

Risk-Assessment, learning and optimization using surrogate models

Paris Perdikaris

Massachusetts Institute of Technology, Department of Mechanical Engineering

Web: <http://web.mit.edu/parisp/www/>

Email: parisp@mit.edu

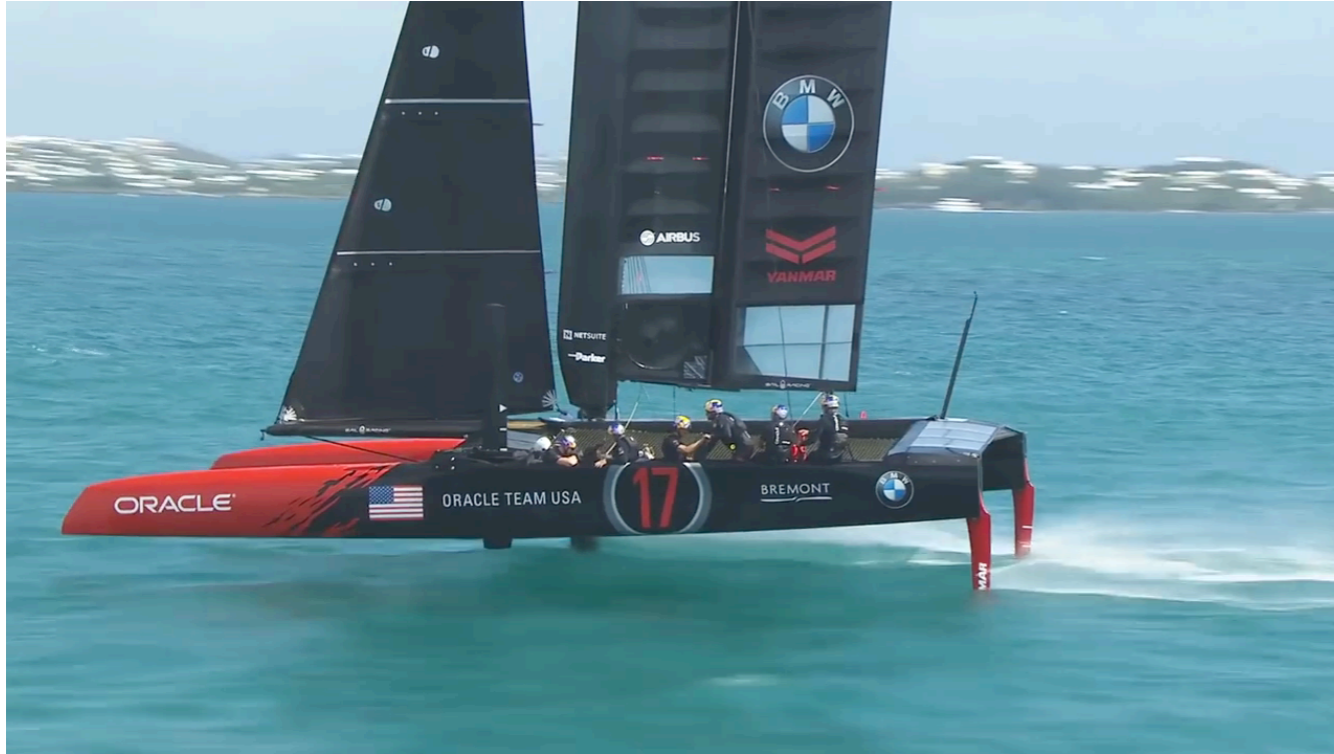
March 1, 2017

SIAM CSE

Atlanta, GA



Motivation

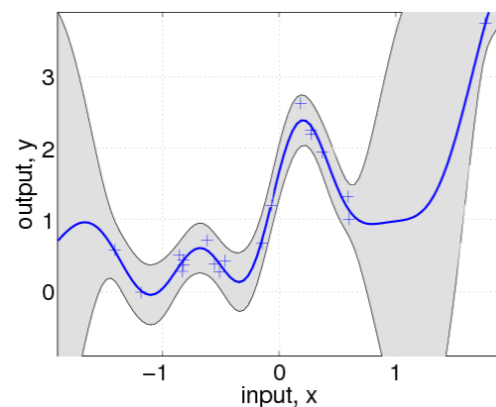


Optimization under uncertainty

$$\boldsymbol{x}, \boldsymbol{\xi} \cdots \rightarrow \blacksquare \cdots \rightarrow Y(\boldsymbol{x}, \boldsymbol{\xi})$$

- ① $L(t) = \{\boldsymbol{x} : \mathcal{R}(Y(\boldsymbol{x}, \boldsymbol{\xi})) \leq t\}$
- ② $\boldsymbol{x}^* = \arg \min_{\boldsymbol{x} \in \mathbb{R}^d} \mathcal{R}(Y(\boldsymbol{x}, \boldsymbol{\xi}))$

- Assess the performance of the system
- Estimate the statistics of the response
- Data acquisition under limited budgets

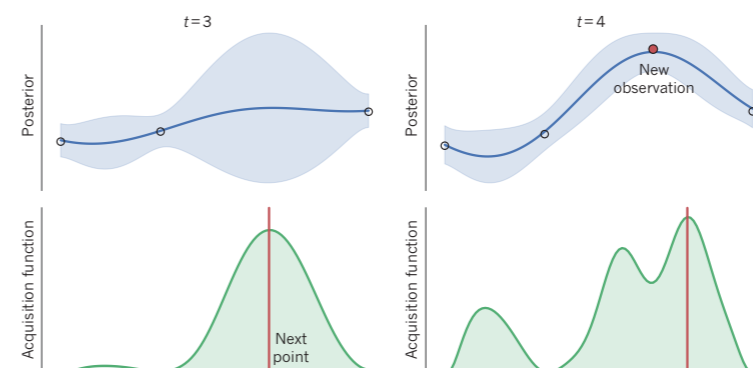


Gaussian process regression

Model & parametric uncertainty

Surrogate uncertainty

Learning uncertainty



Active learning & Bayesian optimization

Gaussian processes

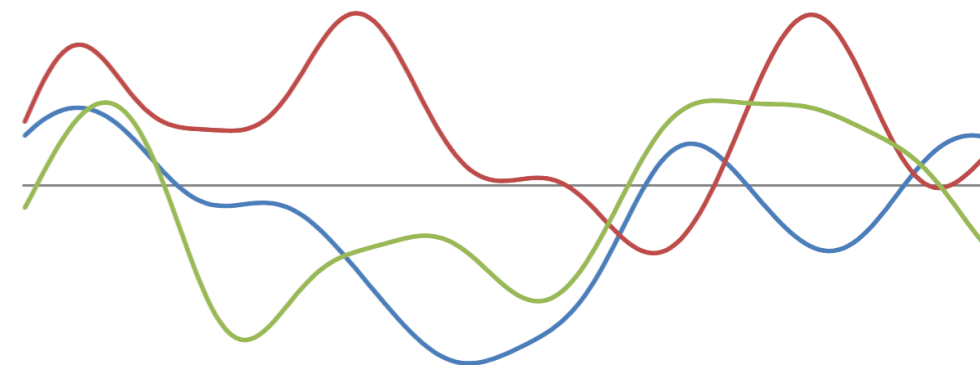
Starting point: The multivariate Gaussian distribution

$$p(\underbrace{f_1, f_2, \dots, f_s}_{\mathbf{f}_A}, \underbrace{f_{s+1}, f_{s+2}, \dots, f_N}_{\mathbf{f}_B}) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}) \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_A \\ \boldsymbol{\mu}_B \end{bmatrix} \quad \text{and} \quad \mathbf{K} = \begin{bmatrix} \mathbf{K}_{AA} & \mathbf{K}_{AB} \\ \mathbf{K}_{BA} & \mathbf{K}_{BB} \end{bmatrix}$$

Generalization: The Gaussian process

$$\boldsymbol{\mu}_\infty = \begin{bmatrix} \mu_f \\ \dots \\ \dots \end{bmatrix} \quad \text{and} \quad \mathbf{K}_\infty = \begin{bmatrix} \mathbf{K}_{ff} & \dots \\ \dots & \dots \end{bmatrix} \quad \mathbf{K}_{i,j} = k(\mathbf{x}_i, \mathbf{x}_j)$$

mean function *covariance function*



Samples from a GP prior

Priors over functions: $f \sim \mathcal{GP}(\mu(x), K(\mathbf{x}, \mathbf{x}'; \theta))$

Infinite dimensional model, but finitely many observations: The marginalization property

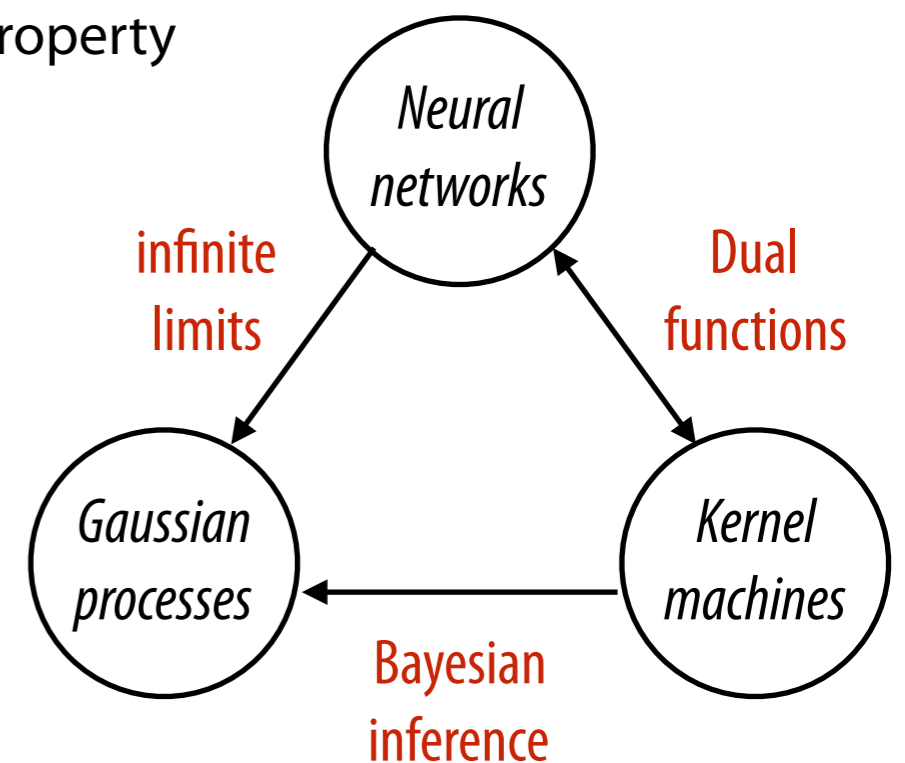
$$p(\mathbf{f}_A, \mathbf{f}_B) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}). \quad \text{Then:}$$

$$p(\mathbf{f}_A) = \int_{\mathbf{f}_B} p(\mathbf{f}_A, \mathbf{f}_B) d\mathbf{f}_B = \mathcal{N}(\boldsymbol{\mu}_A, \mathbf{K}_{AA})$$

Posterior is also Gaussian:

$$p(\mathbf{f}_A, \mathbf{f}_B) \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}). \quad \text{Then:}$$

$$p(\mathbf{f}_A | \mathbf{f}_B) = \mathcal{N}(\boldsymbol{\mu}_A + \mathbf{K}_{AB} \mathbf{K}_{BB}^{-1} (\mathbf{f}_B - \boldsymbol{\mu}_B), \mathbf{K}_{AA} - \mathbf{K}_{AB} \mathbf{K}_{BB}^{-1} \mathbf{K}_{BA})$$



Nonlinear regression with Gaussian processes

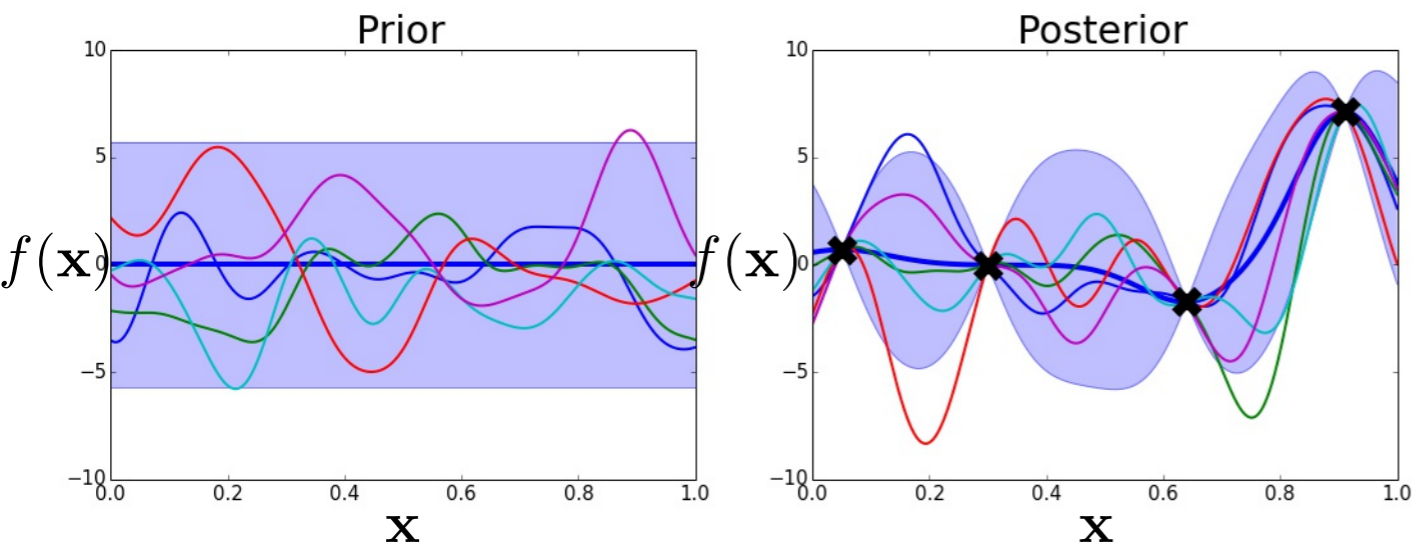
$$y = f(\mathbf{x}) + \epsilon, \quad f \sim \mathcal{GP}(\mu(x), K(\mathbf{x}, \mathbf{x}'; \theta))$$

History:

- Wiener–Kolmogorov filtering (1940)
- Kriging (spatial statistics, 1970)
- GP regression (machine learning, 1996)

Workflow:

- Assign a Gaussian process (GP) prior over functions
- Given a training set of observations (\mathbf{x}, y) calibrate the GP hyper-parameters
- Use the conditional posterior $[f|\mathbf{y}]$ to infer predictions for unobserved \mathbf{x} 's with quantified uncertainty



covariance function

expression

constant

$$\sigma_0^2$$

linear

$$\sum_{d=1}^D \sigma_d^2 x_d x'_d$$

polynomial

$$(\mathbf{x} \cdot \mathbf{x}' + \sigma_0^2)^p$$

squared exponential

$$\exp\left(-\frac{r^2}{2\ell^2}\right)$$

Matérn

$$\frac{1}{2^{\nu-1}\Gamma(\nu)} \left(\frac{\sqrt{2\nu}}{\ell} r\right)^{\nu} K_{\nu}\left(\frac{\sqrt{2\nu}}{\ell} r\right)$$

exponential

$$\exp\left(-\frac{r}{\ell}\right)$$

γ -exponential

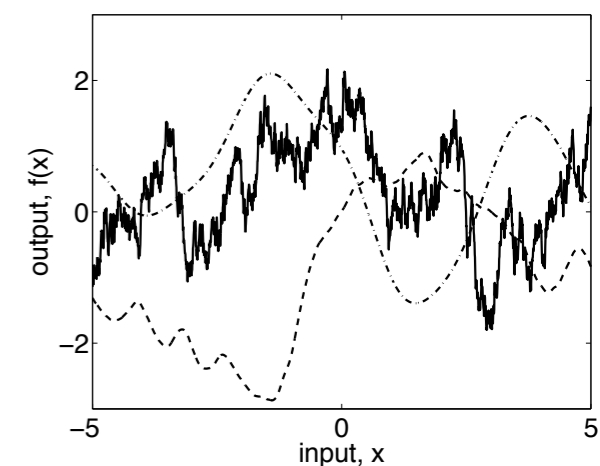
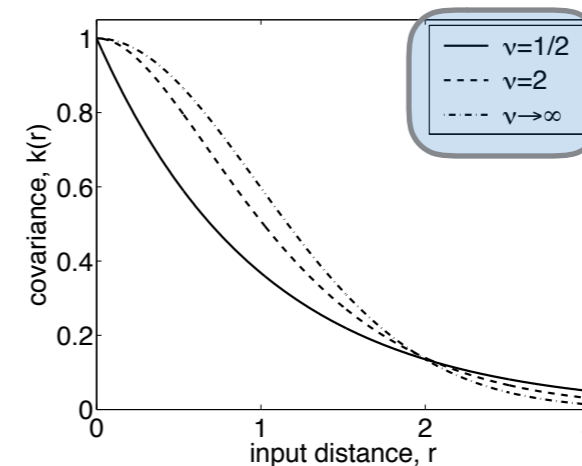
$$\exp\left(-\left(\frac{r}{\ell}\right)^{\gamma}\right)$$

rational quadratic

$$\left(1 + \frac{r^2}{2\alpha\ell^2}\right)^{-\alpha}$$

neural network

$$\sin^{-1}\left(\frac{2\tilde{\mathbf{x}}^{\top}\Sigma\tilde{\mathbf{x}}'}{\sqrt{(1+2\tilde{\mathbf{x}}^{\top}\Sigma\tilde{\mathbf{x}})(1+2\tilde{\mathbf{x}}'^{\top}\Sigma\tilde{\mathbf{x}}')}}\right)$$



Training & prediction

Hyper-parameter estimation:

frequentist approach

The vector of hyper-parameters $\boldsymbol{\theta}$ is determined by maximizing the marginal log-likelihood of the observed data (the so called model evidence), i.e.,

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = -\frac{1}{2} \log |\mathbf{K} + \sigma_\epsilon^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^T (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N}{2} \log 2\pi \quad (8)$$

Bayesian approach

Assign priors over the hyper parameters and marginalize them out using MCMC.

Prediction:

If we consider a Gaussian likelihood $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_\epsilon^2 \mathbf{I})$ then the posterior distribution $p(\mathbf{f}|\mathbf{y}, \mathbf{X})$ is tractable and can be used to perform predictive inference for a new output f_* , given a new input \mathbf{x}_* as

$$p(f_*|\mathbf{y}, \mathbf{X}, \mathbf{x}_*) = \mathcal{N}(f_*|\mu_*, \sigma_*^2), \quad (5)$$

$$\mu_*(\mathbf{x}_*) = \mathbf{k}_{*N} (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y}, \quad (6)$$

$$\sigma_*^2(\mathbf{x}_*) = \mathbf{k}_{**} - \mathbf{k}_{*N} (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{k}_{N*}, \quad (7)$$

where $\mathbf{k}_{*N} = [k(\mathbf{x}_*, \mathbf{x}_1), \dots, k(\mathbf{x}_*, \mathbf{x}_N)]$, $\mathbf{k}_{N*} = \mathbf{k}_{*N}^T$, and $\mathbf{k}_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$. Predictions are computed using the posterior mean μ_* , while prediction uncertainty is quantified through the posterior variance σ_*^2 .

