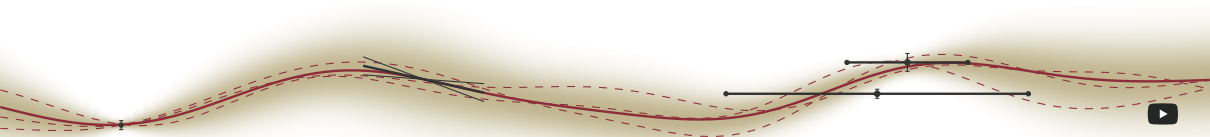# Numerics of Machine Learning
## Lecture 04
## Computation-Aware GP Inference

Jonathan Wenger

10 November 2022

**EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN**

Faculty of Science
Department of Computer Science

Chair for the Methods of Machine Learning

Where are we in the course?

▶ Last week: Contemporary way of solving linear systems for GP regression on large datasets

▶ This week: Probabilistic numerics approach to (approximate) GP regression

Today

▶ **Learning to approximate GPs** with probabilistic numerics.

▶ Quantifying approximation error probabilistically.

▶ Iterative numerical methods for GPs as active learning agents.

▶ Philosophical connections between data and computation.

▶ **Exact uncertainty quantification for GPs in (sub-)quadratic time.**

Recap:
Scalable GP Approximations

**Goal:** Learn an unknown function $f_* : \mathbb{R}^d \rightarrow \mathbb{R}$ from a training dataset of example input-output pairs.
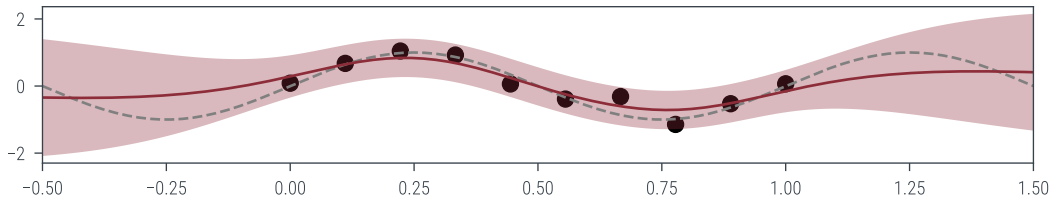
$$f \sim \mathcal{GP}(\mu, k)$$
$$y \mid f(X) \sim \mathcal{N}(f(X), \sigma^2 I)$$
$$f \mid X, y \sim \mathcal{GP}(\mu_{\text{post}}, k_{\text{post}})$$

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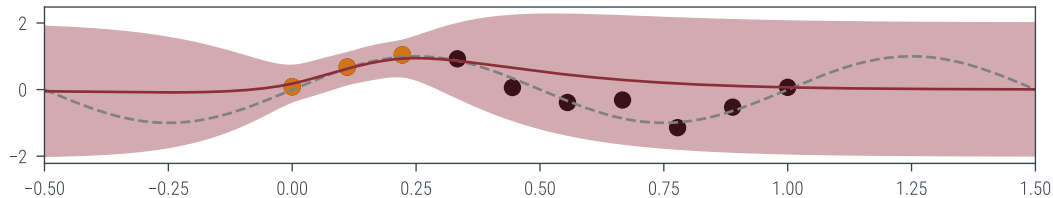


3

$$f \sim \mathcal{GP}(\mu, k)$$
$$y \mid f(X) \sim \mathcal{N}(f(X), \sigma^2 I)$$
$$f \mid X, y \sim \mathcal{GP}(\mu_{\text{post}}, k_{\text{post}})$$

$$\mu_{\text{post}}(x) = \mu(x) + k(x, X)C_i(y - \mu(X))$$
$$k_{\text{post}}(x_0, x_1) = k(x_0, x_1) - k(x_0, X)C_i k(X, x_1)$$

# Recap: Learning to Invert the Kernel Matrix

The Cholesky decomposition as a learning algorithm for the inverse kernel matrix.

UNIVERSITÄT
TÜBINGEN
EBERHARD KARLS

## Algorithm Cholesky with Inverse Approximation

Input: spd matrix $A$
Output: lower triangular $L_i$, s.t. $L_i L_i^\mathsf{T} \approx A$, low-rank $C_i \approx A^{-1}$

1   procedure CHOLESKY($A$)
2     $A' \leftarrow A, C_0 = 0$
3     for $i \in \{1, \ldots, n\}$ do
4       $s_i \leftarrow e_i$                  // Action
5       $d_i \leftarrow (I - C_{i-1}A)s_i$
6       $\eta_i \leftarrow s_i^\mathsf{T} A d_i = e_i^\mathsf{T} A' e_i = \|e_i\|_{A'}^2$,    // Norm. constant
7       $l_i \leftarrow A \frac{1}{\sqrt{\eta_i}} d_i$           // Matrix observation
8       $C_i \leftarrow C_{i-1} + \frac{1}{\eta_i} d_i d_i^\mathsf{T}$      // Inverse estimate
9       $A' \leftarrow A - L_i L_i^\mathsf{T} = A(A^{-1} - C_i)A = A(I - C_i A)$
10      $L_i = \begin{pmatrix} L_{i-1} & l_i \end{pmatrix}$
11     end for
12     return $L_i, C_i$
13 end procedure

Goal: (Low-rank) Approximation $C_i \approx A^{-1}$

Observation: Matrix approx. $\rightarrow$ inverse approx.?

$$L_i L_i^\mathsf{T} \approx A$$

$$\underbrace{(A^{-1}L_i)(A^{-1}L_i)^\mathsf{T}}_{=C_i} \approx A^{-1}$$
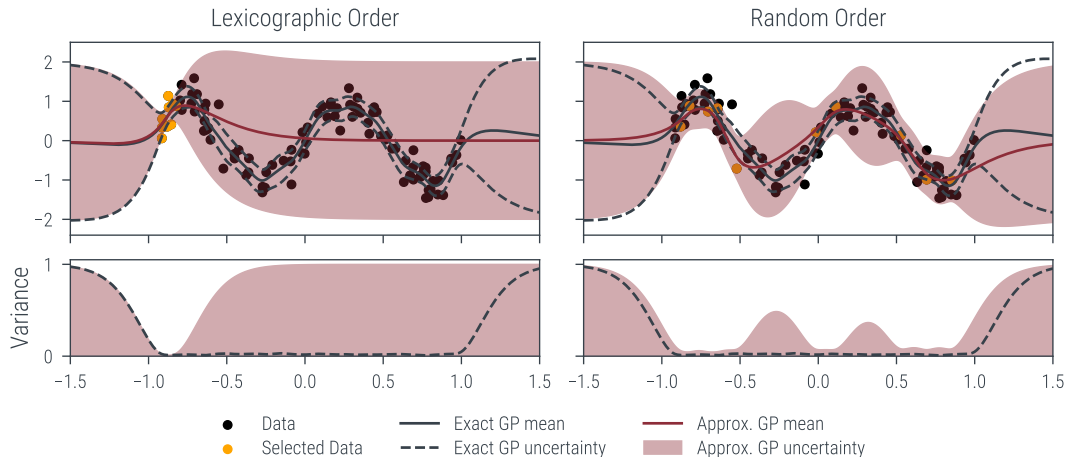
Computational complexity: #flops $\in \mathcal{O}(in^2)$

Cholesky can be seen as an iterative learning algorithm for the kernel matrix and its inverse.

Lexicographic Order — Random Order

Variance

| Data | Exact GP mean |
| Selected Data | Exact GP uncertainty |
| Approx. GP mean | Approx. GP uncertainty |

The selection of datapoints, i.e. choice of actions $s_i$, matters a lot for convergence.

## Partial Cholesky

$$A'e_i = A(I - C_{i-1}A)s_i = Ad_i$$

## Other Method?

$$A'e_i = A(I - C_{i-1}A)s_i = Ad_i$$

Can we learn the kernel matrix (inverse) in a more efficient way via different actions?

# Recap: Method of Conjugate Gradients

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

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

**Goal:** Approximately solve linear system $Ax = b$ with few matrix-vector multiplies.
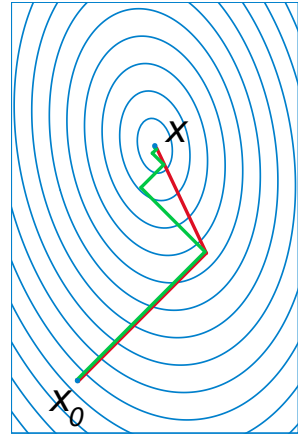
**Idea:** Rephrase as quadratic optimization problem and optimize. Let

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

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. $\langle d_i^\mathsf{T} d_j \rangle_A = d_i^\mathsf{T} A d_j = 0$ for $i \neq j$.
   $\implies$ convergence in at most $n$ steps.
3. Conjugate gradient method: First step $d_0 = r(x_0)$.



