Convergence Rate of CG<sub>[TB97]</sub>)

$$\|x - x_i\|_A \leq 2 \left( \frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \right)^i \|x - x_0\|_A$$

CG converges fast for a small condition number.

# Prior Choice

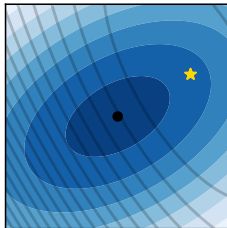
Comparing different choices of prior for BayesCG.



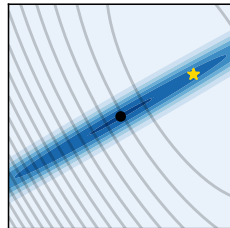
[Coc+19b]

Prior

$$x_* \sim \mathcal{N}(x_0, \Sigma_0)$$



$\Rightarrow$



# Prior Choice

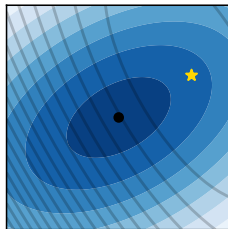
Comparing different choices of prior for BayesCG.



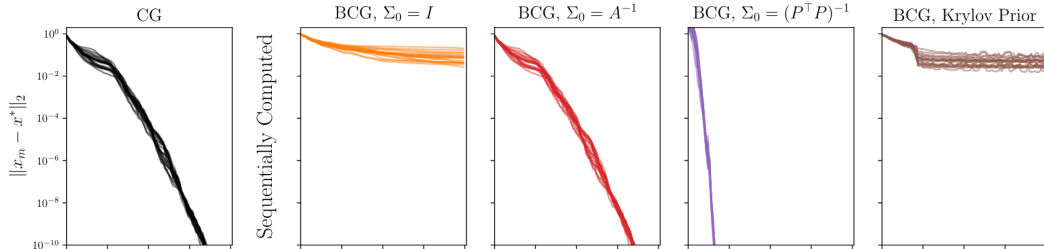
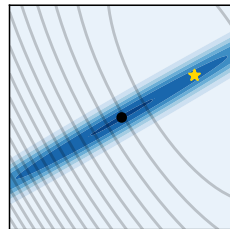
[Coc+19b]

Prior

$$x_* \sim \mathcal{N}(x_0, \Sigma_0)$$



$\Rightarrow$



# Algorithm: Probabilistic Linear Solver



Sequential formulation.

		Time	Space
1	<b>procedure</b> PROBABILISTICLINEARSOLVER( $\mathbf{A}, \mathbf{b}, \mathbf{x}_0 = \mathbf{0}, \Sigma_0$ )		
2	<b>while not</b> STOPPINGCRITERION() <b>do</b>		
3	$\mathbf{s}_i \leftarrow \text{POLICY}()$	Select action via policy.	
4	$\mathbf{r}_{i-1} \leftarrow \mathbf{b} - \mathbf{A}\mathbf{x}_{i-1}$	Residual.	$\mathcal{O}(n)$
5	$\alpha_i \leftarrow \mathbf{s}_i^\top \mathbf{r}_{i-1}$	Observation.	$\mathcal{O}(1)$
6	$\mathbf{z}_i \leftarrow \mathbf{A}\mathbf{s}_i$		$\mathcal{O}(n)$
7	$\mathbf{d}_i \leftarrow \Sigma_{i-1}\mathbf{A}\mathbf{s}_i = \Sigma_{i-1}\mathbf{z}_i$	Search direction.	$\mathcal{O}(n)$
8	$\eta_i \leftarrow \mathbf{s}_i^\top \mathbf{A} \Sigma_{i-1} \mathbf{A} \mathbf{s}_i = \mathbf{z}_i^\top \mathbf{d}_i$		$\mathcal{O}(1)$
9	$\mathbf{C}_i \leftarrow \mathbf{C}_{i-1} + \frac{1}{\eta_i} \mathbf{d}_i \mathbf{d}_i^\top$		$\mathcal{O}(ni)$
10	$\mathbf{x}_i \leftarrow \mathbf{x}_{i-1} + \frac{\alpha_i}{\eta_i} \mathbf{d}_i$	Solution estimate.	$\mathcal{O}(n)$
11	$\Sigma_i \leftarrow \Sigma_0 - \mathbf{C}_i$	Uncertainty.	
12	<b>return</b> $\mathcal{N}(\mathbf{x}_i, \Sigma_i)$		

# Application: Gaussian Processes

Scaling Gaussian processes via probabilistic linear solvers.

# Gaussian Process Regression



Learning an unknown function from data.

**Goal:** Supervised learning from  $n$  data points  $(\mathbf{X}, \mathbf{y})$

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

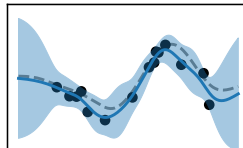
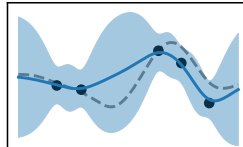
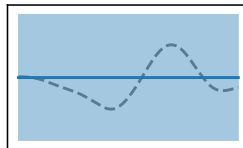
**Likelihood:** Observations  $\mathbf{y} = f(\mathbf{X}) + \epsilon \sim \mathcal{N}(f(\mathbf{X}), \sigma^2 I)$

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

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

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

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



# Gaussian Process Regression



Learning an unknown function from data.

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

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

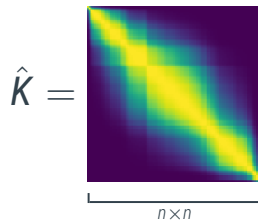
**Likelihood:** Observations  $y = f(X) + \varepsilon \sim \mathcal{N}(f(X), \sigma^2 I)$

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

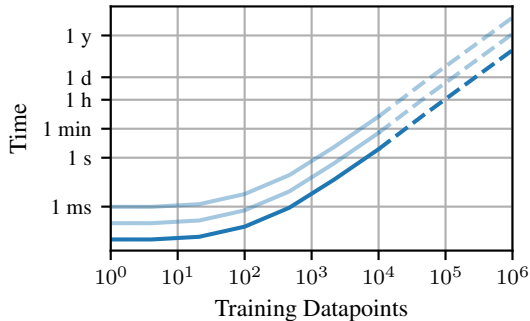


# Computational Cost of Gaussian Processes

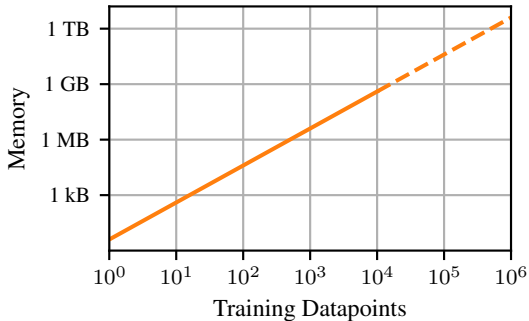


Gaussian processes scale prohibitively with the size  $n$  of the dataset.

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



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



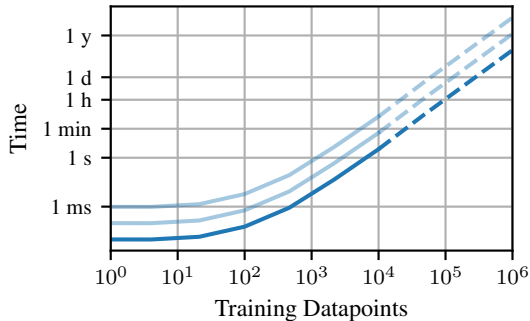


# Computational Cost of Gaussian Processes

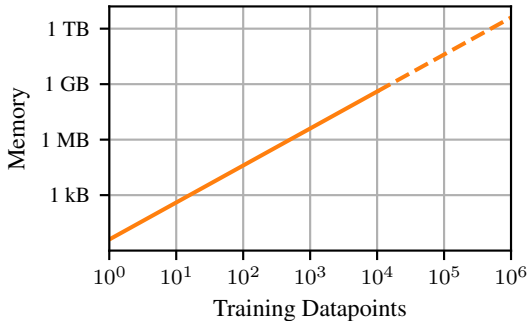


Gaussian processes scale prohibitively with the size  $n$  of the dataset.

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



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



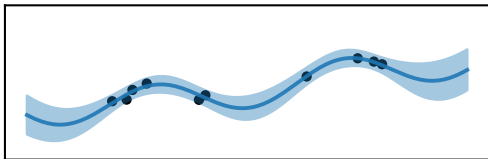
We need to **approximate** the posterior.

# Approximate Gaussian Process Inference

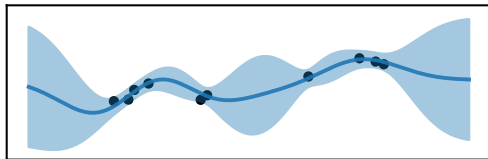


Impact of approximations on uncertainty quantification and decision-making.

