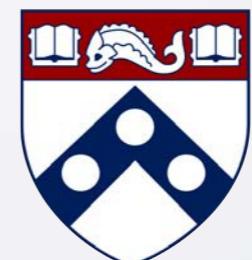


# ENM 540: Data-driven modeling and probabilistic scientific computing

*Scaling up Gaussian processes*

*Active learning & Bayesian optimization*

Paris Perdikaris  
March 13, 2018



**Penn**  
UNIVERSITY of PENNSYLVANIA

# Sparse Gaussian processes

---

## Sparse Gaussian Processes using Pseudo-inputs

---

Edward Snelson

Zoubin Ghahramani

Inference in a GP has the following demands:

Complexity:  $\mathcal{O}(n^3)$   
Storage:  $\mathcal{O}(n^2)$

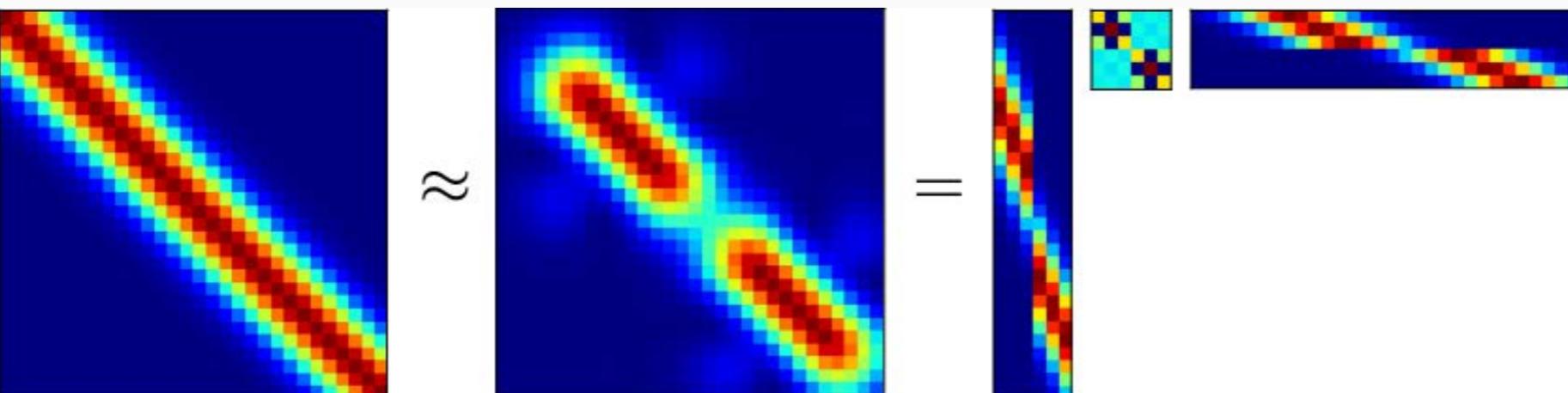
Inference in a *sparse* GP has the following demands:

Complexity:  $\mathcal{O}(nm^2)$   
Storage:  $\mathcal{O}(nm)$

where we get to pick m!

[http://gpss.cc/gpss15/talks/talk\\_james.pdf](http://gpss.cc/gpss15/talks/talk_james.pdf)

# Sparse Gaussian processes



$$\mathbf{K}_{nn} \approx \mathbf{Q}_{nn} = \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn}$$

Instead of inverting  $\mathbf{K}_{nn}$ , we make a low rank (or Nyström) approximation, and invert  $\mathbf{K}_{mm}$  instead.

[http://gpss.cc/gpss15/talks/talk\\_james.pdf](http://gpss.cc/gpss15/talks/talk_james.pdf)

# Sparse Gaussian processes

## CHRONOLOGY

### Subset of data

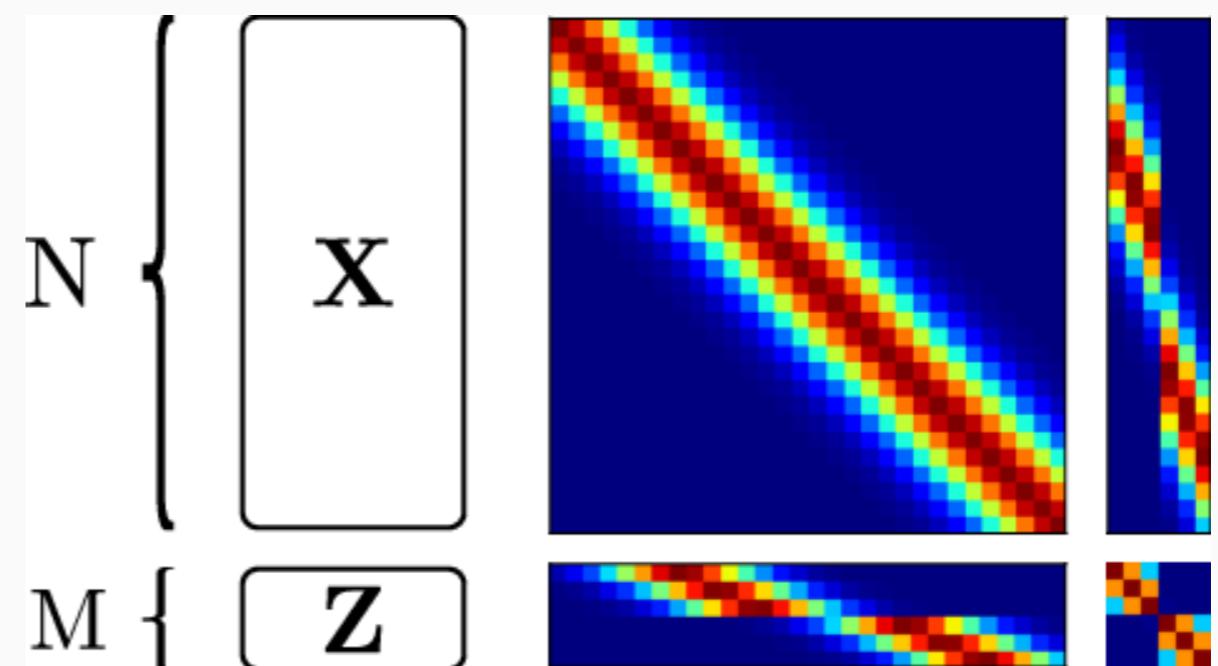
- Silverman 1985 (subset of regressors)
- Smola and Bartlett 2001 (greedy selection)

### Pseudo-input approximations

- Snelson and Ghahramani (2005), Snelson (2007)

### Variational approximations

- Titsias (2009) – derived a variational bound
- Matthews et al. (2015) – showed this minimised KL between processes



[http://gpss.cc/gpss15/talks/talk\\_james.pdf](http://gpss.cc/gpss15/talks/talk_james.pdf)

*Snelson, E., & Ghahramani, Z. (2006). Sparse Gaussian processes using pseudo-inputs. In Advances in neural information processing systems (pp. 1257-1264).*

# Sparse Gaussian processes

Take and extra  $M$  points on the function,  $\mathbf{u} = f(\mathbf{Z})$ .

$$p(y, f, u) = p(y | f)p(f | u)p(u)$$

$$p(y | f) = \mathcal{N}(y | f, \sigma^2 I)$$

$$p(f | u) = \mathcal{N}\left(f | K_{nm}K_{mm}u, \tilde{K}\right)$$

$$p(u) = \mathcal{N}(u | 0, K_{mm})$$

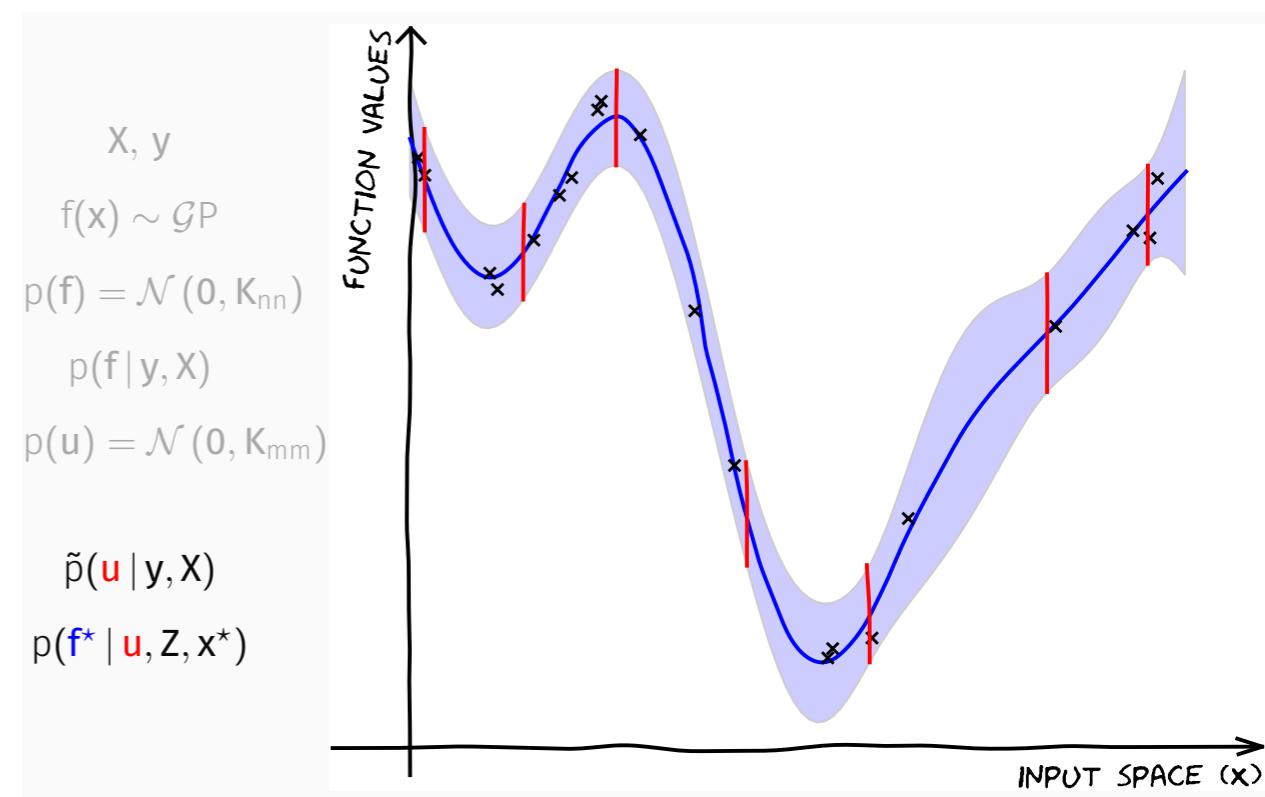
Instead of doing

$$p(f | y, X) = \frac{p(y | f)p(f | X)}{\int p(y | f)p(f | X) df}$$

$$\begin{aligned} x, y \\ f(x) &\sim \mathcal{GP} \\ p(f) &= \mathcal{N}(0, K_{nn}) \end{aligned}$$

$$\begin{aligned} p(f | y, X) \\ p(u) &= \mathcal{N}(0, K_{mm}) \end{aligned}$$

$$\begin{aligned} \tilde{p}(u | y, X) \\ p(f^* | u, Z, x^*) \end{aligned}$$



We'll do

$$p(u | y, Z) = \frac{p(y | u)p(u | Z)}{\int p(y | u)p(u | Z) du}$$

# Sparse Gaussian processes

$$Q_{\mathbf{ff}} = K_{\mathbf{fu}} K_{\mathbf{uu}}^{-1} K_{\mathbf{uf}}$$

$$\mathcal{F} = \frac{N}{2} \log(2\pi) + \underbrace{\frac{1}{2} \log |Q_{\mathbf{ff}} + G|}_{\text{complexity penalty}} + \underbrace{\frac{1}{2} \mathbf{y}^\top (Q_{\mathbf{ff}} + G)^{-1} \mathbf{y}}_{\text{data fit}} + \underbrace{\frac{1}{2\sigma_n^2} \text{tr}(T)}_{\text{trace term}},$$

$$G_{\text{FITC}} = \text{diag}[K_{\mathbf{ff}} - Q_{\mathbf{ff}}] + \sigma_n^2 I$$

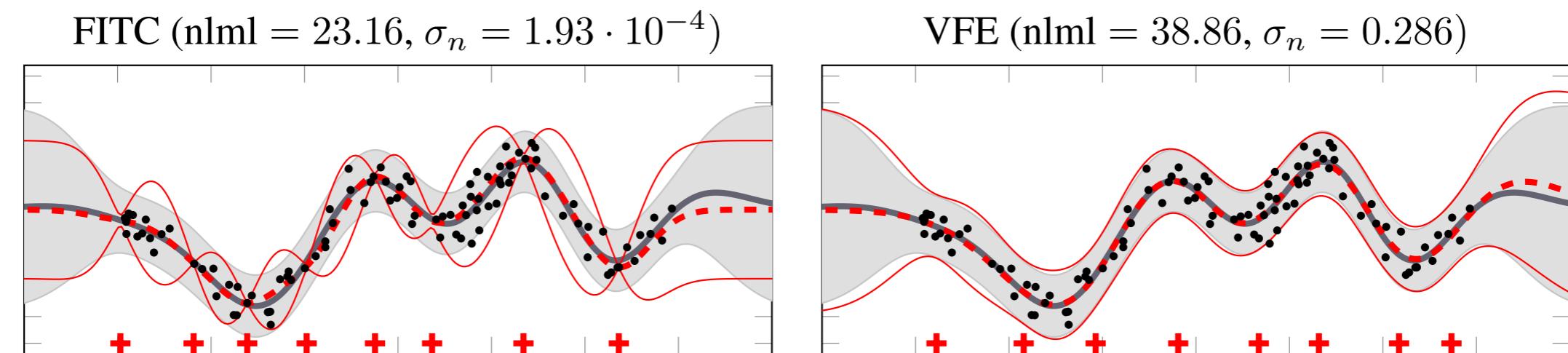
$$T_{\text{FITC}} = 0$$

*Fully Independent Training Conditional (FITC)*

$$G_{\text{VFE}} = \sigma_n^2 I$$

$$T_{\text{VFE}} = K_{\mathbf{ff}} - Q_{\mathbf{ff}}.$$

*Variational Free Energy (VFE)*



*Snelson, E., & Ghahramani, Z. (2006). Sparse Gaussian processes using pseudo-inputs. In Advances in neural information processing systems (pp. 1257-1264).*