Oleg Alexandrov, com-
mons.wikimedia.org/w/index.php?curid=2267598

# Recap: Algorithm: Method of Conjugate Gradients

We can interpret CG as a learning algorithm for the matrix inverse as well.
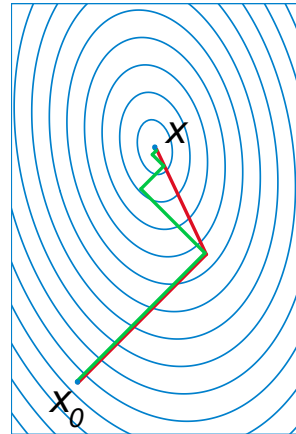
## Algorithm CG with Inverse Approximation

**Input:** spd matrix $A$, vector $b$, initial guess $x_0$
**Output:** approximate solution $x_i \approx A^{-1}b$, low-rank $C_i \approx A^{-1}$

1  **procedure** $CG(A, b, x_0)$
2    **while** $\|r_i\|_2 > \max(\delta_{\text{rtol}}\|b\|_2, \delta_{\text{atol}})$ **do**
3      $r_{i-1} \leftarrow b - Ax_{i-1}$      // Residual
4      $s_i \leftarrow r_{i-1}$      // Action
5      $\alpha_i \leftarrow s_i^\mathsf{T} r_{i-1}$      // Observation
6      $d_i \leftarrow (I - C_{i-1}A)s_i$      // Search direction
7      $\eta_i \leftarrow s_i^\mathsf{T} A d_i = d_i^\mathsf{T} A d_i$      // Norm. constant
8      $C_i \leftarrow C_{i-1} + \frac{1}{\eta_i} d_i d_i^\mathsf{T}$      // Inverse estimate
9      $x_i \leftarrow x_{i-1} + \frac{\alpha_i}{\eta_i} d_i = C_i b$      // Solution estimate
10    **end while**
11    **return** $x_i, C_i$
12 **end procedure**



Oleg Alexandrov, com-
mons.wikimedia.org/w/index.php?curid=2267598

# Recap: Algorithm: Method of Conjugate Gradients

We can interpret CG as a learning algorithm for the matrix inverse as well.

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

## Algorithm CG with Inverse Approximation

**Input:** spd matrix $A$, vector $b$, initial guess $x_0$
**Output:** approximate solution $x_i \approx A^{-1}b$, low-rank $C_i \approx A^{-1}$

1  **procedure** CG($A, b, x_0$)
2      **while** $\|r_i\|_2 > \max(\delta_{\text{rtol}}\|b\|_2, \delta_{\text{atol}})$ **do**
3          $r_{i-1} \leftarrow b - Ax_{i-1}$                    // Residual
4          $s_i \leftarrow r_{i-1}$                              // Action
5          $\alpha_i \leftarrow s_i^\mathsf{T} r_{i-1}$          // Observation
6          $d_i \leftarrow (I - C_{i-1}A)s_i$                    // Search direction
7          $\eta_i \leftarrow s_i^\mathsf{T} A d_i = d_i^\mathsf{T} A d_i$   // Norm. constant
8          $C_i \leftarrow C_{i-1} + \frac{1}{\eta_i}d_i d_i^\mathsf{T}$     // Inverse estimate
9          $x_i \leftarrow x_{i-1} + \frac{\alpha_i}{\eta_i}d_i = C_i b$      // Solution estimate
10     **end while**
11     **return** $x_i, C_i$
12  **end procedure**

## Algorithm Cholesky with Inverse Approximation

**Input:** spd matrix $A$
**Output:** lower triangular $L_i$, s.t. $L_i L_i^\mathsf{T} \approx A$, low-rank $C_i \approx A^{-1}$

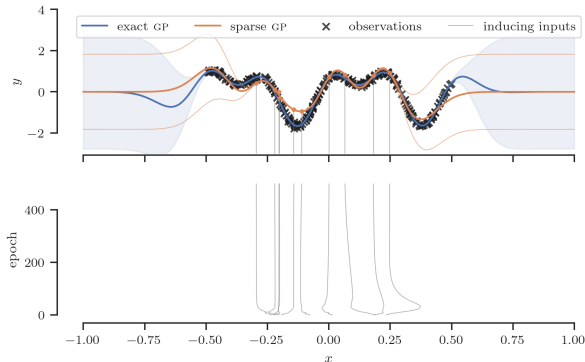1  **procedure** CHOLESKY($A$)
2      $A' \leftarrow A, C_0 = 0$
3      **for** $i \in \{1, \ldots, n\}$ **do**
4          $s_i \leftarrow e_i$                                  // Action
5          $d_i \leftarrow (I - C_{i-1}A)s_i$
6          $\eta_i \leftarrow s_i^\mathsf{T} A d_i = e_i^\mathsf{T} A' e_i = \|e_i\|_{A'}^2$   // Norm. constant
7          $l_i \leftarrow A\frac{1}{\sqrt{\eta_i}}d_i$           // Matrix observation
8          $C_i \leftarrow C_{i-1} + \frac{1}{\eta_i}d_i d_i^\mathsf{T}$   // Inverse estimate
9          $A' \leftarrow A - L_i L_i^\mathsf{T} = A(A^{-1} - C_i)A = A(I - C_i A)$
10         $L_i = (L_{i-1} \quad l_i)$
11     **end for**
12     **return** $L_i, C_i$
13  **end procedure**

# Recap: Stochastic Variational Gaussian Processes

**Idea**: Linear time GP approximation via inducing points.



Source: https://tiao.io/post/sparse-variational-gaussian-processes/

Can we design a method where we can trust the UQ *no matter how much computation we've done*?

Exact UQ for GP approximation with arbitrary amounts of compute.
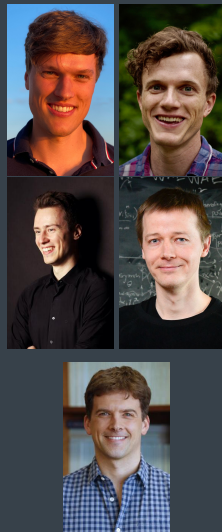Computation-aware GP Inference

## Posterior and Computational Uncertainty in Gaussian Processes

Jonathan Wenger, Geoff Pleiss, Marvin Pförtner, Philipp Hennig and John Cunningham

▶ IterGP: new class of GP approximations *accounting for computational uncertainty*.

▶ IterGP instances extend classic methods (Cholesky, CG, Nyström, …).

▶ Strong theoretical guarantees.

▶ Modeling computational uncertainty either *saves computation* or *improves generalization*.

Paper   arXiv   `https://arxiv.org/abs/2107.00243`

Implementation   🐙   `https://github.com/JonathanWenger/itergp`

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

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

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

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

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

$$= \mu(x) + k(x, X) \underbrace{v_*}_{\text{representer weights}} = \mu(x) + \sum_{j=1}^{n} k(x, x_j)(v_*)_j$$



— Exact GP Mean  ● Data  — Kernel Function(s) × Representer Weights

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

**Observation:** Iterative linear solvers (e.g. CG) approximate the representer weights $v_i \approx v_* = \hat{K}^{-1} y$.

$$\mu_{\text{post}}(x) = \mu(x) + k(x, X) \underbrace{v_*}_{\text{representer weights}} = \mu(x) + \sum_{j=1}^{n} k(x, x_j)(v_*)_j$$

$$\approx \mu(x) + k(x, X) v_i$$



| — Exact GP Mean | ● Data | — Kernel Function(s) × Approx. Representer Weights |

Can we quantify the approximation error $\|v_* - v_i\|$ in the representer weights *probabilistically*?

# Interlude: Gaussians provide the **Linear Algebra of Inference**

If all joints are Gaussian and all observations are linear, all posteriors are Gaussian.

▶ products of Gaussians are Gaussians

$$\mathcal{N}(x; a, A)\mathcal{N}(x; b, B)$$
$$= \mathcal{N}(x; c, C)\mathcal{N}(a; b, A + B)$$
$$C := (A^{-1} + B^{-1})^{-1} \quad c := C(A^{-1}a + B^{-1}b)$$

▶ linear projections of Gaussians are Gaussians

$$p(z) = \mathcal{N}(z; \mu, \Sigma)$$
$$\Rightarrow \quad p(Az) = \mathcal{N}(Az, A\mu, A\Sigma A^{\mathsf{T}})$$

▶ marginals of Gaussians are Gaussians

$$\int \mathcal{N}\left[\begin{pmatrix} x \\ y \end{pmatrix}; \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix}\right] dy = \mathcal{N}(x; \mu_x, \Sigma_{xx})$$

▶ (linear) conditionals of Gaussians are Gaussians

$$p(x \mid y) = \frac{p(x, y)}{p(y)}$$
$$= \mathcal{N}\left(x; \mu_x + \Sigma_{xy}\Sigma_{yy}^{-1}(y - \mu_y), \Sigma_{xx} - \Sigma_{xy}\Sigma_{yy}^{-1}\Sigma_{yx}\right)$$

### Bayesian inference becomes linear algebra

If $p(x) = \mathcal{N}(x; \mu, \Sigma)$ and $p(y \mid x) = \mathcal{N}(y; A^{\mathsf{T}}x + b, \Lambda)$, then

$$p(B^{\mathsf{T}}x + c \mid y) = \mathcal{N}[B^{\mathsf{T}}x + c; B^{\mathsf{T}}\mu + c + B^{\mathsf{T}}\Sigma A(A^{\mathsf{T}}\Sigma A + \Lambda)^{-1}(y - A^{\mathsf{T}}\mu - b), B^{\mathsf{T}}\Sigma B - B^{\mathsf{T}}\Sigma A(A^{\mathsf{T}}\Sigma A + \Lambda)^{-1}A^{\mathsf{T}}\Sigma B]$$

# Probabilistic Linear Solvers

Learning the solution of linear systems while quantifying computational uncertainty.

UNIVERSITÄT
TÜBINGEN
EBERHARD KARLS

**Goal:** Quantify approximation error when solving $v_* = \hat{K}^{-1}y$ probabilistically, i.e. $v_* \sim \mathcal{N}(v_i, \Sigma_i)$.

