

# Bayesian Optimization

Evgeny Burnaev, Alexey Artemov

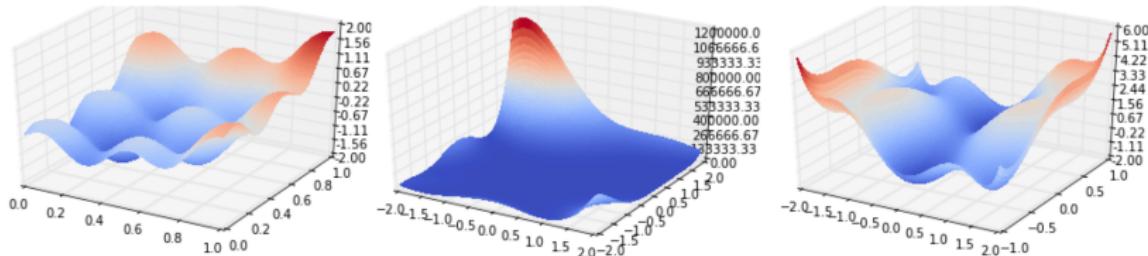
Skoltech, Moscow, Russia

## 1 Bayesian Optimization

# 1 Bayesian Optimization

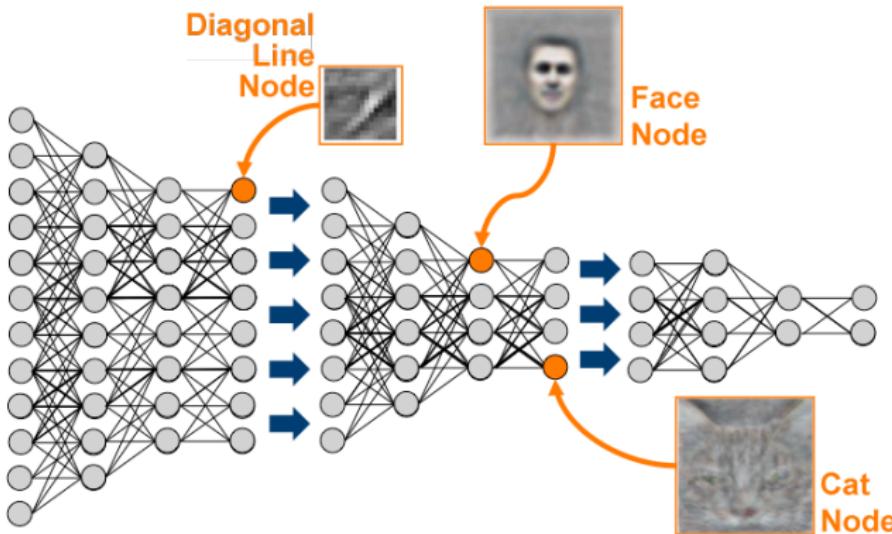
Consider a “well behaved” function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , with  $\mathcal{X} \subseteq \mathbb{R}^d$  being a compact set

$$x_{\min} = \arg \min_{x \in \mathcal{X}} f(x)$$



- $f$  is explicitly unknown and multimodal
- Evaluations of  $f$  may be perturbed
- Evaluations of  $f$  are expensive ⇒
  - Gradient and Hessian are not computable
  - Grid search is not possible

## Parameter tuning in ML algorithms



- Number of layers/units per layer
- Types of each layer
- Regularization coefficients
- Learning rates, etc.

## Parameter tuning in ML algorithms: Example of DNN

Input  $x$ :

- Number of layers/units per layer
- Types of each layer
- Regularization coefficients
- Learning rates, etc.

Output:  $f(x)$

- Deep Neural Network accuracy
- Estimated using cross-validation and/or test set
- Very time-consuming

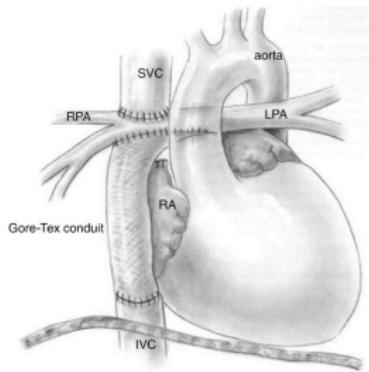


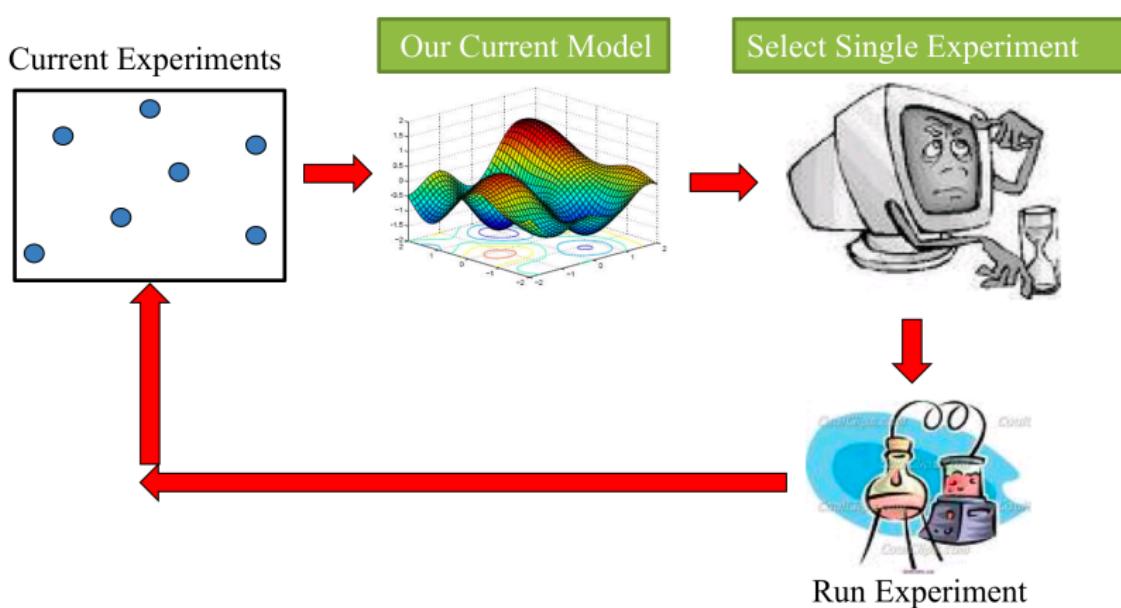
Fig. 1. Extracardiac total cavopulmonary connection. The IVC is disconnected from the right atrium (RA) and connected to the PAs via a Gore-Tex conduit. Figure taken from Reddy et al. [13].

- Design of grafts to be used in heart surgery
- Design of aerodynamic structures, e.g., cars, airplanes
- Calibrating parameters of complex physical models to experimental data
- A/B testing data to optimize the web design to maximize sign-ups, downloads, purchases, etc.

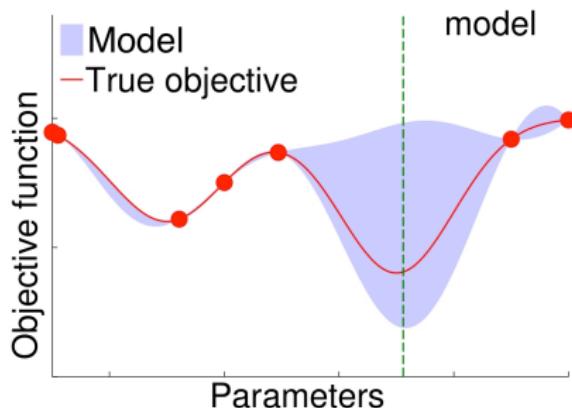
**NB!** There exists commercial services for optimizing black-box functions: SIGOPT, Google Vizier, etc.

