# Energetic Variational Gaussian Process Regression for Computer Experiments

Lulu Kang, Yuanxing Cheng, Yiwei Wang, Chun Liu

Department of Applied Mathematics, Illinois Institute of Technology, Chicago, IL 60616

## Abstract

Gaussian process (GP) regression model is a popular surrogate modeling approach for computer experiments. It provides both accurate prediction and inference for the computer experimental data. Both frequentist and Bayesian framework can be used to perform the estimation and inference tasks for GP. In this chapter, we build the GP model via variational inference, or more specifically, the newly proposed energetic variational inference by [1]. Following the GP model assumption, we first derive the posterior distributions of the parameters for GP model with non-zero mean functions. The energetic variational inference method is used to generate samples of the posterior distributions. More importantly, it can bridge the Bayesian sampling and optimization and provides a much more computationally efficient solution. With a normal prior on the mean component of the GP model, shrinkage estimation is also applied to the parameters to achieve variable selection of the mean function. The performance of the proposed GP model is illustrated via several examples.

**Keyword:** Gaussian process regression, Energetic variational approach, Bayesian inference

## Gaussian Process Regression: Bayesian Approach

▶ Denote $\{x_i, y_i\}_{i=1}^{n}$ as the data from computer experiments, where $x_i \in \Omega \subseteq \mathbb{R}^d$ are the $i$th experimental input values and $y_i \in \mathbb{R}$ is the corresponding output. Here we consider the case of univariate response.

▶ Gaussian process regression is built on the following model assumption of the response,

$$y_i = g(x_i)^\top \beta + Z(x_i) + \epsilon_i, \quad i = 1, \ldots, n, \tag{1}$$

- $g(x)$ is a $p$-dim vector of user-specified basis functions and $\beta$ are the linear coefficients combining the basis functions.
- The noise $\epsilon_i$'s are IID following $N(0, \sigma^2)$ and also independent of the other stochastic ingredients of (1).
- The Gaussian process prior means that the stochastic function $Z(x)$ is a Gaussian process $Z(\cdot) \sim GP(0, \tau^2 K)$, which indicates $\mathbb{E}[Z(x)] = 0$ and

$$\text{Cov}\left[Z(x_1), Z(x_2)\right] = \tau^2 K(x_1, x_2).$$

- Here $\tau^2$ 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 $\omega$.
- For any Gaussian process, the correlation function should be a symmetric positive definite kernel function. For instance, if we use the Gaussian kernel function, then

$$K(x_1, x_2; \omega) = \exp\left\{-\sum_{j=1}^{d} \omega_j (x_{1j} - x_{2j})^2\right\},$$

with $\omega \in \mathbb{R}^d$ and $\omega \geq 0$.

▶ In terms of response $y$, it follows a Gaussian Process with the following mean and covariance,

$$\mathbb{E}[y(x)] = g(x)^\top \beta, \quad \forall x \in \Omega \tag{2}$$
$$\text{Cov}\left[y(x_1), y(x_2)\right] = \tau^2 K(x_1, x_2; \omega) + \sigma^2 \delta(x_1, x_2), \quad \forall x_1, x_2 \in \Omega, \tag{3}$$
$$= \tau^2 \left[K(x_1, x_2; \omega) + \eta \delta(x_1, x_2)\right], \tag{4}$$

where $\delta(x_1, x_2) = 1$ if $x_1 = x_2$ and 0 otherwise, and $\eta = \sigma^2/\tau^2$.

▶ $\eta$ is interpreted as the noise-to-signal ratio. For deterministic computer experiments, the noise component is not part of the model, i.e., $\sigma^2 = 0$. However, a nugget effect, which is a small $\eta$ value, is usually included in the covariance function to avoid singularity of the covariance matrix. The unknown parameter values are $\theta = (\beta, \omega, \tau^2, \eta)$.