Active learning of level sets

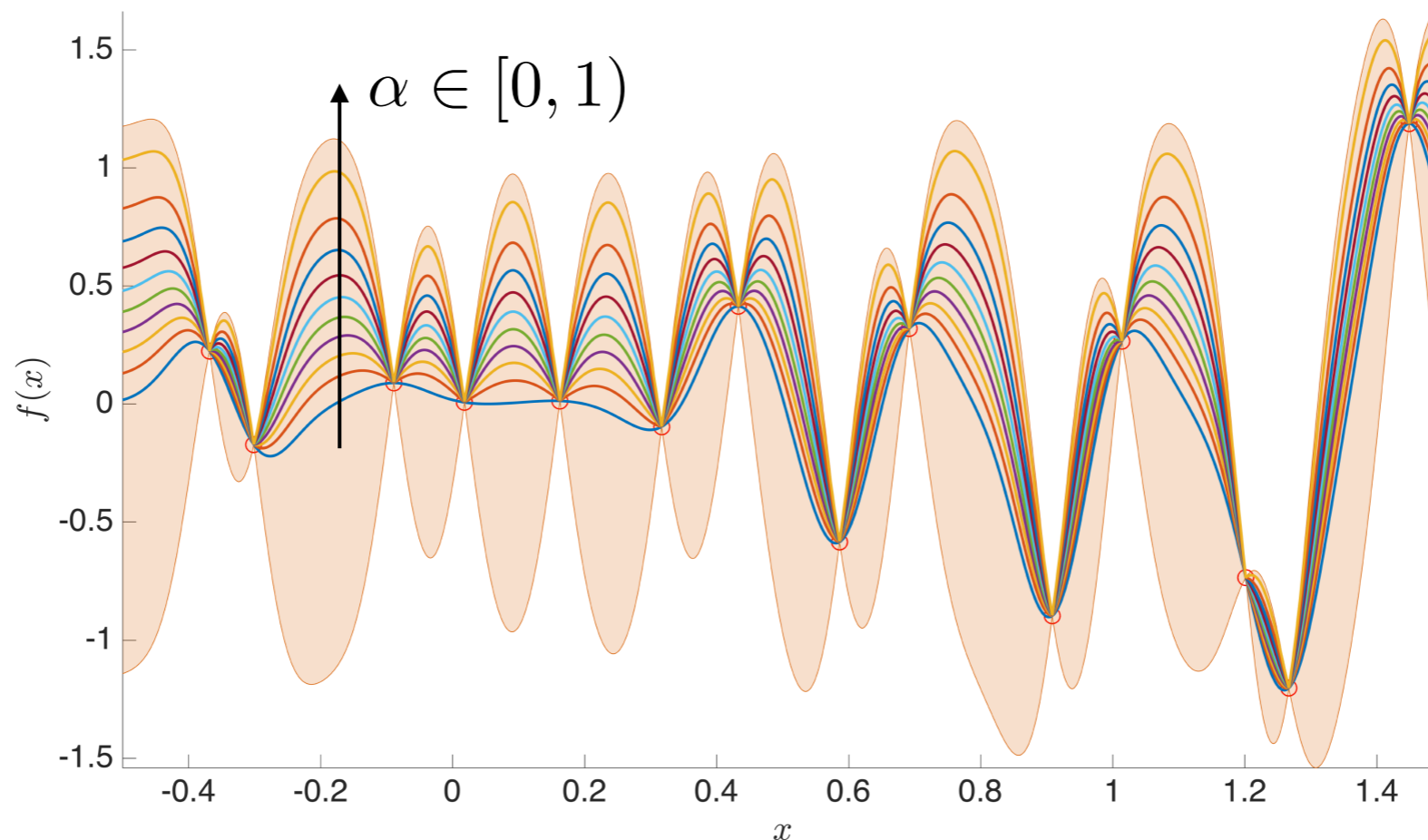
Goal: Identify the sets $L(t) = \{\mathbf{x} : g(\mathbf{x}) \leq t\}$, $g(\mathbf{x}) := \mathcal{R}(Y(\mathbf{x}, \boldsymbol{\xi}))$

Approximate the true objective with a GP surrogate: $g(\mathbf{x}) \approx f(\mathbf{x}) \sim \mathcal{GP}(f|0, k(\mathbf{x}, \mathbf{x}'; \theta))$

$f(x)$ is random so it becomes natural to quantify “surrogate” uncertainty using for e.g. α -superquantile risk measure:

$$\mathcal{R}_\alpha(f(\mathbf{x})) = \min_{c \in \mathbb{R}} \left\{ c + \frac{1}{1-\alpha} \mathbb{E}[\max\{0, f(\mathbf{x}) - c\}] \right\} \quad \begin{array}{l} \text{average of the worst} \\ (1-\alpha)\% \text{ outcomes of } f(\mathbf{x}) \end{array}$$

For $f(x)$ being a GP this can be simplified to: $\mathcal{R}_\alpha(f(\mathbf{x})) = \mu(\mathbf{x}) + \frac{\phi(\Phi^{-1}(\alpha))}{1-\alpha} \sigma(\mathbf{x})$



Active learning of level sets

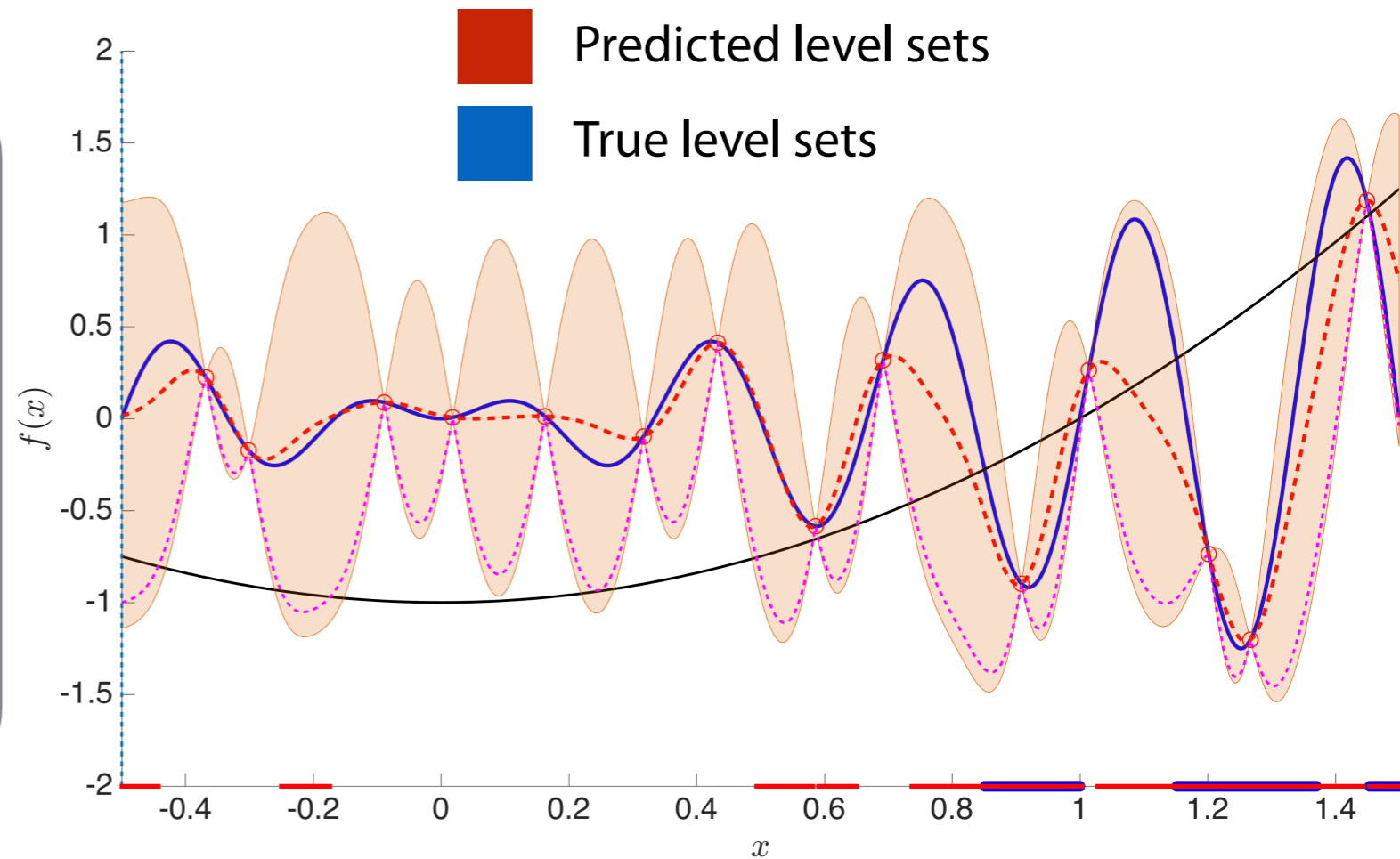
Goal: Identify the sets $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t\}$

Utilize the posterior to guide a sequential sampling policy by optimizing a chosen expected utility function

$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_\theta \mathbb{E}_v |_{\mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

e.g. sample at the locations that maximize the posterior variance in $L(t)$

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x} \in L_\alpha(t)} V(f(\mathbf{x}))$$



Terminate iteration when the “volume” of the predicted level sets is below a given threshold:

$$|V_{n+1}(t) - V_n(t)| < \epsilon, \quad V_n(t) = \int_{L_\alpha(t)} \mathbf{1}_{[-\infty, t]} d\mathbf{x}$$

Remarks:

- The choice of risk-averseness level $\alpha \in [0, 1)$ controls the exploration vs exploitation trade-off.
- Upon convergence the predicted levels sets are guaranteed to be a subset of the true level sets.

Active learning of level sets

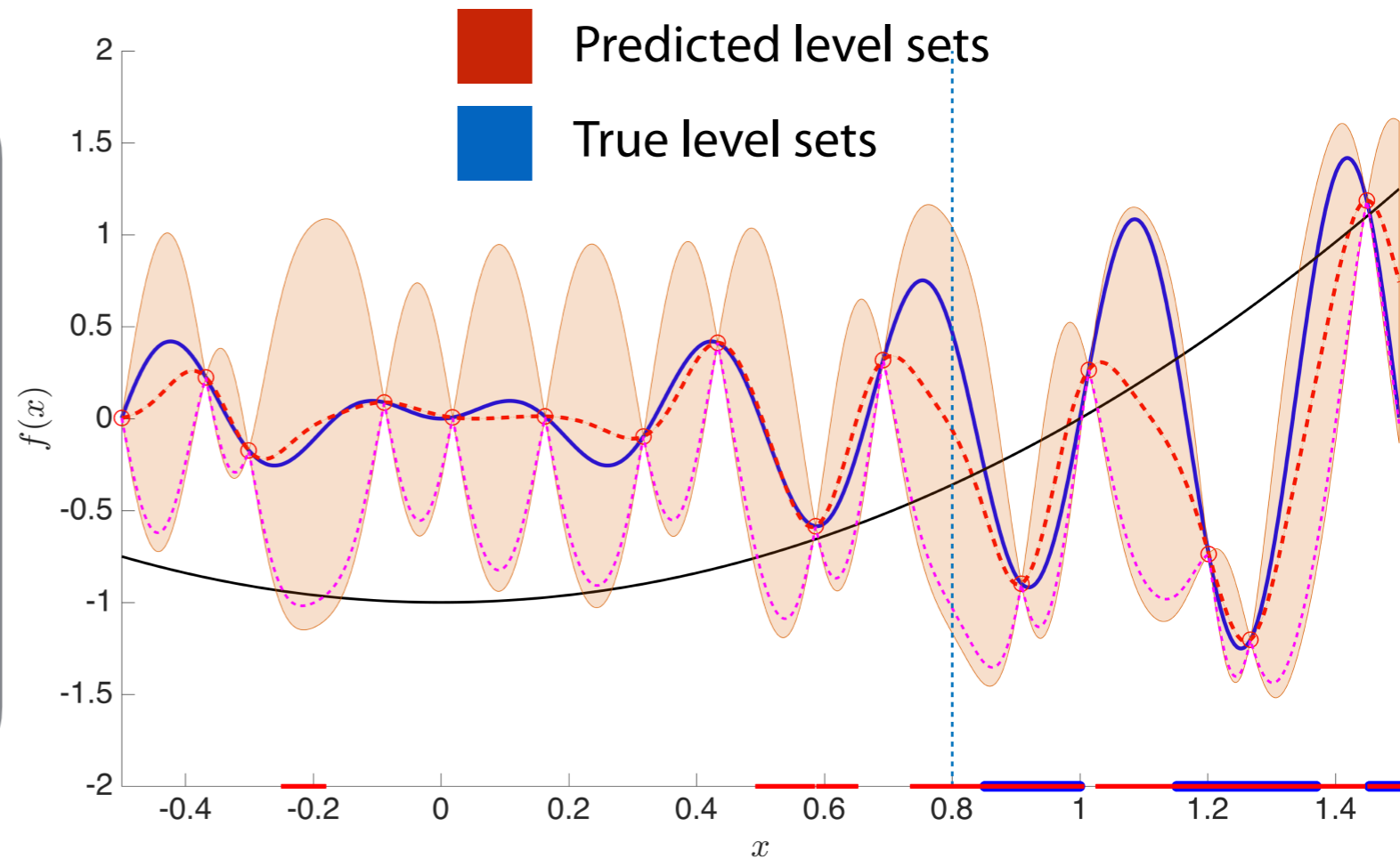
Goal: Identify the sets $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t\}$

Utilize the posterior to guide a sequential sampling policy by optimizing a chosen expected utility function

$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_\theta \mathbb{E}_v |_{\mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

e.g. sample at the locations that maximize the posterior variance in $L(t)$

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x} \in L_\alpha(t)} V(f(\mathbf{x}))$$



Terminate iteration when the “volume” of the predicted level sets is below a given threshold:

$$|V_{n+1}(t) - V_n(t)| < \epsilon, \quad V_n(t) = \int_{L_\alpha(t)} \mathbf{1}_{[-\infty, t]} d\mathbf{x}$$

Remarks:

- The choice of risk-averseness level $\alpha \in [0, 1)$ controls the exploration vs exploitation trade-off.
- Upon convergence the predicted levels sets are guaranteed to be a subset of the true level sets.

Active learning of level sets

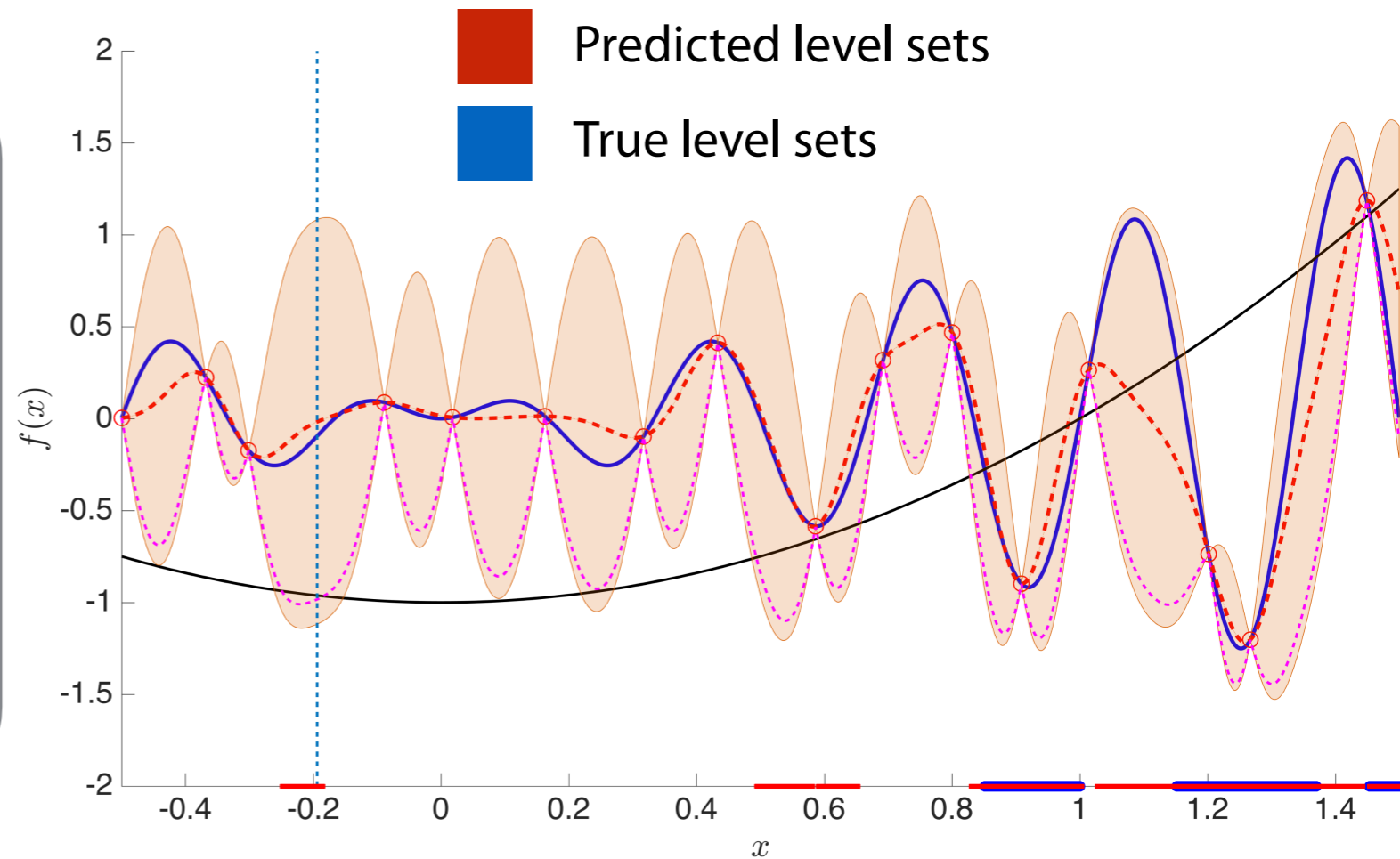
Goal: Identify the sets $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t\}$

Utilize the posterior to guide a sequential sampling policy by optimizing a chosen expected utility function

$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_\theta \mathbb{E}_v |_{\mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

e.g. sample at the locations that maximize the posterior variance in $L(t)$

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x} \in L_\alpha(t)} V(f(\mathbf{x}))$$



Terminate iteration when the “volume” of the predicted level sets is below a given threshold:

$$|V_{n+1}(t) - V_n(t)| < \epsilon, \quad V_n(t) = \int_{L_\alpha(t)} \mathbf{1}_{[-\infty, t]} d\mathbf{x}$$

Remarks:

- The choice of risk-averseness level $\alpha \in [0, 1)$ controls the exploration vs exploitation trade-off.
- Upon convergence the predicted levels sets are guaranteed to be a subset of the true level sets.

Active learning of level sets

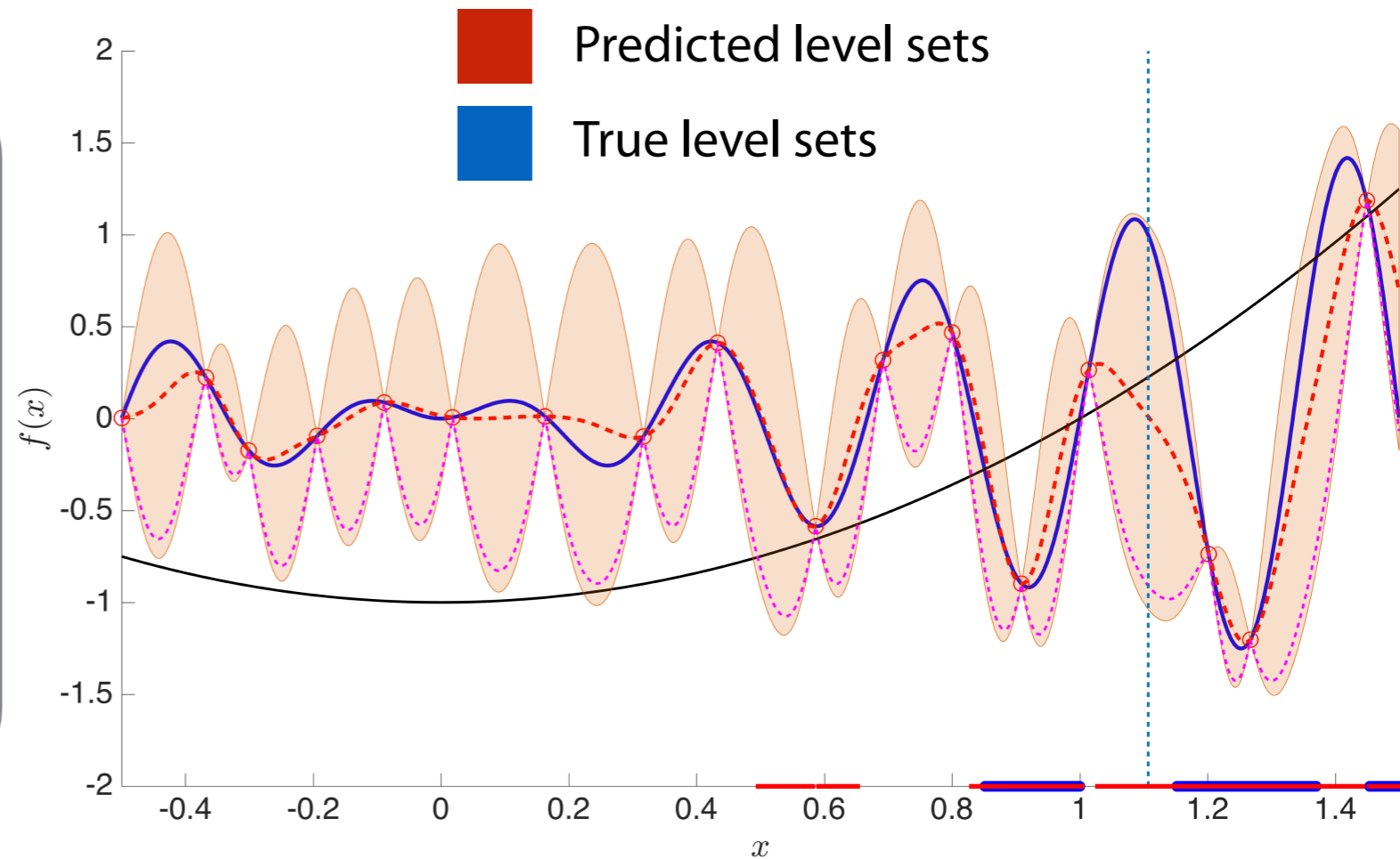
Goal: Identify the sets $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t\}$

Utilize the posterior to guide a sequential sampling policy by optimizing a chosen expected utility function

$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_\theta \mathbb{E}_v |_{\mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

e.g. sample at the locations that maximize the posterior variance in $L(t)$

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x} \in L_\alpha(t)} V(f(\mathbf{x}))$$



Terminate iteration when the “volume” of the predicted level sets is below a given threshold:

$$|V_{n+1}(t) - V_n(t)| < \epsilon, \quad V_n(t) = \int_{L_\alpha(t)} \mathbf{1}_{[-\infty, t]} d\mathbf{x}$$

Remarks:

- The choice of risk-averseness level $\alpha \in [0, 1)$ controls the exploration vs exploitation trade-off.
- Upon convergence the predicted levels sets are guaranteed to be a subset of the true level sets.