# Big Picture

- Since Running experiment is very expensive we use BO
- Select one experiment to run at a time based on results of previous experiments



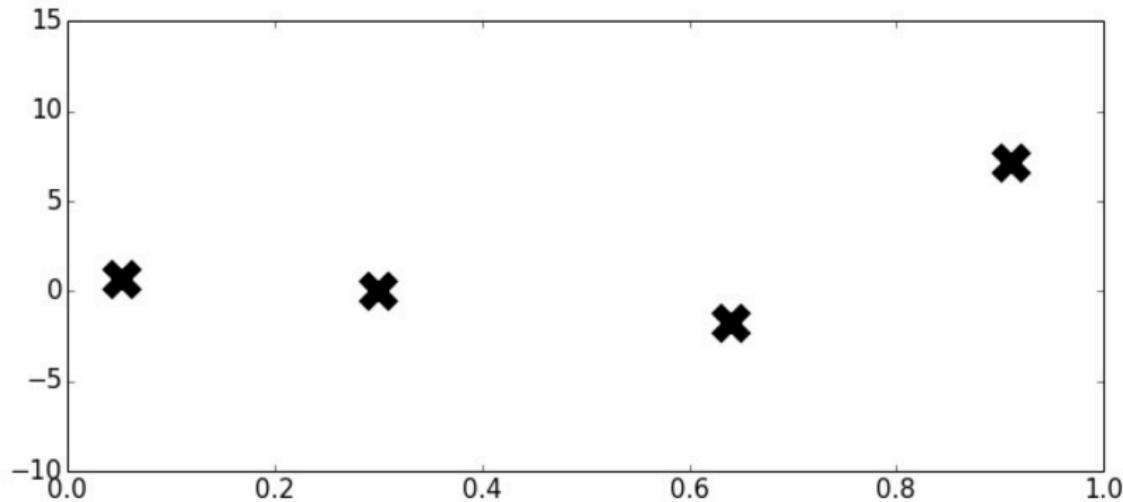
- In high-dimensional case we need many functions evaluations
- Often each evaluation is costly, e.g. in case of experiments



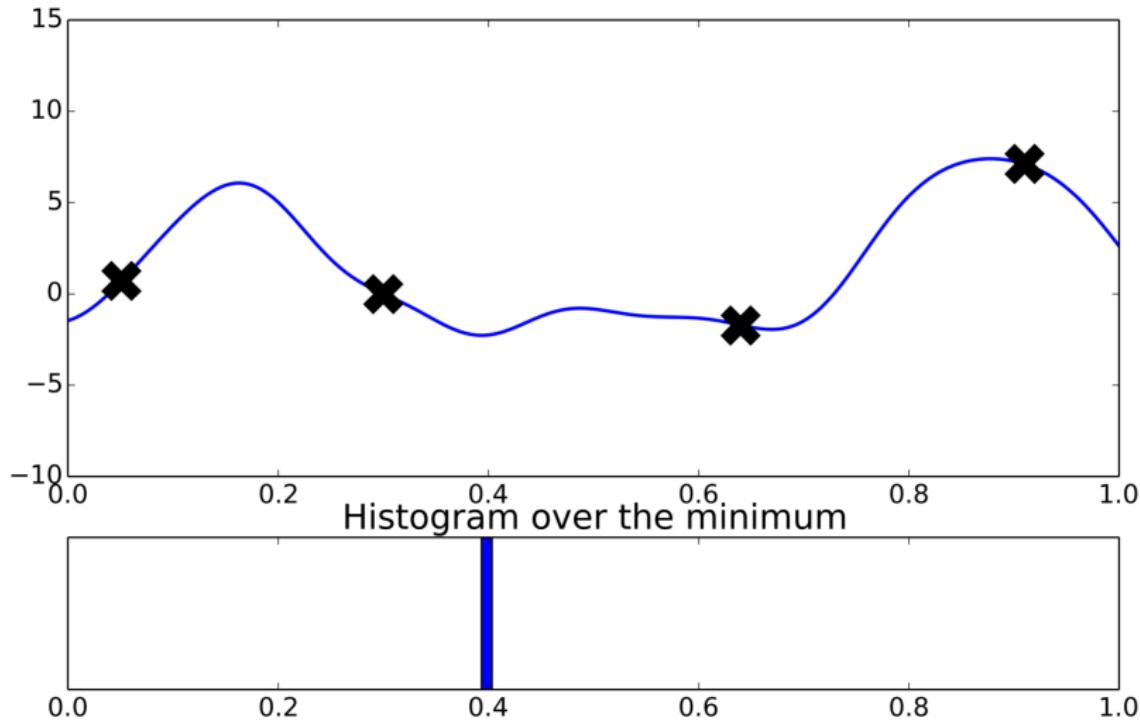
- Error bars are needed to see if a region is still promising

## Typical situation

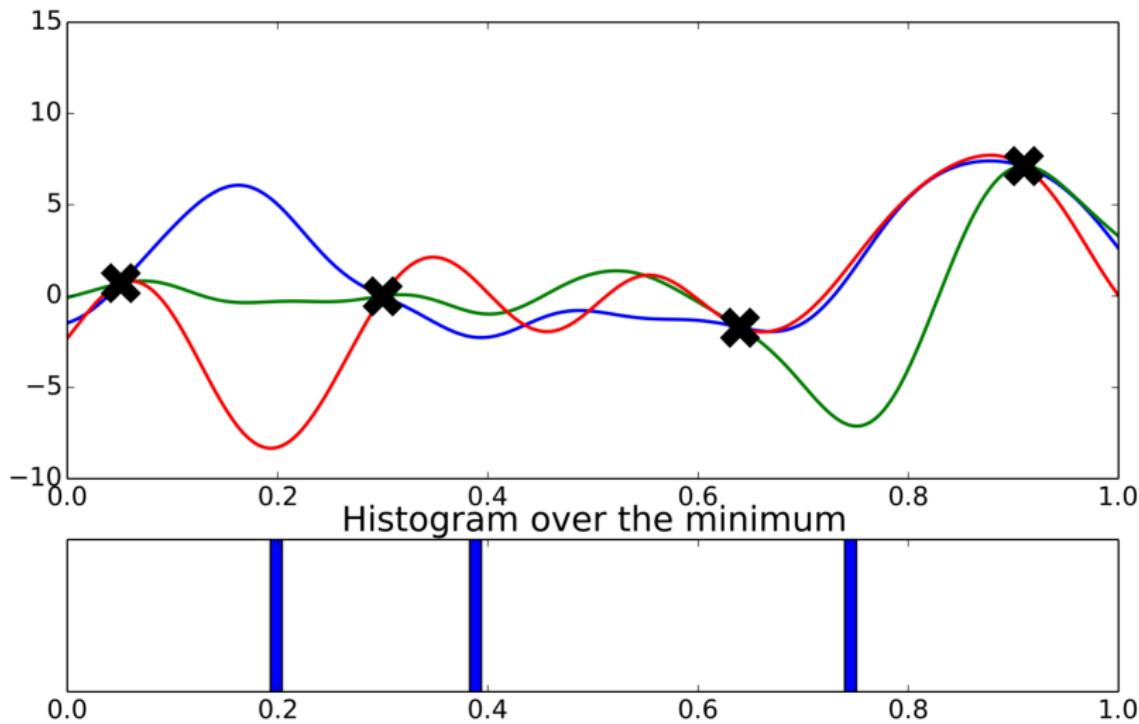
---



Where is the minimum of  $f$ ?  
Where should we evaluate the function next?