*Titsias, M. (2009, April). Variational learning of inducing variables in sparse Gaussian processes. In Artificial Intelligence and Statistics (pp. 567-574).*

# Gaussian processes for big data

---

## Gaussian Processes for Big Data

---

**James Hensman\***  
Dept. Computer Science  
The University of Sheffield  
Sheffield, UK

**Nicolò Fusi\***  
Dept. Computer Science  
The University of Sheffield  
Sheffield, UK

**Neil D. Lawrence\***  
Dept. Computer Science  
The University of Sheffield  
Sheffield, UK

$$\log p(\mathbf{y} \mid \mathbf{X}) \geq \langle \mathcal{L}_1 + \log p(\mathbf{u}) - \log q(\mathbf{u}) \rangle_{q(\mathbf{u})} \triangleq \mathcal{L}_3$$

$$\begin{aligned} \mathcal{L}_3 = \sum_{i=1}^n & \left\{ \log \mathcal{N} \left( y_i \mid \mathbf{k}_i^\top \mathbf{K}_{mm}^{-1} \mathbf{m}, \beta^{-1} \right) \right. \\ & - \frac{1}{2} \beta \tilde{k}_{i,i} - \frac{1}{2} \text{tr} (\mathbf{S} \boldsymbol{\Lambda}_i) \Big\} \\ & - \text{KL} (q(\mathbf{u}) \parallel p(\mathbf{u})) \end{aligned}$$

# Gaussian processes for big data

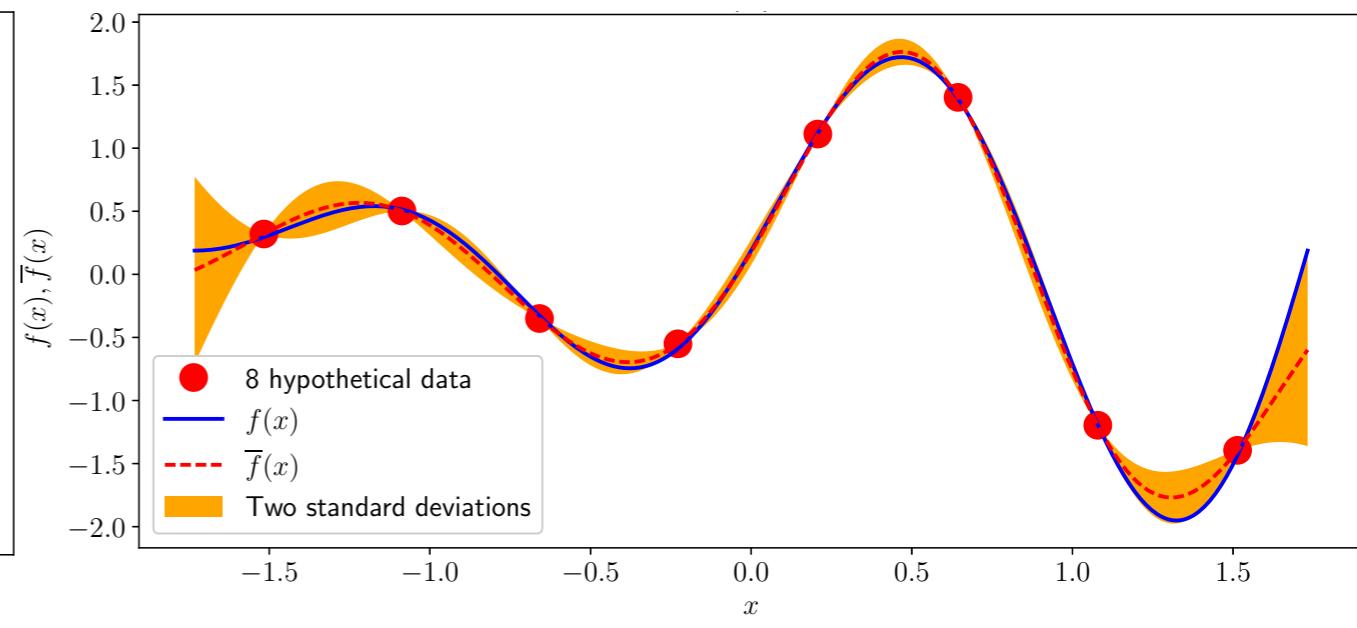
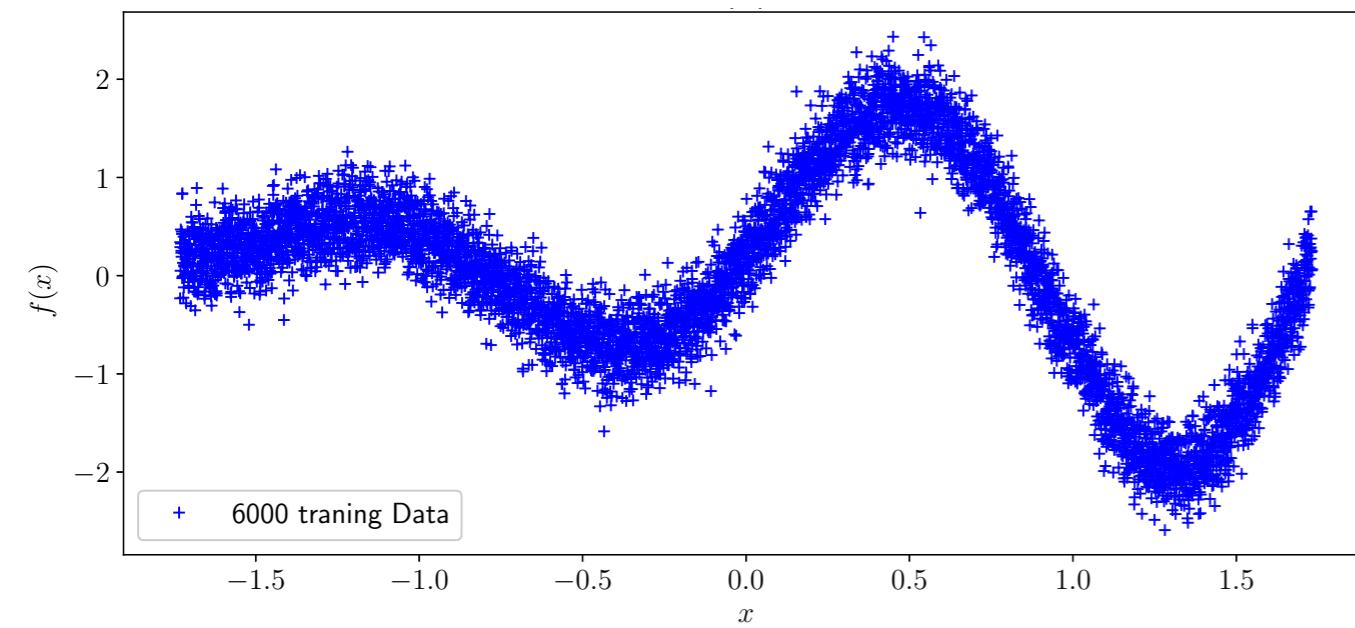
## Parametric Gaussian Process Regression for Big Data

Maziar Raissi\*

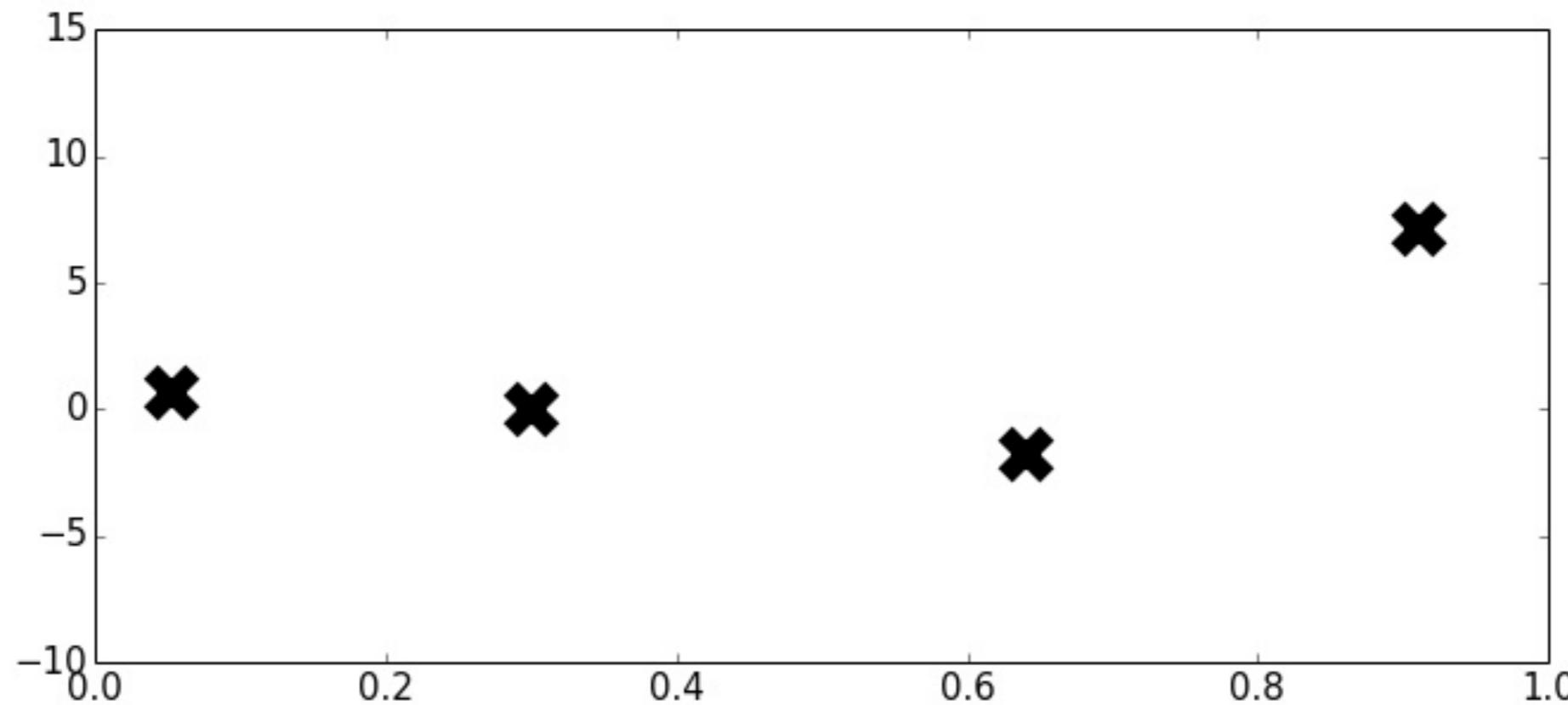
Division of Applied Mathematics  
Brown University  
Providence, RI 02912  
maziar\_raissi@brown.edu

$$\mathbf{u} \sim \mathcal{N}(\mathbf{m}, \mathbf{S})$$

$$\begin{aligned}\mathbf{m} &\leftarrow \mu(\mathbf{Z}; \boldsymbol{\theta}, \mathbf{m}) + \Sigma(\mathbf{Z}, \widetilde{\mathbf{X}}; \boldsymbol{\theta}, \mathbf{S}) \left( \Sigma(\widetilde{\mathbf{X}}, \widetilde{\mathbf{X}}; \boldsymbol{\theta}, \mathbf{S}) + \sigma_\epsilon^2 \mathbf{I} \right)^{-1} \left[ \widetilde{\mathbf{y}} - \mu(\widetilde{\mathbf{X}}; \boldsymbol{\theta}, \mathbf{m}) \right], \\ \mathbf{S} &\leftarrow \Sigma(\mathbf{Z}, \mathbf{Z}; \boldsymbol{\theta}, \mathbf{S}) - \Sigma(\mathbf{Z}, \widetilde{\mathbf{X}}; \boldsymbol{\theta}, \mathbf{S}) \left( \Sigma(\widetilde{\mathbf{X}}, \widetilde{\mathbf{X}}; \boldsymbol{\theta}, \mathbf{S}) + \sigma_\epsilon^2 \mathbf{I} \right)^{-1} \Sigma(\widetilde{\mathbf{X}}, \mathbf{Z}; \boldsymbol{\theta}, \mathbf{S}).\end{aligned}$$