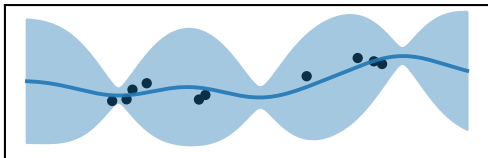
RFFGP



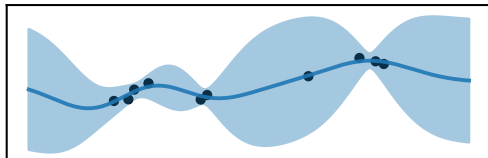
CGGP



SVGP-fixed



SVGP-opt



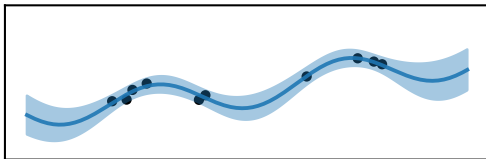
● Data    — Approx. Posterior Mean    ■ Approx. Posterior Uncertainty

# Approximate Gaussian Process Inference

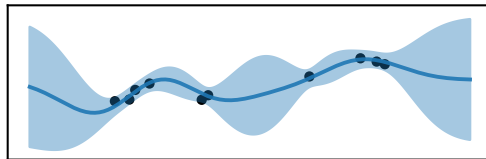


Impact of approximations on uncertainty quantification and decision-making.

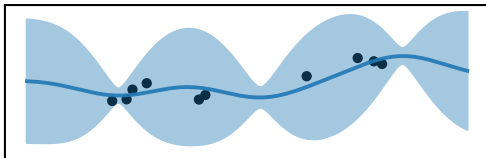
RFFGP



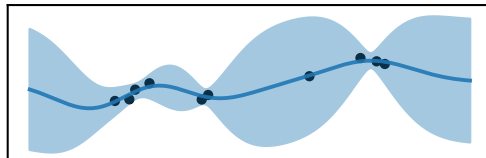
CGGP



SVGP-fixed



SVGP-opt



● Data    — Approx. Posterior Mean    ■ Approx. Posterior Uncertainty

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.

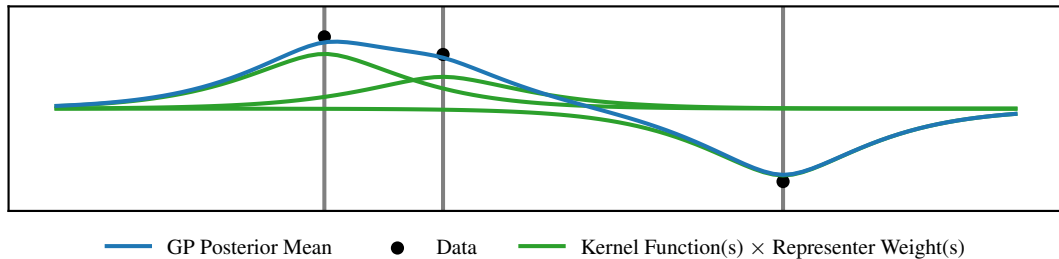
# Representer Weights



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

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

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



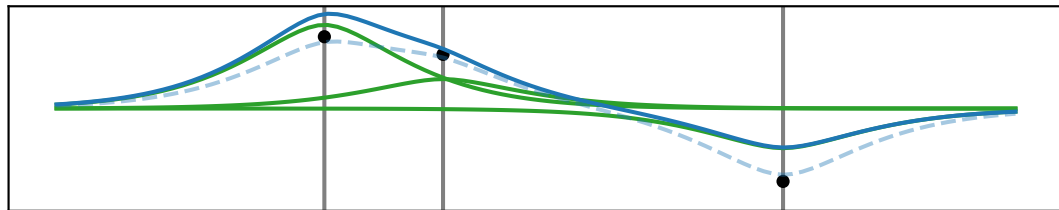
# Approximating Representer Weights



Iterative linear solvers can be used to approximate the representer weights.

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

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



— Approx. GP Posterior Mean

● Data

— Kernel Function(s)  $\times$  Approx. Representer Weight(s)



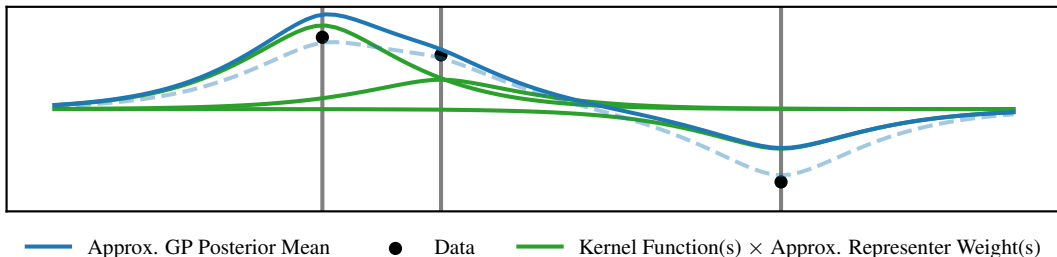
# Approximating Representer Weights



Iterative linear solvers can be used to approximate the representer weights.

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

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



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

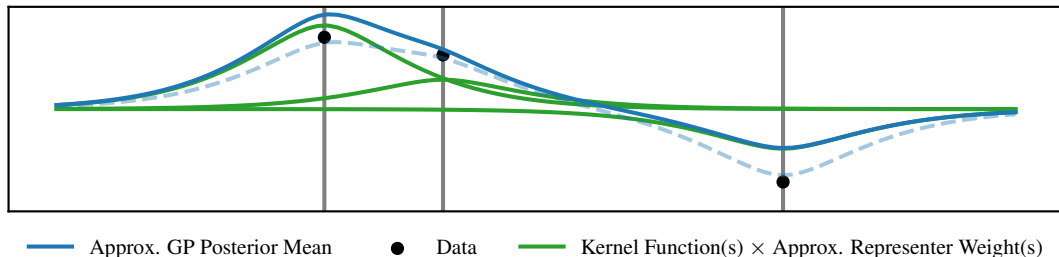
# Approximating Representer Weights



Iterative linear solvers can be used to approximate the representer weights.

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

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



**Problem:** Approximation error of the linear solve.

## Q2: Can We Quantify Approximation Error?

Probabilistic error quantification at prediction time using probabilistic linear solvers.

# Linear Solver Prior for GP Inference



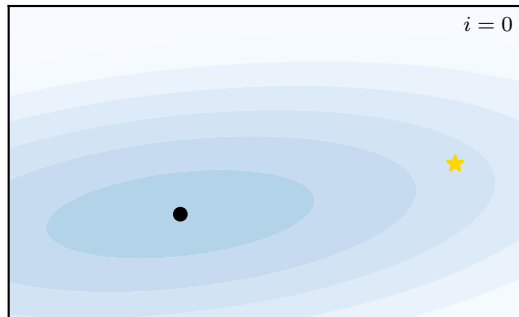
The Gaussian process prior makes assumptions about the representer weights.

[Wen+22]

## Observation:

GP prior induces representer weights prior:

$$y - \mu \sim \mathcal{N}(\mathbf{0}, \hat{K})$$



- ★ Solution  $\mathbf{x}_*$
- Approximation  $\mathbf{x}_i$
- Belief  $p(\mathbf{x}_*)$

# Linear Solver Prior for GP Inference



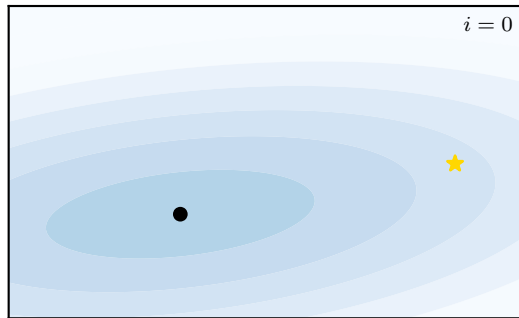
The Gaussian process prior makes assumptions about the representer weights.

[Wen+22]

## Observation:

GP prior induces representer weights prior:

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



- ★ Solution  $\mathbf{x}_*$
- Approximation  $\mathbf{x}_i$
- Belief  $p(\mathbf{x}_*)$

# Linear Solver Posterior for GP Inference

Estimation of representer weights with a probabilistic linear solver.

