

Outline

1 Introduction

2 EVIGP

- Gaussian Process Regression
- Bayesian Approach
- Energetic Variational Inference
- Numerical results

3 Active Learning Manifold Gaussian Process

- Manifold Gaussian Process
- Active Learning Cohn
- Numerical results

Introduction: Computer Experiments

What is a computer experiment?

Introduction: Computer Experiments

What is a computer experiment? It is a simulation using computer models.

- physical,
- biological,
- engineering, etc.

Challenges

- dimension
- cost of data acquisition

Introduction: Gaussian Process

Why Gaussian Process (GP)?

- simple assumption with an explicit prediction formula
- can perform Uncertainty Quantification (UQ) easily

When using Bayesian Approach to do Gaussian Process regression,

- 1 assume prior distributions with hyperparameters
- 2 obtain the posterior distribution
- 3 maximize the posterior to find the best value for the hyperparameters
- 4 find posterior predictive distribution at untried x

Energetic Variational Inference Gaussian Process

The first method uses Energetic Variational Inference to find the maximum point of the posterior distribution function.

- Energetic Variational Approach (Wang Y., Liu C., 2022)
- Variational Inference (Wang, Y., Chen, J., Liu, C. et al., 2021)

We have

- stable posterior approximation
- variable selection

Introduction: Active Learning Manifold Gaussian Process

The second method uses a quick active learning method on a manifold Gaussian process. We aim to achieve the following:

- Dimension reduction
- Active learning

We use:

- Manifold learning: Neural network
- Active learning criteria: Active Learning Cohn

Outline

1 Introduction

2 EVIGP

- Gaussian Process Regression
- Bayesian Approach
- Energetic Variational Inference
- Numerical results

3 Active Learning Manifold Gaussian Process

- Manifold Gaussian Process
- Active Learning Cohn
- Numerical results

Gaussian Process Regression

The Gaussian Process model.

$$y_i = y(\mathbf{x}_i) = \mathbf{g}(\mathbf{x}_i)^\top \boldsymbol{\beta} + Z(\mathbf{x}_i) + \epsilon_i, \quad i = 1, \dots, n, \quad (1)$$

- $\{\mathbf{x}_i, y_i\}_{i=1}^n$ are the data, $\mathbf{x}_i \in \Omega \subseteq \mathbb{R}^d$, $y_i \in \mathbb{R}$
- $\mathbf{g}(\mathbf{x})$ is a p -dimensional vector of basis function, and $\boldsymbol{\beta}$ are their linear coefficients
 - zero order effect (intercept), $\mathbf{g}(\mathbf{x}) = 1$, $\boldsymbol{\beta} = \beta_0$, $p = 1$
 - first order effect, $\mathbf{g}(\mathbf{x}) = [1, x_1, x_2, \dots, x_d]$,
 $\boldsymbol{\beta} = [\beta_0, \beta_1, \beta_2, \dots, \beta_d]$
 - second order effect, etc.
- The noise $\epsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_\omega^2)$ and is also independent of the other stochastic ingredients of (1).

About Gaussian Process $Z(\cdot)$

$Z(\cdot) \sim GP(0, \tau^2 \mathbf{K})$, which indicates $\mathbb{E}[Z(\mathbf{x})] = 0$ and $\text{Cov}[Z(\mathbf{x}_1), Z(\mathbf{x}_2)] = \tau^2 k(\mathbf{x}_1, \mathbf{x}_2)$.

- τ^2 is a constant variance and $k(\cdot, \cdot; \omega) : \Omega \times \Omega \mapsto \mathbb{R}_+$ is the correlation function of the stochastic process with hyperparameters ω .
- Here we use *Gaussian kernel* function, symmetric, and positive definite. $k(x_1, x_2; \omega) = \exp \left\{ - \sum_{j=1}^d \omega_j (x_{1j} - x_{2j})^2 \right\}$, with $\omega \in \mathbb{R}^d$ and $\omega \geq 0$.
- Here we also define the noise-to-signal ratio $\eta = \sigma_\omega^2 / \tau^2$. It's expected to be small

About response y

With above assumption, \mathbf{y}_n is also a gaussian process where $\mathbb{E}[y(\mathbf{x})] = \mathbf{g}(\mathbf{x})^\top \boldsymbol{\beta}, \forall \mathbf{x} \in \Omega$ and

$$\begin{aligned} \text{Cov}[y(\mathbf{x}_1), y(\mathbf{x}_2)] &= \tau^2 k(\mathbf{x}_1, \mathbf{x}_2; \boldsymbol{\omega}) + \sigma^2 \delta(\mathbf{x}_1, \mathbf{x}_2), \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \Omega, \\ &= \tau^2 [k(\mathbf{x}_1, \mathbf{x}_2; \boldsymbol{\omega}) + \eta \delta(\mathbf{x}_1, \mathbf{x}_2)], \end{aligned}$$

where $\delta(\mathbf{x}_1, \mathbf{x}_2) = 1$ if $\mathbf{x}_1 = \mathbf{x}_2$ and 0 otherwise.

In summary, all the unknown parameters are $\boldsymbol{\beta}, \boldsymbol{\omega}, \tau^2$, and η , and we define $\boldsymbol{\theta} = (\boldsymbol{\beta}, \boldsymbol{\omega}, \tau^2, \eta)$.