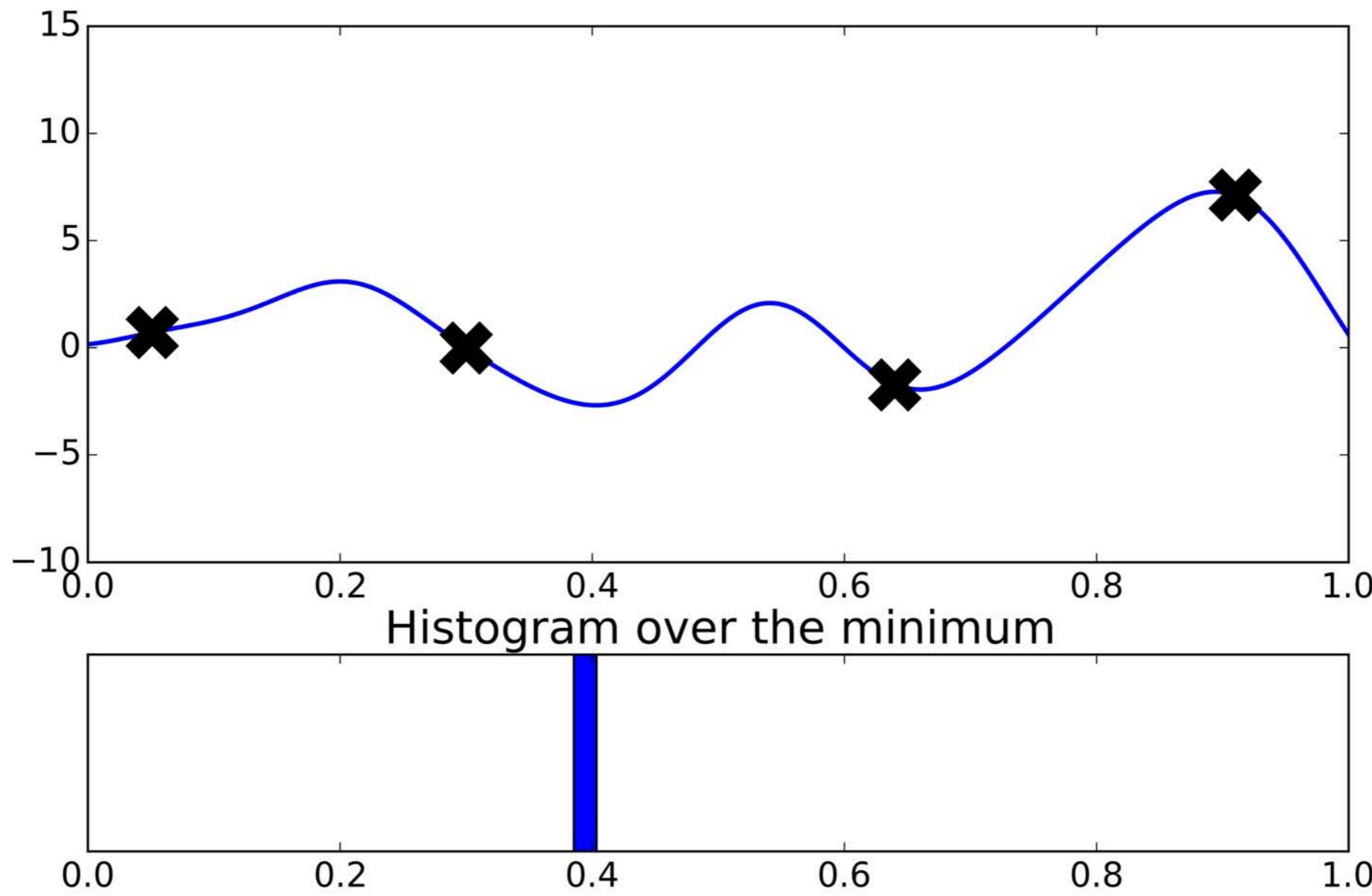


# Bayesian optimization

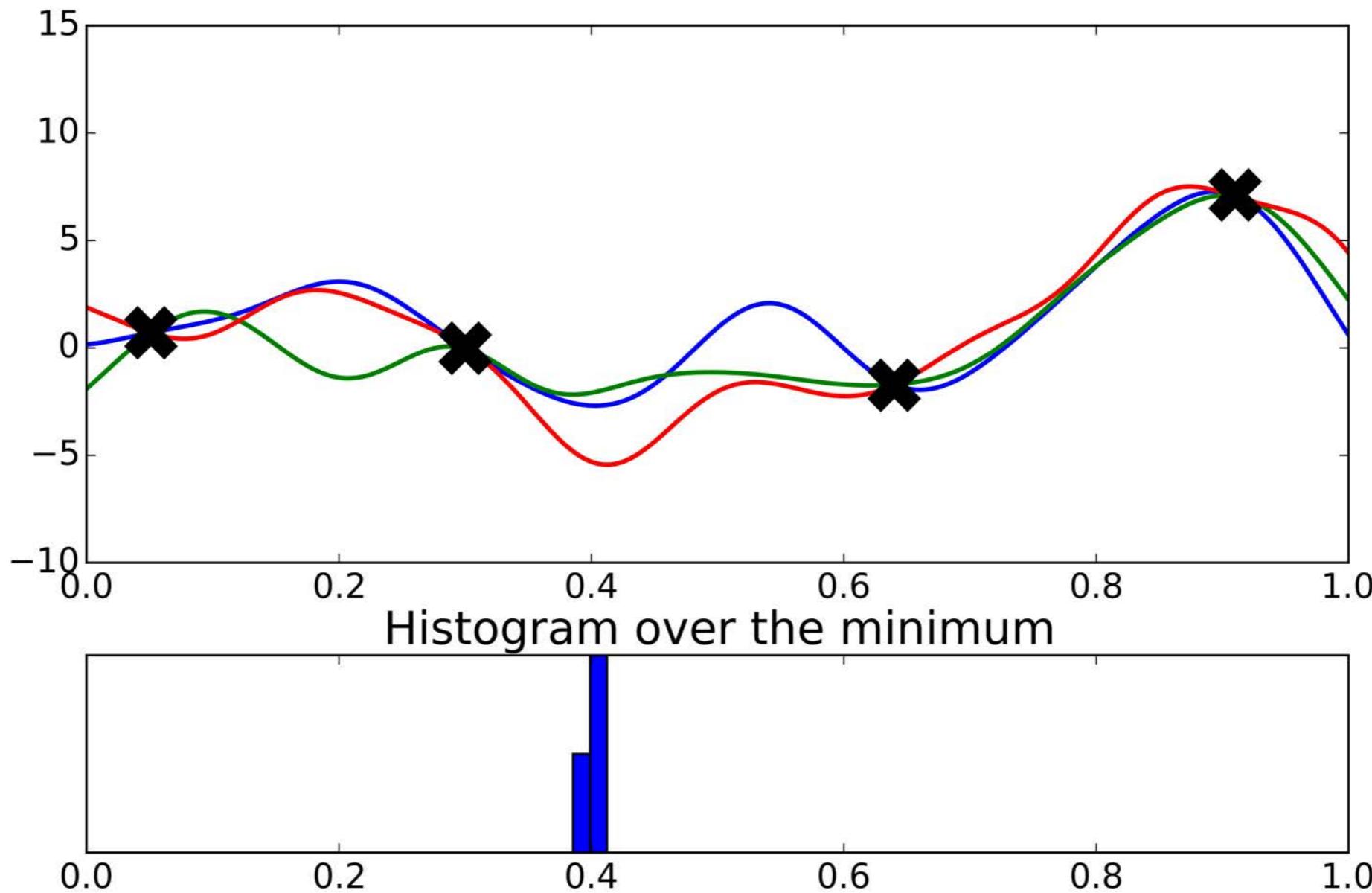


**Where is the minimum of  $f$ ?**  
**Where should we take the next evaluation?**

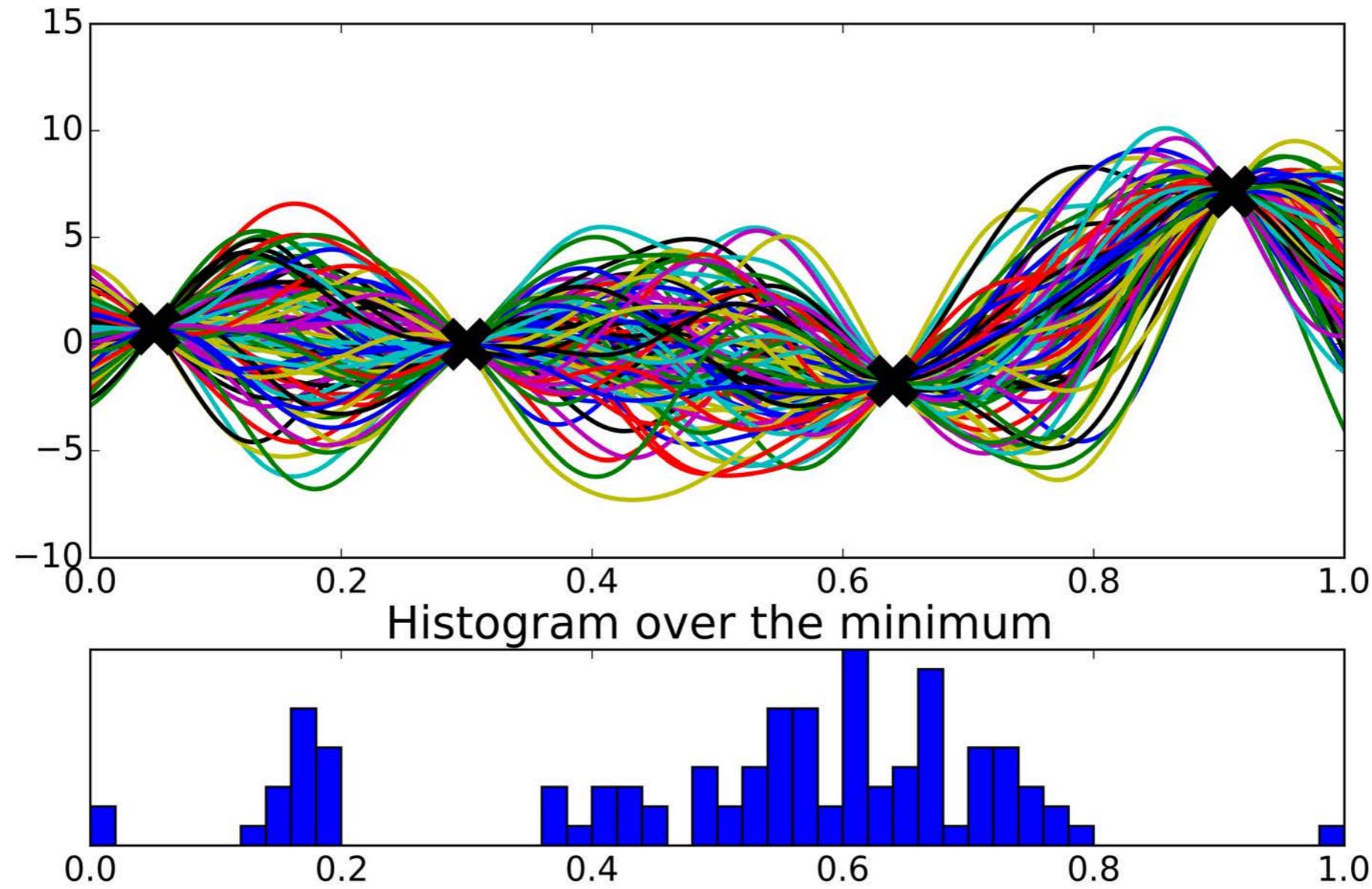
# Bayesian optimization



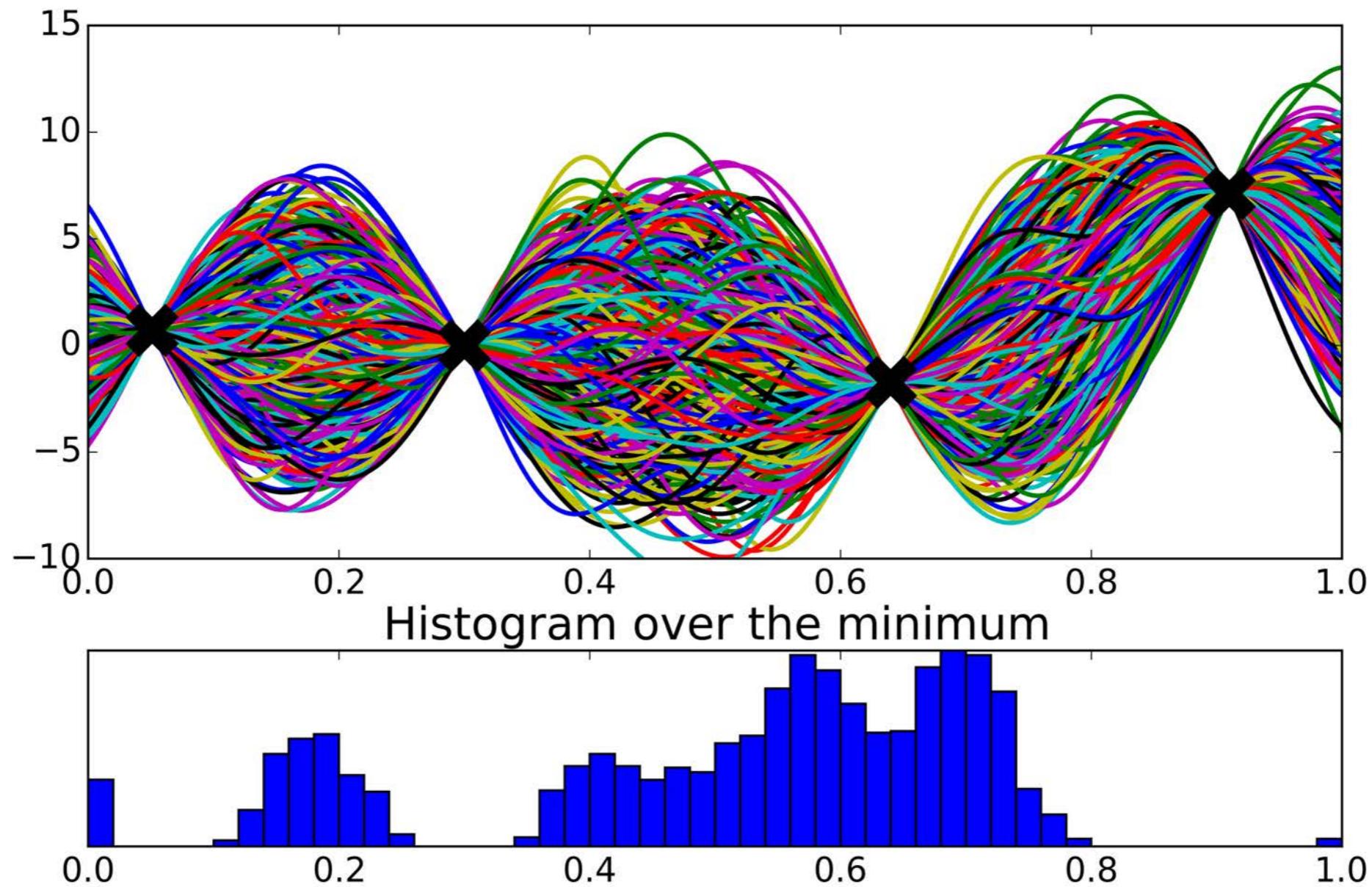
# Bayesian optimization



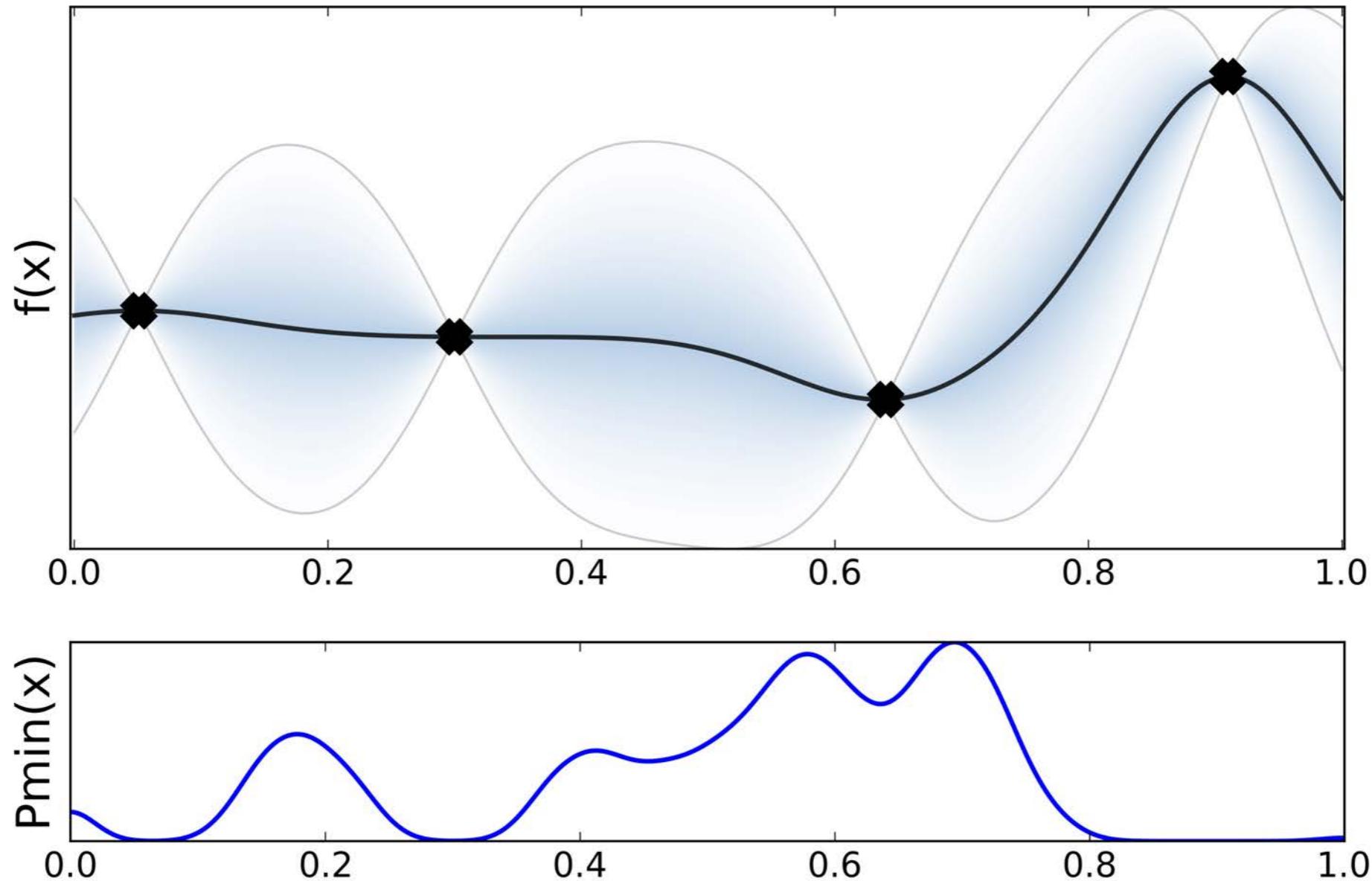
# Bayesian optimization



# Bayesian optimization



# Bayesian optimization



$x^* = \arg \min f(x)$  ..but  $p(x^* | \mathcal{D})$  is not tractable.

# 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:**  $g(\mathbf{x})$  is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

**Idea:** Approximate  $g(\mathbf{x})$  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

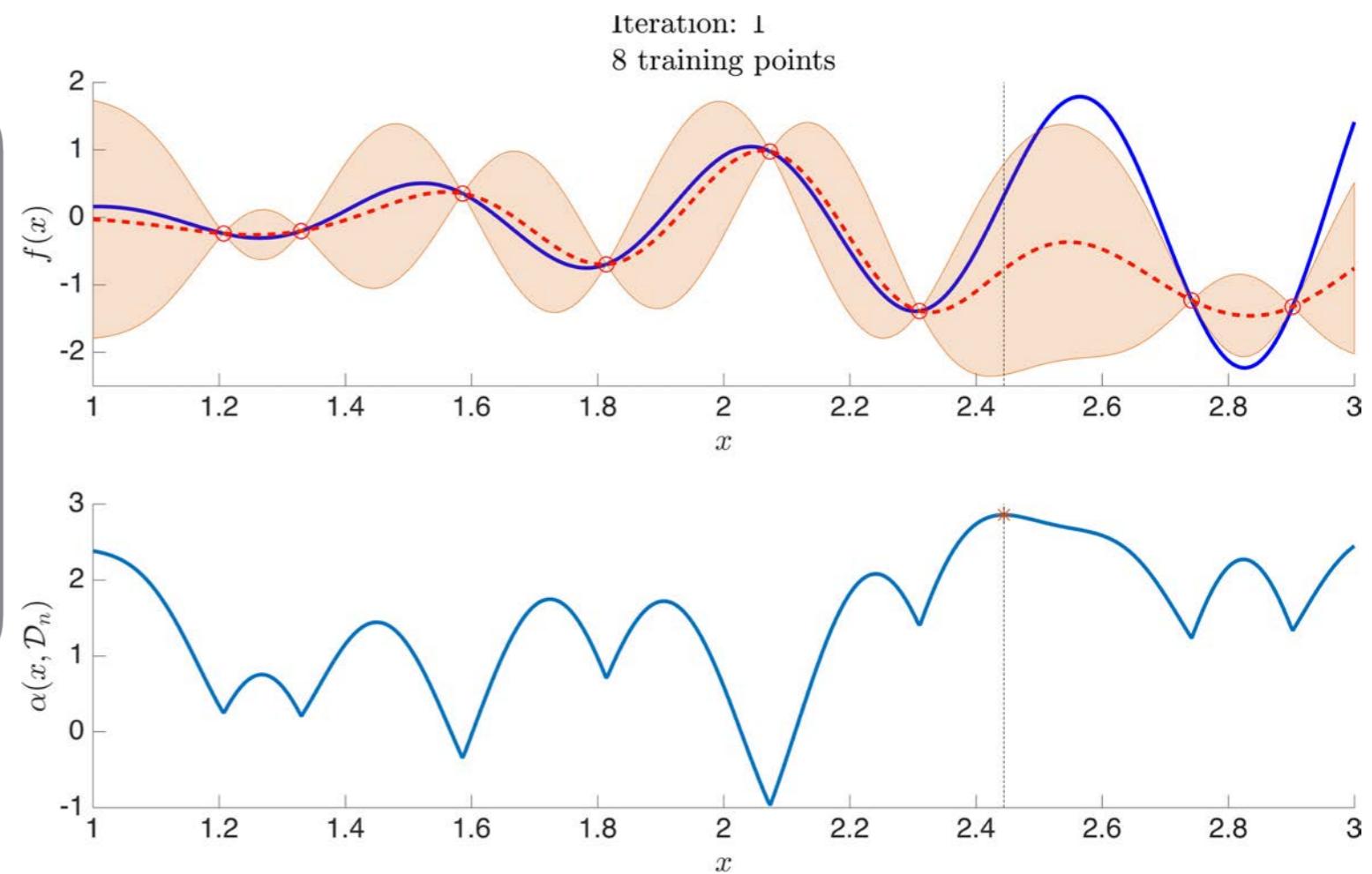