**Prior:** $v_* \sim \mathcal{N}(v_0, \Sigma_0)$
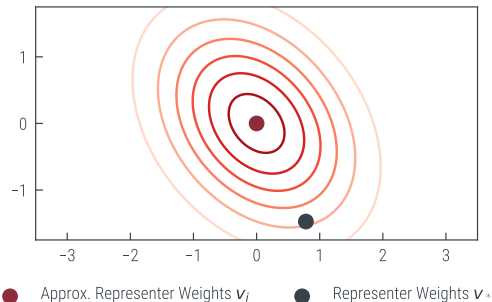


● Approx. Representer Weights $v_i$          ● Representer Weights $v_*$

# Probabilistic Linear Solvers

Learning the solution of linear systems while quantifying computational uncertainty.

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

**Goal:** Quantify approximation error when solving $v_* = \hat{K}^{-1}y$ probabilistically, i.e. $v_* \sim \mathcal{N}(v_i, \Sigma_i)$.

**Prior:** $v_* \sim \mathcal{N}(v_0, \Sigma_0)$

**Likelihood:** Observe representer weights indirectly via matrix-vector multiplication with the residual:

$$\alpha_i := s_i^\mathsf{T} r_{i-1} = s_i^\mathsf{T}((y - \mu) - \hat{K}v_{i-1}) = s_i^\mathsf{T}\hat{K}(v_* - v_{i-1})$$



● Approx. Representer Weights $v_i$   ● Representer Weights $v_*$

# Probabilistic Linear Solvers

Learning the solution of linear systems while quantifying computational uncertainty.

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

**Goal:** Quantify approximation error when solving $v_* = \hat{K}^{-1}y$ probabilistically, i.e. $v_* \sim \mathcal{N}(v_i, \Sigma_i)$.

**Prior:** $v_* \sim \mathcal{N}(v_0, \Sigma_0)$

**Likelihood:** Observe representer weights indirectly via matrix-vector multiplication with the residual:

$$\alpha_i := s_i^\mathsf{T} r_{i-1} = s_i^\mathsf{T}((y - \mu) - \hat{K}v_{i-1}) = s_i^\mathsf{T}\hat{K}(v_* - v_{i-1})$$

**Posterior:** Affine Gaussian inference!



●  Approx. Representer Weights $v_i$          ●  Representer Weights $v_*$

# Probabilistic Linear Solvers

Learning the solution of linear systems while quantifying computational uncertainty.

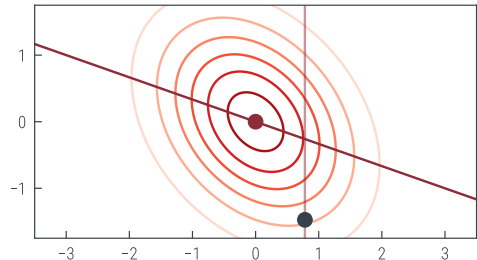**Goal:** Quantify approximation error when solving $v_* = \hat{K}^{-1} y$ probabilistically, i.e. $v_* \sim \mathcal{N}(v_i, \boldsymbol{\Sigma}_i)$.
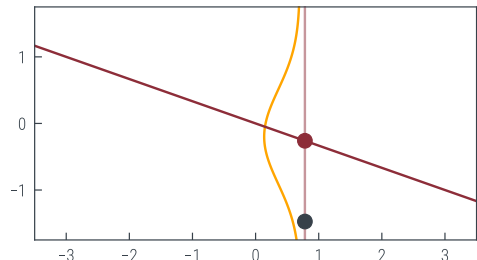
**Prior:** $v_* \sim \mathcal{N}(v_0, \boldsymbol{\Sigma}_0)$

**Likelihood:** Observe representer weights indirectly via matrix-vector multiplication with the residual:

$$\alpha_i := s_i^\mathsf{T} r_{i-1} = s_i^\mathsf{T}((y - \boldsymbol{\mu}) - \hat{K} v_{i-1}) = s_i^\mathsf{T} \hat{K}(v_* - v_{i-1})$$

**Posterior:** $v_* \mid \alpha_i \sim \mathcal{N}(v_i, \boldsymbol{\Sigma}_i)$, where

$$v_i = v_{i-1} + \overbrace{\boldsymbol{\Sigma}_{i-1} \hat{K} s_i}^{=:d_i} \overbrace{(s_i^\mathsf{T} \hat{K} \boldsymbol{\Sigma}_{i-1} \hat{K} s_i)^{-1}}^{=:\eta_i} \overbrace{s_i^\mathsf{T} \hat{K}(v_* - v_{i-1})}^{=\alpha_i}$$

$$= C_i(y - \boldsymbol{\mu})$$

$$\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}_{i-1} - \overbrace{\boldsymbol{\Sigma}_{i-1} \hat{K} s_i}^{=d_i} \overbrace{(s_i^\mathsf{T} \hat{K} \boldsymbol{\Sigma}_{i-1} \hat{K} s_i)^{-1}}^{=\eta_i} \overbrace{s_i^\mathsf{T} \hat{K} \boldsymbol{\Sigma}_{i-1}}^{=d_i^\mathsf{T}}$$

$$= \boldsymbol{\Sigma}_0 - C_i = \boldsymbol{\Sigma}_0 - \sum_{j=1}^{i} \frac{1}{\eta_j} d_j d_j^\mathsf{T} = \boldsymbol{\Sigma}_0 - S_i(S_i^\mathsf{T} \hat{K} S_i)^{-1} S_i^\mathsf{T}$$



● Approx. Representer Weights $v_i$ ● Representer Weights $v_*$

**Problem**: How to choose the linear solver prior?



● Approx. Representer Weights $v_i$     ● Representer Weights $v_*$

**Problem**: How to choose the linear solver prior?

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

Gaussian process prior $f \sim \mathcal{GP}(0, k)$ gives:

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



● Approx. Representer Weights $v_i$    ● Representer Weights $v_*$

**Problem**: How to choose the linear solver prior?

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

Gaussian process prior $f \sim \mathcal{GP}(0, k)$ gives:

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

$$v_* = \hat{K}^{-1} y \sim \mathcal{N}(0, \hat{K}^{-1}) = \mathcal{N}(v_0, \Sigma_0)$$



● Approx. Representer Weights $v_i$        ● Representer Weights $v_*$

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

### Gaussian Processes

Mathematical posterior: $f_\diamond \mid v_* \sim \mathcal{N}(\mu_*(X_\diamond), k_*(X_\diamond, X_\diamond))$, s.t.

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

### Learning the Representer Weights

Infer representer weights via probabilistic linear solver: $p(v_*) = \mathcal{N}(v_*; v_i, \Sigma_i)$, s.t.

$$v_i = C_i(y - \mu) \quad \text{and} \quad \Sigma_i = \Sigma_0 - C_i = \hat{K}^{-1} - C_i$$

### Combined Uncertainty

Marginal distribution: $p(f_\diamond) = \int p(f_\diamond \mid v_*)p(v_*) \, dv_* = \mathcal{N}(f_\diamond; \mu_i(X_\diamond), k_i(X_\diamond, X_\diamond))$, s.t.

$$\mu_i(\cdot) = \mu(\cdot) + k(\cdot, X)v_i$$

$$k_i(\cdot, \cdot) = \underbrace{k(\cdot, \cdot) - k(\cdot, X)\hat{K}^{-1}k(X, \cdot)}_{\text{mathematical uncertainty} \;\bigcirc} + \underbrace{k(\cdot, X)\Sigma_i k(X, \cdot)}_{\text{computational uncertainty} \;\bigcirc} = \underbrace{k(\cdot, \cdot) - k(\cdot, X)C_i k(X, \cdot)}_{\text{combined uncertainty} \;\bigcirc}$$

The covariance can be interpreted as a squared error.

## Combined Uncertainty

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

$$\mu_i(\cdot) = \mu(\cdot) + k(\cdot, X)v_i$$

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

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

### Combined Uncertainty

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

$$\mu_i(\cdot) = \mu(\cdot) + k(\cdot, X)v_i$$

$$k_i(\cdot, \cdot) = \underbrace{k(\cdot, \cdot) - k(\cdot, X)\hat{K}^{-1}k(X, \cdot)}_{\text{mathematical uncertainty} \, \bigcirc} + \underbrace{k(\cdot, X)\Sigma_i k(X, \cdot)}_{\text{computational uncertainty} \, \bigcirc} = \underbrace{k(\cdot, \cdot) - k(\cdot, X)C_i k(X, \cdot)}_{\text{combined uncertainty} \, \bigcirc}$$

**Remember**: $\mathrm{Cov}(f(x), f(x)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2] = \mathbb{E}[(f(x) - \mu(x))^2]$

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

### Combined Uncertainty

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

$$\mu_i(\cdot) = \mu(\cdot) + k(\cdot, X)v_i$$

$$k_i(\cdot, \cdot) = \underbrace{k(\cdot, \cdot) - k(\cdot, X)\hat{K}^{-1}k(X, \cdot)}_{\text{mathematical uncertainty} \; \bullet} + \underbrace{k(\cdot, X)\Sigma_i k(X, \cdot)}_{\text{computational uncertainty} \; \bullet} = \underbrace{k(\cdot, \cdot) - k(\cdot, X)C_i k(X, \cdot)}_{\text{combined uncertainty} \; \bullet}$$

**Remember**: $\text{Cov}(f(x), f(x)) = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2] = \mathbb{E}[(f(x) - \mu(x))^2]$

$$k_{\text{post}}(x, x) = \underbrace{k(x, x) - k(x, X)\hat{K}^{-1}k(X, x)}_{\text{mathematical uncertainty} \; \bullet} = \mathbb{E}[(f(x) - \mu_*(x))^2]$$

$$k_i^{\text{comp}}(x, x) = \underbrace{k(x, X)\Sigma_i k(X, x)}_{\text{computational uncertainty} \; \bullet} \underset{\Sigma_i = \text{Cov}(v_*) = \mathbb{E}[(v_* - v_i)^2]}{=} \mathbb{E}[(\mu_*(x) - \mu_i(x))^2]$$

# Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.