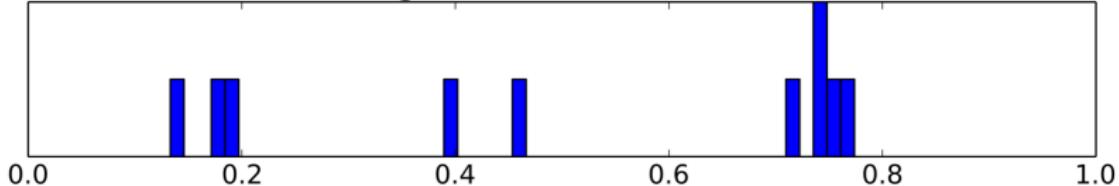
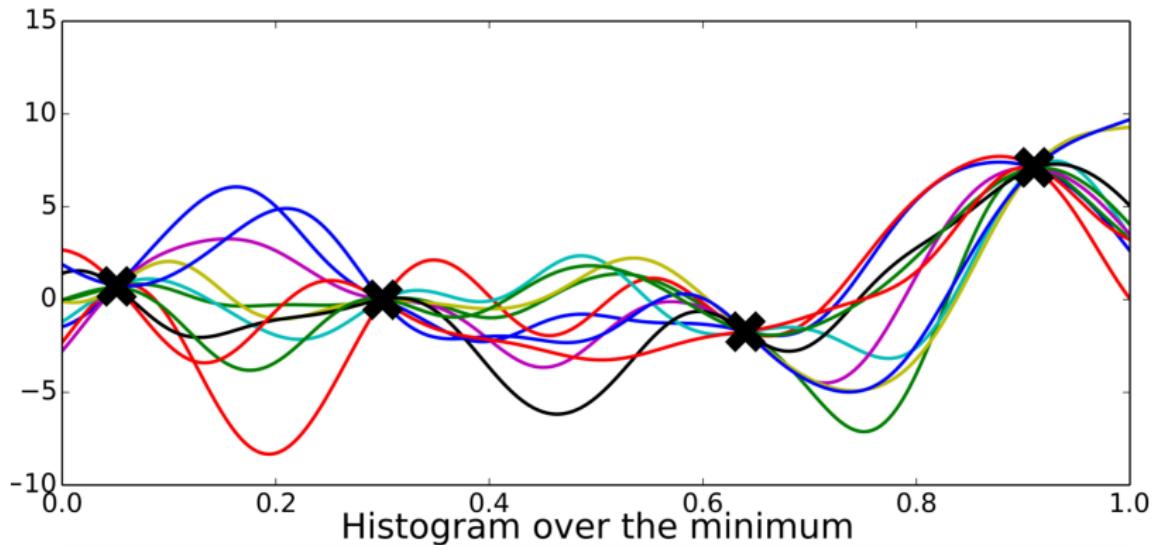