Hyperparameters learning objective

log-marginal likelihood:

$$\begin{aligned} \ell(\tau^2, \eta) = & -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \tau^2 - \frac{1}{2} \log \det(\mathbf{K}_n + \eta \mathbf{I}_n) \\ & - \frac{1}{2\tau^2} \mathbf{Y}_n^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{Y}_n. \end{aligned} \quad (2)$$

differentiating ℓ with respect to τ^2 , solve

$$\hat{\tau}^2 = \frac{1}{n} \mathbf{Y}_n^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{Y}_n. \quad (3)$$

Final form

$$\ell(\eta) = -\frac{1}{2}\mathbf{Y}_n^\top(\mathbf{K}_n + \eta\mathbf{I}_n)^{-1}\mathbf{Y}_n - \frac{1}{2}\log\det(\mathbf{K}_n + \eta\mathbf{I}_n) + C, \quad (4)$$

Outline

1 Introduction

2 EVIGP

- Gaussian Process Regression
- **Bayesian Approach**
- Energetic Variational Inference
- Numerical results

3 Active Learning Manifold Gaussian Process

- Manifold Gaussian Process
- Active Learning Cohn
- Numerical results

Priors

We assume the following prior distributions for the parameters

- $\beta \sim \mathcal{MVN}_p(\mathbf{0}, \nu^2 \mathbf{R})$ or $\rho(\beta) \propto 1$ for a non-informative prior.
- $\omega_i \stackrel{\text{i.i.d.}}{\sim} \text{Gamma}(a_\omega, b_\omega)$, for $i = 1, \dots, d$
- $\tau^2 \sim \text{Inverse-}\chi^2(df_{\tau^2})$
- $\eta \sim \text{Gamma}(a_\eta, b_\eta)$.

Then we have the sampling distribution

$$\mathbf{y}_n | \boldsymbol{\theta} \sim \mathcal{MVN}_n(\mathbf{G}\boldsymbol{\beta}, \tau^2(\mathbf{K}_n + \eta \mathbf{I}_n)), \quad (5)$$

where \mathbf{y}_n is the vector of y_i 's, \mathbf{G} is a matrix with each row being $\mathbf{g}(\mathbf{x}_i)^\top$. The matrix \mathbf{K}_n is the $n \times n$ kernel matrix and $\mathbf{K}_n[i, j] = K(\mathbf{x}_i, \mathbf{x}_j)$ and is symmetric and positive definite. \mathbf{I}_n is an $n \times n$ identity matrix.

Conditional posterior distributions of β

We start with the informative prior case. Using the Bayes formula we obtain

$$\begin{aligned} p(\boldsymbol{\theta} \mid \mathbf{y}_n) &\propto p(\mathbf{y}_n \mid \boldsymbol{\theta}) \cdot p(\boldsymbol{\theta}) \\ p(\boldsymbol{\beta} \mid \mathbf{y}_n, \boldsymbol{\omega}, \tau^2, \eta) &\propto p(\mathbf{y}_n \mid \boldsymbol{\beta}, \boldsymbol{\omega}, \tau^2, \eta) \cdot p(\boldsymbol{\beta} \mid \boldsymbol{\omega}, \tau^2, \eta) \\ &\propto p(\mathbf{y}_n \mid \boldsymbol{\theta}) \cdot p(\boldsymbol{\beta}) \end{aligned}$$

Then, $\beta|\mathbf{y}_n, \omega, \tau^2, \eta \sim \mathcal{MVN}_d(\hat{\beta}_n, \Sigma_{\beta|n})$, where

$$\begin{aligned}\Sigma_{\beta|n} &= \left[\frac{1}{\tau^2} \mathbf{G}^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{G} + \frac{1}{\nu^2} \mathbf{R}^{-1} \right]^{-1} \\ \hat{\beta}_n &= \Sigma_{\beta|n} \frac{(\mathbf{G}^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{y}_n)}{\tau^2}\end{aligned}$$

$\Sigma_{\beta|n}$ and $\hat{\beta}_n$ are used to form an interval estimate. When the confidence interval contains 0, the corresponding effect is considered nonsignificant and is therefore discarded.

$$\Sigma_{\beta|n} = \left[\frac{1}{\tau^2} \mathbf{G}^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{G} + \frac{1}{\nu^2} \mathbf{R}^{-1} \right]^{-1}$$

Thus, we use \mathbf{R} for regularization and ν as the shrinking coefficient. We define the matrix \mathbf{R} to be a diagonal matrix and follow the Hierarchy principle,

- lower order effects are more likely to be significant than higher order effects
- effects of the same order are equally likely to be significant

Structure of Basis Function and Regularization Matrix

$$\mathbf{g}(\mathbf{x}_i)^\top = [1 \quad x_1 \quad x_2 \quad \cdots \quad x_d \quad x_1^2 \quad x_1 x_2 \quad \cdots \quad x_d^2]$$

Order : 0th 1st 2nd

$$\mathbf{R} = \begin{bmatrix} 1 & & & & & & \\ & r & & & & & \\ & & r & & & & \\ & & & \ddots & & & \\ & & & & r & & \\ & & & & & r^2 & \\ & & & & & & r^2 \\ & & & & & & & \ddots \\ & & & & & & & & r^2 \end{bmatrix}, \quad r > 1$$

Variable Selection step

- 1 use cross validation to find the best ν
- 2 with fixed ν , do regression and calculate the interval estimate of $\hat{\beta}_n$
- 3 if a confidence interval contains 0, it's nonsignificant and will be discarded
- 4 according to the weak Heredity principle, *higher-order effects are only considered if one or more of its parent effects are significant*; otherwise, they are removed.
 - e.g. x_1x_2 will be discarded if neither x_1 nor x_2 is significant
- 5 The new model has a new q and β , repeat the above steps.

Posterior, informative prior case

With the informative prior, the posterior distribution is

$$\begin{aligned}
p(\boldsymbol{\omega}, \tau^2, \eta | \mathbf{y}_n) &\propto \int p(\boldsymbol{\theta} | \mathbf{y}_n) \mathrm{d}\boldsymbol{\beta} \\
&\propto \det(\boldsymbol{\Sigma}_{\boldsymbol{\beta} | n})^{1/2} \\
&\times \exp \left[-\frac{1}{2} \hat{\boldsymbol{\beta}}_n^\top \boldsymbol{\Sigma}_{\boldsymbol{\beta} | n}^{-1} \hat{\boldsymbol{\beta}}_n + \frac{1}{2\tau^2} \mathbf{y}_n^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{y}_n \right] \\
&\times \det(\mathbf{K}_n + \eta \mathbf{I}_n)^{1-1/2} (\tau^2)^{-n/2} p(\tau^2) p(\boldsymbol{\omega}) p(\eta) \quad (6)
\end{aligned}$$

And if using the noninformative prior, we can further integrate (6) with respect to τ^2 and obtain $p(\boldsymbol{\omega}, \eta | \mathbf{y}_n)$.

