

Bayesian Optimization

CSC 412/2506 Tutorial

Geoffrey Roeder
Mar 31, 2017

Slides from Kevin Swersky, Nando de Freitas

Course of tutorial

- What problem are we solving with BayesOpt?
- Gaussian Process review: notebook
- Acquisition functions
- EI Example

What problem are we solving?

Machine Learning Basics

- Given: input/output pairs $\{(x_n, y_n)\}_{n=1}^N$
- Goal: find a function that maps inputs to outputs

$$f(x) = \hat{y}$$

- Want to predict correct outputs for *unseen* inputs

Example: Linear Regression

- Problem: fit a curve to one-dimensional data
- Model: linear weights with polynomial basis

$$f(x; a) = a_0 + \sum_{i=1}^p a_i x_i^p$$

- Training algorithm: least squares

$$\hat{a} = \operatorname{argmin}_a -\frac{1}{2} \sum_{n=1}^N (y_n - f(x_n; a))^2$$

- Validation: 5-fold cross-validation

Meta-Parameters

- The model parameters are the regression coefficients.
- We train these with least squares.
- Are there any other parameters?
 - Yes!
 - Polynomial degree
 - Choice of basis (Polynomial, Fourier, Wavelet, ...)
 - Learning algorithm (Least squares, gradient descent, ...)
 - Regularization strength
 - Regularizer (L1, L2, ...)

Meta-Parameters

- Modeling decisions or free variables that cannot be trained using gradient (or other principled) methods.
- Can only evaluate a setting of the meta-parameters by training the model.
- This is very expensive, we want to do this as little as possible.

Typical Search Strategies

- Expert Intuition
- Grid Search
- Random Search
- Grad Student Search

Optimization Framework

- Meta-parameter search is an optimization problem!
 - There is some latent, potentially noisy function that maps meta-parameter settings to a score.
 - The input domain is bounded to some reasonable range.
 - Find the setting that minimizes the score.
- Each function evaluation is expensive, so we need to be clever about how often we query it.

Uncertainty

- In order to perform global optimization we need to characterize our *uncertainty*.
 - Explore places we are unsure about.
 - Exploit when we are sure we can improve.
- Two major sources of uncertainty:
 - Process noise - the observations are not perfectly accurate.
 - Model uncertainty - the response surface is one of many sensible possibilities.

Bayesian Optimization

- Mockus, 1978
 1. Incorporate a prior over the space of possible objective functions.
 2. Combine the prior and likelihood (model fit to data) to get a posterior over function values given observations.
 3. Select the next input to evaluate based on the posterior.
- According to what strategy?

Gaussian Processes

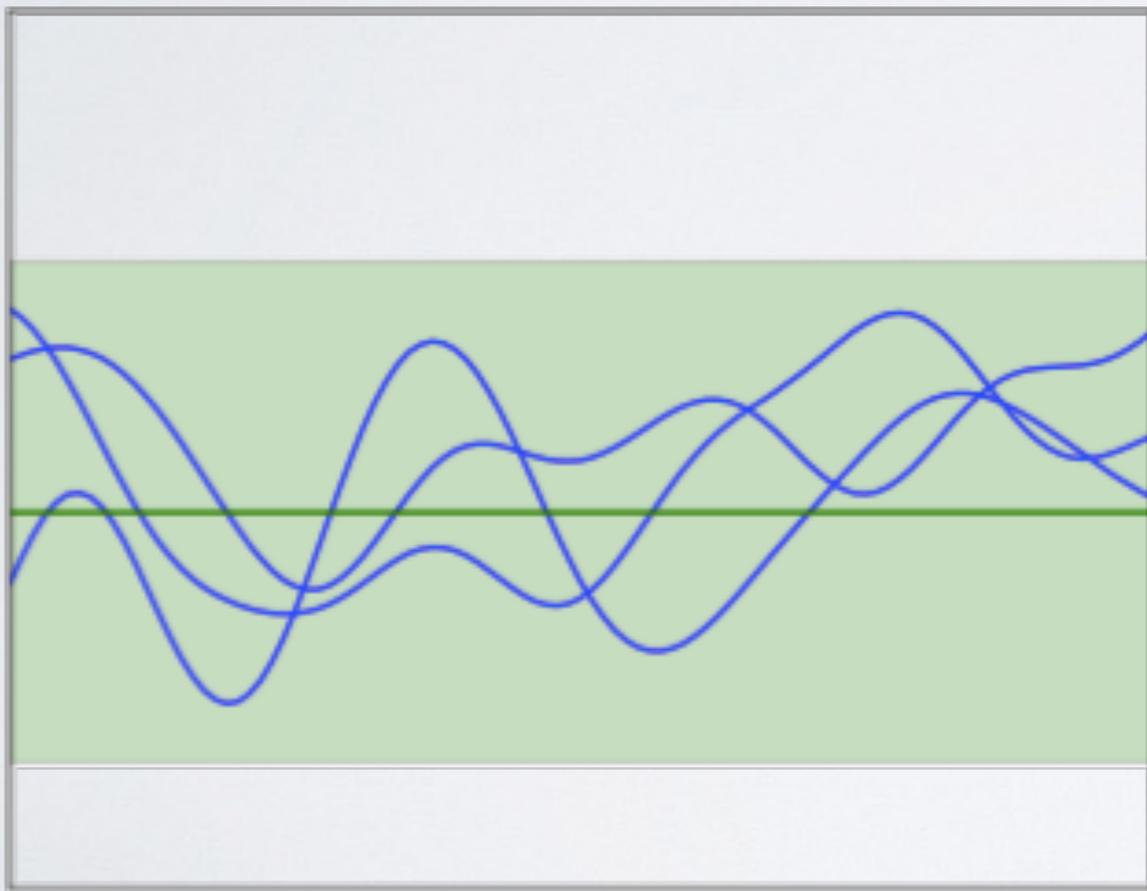
- Distribution over functions $f : \mathcal{X} \rightarrow \mathbb{R}$
- The observations at points $\mathbf{X} = \{x_n \in \mathcal{X}\}_{n=1}^N$ are jointly Gaussian
- Specified by a mean $m : \mathcal{X} \rightarrow \mathbb{R}$ and covariance $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$
- Predictive mean and covariance given observations $(\mathbf{X}, \mathbf{y}) = \{x_n, y_n\}_{n=1}^N$

$$\mu(x; \mathbf{X}, \mathbf{y}, \theta) = K(\mathbf{X}, x)^\top K(\mathbf{X}, \mathbf{X})^{-1}(\mathbf{y} - m(\mathbf{X}))$$
$$\Sigma(x, x'; \mathbf{X}, \mathbf{y}, \theta) = K(x, x') - K(\mathbf{X}, x)^\top K(\mathbf{X}, \mathbf{X})^{-1}K(\mathbf{X}, x')$$