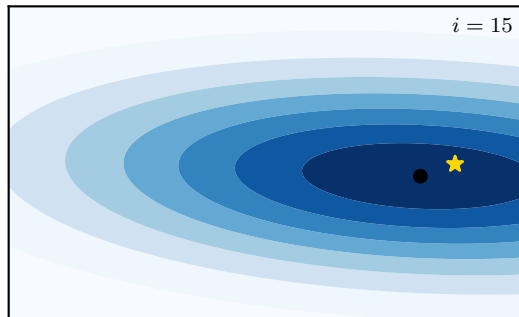
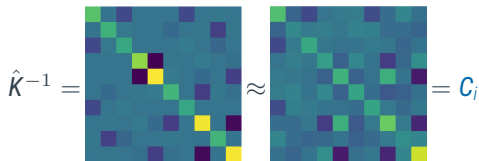


[Wen+22]

Representer weights posterior  $\mathbf{v}_* \sim \mathcal{N}(\mathbf{v}_i, \Sigma_i)$ , s.t.

$$\mathbf{v}_i = \mathbf{C}_i(\mathbf{y} - \boldsymbol{\mu})$$

$$\Sigma_i = \hat{K}^{-1} - \mathbf{C}_i$$



- ★ Solution  $\mathbf{x}_*$
- Approximation  $\mathbf{x}_i$
- Belief  $p(\mathbf{x}_*)$

# Linear Solver Posterior for GP Inference

Estimation of representer weights with a probabilistic linear solver.

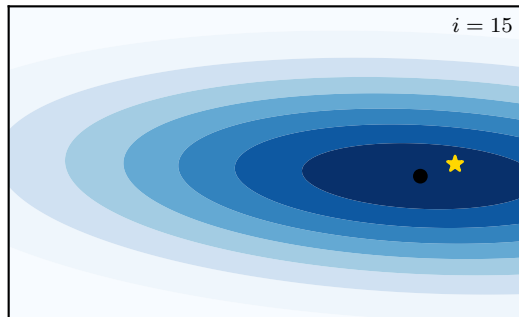
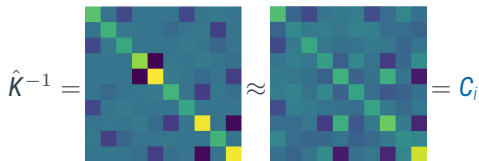


[Wen+22]

Representer weights posterior  $\mathbf{v}_* \sim \mathcal{N}(\mathbf{v}_i, \Sigma_i)$ , s.t.

$$\mathbf{v}_i = \mathbf{C}_i(\mathbf{y} - \boldsymbol{\mu})$$

$$\Sigma_i = \hat{\mathbf{K}}^{-1} - \mathbf{C}_i$$



- ★ Solution  $\mathbf{x}_*$
- Approximation  $\mathbf{x}_i$
- Belief  $p(\mathbf{x}_*)$

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

# IterGP: Computation-Aware Gaussian Process Inference



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

[Wen+22]

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



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

**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.

[Wen+22]

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

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

**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.

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

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

**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.

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

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

**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:  $(f \mid \mathbf{y})(\cdot) \stackrel{d}{=} f(\cdot) + k(\cdot, \mathbf{X}) \underbrace{\hat{\mathbf{K}}^{-1}(\mathbf{y} - \boldsymbol{\mu})}_{=\mathbf{v}_*} \stackrel{d}{=} (f \mid \mathbf{v}_*)(\cdot)$

2 Marginalize representer weights belief:  $p(f(\cdot)) = \int p(f(\cdot) \mid \mathbf{v}_*) p(\mathbf{v}_*) d\mathbf{v}_* = \mathcal{GP}(f; \mu_i, k_i)$

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

**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.

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

2 Marginalize representer weights belief:  $p(f(\cdot)) = \int p(f(\cdot) \mid \mathbf{v}_*) p(\mathbf{v}_*) d\mathbf{v}_* = \mathcal{GP}(f; \mu_i, k_i)$

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

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

**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.

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

2 Marginalize representer weights belief:  $p(f(\cdot)) = \int p(f(\cdot) | \mathbf{v}_*) p(\mathbf{v}_*) d\mathbf{v}_* = \mathcal{GP}(f; \mu_i, k_i)$

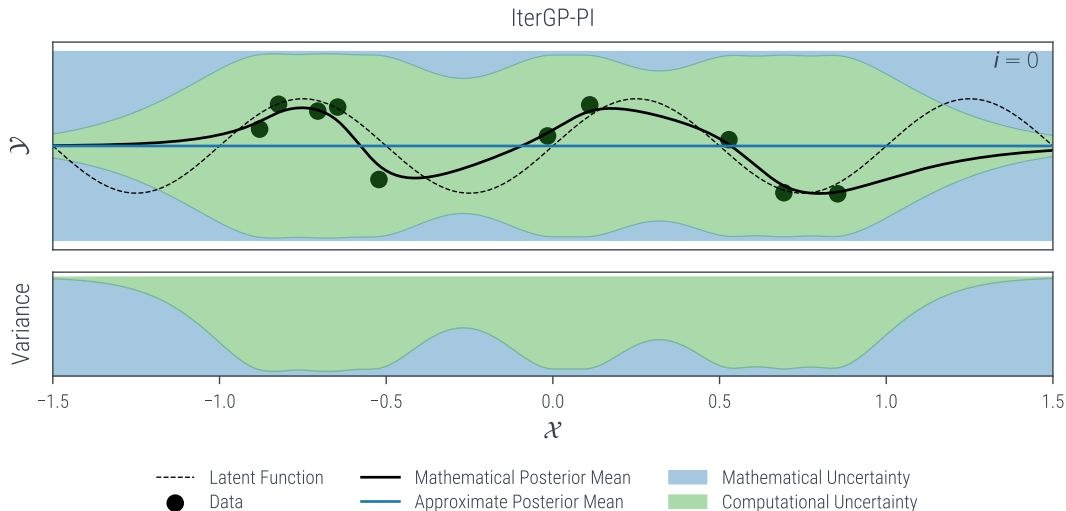
$$\mu_i(\cdot) = \mu(\cdot) + K(\cdot, \mathbf{X}) \mathbf{v}_i$$

$$k_i(\cdot, \cdot) = \underbrace{K(\cdot, \cdot) - K(\cdot, \mathbf{X}) \hat{K}^{-1} K(\mathbf{X}, \cdot)}_{\substack{= \mathbb{E}((f(\cdot) - \mu_*(\cdot))^2) \\ \text{mathematical uncertainty} \quad \text{blue circle}}} + \underbrace{K(\cdot, \mathbf{X}) \boldsymbol{\Sigma}_i K(\mathbf{X}, \cdot)}_{\substack{= \mathbb{E}((\mu_*(\cdot) - \mu_i(\cdot))^2) \\ \text{computational uncertainty} \quad \text{green circle}}} = \underbrace{K(\cdot, \cdot) - K(\cdot, \mathbf{X}) \mathbf{C}_i K(\mathbf{X}, \cdot)}_{\substack{= \mathbb{E}((f(\cdot) - \mu_i(\cdot))^2) \\ \text{combined uncertainty} \quad \text{purple circle}}}$$

# Computation-Aware GP Inference Illustrated



Interpreting computational and combined uncertainty as error quantification.