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v \mid \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

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 super-quintile 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:**  $g(\mathbf{x})$  is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

**Idea:** Approximate  $g(\mathbf{x})$  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

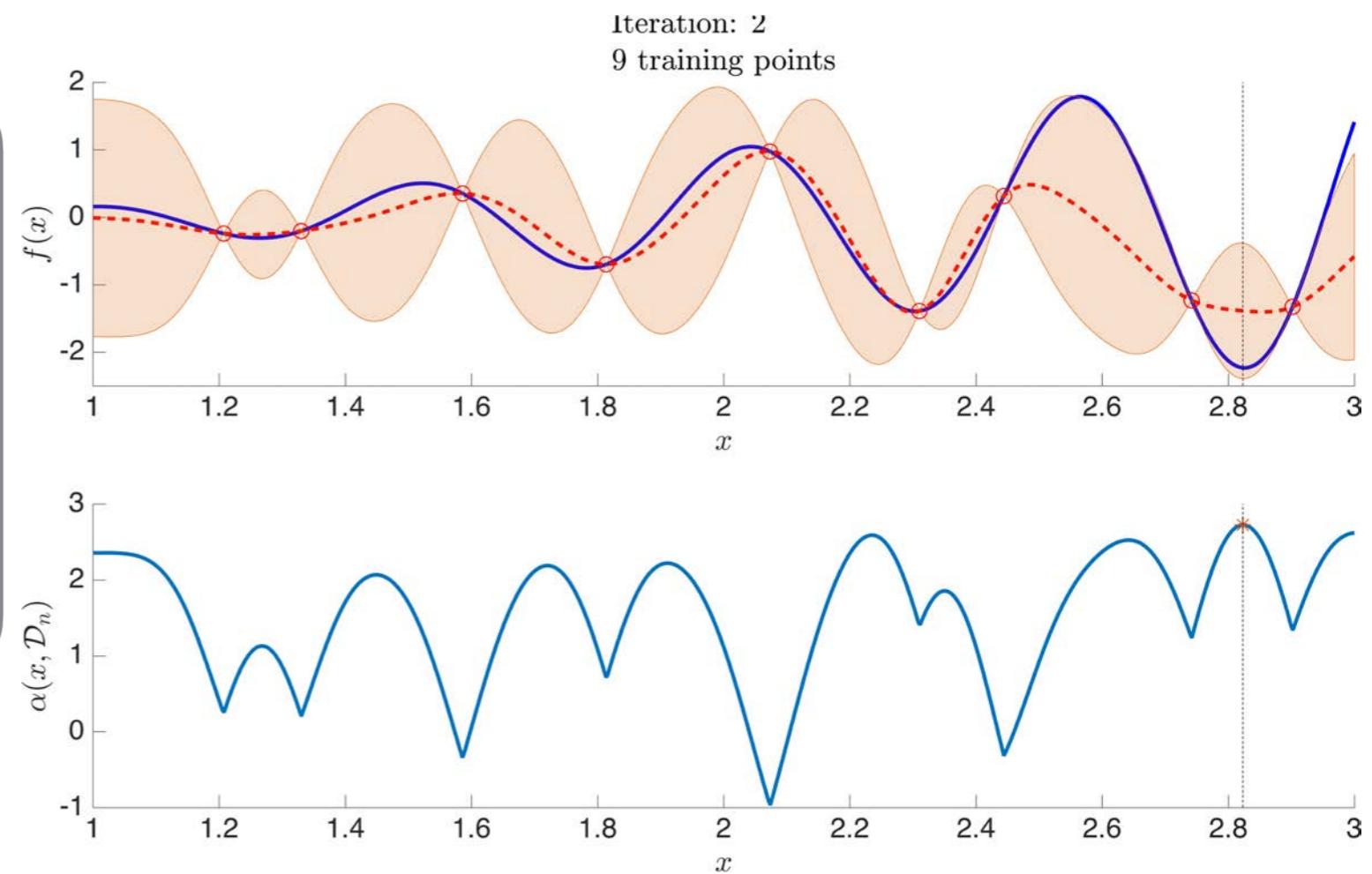
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v \mid \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

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 super-quintile 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:**  $g(\mathbf{x})$  is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