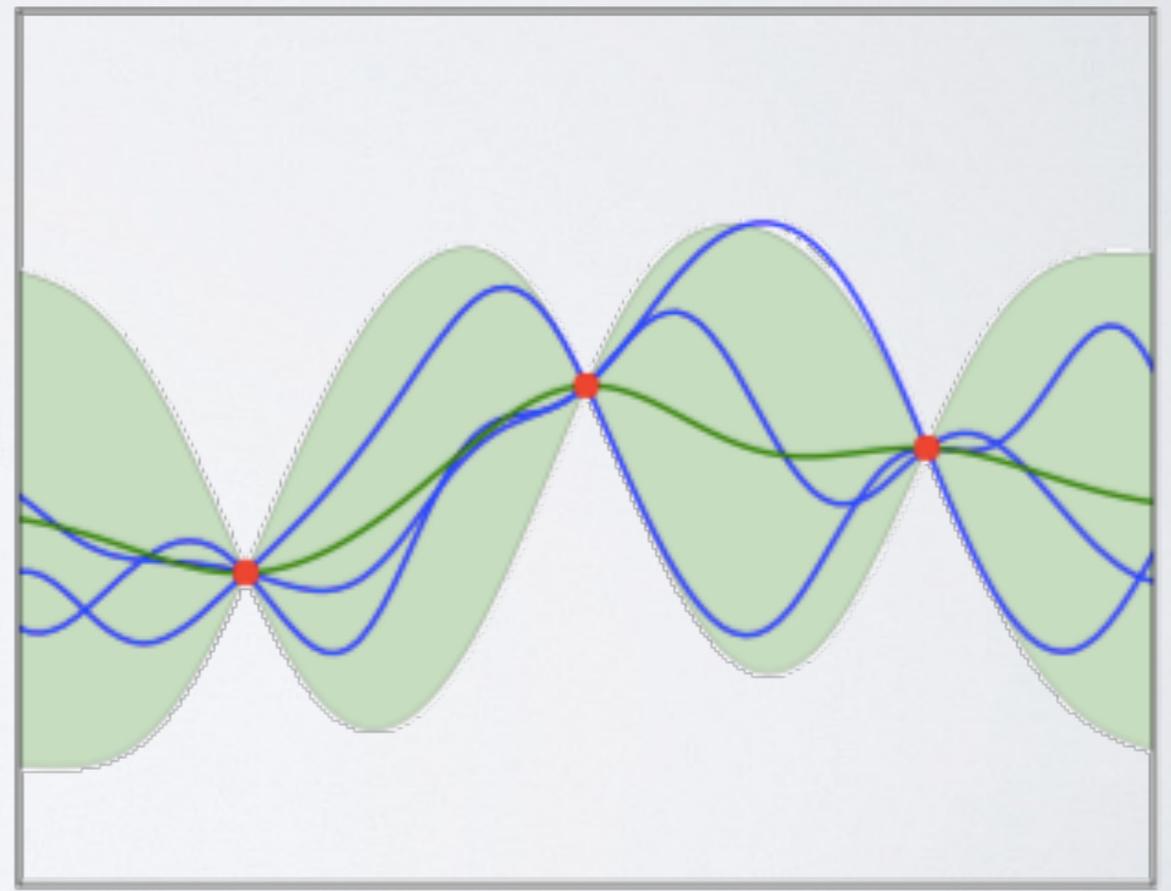
- Intuition:
 - A prior for smooth functions
 - Similar inputs (high covariance) have similar outputs
- Can compute expected value and uncertainty for test inputs easily

GPs as Distributions over Functions

Prior



Posterior

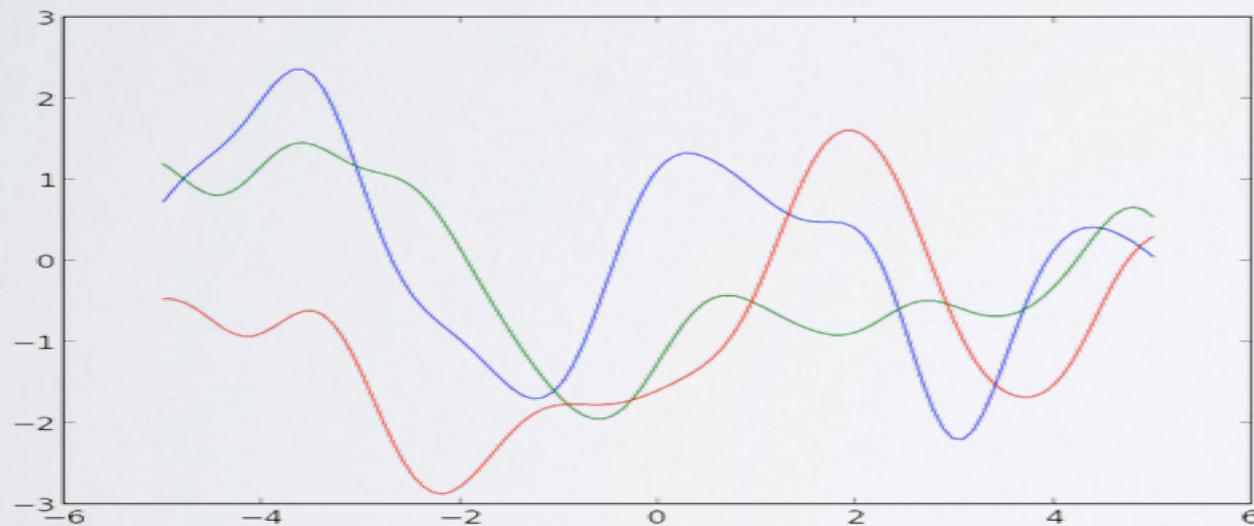


*Samples in blue

GPs Allow High Level Specification

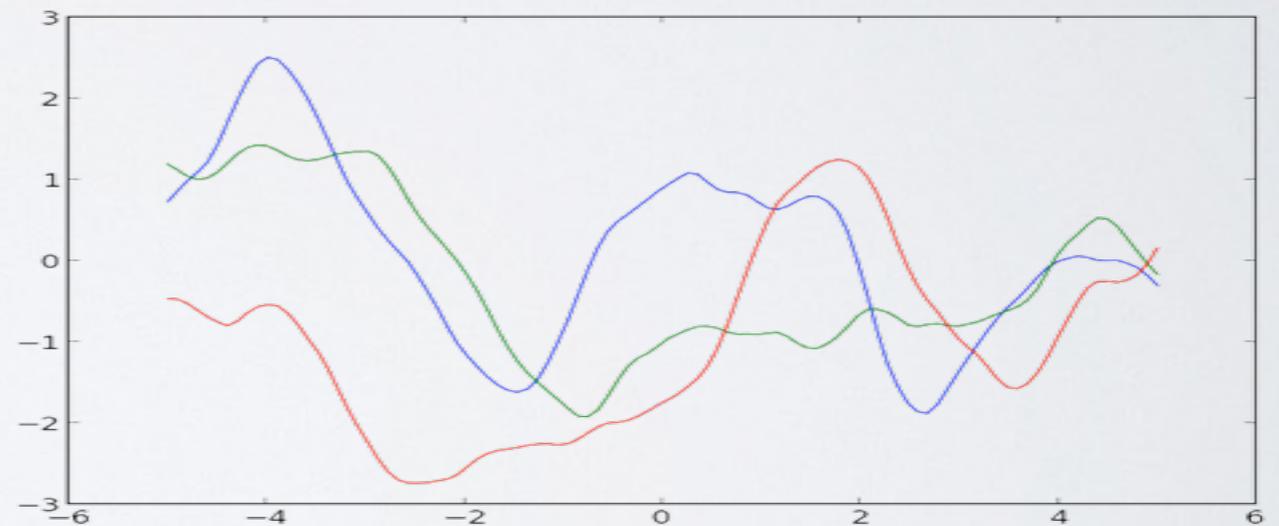
- Gaussian processes are nonparametric, their behaviour is specified at a high level by the choice of *kernel*.

Infinitely differentiable



$$K_{SE}(x, x') = \theta_0 \exp(-r^2(x, x'))$$

Twice differentiable



$$K_{M52}(x, x') = \theta_0 \left(1 + \sqrt{5r^2(x, x')} + \frac{5}{3}r^2(x, x') \right) \exp \left(-\sqrt{5r^2(x, x')} \right)$$