IterGP-Cholesky

$i = 0$

Legend:
- - - - Latent Function
- ● Data
- —— Mathematical Posterior Mean
- —— Approximate Posterior Mean
- Mathematical Uncertainty
- Computational Uncertainty

# Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.

IterGP-Cholesky

$i = 1$

- - - - Latent Function
● Data
——— Mathematical Posterior Mean
——— Approximate Posterior Mean
Mathematical Uncertainty
Computational Uncertainty

# Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.

IterGP-Cholesky

$i = 2$

Legend:
- - - - Latent Function
- ● Data
- Mathematical Posterior Mean
- Approximate Posterior Mean
- Mathematical Uncertainty
- Computational Uncertainty

# Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.

IterGP-Cholesky

$i = 3$

Legend:
- - - - Latent Function
- ● Data
- —— Mathematical Posterior Mean
- —— Approximate Posterior Mean
- Mathematical Uncertainty
- Computational Uncertainty

# Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.

IterGP-Cholesky

$i = 5$

# Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.

IterGP-CG

$i = 0$

Latent Function — — — —
Data ●
Mathematical Posterior Mean ———
Approximate Posterior Mean ———
Mathematical Uncertainty
Computational Uncertainty

# Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.

IterGP-CG

$i = 1$

$y$

Variance

$x$

-1.5    -1.0    -0.5    0.0    0.5    1.0    1.5

- - - -  Latent Function
●  Data
——  Mathematical Posterior Mean
——  Approximate Posterior Mean
    Mathematical Uncertainty
    Computational Uncertainty

# Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.

IterGP-CG

$i = 2$

$\mathcal{Y}$

Variance

$\mathcal{X}$

- - - - Latent Function
● Data
⸺ Mathematical Posterior Mean
⸺ Approximate Posterior Mean
▬ Mathematical Uncertainty
▬ Computational Uncertainty

# Computation-Aware GP Inference Illustrated

Interpreting computational and combined uncertainty as error quantification.

IterGP-CG

$i = 3$

Legend:
- - - - Latent Function
- ● Data
- —— Mathematical Posterior Mean
- —— Approximate Posterior Mean
- Mathematical Uncertainty
- Computational Uncertainty

# Computation-Aware GP Inference Illustrated

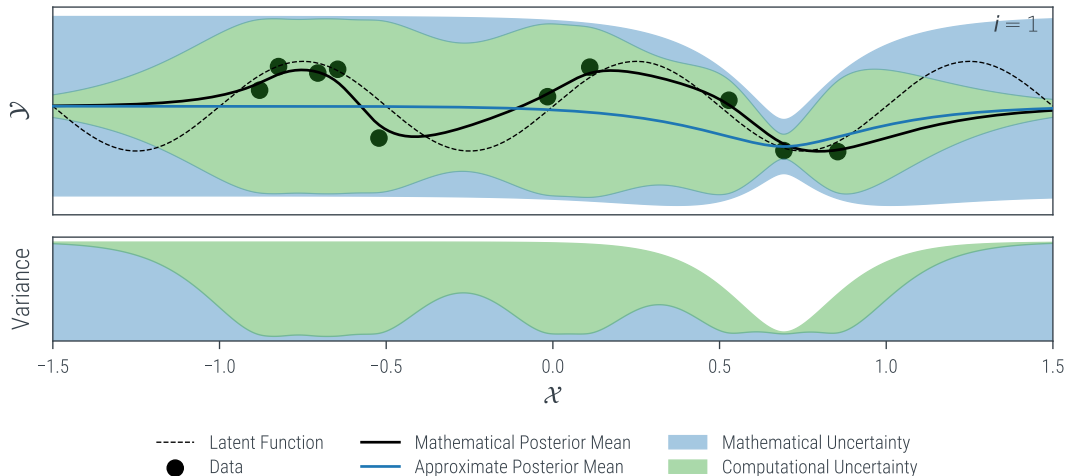Interpreting computational and combined uncertainty as error quantification.

