

Lecture 25

How to Optimize Expensive Functions

Objectives

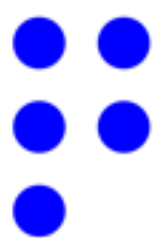
- Quantify the value of the information extracted from an experiment/simulation.
- Optimize an expensive black-box function under a limited budget.

The problem

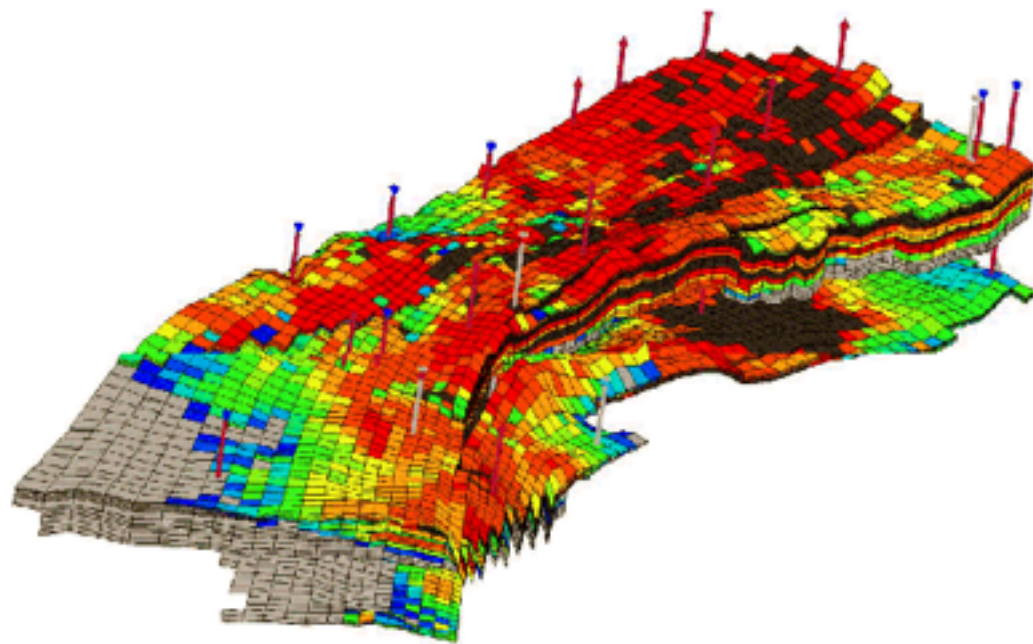
- Problem:

$$\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}} f(\mathbf{x})$$

- when the objective is:
 - very expensive to evaluate
 - you don't have gradients
 - dimensionality < 30 parameters

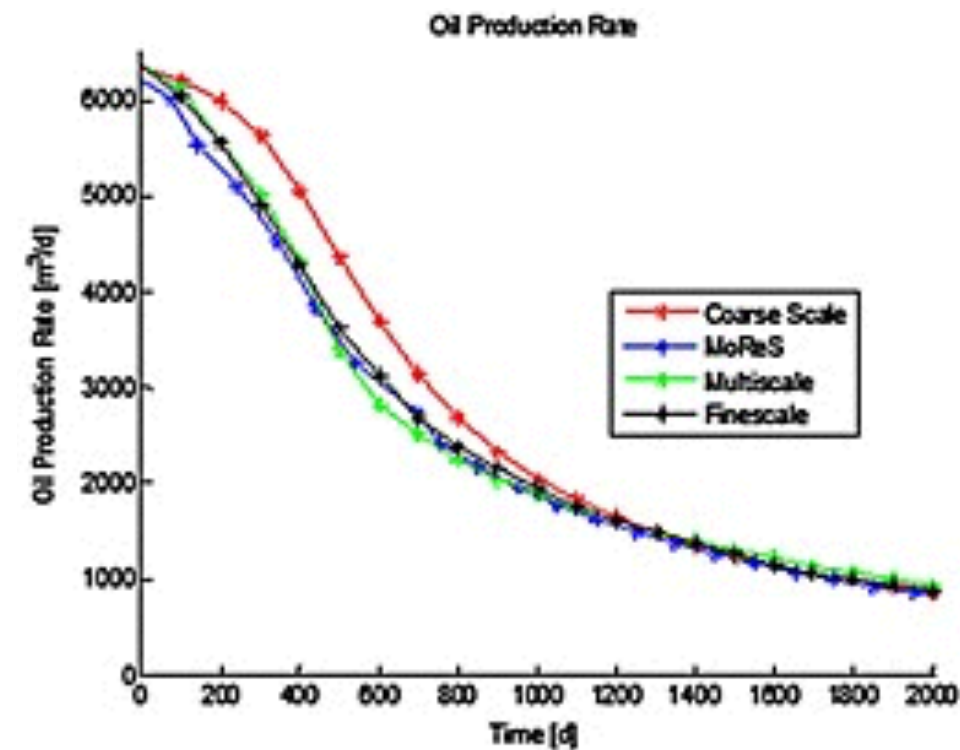


Example 1: Oil-well placement problem



\mathbf{x} = well locations

Simulation



$f(\mathbf{x})$ = expected net
present cost of
investment

Pandita, Billionis, and Panchal, 2016
<http://arxiv.org/abs/1604.01147>

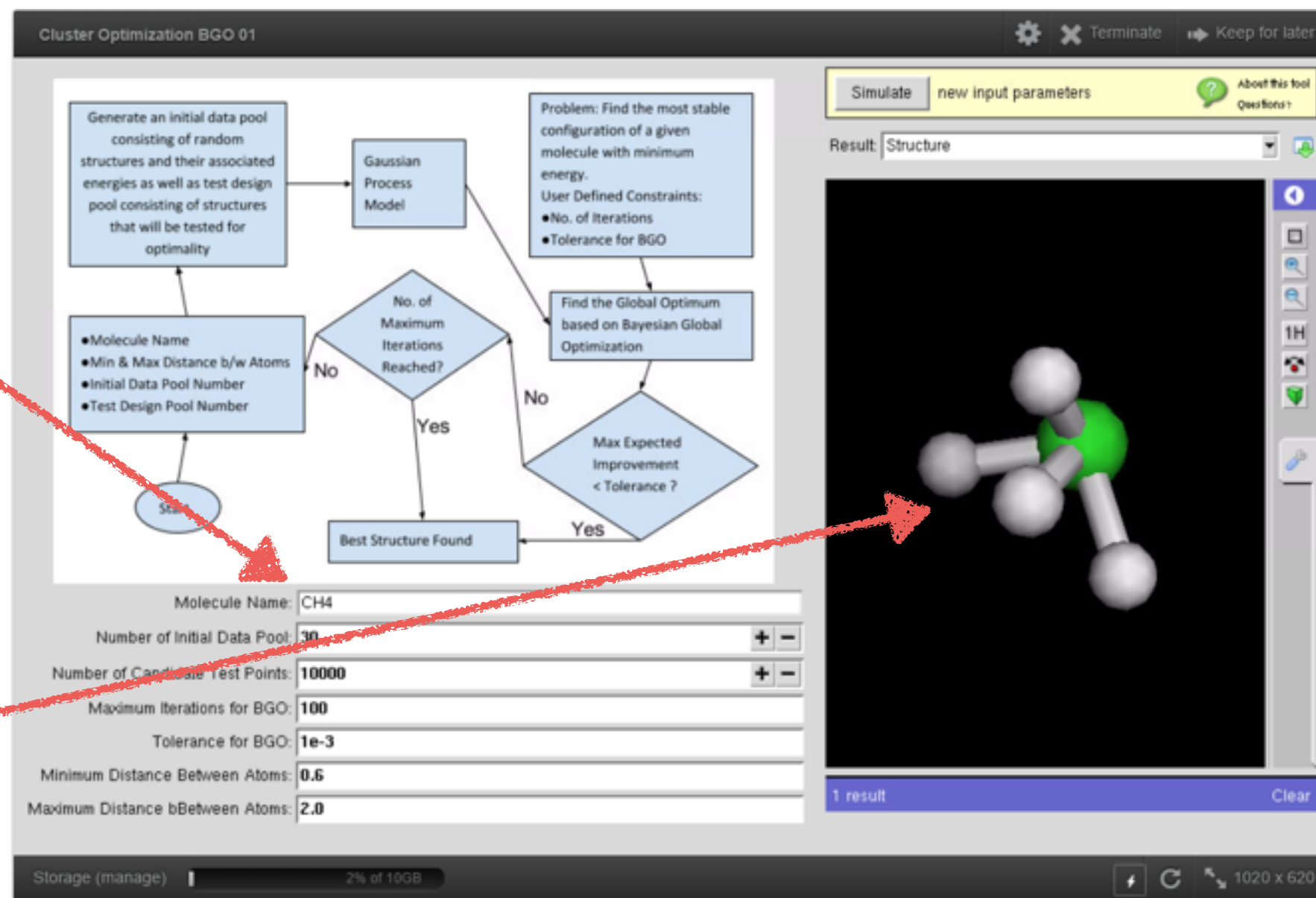
PREDICTIVE
SCIENCE LABORATORY

Example 2: Find stable structures

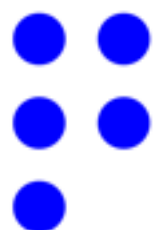
Surf student project

Chemical formula


Geometry (**x**)
with minimum
energy (**y**)



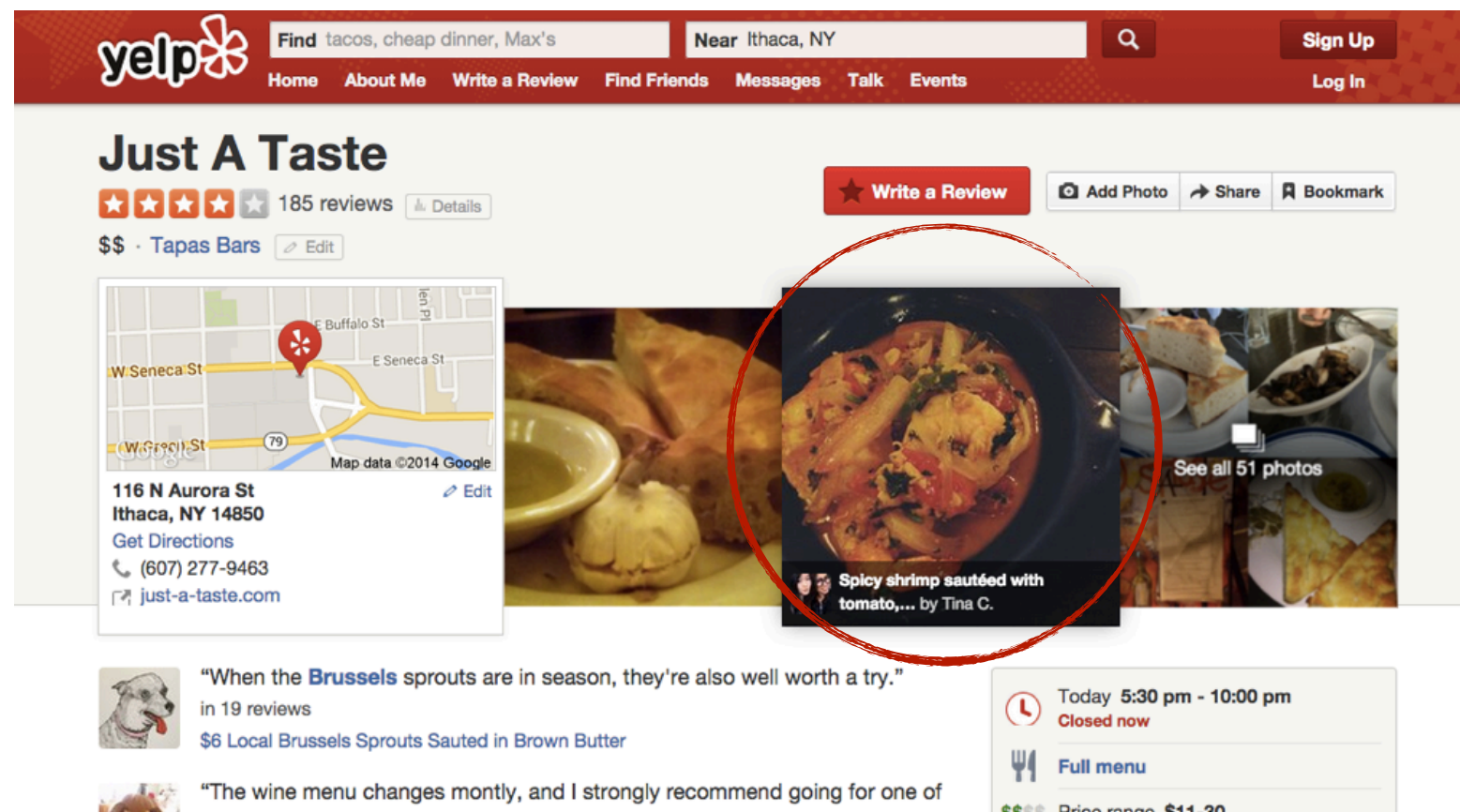
<https://nanohub.org/tools/clusterbgopro>



Example 3: Web site optimization

x = web design  measurement

-y = Number of views, seconds per view, etc.



Example 4: Training a robot to walk

<https://www.youtube.com/watch?v=uainbKfkc3Q>

<https://www.youtube.com/watch?v=GiqNQdzc5TI>

Other examples

- Model calibration (if posed as an optimization problem).
- Maximize efficiency in solar cells.
- Drug development.
- ...

Startup idea

<https://sigopt.com/>

[Solutions](#)[Pricing](#)[FAQ](#)[About](#)[Docs](#)[Login](#)[Sign Up](#)

Pricing

We have plans for teams both large and small. Start for free with no commitment or [contact us](#) and we'll help find the plan that is right for you.

SigOpt for Individuals

\$999/month

30 day free trial, no credit card required

10 experiments

Email support

[Sign Up](#)

SigOpt for Enterprise

Contact for pricing

Free pilot

Customized number of experiments

Premium support

[Contact Us](#)

We love students here at SigOpt. So we're making our individual plan available to academics, free of charge. We also offer discounted rates on our Enterprise plans for educational usage. Contact education@sigopt.com for more information.