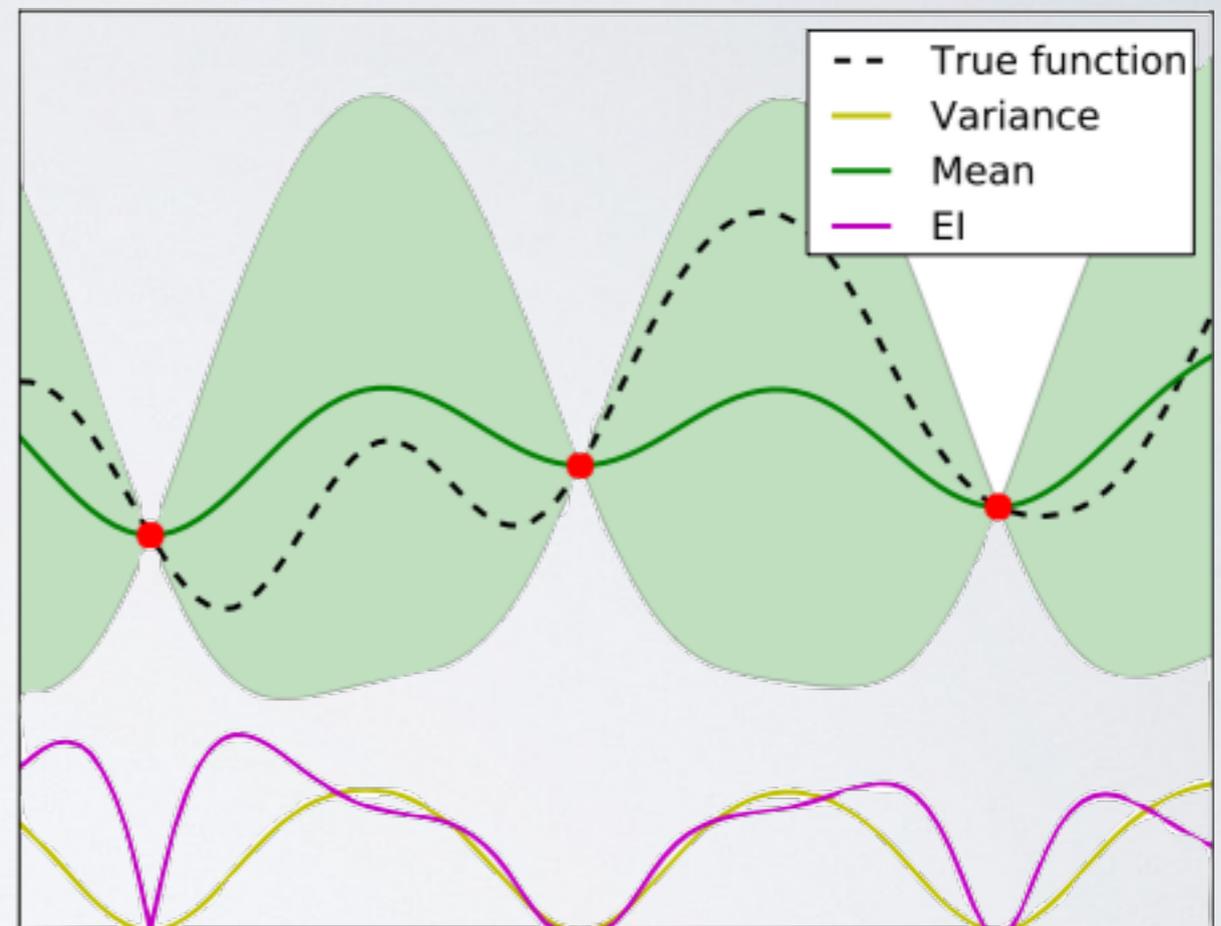
$$r^2(x, x') = \frac{(x-x')^2}{\theta_1^2}$$

Gaussian Process Notebook

Acquisition Functions

Choosing where to Search Next

- The GP gives us a mean and variance for each input.
- Minimum expected value: purely exploitative
- Maximum uncertainty: purely explorative
- Expected improvement: trade-off (Mockus, 1978)
- Many other *acquisition functions* have been proposed in the literature



$$EI(x; \mathbf{X}, \mathbf{y}, \theta) = \int_{y_{\text{best}}}^{\infty} \max(0, y - y_{\text{best}}) P(y|x; \mathbf{X}, \mathbf{y}, \theta) dy$$

Exploration-exploitation tradeoff

Recall the expressions for GP prediction:

$$P(y_{t+1} | \mathcal{D}_{1:t}, \mathbf{x}_{t+1}) = \mathcal{N}(\mu_t(\mathbf{x}_{t+1}), \sigma_t^2(\mathbf{x}_{t+1}) + \sigma_{\text{noise}}^2)$$

$$\mu_t(\mathbf{x}_{t+1}) = \mathbf{k}^T [\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I}]^{-1} \mathbf{y}_{1:t}$$

$$\sigma_t^2(\mathbf{x}_{t+1}) = k(\mathbf{x}_{t+1}, \mathbf{x}_{t+1}) - \mathbf{k}^T [\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I}]^{-1} \mathbf{k}$$

Exploration-exploitation tradeoff

Recall the expressions for GP prediction:

$$P(y_{t+1} | \mathcal{D}_{1:t}, \mathbf{x}_{t+1}) = \mathcal{N}(\mu_t(\mathbf{x}_{t+1}), \sigma_t^2(\mathbf{x}_{t+1}) + \sigma_{\text{noise}}^2)$$

$$\mu_t(\mathbf{x}_{t+1}) = \mathbf{k}^T [\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I}]^{-1} \mathbf{y}_{1:t}$$

$$\sigma_t^2(\mathbf{x}_{t+1}) = k(\mathbf{x}_{t+1}, \mathbf{x}_{t+1}) - \mathbf{k}^T [\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I}]^{-1} \mathbf{k}$$

We should choose the next point \mathbf{x} where the mean is high (**exploitation**) and the variance is high (**exploration**).

Exploration-exploitation tradeoff

Recall the expressions for GP prediction:

$$P(y_{t+1} | \mathcal{D}_{1:t}, \mathbf{x}_{t+1}) = \mathcal{N}(\mu_t(\mathbf{x}_{t+1}), \sigma_t^2(\mathbf{x}_{t+1}) + \sigma_{\text{noise}}^2)$$
$$\mu_t(\mathbf{x}_{t+1}) = \mathbf{k}^T [\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I}]^{-1} \mathbf{y}_{1:t}$$
$$\sigma_t^2(\mathbf{x}_{t+1}) = k(\mathbf{x}_{t+1}, \mathbf{x}_{t+1}) - \mathbf{k}^T [\mathbf{K} + \sigma_{\text{noise}}^2 \mathbf{I}]^{-1} \mathbf{k}$$

We should choose the next point \mathbf{x} where the mean is high (**exploitation**) and the variance is high (**exploration**).

We could balance this tradeoff with an acquisition function as follows:

$$\mu(\mathbf{x}) + \kappa \sigma(\mathbf{x})$$

Probability of Improvement

An acquisition function: Probability of Improvement

$$\text{PI}(\mathbf{x}) = P(f(\underline{\mathbf{x}}) \geq \mu^+ + \xi)$$

μ^+ best observed value

An acquisition function: Probability of Improvement

$$\begin{aligned} \text{PI}(\mathbf{x}) &= P(f(\mathbf{x}) \geq \mu^+ + \xi) \\ &= \Phi\left(\frac{\mu(\mathbf{x}) - \mu^+ - \xi}{\sigma(\mathbf{x})}\right) \end{aligned}$$

cumulative

y^+ best observed value

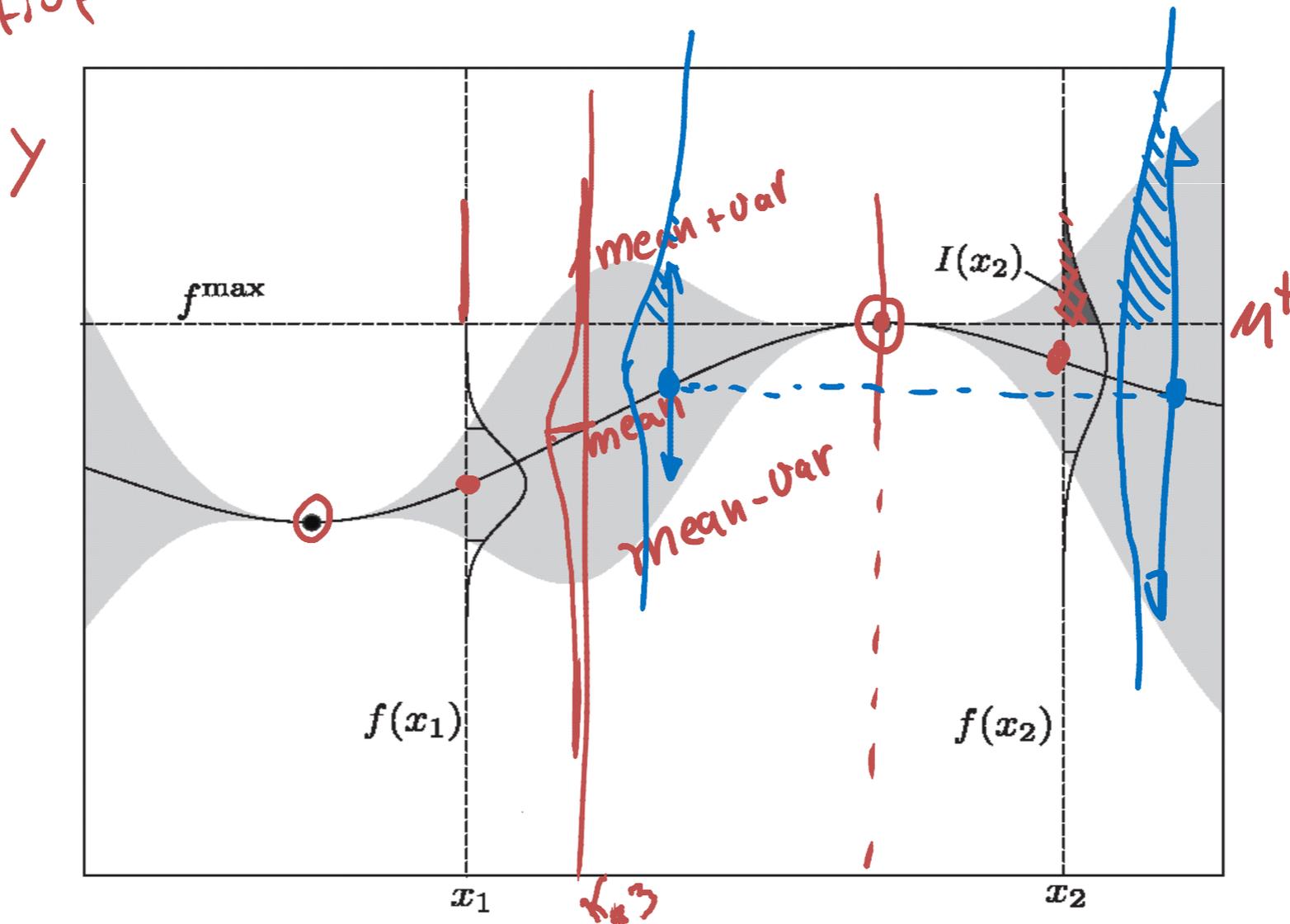
An acquisition function: Probability of Improvement