Noninformative prior case

In the noninformative prior case, we have the estimate of $\hat{\beta}_n$ being normal with its variance and mean being

$$\begin{aligned}\Sigma_{\beta|n} &= \tau^2 [\mathbf{G}^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{G}]^{-1} \\ \hat{\beta}_n &= [\mathbf{G}^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{G}]^{-1} [\mathbf{G}^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1}] \mathbf{y}_n\end{aligned}$$

In this case, we can define

$$\begin{aligned} s_n^2 &= \tau^{-2} \hat{\boldsymbol{\beta}}_n^\top \boldsymbol{\Sigma}_{\boldsymbol{\beta}|n}^{-1} \hat{\boldsymbol{\beta}}_n + \mathbf{y}_n^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{y}_n \\ &= \mathbf{y}_n^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \left[\mathbf{G} (\mathbf{G}^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{G})^{-1} \mathbf{G}^\top + (\mathbf{K}_n + \eta \mathbf{I}_n) \right] \\ &\quad \times (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{y}_n. \end{aligned}$$

so that above posterior (6) can be simplified.

$$\hat{A}^2 = (1 + \hat{A}^2) / (1 + \hat{A}^2) \quad \text{and} \quad \hat{A}^2 = (1 + \hat{A}^2) / (1 + \hat{A}^2)$$
$$-2|_{\text{Coulomb}} = \text{Coulomb Inverse} = 2/(1f - 1 - \text{Coulomb}^2)$$

Posterior predictive distribution

Given the parameters (ω, τ^2, η) , the posterior predictive distribution of $y(x)$ at query point $x \in \Omega$ is the following

$$y(x)|y_n, \omega, \tau^2, \eta \sim \mathcal{N}(\hat{\mu}(x), \sigma_n^2(x)), \quad (9)$$

where

$$\hat{\mu}(x) = g(x)^\top \hat{\beta}_n + K(x, \mathcal{X}_n)(K_n + \eta I_n)^{-1}(y_n - G\hat{\beta}_n), \quad (10)$$

$$\sigma_n^2(\mathbf{x}) = \tau^2 \left\{ 1 - K(\mathbf{x}, \mathcal{X}_n)(\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} K(\mathcal{X}_n, \mathbf{x}) + \mathbf{c}(\mathbf{x})^\top [\mathbf{G}^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{G}]^{-1} \mathbf{c}(\mathbf{x}) \right\} \quad (11)$$

$$\begin{aligned} c(x) &= g(x) - G^\top (K_n + \eta I_n)^{-1} K(\mathcal{X}_n, x), \\ K(x, \mathcal{X}_n) &= K(\mathcal{X}_n, x)^\top = [K(x, x_1), \dots, K(x, x_n)] \end{aligned}$$

Computational Issue

To avoid computational issues caused by ill-conditioned matrices, here we plug in $\hat{\beta}_n$ into $\hat{\mu}(x)$ and obtain the alternative formula.

$$\hat{\mu}(\mathbf{x}) = \left((\mathbf{g}(\mathbf{x})^\top - K(\mathbf{x}, \mathcal{X}_n)(\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{G}) \right. \\ \left. \times (\mathbf{G}^\top (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{G})^{-1} \mathbf{G}^\top + K(\mathbf{x}, \mathcal{X}_n) \right) (\mathbf{K}_n + \eta \mathbf{I}_n)^{-1} \mathbf{y}_n$$

Sampling Procedure

The sampling procedure is

- 1 Generate posterior samples from $p(\boldsymbol{\omega}, \tau^2, \eta | \mathbf{y}_n)$ or $p(\boldsymbol{\omega}, \eta | \mathbf{y}_n)$, depending on which prior is used for $\boldsymbol{\beta}$, and then select the mode.
- 2 Based on the posterior samples of $(\boldsymbol{\omega}, \tau^2, \eta)$, generate the conditional posterior distribution for $\boldsymbol{\beta}$. If non-informative prior is used, generate the conditional distribution of $\tau^2 | \boldsymbol{\omega}, \eta$.
- 3 Based on the posterior samples of $(\boldsymbol{\omega}, \tau^2, \eta)$ (or $(\boldsymbol{\omega}, \eta)$ for non-informative prior), generate the posterior predictive distribution of any query point $y(x) | \mathbf{y}_n, \boldsymbol{\omega}, \tau^2, \eta$. For non-informative prior, the posterior predictive distribution $y(x) | \mathbf{y}_n, \boldsymbol{\omega}, \eta$ is the same except τ^2 is replaced by $\hat{\tau}^2$.

Outline

1 Introduction

2 EVIGP

- Gaussian Process Regression
- Bayesian Approach
- Energetic Variational Inference
- Numerical results

3 Active Learning Manifold Gaussian Process

- Manifold Gaussian Process
- Active Learning Cohn
- Numerical results

Preliminary

We start with the continuous formulation.

- denote the posterior, either (6) or (8) to be ρ^*
- define dynamic flow map $\mathbf{x} = \phi(\mathbf{z}, t) : \mathbb{R}^d \rightarrow \mathbb{R}^d$, smooth and one-to-one.
 - It is designed to transform an initial pdf $\rho^0(\mathbf{z})$ to $\rho^*(\mathbf{x})$ as $t \rightarrow \infty$
 - velocity: $\dot{\phi}(\mathbf{z}, t) := \mathbf{u}(\phi(\mathbf{z}, t), t)$, $\mathbf{z} \in \mathbb{R}^d$, $t > 0$
 - transportation equation: $\dot{\rho} + \nabla \cdot (\rho \mathbf{u}) = 0$, $\rho(\mathbf{x}, 0) = \rho^0(\mathbf{x})$ and $\rho(\mathbf{x}, \infty) = \rho^*(\mathbf{x})$
- the energy functional $\mathcal{F}[\phi_t]$ is the KL divergence of ρ and ρ^* :
$$\mathcal{F}[\phi_t] = \int \rho(\mathbf{x}) \ln \left(\frac{\rho}{\rho^*} \right) d\mathbf{x}.$$
 And we expect $\mathcal{F}[\phi_t] \rightarrow 0$ as $t \rightarrow \infty$.