## Intuition: three curves



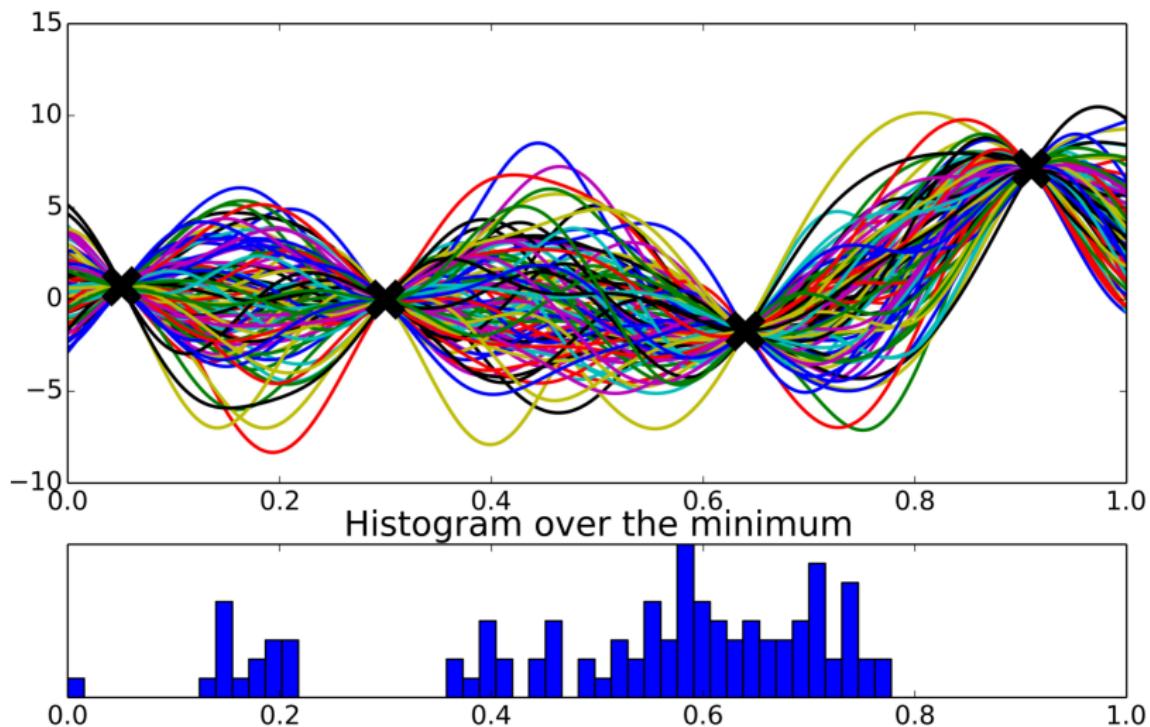
## Intuition: ten curves

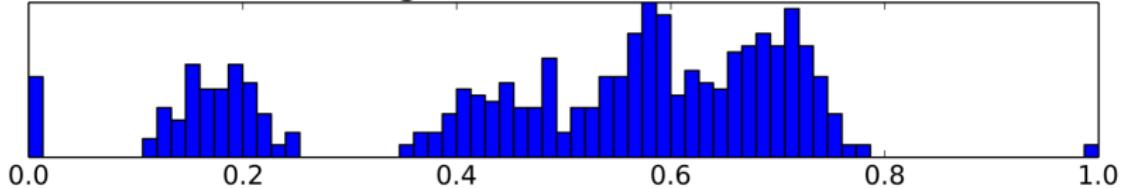
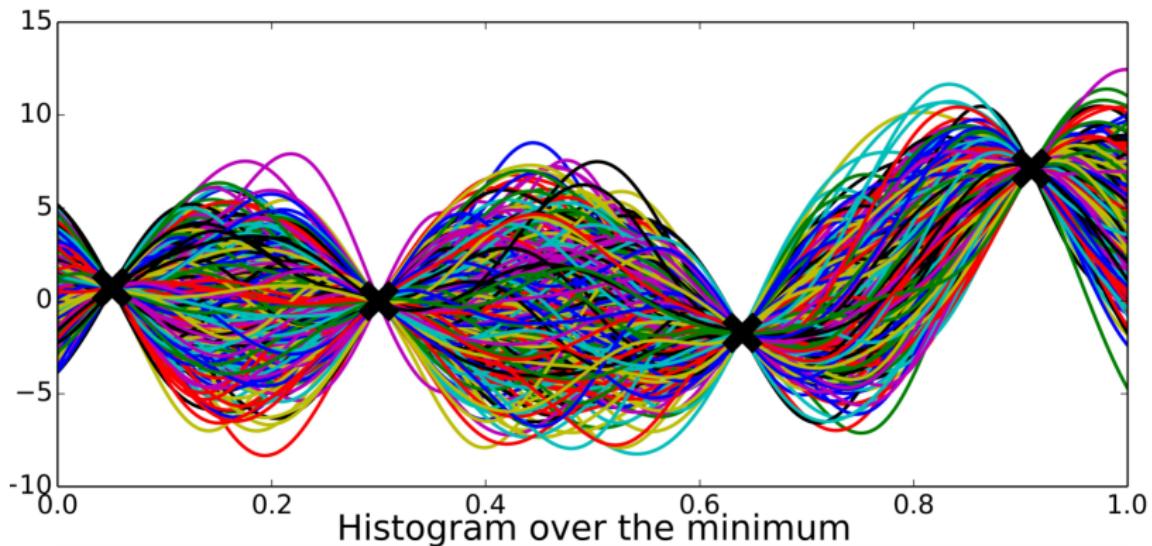
---



# Intuition: hundred curves

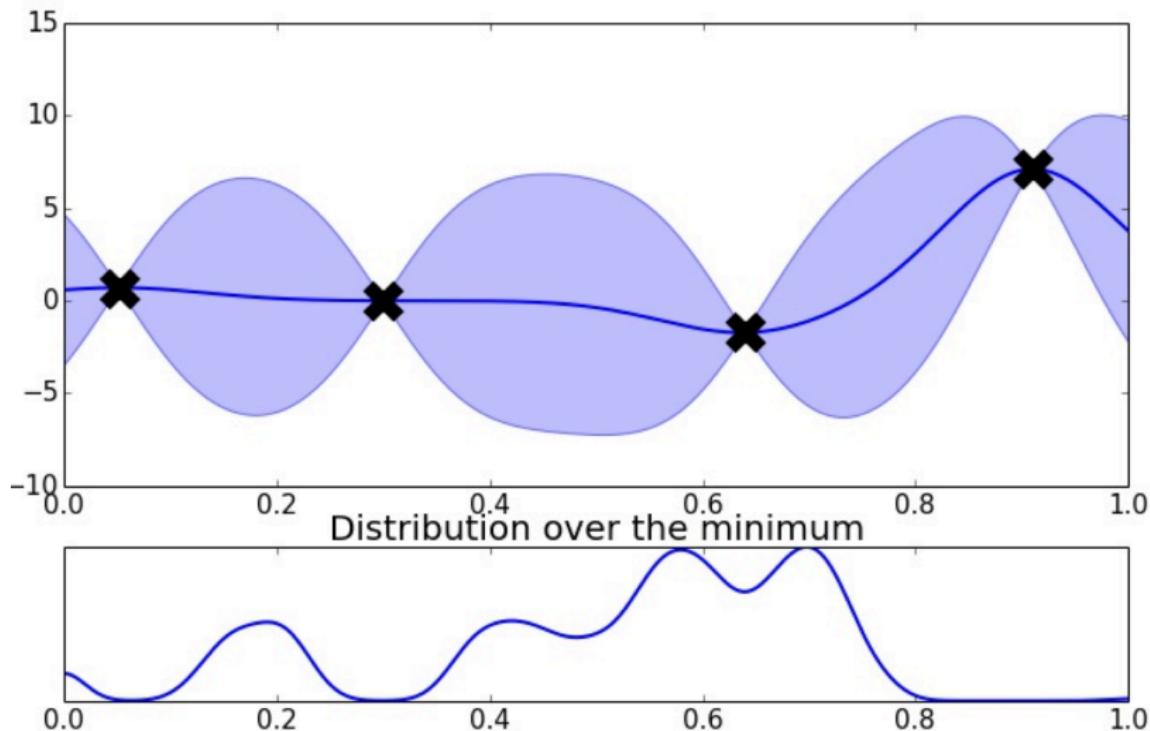
---





## Intuition: infinite number of curves

---



- We made prior assumption about  $f$
- Information about the minimum is now encoded in a new function (the probability distribution  $p_{\min}$  over the minimum in this case)
- We can use  $p_{\min}$  (or a functional of it) to decide where to sample next
- Other functions to encode relevant information about the minimum are possible, e. g. the “marginal expected gain” at each location.

## Using GP as a prior

---

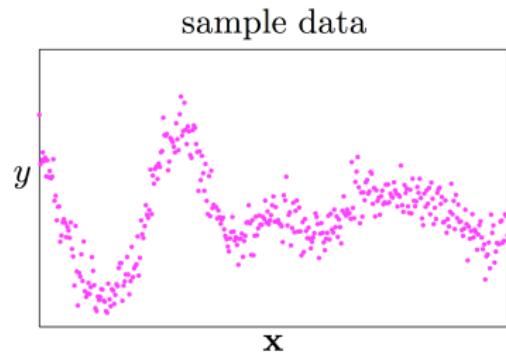
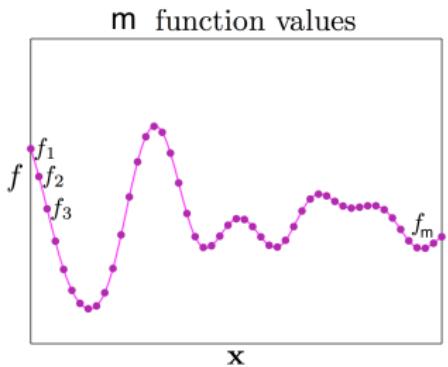
- Training data set  $S_m = \{X, y\} = \{(x_i, y_i)\}_{i=1}^m$

- Model:

$$y_i = f(x_i) + \varepsilon_i,$$

$f \sim \mathcal{GP}(\cdot | 0, K)$ , with  $K(x, x') = \text{cov}(f(x), f(x'))$ ,

$\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$  is a white noise



- Training data set  $S_m = \{X, y\} = \{(x_i, y_i)\}_{i=1}^m$
- Model:

$$y_i = f(x_i) + \varepsilon_i$$

$$f \sim \mathcal{GP}(\cdot | 0, K)$$

$\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$  is a white noise

- The prior is

$$p(f) = \mathcal{N}(f | 0, K)$$

- The noise model, or likelihood is

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

- Integrating over the function variables  $f$  we get the marginal likelihood

$$p(y) = \int p(y|f)p(f)df = \mathcal{N}(y|0, K + \sigma^2 I_m)$$

- Let us denote input test point as  $x_*$ , and output

$$y_* = f_* + \varepsilon_*, \quad f_* = f(x_*)$$

- Consider joint training and test marginal likelihood

$$p(y, f_*) = \mathcal{N} \left( \begin{array}{c} y \\ f(x_*) \end{array} \middle| 0, \begin{bmatrix} K + \sigma^2 I_m & k_* \\ k_*^\top & K_{**} \end{bmatrix} \right),$$

where  $k_* = \{K(x_*, x_i)\}_{i=1}^m$  and  $K_{**} = K(x_*, x_*)$

- What we know about noiseless value  $f(x_*)$ ?