Active learning of level sets

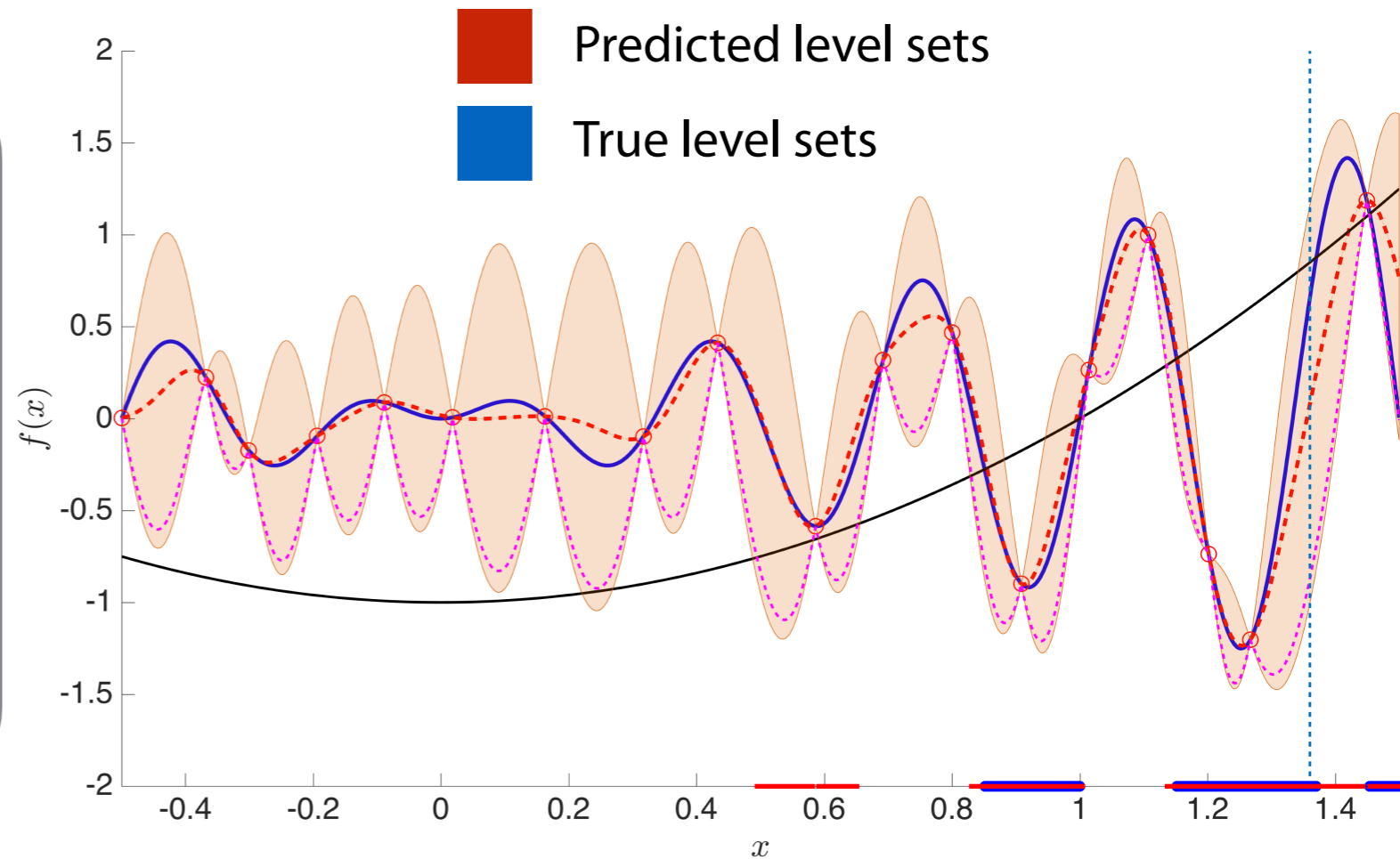
Goal: Identify the sets $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t\}$

Utilize the posterior to guide a sequential sampling policy by optimizing a chosen expected utility function

$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_\theta \mathbb{E}_v |_{\mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

e.g. sample at the locations that maximize the posterior variance in $L(t)$

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x} \in L_\alpha(t)} V(f(\mathbf{x}))$$



Terminate iteration when the “volume” of the predicted level sets is below a given threshold:

$$|V_{n+1}(t) - V_n(t)| < \epsilon, \quad V_n(t) = \int_{L_\alpha(t)} \mathbf{1}_{[-\infty, t]} d\mathbf{x}$$

Remarks:

- The choice of risk-averseness level $\alpha \in [0, 1)$ controls the exploration vs exploitation trade-off.
- Upon convergence the predicted levels sets are guaranteed to be a subset of the true level sets.

Active learning of level sets

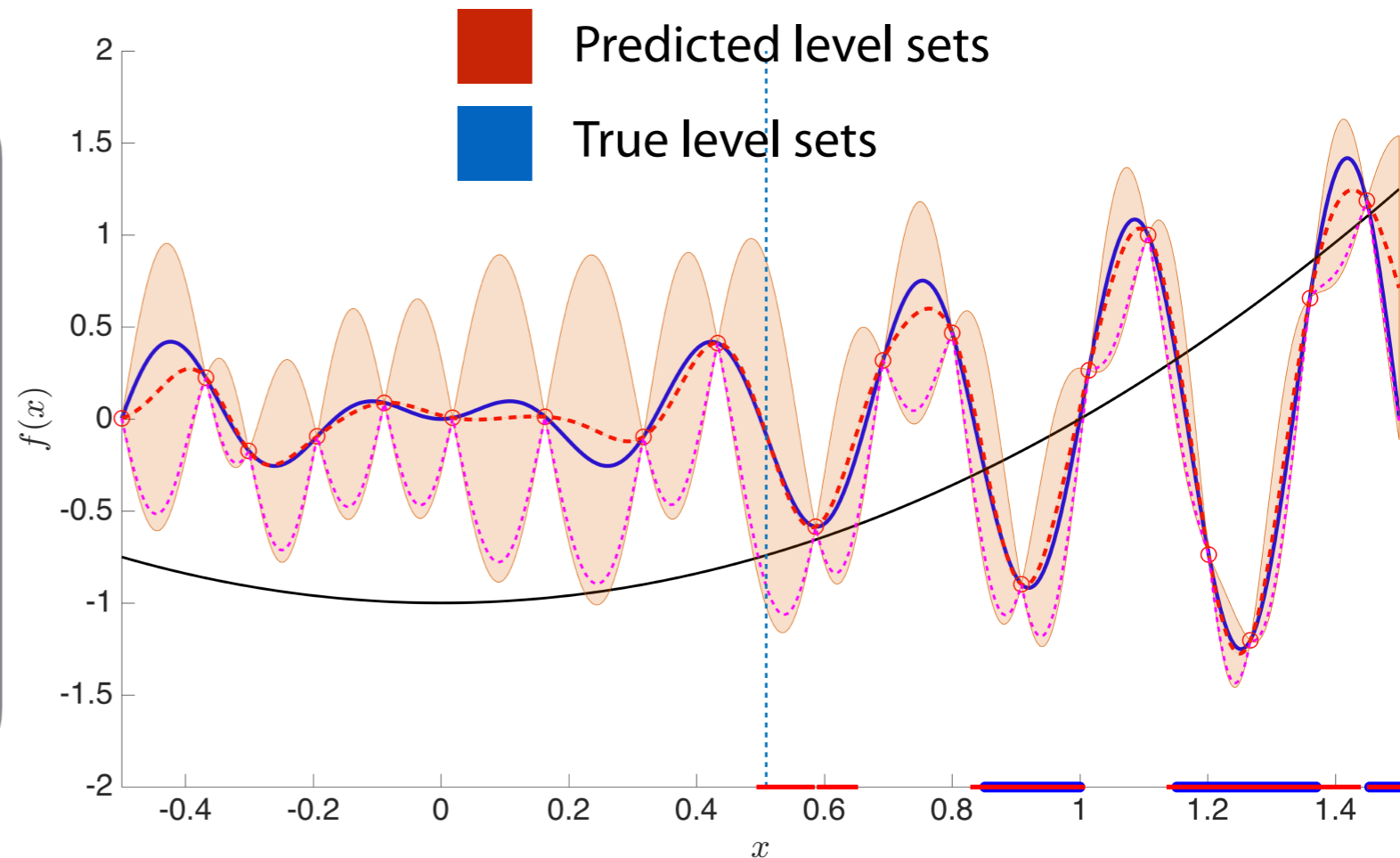
Goal: Identify the sets $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t\}$

Utilize the posterior to guide a sequential sampling policy by optimizing a chosen expected utility function

$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_\theta \mathbb{E}_v |_{\mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

e.g. sample at the locations that maximize the posterior variance in $L(t)$

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x} \in L_\alpha(t)} V(f(\mathbf{x}))$$



Terminate iteration when the “volume” of the predicted level sets is below a given threshold:

$$|V_{n+1}(t) - V_n(t)| < \epsilon, \quad V_n(t) = \int_{L_\alpha(t)} \mathbf{1}_{[-\infty, t]} d\mathbf{x}$$

Remarks:

- The choice of risk-averseness level $\alpha \in [0, 1)$ controls the exploration vs exploitation trade-off.
- Upon convergence the predicted levels sets are guaranteed to be a subset of the true level sets.

Active learning of level sets

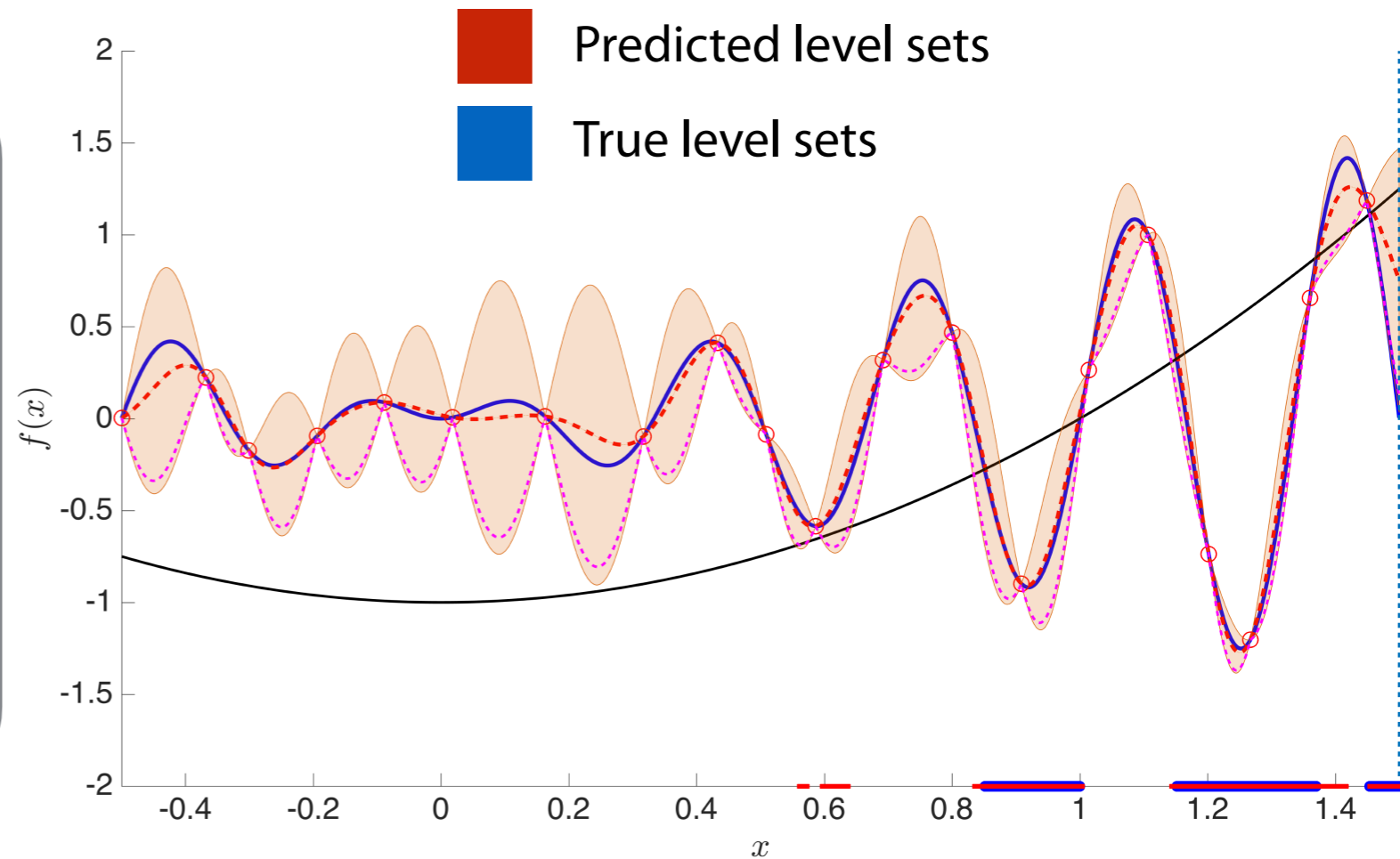
Goal: Identify the sets $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t\}$

Utilize the posterior to guide a sequential sampling policy by optimizing a chosen expected utility function

$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_\theta \mathbb{E}_v |_{\mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

e.g. sample at the locations that maximize the posterior variance in $L(t)$

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x} \in L_\alpha(t)} V(f(\mathbf{x}))$$



Terminate iteration when the “volume” of the predicted level sets is below a given threshold:

$$|V_{n+1}(t) - V_n(t)| < \epsilon, \quad V_n(t) = \int_{L_\alpha(t)} \mathbf{1}_{[-\infty, t]} d\mathbf{x}$$

Remarks:

- The choice of risk-averseness level $\alpha \in [0, 1)$ controls the exploration vs exploitation trade-off.
- Upon convergence the predicted levels sets are guaranteed to be a subset of the true level sets.

Active learning of level sets

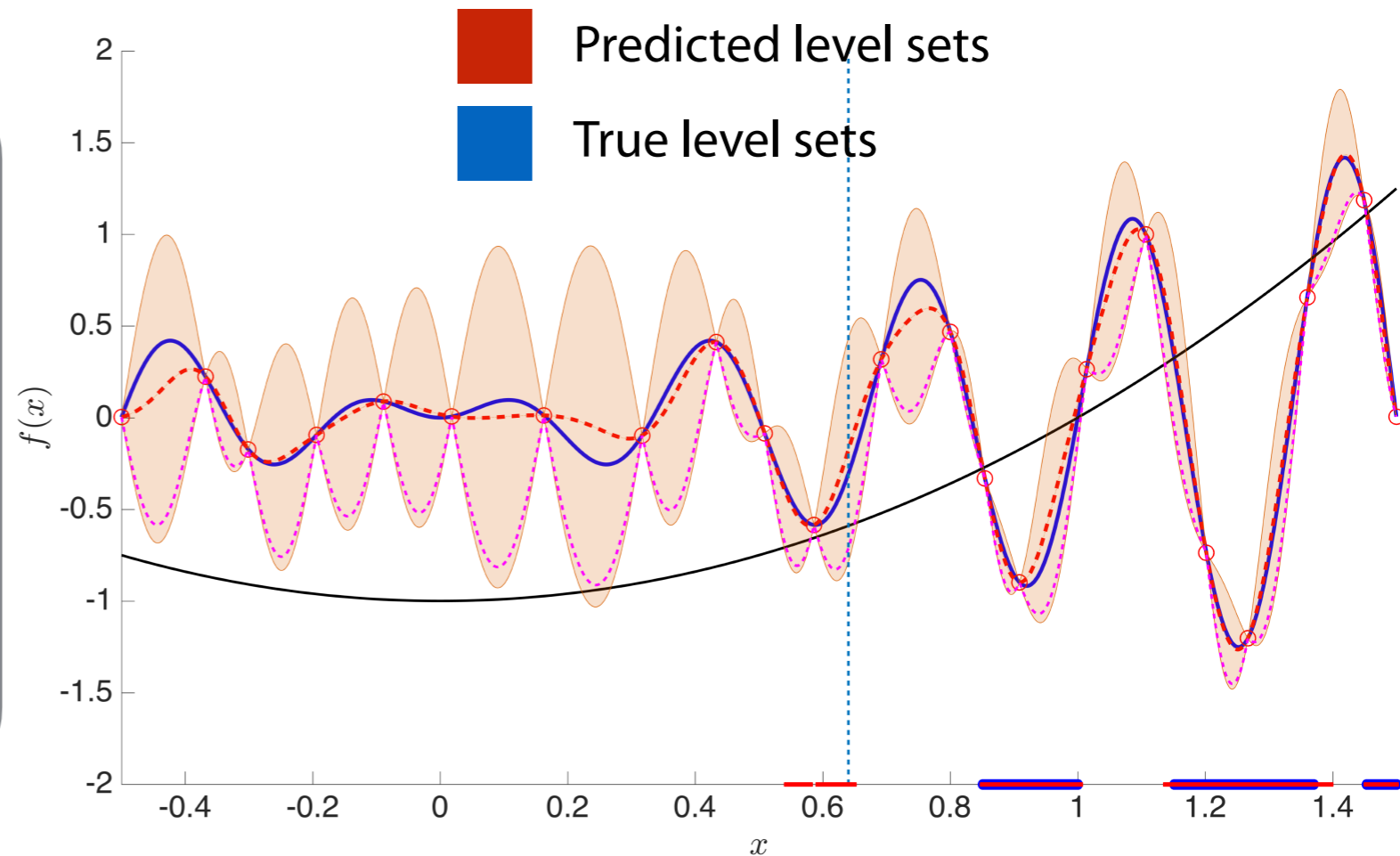
Goal: Identify the sets $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t\}$

Utilize the posterior to guide a sequential sampling policy by optimizing a chosen expected utility function

$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_\theta \mathbb{E}_v |_{\mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

e.g. sample at the locations that maximize the posterior variance in $L(t)$

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x} \in L_\alpha(t)} V(f(\mathbf{x}))$$



Terminate iteration when the “volume” of the predicted level sets is below a given threshold:

$$|V_{n+1}(t) - V_n(t)| < \epsilon, \quad V_n(t) = \int_{L_\alpha(t)} \mathbf{1}_{[-\infty, t]} d\mathbf{x}$$

Remarks:

- The choice of risk-averseness level $\alpha \in [0, 1)$ controls the exploration vs exploitation trade-off.
- Upon convergence the predicted levels sets are guaranteed to be a subset of the true level sets.