**PREDICTIVE
SCIENCE LABORATORY**

Free options

- <https://github.com/PredictiveScienceLab/py-bgo> (features stochastic and multi-objective optimization)
- <https://github.com/SheffieldML/GPyOpt> (features parallel optimization)

```
pip install GPyOpt
```

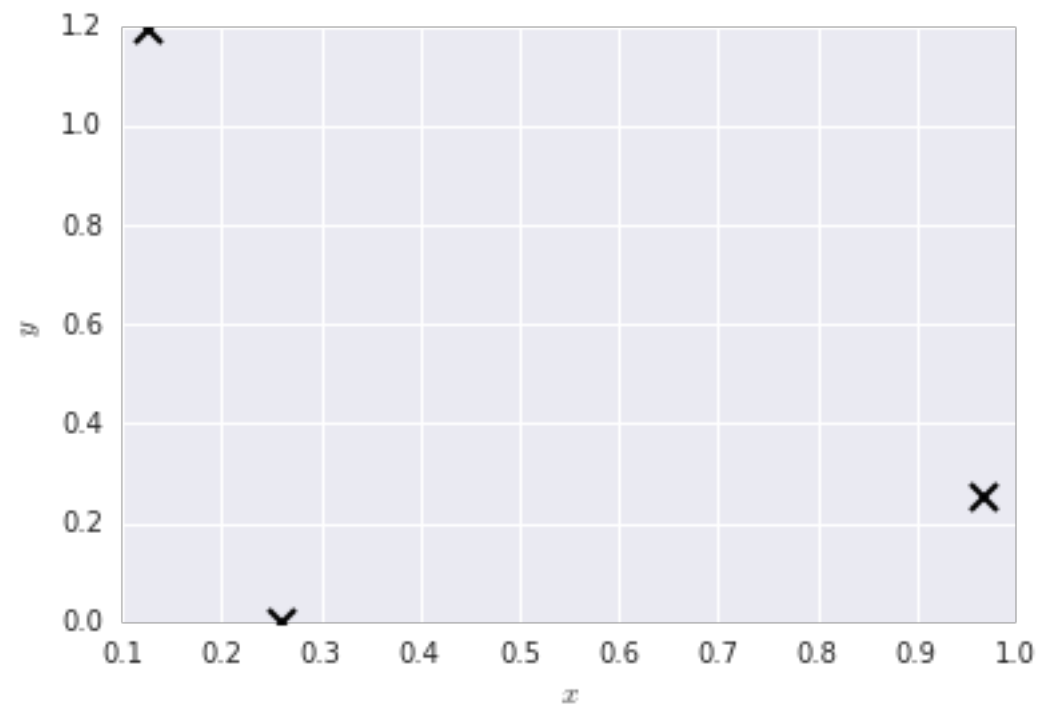
Idea

- 
1. We have some observed data (designs \mathbf{x} - vs - objectives \mathbf{y}).
 2. We fit a **statistical regression model** to the data.
 3. For each candidate design, compute the **value of information (Vol)**.
 4. We find the design with the **maximum Vol**.
 5. We compute the objective for this design.

Repeat until:

- budget is exhausted;
- Vol low.

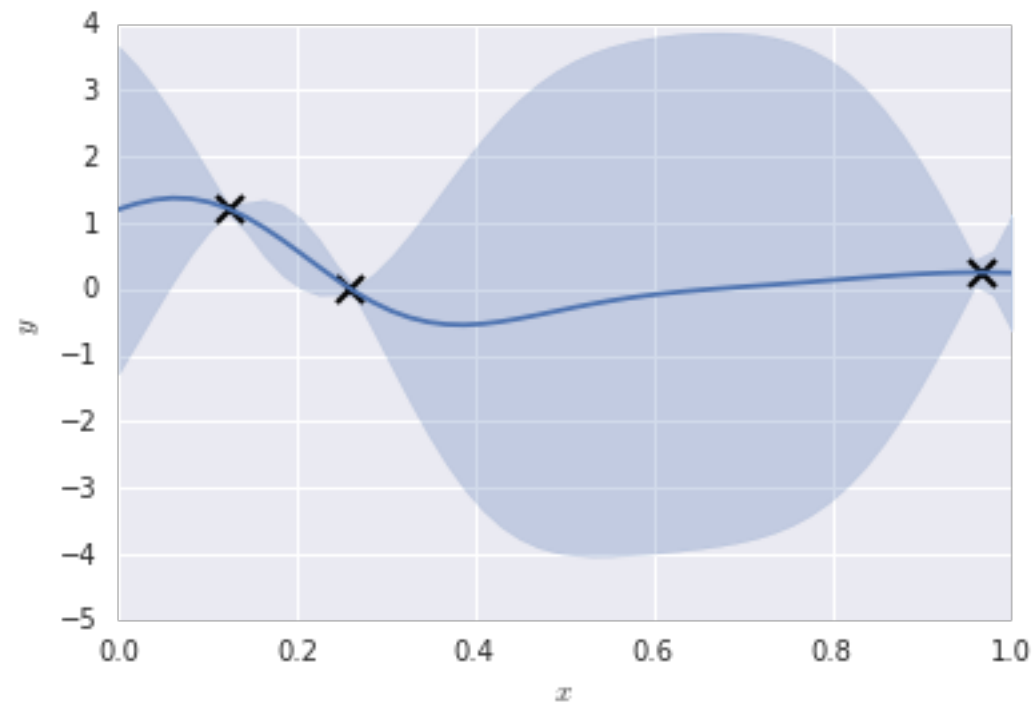
We have some data



$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

$$y_{1:n} = \{y_1, \dots, y_n\}$$

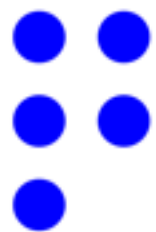
We fit a statistical model



$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$



Gaussian process regression

- Assume that we have observed:

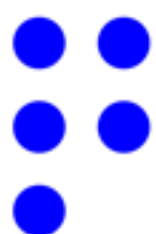
$$\mathbf{X} = \{x_1, \dots, x_N\},$$

$$\mathbf{f} = \{f(x_1), \dots, f(x_N)\}$$

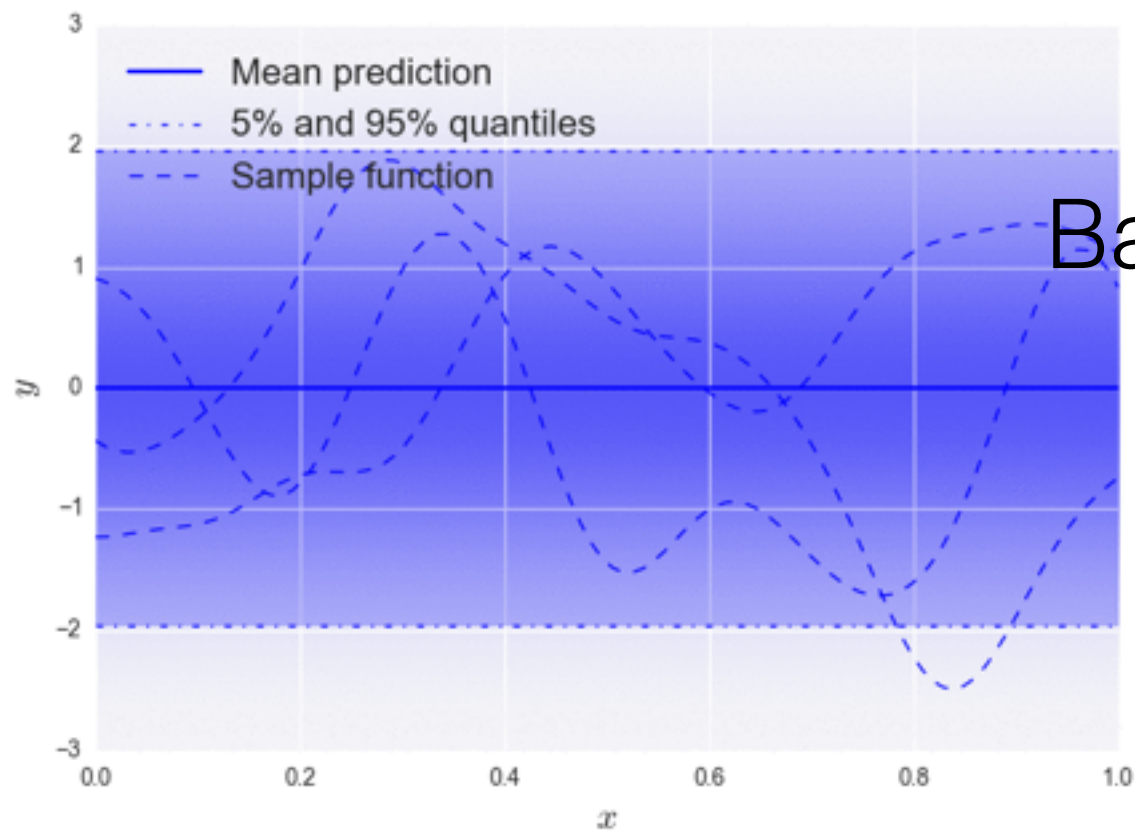
- and that we want to make predictions at an arbitrary set of *test* inputs:

$$\mathbf{X}^* = \{x_1^*, \dots, x_{N^*}^*\}$$

$$\mathbf{f}^* = \{f(x_1^*), \dots, f(x_{N^*}^*)\}$$

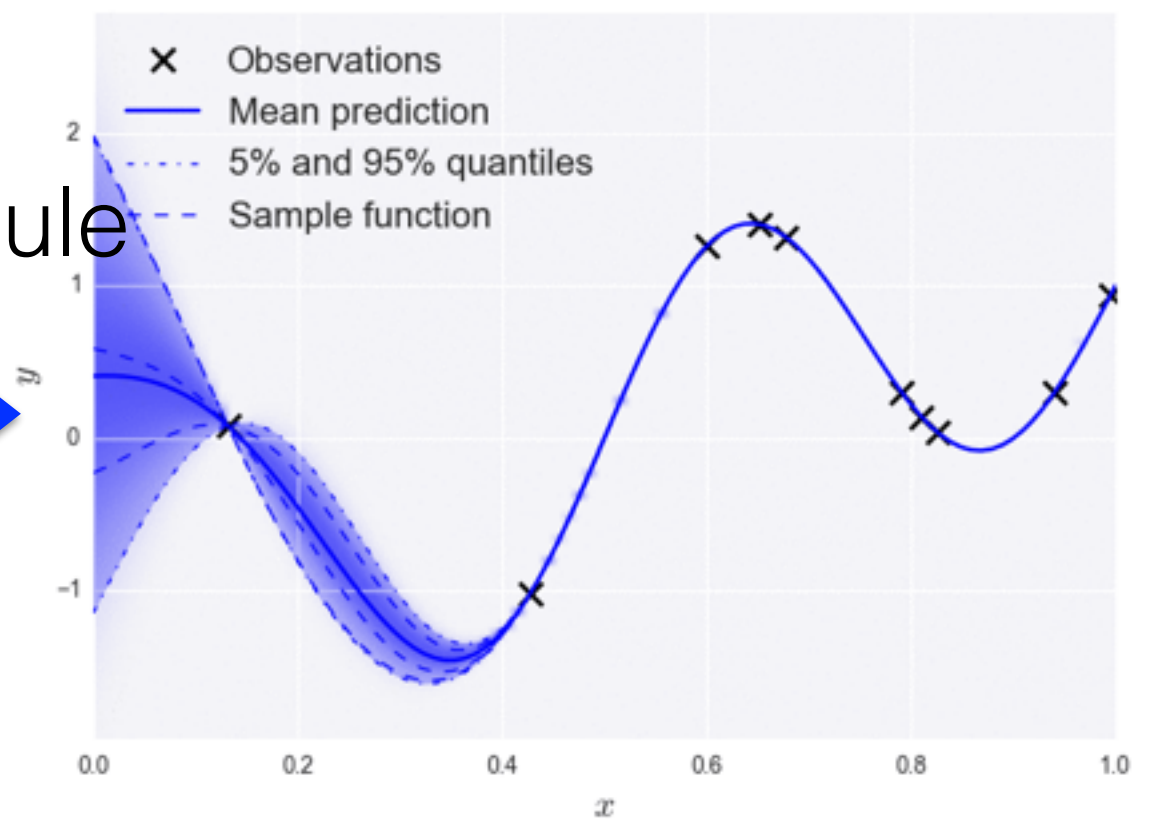
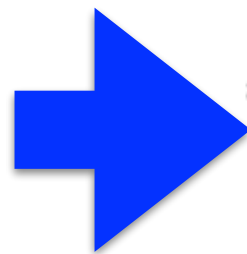


Gaussian process regression



Prior GP

Bayes rule



Posterior GP

The point predictive distribution

- Posterior GP:

$$f(\cdot) | \mathbf{X}, \mathbf{f} \sim \text{GP}(f(\cdot) | \tilde{m}(\cdot), \tilde{k}(\cdot, \cdot)),$$

- Looking at just one point, we get the *point predictive distribution*:

$$y | \mathbf{x}, \mathbf{X}, \mathbf{f} \sim \mathcal{N}(y | \tilde{m}(\mathbf{x}), \tilde{\sigma}^2(\mathbf{x})),$$