**Idea:** Approximate  $g(\mathbf{x})$  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

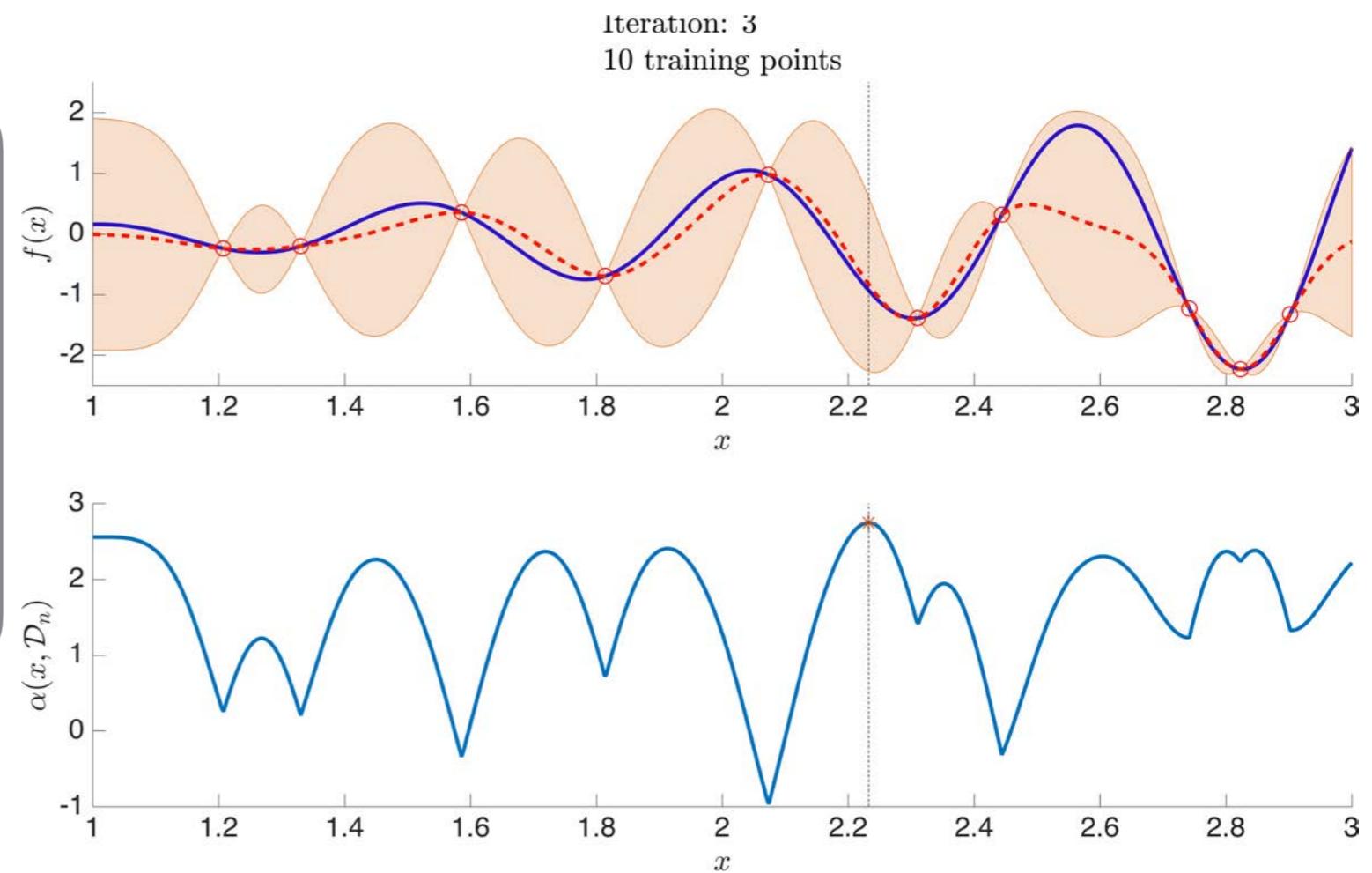
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v \mid \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

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 super-quintile 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:**  $g(\mathbf{x})$  is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

**Idea:** Approximate  $g(\mathbf{x})$  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

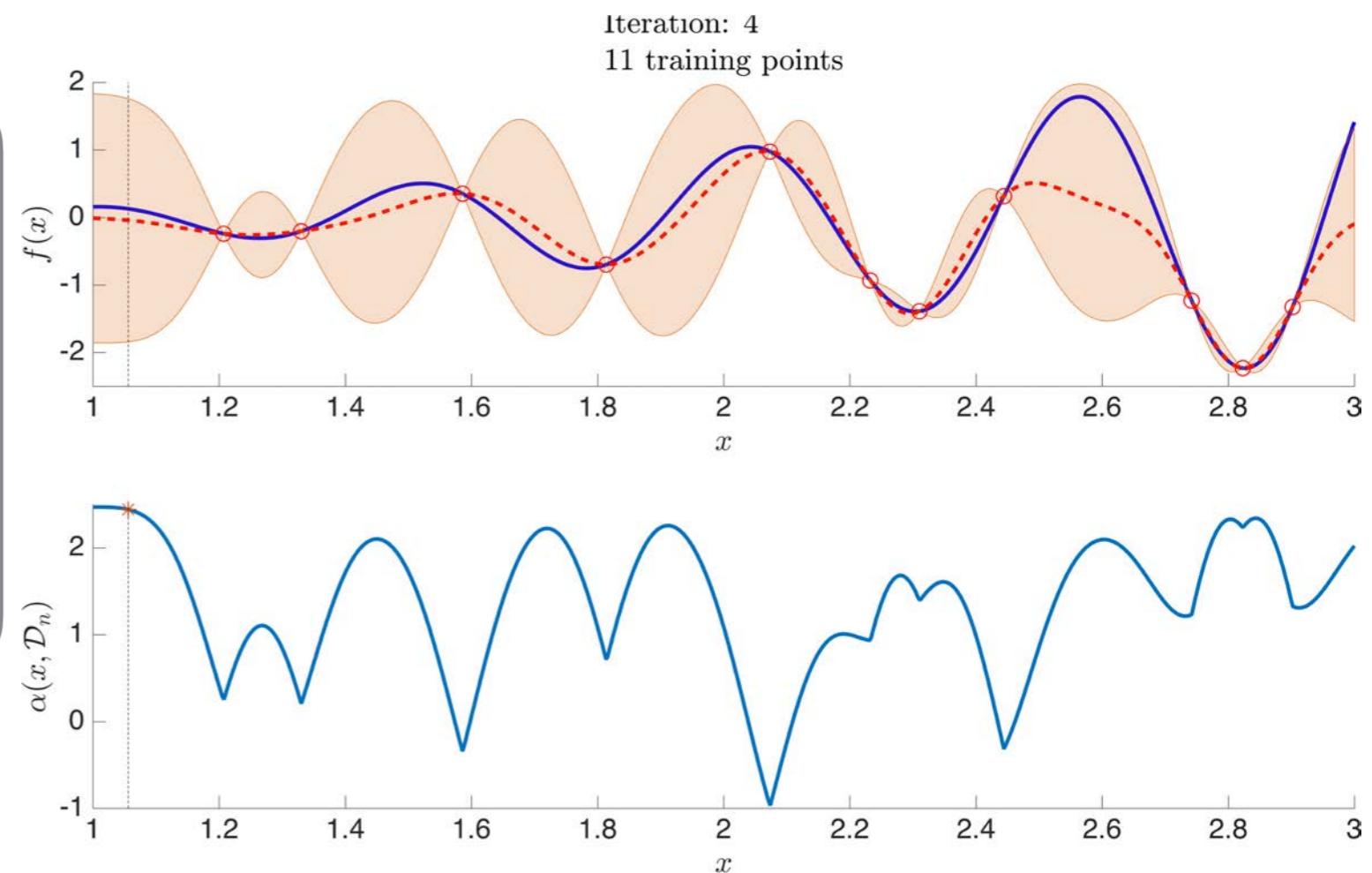
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v \mid \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

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 super-quintile 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:**  $g(\mathbf{x})$  is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

**Idea:** Approximate  $g(\mathbf{x})$  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

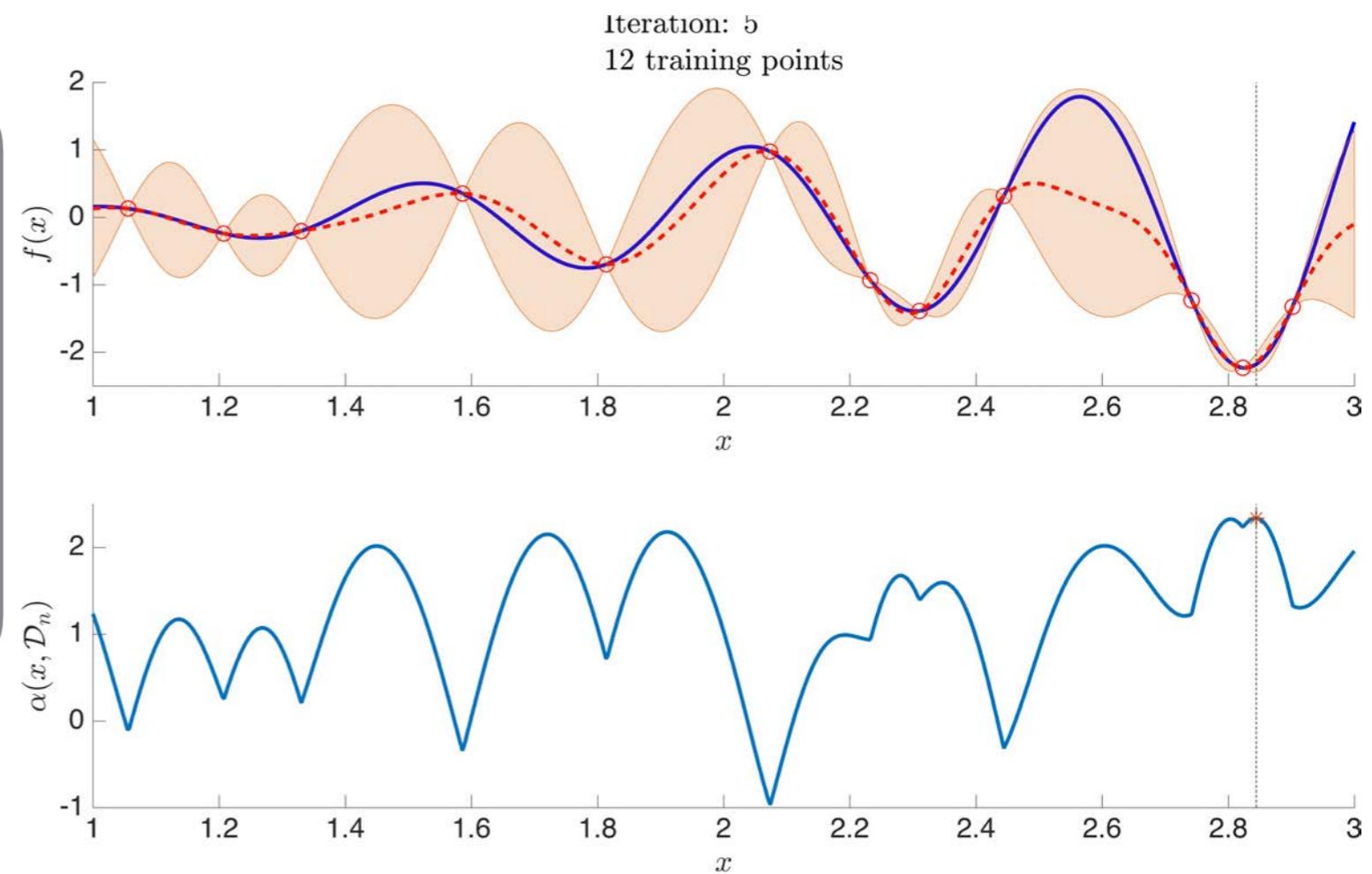
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v \mid \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

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 super-quintile 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:**  $g(\mathbf{x})$  is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

**Idea:** Approximate  $g(\mathbf{x})$  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

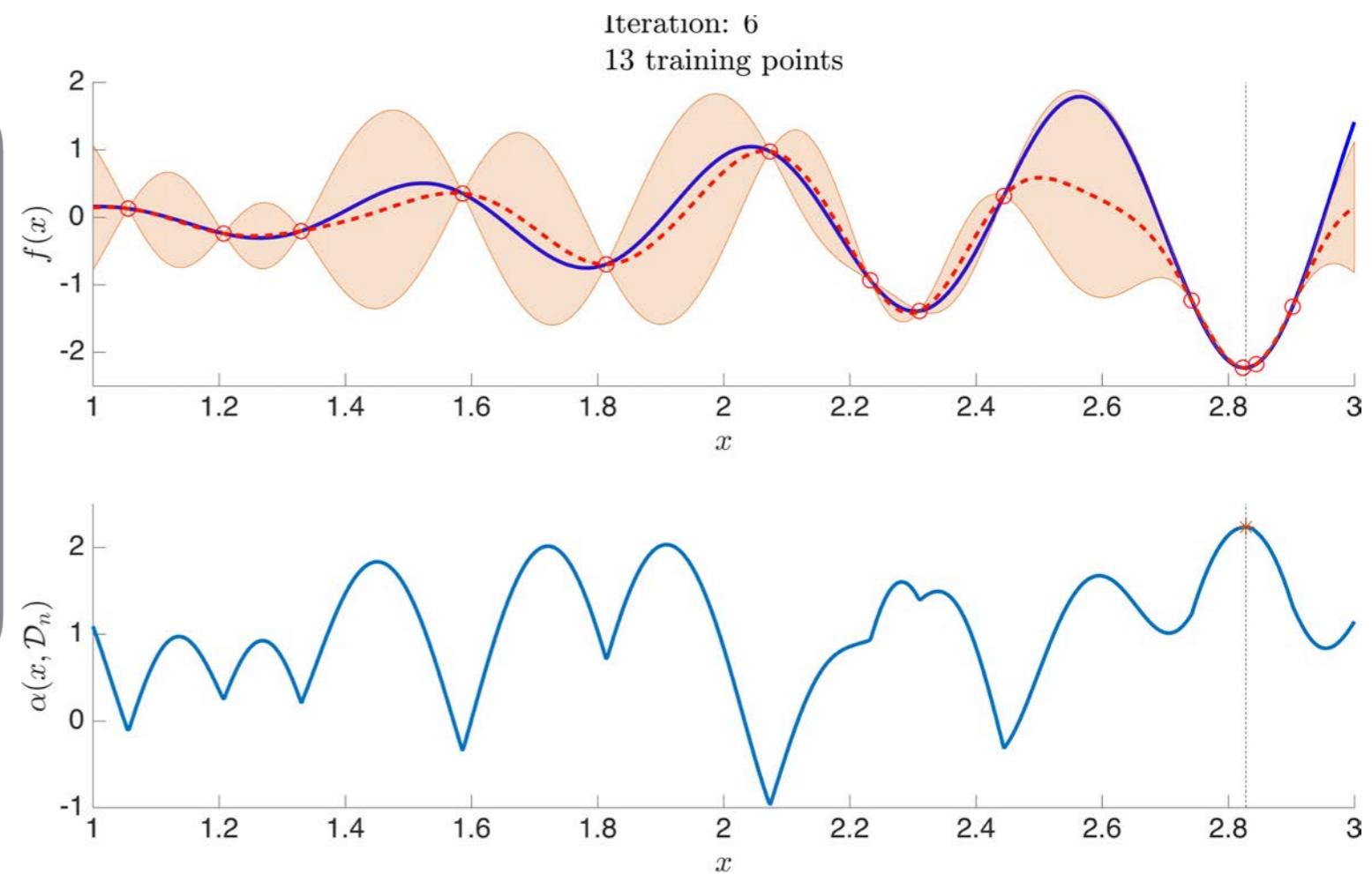
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v \mid \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

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 super-quintile 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:**  $g(\mathbf{x})$  is a black-box and expensive to evaluate objective function, noisy observations, no gradients.