IterGP-CG

$i = 5$

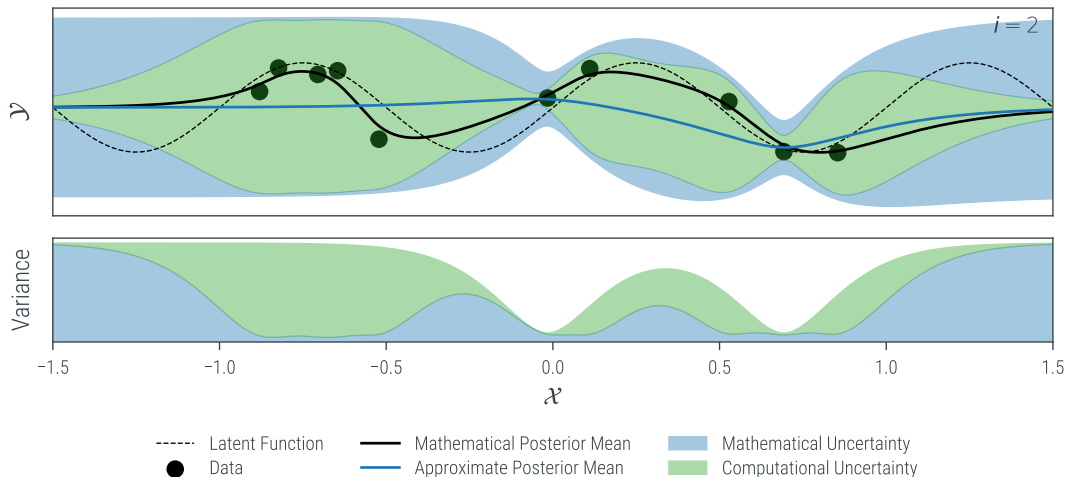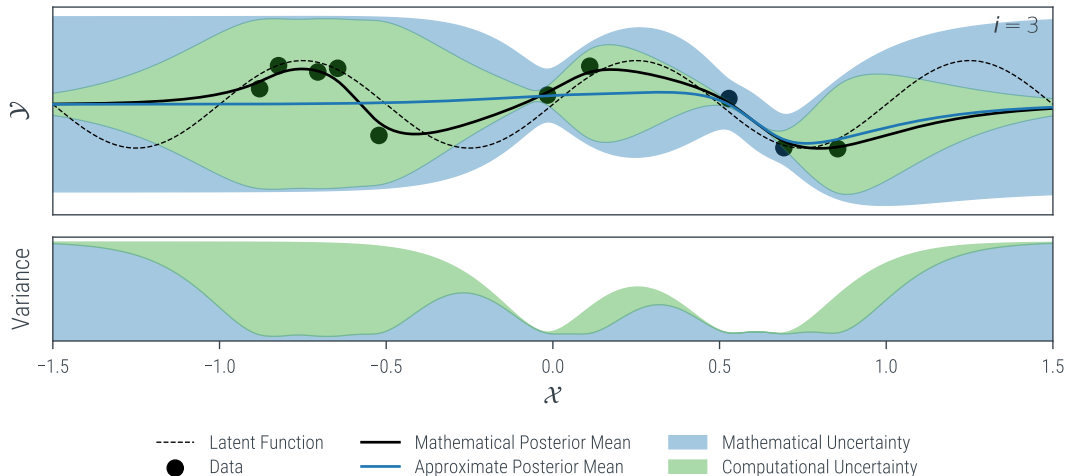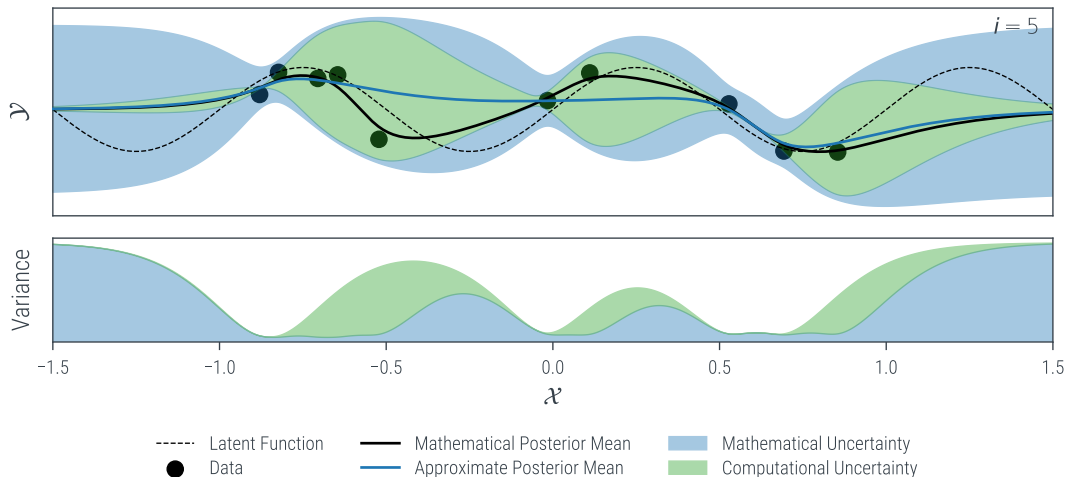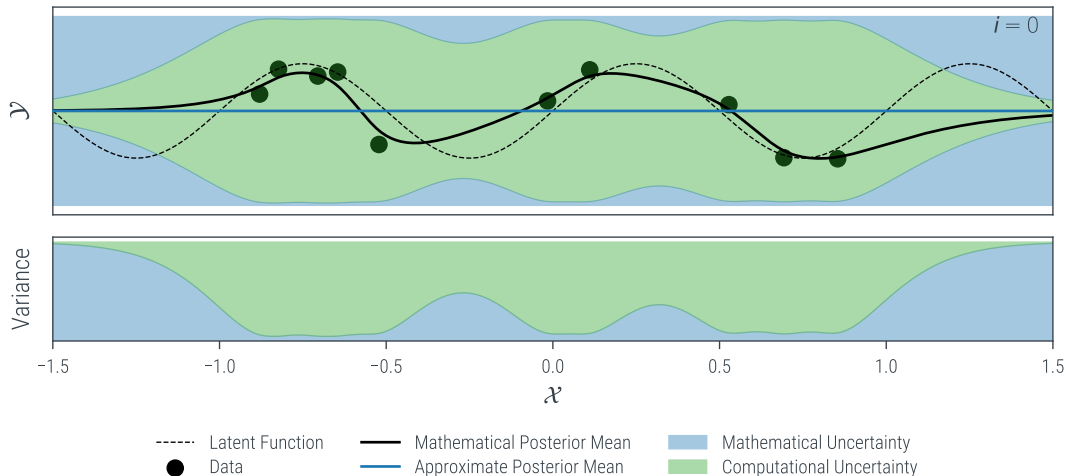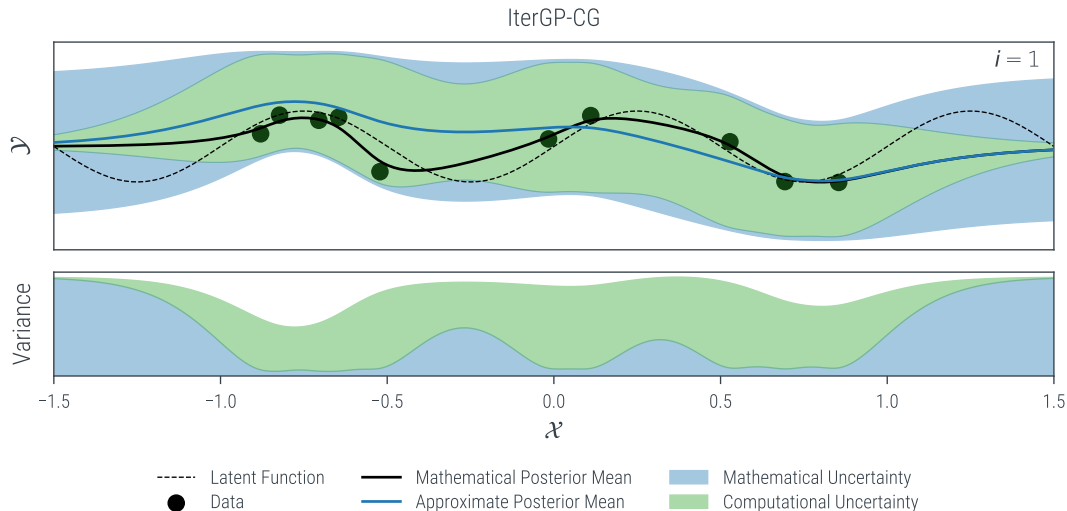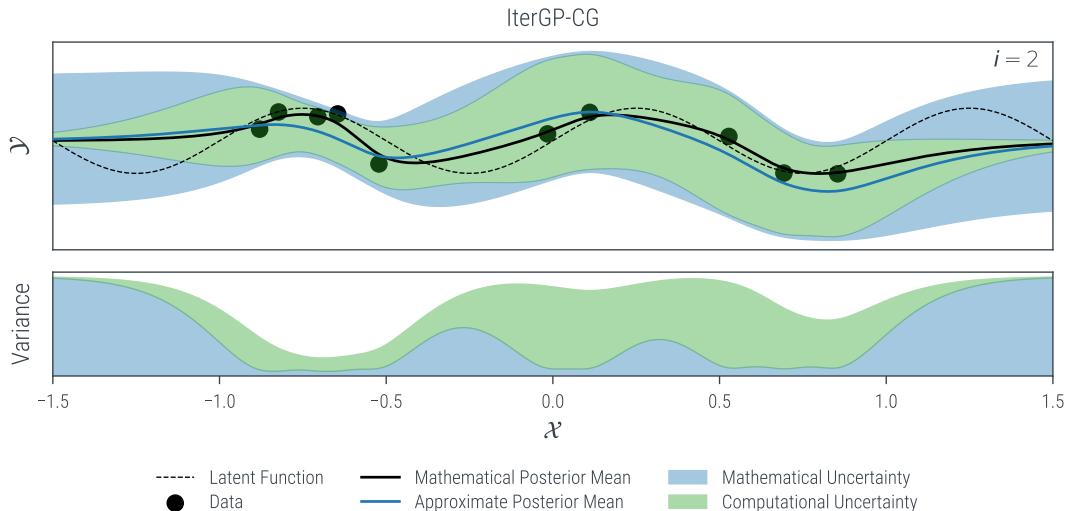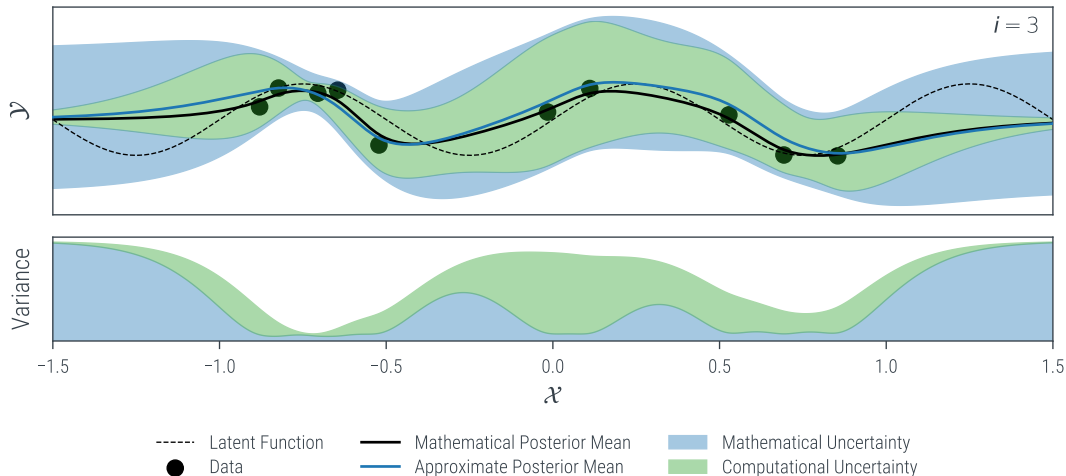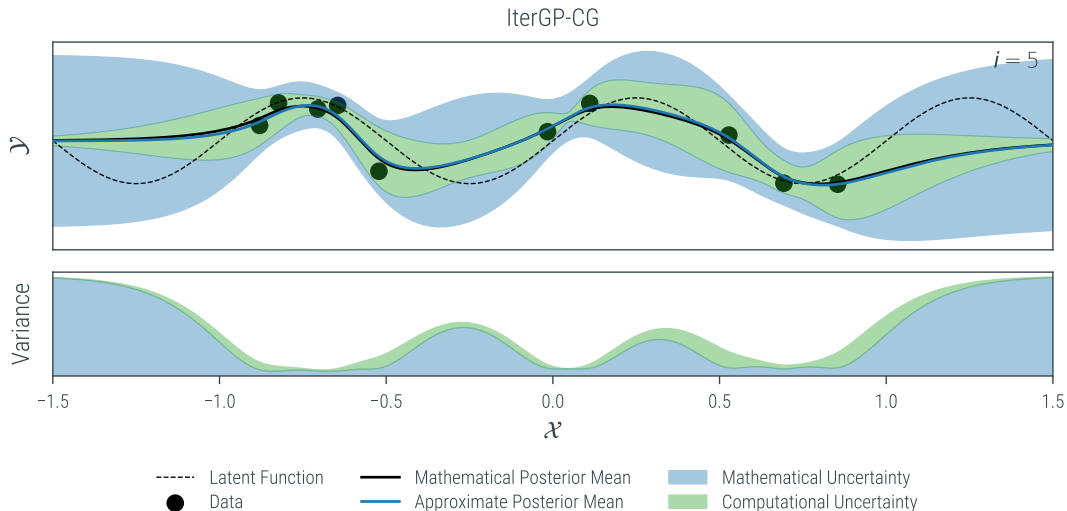Legend:
- - - - Latent Function
- ● Data
- —— Mathematical Posterior Mean
- —— Approximate Posterior Mean
- Mathematical Uncertainty
- Computational Uncertainty

A class of computation-aware iterative methods for GP approximation.