Energy Dissipation Law

$$\frac{d}{dt}\mathcal{F}[\phi_t] = -\Delta(\phi_t, \dot{\phi}_t), \quad (12)$$

The $\Delta(\phi_t, \dot{\phi}_t)$ represents the rate of energy dissipation,

- $\Delta(\phi_t, \dot{\phi}_t) \geq 0$ so that $\mathcal{F}(\phi_t)$ decreases with time
- a simple quadratic functional in terms of $\dot{\phi}_t$

$$\Delta(\phi_t, \dot{\phi}_t) = \int_{\Omega_t} \rho_{[\phi_t]} \|\dot{\phi}_t\|_2^2 dx,$$

where

- $\rho_{[\phi_t]}$ denotes the pdf of the distribution at t
- Ω_t is the current support
- $\|\mathbf{a}\|_2 = \mathbf{a}^\top \mathbf{a}, \forall \mathbf{a} \in \mathbb{R}^d$.

Solving, step 1

As stated in *Wang, Y., Chen, J., Liu, C. et al., 2021*

- 1 perform discretization on space parameter, which is called the "approximation-then-variation" approach and obtain the particle version of the PDE.

$$\frac{d}{dt} \mathcal{F}_h(\{\mathbf{x}_i(t)\}_{i=1}^N) = -\Delta_h(\{\mathbf{x}_i(t)\}_{i=1}^N, \{\dot{\mathbf{x}}_i(t)\}_{i=1}^N). \quad (13)$$

Here $\mathbf{x}_i(t) = \phi(\mathbf{x}_i(0), t)$ where $\mathbf{x}_i(0)$ is sampled from the initial distribution ρ^0 , and

$$\mathcal{F}_h(\mathbf{x}_i) = \frac{1}{N} \sum_{i=1}^N \left(\ln \left(\frac{1}{N} \sum_{j=1}^N K_h(\mathbf{x}_i, \mathbf{x}_j) \right) + V(\mathbf{x}_i) \right)$$

where $K_h(\mathbf{x}_i, \mathbf{x}_j) \propto \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/h^2)$, h is kernel bandwidth and $V = -\log \rho^*$.

Solving, step 3

- 3** using implicit Euler for time discretization, the numerical scheme is

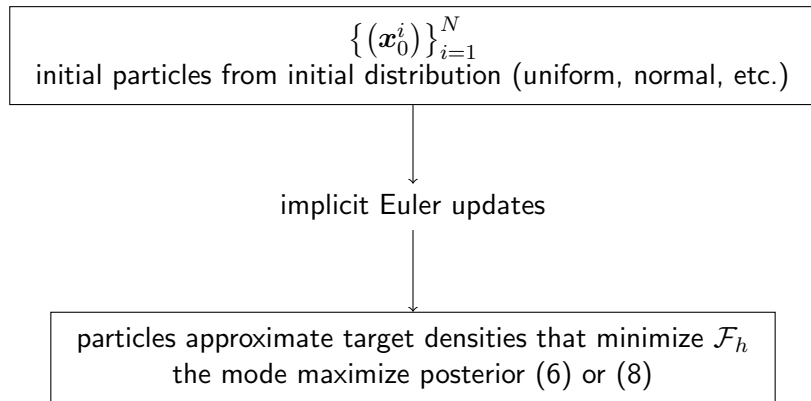
$$\frac{1}{N} \frac{\mathbf{x}_i^{n+1} - \mathbf{x}_i^n}{\tau} = -\frac{\delta \mathcal{F}_h}{\delta \mathbf{x}_i} \left(\{\mathbf{x}_i^{n+1}\}_{i=1}^N \right)$$

Using Proximal Point Algorithm, we need to solve

$$\{\mathbf{x}_i^{n+1}\}_{i=1}^N = \arg \min_{\{\mathbf{x}_i\}_{i=1}^N} J(\{\mathbf{x}_i\}_{i=1}^N)$$

- $J(\{\mathbf{x}_i\}_{i=1}^N) := \frac{1}{2\tau} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{x}_i^n\|^2 / N + \mathcal{F}_h(\{\mathbf{x}_i\}_{i=1}^N)$
- $\{\mathbf{x}(t)\}_{i=1}^N$ is the locations of particles at time t
- $\dot{\mathbf{x}}_i(t)$ is the derivative of \mathbf{x}_i with t
- subscript h is the bandwidth parameter of the kernel

EVI as sampling method road map



About the diffusion term

With KL divergence as the energy functional, EnVarA gives

$$\mathcal{F}[\phi_t] = \int \rho \ln \rho + \rho V \, \mathrm{d}\mathbf{x} \implies \rho \mathbf{u} = -(\nabla \rho + \rho \nabla V)$$

Combine with the transportation equation we have

$$\dot{\rho} = \nabla \cdot (\nabla \rho + \rho \nabla V)$$

We can set the diffusion term to 0. And this correspond to use

$$\mathcal{F}[\phi_t] = \int \rho V \, \mathrm{d}\mathbf{x} \text{ and thus}$$

$$\mathcal{F}_h(\mathbf{x}_i) = V(\mathbf{x}_i) = -\log \rho^*(\mathbf{x}_i)$$

In this case, the scheme ends up a simple gradient flow that should converges to the local maximum points of $\rho^*(\mathbf{x}_i)$.

Two variants of EVI-GP

- EVI-Post: Generate N particles and they will approximate the posterior samples through implicit Euler updates
- EVI-MAP: This variant corresponds to Maximum A Posteriori (MAP), and is thus named EVI-MAP. Simply set $\mathcal{F}(\mathbf{x}) = V(\mathbf{x})$ and $N = 1$. This is equivalent to proximal point algorithm.