$$\begin{aligned}
 \text{PI}(\mathbf{x}) &= P(f(\mathbf{x}) \geq \mu^+ + \xi) \\
 &= \Phi\left(\frac{\mu(\mathbf{x}) - \mu^+ - \xi}{\sigma(\mathbf{x})}\right)
 \end{aligned}$$

cumulative

μ^+ best observed value

$$f^{\max} = \mu^+ + \xi \text{ (small)}$$



Expected Improvement

Bayes and decision theory

Utilitarian view: We need models to make the right decisions under uncertainty. Inference and decision making are intertwined.

Bayes and decision theory

Utilitarian view: We need models to make the right decisions under uncertainty. Inference and decision making are intertwined.

Learned posterior

$$\left\{ \begin{array}{l} P(x=\text{healthy}|\text{data}) = 0.9 \\ P(x=\text{cancer}|\text{data}) = 0.1 \end{array} \right.$$

Cost/Reward model $u(x,a)$

	$a = \text{no treatment}$	$a = \text{treatment}$
$x = \text{healthy}$	0	-30
$x = \text{cancer}$	-100	-20

Handwritten annotations in blue ink: circles around the values 0, -30, -100, and -20. A vertical line points to the 0. A horizontal line connects the 0 and -30. A horizontal line connects the -100 and -20. A diagonal line is drawn through the -20.

Bayes and decision theory

Utilitarian view: We need models to make the right decisions under uncertainty. Inference and decision making are intertwined.

Learned posterior

$$\begin{cases} P(x=\textit{healthy}|\textit{data}) = 0.9 \\ P(x=\textit{cancer}|\textit{data}) = 0.1 \end{cases}$$

Cost/Reward model $u(x,a)$

	$a = \textit{no treatment}$	$a = \textit{treatment}$
$x = \textit{healthy}$	0	-30
$x = \textit{cancer}$	-100	-20

We choose the action that maximizes the **expected utility**, or equivalently, which minimizes the **expected cost**.

$$EU(a) = \sum u(x,a) P(x|\textit{data})$$

Bayes and decision theory

Utilitarian view: We need models to make the right decisions under uncertainty. Inference and decision making are intertwined.

Learned posterior

$$\begin{cases} P(x=\text{healthy}|\text{data}) = 0.9 \\ P(x=\text{cancer}|\text{data}) = 0.1 \end{cases}$$

Cost/Reward model $u(x,a)$

	$a = \text{no treatment}$	$a = \text{treatment}$
$x = \text{healthy}$	0	-30
$x = \text{cancer}$	-100	-20

We choose the action that maximizes the **expected utility**, or equivalently, which minimizes the **expected cost**.

$$EU(a) = \sum_x u(x,a) P(x|\text{data})$$

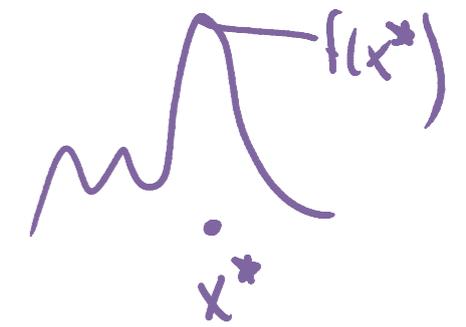
$$EU(a=\text{treatment}) = u(\text{healthy, treatment}) P(x=\text{healthy}|\text{data}) + u(\text{cancer, treatment}) P(x=\text{cancer}|\text{data})$$

$$= (-30)(0.9) + (-20)(0.1) =$$

$$EU(a=\text{no treatment}) =$$

An expected utility criterion

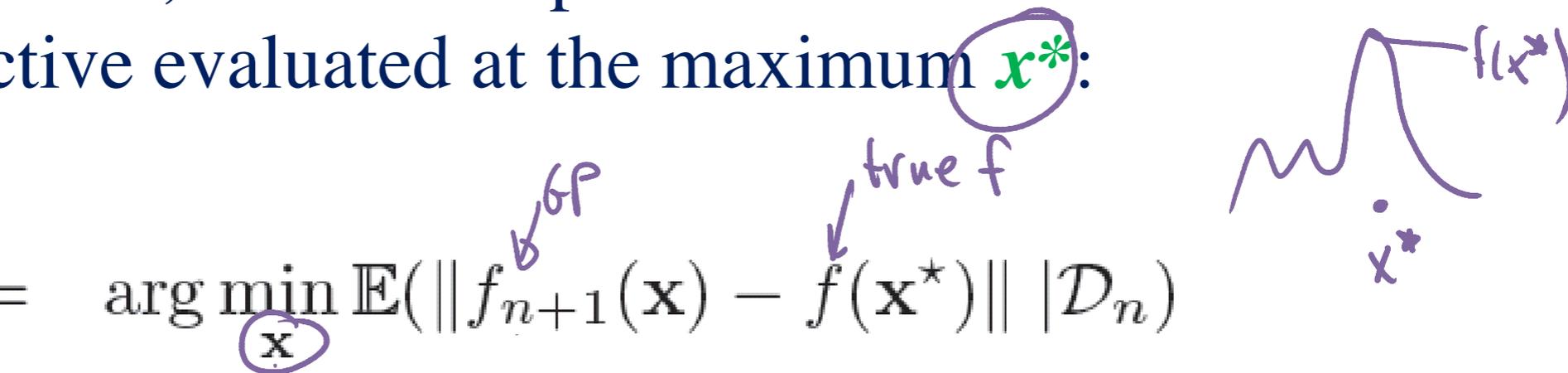
At iteration $n+1$, choose the point that minimizes the distance to the objective evaluated at the maximum \mathbf{x}^* :

$$\begin{aligned} \mathbf{x}_{n+1} &= \arg \min_{\mathbf{x}} \mathbb{E}(\|f_{n+1}(\mathbf{x}) - f(\mathbf{x}^*)\| | \mathcal{D}_n) \\ &= \arg \min_{\mathbf{x}} \int \|f_{n+1}(\mathbf{x}) - f(\mathbf{x}^*)\| p(f_{n+1} | \mathcal{D}_n) df_{n+1} \end{aligned}$$


The diagram shows a hand-drawn purple curve representing a function $f(x)$. The curve has a jagged, oscillatory appearance on the left side and then rises to a peak at a point labeled x^* . The value of the function at this peak is labeled $f(x^*)$. A purple dot marks the location of x^* on the horizontal axis.