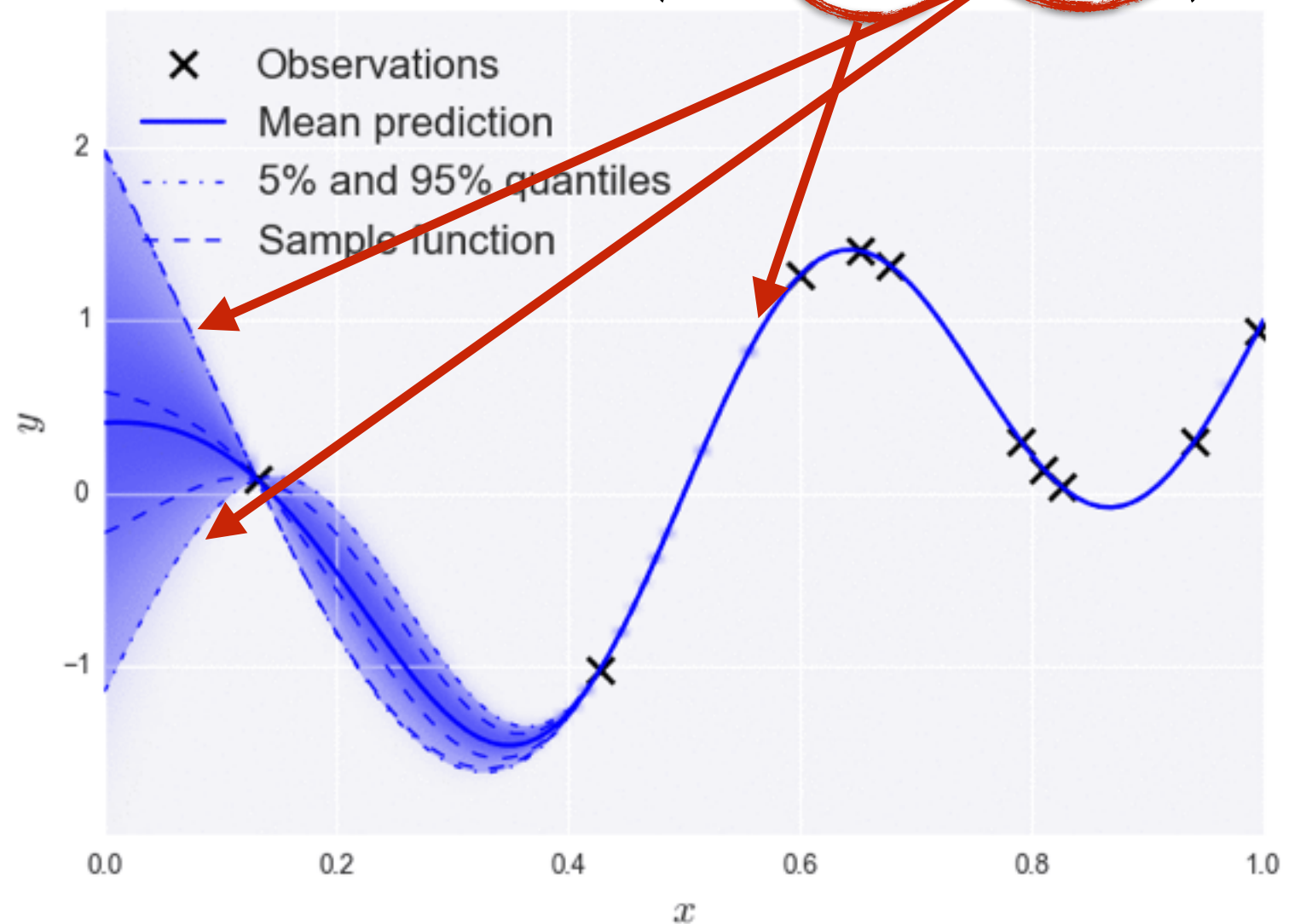
$$\tilde{\sigma}^2(\mathbf{x}) = \tilde{k}(\mathbf{x}, \mathbf{x}).$$

- You may use the mean as a surrogate.

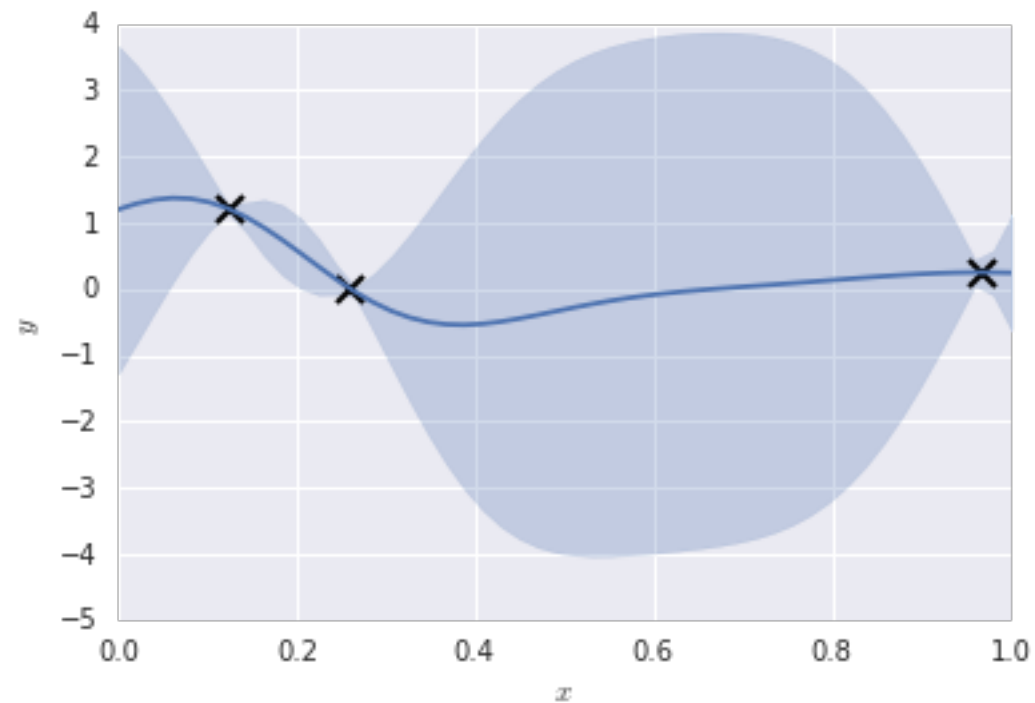
Gaussian process regression

$$y | \mathbf{x}, \mathbf{X}, \mathbf{f} \sim \mathcal{N}(y | \tilde{m}(\mathbf{x}), \tilde{\sigma}^2(\mathbf{x})),$$

$$f(\mathbf{x}) = \tilde{m}(\mathbf{x}) \pm 2\tilde{\sigma}(\mathbf{x})$$



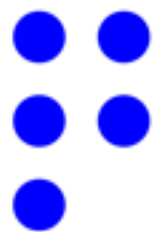
We fit a statistical model



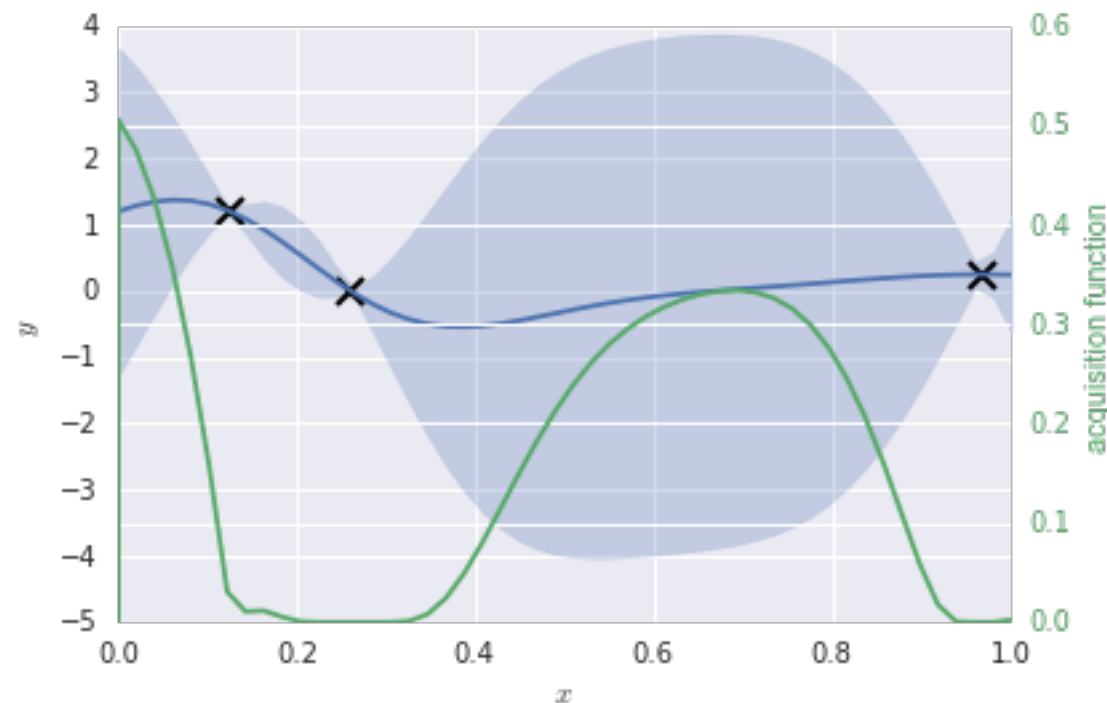
$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$



Quantify the value of information via an acquisition function

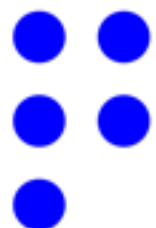


$a(\mathbf{x})$

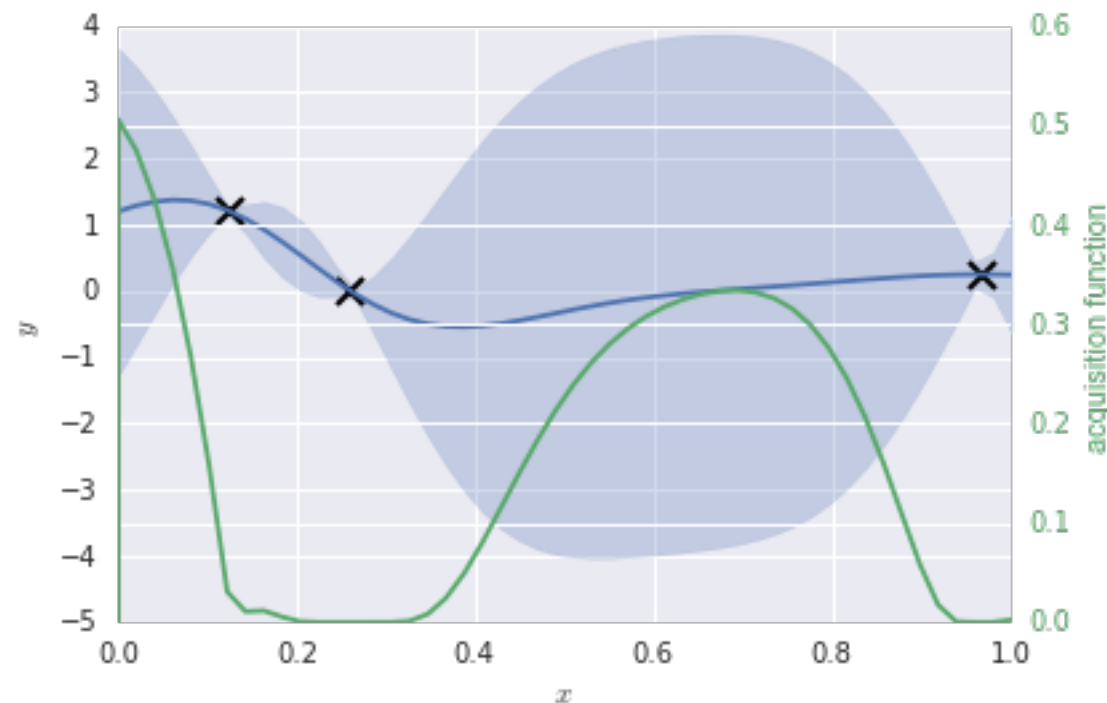
$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$



Quantify the value of information via an acquisition function



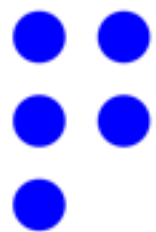
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

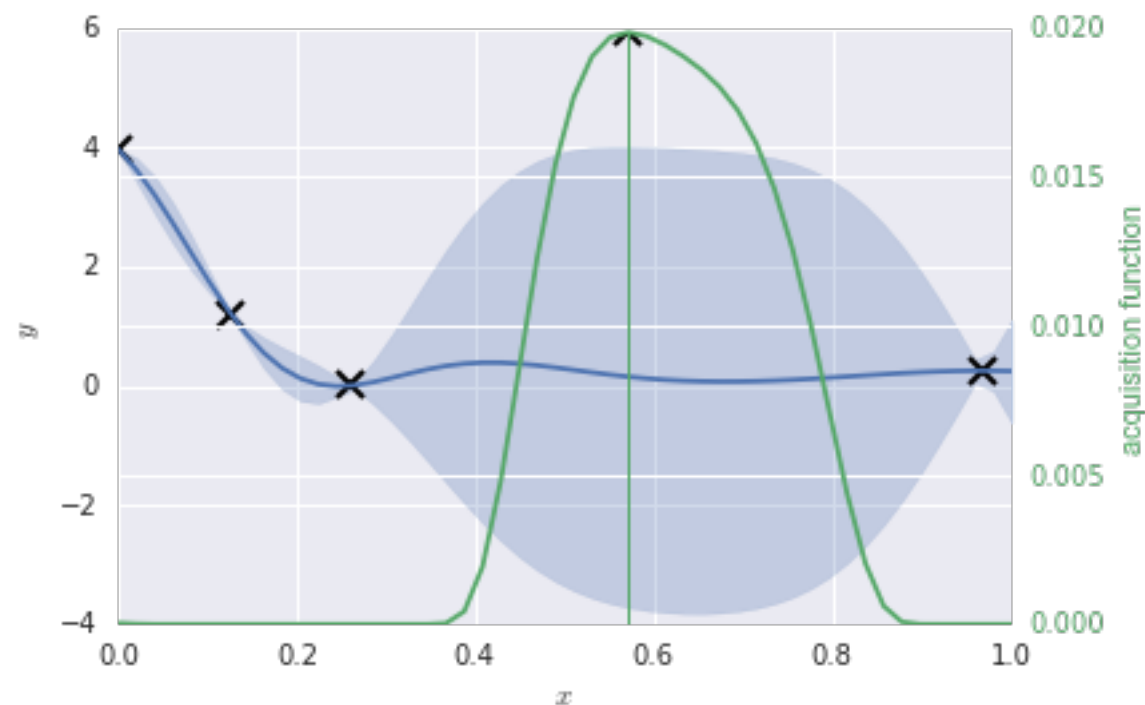
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$



Repeat (Iteration 2)



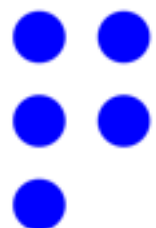
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