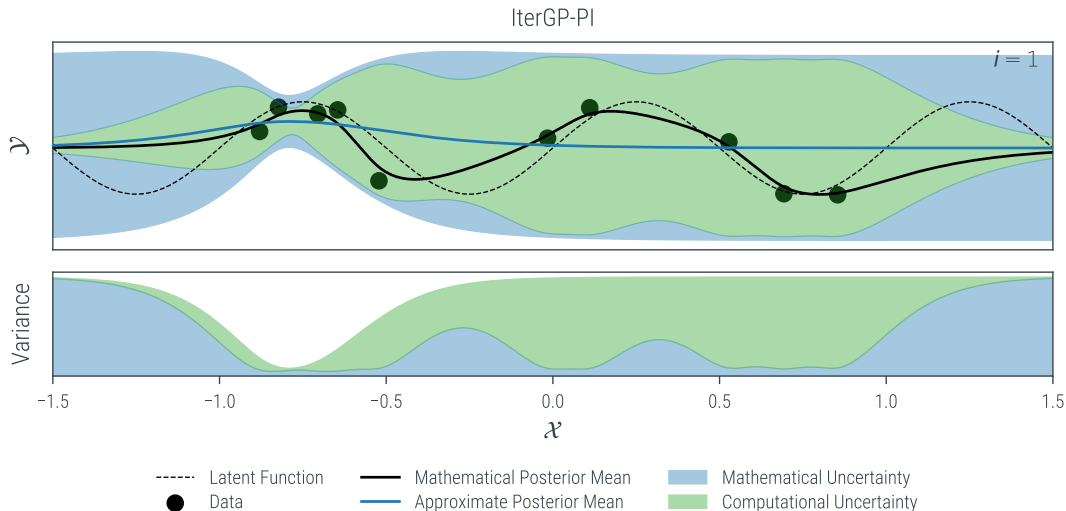




# Computation-Aware GP Inference Illustrated



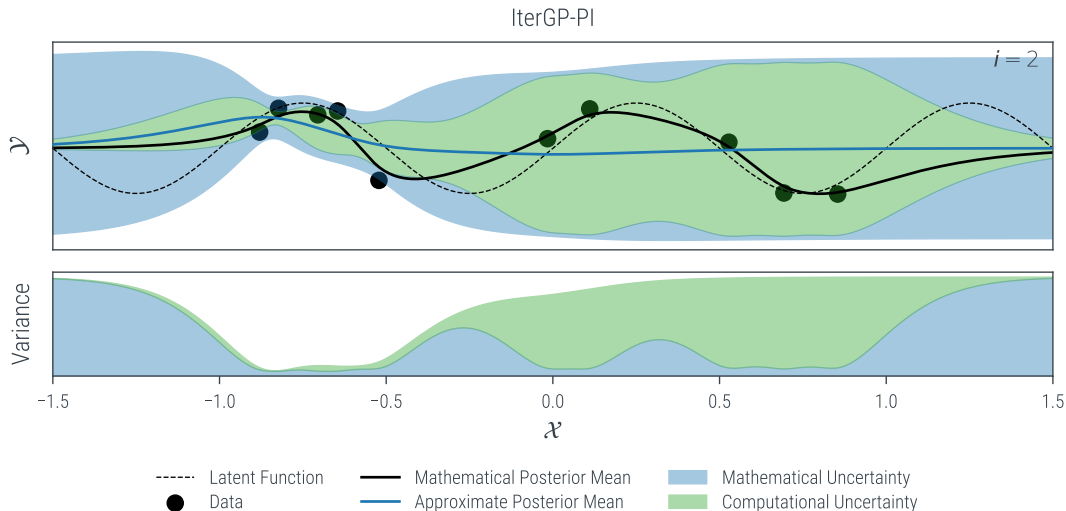
Interpreting computational and combined uncertainty as error quantification.



# Computation-Aware GP Inference Illustrated



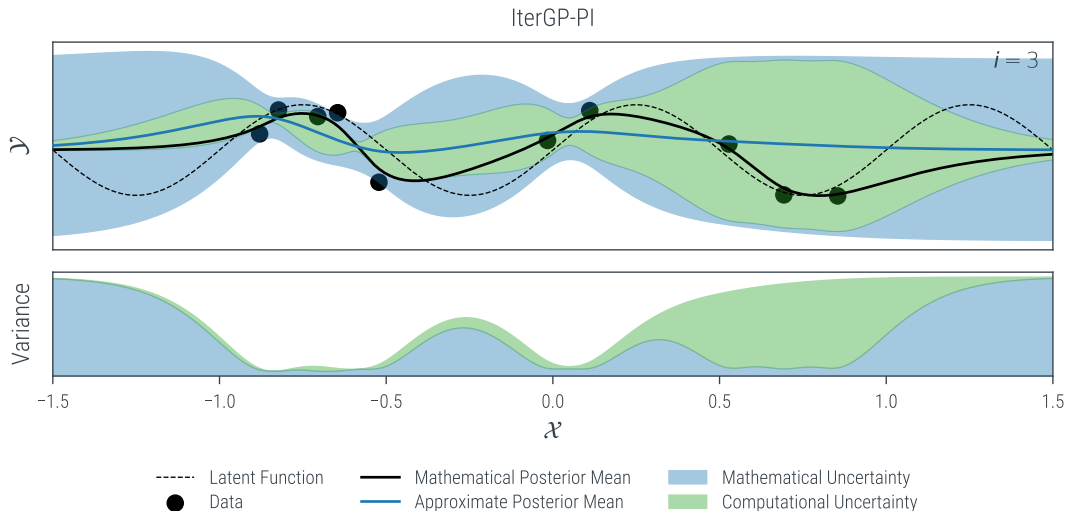
Interpreting computational and combined uncertainty as error quantification.



# Computation-Aware GP Inference Illustrated



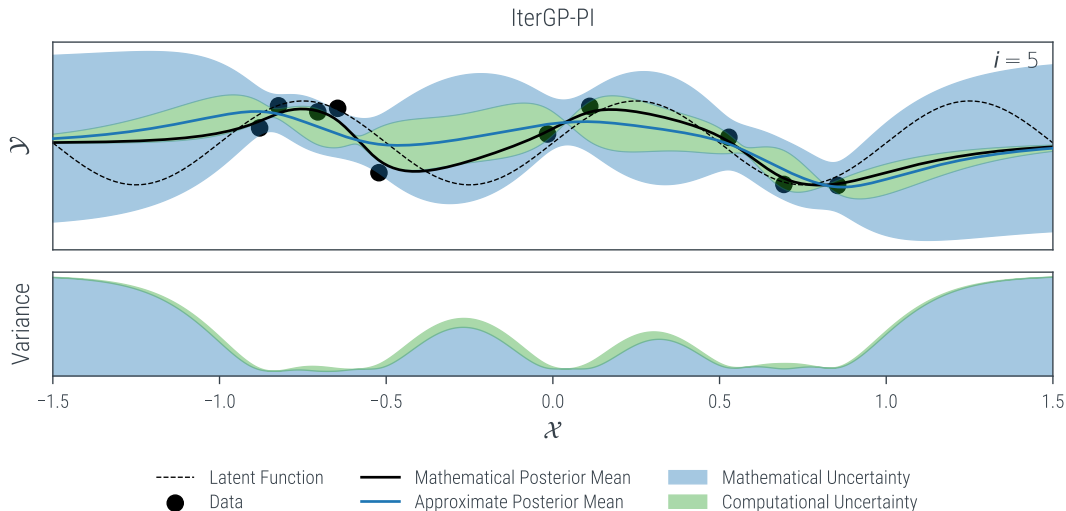
Interpreting computational and combined uncertainty as error quantification.



# Computation-Aware GP Inference Illustrated



Interpreting computational and combined uncertainty as error quantification.



		Time	Space
1	<b>procedure</b> ITERGP( $\mu, K, X, y, C_0 = \mathbf{0}$ )		
2	<b>while not</b> STOPPINGCRITERION() <b>do</b>		
3	$s_i \leftarrow \text{POLICY}()$	Select action via policy.	
4	$r_{i-1} \leftarrow (y - \mu) - \hat{K}v_{i-1}$	Residual.	$\mathcal{O}(n^2)$
5	$\alpha_i \leftarrow s_i^\top r_{i-1}$	Observation.	$\mathcal{O}(n)$
6	$z_i \leftarrow \hat{K}s_i$		$\mathcal{O}(1)$
7	$d_i \leftarrow \Sigma_{i-1}\hat{K}s_i = s_i - C_{i-1}z_i$	Search direction.	$\mathcal{O}(n)$
8	$\eta_i \leftarrow s_i^\top \hat{K}\Sigma_{i-1}\hat{K}s_i = z_i^\top d_i$		$\mathcal{O}(ni)$
9	$C_i \leftarrow C_{i-1} + \frac{1}{\eta_i}d_id_i^\top$	Precision matrix approx. $C_i \approx \hat{K}^{-1}$ .	$\mathcal{O}(n)$
10	$v_i \leftarrow v_{i-1} + \frac{\alpha_i}{\eta_i}d_i$	Representer weights estimate.	$\mathcal{O}(n)$
11	$\Sigma_i \leftarrow \Sigma_0 - C_i$	Representer weights uncertainty.	$\mathcal{O}(n)$
12	$\mu_i(\cdot) \leftarrow \mu(\cdot) + K(\cdot, X)v_i$	Approximate posterior mean.	$\mathcal{O}(n_\diamond n)$
13	$K_i(\cdot, \cdot) \leftarrow K(\cdot, \cdot) - K(\cdot, X)C_iK(X, \cdot)$	Combined covariance function.	$\mathcal{O}(n_\diamond ni)$
14	<b>return</b> $\mathcal{GP}(\mu_i, K_i)$		$\mathcal{O}(n_\diamond^2)$

Theorem (Kanagawa et al. [Kan+18])

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

## Theorem (Kanagawa et al. [Kan+18])

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

## Theorem (Wenger et al. [Wen+22])

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

## Theorem (Kanagawa et al. [Kan+18])

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

## Theorem (Wenger et al. [Wen+22])

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

Exact quantification of uncertainty from **limited data** and **limited computation**.