An expected utility criterion

At iteration $n+1$, choose the point that minimizes the distance to the objective evaluated at the maximum \mathbf{x}^* :

$$\begin{aligned}\mathbf{x}_{n+1} &= \arg \min_{\mathbf{x}} \mathbb{E}(\|f_{n+1}(\mathbf{x}) - f(\mathbf{x}^*)\| | \mathcal{D}_n) \\ &= \arg \min_{\mathbf{x}} \int \|f_{n+1}(\mathbf{x}) - f(\mathbf{x}^*)\| p(f_{n+1} | \mathcal{D}_n) df_{n+1}\end{aligned}$$


We don't know the true objective at the maximum. To overcome this, Mockus proposed the following acquisition function:

$$\underline{\mathbf{x}} = \arg \max_{\mathbf{x}} \mathbb{E}(\max\{0, f_{n+1}(\mathbf{x}) - f^{\max}\} | \mathcal{D}_n)$$

Expected improvement

$$\mathbf{x} = \arg \max_{\mathbf{x}} \mathbb{E}(\max\{0, f_{n+1}(\mathbf{x}) - \underbrace{f^{\max}}_{\mu^+ + \xi}\} | \mathcal{D}_n)$$

For this acquisition, we can obtain an analytical expression:

$$\text{EI}(\mathbf{x}) = \begin{cases} (\mu(\mathbf{x}) - \mu^+ - \xi)\Phi(Z) + \sigma(\mathbf{x})\phi(Z) & \text{if } \sigma(\mathbf{x}) > 0 \\ 0 & \text{if } \sigma(\mathbf{x}) = 0 \end{cases}$$

$$Z = \frac{\mu(\mathbf{x}) - \mu^+ - \xi}{\sigma(\mathbf{x})}$$

where $\phi(\cdot)$ and $\Phi(\cdot)$ denote the PDF and CDF of the standard Normal

$$\mu^+ = \operatorname{argmax}_{\mathbf{x}_i \in \mathbf{x}_{1:t}} \mu(\mathbf{x}_i)$$

- Probability of Improvement

$$\text{PI}(\mathbf{x}) = \Phi\left(\frac{\mu(\mathbf{x}) - \mu^+ - \xi}{\sigma(\mathbf{x})}\right)$$

Kushner 1964

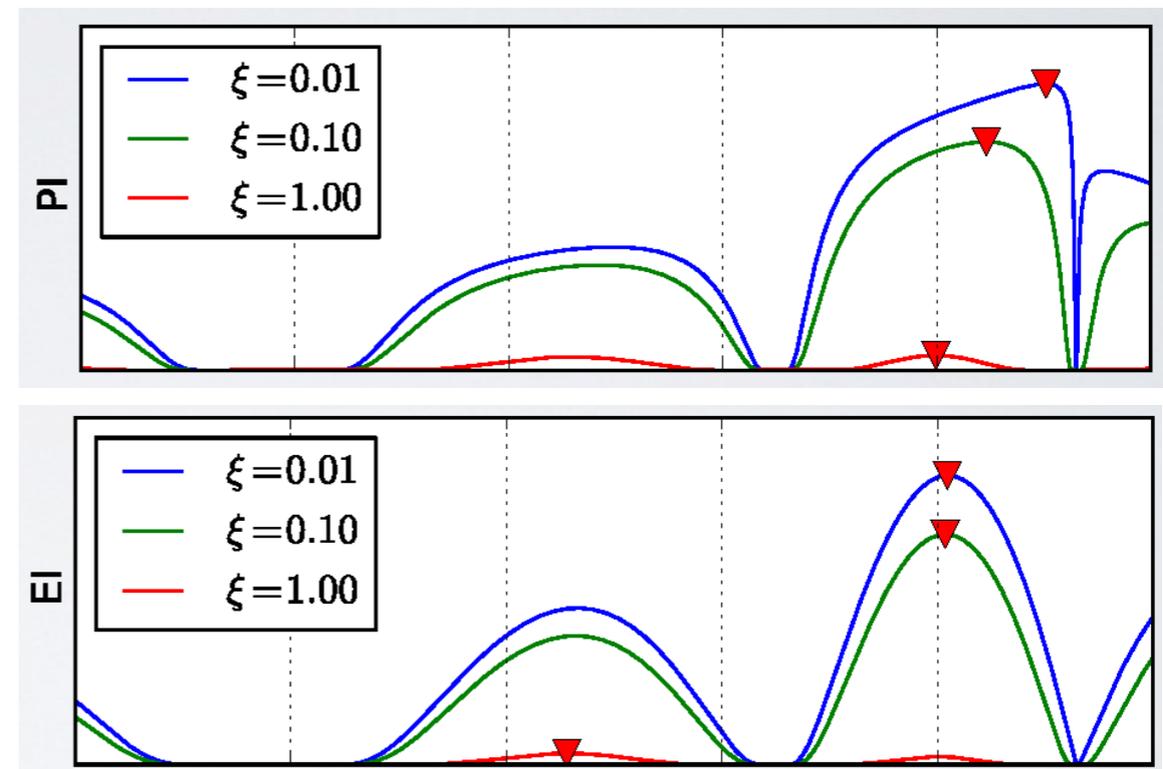
- Expected Improvement

$$\text{EI}(\mathbf{x}) = (\mu(\mathbf{x}) - \mu^+ - \xi)\Phi(Z) + \sigma(\mathbf{x})\phi(Z)$$