- Joint training and test marginal likelihood

$$p(y, f_*) = \mathcal{N} \left( \begin{bmatrix} y \\ f(x_*) \end{bmatrix} \mid 0, \begin{bmatrix} K + \sigma^2 I_m & k_* \\ k_*^\top & K_{**} \end{bmatrix} \right),$$

where  $k_* = \{K(x_*, x_i)\}_{i=1}^m$  and  $K_{**} = K(x_*, x_*)$

- Condition on training outputs  $y$  we get

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

where

$$\begin{aligned} \mu_*(x_*) &= k_*^\top [K + \sigma^2 I_m]^{-1} y = \\ &= \sum_{i=1}^m \alpha_i K(x_*, x_i) \text{ with } \boldsymbol{\alpha} = [K + \sigma^2 I_m]^{-1} y \text{ (aka KRR)} \end{aligned}$$

$$\sigma_*^2(x_*) = K_{**} - k_*^\top [K + \sigma^2 I_m]^{-1} k_*$$

- $p(y_*|y) = \mathcal{N}(y_* | \mu_*, \sigma_*^2 + \sigma^2)$  predicts what we'll see next

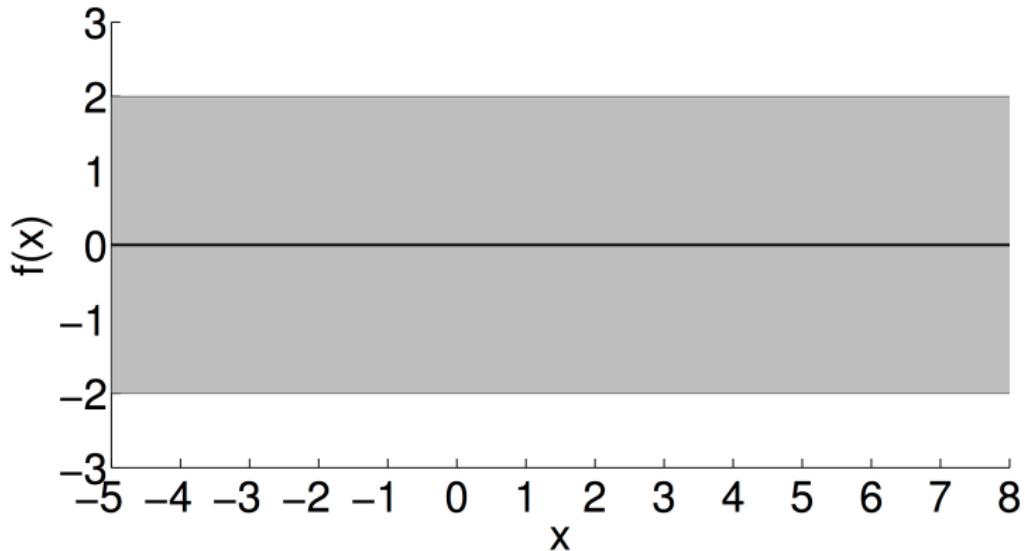


Figure – Prior belief about the function

Predictive (marginal) mean and variance

$$\mathbb{E}[f(x_*)|x_*, \emptyset] = \mu_*(x_*) = 0$$

$$\text{Var}[f(x_*)|x_*, \emptyset] = \sigma_*^2(x_*) = K(x_*, x_*)$$

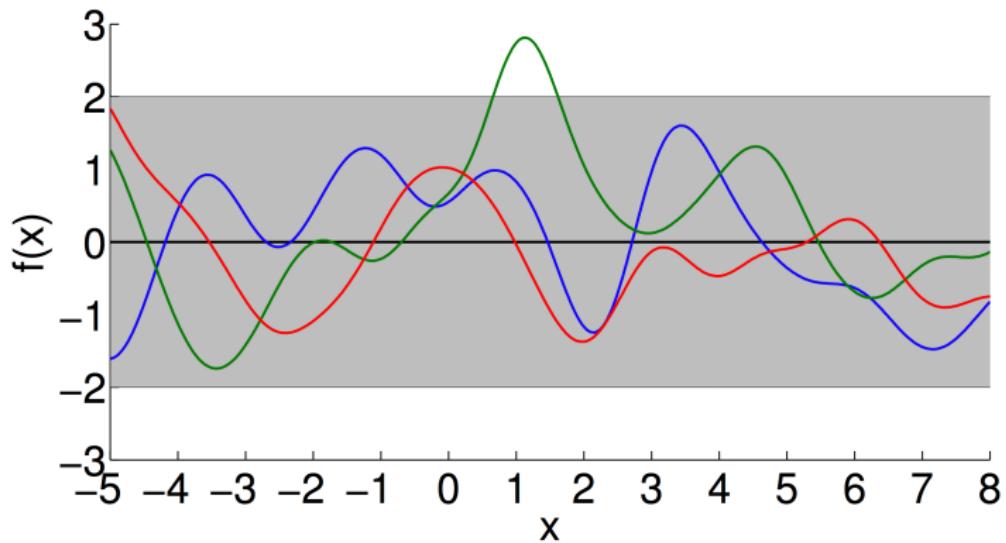


Figure – Prior belief about the function

Predictive (marginal) mean and variance

$$\mathbb{E}[f(x_*)|x_*, \emptyset] = \mu_*(x_*) = 0$$

$$\text{Var}[f(x_*)|x_*, \emptyset] = \sigma_*^2(x_*) = K(x_*, x_*)$$

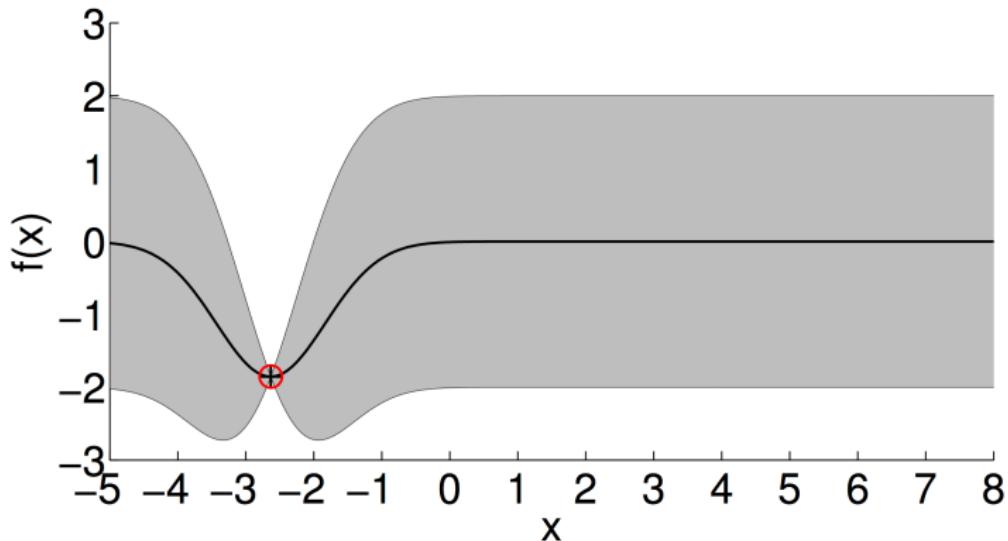


Figure – Posterior belief about the function

Predictive (marginal) mean and variance

$$\mathbb{E}[f(x_*)|x_*, X, y] = \mu_*(x_*) = k_*^\top [K + \sigma^2 I_m]^{-1} y$$

$$\text{Var}[f(x_*)|x_*, X, y] = \sigma_*^2(x_*) = K_{**} - k_*^\top [K + \sigma^2 I_m]^{-1} k_*$$