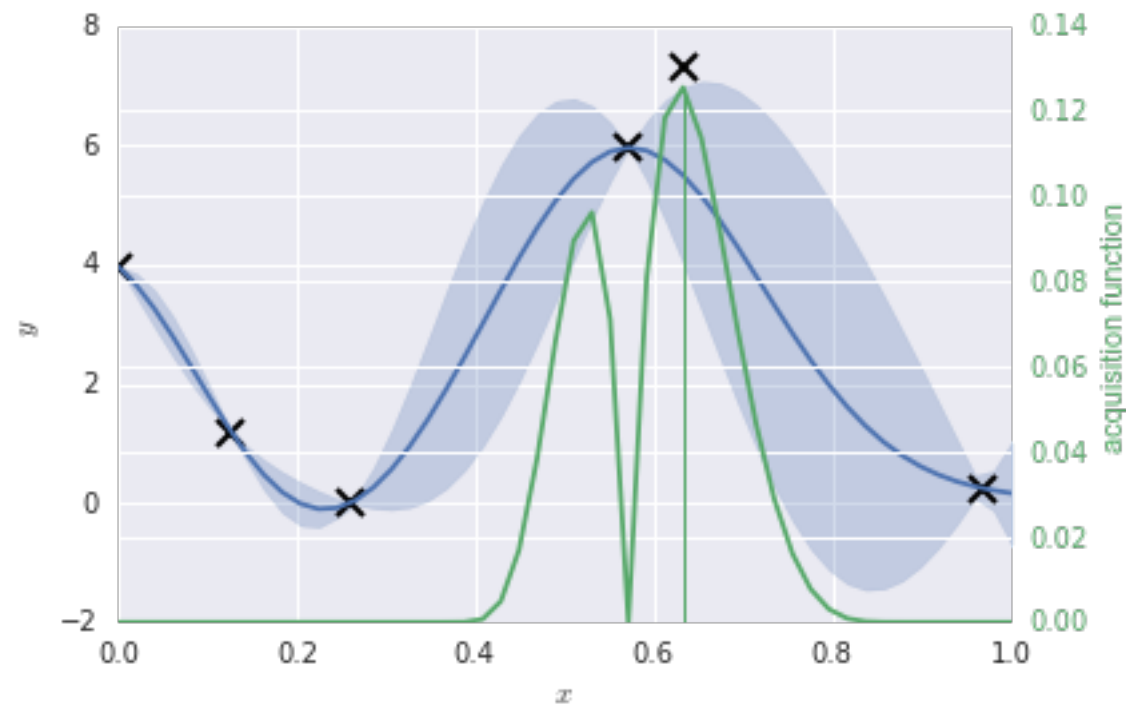
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$



Repeat (Iteration 3)



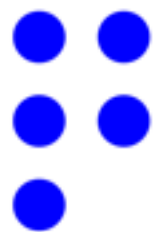
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

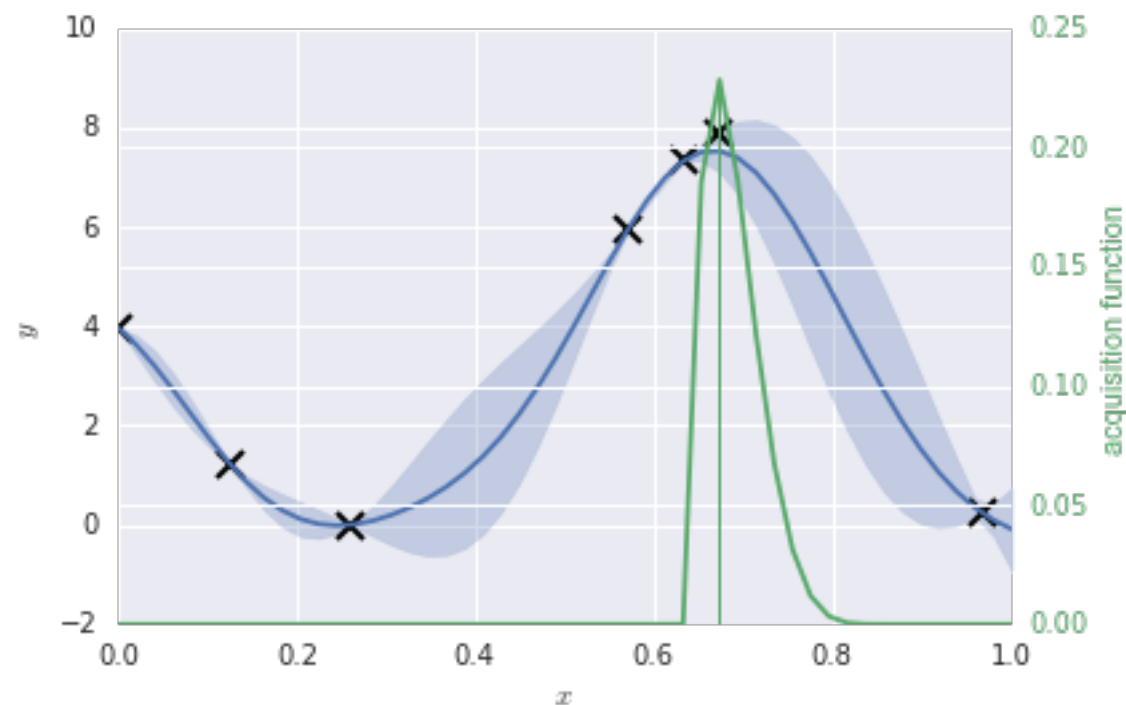
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$



Repeat (Iteration 3)



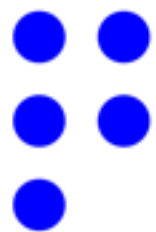
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

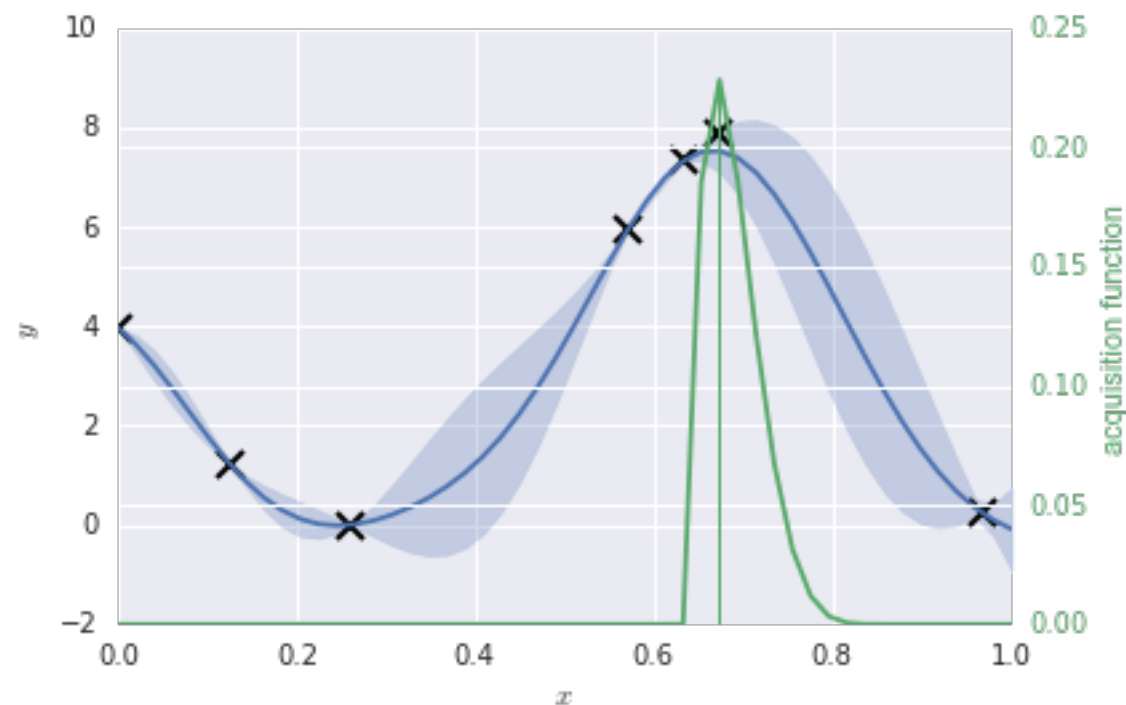
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$



Repeat (Iteration 4)



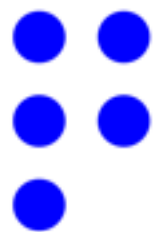
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

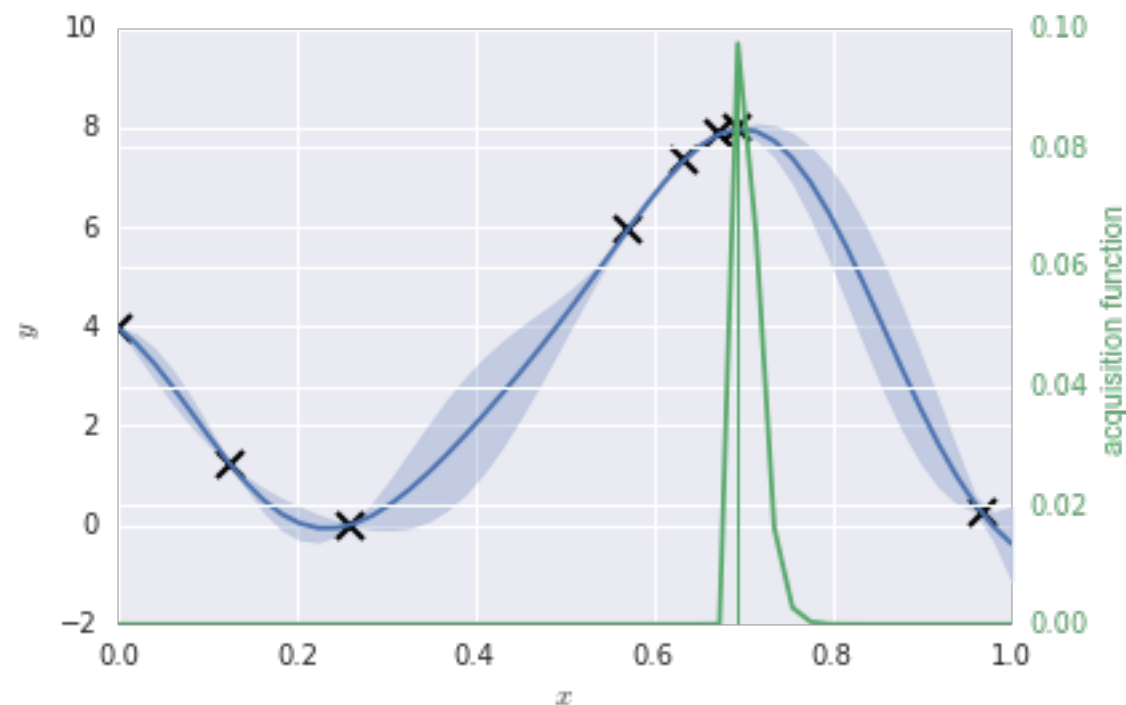
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$



Repeat (Iteration 5)



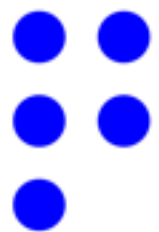
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$



The problem

- Problem:

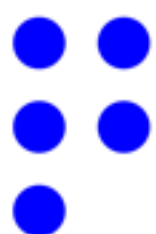
$$\mathbf{x}^* = \operatorname{argmax}_{\mathbf{x}} f(\mathbf{x})$$

- when the objective is:

- very expensive to evaluate

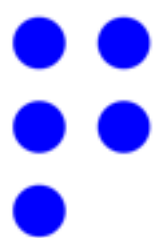
- you don't have gradients

- dimensionality < 30 parameters



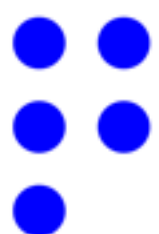
Dynamic vs Myopic Information Acquisition

- Optimal information acquisition policies...
- \Rightarrow Dynamic programming/control theory.
- Too hard mathematical/computational problems.
- What if, we just pick one piece of information at a time?
- Myopic (one-step-look-ahead) policies.



The value of information

- The value of information (Vol) depends on what you want to do.
- Can be quantified objectively if:
 - you have assigned probabilities over all possibilities.
 - you can quantify your profit/loss if any of the possibilities happen.