## Prior and Posterior Distribution

We assume the following prior distributions for the parameters,

$$\beta \sim \mathcal{MVN}_p(0, v^2 R) \tag{5}$$
$$\omega_i \overset{\text{i.i.d.}}{\sim} \text{Gamma}(a_\omega, b_\omega), \text{ for } i = 1, \ldots, d \tag{6}$$
$$\tau^2 \sim \text{Inverse-}\chi^2(df_{\tau^2}) \tag{7}$$
$$\eta \sim \text{Gamma}(a_\eta, b_\eta). \tag{8}$$

In practice, using non-informative prior $p(\beta) \propto 1$ also gives promising result.

▶ Based on the data, the sampling distribution is

$$y_n | \theta \sim \mathcal{MVN}_n\left(G\beta, \tau^2(K_n + \eta I_n)\right), \tag{9}$$

where $y_n$ is the vector of $y_i$'s and $G$ is a matrix of rows and each row is $g(x_i)^\top$. The matrix $K_n$ is the $n \times n$ kernel matrix and $K_n[i, j] = K(x_i, x_j)$ and is symmetric and positive definite, and $I_n$ is an $n \times n$ identity matrix.

▶ The pdf of $y_n | \theta$ is

$$p(y_n | \theta) \propto (\tau^2)^{-\frac{n}{2}} \det(K_n + \eta I_n)^{-1/2} \exp\left(-\frac{1}{2\tau^2}(y_n - G\beta)^\top (K_n + \eta I_n)^{-1}(y_n - G\beta)\right). \tag{10}$$

▶ The joint posterior distribution of all parameters is

$$p(\theta | y_n) \propto p(\theta) p(y_n | \theta) \tag{11}$$
$$\propto p(\beta)\left(\prod_{j=1}^{d} p(\omega_i)\right) p(\tau^2) p(\eta) p(y_n | \theta). \tag{12}$$

▶ Due to the conditional conjugacy,

$$\beta | y_n, \omega, \tau^2, \eta \sim \mathcal{MVN}_d\left(\hat{\beta}_n, \Sigma_{\beta|n}\right),$$

where

$$\Sigma_{\beta|n} = \left[\frac{1}{\tau^2} G^\top (K_n + \eta I_n)^{-1} G + \frac{1}{v^2} R^{-1}\right]^{-1}$$
$$\hat{\beta}_n = \Sigma_{\beta|n} \frac{(G^\top (K_n + \eta I_n)^{-1} y_n)}{\tau^2}$$

▶ if with non-informative prior for $\beta$, we have

$$\Sigma_{\beta|n} = \tau^2 \left[G^\top (K_n + \eta I_n)^{-1} G\right]^{-1}$$
$$\hat{\beta}_n = \left[G^\top (K_n + \eta I_n)^{-1} G\right]^{-1} \left[G^\top (K_n + \eta I_n)^{-1}\right] y_n$$

## Estimate $\tau$

$$p(\omega, \tau^2, \eta | y_n) \propto \int p(\beta) p(y_n | \theta) p(\omega) p(\tau^2) p(\eta) \, d\beta \tag{13}$$
$$\propto \det(\Sigma_{\beta|n})^{1/2} \exp\left[-\frac{1}{2}\hat{\beta}_n^\top \Sigma_{\beta|n}^{-1} \hat{\beta}_n - \frac{1}{2\tau^2} y_n^\top (K_n + \eta I_n)^{-1} y_n\right] (\tau^2)^{-n/2} \det(K_n + \eta I_n)^{-1/2} p(\tau^2) p(\omega) p(\eta)$$

Define

$$s_n^2 = \tau^{-2} \hat{\beta}_n^\top \Sigma_{\beta|n}^{-1} \hat{\beta}_n + y_n^\top (K_n + \eta I_n)^{-1} y_n$$
$$= y_n^\top (K_n + \eta I_n)^{-1} \left[G\left(G^\top (K_n + \eta I_n)^{-1} G\right)^{-1} G^\top + (K_n + \eta I_n)\right] (K_n + \eta I_n)^{-1} y_n.$$