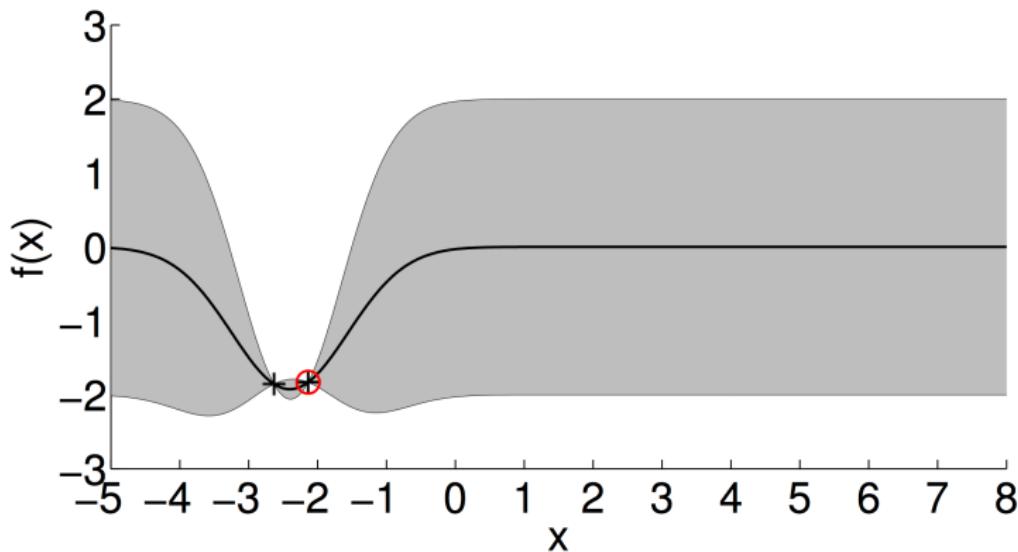


Figure – Posterior belief about the function

Predictive (marginal) mean and variance

$$\mathbb{E}[f(x_*)|x_*, X, y] = \mu_*(x_*) = k_*^\top [K + \sigma^2 I_m]^{-1} y$$

$$\text{Var}[f(x_*)|x_*, X, y] = \sigma_*^2(x_*) = K_{**} - k_*^\top [K + \sigma^2 I_m]^{-1} k_*$$

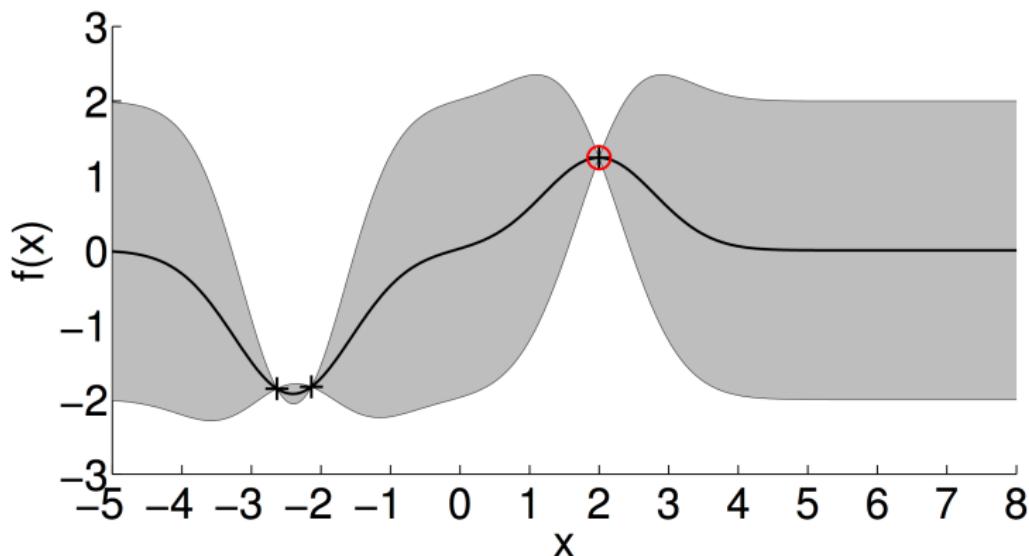


Figure – Posterior belief about the function

Predictive (marginal) mean and variance

$$\mathbb{E}[f(x_*)|x_*, X, y] = \mu_*(x_*) = k_*^\top [K + \sigma^2 I_m]^{-1} y$$

$$\text{Var}[f(x_*)|x_*, X, y] = \sigma_*^2(x_*) = K_{**} - k_*^\top [K + \sigma^2 I_m]^{-1} k_*$$

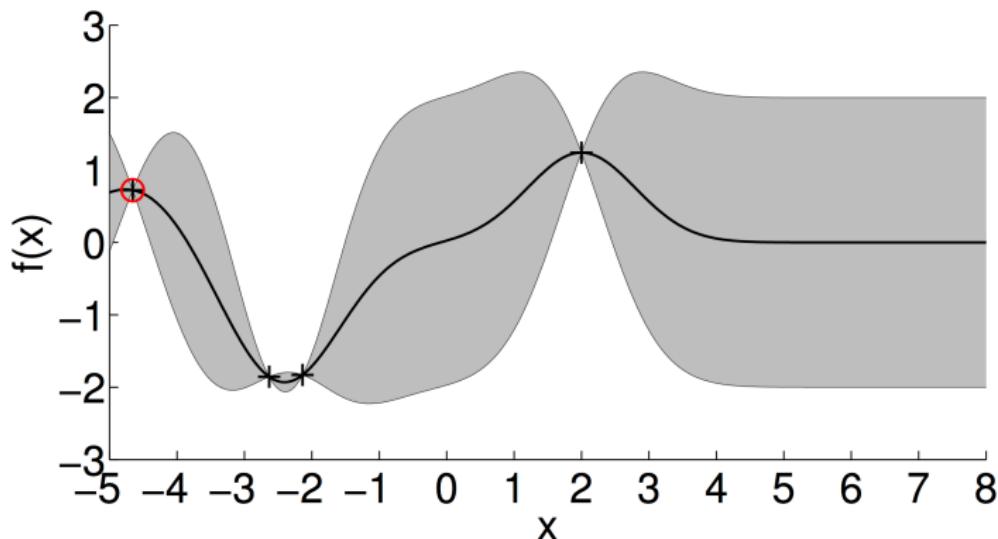


Figure – Posterior belief about the function

Predictive (marginal) mean and variance

$$\mathbb{E}[f(x_*)|x_*, X, y] = \mu_*(x_*) = k_*^\top [K + \sigma^2 I_m]^{-1} y$$

$$\text{Var}[f(x_*)|x_*, X, y] = \sigma_*^2(x_*) = K_{**} - k_*^\top [K + \sigma^2 I_m]^{-1} k_*$$