The value of information

Vol of **x** = how much expected gain if I measure at **x**

= expected profit if I measure at **x**

- current best alternative

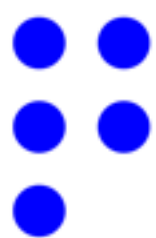
= expected income if I measure at **x**

- cost of measuring **x**

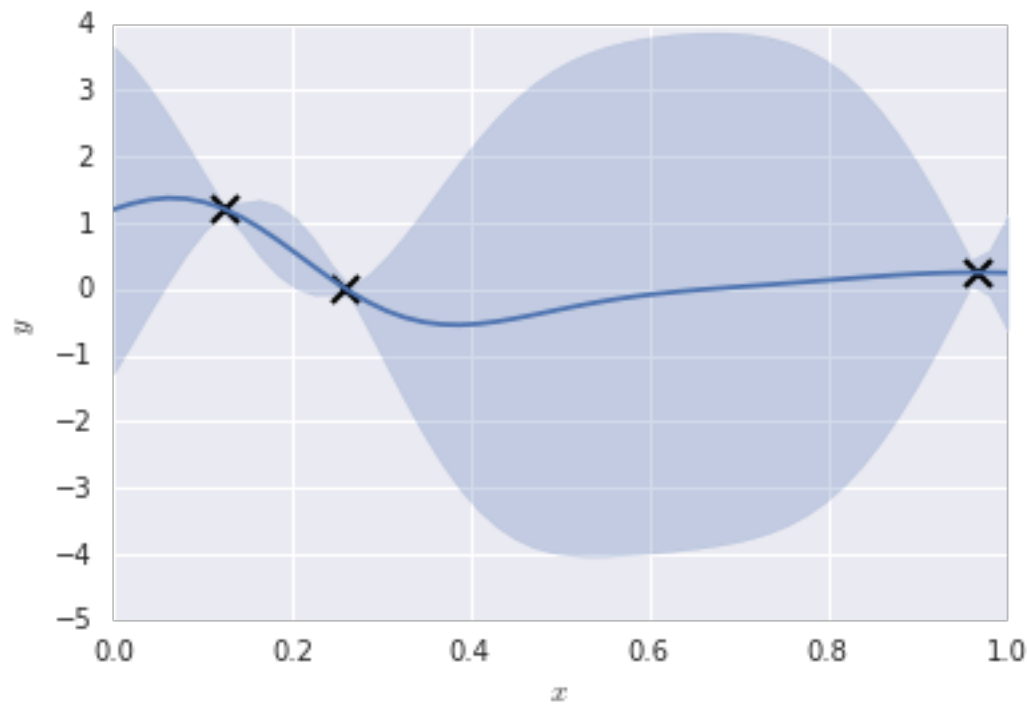
- current best alternative

Acquisition function

- Most of the times, we don't have the details to find the Vol.
- We use heuristic approximations to Vol such as:
 - the probability of improvement
 - the expected improvement
 - the knowledge gradient
 - the expected information gain



Maximum Mean



Add the point with the maximum expected mean:

$$a(\mathbf{x}) = m(\mathbf{x})$$

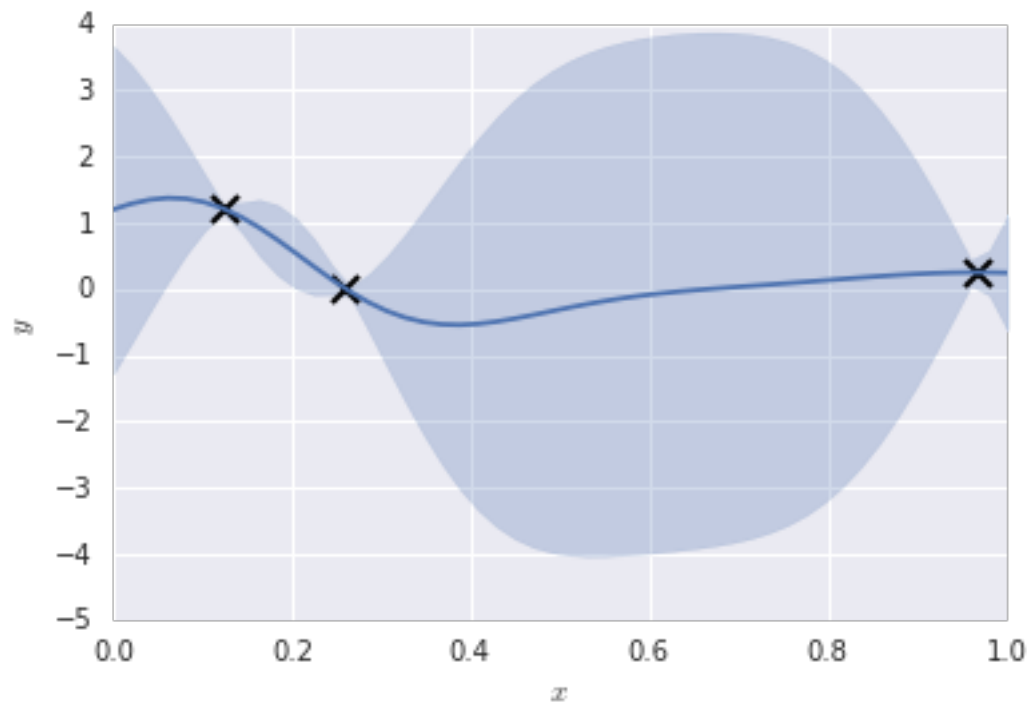
Tries to exploit what we know.

It does not converge.

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$
$$y_{1:n} = \{y_1, \dots, y_n\}$$
$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

How can we add some elements of exploration?

Maximum Upper Interval



Use the variance to explore:

$$a(\mathbf{x}) = m(\mathbf{x}) + \psi\sigma(\mathbf{x})$$

Adds upper quantile.

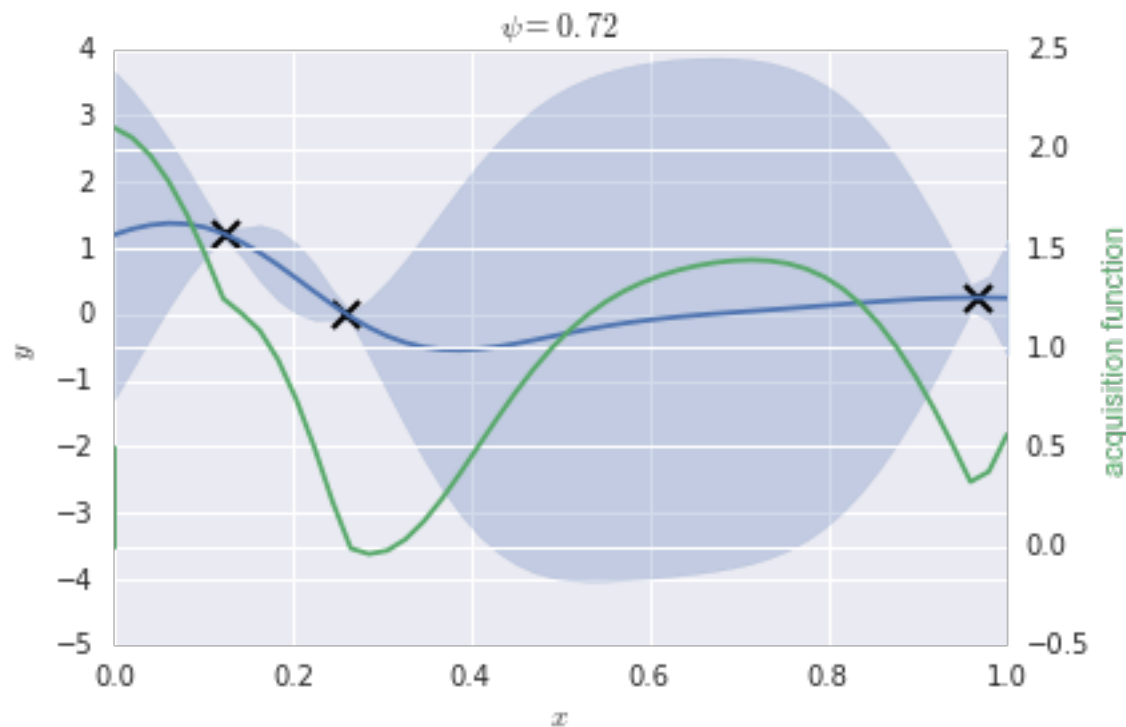
Provable convergence to the global maximum!

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

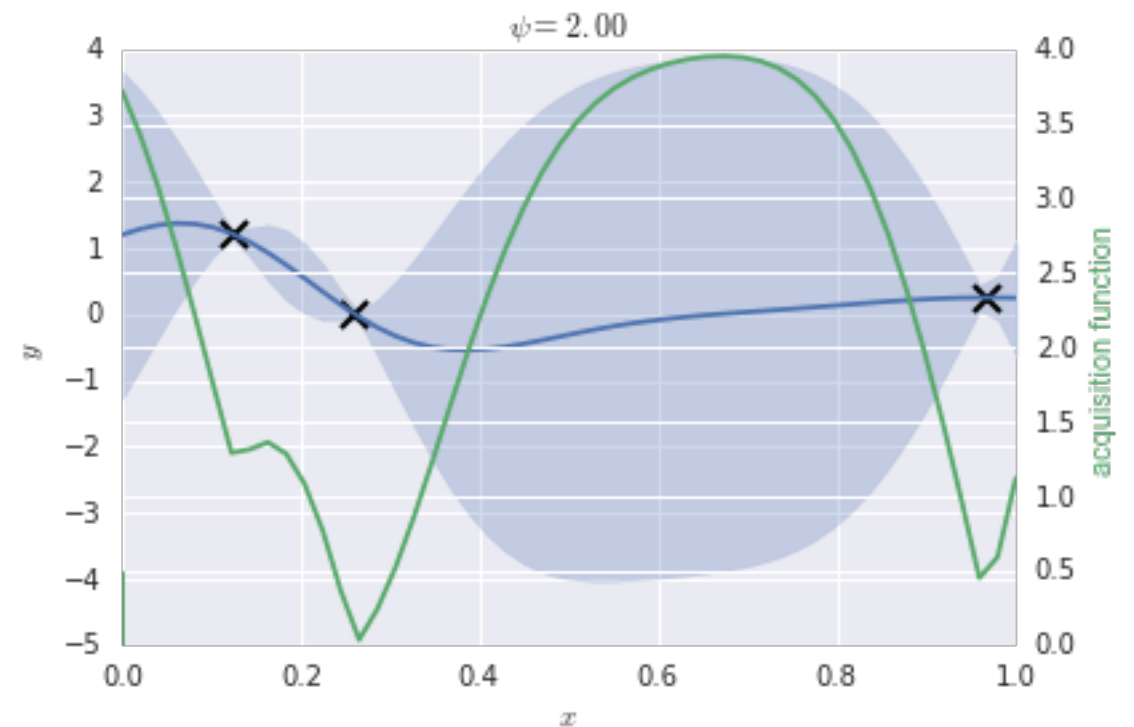
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

Maximum Upper Interval



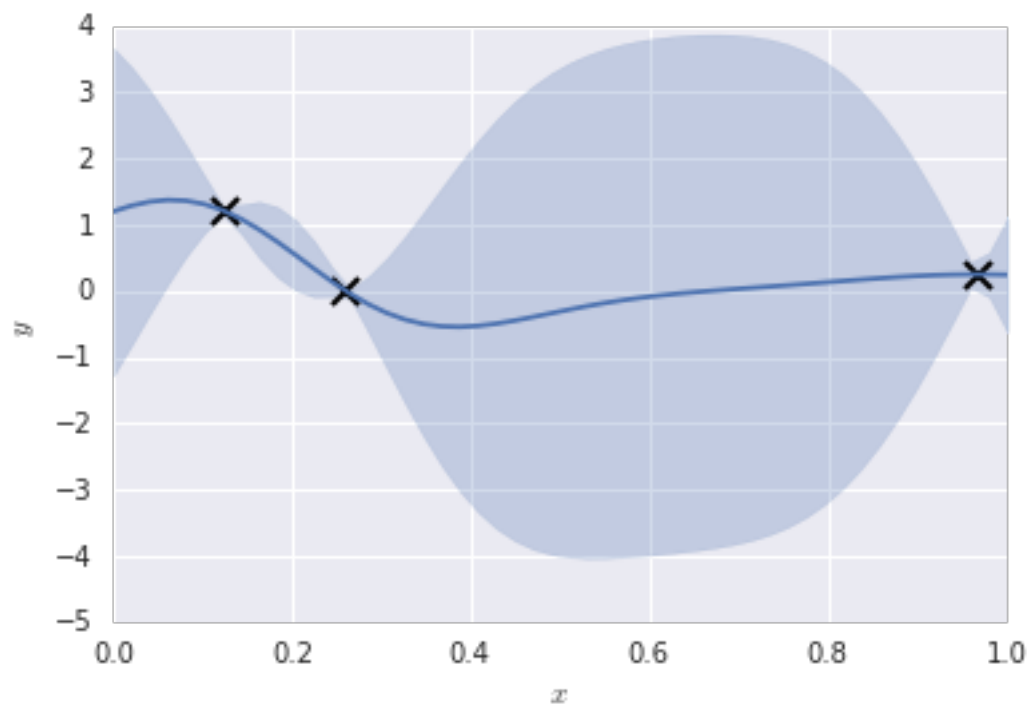
Exploits more



Explores more

Too much exploration...

Probability of Improvement



Current best:

$$\tilde{y}_n = \max_{1 \leq i \leq n} y_i$$

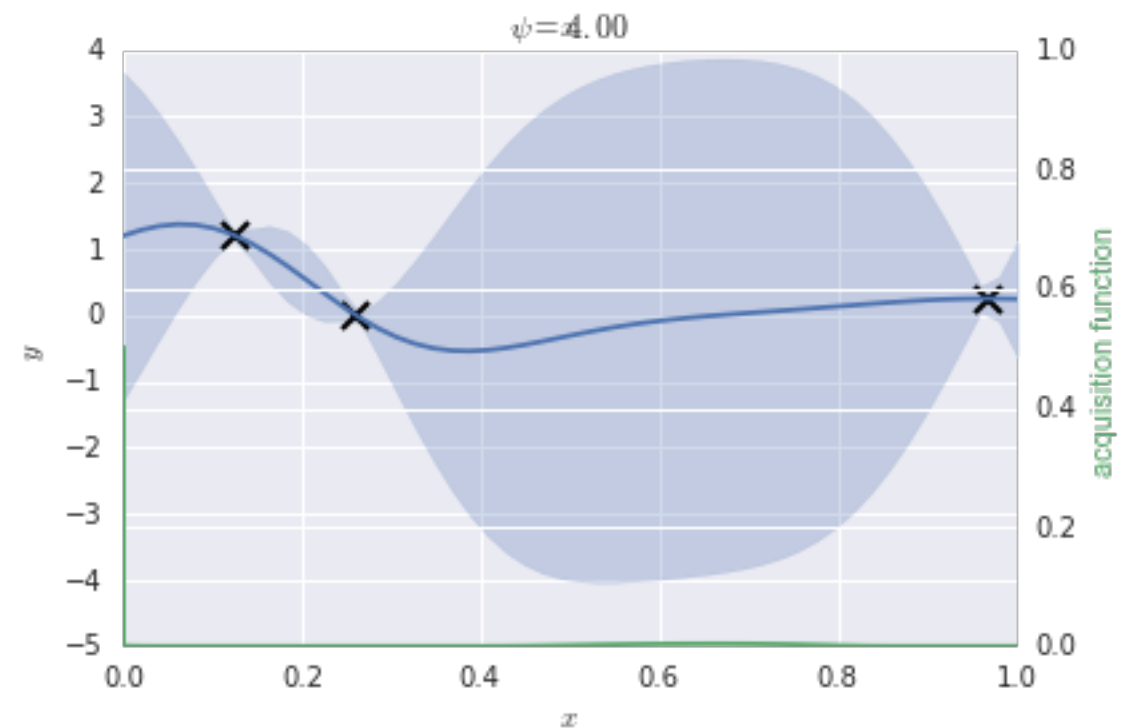
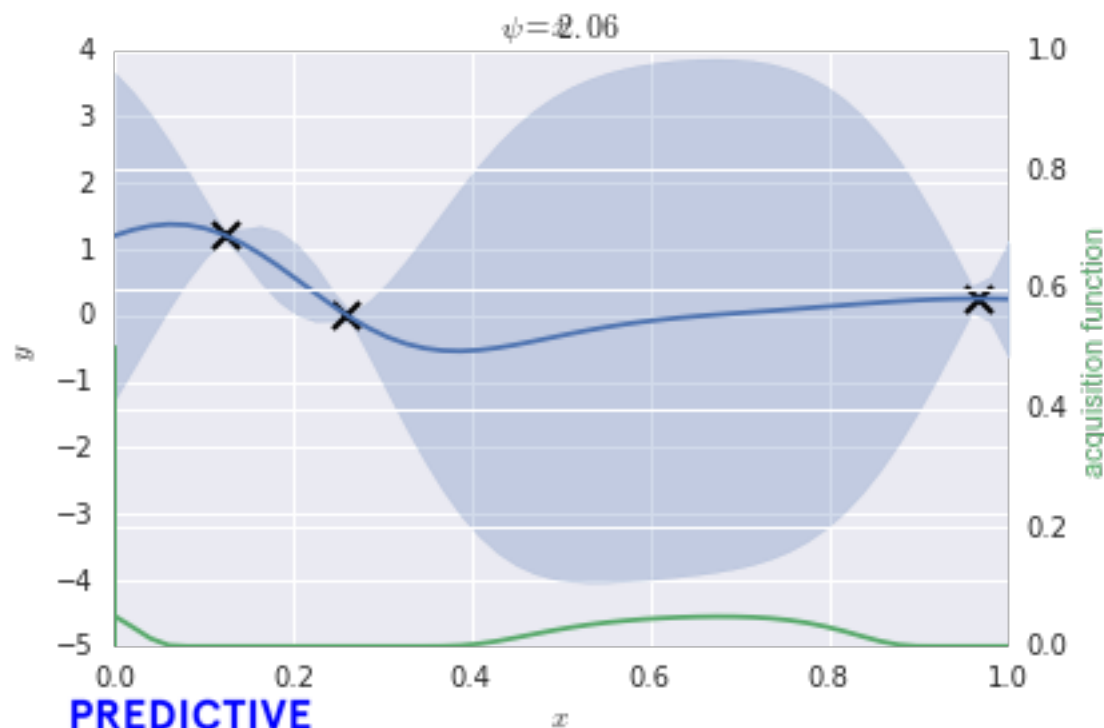
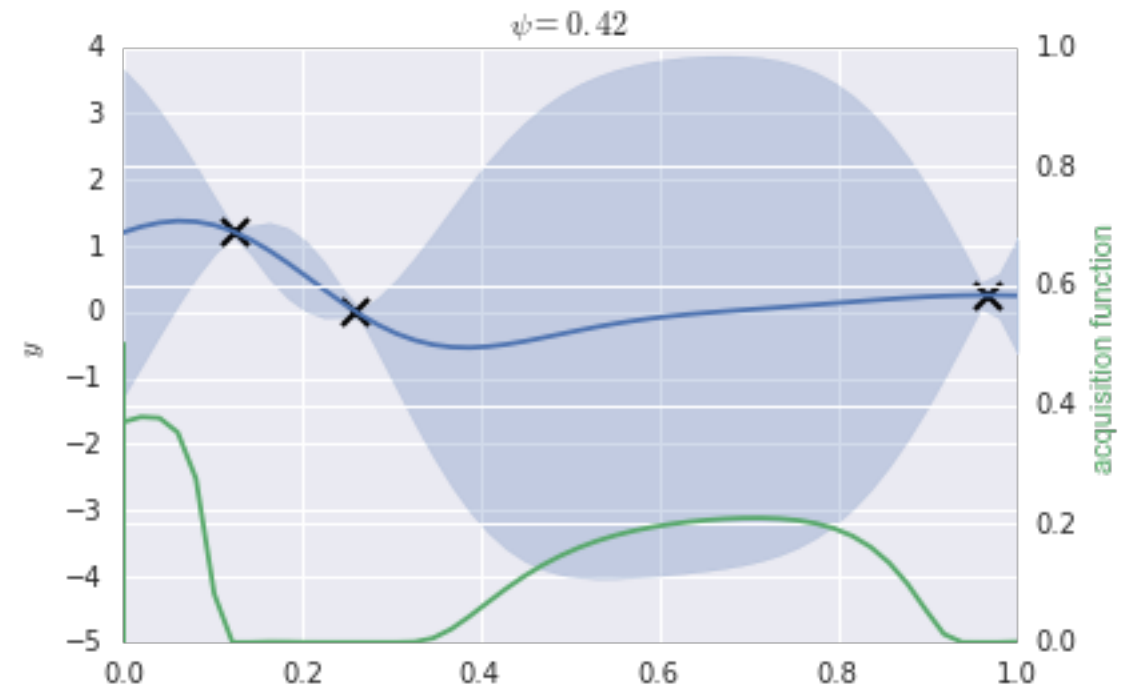
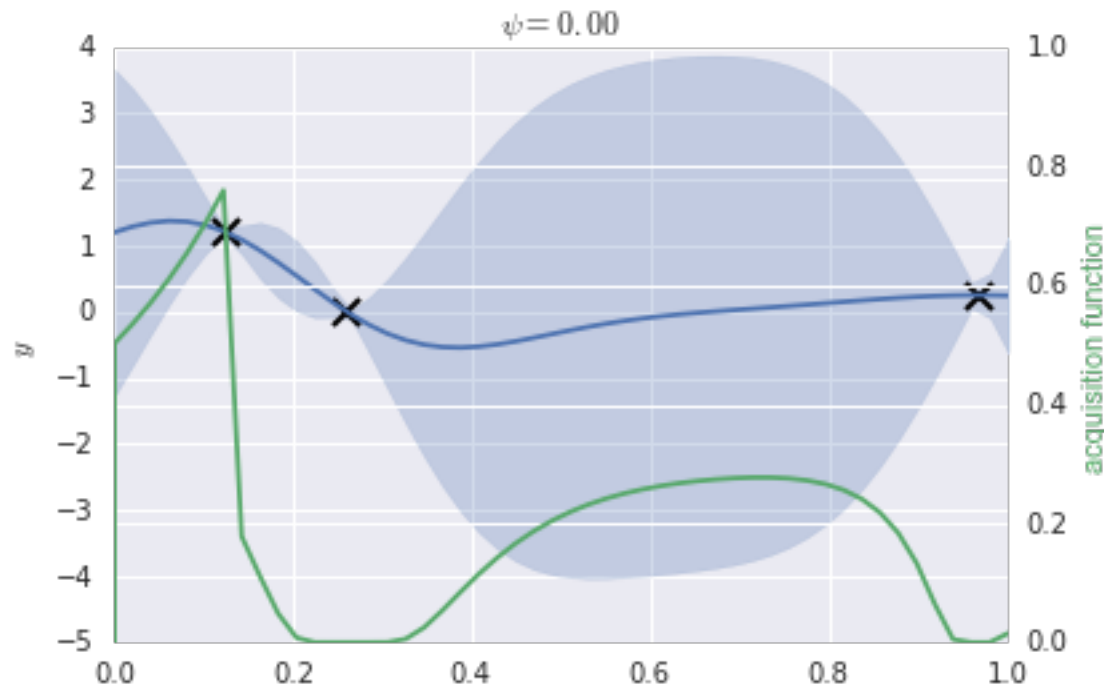
Hypothetical simulation at \mathbf{x} yields y

The probability of improvement is:

$$\begin{aligned}\mathbf{x}_{1:n} &= \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \\ y_{1:n} &= \{y_1, \dots, y_n\} \\ p(y | \mathbf{x}) &\approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))\end{aligned}$$

$$\begin{aligned}a(\mathbf{x}) &= P[y > \tilde{y}_n + \psi | \mathbf{x}] \\ &= \int_{\tilde{y}_n + \psi}^{\infty} p(y | \mathbf{x}) dy \\ &= 1 - \Phi\left(\frac{\tilde{y}_n + \psi - m(\mathbf{x})}{\sigma(\mathbf{x})}\right)\end{aligned}$$

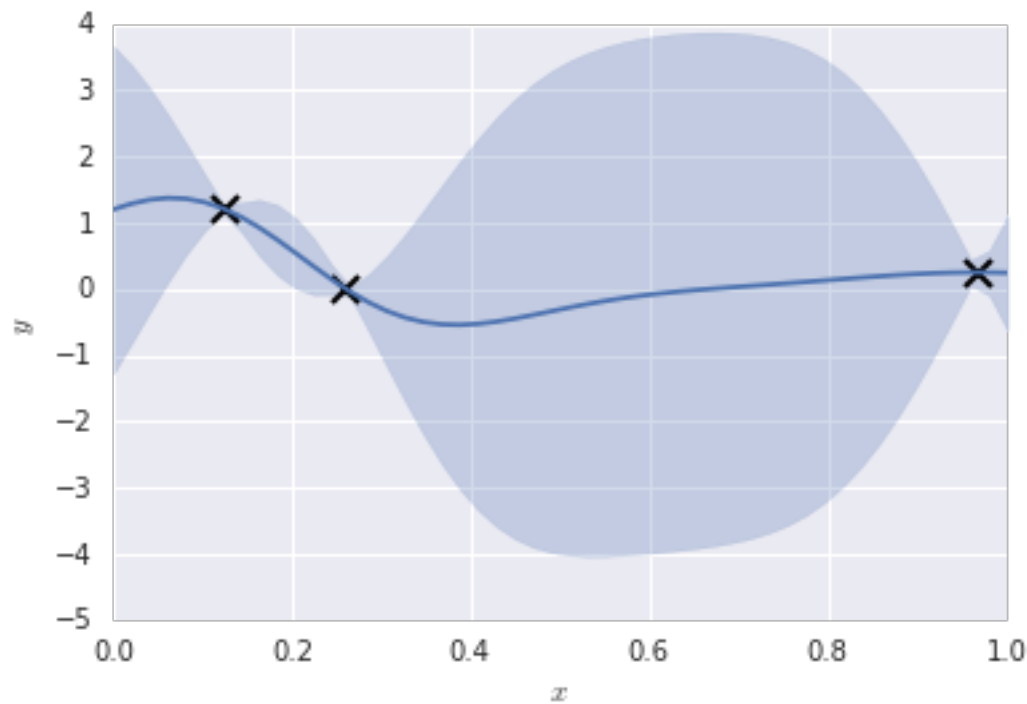
Probability of Improvement



Probability of Improvement - Why not use it?

- Large value of ψ -> exploration.
- Small value of ψ -> exploitation.
- But how to you pick it?

Expected Improvement



Current best:

$$\tilde{y}_n = \max_{1 \leq i \leq n} y_i$$

Hypothetical simulation at \mathbf{x} yields y

The improvement is:

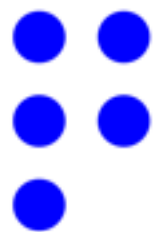
$$a(\mathbf{x}, y) = \max\{0, y - \tilde{y}_n\}$$

Integrate over y :

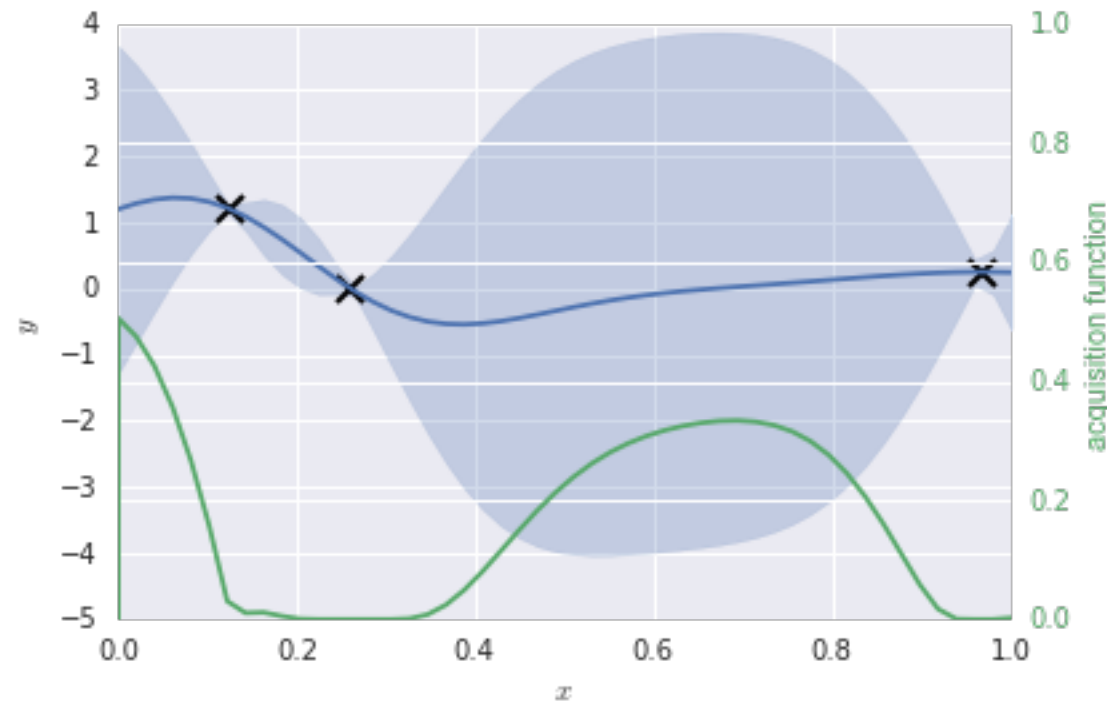
$$\begin{aligned}\mathbf{x}_{1:n} &= \{\mathbf{x}_1, \dots, \mathbf{x}_n\} \\ y_{1:n} &= \{y_1, \dots, y_n\} \\ p(y | \mathbf{x}) &\approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))\end{aligned}$$

$$a(\mathbf{x}) = \int_{-\infty}^{\infty} \max\{0, y - \tilde{y}_n\} p(y | \mathbf{x}) dy$$

$$= (m(\mathbf{x}) - \tilde{y}_n) \Phi\left(\frac{m(\mathbf{x}) - \tilde{y}_n}{\sigma(\mathbf{x})}\right) + \sigma(\mathbf{x}) \phi\left(\frac{m(\mathbf{x}) - \tilde{y}_n}{\sigma(\mathbf{x})}\right)$$