Active learning of level sets

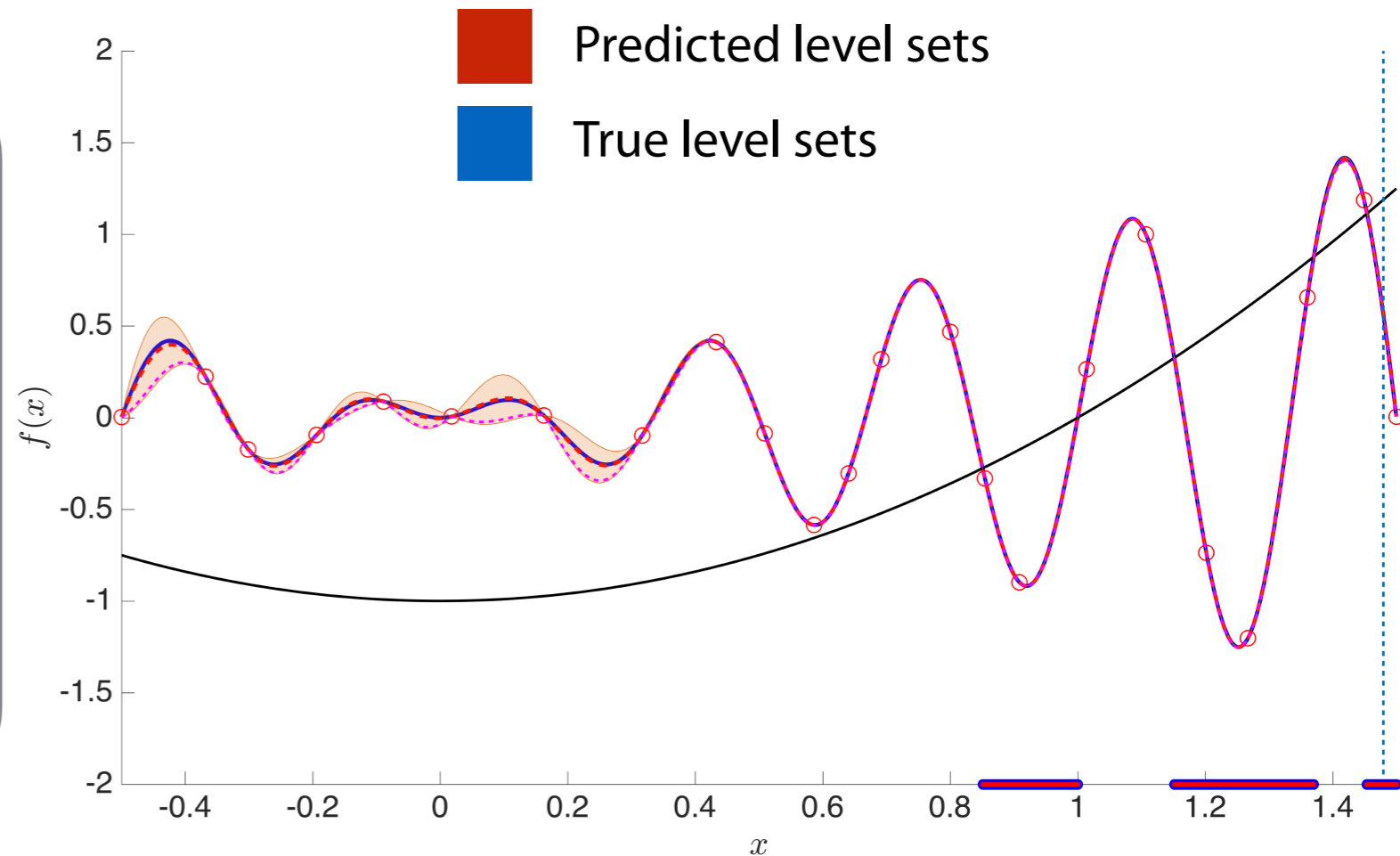
Goal: Identify the sets $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t\}$

Utilize the posterior to guide a sequential sampling policy by optimizing a chosen expected utility function

$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_\theta \mathbb{E}_v |_{\mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

e.g. sample at the locations that maximize the posterior variance in $L(t)$

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x} \in L_\alpha(t)} V(f(\mathbf{x}))$$



Terminate iteration when the “volume” of the predicted level sets is below a given threshold:

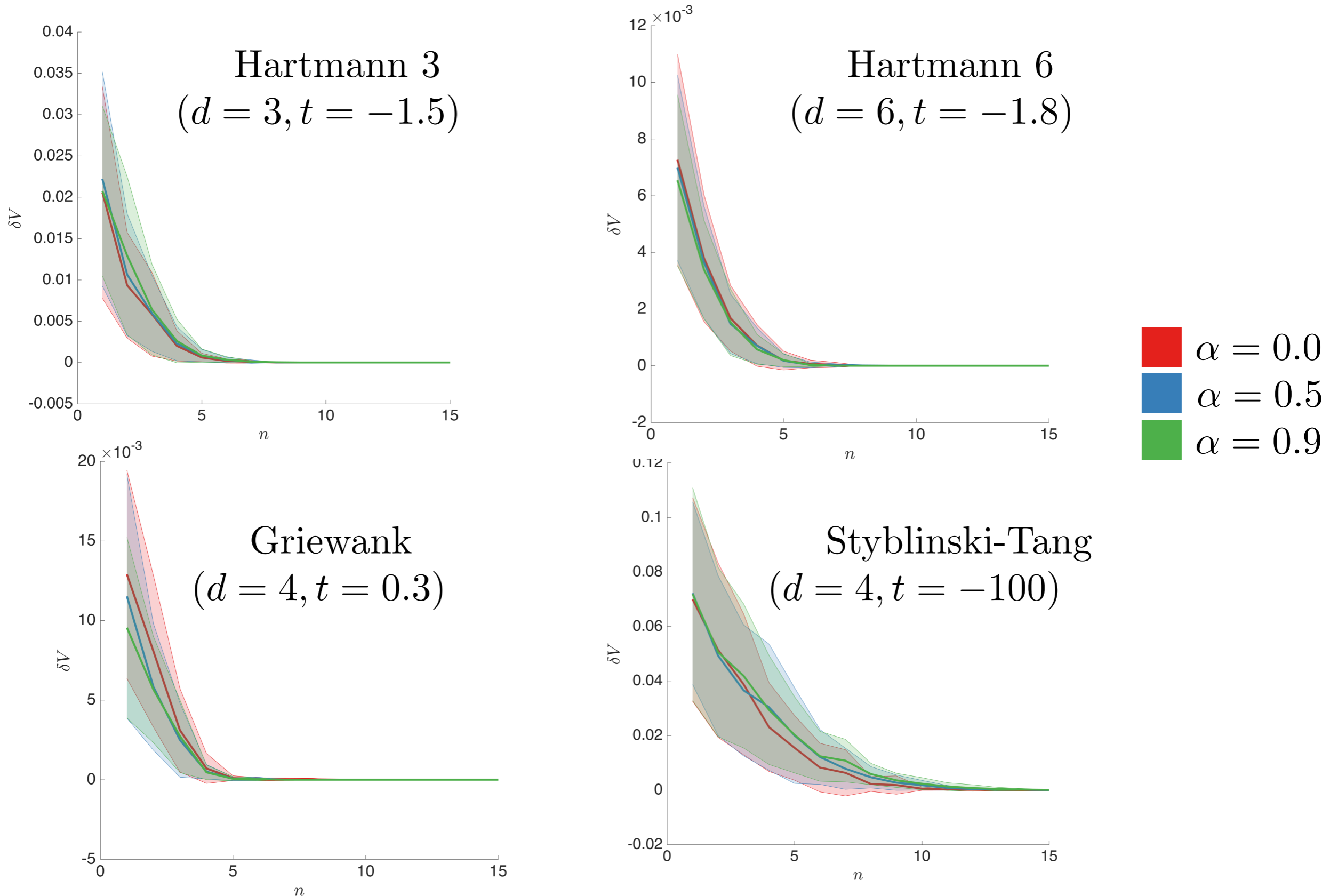
$$|V_{n+1}(t) - V_n(t)| < \epsilon, \quad V_n(t) = \int_{L_\alpha(t)} \mathbf{1}_{[-\infty, t]} d\mathbf{x}$$

Remarks:

- The choice of risk-averseness level $\alpha \in [0, 1)$ controls the exploration vs exploitation trade-off.
- Upon convergence the predicted levels sets are guaranteed to be a subset of the true level sets.

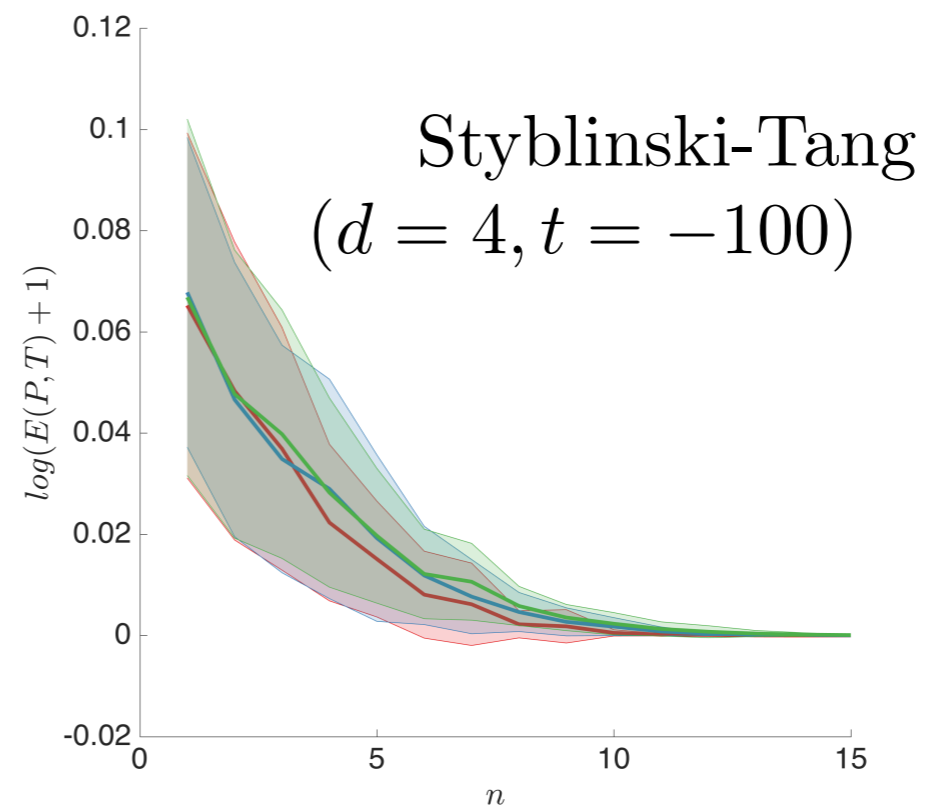
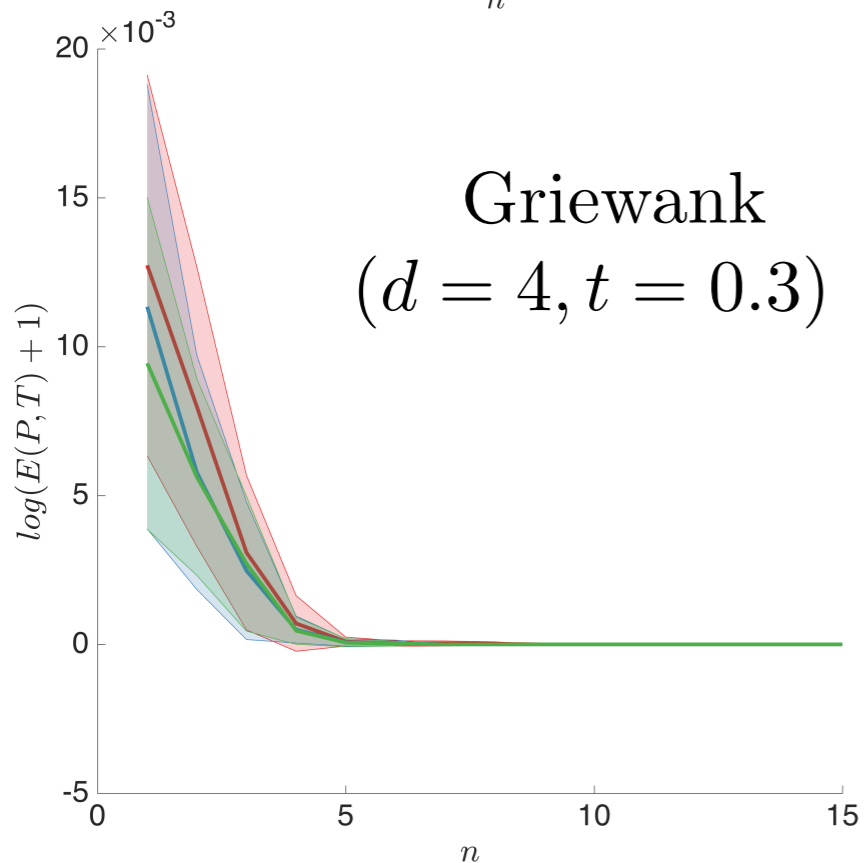
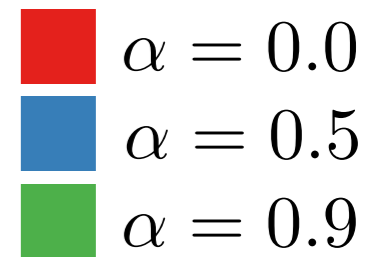
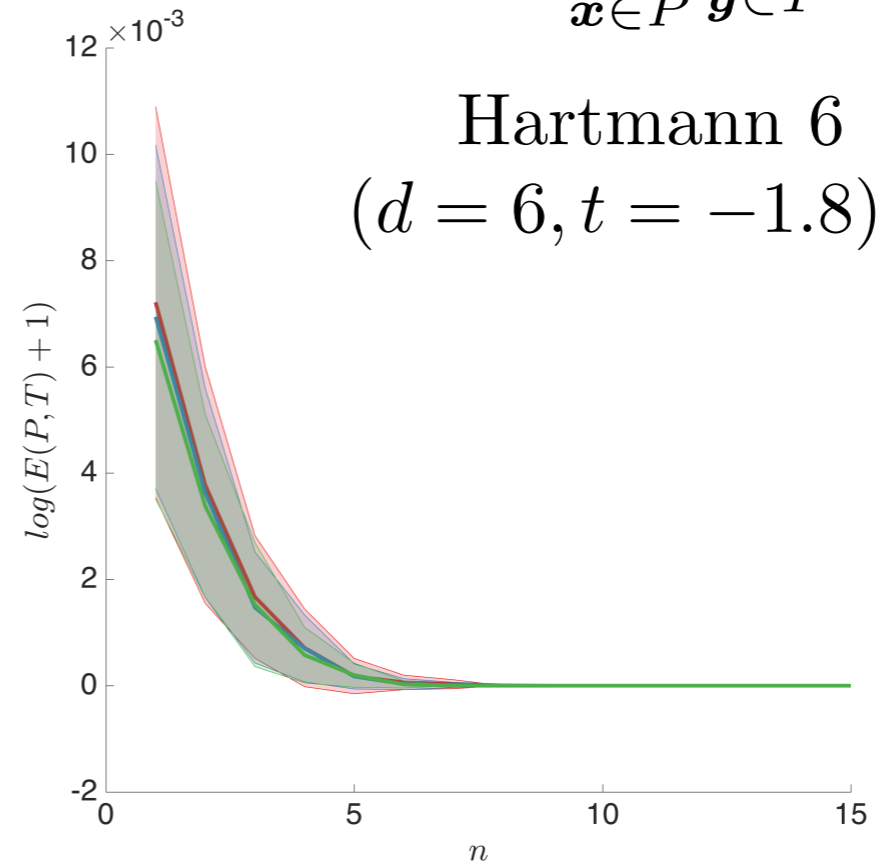
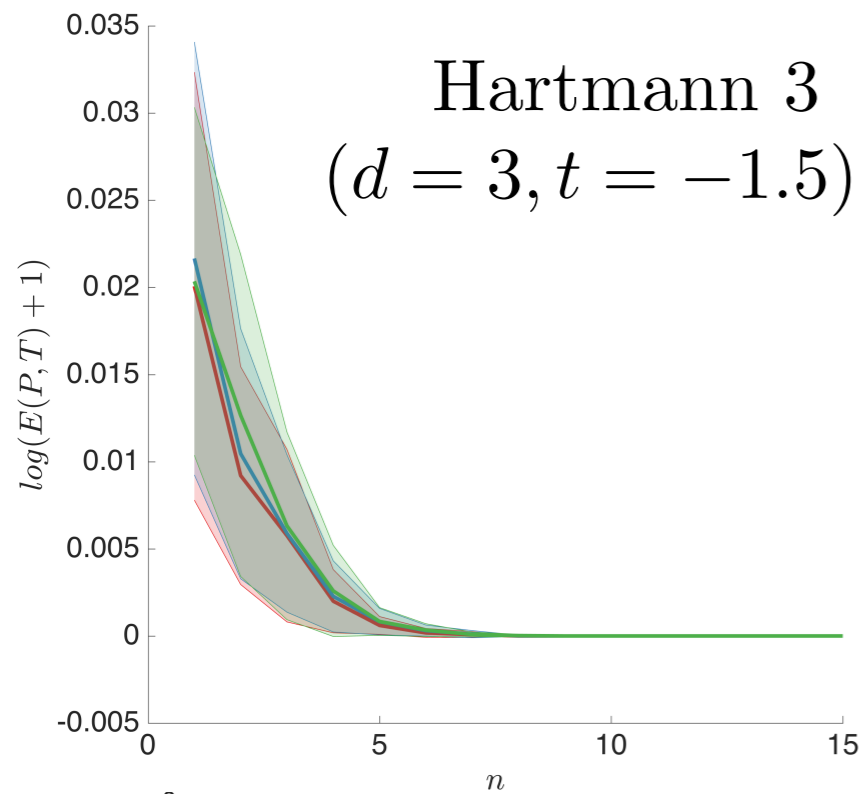
Active learning of level sets

Change in volume after between two consecutive iterations: $\delta V = |V_{n+1} - V_n|$