Outline

1 Introduction

2 EVIGP

- Gaussian Process Regression
- Bayesian Approach
- Energetic Variational Inference
- Numerical results

3 Active Learning Manifold Gaussian Process

- Manifold Gaussian Process
- Active Learning Cohn
- Numerical results

Performance Metric

- Performance is measured by the **Standardized Root Mean Squared Predictive Error (RMSPE)**.

$$\text{standardized RMSPE} = \left(\sqrt{\frac{1}{N} \sum_i (\hat{y}_i - y_i)^2} \right) / \text{std. dev.}(y_i)$$

- \hat{y}_i is the predicted value at x_i
- y_i is the true value at x_i
- $\sigma(y_i)$ is the standard deviation of all y_i
- The closer the RMSPE is to **0**, the better the model performance.
- Our method (EVIGP) is compared against the R packages: **gpfit**, **mlegp**, and **laGP**.

Example 1: $x \sin(x)$ (Setup)

- **Data:** training data size $n_{train} = 11$, testing data size $n_{test} = 100$, from latin hypercube sampling on $[0, 1]$ and then scaled to $[0, 10]$
- **Noise:** $\epsilon \sim \mathcal{N}(0, 0.5^2)$
- **Priors:** $a_\omega = a_\eta = 1$, $b_\omega = b_\eta = 5$.
- **EVI Settings:** Time step $h = 0.02$, scale $\tau = 1$. Initially, 100 EVI particle uniformly distributed at $[0, 0.1] \times [0.1, 0.4]$
- **Solver:** L-BFGS (PyTorch) with history size 50, max inner iterations 100, and max epoch 500.

Example 1: $x \sin(x)$ (Result)

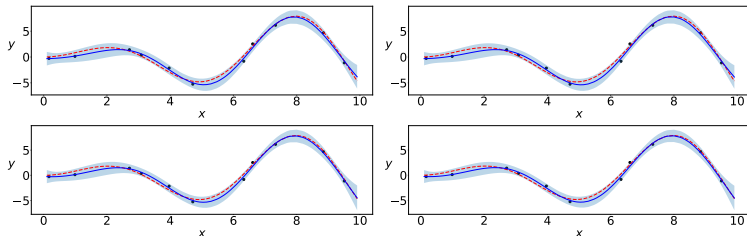


Figure: The gaussian process prediction. Black dots are training data, red curve is the true $x \sin(x)$ and blue curve is the predicted, light blue area shows the predictive confidence interval

Example 1: $x \sin(x)$ (Robustness)

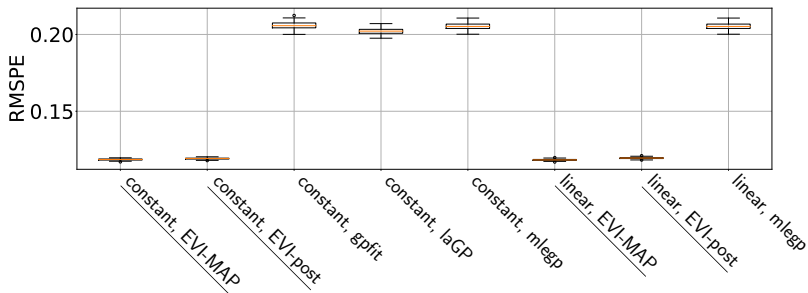


Figure: The mean standardized RMSPE from all different cases and methods as indicated in the label.

Example 2: OTL Circuit (Setup)

- **Function:** A 6-dimensional electronic circuit function.

$$V_m(\mathbf{x}) = \frac{(V_{b1} + 0.74)\beta(R_{c2} + 9)}{\beta(R_{c2} + 9) + R_f} + \frac{11.35R_f}{\beta(R_{c2} + 9) + R_f} \\ + \frac{0.74R_f\beta(R_{c2} + 9)}{(\beta(R_{c2} + 9) + R_f)R_{c1}}, V_{b1} = \frac{12R_{b2}}{R_{b1} + R_{b2}}$$

- **Data:** $n_{\text{train}} = 200$, $n_{\text{test}} = 1000$.
- **Noise:** $\epsilon \sim \mathcal{N}(0, 0.02^2)$.
- **Priors:** $a_{\omega} = a_{\eta} = 1.1$, $b_{\omega} = b_{\eta} = 2.5$.
- **EVI Settings:** Time step $h = 0.001$, scale $\tau = 1$
- **Informative Prior:** A 5-fold cross-validation was used to select a shrinkage parameter $\nu = 2.65$ for variable selection.
- **Solver:** L-BFGS (PyTorch) with history size 50, max inner iterations 100, and max epoch 500

informative prior and shrinkage effect

With an informative prior, we first use cross-validation to choose the parameter ν and then select parameters based on their significance.

- \mathbf{R} is a *diagonal matrix*. And on the diagonal, the values are
 - 1, corresponds to the intercept term in \mathbf{g}
 - $r = 4/3$, linear effect
 - r^2 , quadratic effect,
- ν is searched in a evenly spaced grid from 0 to 5 with 0.05 grid size. A 5-fold cross validation is used to find ν with the smallest standardized RMSPE. The final result is $\nu = 2.65$.

To better illustrate, we rename $R_{b1}, R_{b2}, R_f, R_{c1}, R_{c2}, \beta$ to be x_1, x_2, \dots, x_6 , respectively.

Example 2: OTL circuit (Comparison)

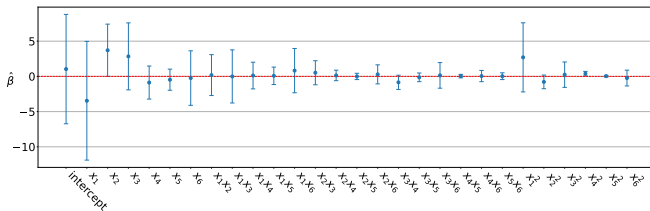
The mean of RMSPE of all methods tested are listed below.