Then $p(\omega, \tau^2, \eta | y_n)$ is

$$p(\omega, \tau^2, \eta | y_n) \propto (\tau^2)^{-\left[\frac{1}{2}(df_{\tau^2} + n - p) + 1\right]} \exp\left(-\frac{1 + s_n^2}{2\tau^2}\right) \det(G^\top (K_n + \eta I_n)^{-1} G)^{-1/2} \det(K_n + \eta I_n)^{-1/2} p(\omega) p(\eta).$$

Due to the conditional conjugacy, we can see the conditional posterior distribution of $\tau^2 | \omega, \eta, y_n$ is
$\tau^2 | \omega, \eta, y_n \sim \text{Scaled Inverse-}\chi^2(df_{\tau^2} + n - p, \hat{\tau}^2)$ where $\hat{\tau}^2 = (1 + s_n^2)/(df_{\tau^2} + n - p)$.

## Posterior Distribution

Using the informative prior $\beta \sim \mathcal{MVN}_p(0, v^2 R)$, we obtain the posterior distribution (14) in which $\beta$ has been integrated out, and thus it saves the computation involved in posterior sampling stage.

$$p(\omega, \tau^2, \eta | y_n) \propto \det(\Sigma_{\beta|n})^{1/2} \exp\left[-\frac{1}{2}\hat{\beta}_n^\top \Sigma_{\beta|n}^{-1} \hat{\beta}_n + \frac{1}{2\tau^2} y_n^\top (K_n + \eta I_n)^{-1} y_n\right] \det(K_n + \eta I_n)^{-1/2}$$
$$\times (\tau^2)^{-n/2} p(\tau^2) p(\omega) p(\eta). \tag{14}$$

Using the non-informative prior $p(\beta) \propto 1$, we obtain the posterior distribution (15) in which both $\beta$ and $\tau^2$ are integrated out.

$$p(\omega, \eta | y_n) \propto (\hat{\tau}^2)^{-\frac{1}{2}(df_{\tau^2} + n - p)} \det(G^\top (K_n + \eta I_n)^{-1} G)^{-1/2} \det(K_n + \eta I_n)^{-1/2} p(\omega) p(\eta). \tag{15}$$

Given the parameters $(\omega, \tau^2, \eta)$, the posterior predictive distribution of $y(x)$ at any query point $x$ is the following normal distribution.

$$y(x) | y_n, \omega, \tau^2, \eta \sim \mathcal{N}(\hat{\mu}(x), \sigma_n^2(x)), \tag{16}$$

where

$$\hat{\mu}(x) = g(x)^\top \hat{\beta}_n + K(x, \mathcal{X}_n)(K_n + \eta I_n)^{-1}(y_n - G\hat{\beta}_n), \tag{17}$$
$$\sigma_n^2(x) = \tau^2 \left\{1 - K(x, \mathcal{X}_n)(K_n + \eta I_n)^{-1} K(\mathcal{X}_n, x) + c(x)^\top \left[G^\top (K_n + \eta I_n)^{-1} G\right]^{-1} c(x)\right\}, \tag{18}$$
$$c(x) = g(x) - G^\top (K_n + \eta I_n)^{-1} K(\mathcal{X}_n, x),$$
$$K(x, \mathcal{X}_n) = K(\mathcal{X}_n, x)^\top = [K(x, x_1), \ldots, K(x, x_n)]$$

For the sake of computational stability, we use an alternative formula for $\hat{\mu}(x)$ by first plugging in $\hat{\beta}_n$ into $\hat{\mu}(x)$.

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

We can use cholesky decomposition on $K_n + \eta I_n$ where we define $LL^\top = K_n + \eta I_n$. Then we can obtain higher accuracy in calculating $(K_n + \eta I_n)^{-1} G$ by solving $\gamma_1$ from the following equations.

$$L^\top \alpha = \eta G$$
$$L\gamma_1 = \alpha/\eta \tag{19}$$

Similar method is used on $(K_n + \eta I_n)^{-1} y_n = \gamma_2$. Then the final formula reads,

$$\hat{\mu}(x) = \left((1 - K(x, \mathcal{X}_n)\gamma_1)\left(G^\top \gamma_1\right)^{-1} G^\top + K(x, \mathcal{X}_n)\right)\gamma_2$$