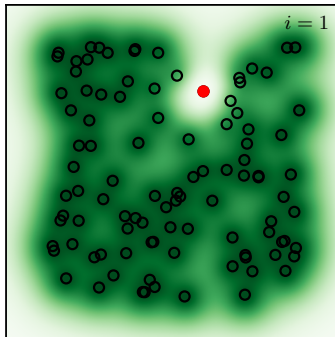
# Policy Choice and Connection to Other Approximations



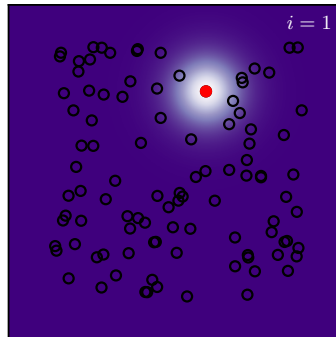
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $s_i$	Classic Analog
IterGP-Cholesky	$e_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{P}^{-1}r_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(X, z_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty



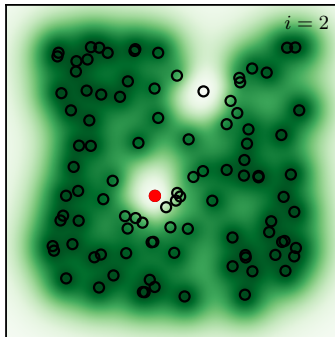
# Policy Choice and Connection to Other Approximations



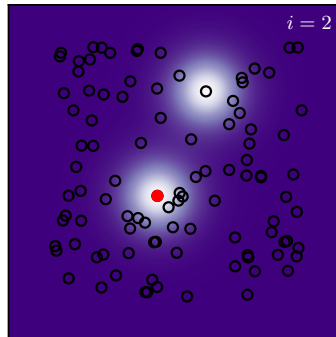
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $s_i$	Classic Analog
IterGP-Cholesky	$e_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{P}^{-1}r_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(X, z_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty



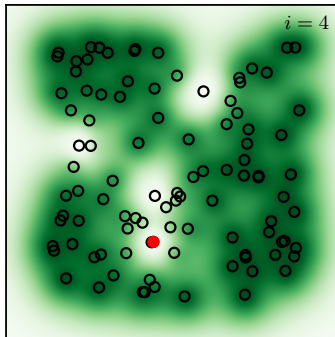
# Policy Choice and Connection to Other Approximations



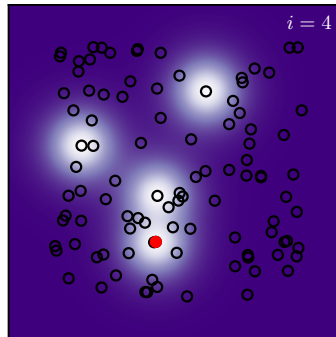
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $s_i$	Classic Analog
IterGP-Cholesky	$e_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{P}^{-1}r_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(X, z_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty



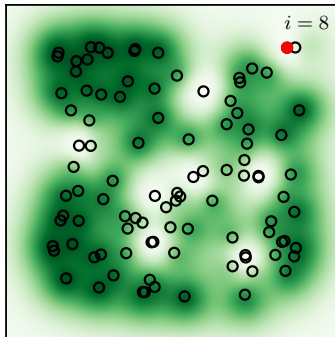
# Policy Choice and Connection to Other Approximations



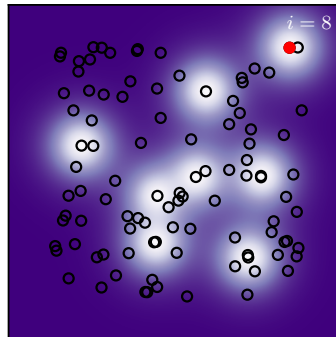
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $s_i$	Classic Analog
IterGP-Cholesky	$e_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{P}^{-1}r_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(X, z_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty



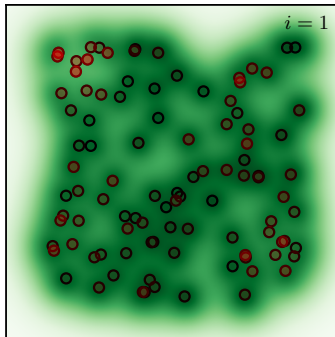
# Policy Choice and Connection to Other Approximations



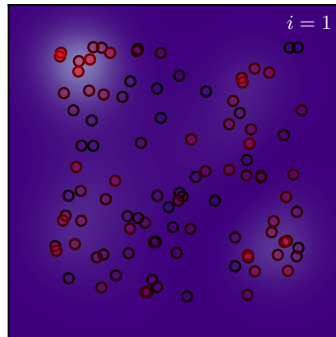
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $\mathbf{s}_i$	Classic Analog
IterGP-Cholesky	$\mathbf{e}_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{\mathbf{P}}^{-1}\mathbf{r}_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(\mathbf{X}, \mathbf{z}_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty



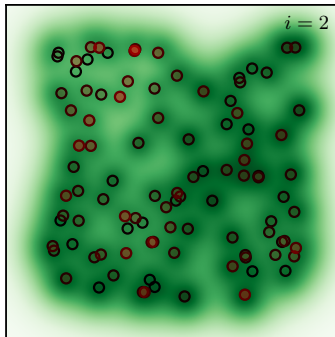
# Policy Choice and Connection to Other Approximations



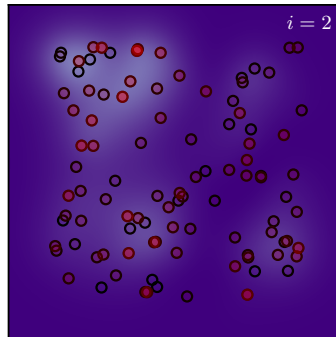
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $\mathbf{s}_i$	Classic Analog
IterGP-Cholesky	$\mathbf{e}_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{\mathbf{P}}^{-1}\mathbf{r}_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(\mathbf{X}, \mathbf{z}_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty



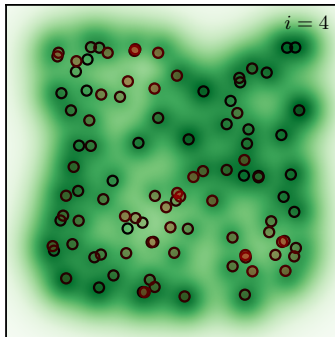
# Policy Choice and Connection to Other Approximations



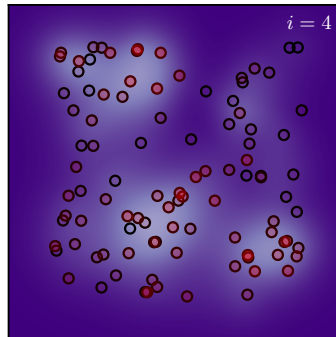
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $\mathbf{s}_i$	Classic Analog
IterGP-Cholesky	$\mathbf{e}_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{\mathbf{P}}^{-1}\mathbf{r}_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(\mathbf{X}, \mathbf{z}_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty



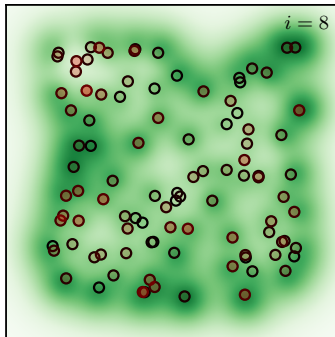
# Policy Choice and Connection to Other Approximations



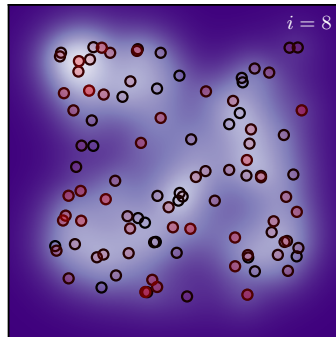
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $s_i$	Classic Analog
IterGP-Cholesky	$e_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{P}^{-1}r_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(X, z_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty





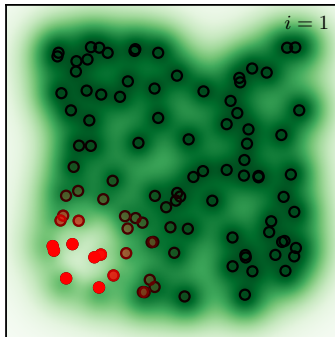
# Policy Choice and Connection to Other Approximations



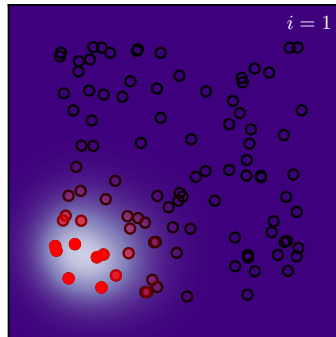
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $\mathbf{s}_i$	Classic Analog
IterGP-Cholesky	$\mathbf{e}_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{\mathbf{P}}^{-1} \mathbf{r}_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(\mathbf{X}, \mathbf{z}_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty



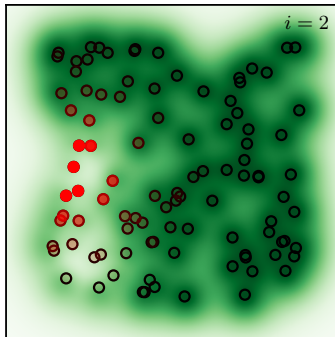
# Policy Choice and Connection to Other Approximations



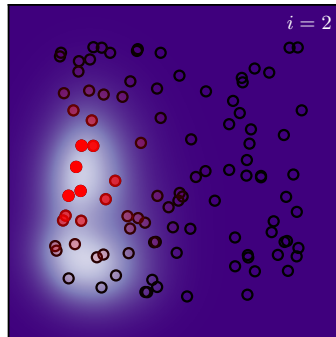
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $s_i$	Classic Analog
IterGP-Cholesky	$e_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{P}^{-1}r_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(X, z_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty



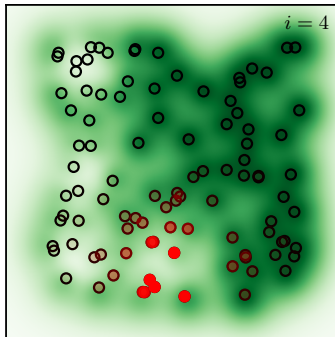
# Policy Choice and Connection to Other Approximations



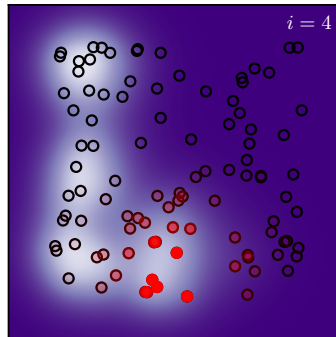
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $s_i$	Classic Analog
IterGP-Cholesky	$e_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{P}^{-1}r_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(X, z_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty



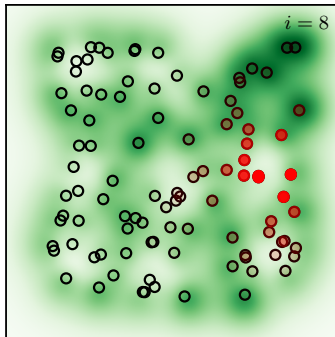
# Policy Choice and Connection to Other Approximations



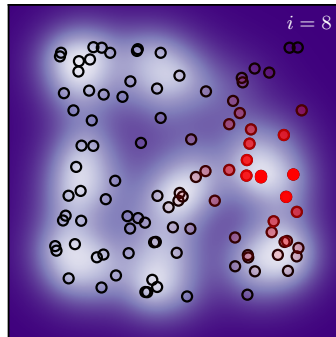
IterGP extends the most commonly used GP approximations to include computational uncertainty, with at most quadratic cost.

	Actions $s_i$	Classic Analog
IterGP-Cholesky	$e_i$	(Partial) Cholesky / Subset of data
IterGP-CG	$\hat{P}^{-1}r_i$	(Preconditioned) CG
IterGP-PseudoInput	$k(X, z_i)$	$\approx$ SVGP

Computational Uncertainty



Combined Uncertainty

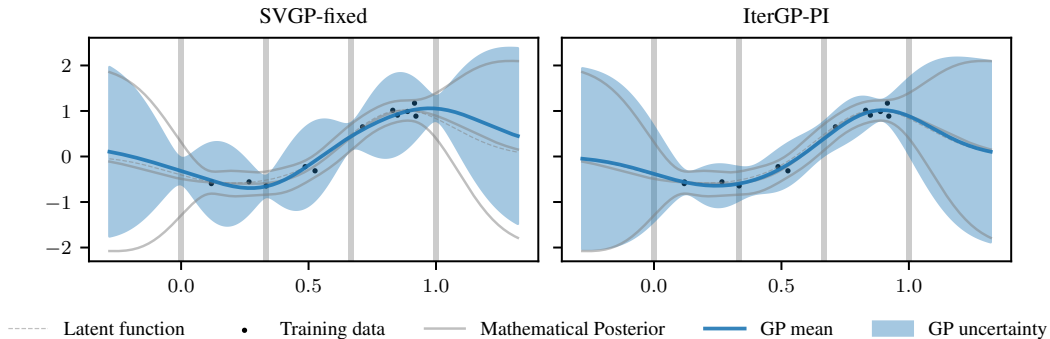


# SVGP versus IterGP-PI



Quantifying computational uncertainty improves generalization of inducing point methods like SVGP.

[Tit09; HFL13]

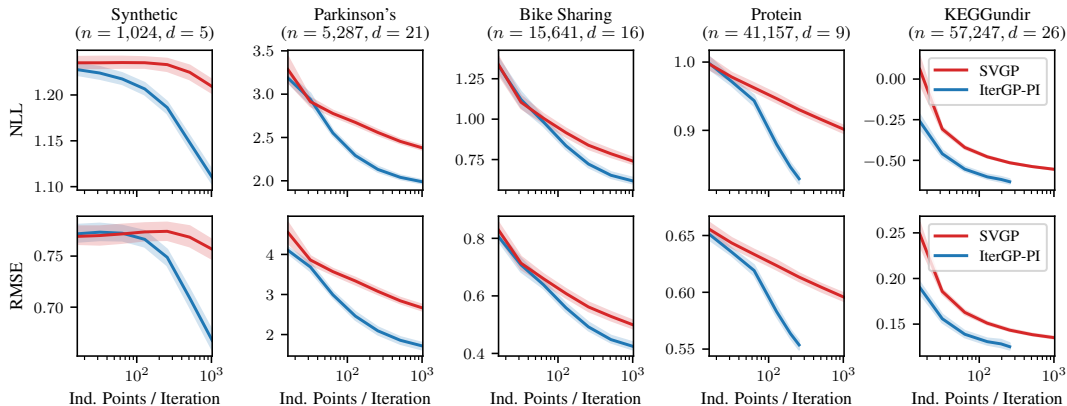


# SVGP versus IterGP-PI



Quantifying computational uncertainty improves generalization of inducing point methods like SVGP.

[Tit09; HFL13]

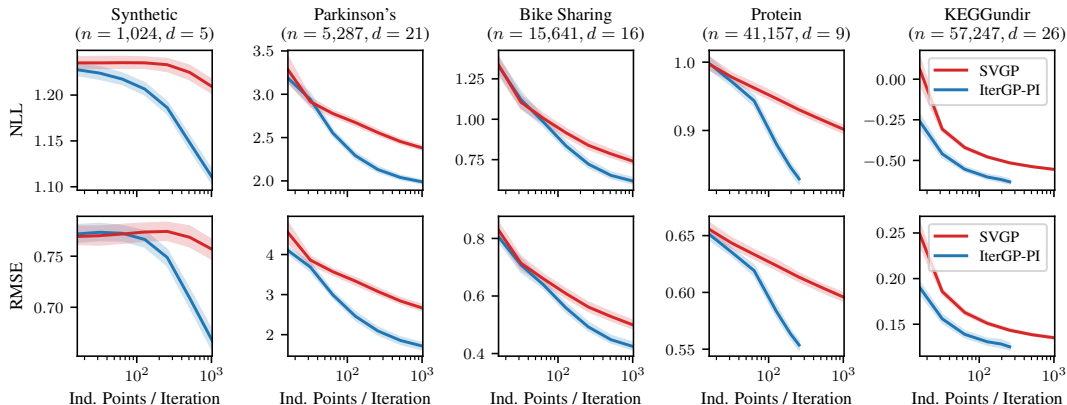


# SVGP versus IterGP-PI



Quantifying computational uncertainty improves generalization of inducing point methods like SVGP.

[Tit09; HFL13]



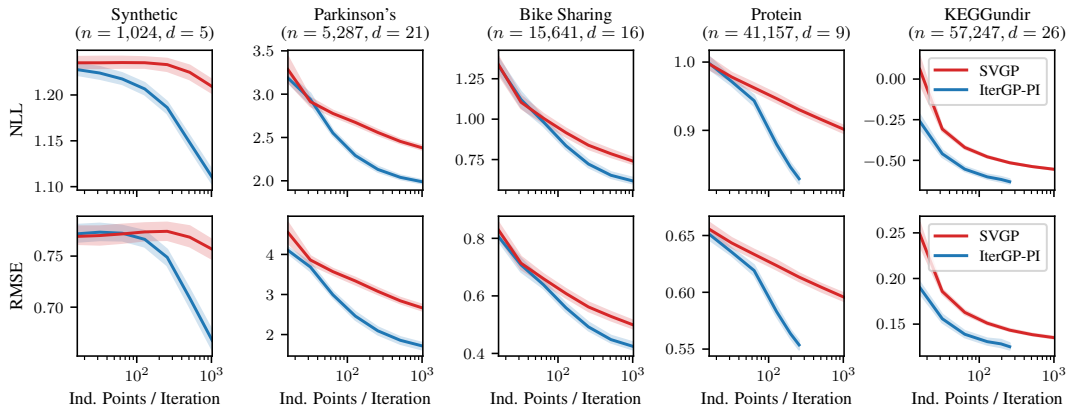
What about optimizing inducing point locations?