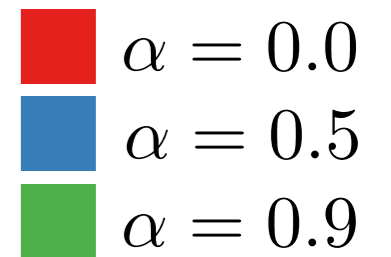
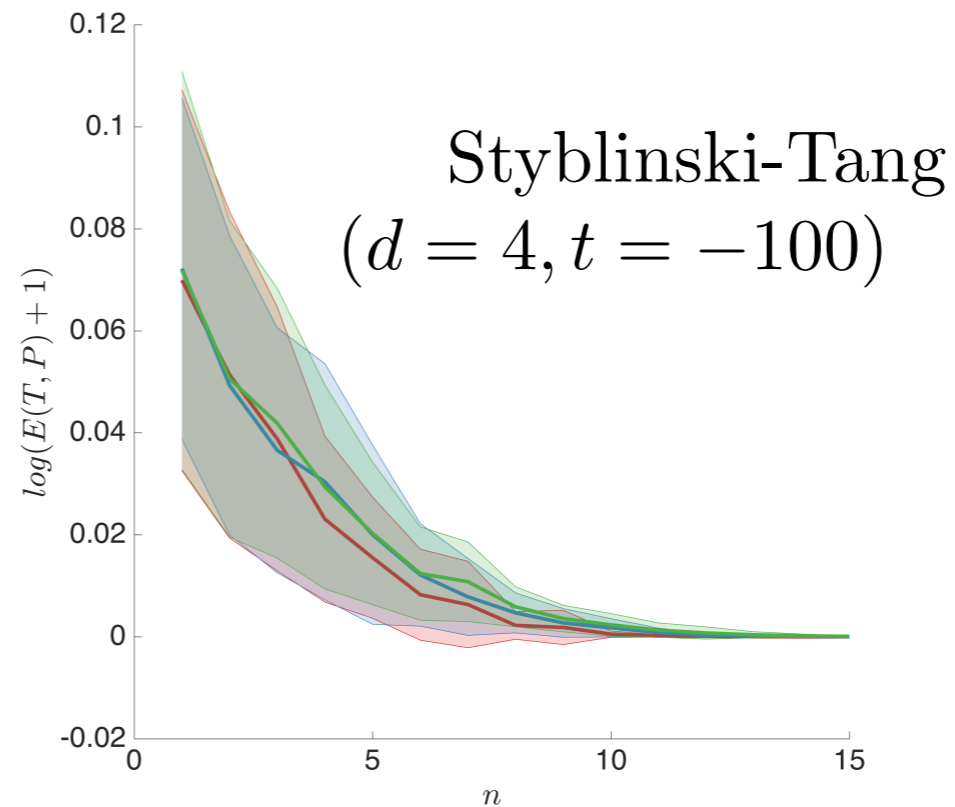
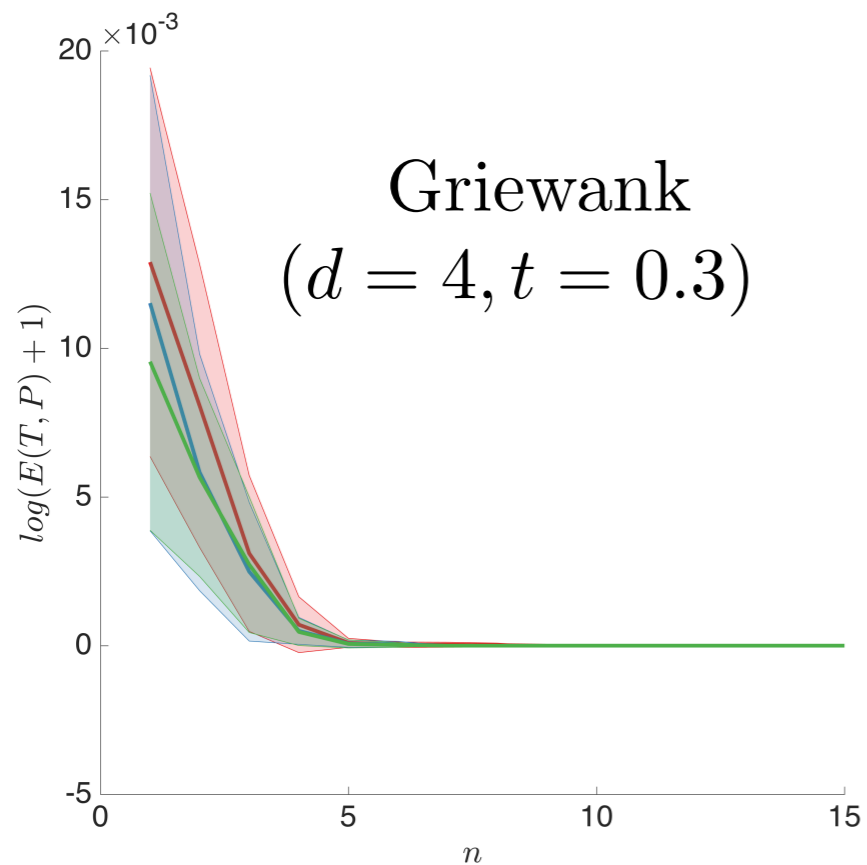
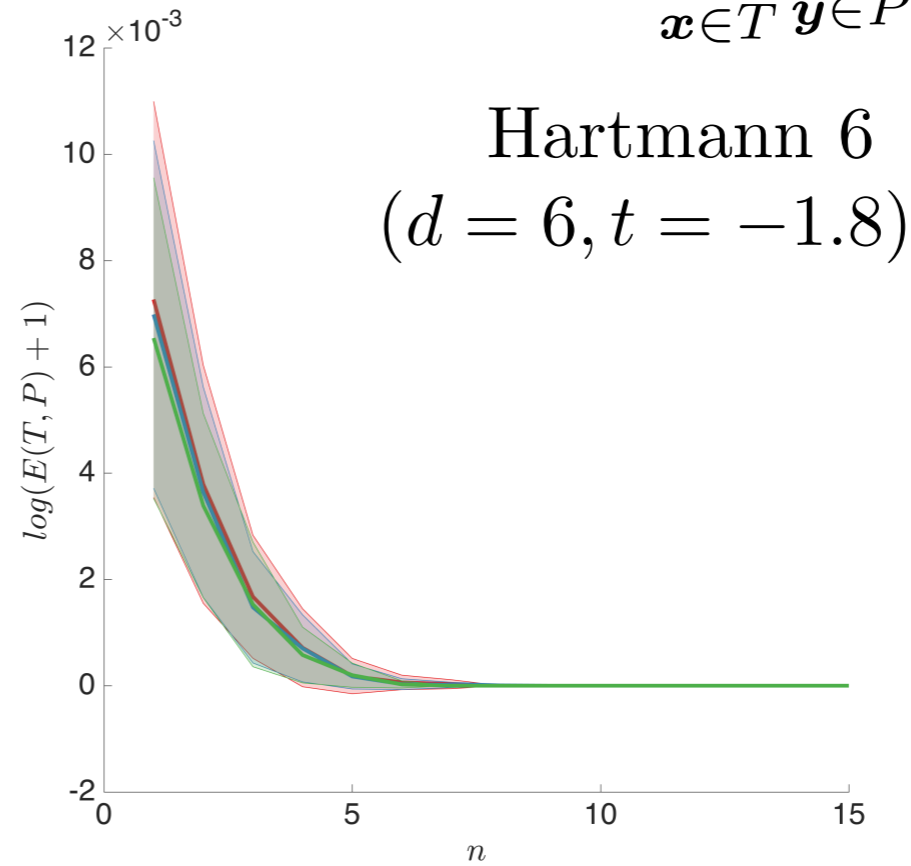
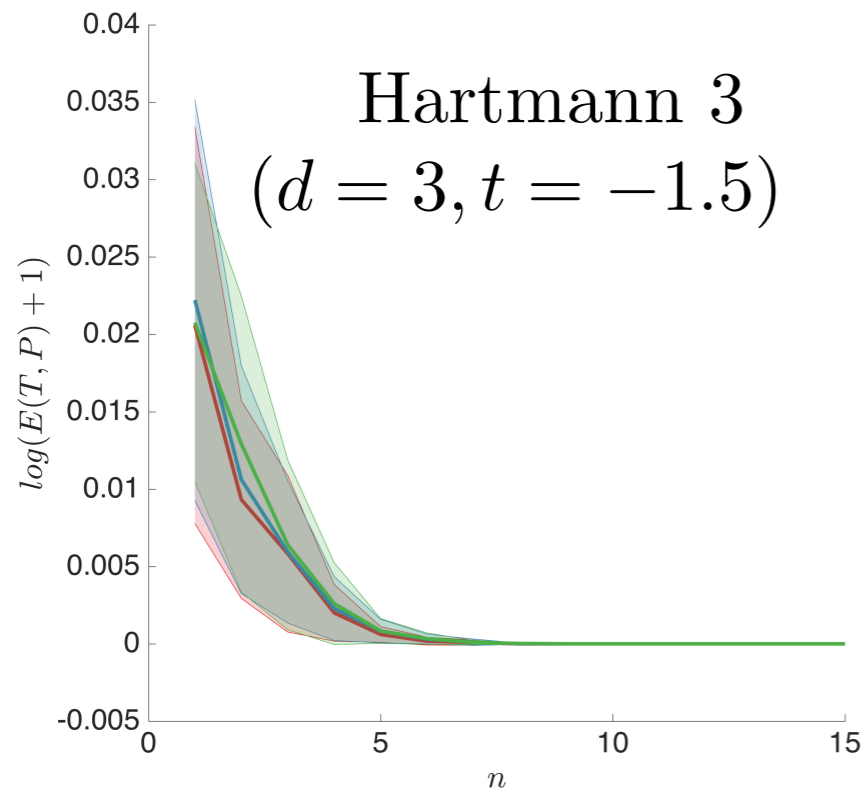
Active learning of level sets

Excess of predicted over the true level sets: $E(P, T) = \sup_{x \in P} \inf_{y \in T} d(x, y)$



Active learning of level sets

Excess of true over the predicted level sets: $E(T, P) = \sup_{x \in T} \inf_{y \in P} d(x, y)$



Active learning of level sets

Case	α	0.0	0.5	0.9
Hartmann ($d = 3$)		6.7	7.3	8.1
Hartmann ($d = 6$)		6.2	6.2	6.5
Griewank ($d = 4$)		6.5	6.5	6.6
Styblinski-Tang ($d = 4$)		9.3	10.6	12.3

Table 1: Average # of iterations to convergence for 100 independent trials.

Bayesian optimization

Goal: Estimate the global minimum of a function: $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} g(\mathbf{x})$ (potentially intractable)

Setup: $g(x)$ is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

Idea: Approximate $g(x)$ using a GP surrogate: $y = f(\mathbf{x}) + \epsilon$, $f \sim \mathcal{GP}(f|0, k(\mathbf{x}, \mathbf{x}'; \theta))$

Utilize the posterior to guide a sequential or parallel sampling policy by optimizing a chosen expected utility function

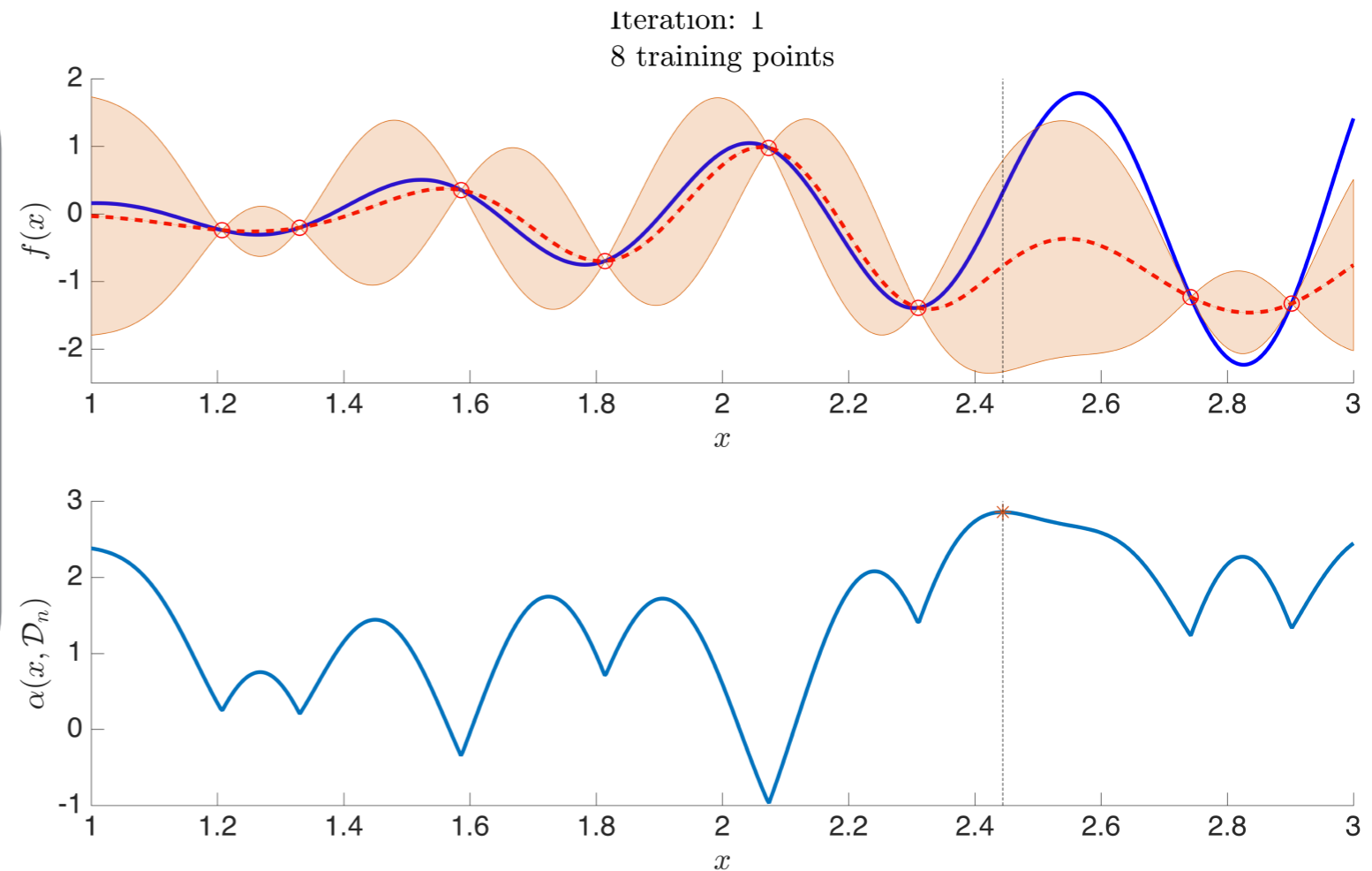
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v | \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

The optimization problem is transformed to:

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x}} \alpha(\mathbf{x}; \mathcal{D}_n)$$

Remark:

Acquisition functions aim to balance the trade-off between exploration and exploitation.



e.g. sample at the locations that minimize the lower superquintile risk confidence bound

$$\mathbf{x}_{n+1} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \mu(\mathbf{x}) - \frac{\phi(\Phi^{-1}(\alpha))}{1 - \alpha} \sigma(\mathbf{x})$$

Bayesian optimization

Goal: Estimate the global minimum of a function: $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} g(\mathbf{x})$ (potentially intractable)

Setup: $g(x)$ is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

Idea: Approximate $g(x)$ using a GP surrogate: $y = f(\mathbf{x}) + \epsilon$, $f \sim \mathcal{GP}(f|0, k(\mathbf{x}, \mathbf{x}'; \theta))$

Utilize the posterior to guide a sequential or parallel sampling policy by optimizing a chosen expected utility function

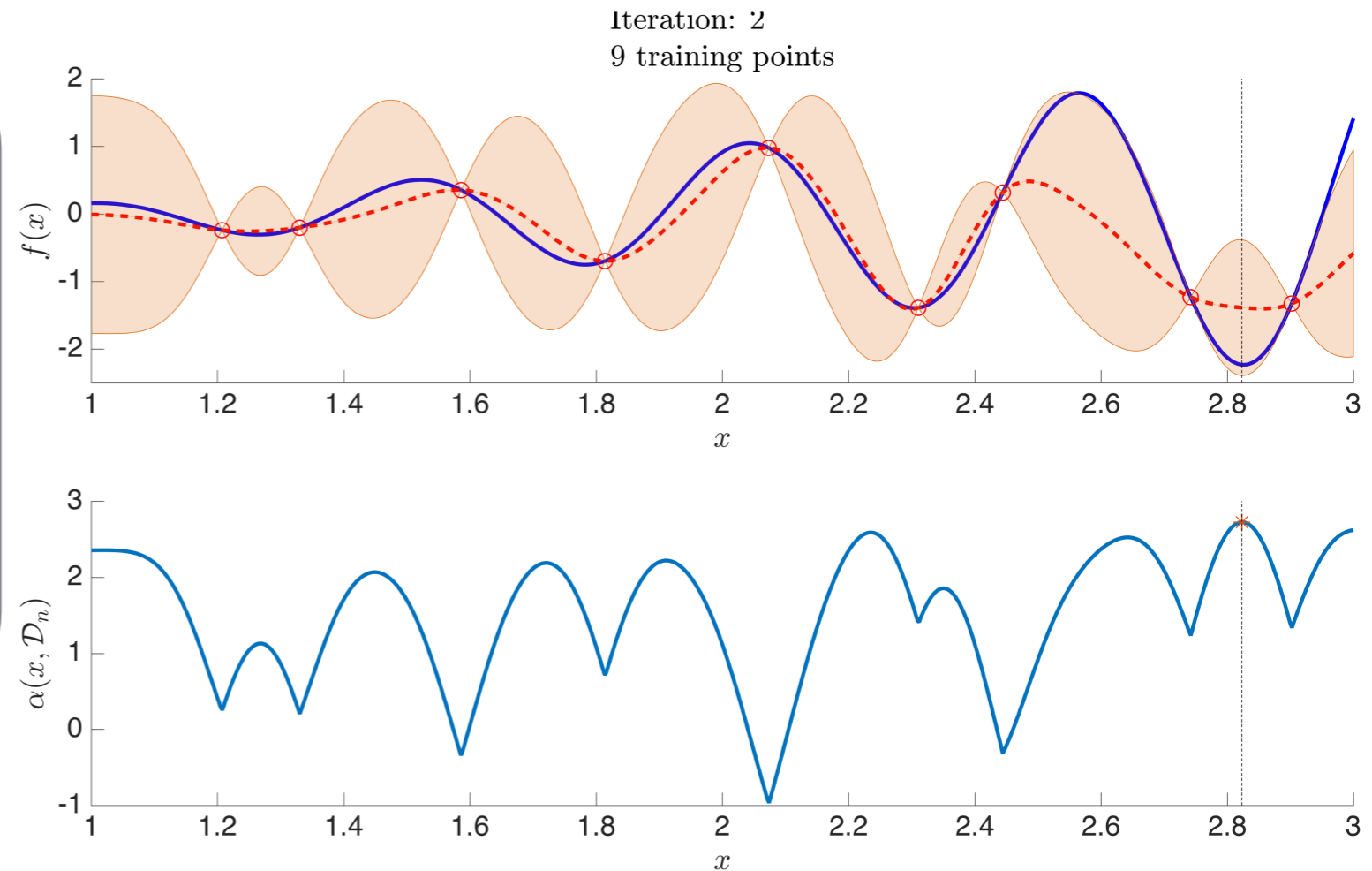
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_\theta \mathbb{E}_v |_{\mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

The optimization problem is transformed to:

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x}} \alpha(\mathbf{x}; \mathcal{D}_n)$$

Remark:

Acquisition functions aim to balance the trade-off between exploration and exploitation.



e.g. sample at the locations that minimize the lower superquintile risk confidence bound

$$\mathbf{x}_{n+1} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \mu(\mathbf{x}) - \frac{\phi(\Phi^{-1}(\alpha))}{1 - \alpha} \sigma(\mathbf{x})$$

Bayesian optimization

Goal: Estimate the global minimum of a function: $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} g(\mathbf{x})$ (potentially intractable)

Setup: $g(x)$ is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

Idea: Approximate $g(x)$ using a GP surrogate: $y = f(\mathbf{x}) + \epsilon$, $f \sim \mathcal{GP}(f|0, k(\mathbf{x}, \mathbf{x}'; \theta))$

Utilize the posterior to guide a sequential or parallel sampling policy by optimizing a chosen expected utility function

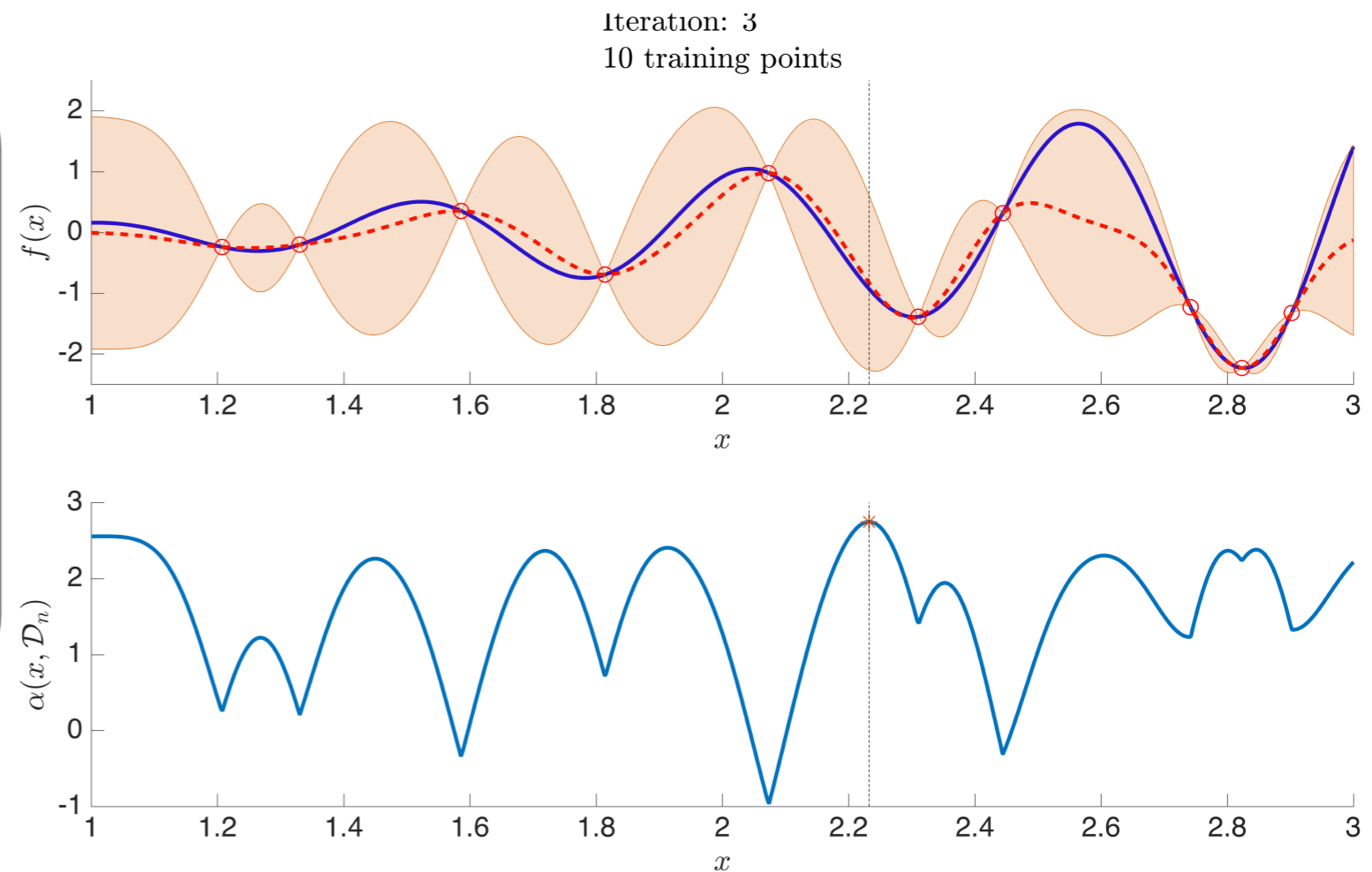
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v | \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