## Sampling procedure and EVI

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

Energetic variational inference method is introduced here. Let $\phi_t$ be the dynamic flow map $\phi_t : \mathbb{R}^d \to \mathbb{R}^d$ that continuously transforms the $d$-dimensional distribution from an initial distribution toward the target one and we require the map $\phi_t$ to be smooth and one-to-one. Then define functional $\mathcal{F}(\phi_t)$ to be (14) or (15). The EVI start with energy dissipation law:

$$\frac{d}{dt}\mathcal{F}(\phi_t) = -\triangle(\phi_t, \dot{\phi}_t), \tag{20}$$

where $\triangle(\phi_t, \dot{\phi}_t)$ is a user-specified functional representing the rate of energy dissipation, and $\dot{\phi}_t$ is the derivative of $\phi_t$ with time $t$. The dissipation functional should satisfy $\triangle(\phi_t, \dot{\phi}_t) \geq 0$ so that $\mathcal{F}(\phi_t)$ decreases with time. As discussed in [?], there are many ways to specify $\triangle(\phi_t, \dot{\phi}_t)$ and the simplest among them is a quadratic functional in terms of $\dot{\phi}_t$,

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

where $\rho_{[\phi_t]}$ denotes the pdf of the current distribution which is the initial distribution transformed by $\phi_t$, $\Omega_t$ is the current support, and $\|a\|_2 = a^\top a$ for $\forall a \in \mathbb{R}^d$. Then the EVI reads:

$$\frac{\delta \frac{1}{2}\triangle}{\delta \dot{\phi}_t} = -\frac{\delta \mathcal{F}}{\delta \phi_t} \implies \rho_{[\phi_t]}\dot{\phi}_t = -\frac{\delta \mathcal{F}}{\delta \phi_t}.$$

## EVI(contin.)

We then do discretization on space parameter, called it "approximation-then-variation" approach and obtain the particle version of the PDE.

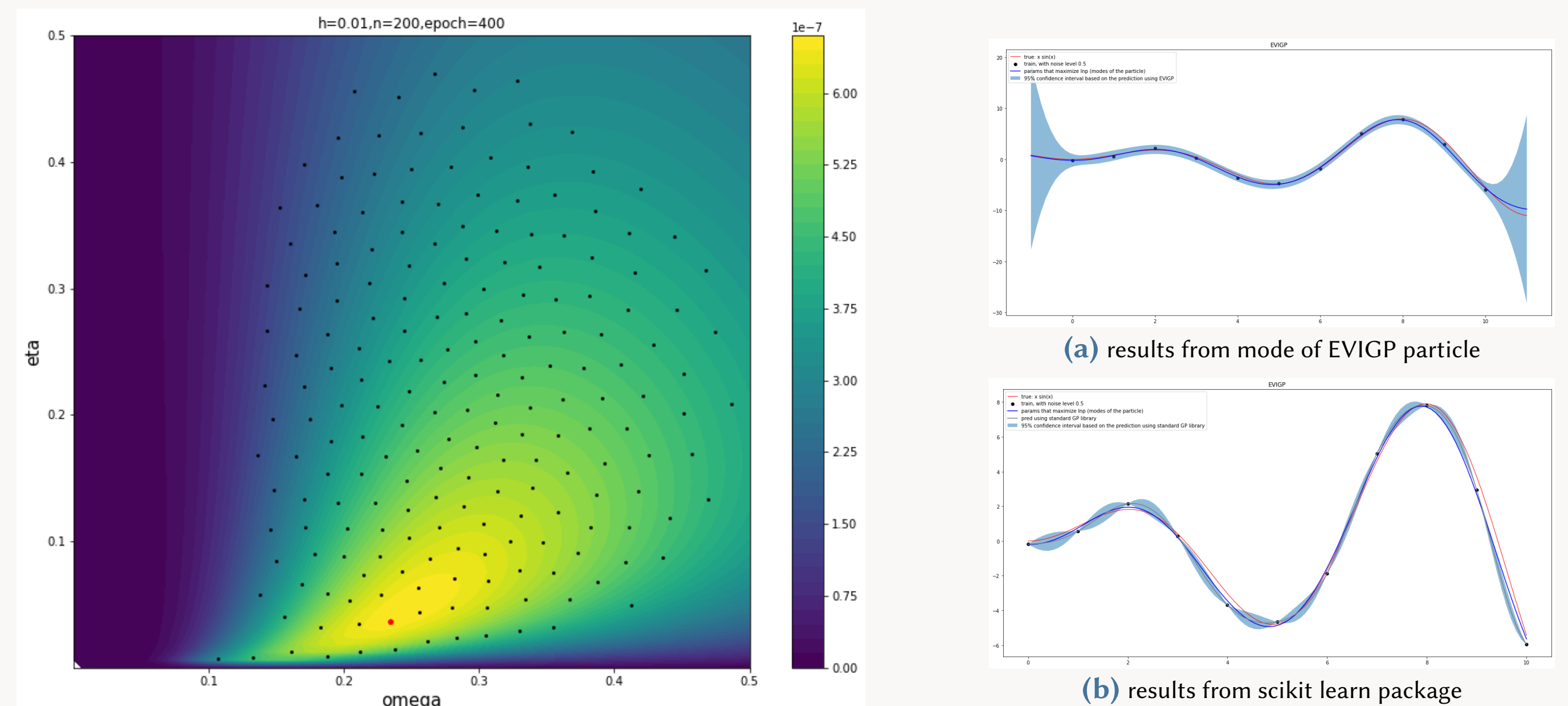$$\frac{d}{dt}\mathcal{F}_h(\{x_i(t)\}_{i=1}^N) = -\triangle_h(\{x_i(t)\}_{i=1}^N, \{\dot{x}_i(t)\}_{i=1}^N). \tag{21}$$