# SVGP versus IterGP-PI



Quantifying computational uncertainty improves generalization of inducing point methods like SVGP.

[Tit09; HFL13]



What about computational cost? SVGP:  $\mathcal{O}(ni^2)$  versus IterGP:  $\mathcal{O}(n^2i)$ .

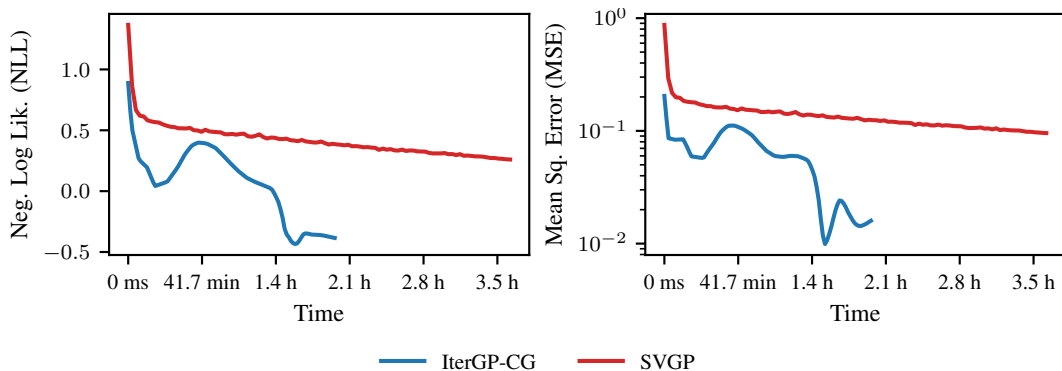


# Training Gaussian Processes on Large-Scale Data



Kernel hyperparameter optimization with SVGP and IterGP on a problem with  $n \approx 4 \cdot 10^5$  data points.

[Wen+24, Unpublished work]

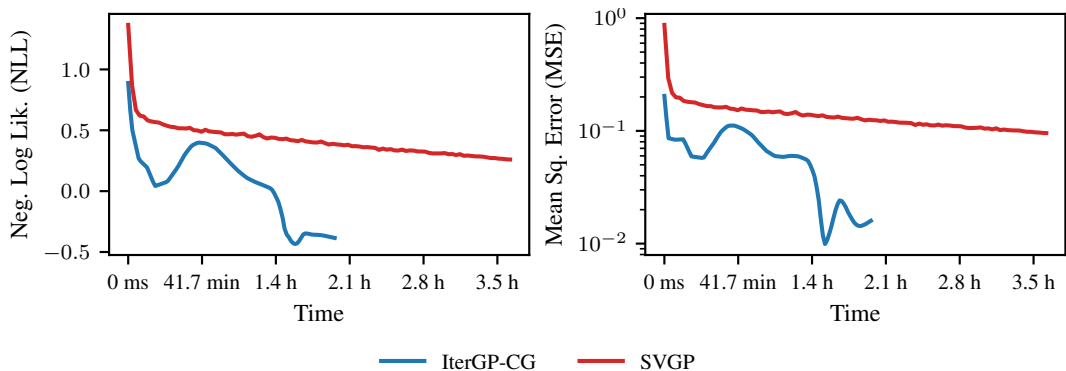


# Training Gaussian Processes on Large-Scale Data



Kernel hyperparameter optimization with SVGP and IterGP on a problem with  $n \approx 4 \cdot 10^5$  data points.

[Wen+24, Unpublished work]



**Faster large-scale Gaussian processes with better generalization!**

# Other Applications

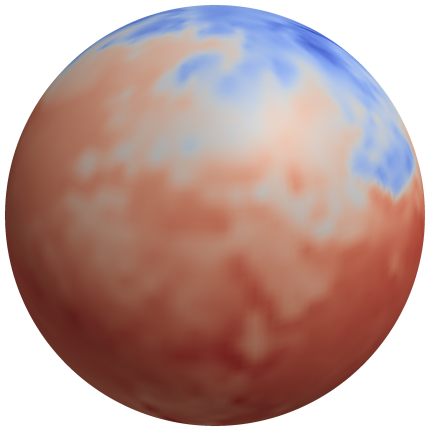
Extending these ideas beyond what we've seen.

# Spatiotemporal Modeling



Spatio-temporal regression of Earth surface temperature via computation-aware filtering and smoothing.

[Pfö+24, Unpublished Work]



(a) Prediction



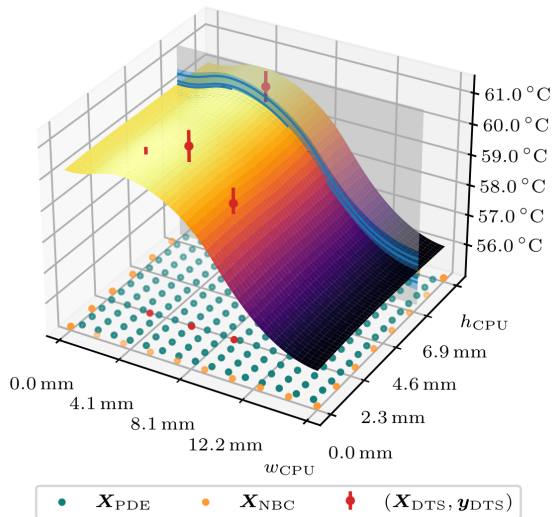
(b) Uncertainty

# Physics-Informed GP Regression

Learning to solve linear partial differential equations.



[Pfö+23]

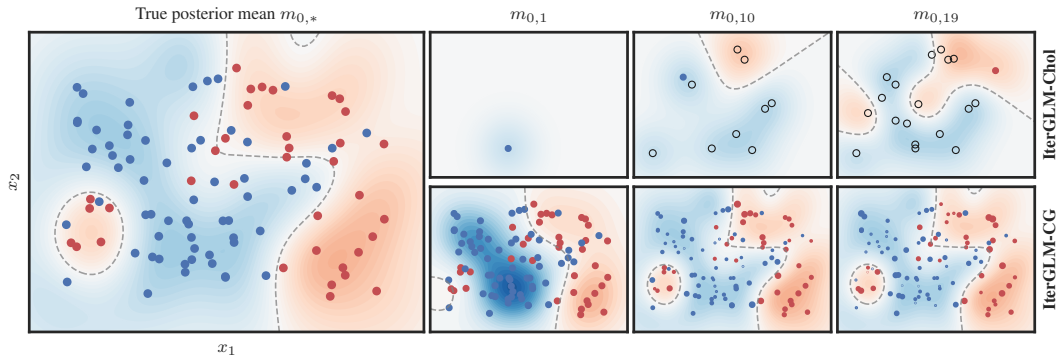


# Generalized Linear Models



Gaussian process classification with IterGLM using two different policies.

[Tat+23]



# Wrapping Up



## Summary

- ▶ Large-scale linear systems are ubiquitous in scientific computation.



## Summary

- ▶ Large-scale linear systems are ubiquitous in scientific computation.
- ▶ Probabilistic linear solvers allow us to explicitly trade off **speed** for **precision**.





## Summary

- ▶ Large-scale linear systems are ubiquitous in scientific computation.
- ▶ Probabilistic linear solvers allow us to explicitly trade off **speed for precision**.

## Application: Gaussian Processes

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



## Summary

- ▶ Large-scale linear systems are ubiquitous in scientific computation.
- ▶ Probabilistic linear solvers allow us to explicitly trade off **speed for precision**.

## Application: Gaussian Processes

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



## Summary

- ▶ Large-scale linear systems are ubiquitous in scientific computation.
- ▶ Probabilistic linear solvers allow us to explicitly trade off **speed** for **precision**.

## Application: Gaussian Processes

- ▶ Large-scale GP models are often as much about the approximation as they are about the data.
- ▶ We can exactly quantify the error from **finite data** and from the **approximation** via a combined uncertainty estimate.  $\Rightarrow$  IterGP



## Summary

- ▶ Large-scale linear systems are ubiquitous in scientific computation.
- ▶ Probabilistic linear solvers allow us to explicitly trade off **speed for precision**.

## Application: Gaussian Processes

- ▶ Large-scale GP models are often as much about the approximation as they are about the data.
- ▶ We can exactly quantify the error from **finite data** and from the **approximation** via a combined uncertainty estimate.  $\Rightarrow$  IterGP
- ▶ **Explicit trade-off** between computation and uncertainty via probabilistic linear solver.