Type of mean	EVIGP	mlegp/lagp
Constant	0.02207	0.04892
Linear	0.01401	0.03593
Quadratic	0.02181	0.03012
Quadratic, after selection	0.01927	N/A

Table: standardized RMSPE

Example 2: OTL circuit (Significance of effects)

We kept only $x_2, x_3, x_4, x_3x_4, x_2^2$ and regress the data again.



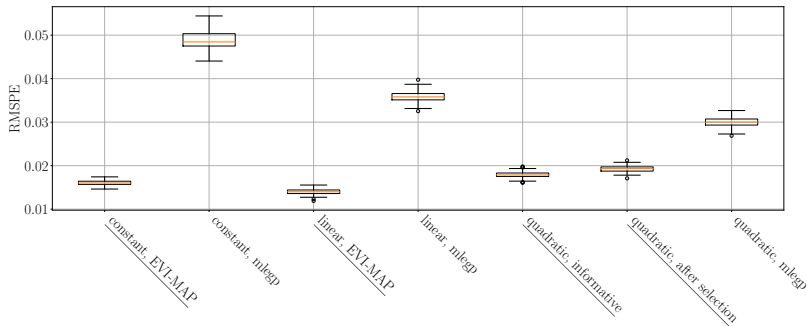


Figure: The mean standardized RMSPE from all different cases and methods as indicated in the label.

Example 3: Borehole Function (Setup)

- **Problem:** An 8-dimensional function modeling water flow through a borehole.

$$f(\mathbf{x}) = \frac{2\pi T_u(H_u - H_l)}{\ln(r/r_\omega) \left(1 + \frac{2LT_u}{\ln(r/r_\omega)r_\omega^2 K_\omega} + \frac{T_u}{T_l}\right)}$$

- **Priors:** $a_\omega = a_\eta = 1$, $b_\omega = b_\eta = 4$.
- **Informative Prior:** Cross-validation resulted in shrinkage parameter $\nu = 0.95$.
- variables are also renamed as x_1, x_2, \dots, x_8 .
- The rest settings are inherited from the OTL example.

Type of mean	EVIGP	mlegp/lagp
Constant	0.0108	0.3349
Linear	0.04190	0.1808
Quadratic	0.03182	0.2022
Quadratic, after selection	0.02310	N/A

Table: standardized RMSPE

Example 3: Borehole Function (Robustness)

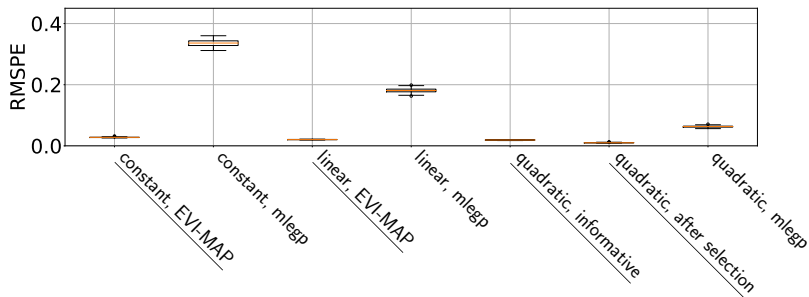


Figure: The mean standardized RMSPE from all different cases and methods as indicated in the label.

Outline

1 Introduction

2 EVIGP

- Gaussian Process Regression
- Bayesian Approach
- Energetic Variational Inference
- Numerical results

3 Active Learning Manifold Gaussian Process

- **Manifold Gaussian Process**
- Active Learning Cohn
- Numerical results

Manifold Gaussian Process

The manifold Gaussian Process (mGP) framework decomposes the regression task into: $F = G \circ M$

- $M : \mathcal{X} \rightarrow \mathcal{H}$ is a deterministic mapping, $\mathcal{H} \subset \mathbb{R}^Q$ latent feature space
- $G : \mathcal{H} \rightarrow \mathcal{Y}$ is a GPR model on \mathcal{H}

Specifically,

$$\mathbf{Y}_n \sim \mathcal{N}(0, \tau^2(\tilde{\mathbf{K}}_n + \eta \mathbf{I}_n)),$$

where $\tilde{K}_{n,ij} = k(M(\mathbf{x}_i), M(\mathbf{x}_j))$. Mapping M can be parameterized via a neural net with parameters θ_M , optimized jointly with GP hyperparameters.

Posterior distribution

Correspondingly,

$$(y(\mathbf{x})|\mathbf{y}, \tau^2, \eta, \boldsymbol{\theta}) = \tilde{\mathbf{k}}_n(\mathbf{x})'(\tilde{\mathbf{K}}_n + \hat{\eta}\mathbf{I}_n)^{-1}\mathbf{y}, \quad (15)$$

$$s_n^2(\mathbf{x}) \triangleq \text{Var} [y(\mathbf{x})|\mathbf{y}, \tau^2, \eta, \boldsymbol{\theta}] \quad (16)$$

$$= \hat{\tau}^2 \left[1 - \tilde{\mathbf{k}}(\mathbf{x})'(\tilde{\mathbf{K}}_n + \hat{\eta}\mathbf{I}_n)^{-1}\tilde{\mathbf{k}}(\mathbf{x}) \right] + \hat{\sigma}^2, \quad (17)$$

And the negative log likelihood is:

$$\text{NLML}(\boldsymbol{\theta}_{\text{mGP}}) = n \log \tau^2 + \log \det(\tilde{\mathbf{K}}_n + \eta\mathbf{I}_n) + \frac{1}{\tau^2} \mathbf{y}^\top (\tilde{\mathbf{K}}_n + \eta\mathbf{I}_n)^{-1} \mathbf{y},$$

Neural Network Architecture

In each layer, the transformation is defined as

$$\mathbf{Z}_i = T_i(\mathbf{X}) = \sigma_M(\mathbf{W}_i \mathbf{Z}_{i-1} + \mathbf{B}_i)$$

- σ_M is the activation function
- \mathbf{W}_i and \mathbf{B}_i are the weight matrix and bias vector for the i -th layer

$$M(\mathbf{X}) = T_l \circ T_{l-1} \circ \cdots \circ T_1(\mathbf{Z}_0)$$

with $\mathbf{Z}_0 = \mathbf{X}$.