Here $\{x(t)\}_{i=1}^N$ is the locations of $N$ particles at time $t$ and $\dot{x}_i(t)$ is the derivative of $x_i$ with $t$, and thus is the velocity of the $i$th particle as it moves toward the target distribution. The subscript $h$ of $\mathcal{F}$ and $\triangle$ denotes the bandwidth parameter of the kernel function used in the kernelization operation. Applying the variational steps to (21), we obtain the dynamics of decreasing $\mathcal{F}$ at the particle level,

$$\frac{\delta \frac{1}{2}\triangle_h}{\delta \dot{x}_i(t)} = -\frac{\delta \mathcal{F}_h}{\delta x_i}, \quad \text{for } i = 1, \ldots, N. \tag{22}$$

Backward-Euler is then applied to solve this ODE, which gives particles representing the samples $\theta$s from the posterior $\log p(\omega, \tau^2, \eta | y_n)$ or $\log p(\omega, \eta | y_n)$.

In accordance with the frequentist approach, we can then pick the mode from those samples that maximize the posterior probability and make predictions according to (17).

## 2D toy example: $y = x\sin(x)$



(a) results from mode of EVIGP particle
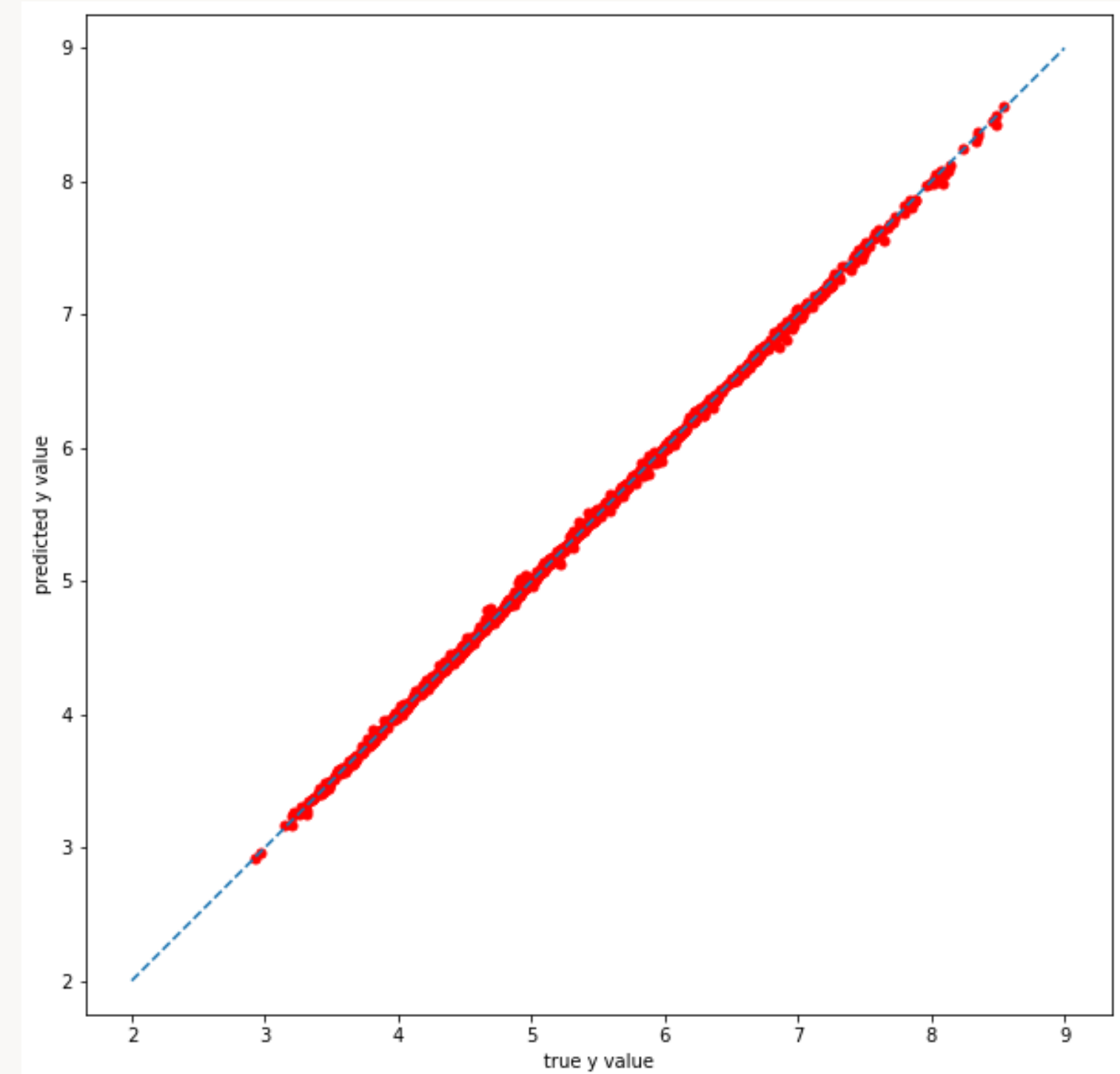


(b) results from scikit learn package

## high dimension function

We compare the regression results for borehole function and OTL-circuit function to those from R package GPfit and mlegp, corresponding to constant mean and linear mean in the GP model, respectively. Here the metric standardized root mean squared percentage error (rmse) is defined as

$$\sqrt{\frac{1}{N_{\text{test}}}\sum\left(y_{\text{pred}}/y_{\text{true}} - 1\right)^2}\Big/\text{std}(y_{\text{true}})$$

And an example of QQ-plot

| function name | borehole | OTC-Circuit |
|---|---|---|
| EVIGP (constant mean) | 0.000267 | 0.00386 |
| gpfit | 0.0224 | 0.0131 |
| EVIGP (linear mean) | 0.000235 | 0.00263 |
| mlegp | 0.0118 | 0.00884 |



## Reference

📄 Wang, Y., Chen, J., Liu, C., and Kang, L.
Particle-based energetic variational inference.
*Statistics and Computing* (04 2020).

Yuanxing Cheng

ycheng46@hawk.iit.edu
Illinoise Institute of Technology

ILLINOIS TECH