$$Z = \frac{\mu(\mathbf{x}) - \mu^+ - \xi}{\sigma(\mathbf{x})}$$

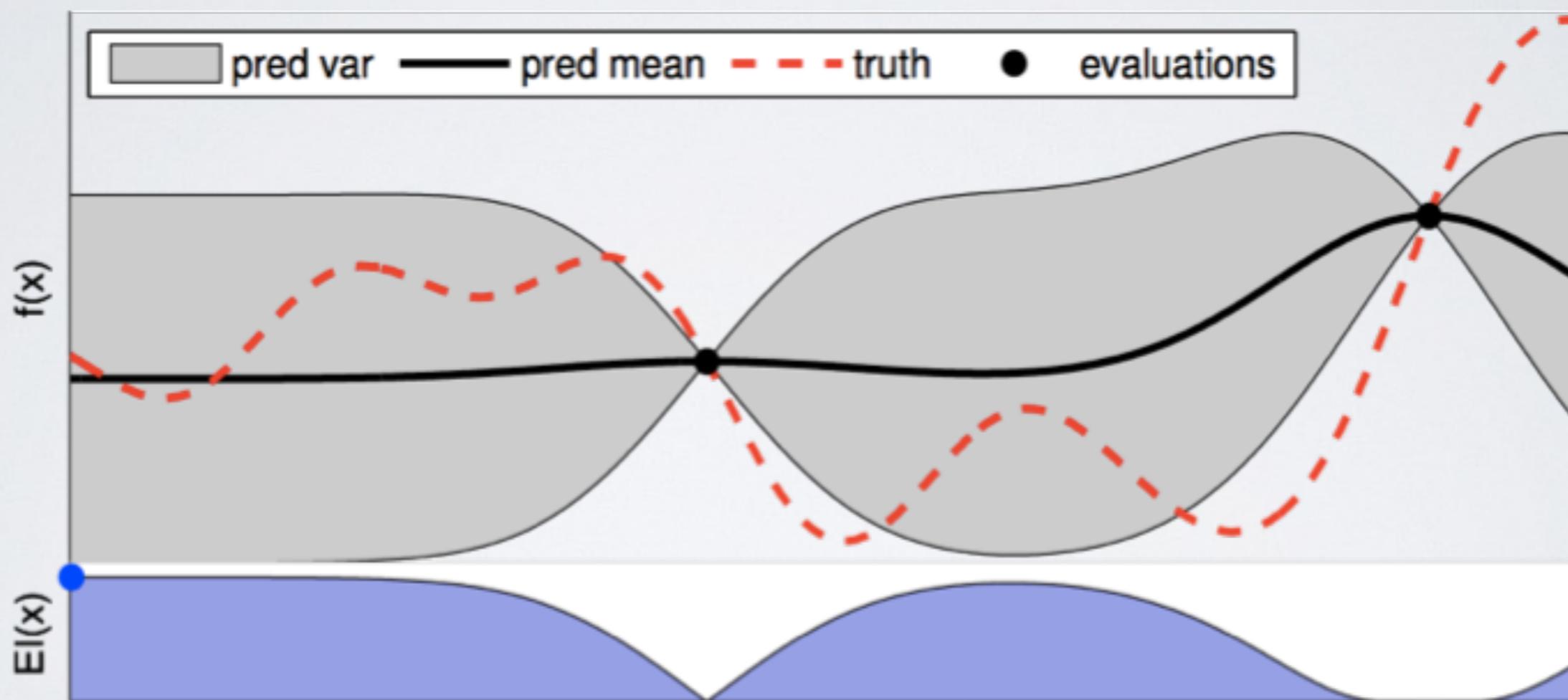
Mockus 1978

Acquisition functions

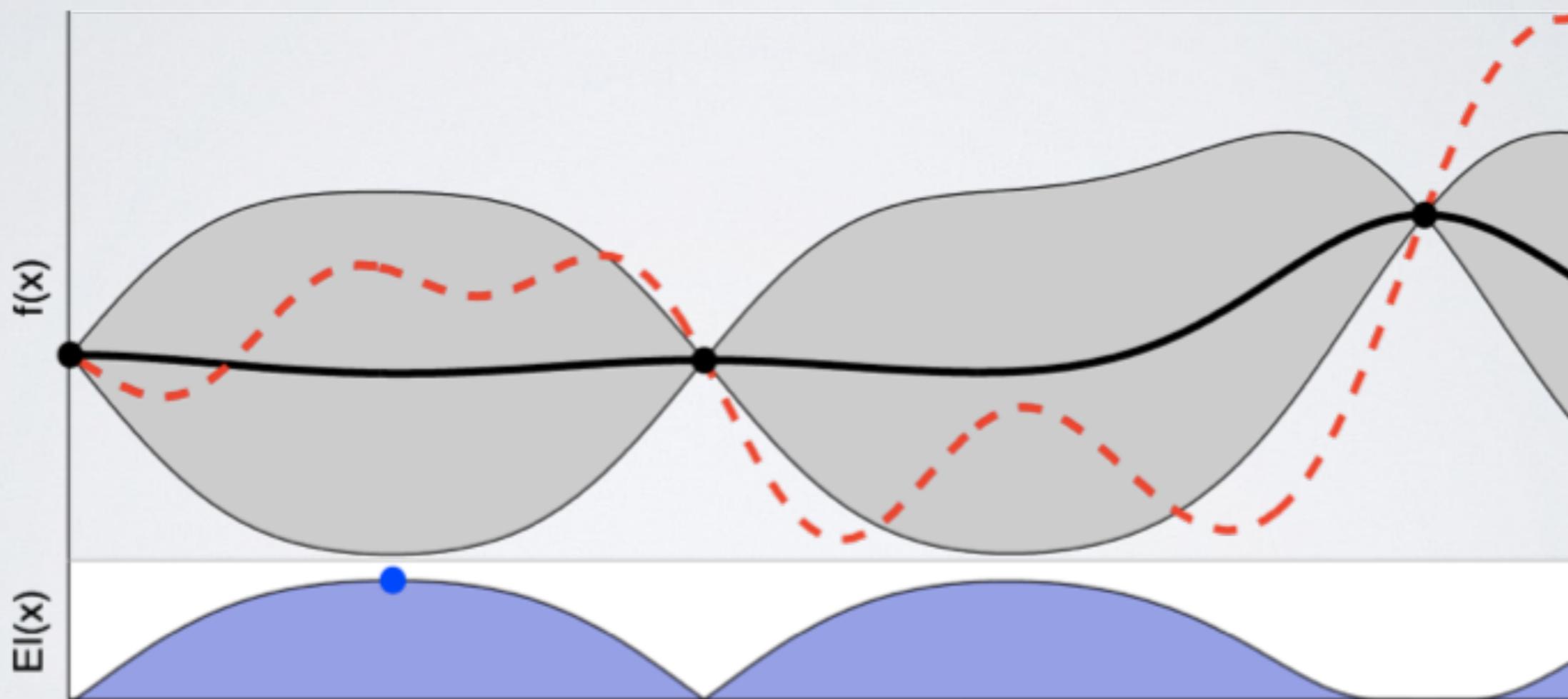


Example: Expected Improvement

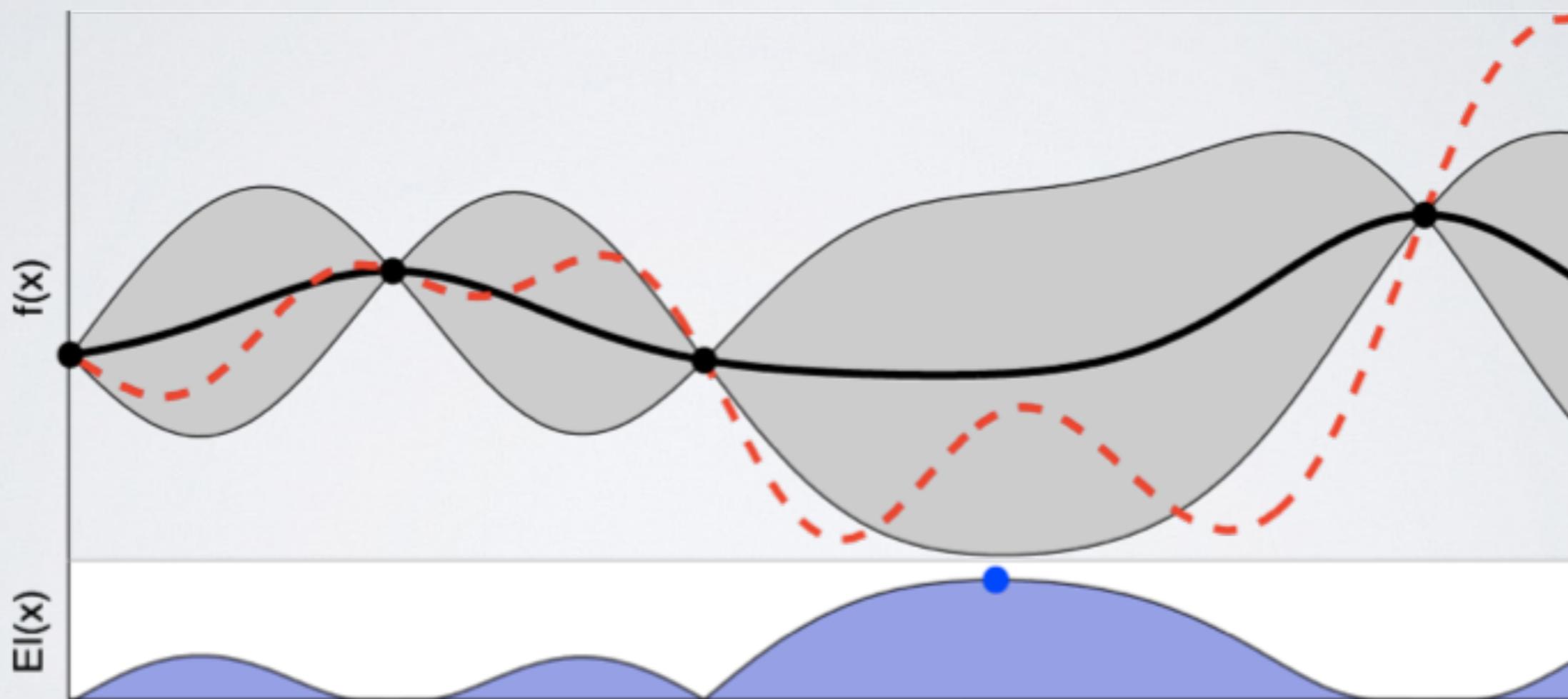