**Input:** prior mean $\mu$, prior kernel $k$, training data $X$, $y$
**Output:** (combined) GP posterior $\mathcal{GP}(\mu_i, k_i)$

1   procedure ITERGP($\mu, k, X, y$)
2     $(\mu_0, k_0) \leftarrow (\mu, k)$
3     $\boldsymbol{\mu} \leftarrow \mu(X)$
4     $\hat{K} \leftarrow k(X, X) + \sigma^2 I$
5     while not STOPPINGCRITERION() do
6       $s_i \leftarrow$ POLICY()        // Select action via policy.
7       $r_{i-1} \leftarrow (y - \boldsymbol{\mu}) - \hat{K} v_{i-1}$      // Predictive residual.
8       $\alpha_i \leftarrow s_i^\mathsf{T} r_{i-1}$        // Observation of linear solver.
9       $d_i \leftarrow \Sigma_{i-1} \hat{K} s_i = (I - C_{i-1} \hat{K}) s_i$    // Search direction.
10      $\eta_i \leftarrow s_i^\mathsf{T} \hat{K} \Sigma_{i-1} \hat{K} s_i = s_i^\mathsf{T} \hat{K} d_i$    // Normalization constant.
11      $C_i \leftarrow C_{i-1} + \frac{1}{\eta_i} d_i d_i^\mathsf{T}$   // Inverse approx. $C_i \approx \hat{K}^{-1}$.
12      $v_i \leftarrow v_{i-1} + \frac{\alpha_i}{\eta_i} d_i$      // Representer weights estimate.
13      $\Sigma_i \leftarrow \Sigma_0 - C_i$       // Computational rep. w. uncertainty.
14     end while
15     $p(v_*) \leftarrow \mathcal{N}(v_*; v_i, \Sigma_i)$   // Belief about representer weights.
16     $\mu_i(\cdot) \leftarrow \mu(\cdot) + k(\cdot, X) v_i$    // Approximate posterior mean.
17     $k_i(\cdot, \cdot) \leftarrow k(\cdot, \cdot) - k(\cdot, X) C_i k(X, \cdot)$   // Combined uncertainty.
18     return $\mathcal{GP}(\mu_i, k_i)$
19 end procedure

## Initialize representer weights belief

$$v_0 = 0$$
$$C_0 = 0$$
$$\Sigma_0 = \Sigma_0 - C_0 = \hat{K}^{-1}$$



●   Approx. Representer Weights $v_i$      ●   Representer Weights $v_*$

# Algorithm: IterGP

A class of computation-aware iterative methods for GP approximation.

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

**Input:** prior mean $\mu$, prior kernel $k$, training data $X$, $y$
**Output:** (combined) GP posterior $\mathcal{GP}(\mu_i, k_i)$
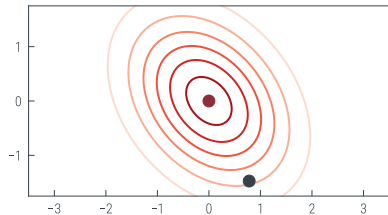
1  **procedure** IterGP$(\mu, k, X, y)$
2      $(\mu_0, k_0) \leftarrow (\mu, k)$
3      $\boldsymbol{\mu} \leftarrow \mu(X)$
4      $\hat{K} \leftarrow k(X, X) + \sigma^2 I$
5      **while not** StoppingCriterion() **do**
6          $s_i \leftarrow$ Policy()                    // Select action via policy.
7          $r_{i-1} \leftarrow (y - \boldsymbol{\mu}) - \hat{K} v_{i-1}$          // Predictive residual.
8          $\alpha_i \leftarrow s_i^{\mathsf{T}} r_{i-1}$              // Observation of linear solver.
9          $d_i \leftarrow \Sigma_{i-1} \hat{K} s_i = (I - C_{i-1}\hat{K}) s_i$     // Search direction.
10          $\eta_i \leftarrow s_i^{\mathsf{T}} \hat{K} \Sigma_{i-1} \hat{K} s_i = s_i^{\mathsf{T}} \hat{K} d_i$      // Normalization constant.
11          $C_i \leftarrow C_{i-1} + \frac{1}{\eta_i} d_i d_i^{\mathsf{T}}$   // Inverse approx. $C_i \approx \hat{K}^{-1}$.
12          $v_i \leftarrow v_{i-1} + \frac{\alpha_i}{\eta_i} d_i$       // Representer weights estimate.
13          $\Sigma_i \leftarrow \Sigma_0 - C_i$       // Computational rep. w. uncertainty.
14      **end while**
15      $p(v_*) \leftarrow \mathcal{N}(v_*; v_i, \Sigma_i)$ // Belief about representer weights.
16      $\mu_i(\cdot) \leftarrow \mu(\cdot) + k(\cdot, X) v_i$   // Approximate posterior mean.
17      $k_i(\cdot, \cdot) \leftarrow k(\cdot, \cdot) - k(\cdot, X) C_i k(X, \cdot)$   // Combined uncertainty.
18      **return** $\mathcal{GP}(\mu_i, k_i)$
19  **end procedure**

## Select action via policy

$$s_i = \textsc{Policy}()$$



● Approx. Representer Weights $v_i$        ● Representer Weights $v_*$

# Algorithm: IterGP

A class of computation-aware iterative methods for GP approximation.

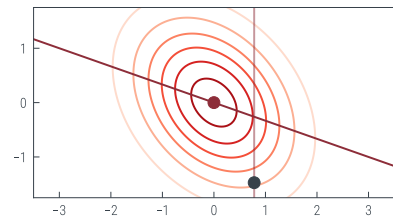**Input:** prior mean $\mu$, prior kernel $k$, training data $X$, $y$
**Output:** (combined) GP posterior $\mathcal{GP}(\mu_i, k_i)$

1  **procedure** ITERGP($\mu, k, X, y$)
2    $(\mu_0, k_0) \leftarrow (\mu, k)$
3    $\boldsymbol{\mu} \leftarrow \mu(X)$
4    $\hat{K} \leftarrow k(X, X) + \sigma^2 I$
5    **while not** STOPPINGCRITERION() **do**
6       $s_i \leftarrow$ POLICY()     // Select action via policy.
7       $r_{i-1} \leftarrow (y - \boldsymbol{\mu}) - \hat{K}v_{i-1}$     // Predictive residual.
8       $\alpha_i \leftarrow s_i^\mathsf{T} r_{i-1}$     // Observation of linear solver.
9       $d_i \leftarrow \Sigma_{i-1} \hat{K} s_i = (I - C_{i-1}\hat{K})s_i$     // Search direction.
10      $\eta_i \leftarrow s_i^\mathsf{T} \hat{K} \Sigma_{i-1} \hat{K} s_i = s_i^\mathsf{T} \hat{K} d_i$     // Normalization constant.
11      $C_i \leftarrow C_{i-1} + \frac{1}{\eta_i} d_i d_i^\mathsf{T}$   // Inverse approx. $C_i \approx \hat{K}^{-1}$.
12      $v_i \leftarrow v_{i-1} + \frac{\alpha_i}{\eta_i} d_i$     // Representer weights estimate.
13      $\Sigma_i \leftarrow \Sigma_0 - C_i$     // Computational rep. w. uncertainty.
14    **end while**
15    $p(v_*) \leftarrow \mathcal{N}(v_*; v_i, \Sigma_i)$   // Belief about representer weights.
16    $\mu_i(\cdot) \leftarrow \mu(\cdot) + k(\cdot, X)v_i$     // Approximate posterior mean.
17    $k_i(\cdot, \cdot) \leftarrow k(\cdot, \cdot) - k(\cdot, X)C_i k(X, \cdot)$   // Combined uncertainty.
18    **return** $\mathcal{GP}(\mu_i, k_i)$
19 **end procedure**

## Observe projected residual

$$\alpha_i = s_i^\mathsf{T} r_{i-1} = (\hat{K}s_i)^\mathsf{T}(v_* - v_{i-1})$$

## Compute search direction

$$d_i = (I - C_{i-1}\hat{K})s_i$$



● Approx. Representer Weights $v_i$      ● Representer Weights $v_*$

# Algorithm: IterGP
A class of computation-aware iterative methods for GP approximation.

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

**Input:** prior mean $\mu$, prior kernel $k$, training data $X$, $y$
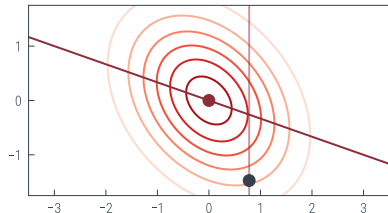**Output:** (combined) GP posterior $\mathcal{GP}(\mu_i, k_i)$
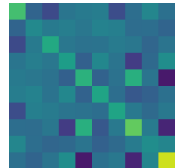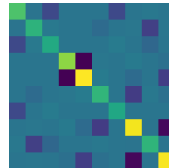
1  **procedure** ITERGP($\mu, k, X, y$)
2    $(\mu_0, k_0) \leftarrow (\mu, k)$
3    $\boldsymbol{\mu} \leftarrow \mu(X)$
4    $\hat{K} \leftarrow k(X, X) + \sigma^2 I$
5    **while not** STOPPINGCRITERION() **do**
6      $s_i \leftarrow$ POLICY()                  // Select action via policy.
7      $r_{i-1} \leftarrow (y - \boldsymbol{\mu}) - \hat{K} v_{i-1}$         // Predictive residual.
8      $\alpha_i \leftarrow s_i^\mathsf{T} r_{i-1}$              // Observation of linear solver.
9      $d_i \leftarrow \Sigma_{i-1} \hat{K} s_i = (I - C_{i-1} \hat{K}) s_i$   // Search direction.
10     $\eta_i \leftarrow s_i^\mathsf{T} \hat{K} \Sigma_{i-1} \hat{K} s_i = s_i^\mathsf{T} \hat{K} d_i$   // Normalization constant.
11     $C_i \leftarrow C_{i-1} + \frac{1}{\eta_i} d_i d_i^\mathsf{T}$   // Inverse approx. $C_i \approx \hat{K}^{-1}$.
12     $v_i \leftarrow v_{i-1} + \frac{\alpha_i}{\eta_i} d_i$       // Representer weights estimate.
13     $\Sigma_i \leftarrow \Sigma_0 - C_i$           // Computational rep. w. uncertainty.
14    **end while**
15    $p(v_*) \leftarrow \mathcal{N}(v_*; v_i, \Sigma_i)$   // Belief about representer weights.
16    $\mu_i(\cdot) \leftarrow \mu(\cdot) + k(\cdot, X) v_i$     // Approximate posterior mean.
17    $k_i(\cdot, \cdot) \leftarrow k(\cdot, \cdot) - k(\cdot, X) C_i k(X, \cdot)$   // Combined uncertainty.
18    **return** $\mathcal{GP}(\mu_i, k_i)$
19  **end procedure**

## Update precision matrix approximation

$$C_i = C_{i-1} + \frac{1}{\eta_i} d_i d_i^\mathsf{T}$$

$$= \begin{pmatrix} | & & | \\ d_1 & \cdots & d_i \\ | & & | \end{pmatrix} \begin{pmatrix} \frac{1}{\eta_1} & & \\ & \ddots & \\ & & \frac{1}{\eta_i} \end{pmatrix} \begin{pmatrix} - d_1^\mathsf{T} - \\ \vdots \\ - d_i^\mathsf{T} - \end{pmatrix}$$

Precision Matrix          Precision Matrix Approx.