**Idea:** Approximate  $g(\mathbf{x})$  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

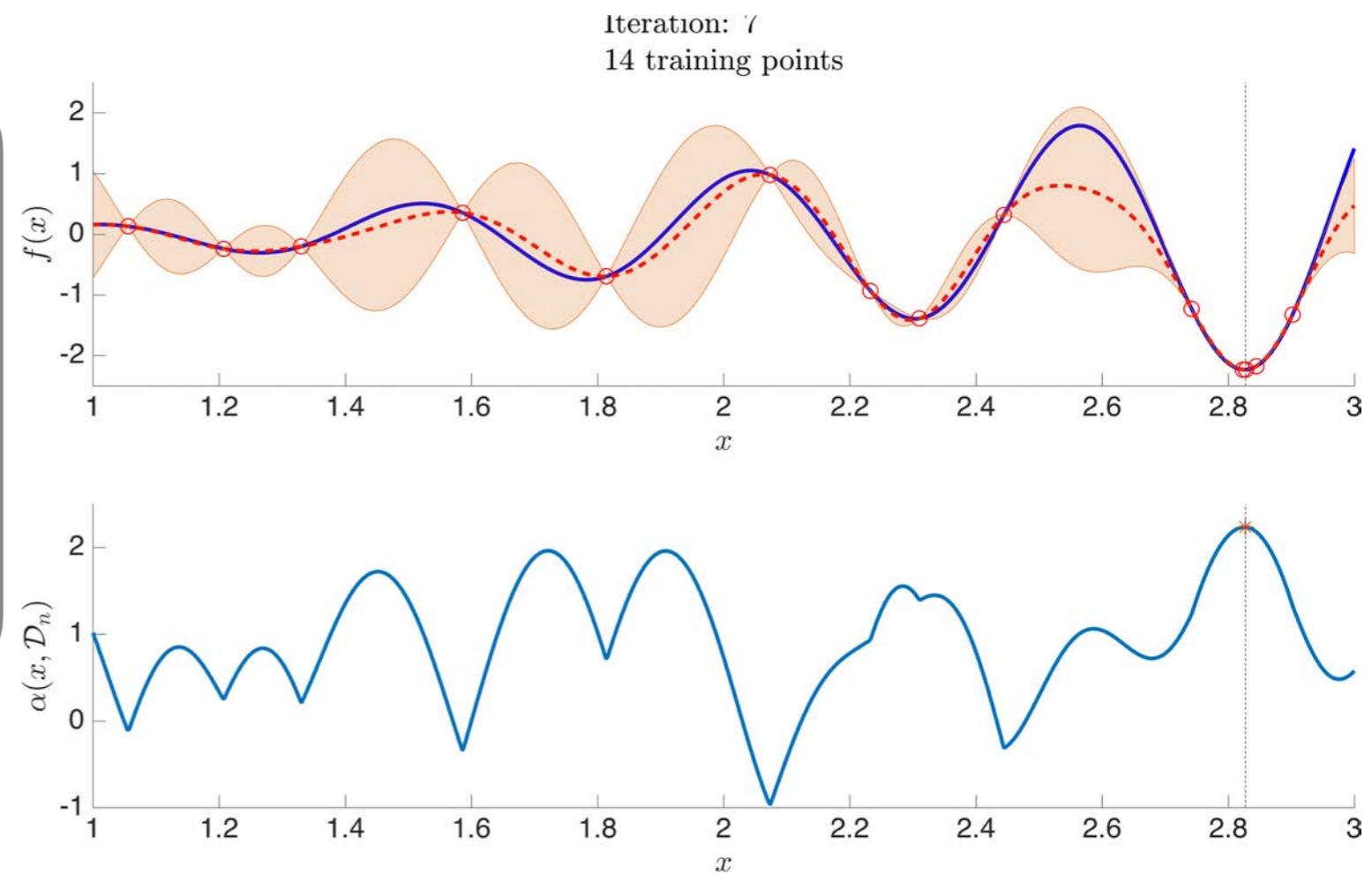
$$\alpha(\mathbf{x}; \mathcal{D}_n) = \mathbb{E}_{\theta} \mathbb{E}_{v \mid \mathbf{x}, \theta} [U(\mathbf{x}, v, \theta)]$$

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 super-quintile risk confidence bound

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

# Learning level sets

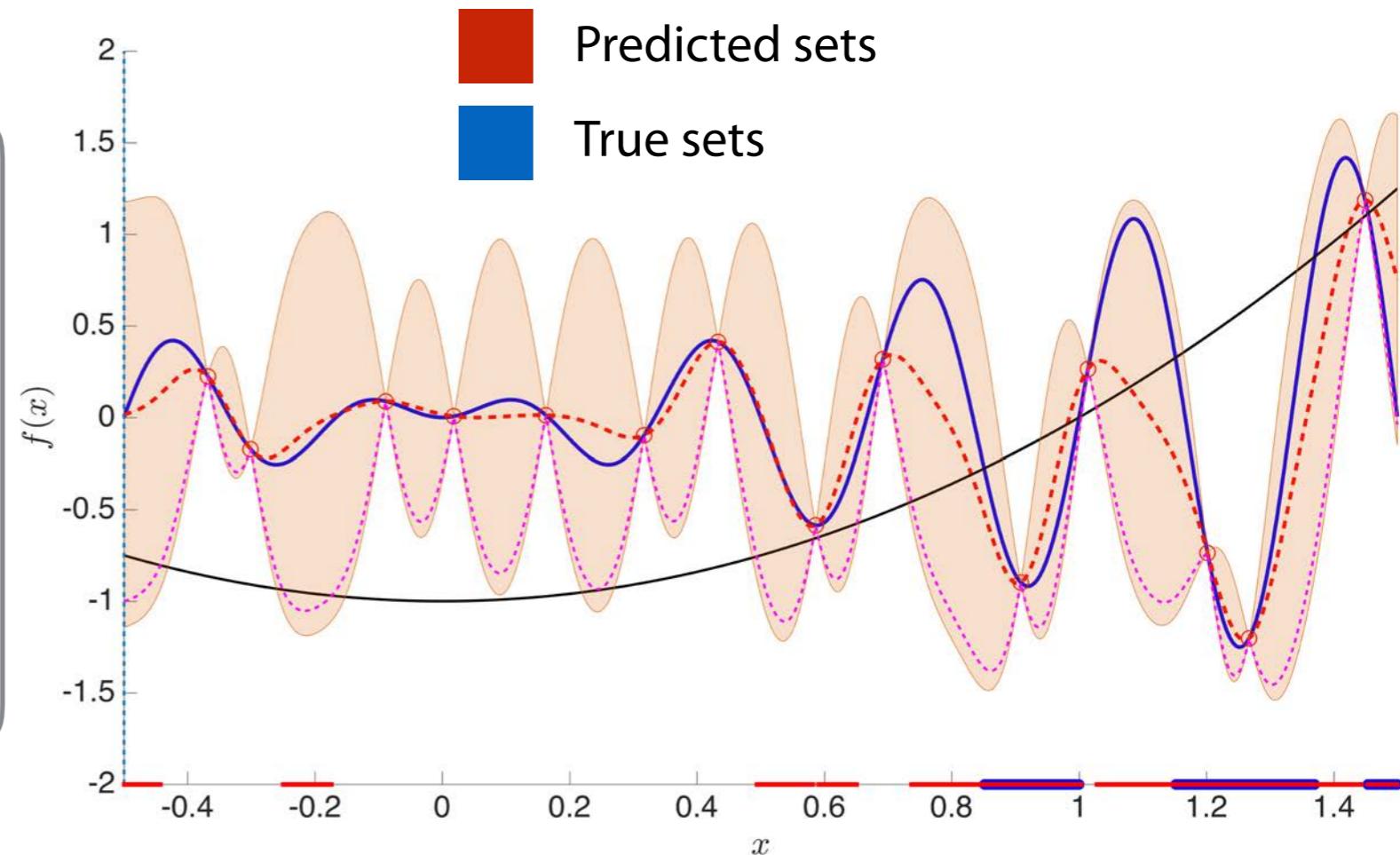
**Goal:** Identify the sets  $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t(\mathbf{x})\}$

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

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

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 sets are guaranteed to be a subset of the true sets.

# Learning level sets

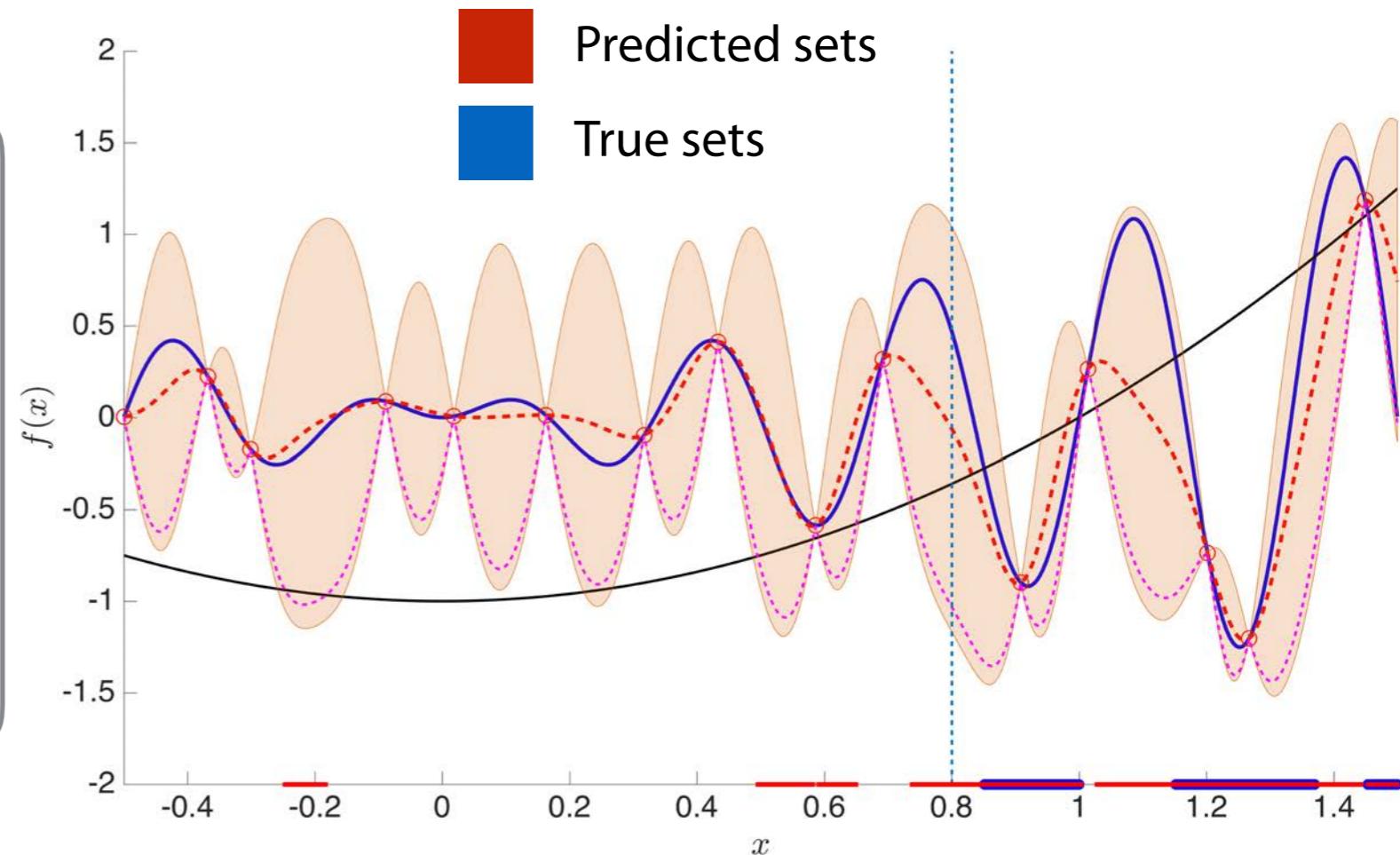
**Goal:** Identify the sets  $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t(\mathbf{x})\}$

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

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

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 sets are guaranteed to be a subset of the true sets.