The optimization problem is transformed to:

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x}} \alpha(\mathbf{x}; \mathcal{D}_n)$$

Remark:

Acquisition functions aim to balance the trade-off between exploration and exploitation.



e.g. sample at the locations that minimize the lower superquintile risk confidence bound

$$\mathbf{x}_{n+1} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \mu(\mathbf{x}) - \frac{\phi(\Phi^{-1}(\alpha))}{1 - \alpha} \sigma(\mathbf{x})$$

Bayesian optimization

Goal: Estimate the global minimum of a function: $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} g(\mathbf{x})$ (potentially intractable)

Setup: $g(x)$ is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

Idea: Approximate $g(x)$ using a GP surrogate: $y = f(\mathbf{x}) + \epsilon$, $f \sim \mathcal{GP}(f|0, k(\mathbf{x}, \mathbf{x}'; \theta))$

Utilize the posterior to guide a sequential or parallel sampling policy by optimizing a chosen expected utility function

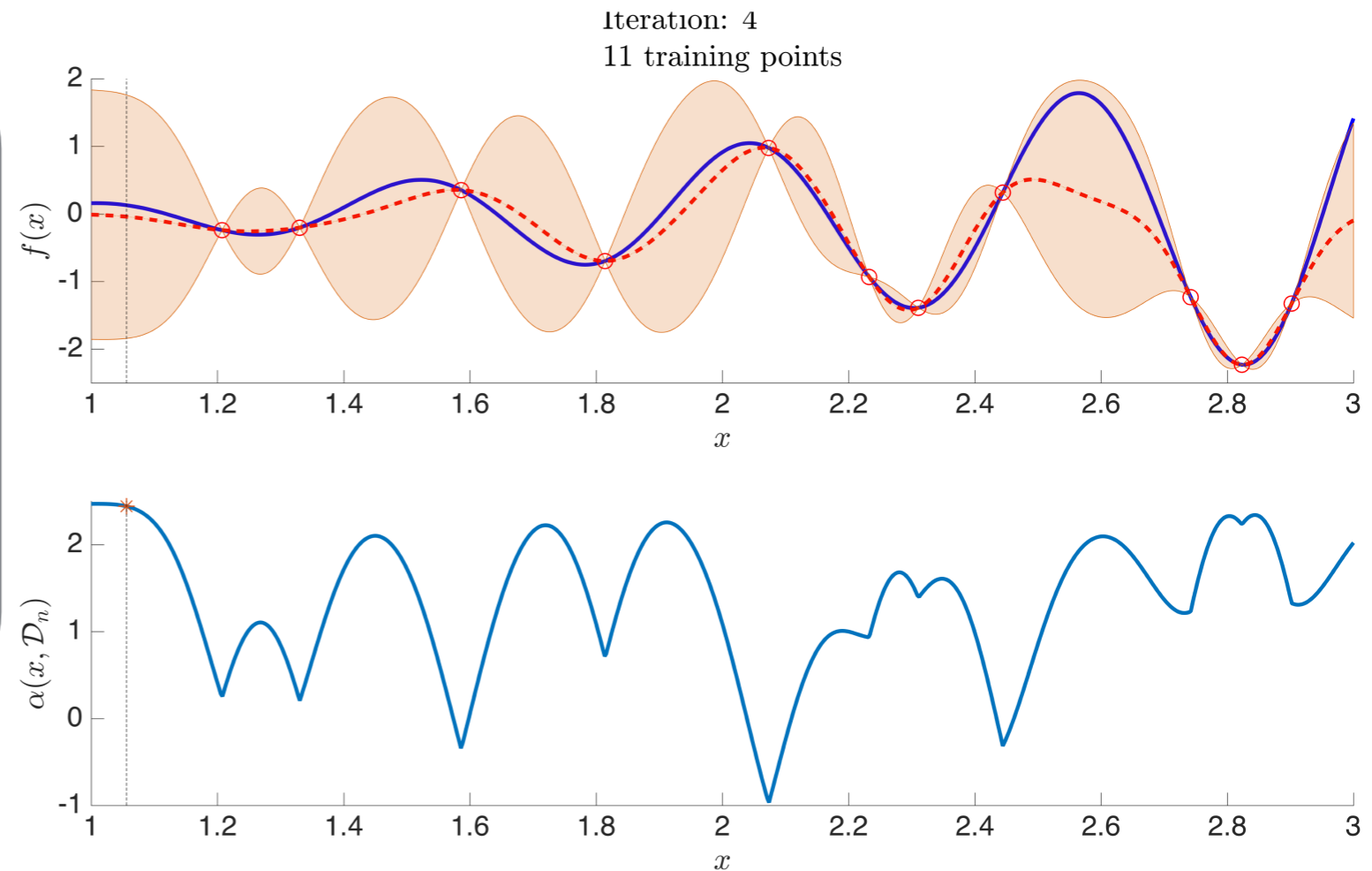
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v | \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

The optimization problem is transformed to:

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x}} \alpha(\mathbf{x}; \mathcal{D}_n)$$

Remark:

Acquisition functions aim to balance the trade-off between exploration and exploitation.



e.g. sample at the locations that minimize the lower superquintile risk confidence bound

$$\mathbf{x}_{n+1} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \mu(\mathbf{x}) - \frac{\phi(\Phi^{-1}(\alpha))}{1 - \alpha} \sigma(\mathbf{x})$$

Bayesian optimization

Goal: Estimate the global minimum of a function: $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} g(\mathbf{x})$ (potentially intractable)

Setup: $g(x)$ is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

Idea: Approximate $g(x)$ using a GP surrogate: $y = f(\mathbf{x}) + \epsilon$, $f \sim \mathcal{GP}(f|0, k(\mathbf{x}, \mathbf{x}'; \theta))$

Utilize the posterior to guide a sequential or parallel sampling policy by optimizing a chosen expected utility function

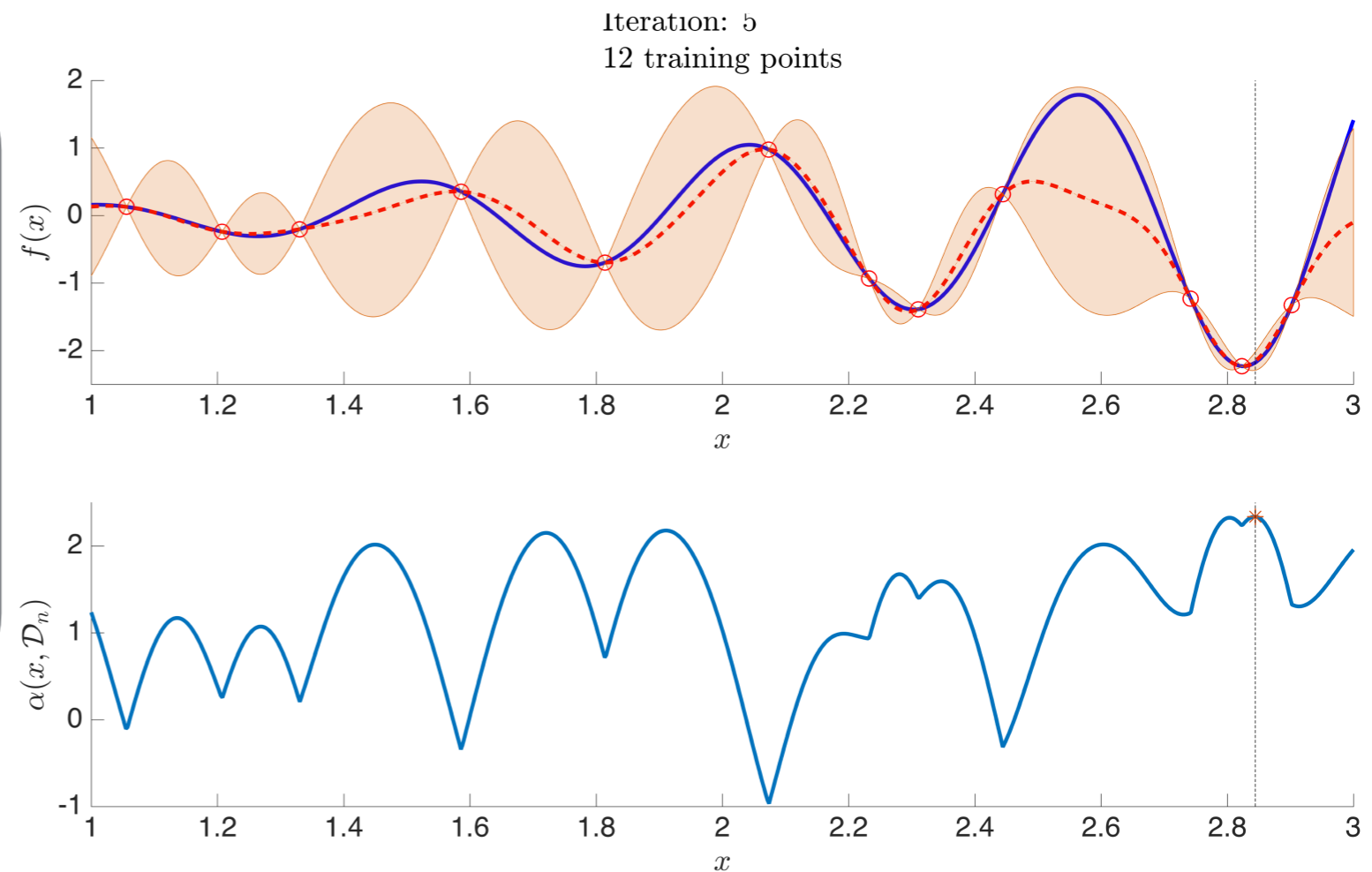
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v | \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

The optimization problem is transformed to:

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x}} \alpha(\mathbf{x}; \mathcal{D}_n)$$

Remark:

Acquisition functions aim to balance the trade-off between exploration and exploitation.



e.g. sample at the locations that minimize the lower superquintile risk confidence bound

$$\mathbf{x}_{n+1} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \mu(\mathbf{x}) - \frac{\phi(\Phi^{-1}(\alpha))}{1 - \alpha} \sigma(\mathbf{x})$$

Bayesian optimization

Goal: Estimate the global minimum of a function: $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} g(\mathbf{x})$ (potentially intractable)

Setup: $g(x)$ is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

Idea: Approximate $g(x)$ using a GP surrogate: $y = f(\mathbf{x}) + \epsilon$, $f \sim \mathcal{GP}(f|0, k(\mathbf{x}, \mathbf{x}'; \theta))$

Utilize the posterior to guide a sequential or parallel sampling policy by optimizing a chosen expected utility function

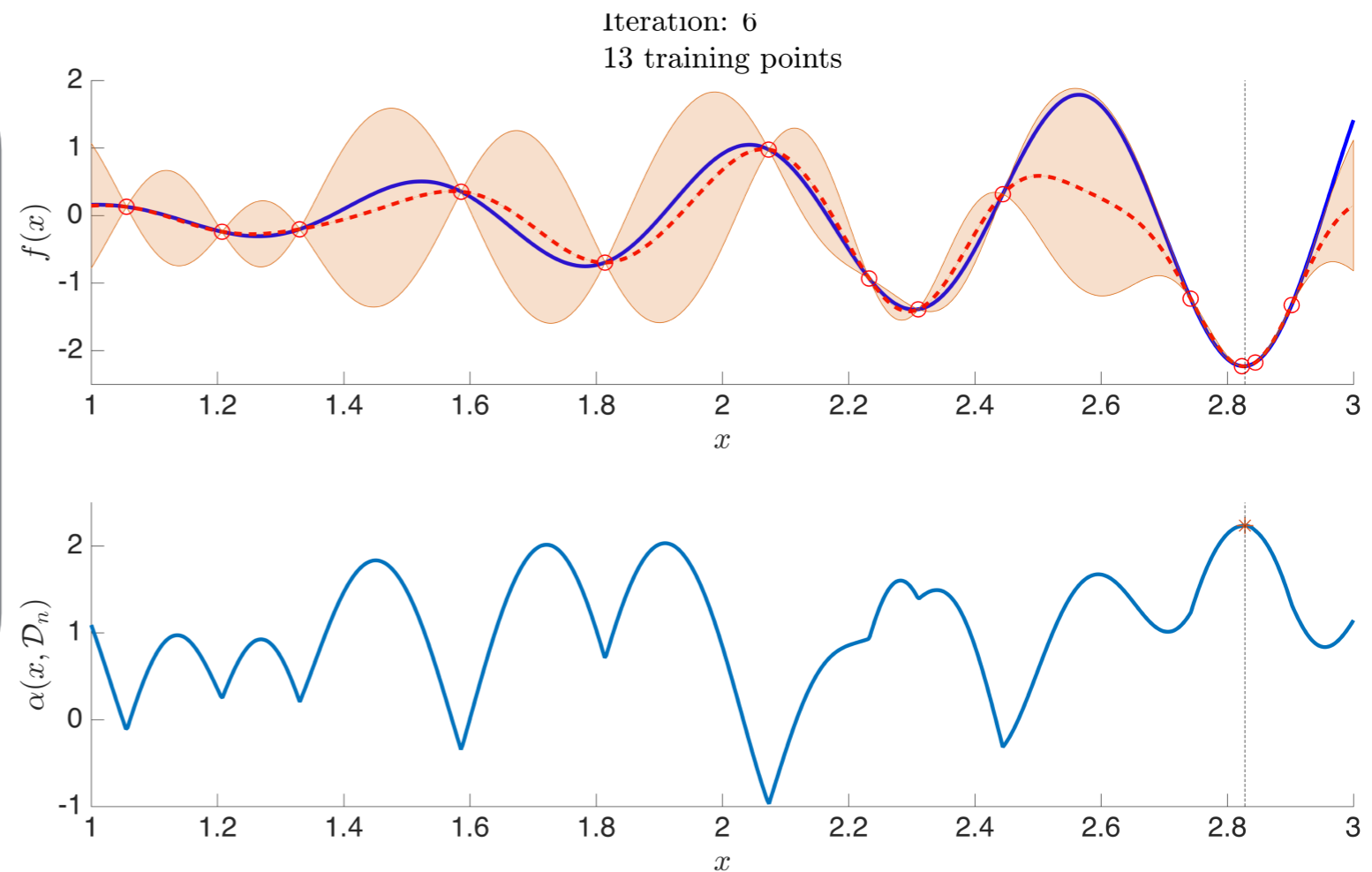
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v | \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

The optimization problem is transformed to:

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x}} \alpha(\mathbf{x}; \mathcal{D}_n)$$

Remark:

Acquisition functions aim to balance the trade-off between exploration and exploitation.



e.g. sample at the locations that minimize the lower superquintile risk confidence bound

$$\mathbf{x}_{n+1} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \mu(\mathbf{x}) - \frac{\phi(\Phi^{-1}(\alpha))}{1 - \alpha} \sigma(\mathbf{x})$$

Bayesian optimization

Goal: Estimate the global minimum of a function: $\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^d} g(\mathbf{x})$ (potentially intractable)

Setup: $g(x)$ is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

Idea: Approximate $g(x)$ using a GP surrogate: $y = f(\mathbf{x}) + \epsilon$, $f \sim \mathcal{GP}(f|0, k(\mathbf{x}, \mathbf{x}'; \theta))$

Utilize the posterior to guide a sequential or parallel sampling policy by optimizing a chosen expected utility function

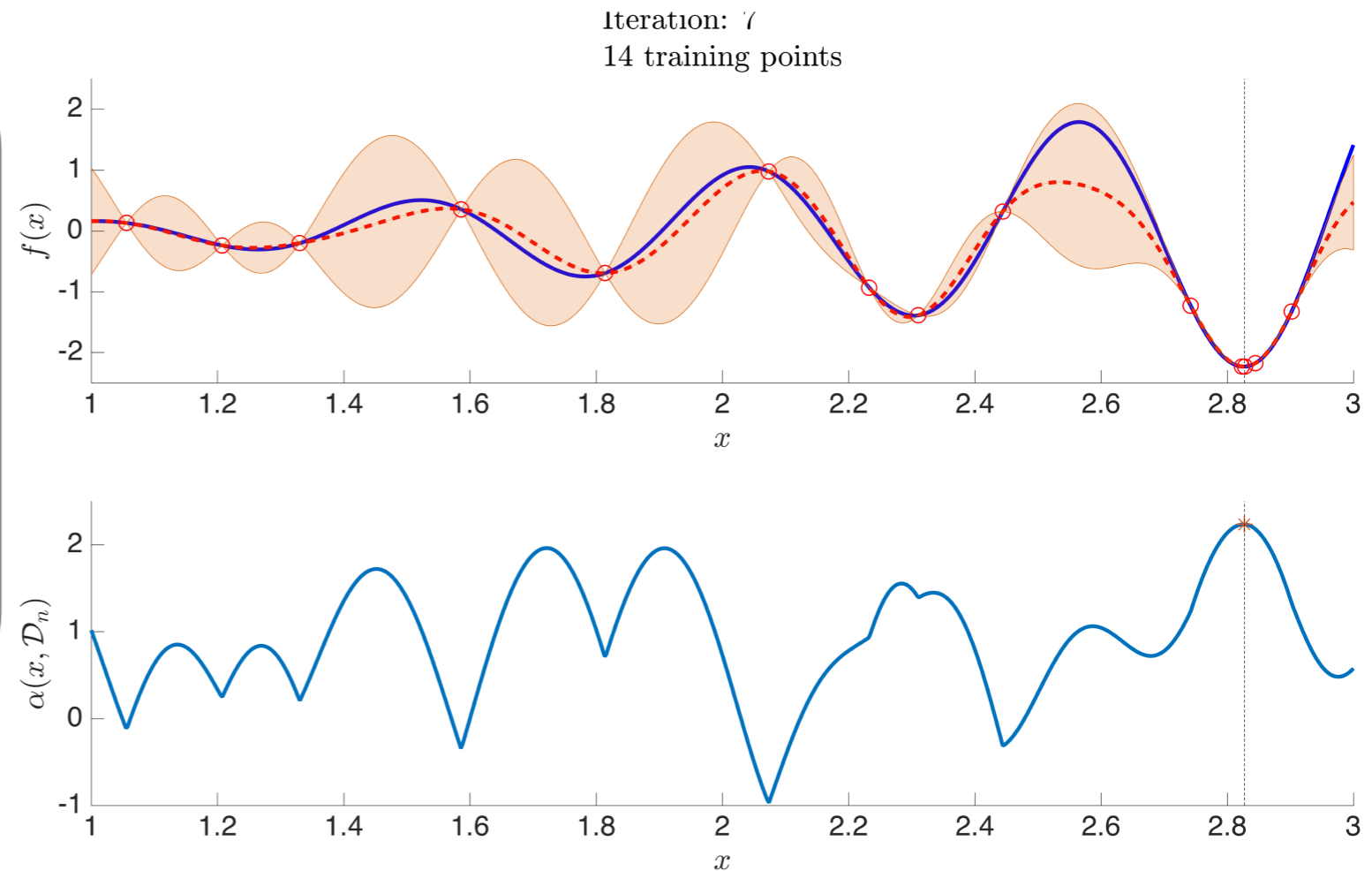
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v | \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

The optimization problem is transformed to:

$$\mathbf{x}_{n+1} = \arg \max_{\mathbf{x}} \alpha(\mathbf{x}; \mathcal{D}_n)$$

Remark:

Acquisition functions aim to balance the trade-off between exploration and exploitation.



e.g. sample at the locations that minimize the lower superquintile risk confidence bound

$$\mathbf{x}_{n+1} = \arg \min_{\mathbf{x} \in \mathbb{R}^d} \mu(\mathbf{x}) - \frac{\phi(\Phi^{-1}(\alpha))}{1 - \alpha} \sigma(\mathbf{x})$$

Summary

Web: <http://web.mit.edu/parisp/www/>

Email: parisp@mit.edu

Questions?

Funding:



Collaborators:

Maziar Raissi (Brown)

Johannes Royset (NPS)

This work received support from DARPA grant HR0011-14-1-0060

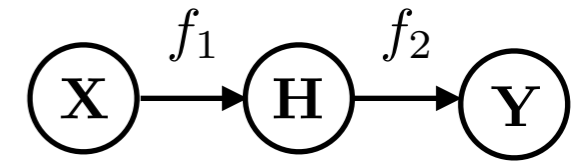
Web: <http://web.mit.edu/parisp/www/>

Email: parisp@mit.edu

Challenges and limitations

Discontinuities and non-stationarity: GPs struggle with discontinuous data

Use warping functions to transform into a jointly stationary input space



- Log, sigmoid, betaCDF —> "Warped GPs" *Snelson, E., C.E. Rasmussen, and Z.Ghahramani. "Warped gaussian processes."*
- Neural networks —> "Manifold GPs" *Calandra, R., et al. "Manifold Gaussian processes for regression."*
- Gaussian processes —> "Deep GPs" *Damianou, A. C., and N.D. Lawrence. "Deep Gaussian processes."*

Theoretical guarantees: Accuracy, convergence rates, posterior consistency, contraction rates, etc.