# Algorithm: IterGP
A class of computation-aware iterative methods for GP approximation.

UNIVERSITÄT TÜBINGEN
EBERHARD KARLS

**Input:** prior mean $\mu$, prior kernel $k$, training data $X$, $y$
**Output:** (combined) GP posterior $\mathcal{GP}(\mu_i, k_i)$

1  **procedure** ITERGP($\mu, k, X, y$)
2  $\quad (\mu_0, k_0) \leftarrow (\mu, k)$
3  $\quad \boldsymbol{\mu} \leftarrow \mu(X)$
4  $\quad \hat{K} \leftarrow k(X, X) + \sigma^2 I$
5  $\quad$ **while not** STOPPINGCRITERION() **do**
6  $\quad\quad s_i \leftarrow$ POLICY()  // Select action via policy.
7  $\quad\quad r_{i-1} \leftarrow (y - \boldsymbol{\mu}) - \hat{K} v_{i-1}$  // Predictive residual.
8  $\quad\quad \alpha_i \leftarrow s_i^\mathsf{T} r_{i-1}$  // Observation of linear solver.
9  $\quad\quad d_i \leftarrow \boldsymbol{\Sigma}_{i-1} \hat{K} s_i = (I - C_{i-1}\hat{K}) s_i$  // Search direction.
10 $\quad\quad \eta_i \leftarrow s_i^\mathsf{T} \hat{K} \Sigma_{i-1} \hat{K} s_i = s_i^\mathsf{T} \hat{K} d_i$  // Normalization constant.
11 $\quad\quad C_i \leftarrow C_{i-1} + \frac{1}{\eta_i} d_i d_i^\mathsf{T}$  // Inverse approx. $C_i \approx \hat{K}^{-1}$.
12 $\quad\quad v_i \leftarrow v_{i-1} + \frac{\alpha_i}{\eta_i} d_i$  // Representer weights estimate.
13 $\quad\quad \boldsymbol{\Sigma}_i \leftarrow \Sigma_0 - C_i$  // Computational rep. w. uncertainty.
14 $\quad$ **end while**
15 $\quad p(v_*) \leftarrow \mathcal{N}(v_*; v_i, \Sigma_i)$  // Belief about representer weights.
16 $\quad \mu_i(\cdot) \leftarrow \mu(\cdot) + k(\cdot, X) v_i$  // Approximate posterior mean.
17 $\quad k_i(\cdot, \cdot) \leftarrow k(\cdot, \cdot) - k(\cdot, X) C_i k(X, \cdot)$  // Combined uncertainty.
18 $\quad$ **return** $\mathcal{GP}(\mu_i, k_i)$
19 **end procedure**

Update precision matrix approximation

$$C_i = C_{i-1} + \frac{1}{\eta_i} d_i d_i^\mathsf{T}$$

$$= \begin{pmatrix} | & & | \\ d_1 & \cdots & d_i \\ | & & | \end{pmatrix} \begin{pmatrix} \frac{1}{\eta_1} & & \\ & \ddots & \\ & & \frac{1}{\eta_i} \end{pmatrix} \begin{pmatrix} - d_1^\mathsf{T} - \\ \vdots \\ - d_i^\mathsf{T} - \end{pmatrix}$$

Update representer weights belief

$$v_i = C_i(y - \boldsymbol{\mu}) = v_{i-1} + \frac{\alpha_i}{\eta_i} d_i$$

$$\Sigma_i = \Sigma_i - C_i$$

What about the partial Cholesky and CG?
Connection to Other GP Approximations

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

| Method | Actions $s_i$ | Classic Analog |
|---|---|---|
| IterGP-Cholesky | $e_i$ | (partial) Cholesky / subset of data |
| IterGP-EVD | $\mathrm{ev}_i(\hat{K})$ | (partial) eigenvalue decomp. |
| IterGP-CG | $s_i^{\mathrm{PCG}}$ or $\hat{P}^{-1} r_i$ | (preconditioned) CG |
| IterGP-PseudoInput | $k(X, z_i)$ | $\approx$ SVGP |



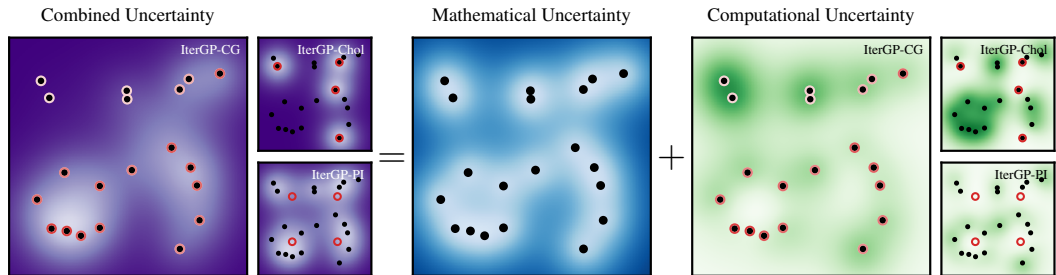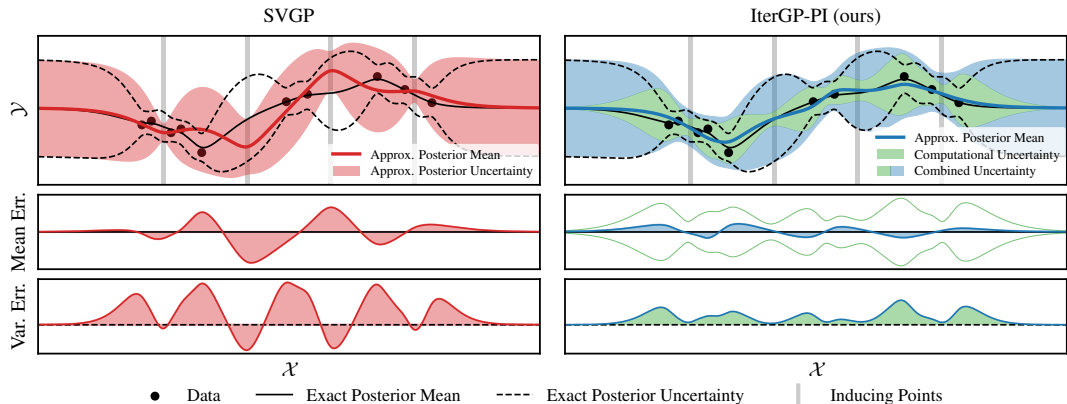Combined Uncertainty   Mathematical Uncertainty   Computational Uncertainty

Figure: Computational uncertainty is small where there either is no data (●) or computation was "targeted" (○).

# Connection to SVGP

One can construct a similar approximation to SVGP with proper uncertainty quantification.

SVGP | IterGP-PI (ours)

Legend:
- Approx. Posterior Mean
- Approx. Posterior Uncertainty
- Computational Uncertainty
- Combined Uncertainty
- $\mathcal{Y}$ (vertical axis)
- Mean Err.
- Var. Err.
- $\mathcal{X}$ (horizontal axis)
- ● Data
- —— Exact Posterior Mean
- - - - Exact Posterior Uncertainty
- | Inducing Points

IterGP-PseudoInput has complexity $\mathcal{O}(n^2 i)$. Are we restricted to quadratic time?

# The Cost of Combined Uncertainty Quantification

Taking a second look at the computational complexity of IterGP.



**Question**: How costly is one iteration of IterGP for a specific policy?

Kernel matrix $\hat{K}$ appears in three ways:

- Observation: $\alpha_i \leftarrow s_i^\mathsf{T}((y - \mu) - \hat{K}v_{i-1})$
- Search direction: $d_i \leftarrow (I - C_{i-1}\hat{K})s_i$
- Normalization const.: $\eta_i \leftarrow s_i^\mathsf{T}\hat{K}d_i$

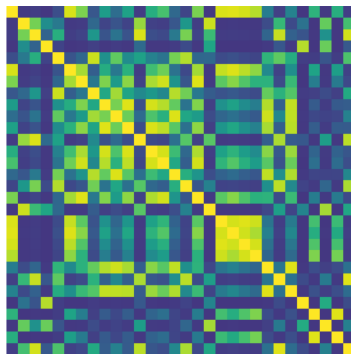# The Cost of Combined Uncertainty Quantification

Taking a second look at the computational complexity of IterGP.