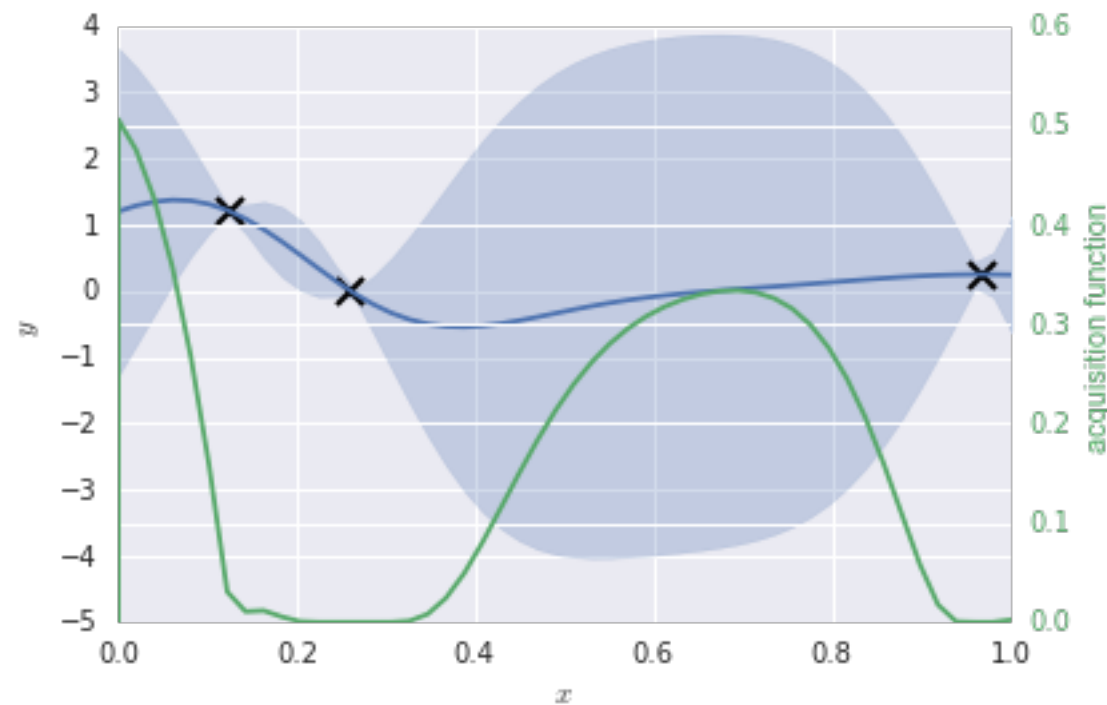


Expected Improvement



Automatic exploration vs exploitation...

Quantify the value of information via an acquisition function

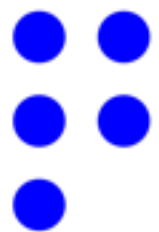


$a(\mathbf{x})$

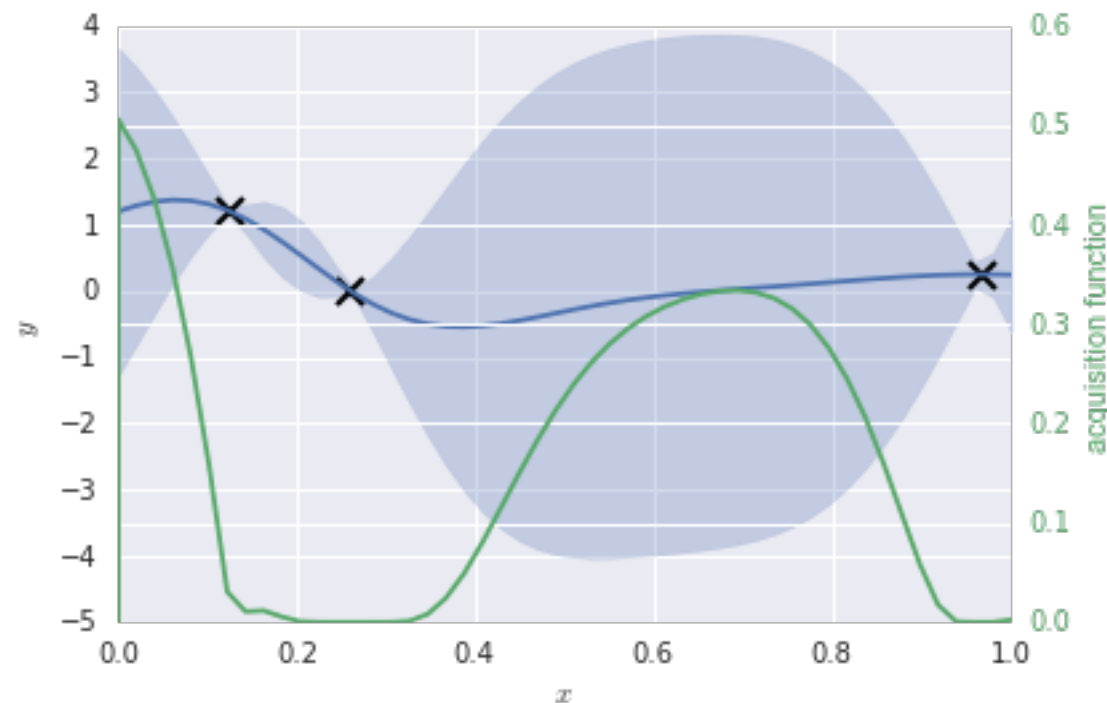
$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$



Quantify the value of information via an acquisition function



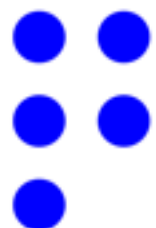
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

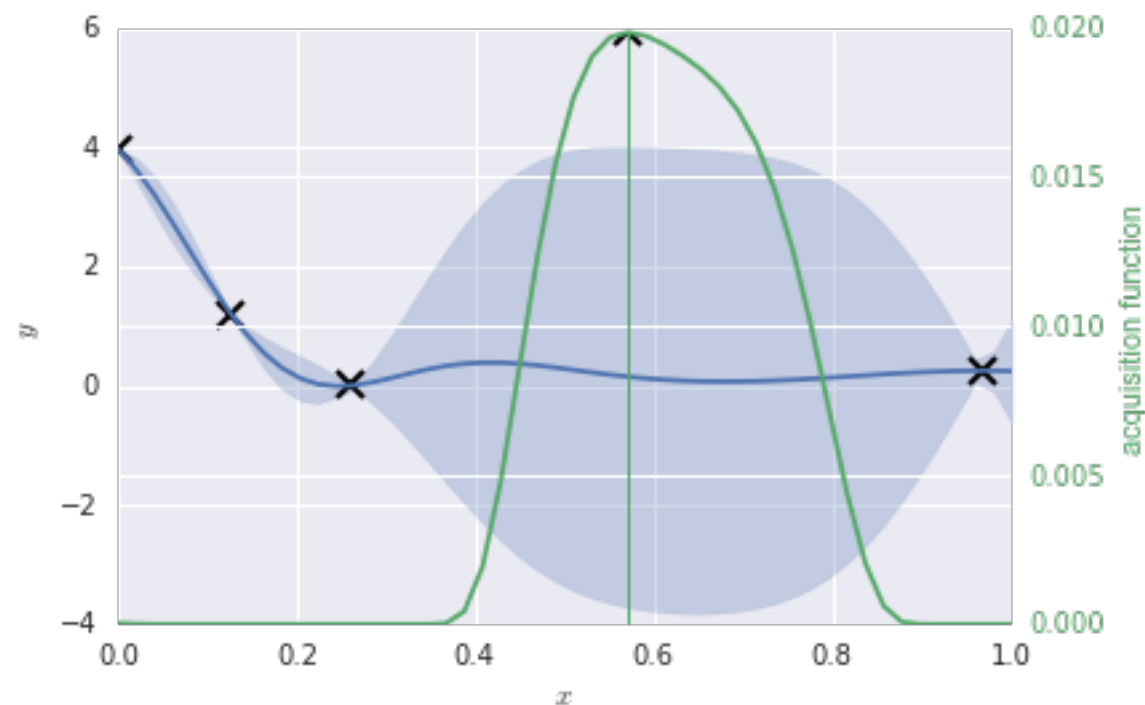
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$



Repeat (Iteration 2)



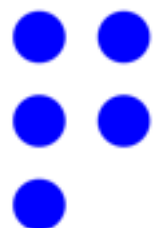
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

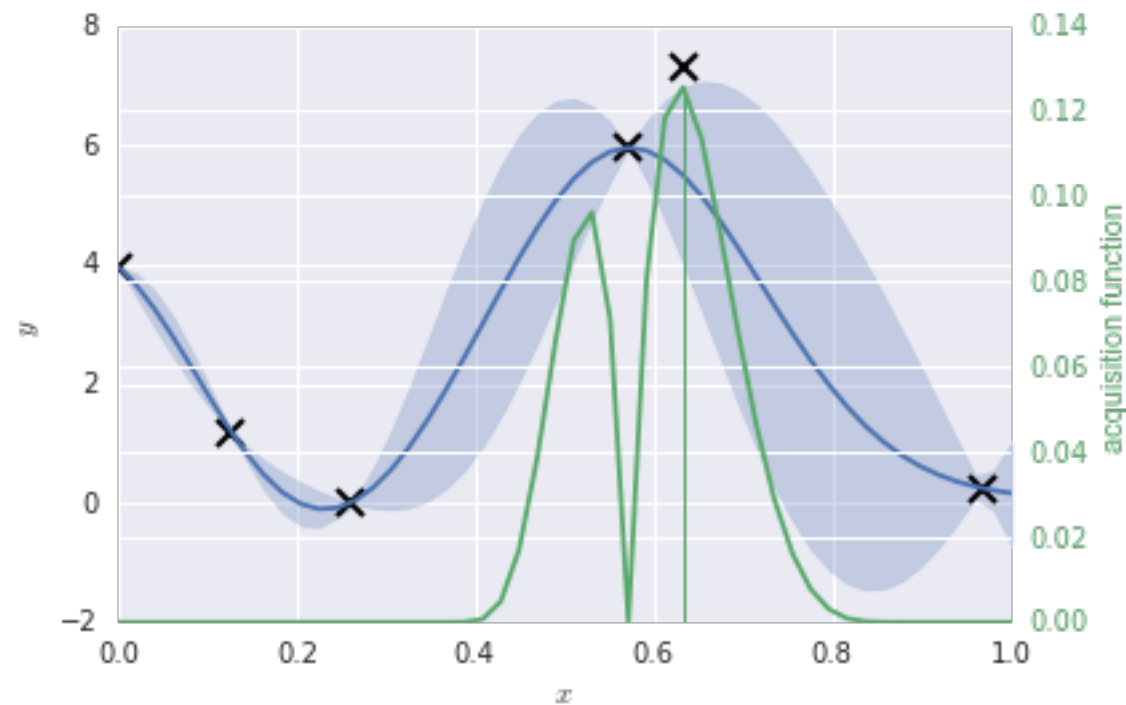
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$



Repeat (Iteration 3)



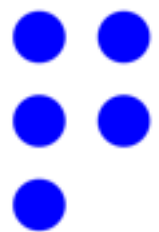
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

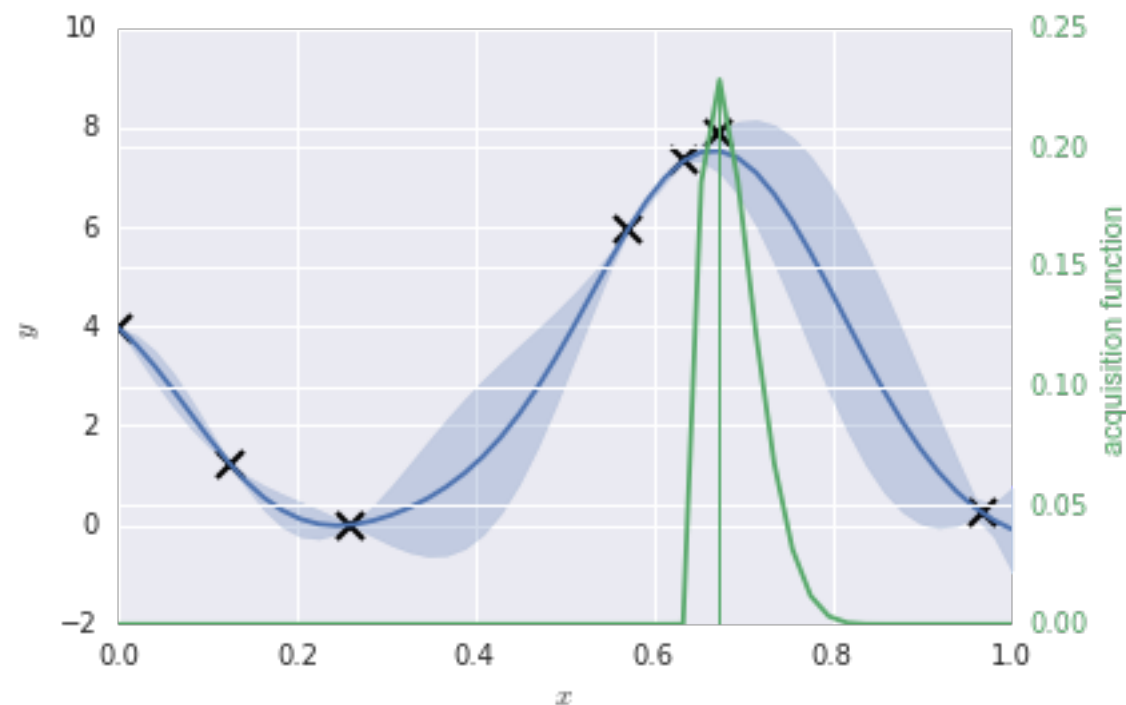
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$



Repeat (Iteration 3)



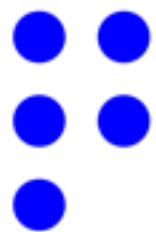
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