Approximation theory in Reproducing Kernel Hilbert Spaces

Stuart, A.M., and A.L. Teckentrup. "Posterior consistency for Gaussian process approximations of Bayesian posterior distributions." arXiv preprint, 2016

Scalability: GPs suffer from a cubic scaling with the data

Low-rank approximations to the covariance

Snelson, E., and Z. Ghahramani. "Sparse Gaussian processes using pseudo-inputs."

Frequency-domain learning algorithms

Perdikaris P., D. Venturi, G.E. Karniadakis "Multi-fidelity information fusion algorithms for high dimensional systems and massive data-sets", SIAM J. Sci. Comput., 2016

Stochastic variational inference

Hensman, J., N. Fusi, and N.D. Lawrence. "Gaussian processes for big data."

High-dimensions: Tensor product kernels suffer from the curse of dimensionality, i.e. they require an exponentially increasing amount of training data

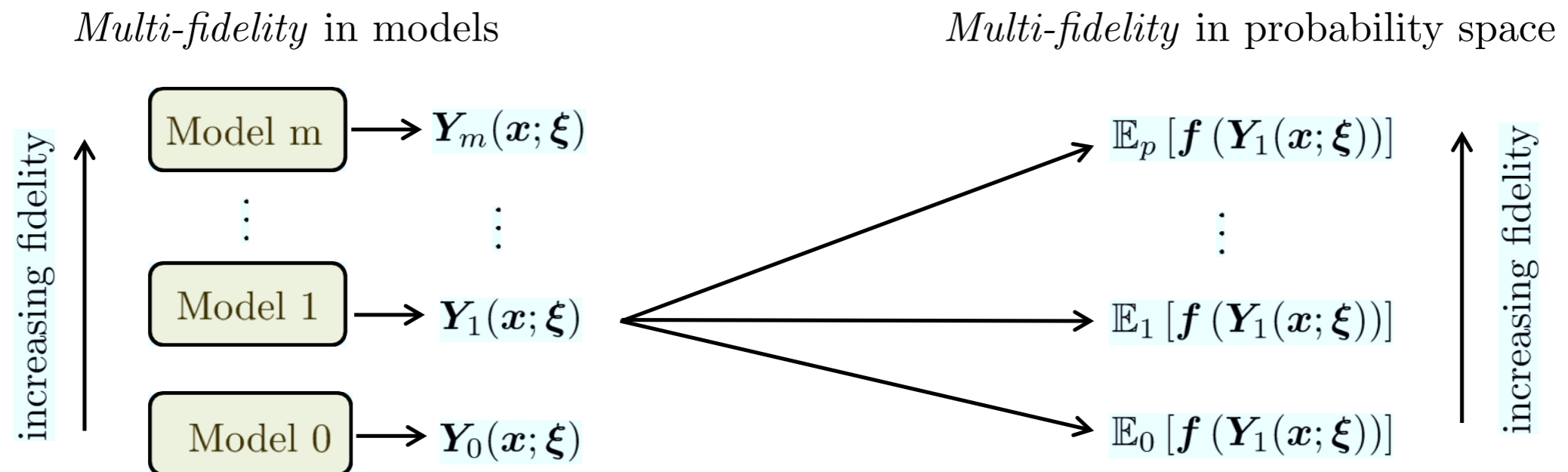
Data-driven additive kernels

Perdikaris P., D. Venturi, G.E. Karniadakis "Multi-fidelity information fusion algorithms for high dimensional systems and massive data-sets", SIAM J. Sci. Comput., 2016

Unsupervised dimensionality-reduction (GPLVM, variational auto-encoders)

Lawrence, N.D. "Gaussian process latent variable models for visualisation of high dimensional data."

Multi-fidelity in physical models and in probability space



PROCEEDINGS

OF THE ROYAL SOCIETY A

rspa.royalsocietypublishing.org

Research



Article submitted to journal

Multi-fidelity modeling via recursive co-kriging and Gaussian Markov random fields

P. Perdikaris¹, D. Venturi¹, J.O. Royset², and G.E. Karniadakis¹

¹Division of Applied Mathematics, Brown University, Providence, RI 02912, USA

²Operations Research Department, Naval Postgraduate School, Monterey, CA 93943, USA

$$\mathbb{E}_{k+1}[f(\mathbf{Y}_l(\mathbf{x}; \xi))] = \rho_{k+1} \mathbb{E}_k[f(\mathbf{Y}_l(\mathbf{x}; \xi))] + \delta_{k+1}(\mathbf{x}), \quad k \leq p, \quad l \leq m$$

$$\begin{pmatrix} \mathbb{E}_1 [f(\mathbf{Y}_1)] & \mathbb{E}_1 [f(\mathbf{Y}_2)] & \cdots & \mathbb{E}_1 [f(\mathbf{Y}_m)] \\ \mathbb{E}_2 [f(\mathbf{Y}_1)] & \mathbb{E}_2 [f(\mathbf{Y}_2)] & \cdots & \mathbb{E}_2 [f(\mathbf{Y}_m)] \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{E}_p [f(\mathbf{Y}_1)] & \mathbb{E}_p [f(\mathbf{Y}_2)] & \cdots & \mathbb{E}_p [f(\mathbf{Y}_m)] \end{pmatrix}$$

Fidelity in probability space

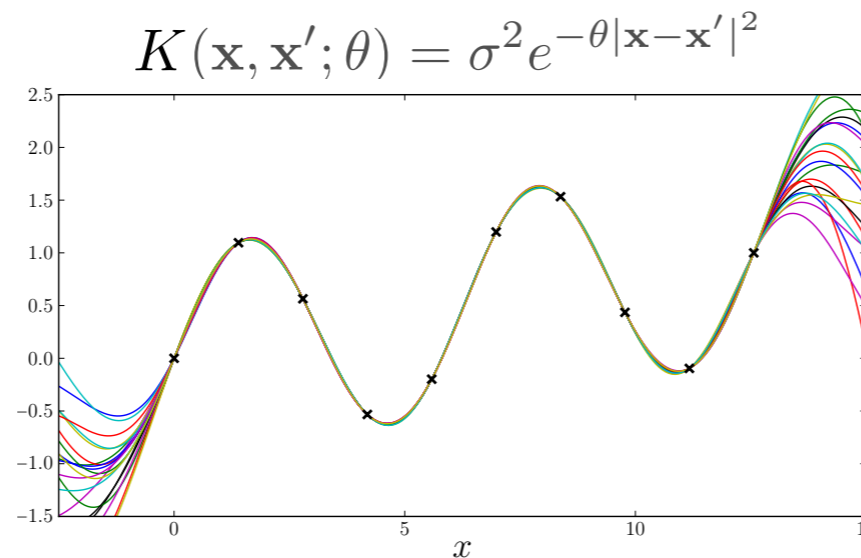
Fidelity in physical models

Importance of the prior

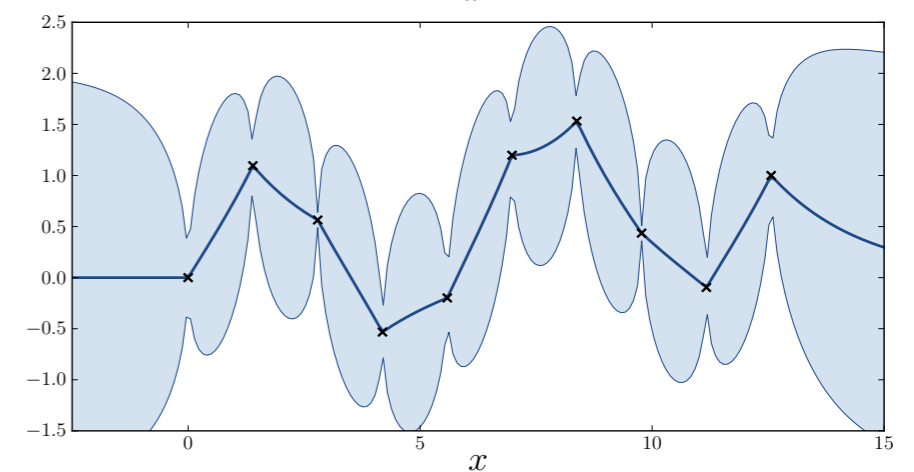
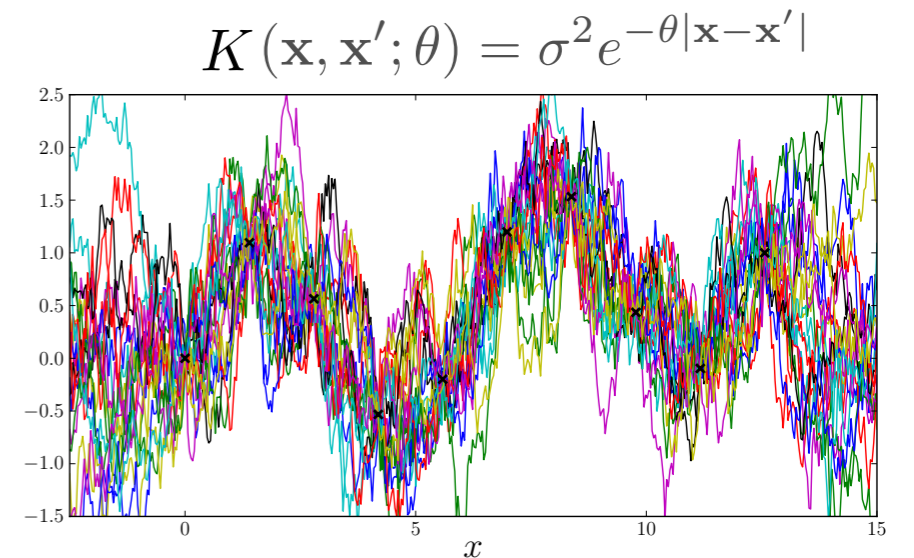
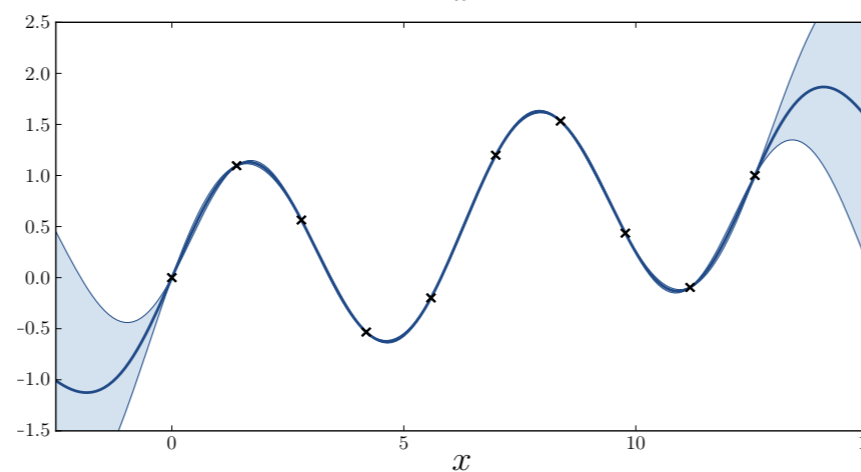
The choice of the covariance kernel has a big impact on the model as it is tightly related to:

- The smoothness of the sample paths, hence the regularity of the predictor.
- The accuracy and uncertainty of the predictor.
- The conditioning of the correlation matrix, hence the efficiency of the learning algorithms.

Samples from a GP posterior



Posterior mean and variance



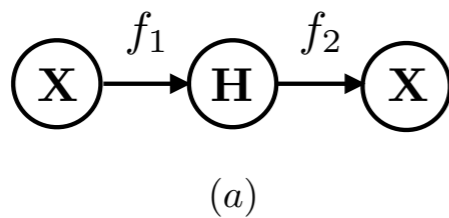
<http://www.automaticstatistician.com/index/>

Model inversion in high-dimensions

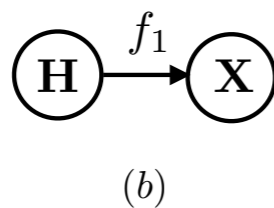
Goal: Developed scalable algorithms for solving high-dimensional inverse problems

Optimization in physical space $\min_{\mathbf{x} \in \mathbb{R}^d} \|g(\mathbf{x}) - y^*\|$ $\xrightarrow{\text{non-linear dim reduction } q \ll d}$ $\min_{\mathbf{h} \in \mathbf{H}^q} \|g(\mathbf{h}) - y^*\|$ Optimization in latent space

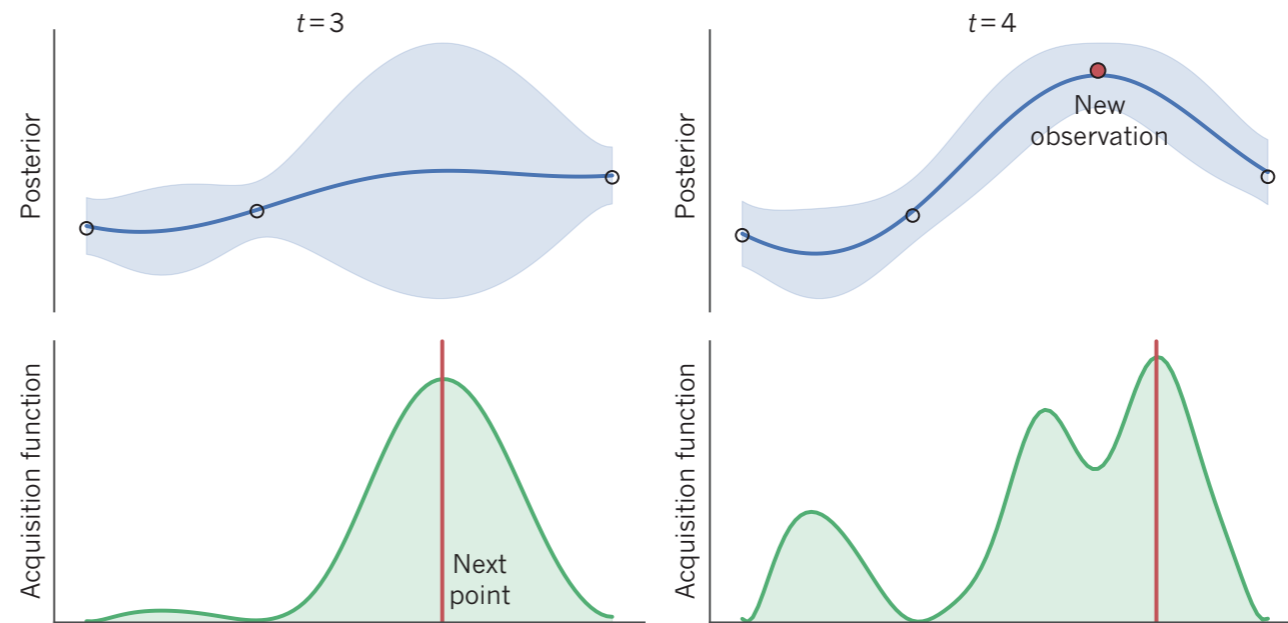
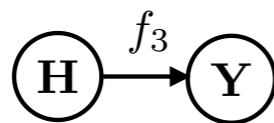
Deep auto-encoder (supervised)



GPLVM (unsupervised)



$$f_i \sim GP$$



Bayesian optimization in latent space

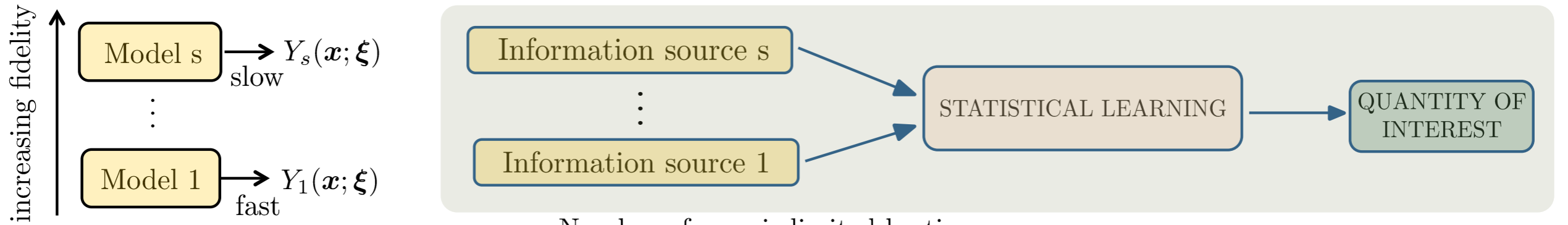
Technical approach:

- Non-linear dimensionality reduction using supervised deep auto-encoders and/or unsupervised GPLVMs
- Bayesian optimization in the low-dimensional latent space

Lawrence, N. D. "Gaussian process latent variable models for visualisation of high dimensional data." *Advances in neural information processing systems* 16.3 (2004): 329-336.

Shahriari, Bobak, et al. "Taking the human out of the loop: A review of bayesian optimization." *Proceedings of the IEEE* 104.1 (2016): 148-175.

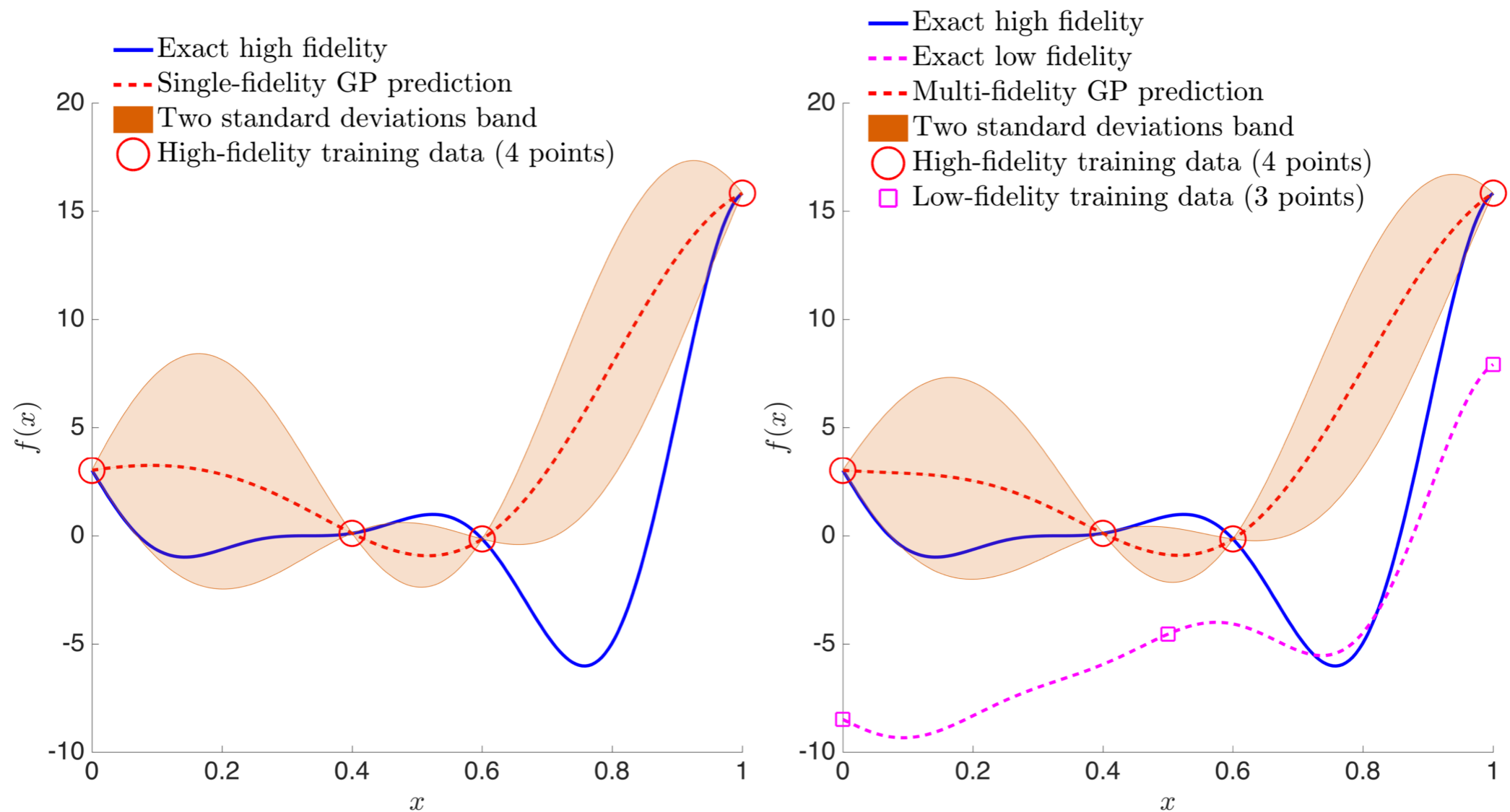
Multi-fidelity modeling



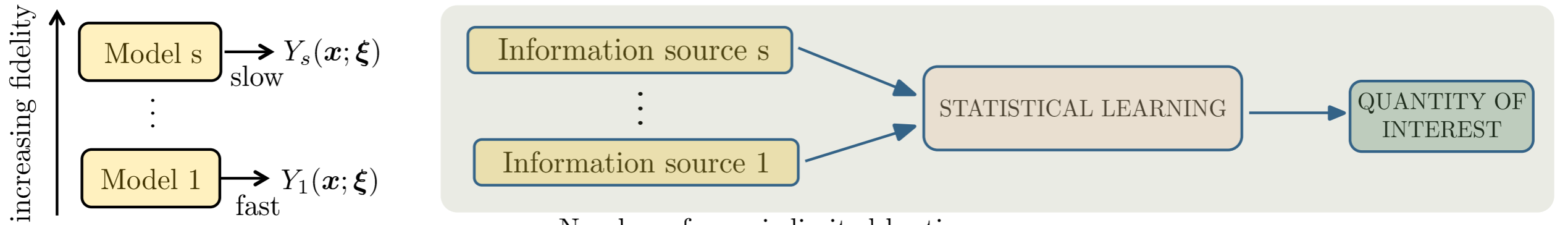
Number of runs is limited by time and computational resources