A batch normalizing layer is also applied so that kernel values have

- kernel values
- well-behaved gradients

Activation function

Log-sigmoid:

$$\log\text{-}\sigma(x) := \log\left(\frac{1}{1+e^{-x}}\right) = -\log(1+e^{-x}).$$

Its derivative:

$$\frac{d}{dx} \log\text{-}\sigma(x) = \frac{1}{1+e^x} = \sigma(-x),$$

- smoother transition
- not magnitude explosion

A combined loss

- reconstruction loss from the autoencoder architecture
- L_2 regularization

Outline

1 Introduction

2 EVIGP

- Gaussian Process Regression
- Bayesian Approach
- Energetic Variational Inference
- Numerical results

3 Active Learning Manifold Gaussian Process

- Manifold Gaussian Process
- Active Learning Cohn**
- Numerical results

Integrated Mean-Squared Error

Active Learning Cohn selects new data points \mathbf{x}_{n+1} that minimize the model's integrated mean-squared error (IMSE) over the input space.

$$\check{s}_{n+1}^2(\mathbf{x}|\mathbf{x}_{n+1}) \triangleq \hat{\tau}^2 \left[1 - \tilde{\mathbf{k}}_{n+1}(\mathbf{x})^\top (\tilde{\mathbf{K}}_{n+1} + \eta \mathbf{I}_{n+1})^{-1} \tilde{\mathbf{k}}_{n+1}(\mathbf{x}) \right],$$

$$\text{IMSE}(\mathbf{x}_{n+1}) \triangleq \int_{\mathcal{X}} \check{s}_{n+1}^2(\mathbf{x}|\mathbf{x}_{n+1}) \mathrm{d}\mathbf{x},$$

- \mathbf{x}_{n+1} is selected from a finite candidate pool $\mathcal{X}_{\text{cand}} \subset \mathcal{X}$.
- $\tilde{\mathbf{k}}_{n+1}(\mathbf{x}) = [\tilde{k}(\mathbf{x}, \mathbf{x}_1), \dots, \tilde{k}(\mathbf{x}, \mathbf{x}_{n+1})]^\top$ $\tilde{\mathbf{K}}_{n+1}$ denote the kernel vector and $(n+1) \times (n+1)$ correlation matrix evaluated in the learned latent space $\mathcal{H} = M(\mathcal{X})$ using the current mGP parameters.

Active Learning Cohn

Minimizing the IMSE is equivalent to maximizing the expected reduction in posterior variance, thus

$$\begin{aligned} \text{ALC}(\mathbf{x}_{n+1}) &= \int_{\mathcal{X}} s_n^2(\mathbf{x}) - \check{s}_{n+1}^2(\mathbf{x}|\mathbf{x}_{n+1}) \, d\mathbf{x} \\ &\propto \int_{\mathcal{X}} \hat{\tau}^2 \tilde{\mathbf{k}}_{n+1}(\mathbf{x})^\top (\tilde{\mathbf{K}}_{n+1} + \eta \mathbf{I}_{n+1})^{-1} \tilde{\mathbf{k}}_{n+1}(\mathbf{x}) \, d\mathbf{x}, \end{aligned}$$

where $s_n^2(\mathbf{x})$ is computed via (16) except that k function is replaced by \tilde{k} , and thus $\mathbf{k}_n(\mathbf{x})$ by $\tilde{\mathbf{k}}_n(\mathbf{x})$ and \mathbf{K}_n by $\tilde{\mathbf{K}}_n$.

ALC Approximation

Numerically, we introduce a fixed reference set \mathcal{X}_{ref} of size m sampled from \mathcal{X} using Latin Hypercube Design. The ALC objective becomes:

$$\text{ALC}(\mathbf{x}_{n+1}) \propto \sum_{\mathbf{x} \in \mathcal{X}_{\text{ref}}} \hat{\tau}^2 \tilde{\mathbf{k}}_{n+1}(\mathbf{x})^\top (\tilde{\mathbf{K}}_{n+1} + \eta \mathbf{I}_{n+1})^{-1} \tilde{\mathbf{k}}_{n+1}(\mathbf{x}).$$

At each active learning iteration, the most informative training point is selected via:

$$\mathbf{x}_{n+1}^* = \arg \max_{\mathbf{x} \in \mathcal{X}_{\text{cand}}} \text{ALC}(\mathbf{x}),$$

and subsequently removed from the candidate pool $\mathcal{X}_{\text{cand}}$.

Improvement

- batch acquisition: acquire B design points per iteration
- pre-screening: shrinking the candidate set, leaving only those with high predictive variance
- 1 among the candidate set $\mathcal{X}_{\text{cand}}$, we identify the top $K > B$ points with the highest predictive variance, denoted \mathcal{X}_r^* :

$$\mathcal{X}_{\text{cand}}^{(r)} = \left\{ \mathbf{x}_1^{(r)}, \mathbf{x}_2^{(r)}, \dots, \mathbf{x}_K^{(r)} \right\},$$

- 2 compute ALC on $\mathcal{X}_{\text{cand}}^{(r)}$.
- 3 select the B points with the highest ALC scores

Algorithm: ALmGP (Part 1: Setup)

Algorithm 1: Active Learning for Manifold Gaussian Process

- 1: **Input:** Initial training dataset $\{\mathbf{x}_i, y_i\}_{i=1}^{n_0}$; reference set \mathcal{X}_{ref} and candidate set $\mathcal{X}_{\text{cand}}$; tuning parameters $K, B, \text{To1}, N_{\text{max}}$; and initial mGP settings.
- 2: **Output:** Trained mGP model and estimated parameters.
- 3: **Step 0:** Fit the mGP model using the initial training set $\{\mathbf{x}_i, y_i\}_{i=1}^{n_0}$.