Question: How costly is one iteration of IterGP for a specific policy?

Kernel matrix $\hat{K}$ appears in three ways:

▶ Observation: $\alpha_i \leftarrow s_i^\mathsf{T}((y - \mu) - \hat{K}v_{i-1}) \implies s_i^\mathsf{T}\hat{K}d_j$

▶ Search direction: $d_i \leftarrow (I - C_{i-1}\hat{K})s_i$

▶ Normalization const.: $\eta_i \leftarrow s_i^\mathsf{T}\hat{K}d_i$

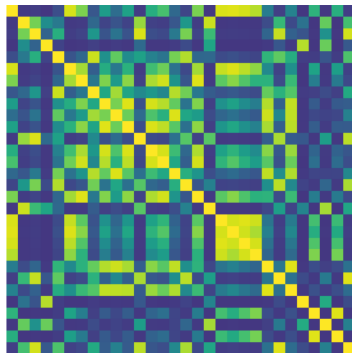# The Cost of Combined Uncertainty Quantification

Taking a second look at the computational complexity of IterGP.

**Question**: How costly is one iteration of IterGP for a specific policy?

Kernel matrix $\hat{K}$ appears in three ways:

- Observation: $\alpha_i \leftarrow s_i^{\mathsf{T}}\left((y - \mu) - \hat{K}v_{i-1}\right) \implies s_i^{\mathsf{T}}\hat{K}d_j$
- Search direction:
  $d_i \leftarrow (I - C_{i-1}\hat{K})s_i = s_i - \sum_{j=1}^{i} \frac{1}{\eta_j} d_j d_j^{\mathsf{T}}\hat{K}s_i \implies s_i^{\mathsf{T}}\hat{K}d_j$
- Normalization const.: $\eta_i \leftarrow s_i^{\mathsf{T}}\hat{K}d_i$

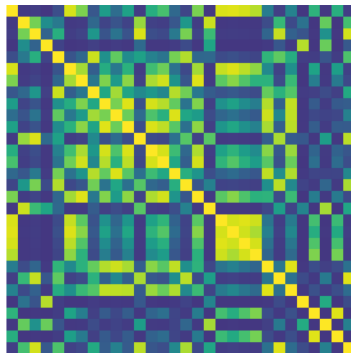# The Cost of Combined Uncertainty Quantification

Taking a second look at the computational complexity of IterGP.

**Question**: How costly is one iteration of IterGP for a specific policy?

Kernel matrix $\hat{K}$ appears in three ways:

▶ Observation: $\alpha_i \leftarrow s_i^\mathsf{T}((y - \mu) - \hat{K}v_{i-1}) \implies s_i^\mathsf{T}\hat{K}d_j$

▶ Search direction:
$d_i \leftarrow (I - C_{i-1}\hat{K})s_i = s_i - \sum_{j=1}^{i} \frac{1}{\eta_j}d_j d_j^\mathsf{T}\hat{K}s_i \implies s_i^\mathsf{T}\hat{K}d_j$

▶ Normalization const.: $\eta_i \leftarrow s_i^\mathsf{T}\hat{K}d_i$

**Idea**: Choose $i$ actions with at most $\ell \ll n$ non-zero entries. $\implies \mathcal{O}(\ell^2)$ per iteration!

We only operate on the data that we target with computation → arbitrary computation cost!

# Working with Infinite Data

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

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

### Theorem (Online GP Approximation with IterGP)

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

$$\tilde{\boldsymbol{s}}_i = \begin{pmatrix} \boldsymbol{s}_i \\ 0 \end{pmatrix} \tag{1}$$

*Then the posterior returned by IterGP for the dataset $(\boldsymbol{X}, \boldsymbol{y})$ using actions $\boldsymbol{s}_i$ is identical to the posterior returned by IterGP for the extended dataset using actions $\tilde{\boldsymbol{s}}_i$, i.e. it holds for any $i \in \{1, \ldots, n\}$, that*

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

# An Approximation or a Better Model?

An alternative view of IterGP as a better model for the way we do inference with a computer.

UNIVERSITÄT TÜBINGEN
EBERHARD KARLS

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



▶ The distinction between data and computation vanishes from this perspective.

# An Approximation or a Better Model?

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?

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

An Approximation or a Better Model?

An alternative view of IterGP as a better model for the way we do inference with a computer.

UNIVERSITÄT
TÜBINGEN
EBERHARD KARLS

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?

$$f \sim \mathcal{GP}(\mu, k)$$
$$\tilde{y} \mid f(X) \sim \mathcal{N}(S_i^{\mathsf{T}} f(X), \sigma^2 S_i^{\mathsf{T}} S_i)$$
$$f \mid X, \tilde{y} \sim \mathcal{GP}(\mu_i, k_i)$$

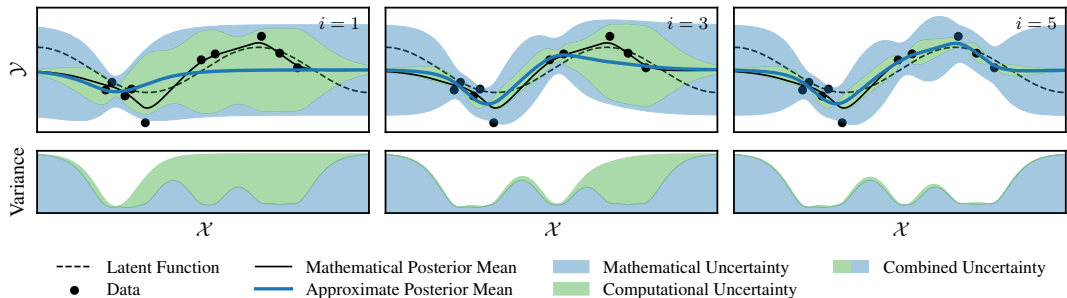▶ IterGP's combined posterior is equivalent exact GP regression for linearly projected data.

How meaningful is computational and combined uncertainty?
Theoretical Analysis

# Combined Uncertainty as Worst Case Error

The combined uncertainty is a tight worst case bound on the relative error to the latent function.

GP:
$$\frac{|\text{Latent Function}(x) - \text{Math. Posterior Mean}(x)|}{\|\text{Latent Function}\|} \leq \text{Posterior Pred. Std. Deviation}(x) \; \bullet$$

IterGP:
$$\frac{|\text{Latent Function}(x) - \text{Approx. Posterior Mean}(x)|}{\|\text{Latent Function}\|} \leq \text{Combined Std. Deviation}(x) \; \bullet + \bullet$$

**Exact uncertainty quantification in quadratic / linear / constant time!**

# Combined Uncertainty as Worst Case Error

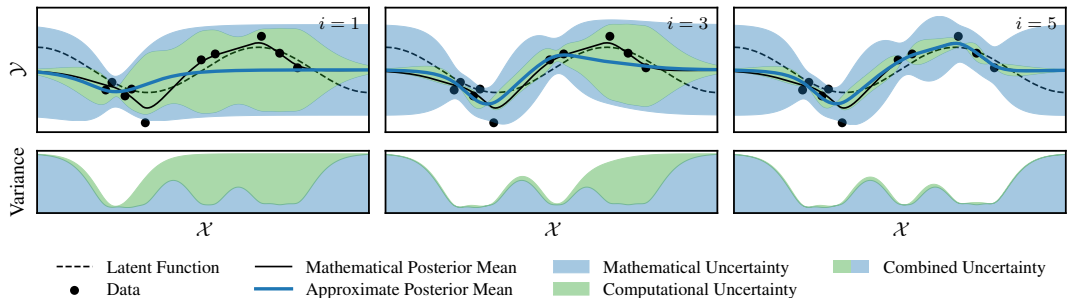The combined uncertainty is a tight worst case bound on the relative error to the latent function.

placeholder

- - - - Latent Function
- Data
—— Mathematical Posterior Mean
—— Approximate Posterior Mean
Mathematical Uncertainty
Computational Uncertainty
Combined Uncertainty

## Theorem

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

error of approximate posterior mean ●

$$\sup_{g \in \mathcal{H}_{k\sigma} : \|g\|_{\mathcal{H}_{k\sigma}} \leq 1} \underbrace{\mu_*^g(\boldsymbol{x}) - \mu_i^g(\boldsymbol{x})}_{\text{computational error} \bullet} = \sqrt{k_i^{\text{comp}}(\boldsymbol{x}, \boldsymbol{x})} \tag{3}$$

Numerics of Machine Learning – Winter 2022/23 – Lecture 04: Computation-Aware GP Inference – © N. Bosch, J. Grosse, P. Hennig, A. Kristiadi, M. Pförtner, J. Schmidt, F. Schneider, L. Tatzel, J. Wenger, 2022 CC BY-NC-SA   ▶  29

## Summary

- ▶ Approximate GPs by learning the representer weights.
- ▶ Can quantify approximation error *probabilistically*.
- ▶ Variants of IterGP defined via the policy learn actively.
- ▶ Distinction between data and computation vanishes.
- ▶ Exact UQ in arbitrary time with strong guarantees.

### Please cite this course, as

```
@techreport{NoML22,
    title = {Numerics of Machine Learning},
    author = {N. Bosch and J. Grosse
    and P. Hennig and A. Kristiadi
    and M. Pförtner and J. Schmidt
    and F. Schneider and L. Tatzel
    and J. Wenger},
    series = {Lecture Notes in Machine Learning},
    year = {2022},
    institution = {Tübingen AI Center},
}
```

## Next Week

- ▶ How to simulate, i.e. *learn* the dynamics of systems that follow (partially-known) physical laws.