## Summary

- ▶ Large-scale linear systems are ubiquitous in scientific computation.
- ▶ Probabilistic linear solvers allow us to explicitly trade off **speed** for **precision**.

## Application: Gaussian Processes

- ▶ Large-scale GP models are often as much about the approximation as they are about the data.
- ▶ We can exactly quantify the error from **finite data** and from the **approximation** via a combined uncertainty estimate.  $\Rightarrow$  IterGP
- ▶ **Explicit trade-off** between computation and uncertainty via probabilistic linear solver.

## Open Research Questions / Future Directions

- ▶ Calibration.
- ▶ Policy design for downstream tasks (Active learning, Bayesian optimization, ...).
- ▶ ...?



## 1 Introduction

## 2 Probabilistic Linear Solvers

### 2.1 Derivation

### 2.2 Policy Choice

### 2.3 Prior Choice

### 2.4 Algorithm

## 3 Application: Large-Scale Gaussian Processes

### 3.1 Gaussian Process Inference at Scale

### 3.2 Quantifying Approximation Error

### 3.3 Algorithm: IterGP

### 3.4 Theoretical Analysis

### 3.5 Policy Choice Illustrated

### 3.6 Experiments

## 4 Summary and Extensions

## 5 Additional Material

### 5.1 Calibration

### 5.2 An approximation method or a better model?

# Additional Material

**Observation:** Uncertainty quantification of probabilistic linear solvers can be conservative!

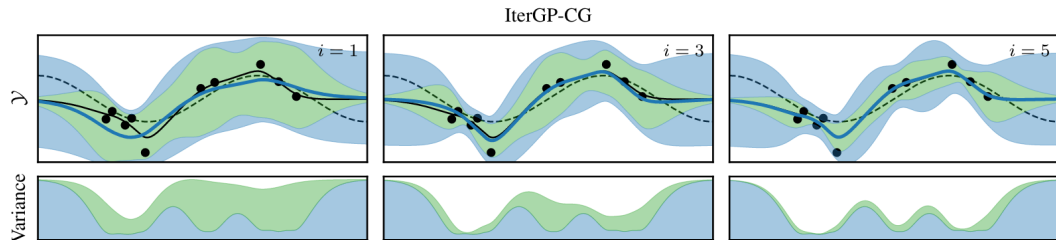


Figure: IterGP using a (conjugate) gradient policy.



**Observation:** Uncertainty quantification of probabilistic linear solvers can be conservative!

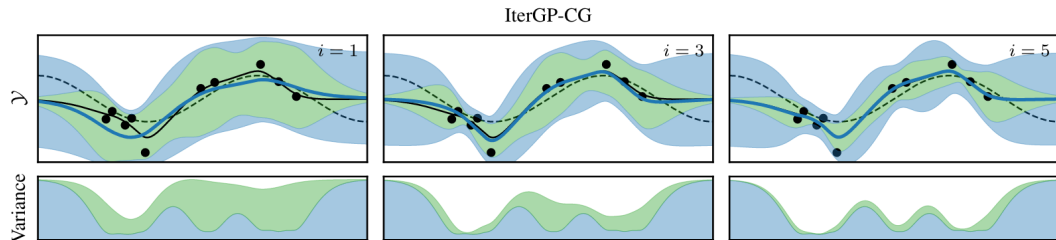


Figure: IterGP using a (conjugate) gradient policy.

**Why is that?** We conditioned on  $\alpha_i = \mathbf{s}_i^\top \mathbf{r}_{i-1} = \mathbf{s}_i^\top \mathbf{A}(\mathbf{x}_* - \mathbf{x}_{i-1})$ .

**But:** We've "cheated" for a gradient policy, since  $\mathbf{s}_i = \mathbf{b} - \mathbf{A}\mathbf{x}_{i-1} = \mathbf{A}(\mathbf{x}_* - \mathbf{x}_{i-1}) = \mathbf{s}_i(\mathbf{x}_*)$ .

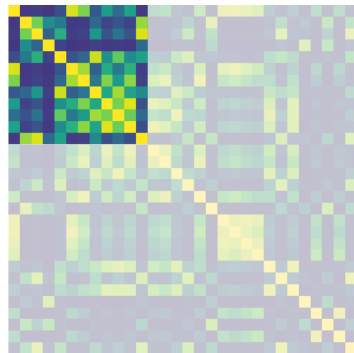
## Theorem (Online GP Approximation with IterGP)

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

$$\tilde{\mathbf{s}}_i = \begin{pmatrix} \mathbf{s}_i \\ \mathbf{0} \end{pmatrix} \quad (1)$$

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  $\tilde{\mathbf{s}}_i$ :

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



# An Approximation Method or a Better Model?



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

◀ ToC

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



# An Approximation Method or a Better Model?



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

◀ ToC

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



The distinction between data and computation vanishes.

# An Approximation Method or a Better Model?



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

◀ ToC

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



The distinction between data and computation vanishes.

What if we modeled this situation with a Gaussian process?

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

# An Approximation Method or a Better Model?



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

◀ ToC

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



The distinction between data and computation vanishes.

What if we modeled 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.



- ▶ E. Daxberger, A. Kristiadi, A. Immer, R. Eschenhagen, M. Bauer, and P. Hennig. *Laplace Redux – Effortless Bayesian Deep Learning*. 2022. DOI: [10.48550/arXiv.2106.14806](https://doi.org/10.48550/arXiv.2106.14806). URL: <http://arxiv.org/abs/2106.14806> (cit. on p. 7).
- ▶ P. Hennig, M. A. Osborne, and M. Girolami. “Probabilistic numerics and uncertainty in computations”. In: *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences* 471.2179 (2015) (cit. on pp. 9–12).
- ▶ J. Cockayne, C. J. Oates, T. J. Sullivan, and M. Girolami. “Bayesian Probabilistic Numerical Methods”. In: *SIAM Review* 61.4 (2019), pp. 756–789. DOI: [10.1137/17M1139357](https://doi.org/10.1137/17M1139357) (cit. on pp. 9–12).
- ▶ P. Hennig, M. A. Osborne, and H. P. Kersting. *Probabilistic Numerics: Computation as Machine Learning*. CUP, 2022. ISBN: 978-1-316-68141-1. DOI: [10.1017/9781316681411](https://doi.org/10.1017/9781316681411) (cit. on pp. 9–12).
- ▶ P. Hennig. “Probabilistic Interpretation of Linear Solvers”. In: *SIAM Journal on Optimization* 25.1 (2015), pp. 234–260 (cit. on pp. 9–21).



- ▶ J. Cockayne, C. J. Oates, I. C. Ipsen, and M. 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. 9–23, 28, 30, 31).
- ▶ J. Wenger and P. Hennig. “Probabilistic Linear Solvers for Machine Learning”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2020 (cit. on pp. 9–12).
- ▶ J. Wenger, G. Pleiss, M. Pförtner, P. Hennig, and J. P. Cunningham. “Posterior and Computational Uncertainty in Gaussian Processes”. In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2022 (cit. on pp. 28, 48–59, 66–68, 102).
- ▶ L. N. Trefethen and D. Bau. *Numerical Linear Algebra*. Society for Industrial and Applied Mathematics (SIAM), 1997 (cit. on p. 29).
- ▶ M. Kanagawa, P. Hennig, D. Sejdinovic, and B. K. Sriperumbudur. *Gaussian processes and kernel methods: A review on connections and equivalences*. 2018. arXiv: [1807.02582](https://arxiv.org/abs/1807.02582) (cit. on pp. 66–68).





- ▶ M. Titsias. “Variational learning of inducing variables in sparse Gaussian processes”. In: *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2009 (cit. on pp. 81–84).
- ▶ J. Hensman, N. Fusi, and N. D. Lawrence. “Gaussian processes for big data”. In: *Conference on Uncertainty in Artificial Intelligence (UAI)*. 2013 (cit. on pp. 81–84).
- ▶ M. Pförtner, I. Steinwart, P. Hennig, and J. Wenger. *Physics-Informed Gaussian Process Regression Generalizes Linear PDE Solvers*. 2023. DOI: [10.48550/arXiv.2212.12474](https://arxiv.org/abs/2212.12474). URL: <http://arxiv.org/abs/2212.12474> (cit. on p. 89).
- ▶ L. Tatzel, J. Wenger, F. Schneider, and P. Hennig. *Accelerating Generalized Linear Models by Trading off Computation for Uncertainty*. 2023. DOI: [10.48550/arXiv.2310.20285](https://arxiv.org/abs/2310.20285). URL: <http://arxiv.org/abs/2310.20285> (cit. on p. 90).