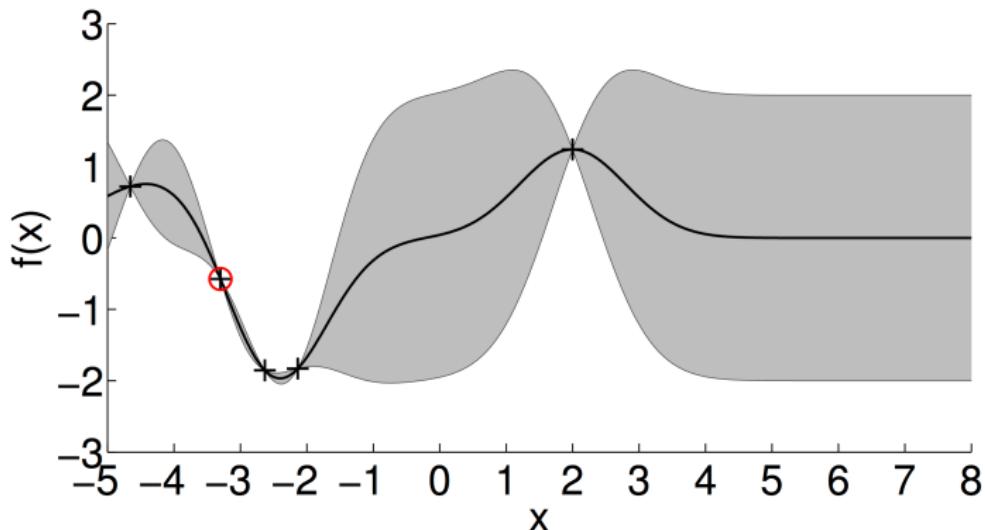


Figure – Posterior belief about the function

Predictive (marginal) mean and variance

$$\mathbb{E}[f(x_*)|x_*, X, y] = \mu_*(x_*) = k_*^\top [K + \sigma^2 I_m]^{-1} y$$

$$\text{Var}[f(x_*)|x_*, X, y] = \sigma_*^2(x_*) = K_{**} - k_*^\top [K + \sigma^2 I_m]^{-1} k_*$$

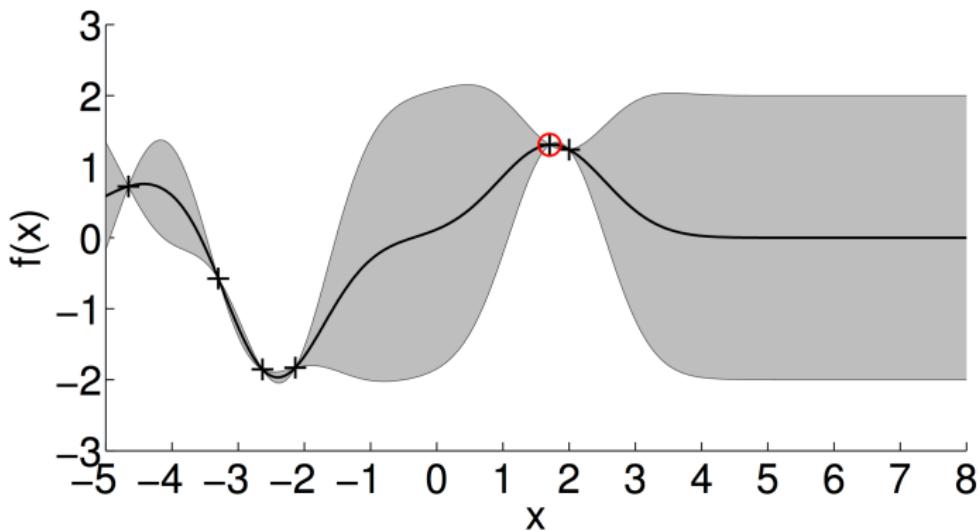


Figure – Posterior belief about the function

Predictive (marginal) mean and variance

$$\mathbb{E}[f(x_*)|x_*, X, y] = \mu_*(x_*) = k_*^\top [K + \sigma^2 I_m]^{-1} y$$

$$\text{Var}[f(x_*)|x_*, X, y] = \sigma_*^2(x_*) = K_{**} - k_*^\top [K + \sigma^2 I_m]^{-1} k_*$$

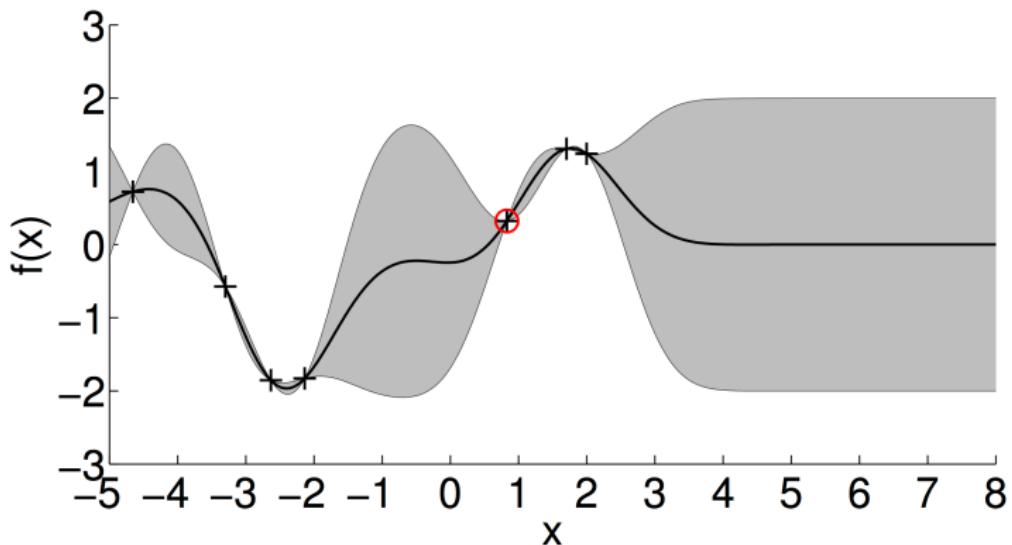


Figure – Posterior belief about the function

Predictive (marginal) mean and variance

$$\mathbb{E}[f(x_*)|x_*, X, y] = \mu_*(x_*) = k_*^\top [K + \sigma^2 I_m]^{-1} y$$

$$\text{Var}[f(x_*)|x_*, X, y] = \sigma_*^2(x_*) = K_{**} - k_*^\top [K + \sigma^2 I_m]^{-1} k_*$$

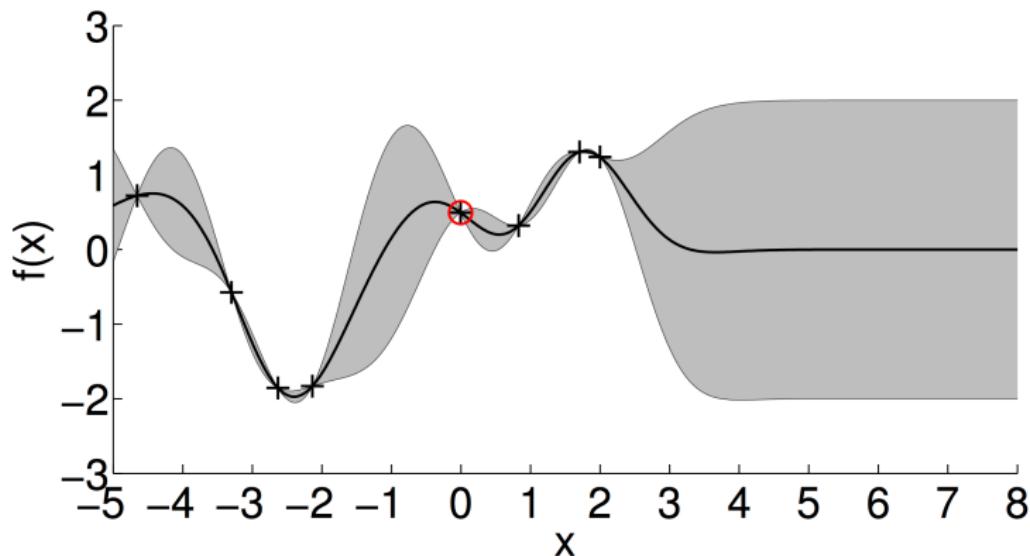


Figure – Posterior belief about the function

Predictive (marginal) mean and variance

$$\mathbb{E}[f(x_*)|x_*, X, y] = \mu_*(x_*) = k_*^\top [K + \sigma^2 I_m]^{-1} y$$

$$\text{Var}[f(x_*)|x_*, X, y] = \sigma_*^2(x_*) = K_{**} - k_*^\top [K + \sigma^2 I_m]^{-1} k_*$$