# Learning level sets

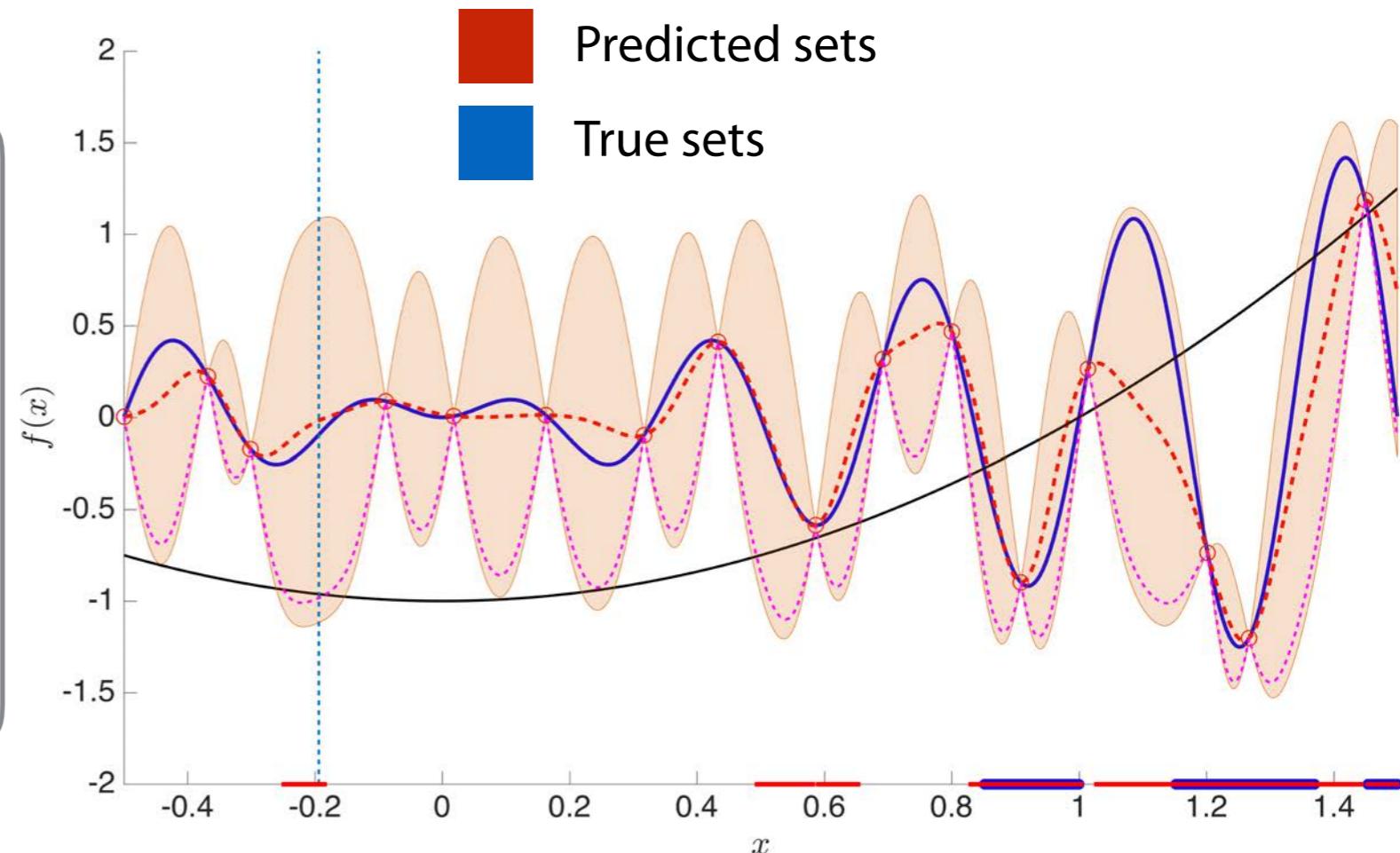
**Goal:** Identify the sets  $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t(\mathbf{x})\}$

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

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

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 sets are guaranteed to be a subset of the true sets.

# Learning level sets

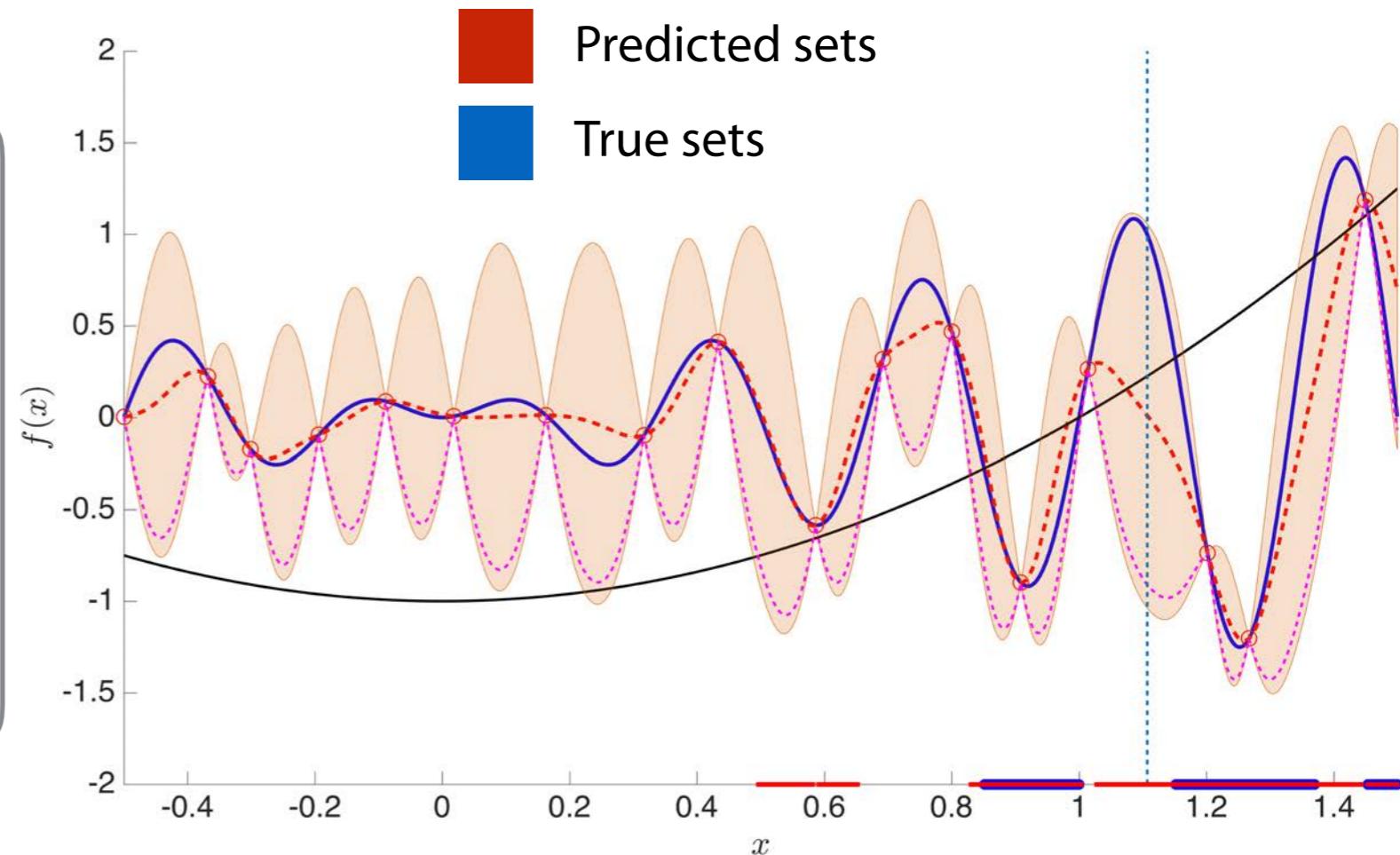
**Goal:** Identify the sets  $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t(\mathbf{x})\}$

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

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

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 sets are guaranteed to be a subset of the true sets.

# Learning level sets

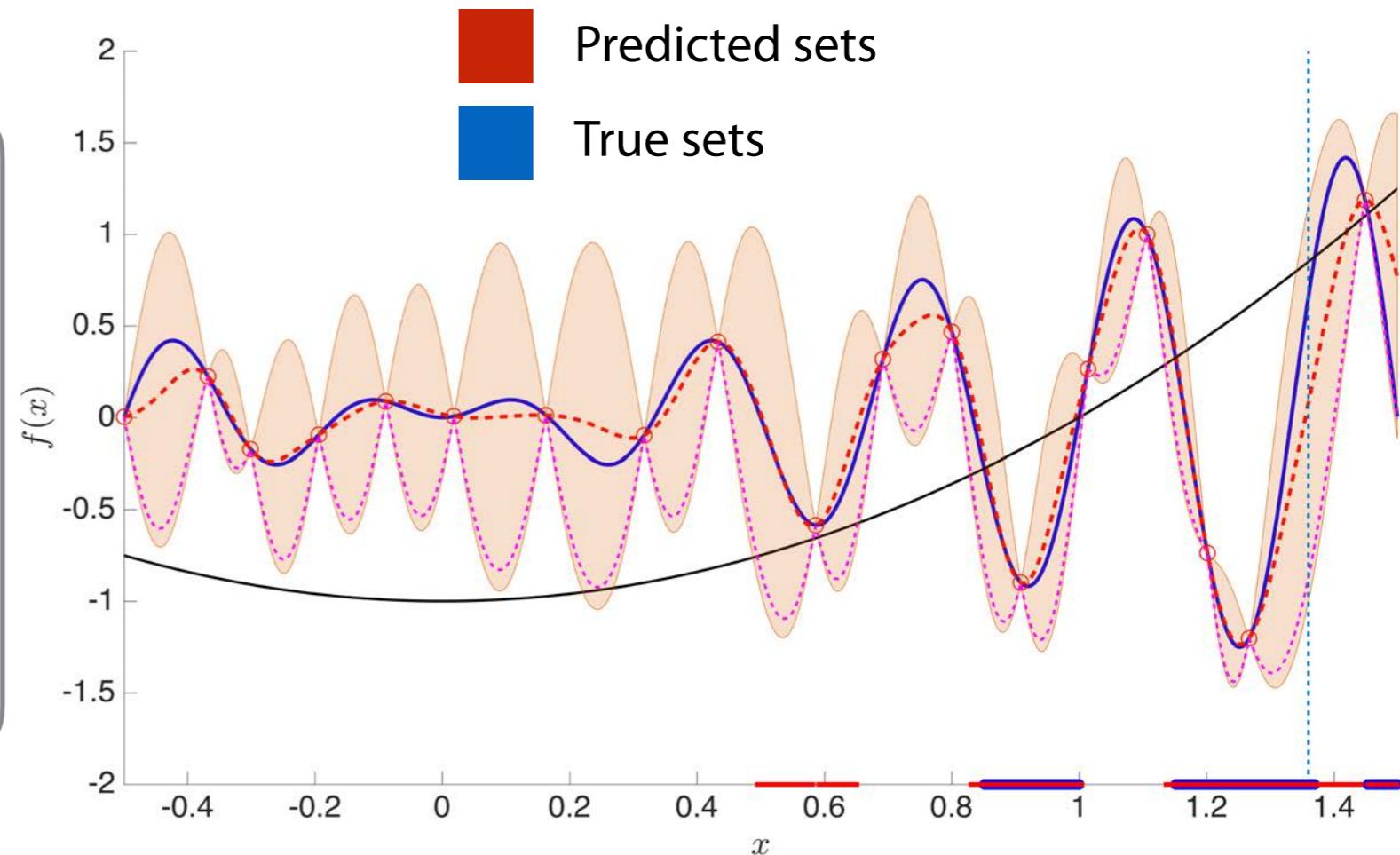
**Goal:** Identify the sets  $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t(\mathbf{x})\}$

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

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

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 sets are guaranteed to be a subset of the true sets.

# Learning level sets

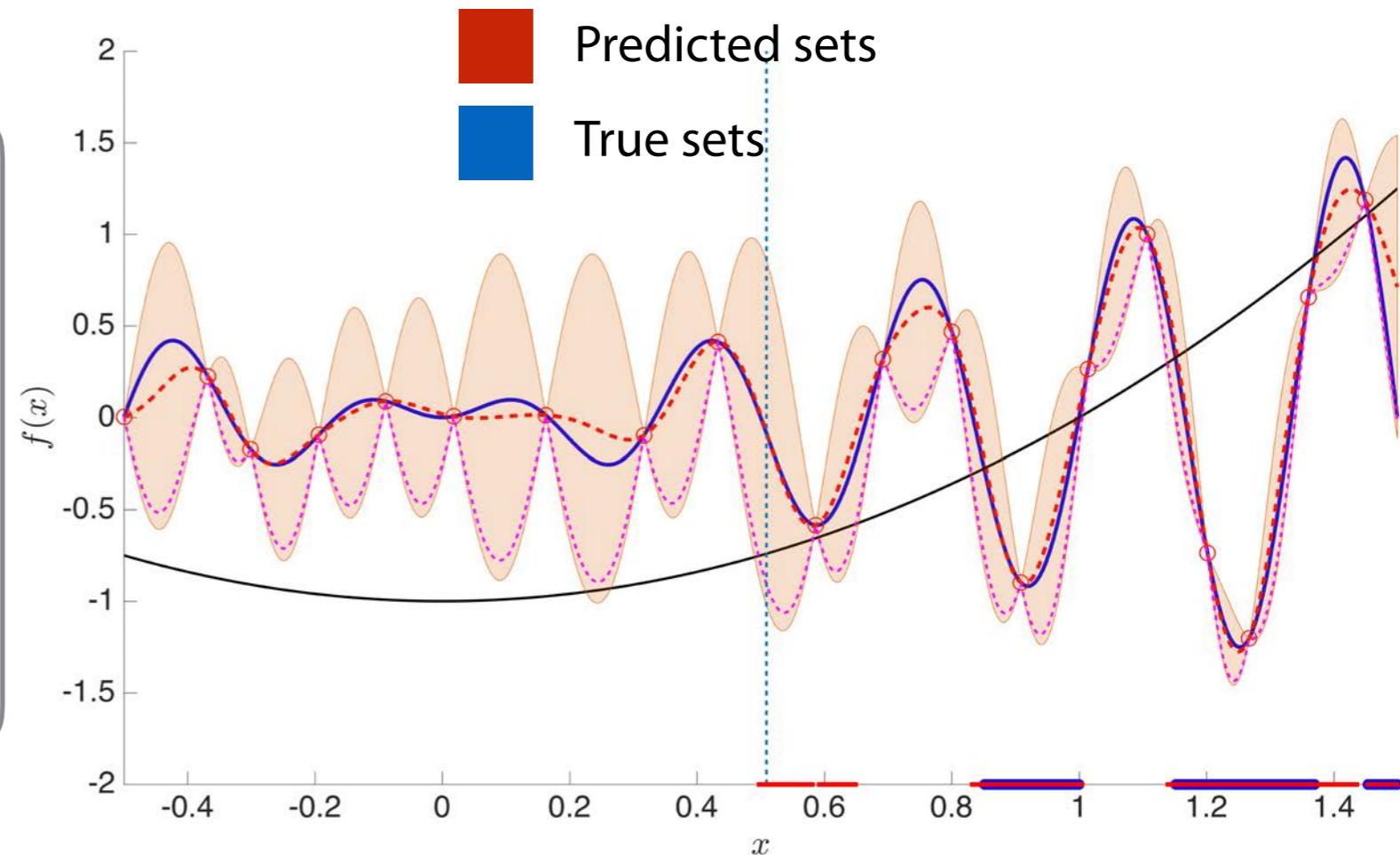
**Goal:** Identify the sets  $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t(\mathbf{x})\}$

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

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

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 sets are guaranteed to be a subset of the true sets.