Algorithm: ALmGP (Part 2: Main Loop)

Algorithm 1: Active Learning for Manifold Gaussian Process (Cont.)

```

4: for  $r = 1, 2, \dots, \lfloor N_{\max}/B \rfloor$  do
5:   Step 1: Screen candidates: Select the top  $K$  points with the
      highest posterior variance from  $\mathcal{X}_{\text{cand}}$ .
6:   Step 2: Score candidates: Compute the ALC acquisition
      score for each screened point.
7:   Step 3: Acquire new data: Select top  $B$  points with the
      highest ALC scores and add to the training set.
8:   Step 4: Update model: Re-train or update the mGP model
      with the expanded training set.
9:   Step 5: Check convergence: If error is below  $\text{To1}$ , terminate.
10: end for

```

Outline

1 Introduction

2 EVIGP

- Gaussian Process Regression
- Bayesian Approach
- Energetic Variational Inference
- Numerical results

3 Active Learning Manifold Gaussian Process

- Manifold Gaussian Process
- Active Learning Cohn
- **Numerical results**

Experimental Setup

- **Optimizer:** L-BFGS with strong Wolfe line search from PyTorch.
- **Stability:** Model hyperparameters are squared to ensure positivity.
- **Metric:** Root Mean Square Error (RMSE) on an independent test set.

$$\text{RMSE} = \sqrt{\frac{1}{n_{\text{test}}} \sum_i (\hat{y}_i - y_i)^2}$$

- **Robustness:** Each experiment is repeated 10 times.
- **Baseline:** Performance is compared against random point acquisition.
- **Data generation:** Latin Hypercube Design

Example 1: Piecewise Trigonometric Function (Setup)

- **Function:** A 1D function with three distinct regions and Gaussian noise $\epsilon \sim \mathcal{N}(0, 0.1^2)$.

$$F(x) = \begin{cases} 1.35 \cos(12\pi x), & x \in [0, 0.33] \\ 1.35, & x \in [0.33, 0.66] \\ 1.35 \cos(6\pi x), & x \in [0.66, 1] \end{cases}$$

- **Parameters:**
 - Initial data: $n_0 = 10$
 - Budget: $N_{\max} = 20$ (with batch size $B = 1$)
 - Neural Network: [1-6-2] architecture

Example 2: 2D Deterministic Function (Setup)

- **Function:** A 2D function composed of Gaussian PDFs and a linear term, rotated by 45° .

$$f(x_1, x_2) = 1 - \phi(x_2; 3, 0.5^2) - \phi(x_2; -3, 0.5^2) + \frac{x_1}{100}$$

- **Parameters:**
 - Initial data: $n_0 = 50$
 - Budget: $N_{\max} = 50$ (with batch size $B = 1$)
 - Neural Network: [2-10-3] architecture

Example 3: Function on the 3D Sphere (Setup)

- **Function:** A function defined on the surface of a 3D unit sphere ($x^2 + y^2 + z^2 = 1$).

$$f(x, y, z) = \cos(x) + y^2 + e^z$$

- **Parameters:**
 - Initial data: $n_0 = 50$
 - Budget: $N_{\max} = 100$ (with batch size $B = 1$)
 - Neural Network: [3-10-2] architecture

Example 3: Function on the 3D Sphere (Results)

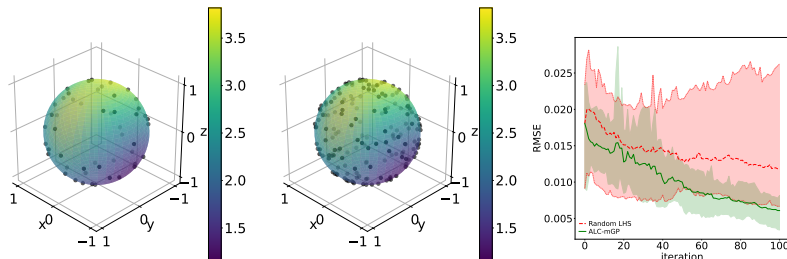


Figure: True function on the sphere, final prediction, and RMSE comparison.

The final average RMSE was 6×10^{-3} .

- **Function:** A well-known 8-dimensional benchmark function that models groundwater flow through a borehole.
- **Challenge:** A complex, high-dimensional, and nonlinear regression problem.
- **Parameters:**
 - Initial data: $n_0 = 50$
 - Budget: $N_{\max} = 150$ (with batch size $B = 1$)
 - Neural Network: [8-30-4] architecture

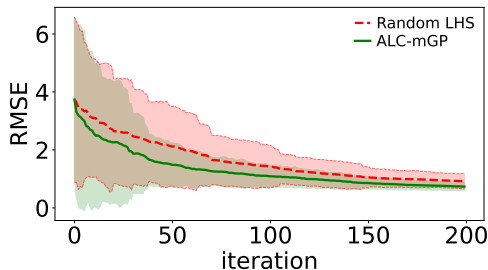


Figure: Test RMSE comparison for the 8D Borehole function.

The final average RMSE was **0.605**.

Conclusion & Future Work

The methods developed in this dissertation make GP modeling more adaptive to structural constraints and more scalable for scientific surrogate modeling. For future research directions:

- **For EVI:** Future research can explore adaptive kernels, hierarchical priors, and connections to Wasserstein gradient flows to improve flexibility.
- **For Active Learning:** We can integrate manifold regularization into the autoencoder or develop an end-to-end framework that co-optimizes for uncertainty and geometry.
- **Unified System:** A powerful next step is to combine both contributions into a single system that dynamically evolves its inference and acquisition strategies, unlocking new applications in autonomous experimentation.

Acknowledgements

I would like to express my sincere gratitude to the following individuals and institutions for their invaluable support throughout my Ph.D. journey:

- **My Advisors:** Dr. Lulu Kang and Dr. Chun Liu, for their exceptional guidance, mentorship, and unwavering support.
- **My Committee:** Thank you to the members of my dissertation committee for their insightful feedback and encouragement.
- **IIT, OSG:** For providing the resources that made this research possible
- **Collaborators & Colleagues:** A special thanks to my peers and collaborators for the stimulating discussions and support.
- **Family & Friends:** Finally, I am deeply grateful to my family and friends for their constant love and encouragement.

Thank You

Questions?