Simple Example



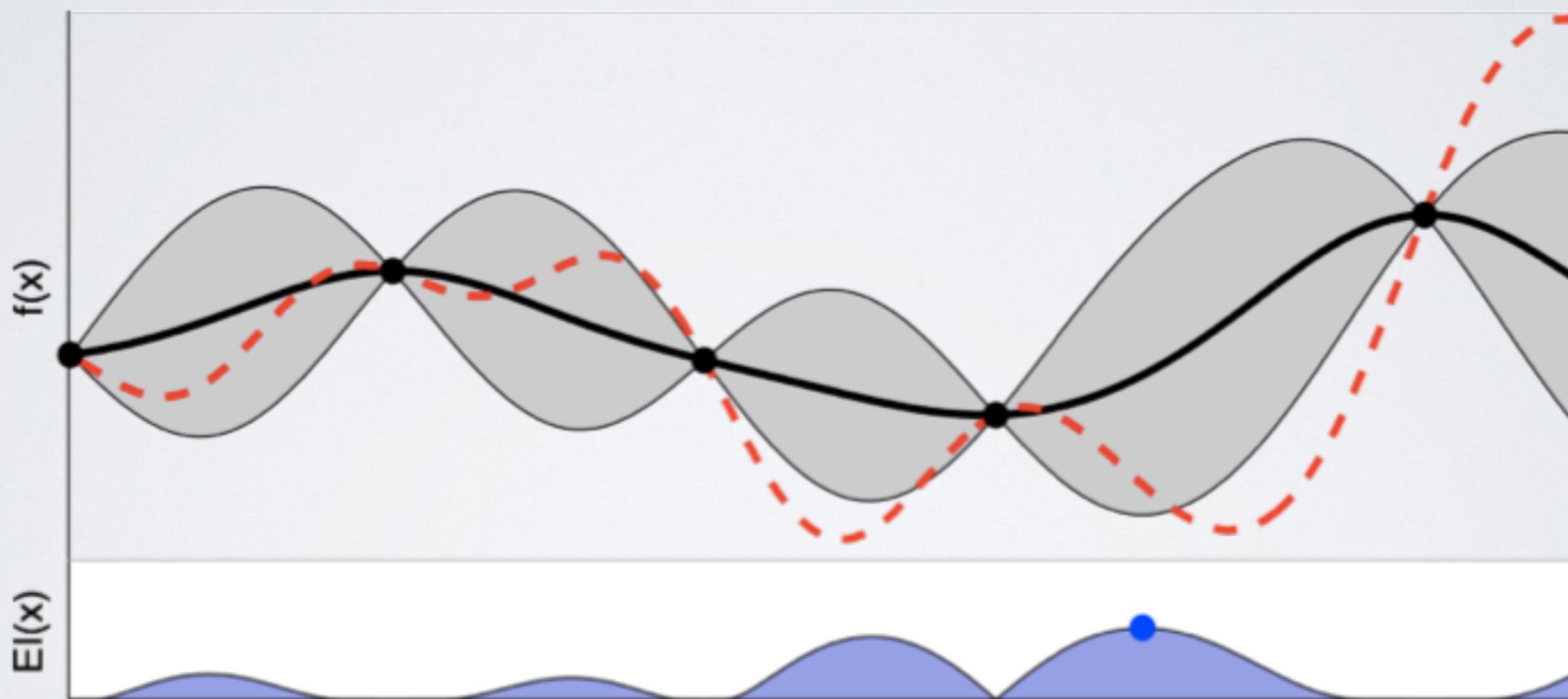
Simple Example



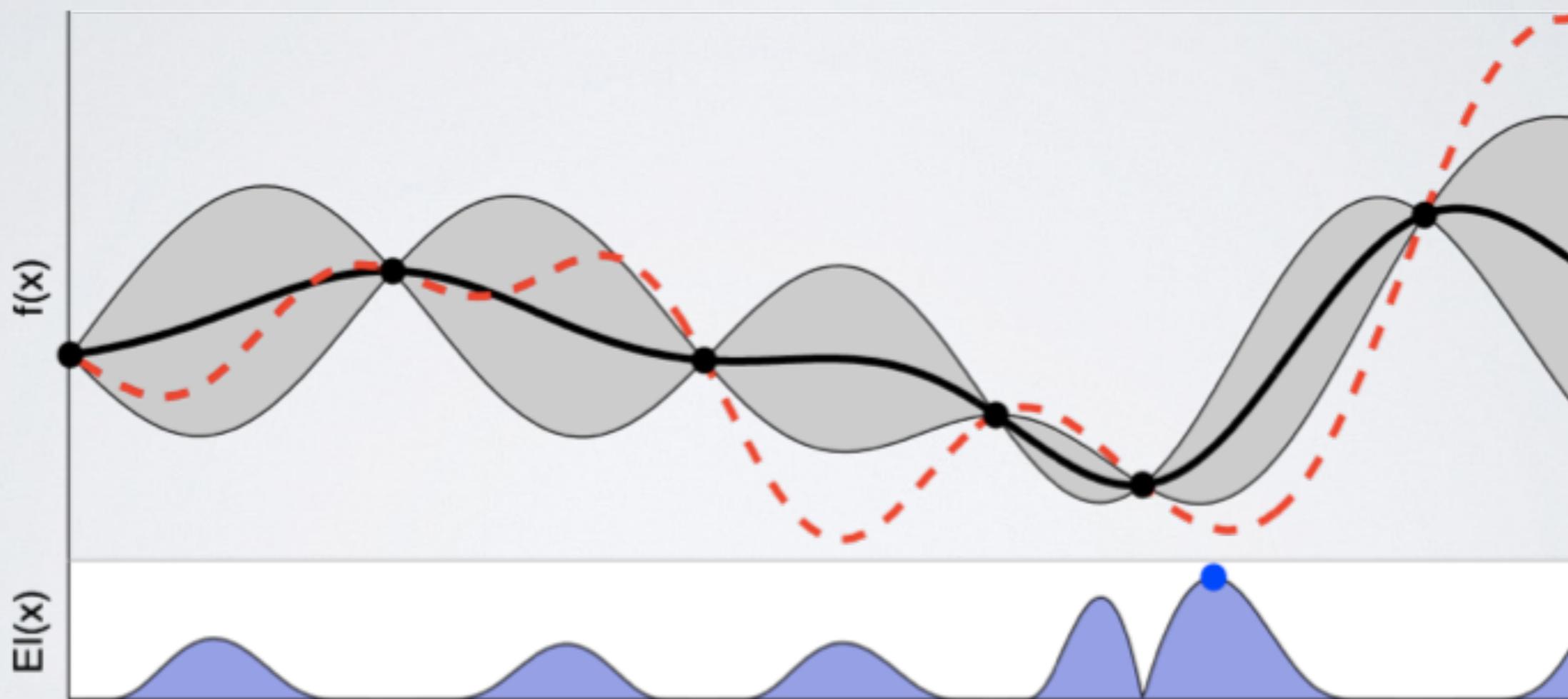
Simple Example



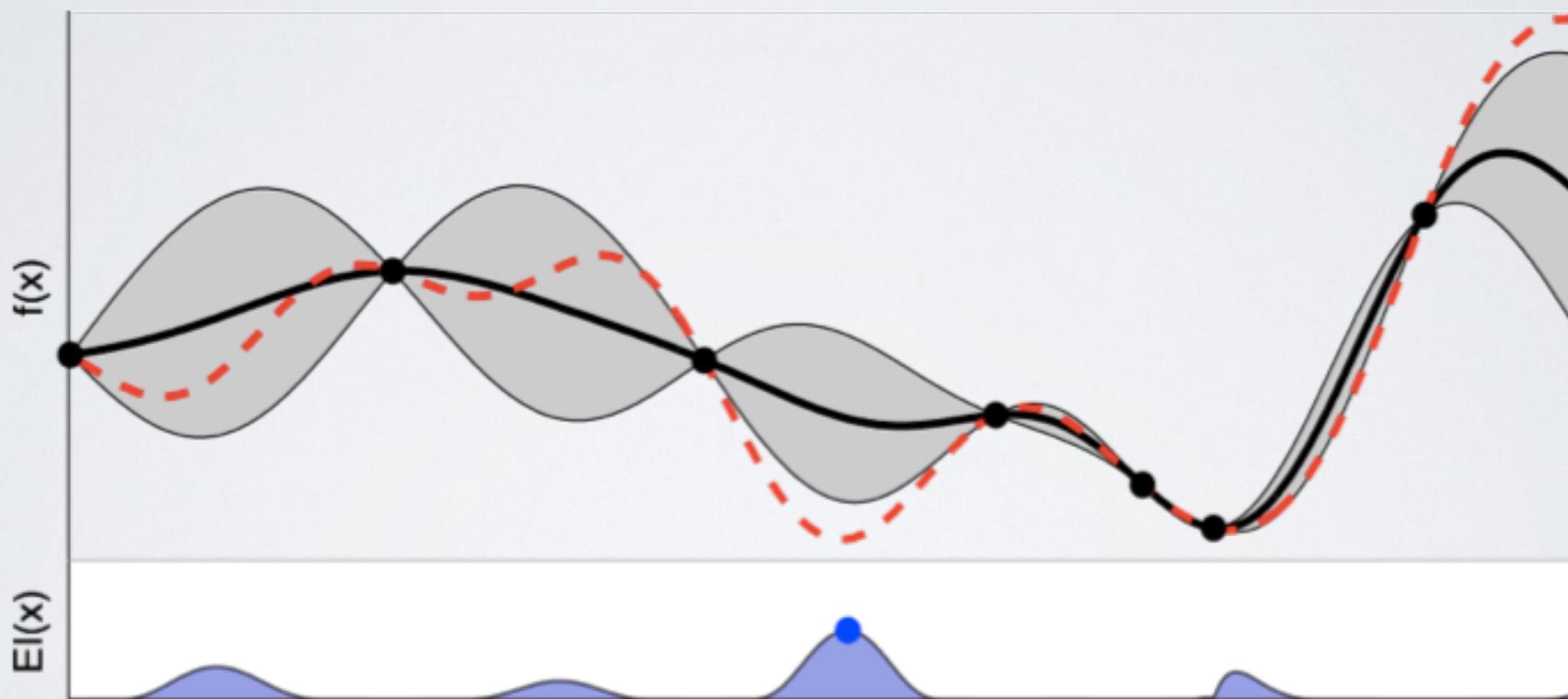
Simple Example



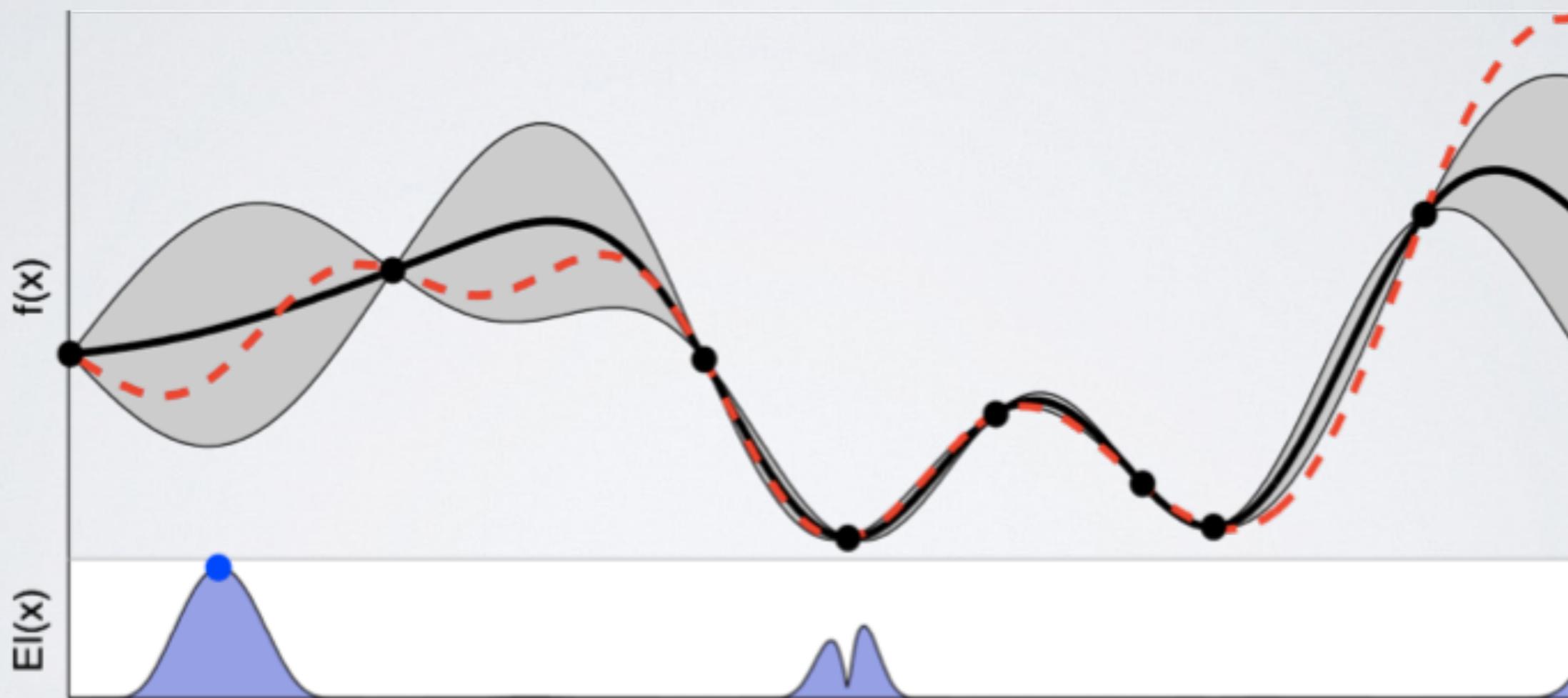