# Learning level sets

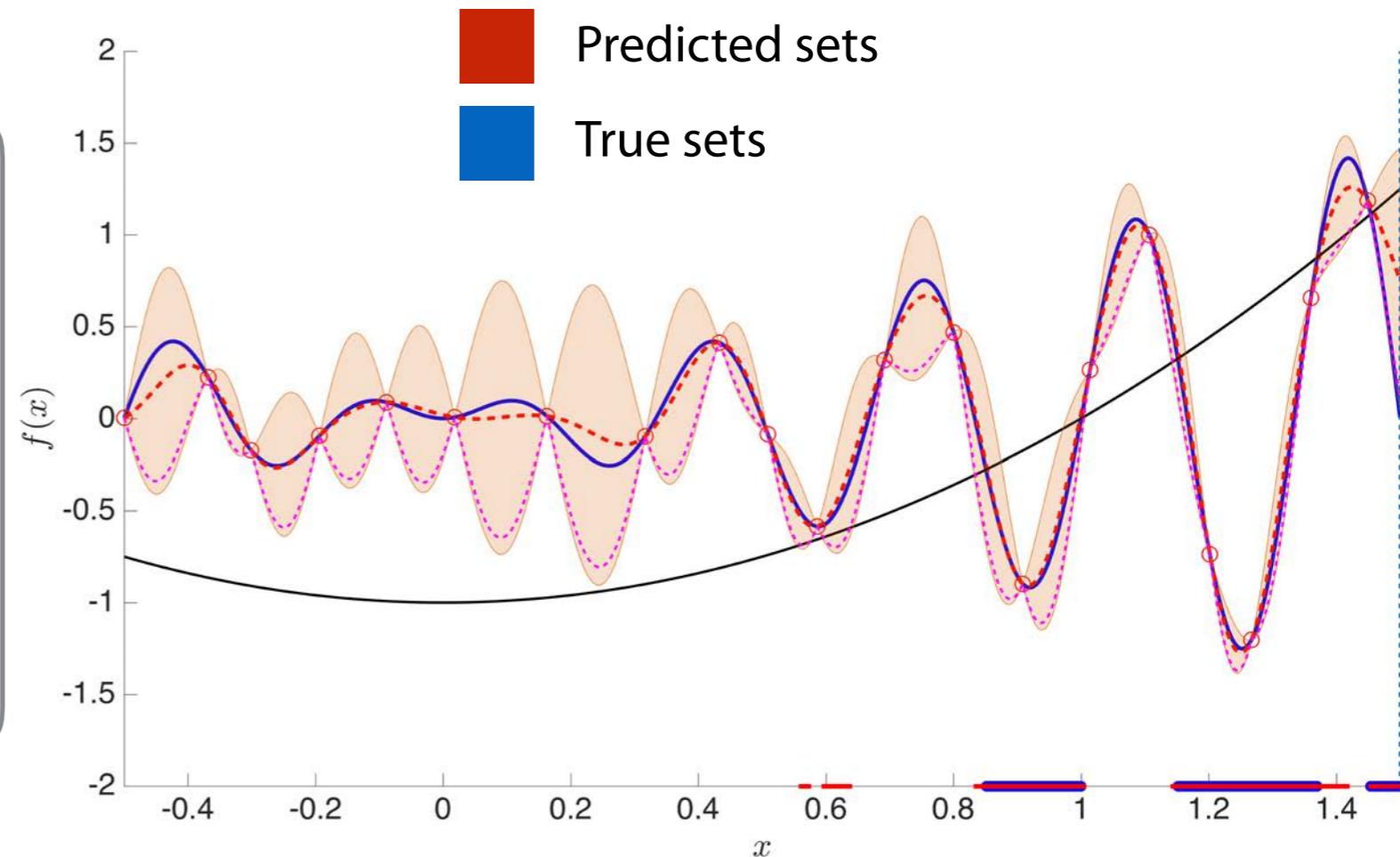
**Goal:** Identify the sets  $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t(\mathbf{x})\}$

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

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

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 sets are guaranteed to be a subset of the true sets.

# Learning level sets

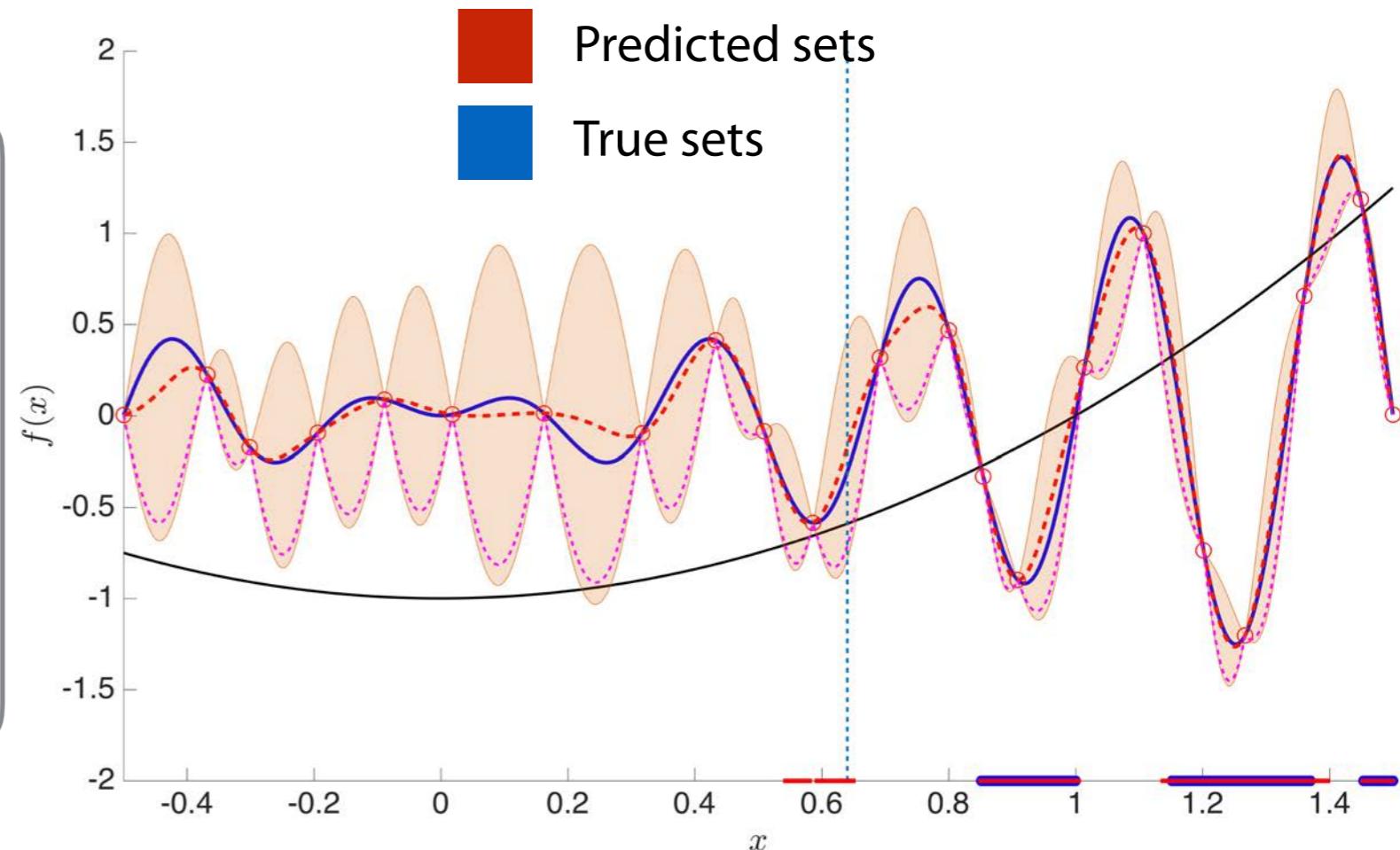
**Goal:** Identify the sets  $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t(\mathbf{x})\}$

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

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

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 sets are guaranteed to be a subset of the true sets.

# Learning level sets

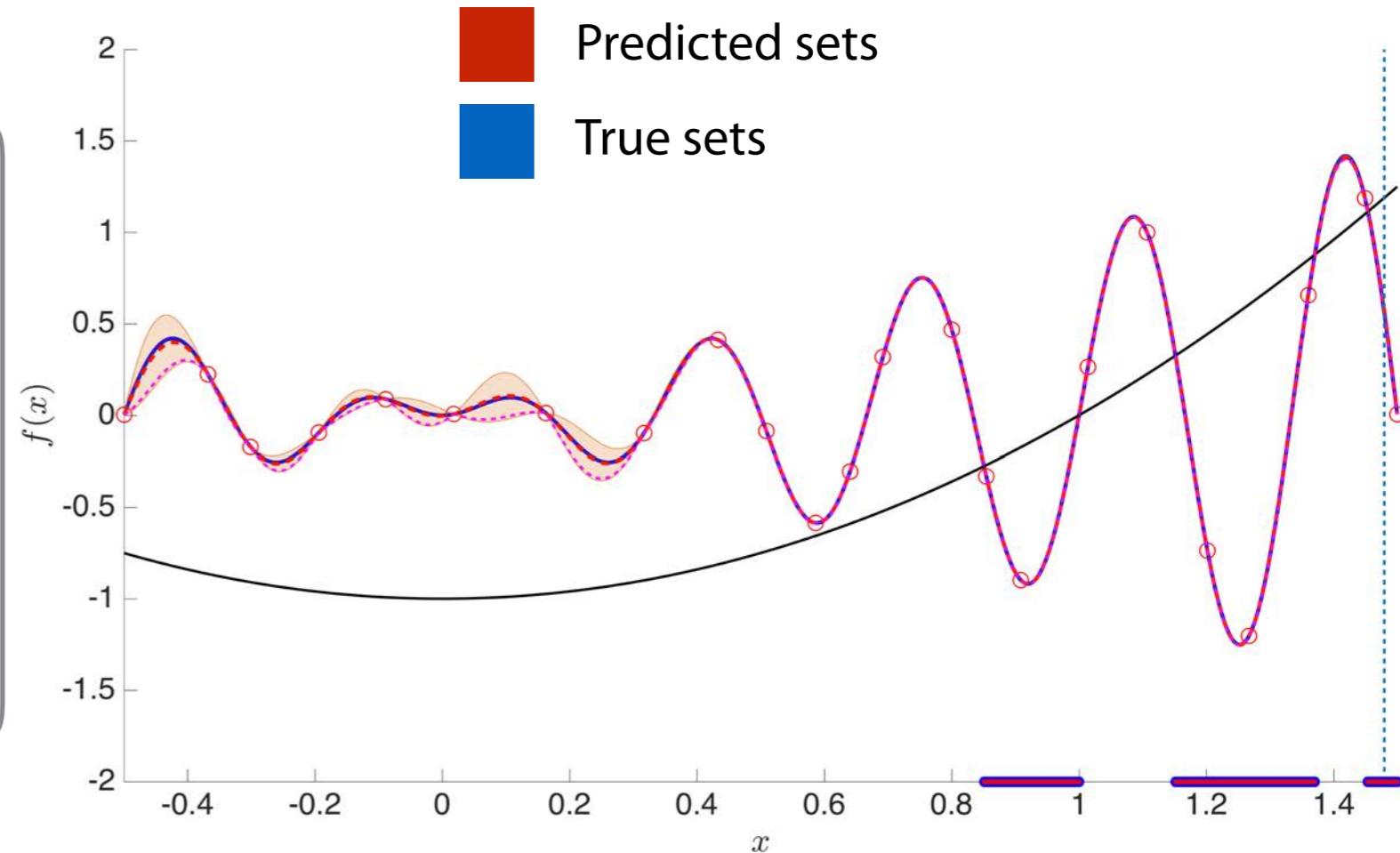
**Goal:** Identify the sets  $L_\alpha(t) = \{\mathbf{x} : \mathcal{R}_\alpha(f(\mathbf{x})) \leq t(\mathbf{x})\}$

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

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

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 sets are guaranteed to be a subset of the true sets.