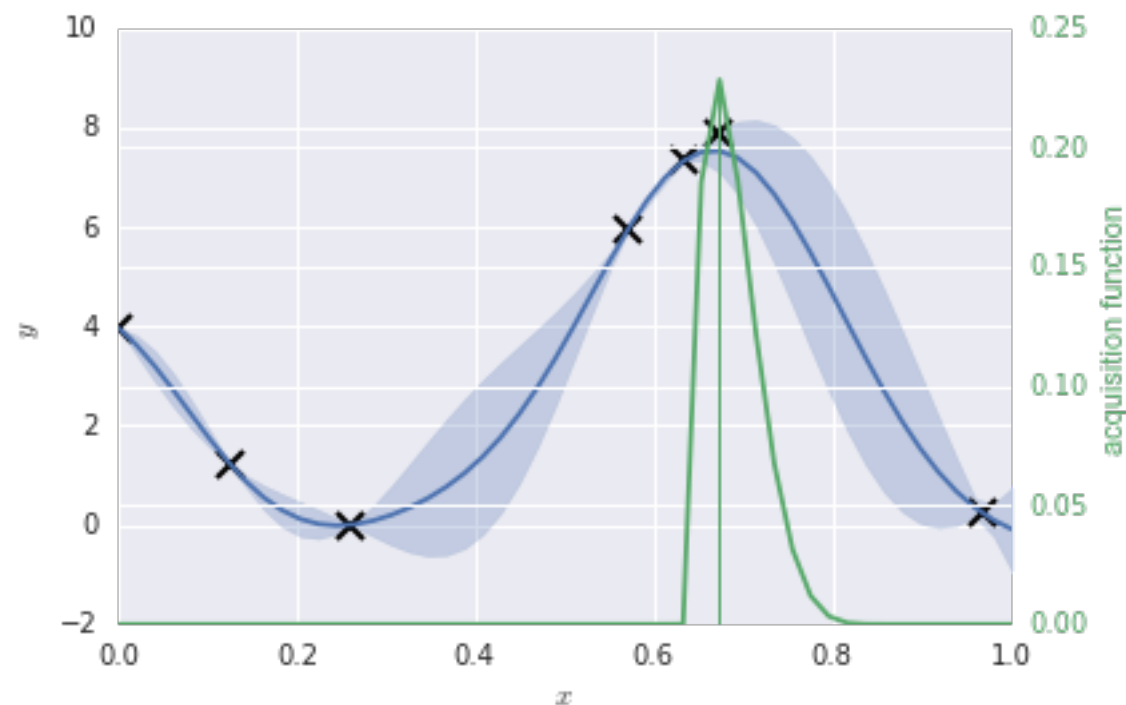
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$



Repeat (Iteration 4)



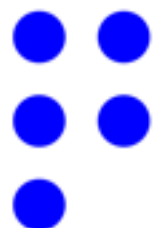
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

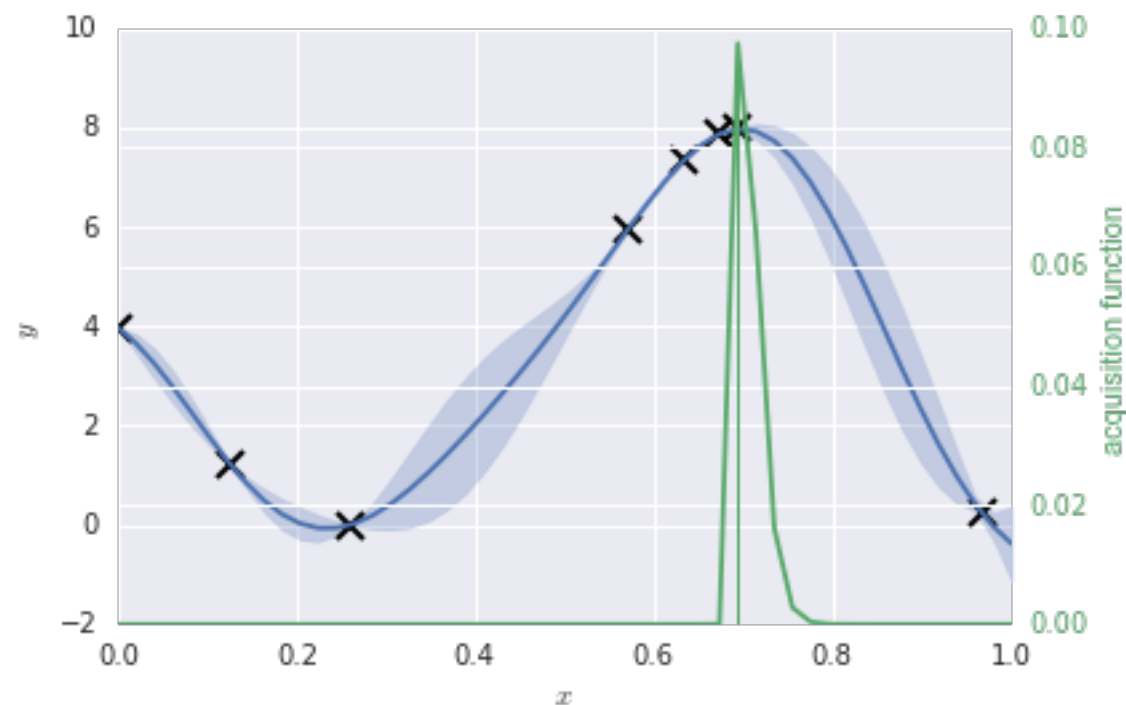
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$



Repeat (Iteration 5)



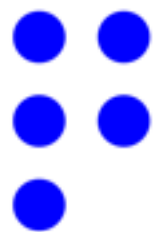
$a(\mathbf{x})$

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$$

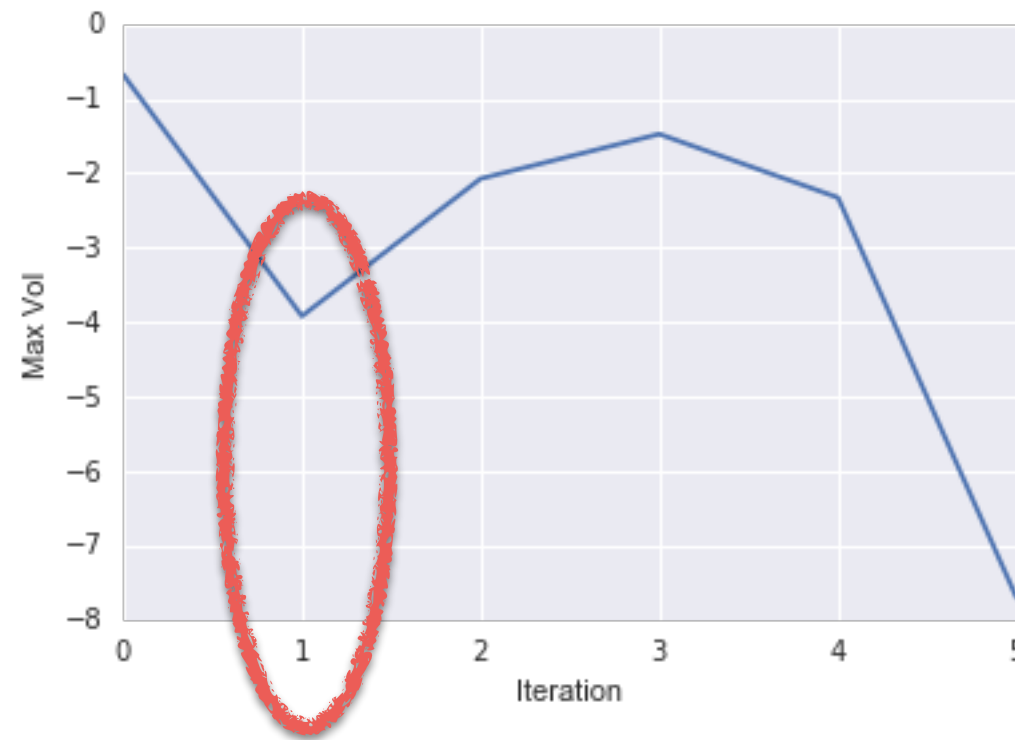
$$y_{1:n} = \{y_1, \dots, y_n\}$$

$$p(y | \mathbf{x}) \approx \mathcal{N}(y | m(\mathbf{x}), \sigma^2(\mathbf{x}))$$

$$\mathbf{x}_{n+1} = \operatorname{argmax}_{\mathbf{x}} a(\mathbf{x})$$

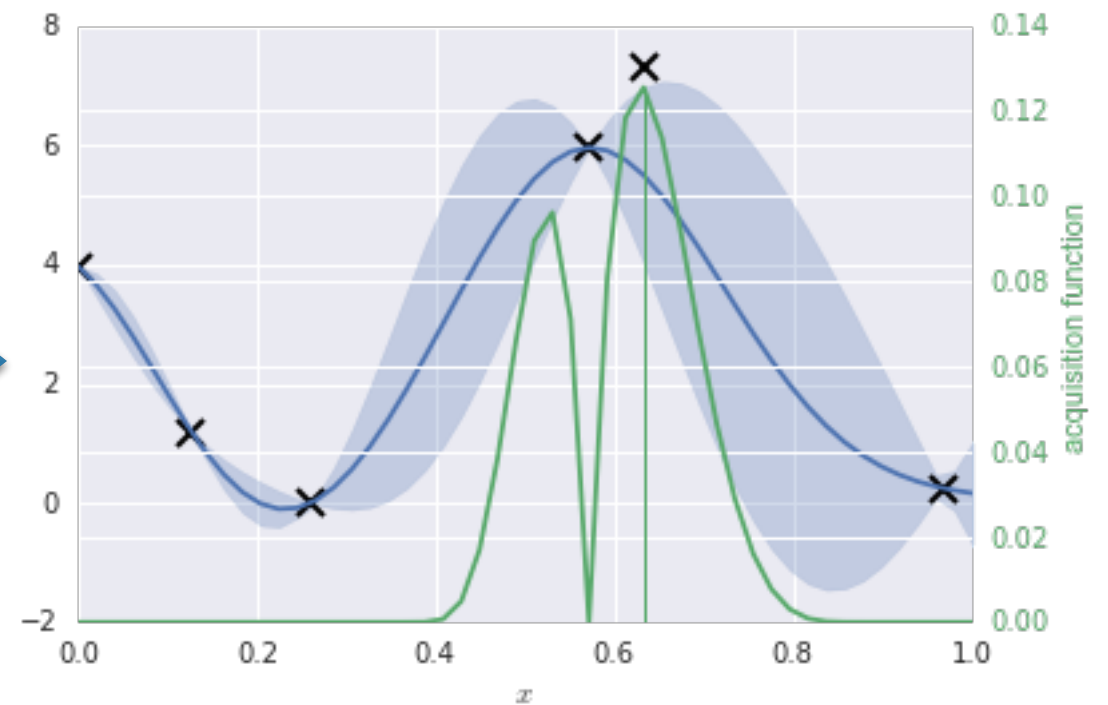
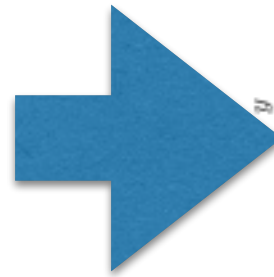
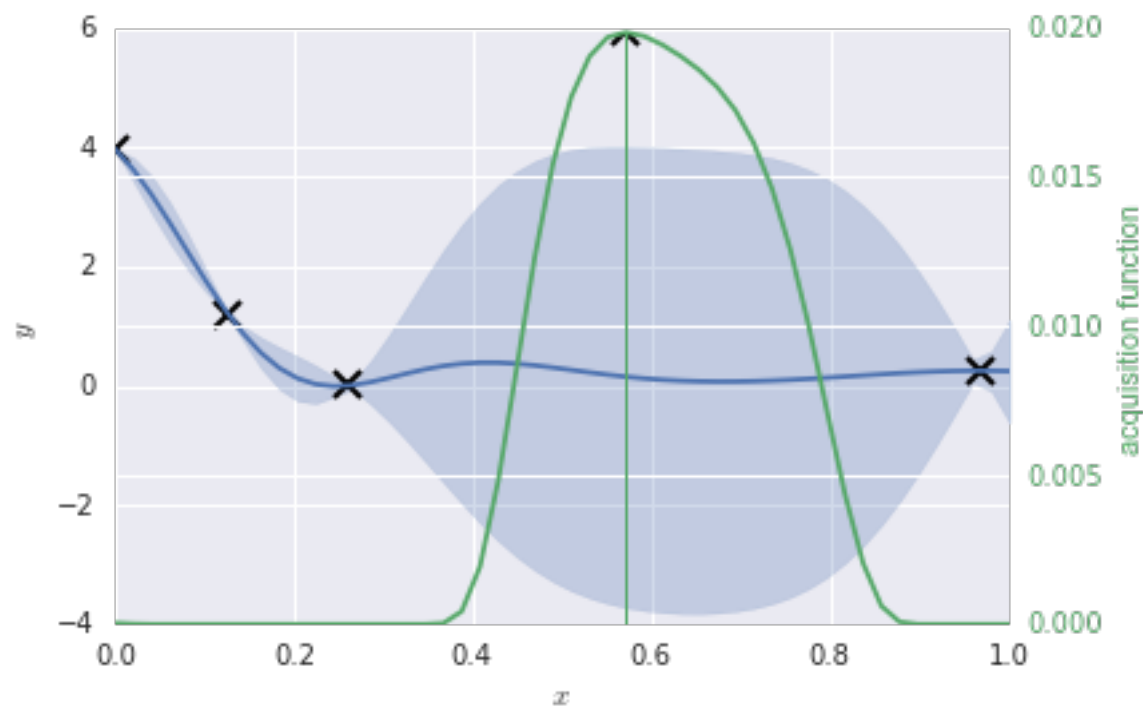


The S-curve Phenomenon of Information



Unexpected information

The S-curve Phenomenon of Information



Advanced Topics for Next Time

- Enforcing optimization constraints.
- Information acquisition for multi-objective optimization.
- Information acquisition under uncertainty.
- Information acquisition from various heterogeneous sources.

Start the notebook of lecture 24