Simple Example



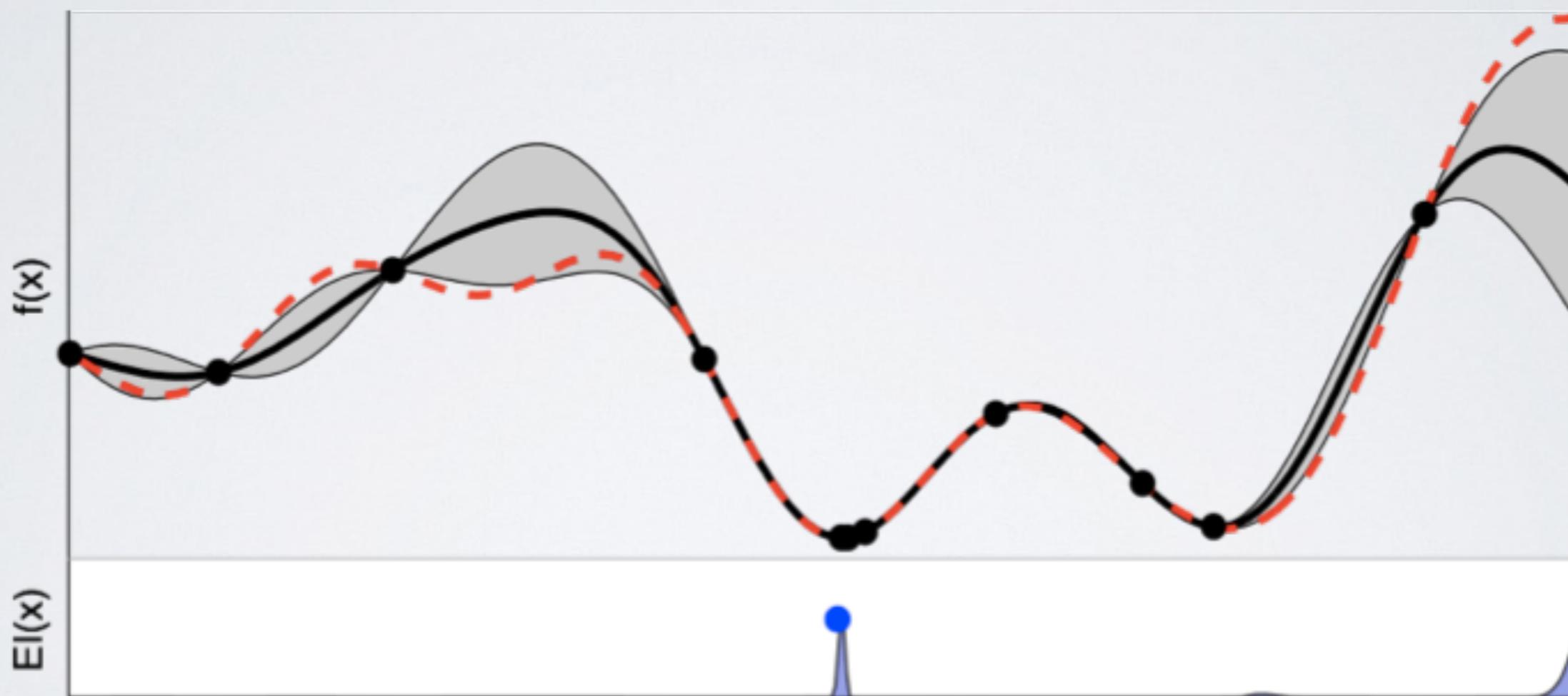
Simple Example



Simple Example



Simple Example



Check out Spearmint:

<https://github.com/HIPS/Spearmint>

Bayesian Optimization package for
Python, Matlab