We cannot compute at all $(\mathbf{x}; \xi)$

Prediction of $Z_i(\mathbf{x}) = \mathbb{E}[f(Y_i(\mathbf{x}; \xi))]$ is a problem of **statistical inference**



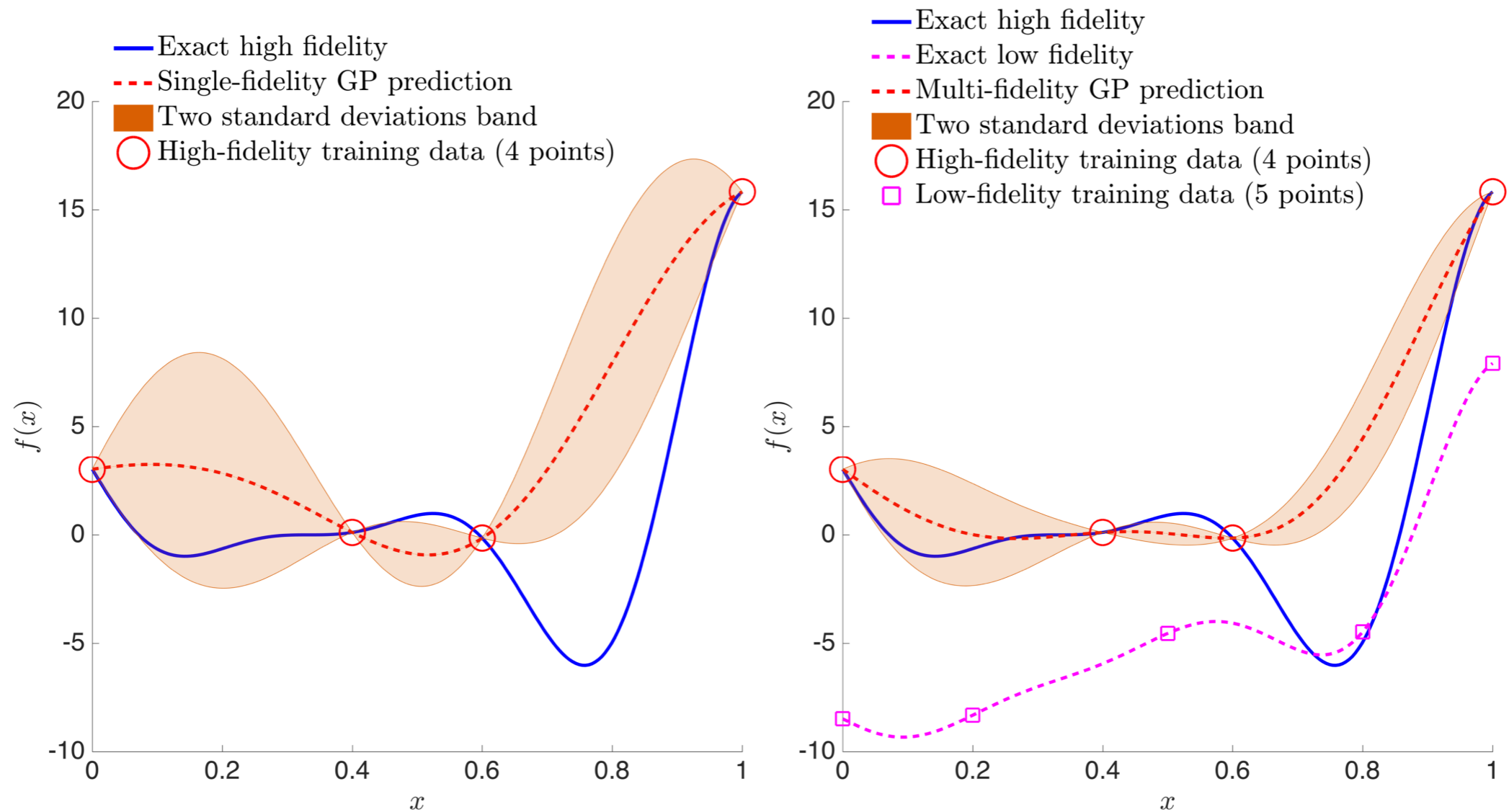
Multi-fidelity modeling



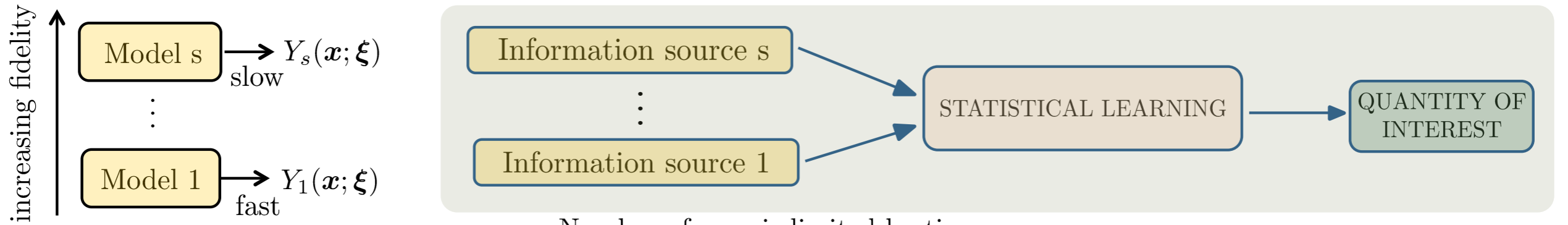
Number of runs is limited by time and computational resources

We cannot compute at all $(\mathbf{x}; \xi)$

Prediction of $Z_i(\mathbf{x}) = \mathbb{E}[f(Y_i(\mathbf{x}; \xi))]$ is a problem of **statistical inference**



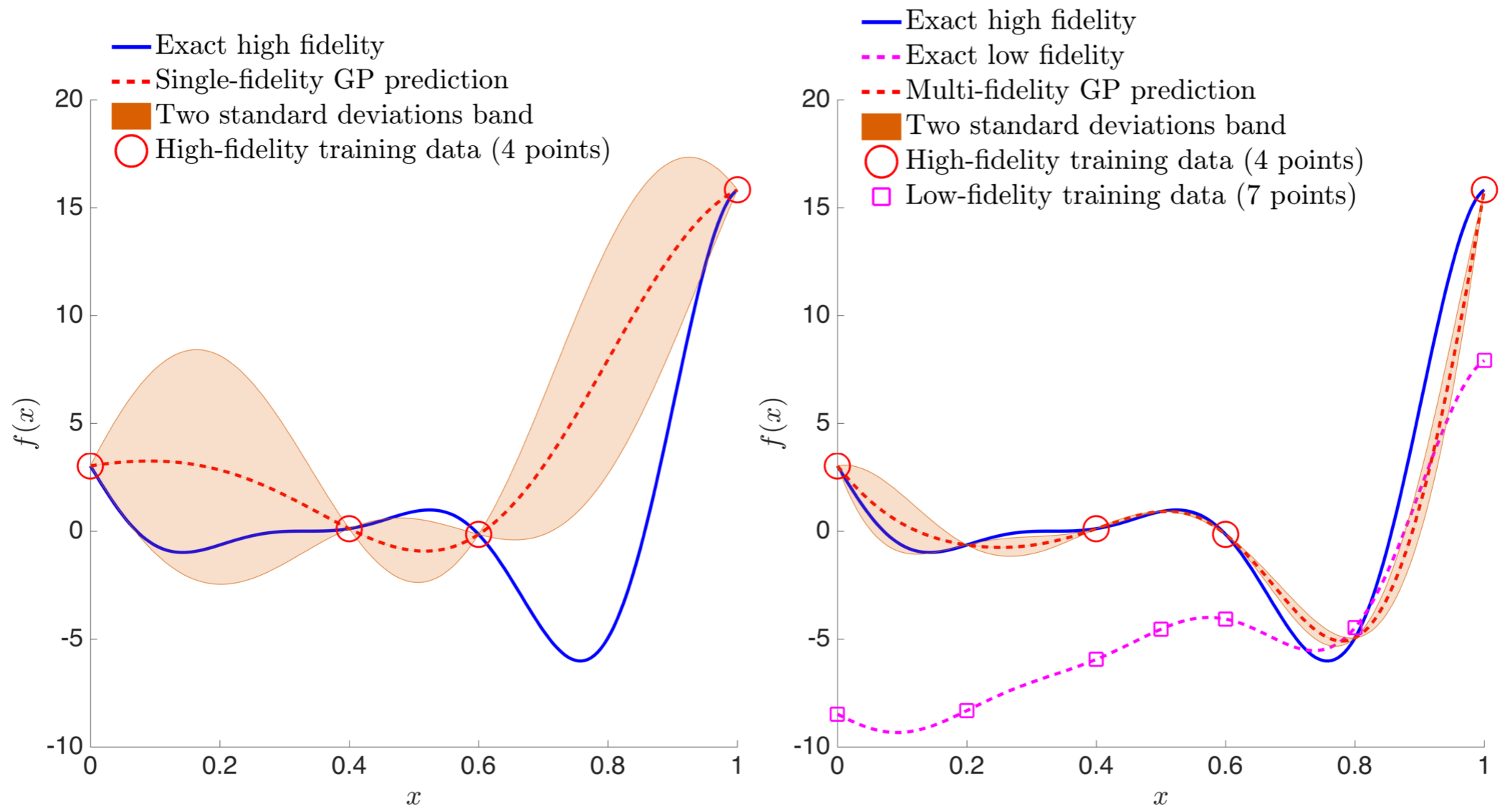
Multi-fidelity modeling



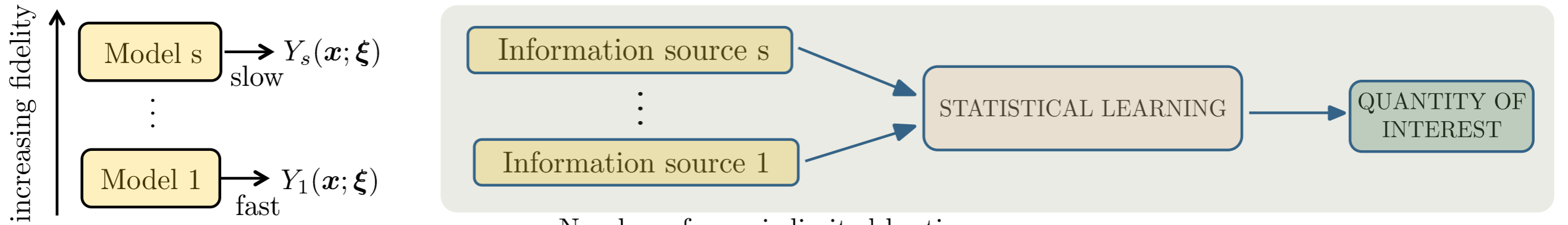
Number of runs is limited by time and computational resources

We cannot compute at all $(\mathbf{x}; \xi)$

Prediction of $Z_i(\mathbf{x}) = \mathbb{E}[f(Y_i(\mathbf{x}; \xi))]$ is a problem of **statistical inference**



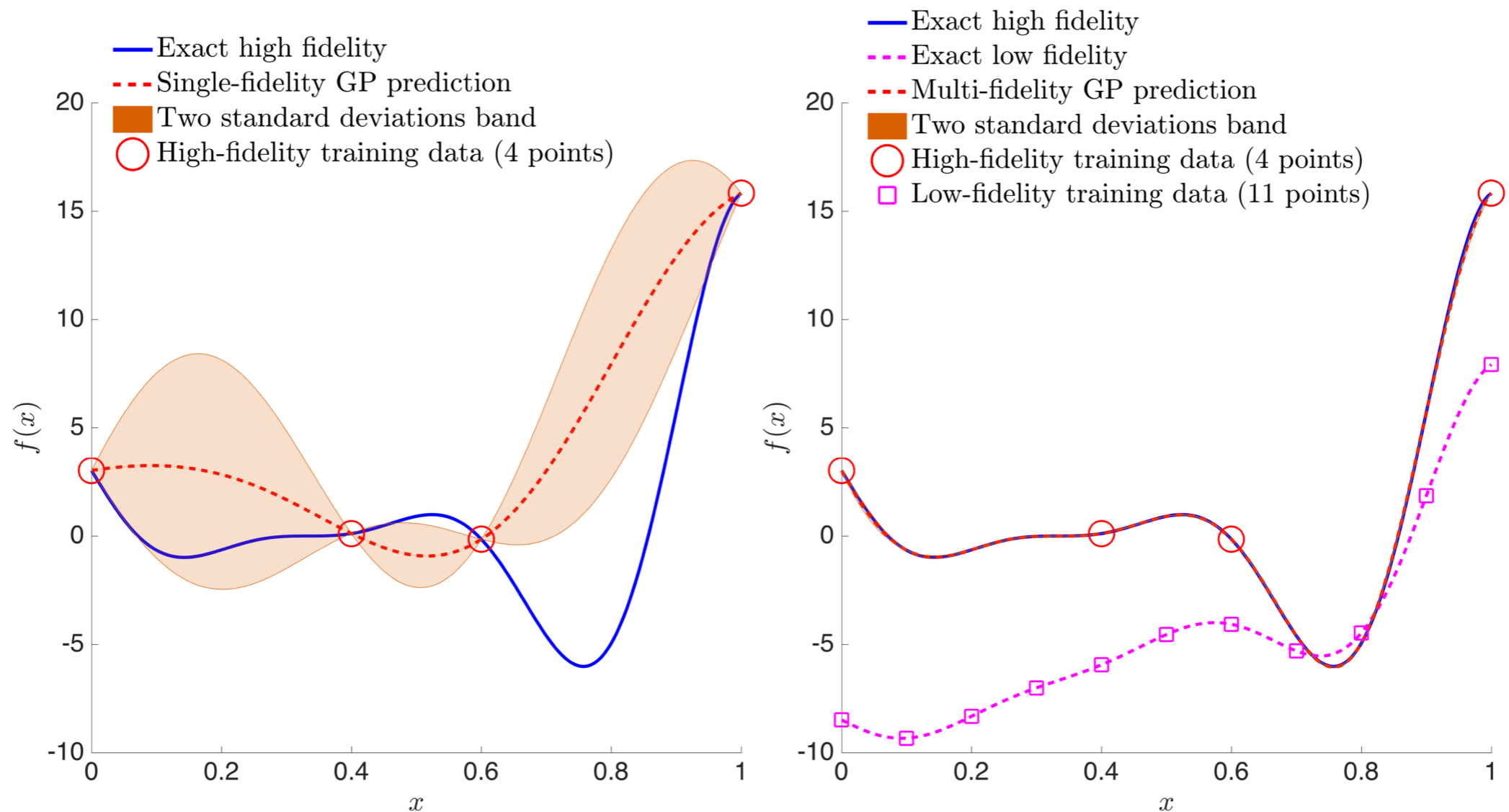
Multi-fidelity modeling



Number of runs is limited by time and computational resources

We cannot compute at all $(\mathbf{x}; \xi)$

Prediction of $Z_i(\mathbf{x}) = \mathbb{E}[f(Y_i(\mathbf{x}; \xi))]$ is a problem of **statistical inference**



Multi-fidelity modeling with GPs

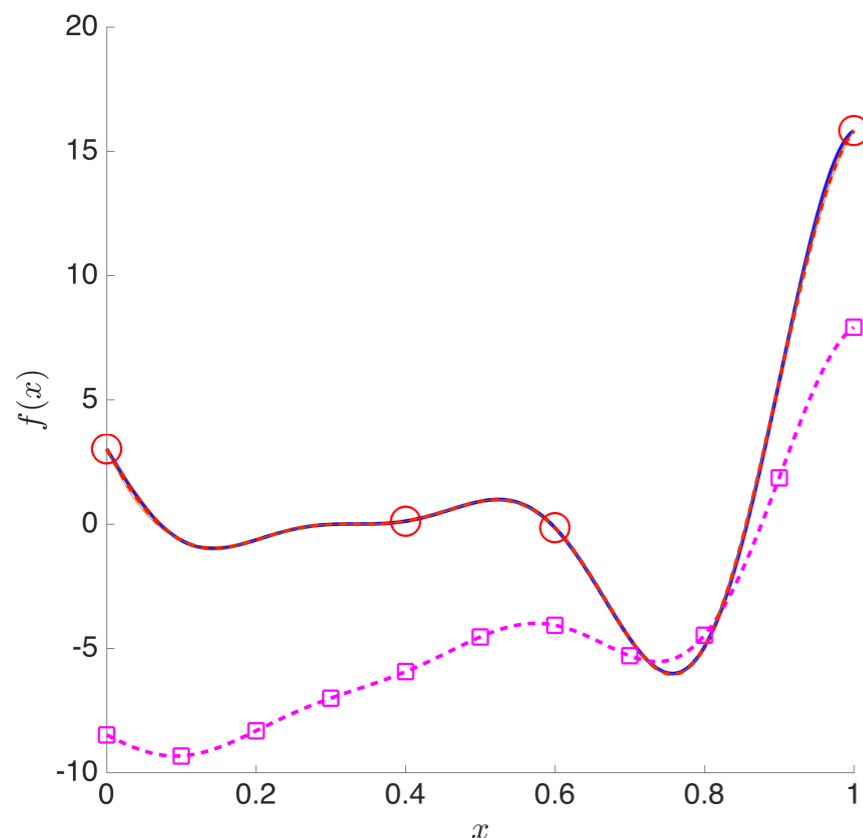


Auto-regressive model:

AR1

$$f_t(\mathbf{x}) = \rho_{t-1}(\mathbf{x})f_{t-1}(\mathbf{x}) + \delta_t(\mathbf{x})$$

$$t = 1, \dots, s$$



Predicting the Output from a Complex Computer Code When Fast Approximations Are Available

M. C. Kennedy; A. O'Hagan

Biometrika, Vol. 87, No. 1. (Mar., 2000), pp. 1-13.

Predictive posterior

$$p(f_* | \mathbf{y}, \mathbf{X}, \mathbf{x}_*) = \mathcal{N}(f_* | \mu_*, \sigma_*^2),$$

$$\mu_*(\mathbf{x}_*) = \mathbf{k}_{*N} (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{y},$$

$$\sigma_*^2(\mathbf{x}_*) = \mathbf{k}_{**} - \mathbf{k}_{*N} (\mathbf{K} + \sigma_\epsilon^2 \mathbf{I})^{-1} \mathbf{k}_{N*},$$

$$\begin{matrix} N_1 & \begin{bmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{bmatrix} \\ N_2 & \end{matrix}$$

Block covariance matrix

Key idea: Replace f_{t-1} with the GP posterior of the previous level \tilde{f}_{t-1}

$$f_t(\mathbf{x}) = \rho_{t-1}(\mathbf{x}) f_{t-1}(\mathbf{x}) + \delta_t(\mathbf{x})$$

$$\tilde{f}_{t-1} \sim f_{t-1} | \mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_{t-1}$$

Theorem (LeGratiet, 2014):

The predictive posterior of the recursive scheme has exactly the same distribution with the the fully coupled model given a nested experimental design. *(under the assumption of nested training sets!)*