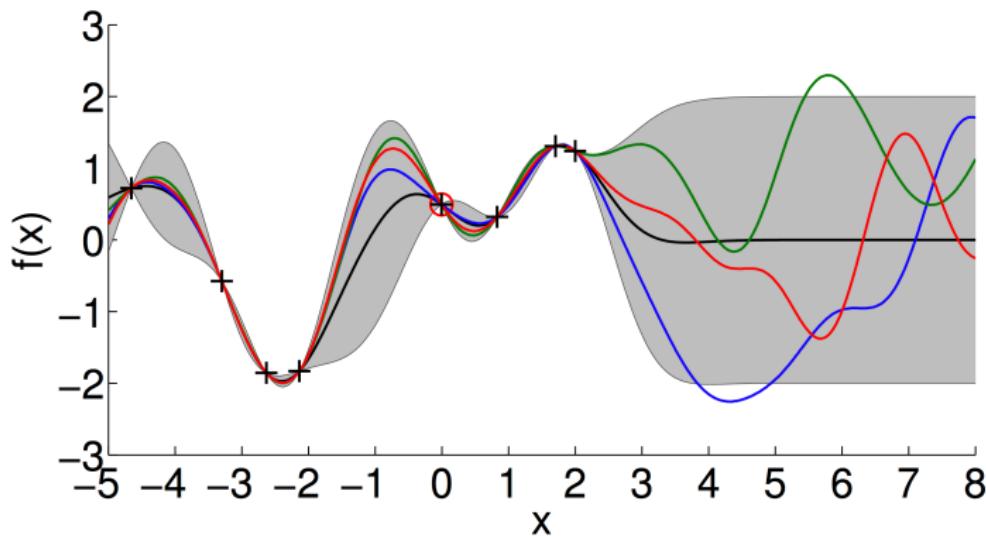


Figure – Posterior belief about the function

Predictive (marginal) mean and variance

$$\mathbb{E}[f(x_*)|x_*, X, y] = \mu_*(x_*) = k_*^\top [K + \sigma^2 I_m]^{-1} y$$

$$\text{Var}[f(x_*)|x_*, X, y] = \sigma_*^2(x_*) = K_{**} - k_*^\top [K + \sigma^2 I_m]^{-1} k_*$$

Methodology to perform global optimization of multimodal black-box functions

1. Choose some **prior measure** over the space of possible objectives  $f$
2. Combine prior and the likelihood to get a **posterior** over the objective given some observations
3. Use the posterior to decide where to take the next evaluation according to some **acquisition function**
4. Augment the data set
5. Iterate between 2 and 4 until the evaluation budget is over

**Comment:** BO can be theoretically formalized in the framework of dynamic programming principle

- Use GP  $\mathcal{GP}(\cdot | \mu(x), K(x, x'))$  as a prior for  $f(\cdot)$
- GP has marginal closed-form for the posterior mean  $\mu_*(x)$  and variance  $\sigma_*^2(x) \Rightarrow$  efficient calculation of acquisition function
  - **Exploration:** Evaluate in places where the variance is large
  - **Exploitation:** Evaluate in places where the mean is low

Acquisition functions balance these two factors to determine where to evaluate next

- BO is an strategy to transform the problem

$$x_{\min} = \arg \min_{x \in \mathcal{X}} f(x)$$

unsolvable!

into a series of problems

$$x_{t+1} = \arg \max_{x \in \mathcal{X}} \alpha(x|S_t),$$

solvable!

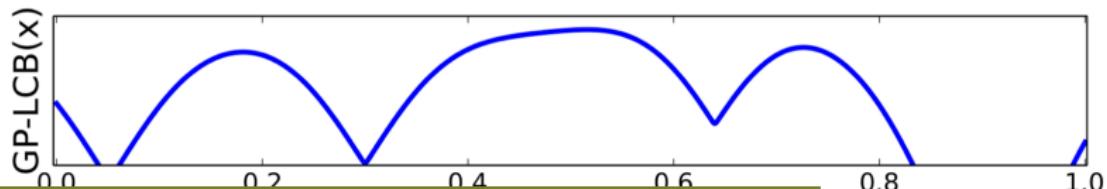
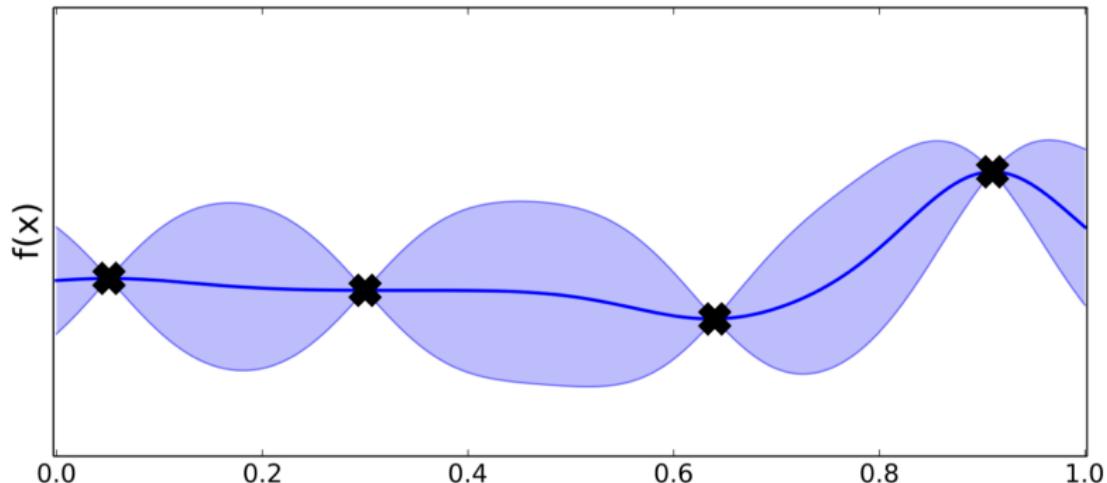
where

- $\alpha(x)$  is not so expensive to evaluate
- Gradients of  $\alpha(x)$  are typically available
- Still need to find  $x_{t+1}$ : DIRECT, gradient methods, SA

## GP Upper (lower) Confidence Band

Direct balance between exploration and exploitation ( $\zeta$  is a user-defined parameter):

$$\alpha_{LCB}(x) = -\mu_*(x) + \zeta \cdot \sigma_*(x)$$

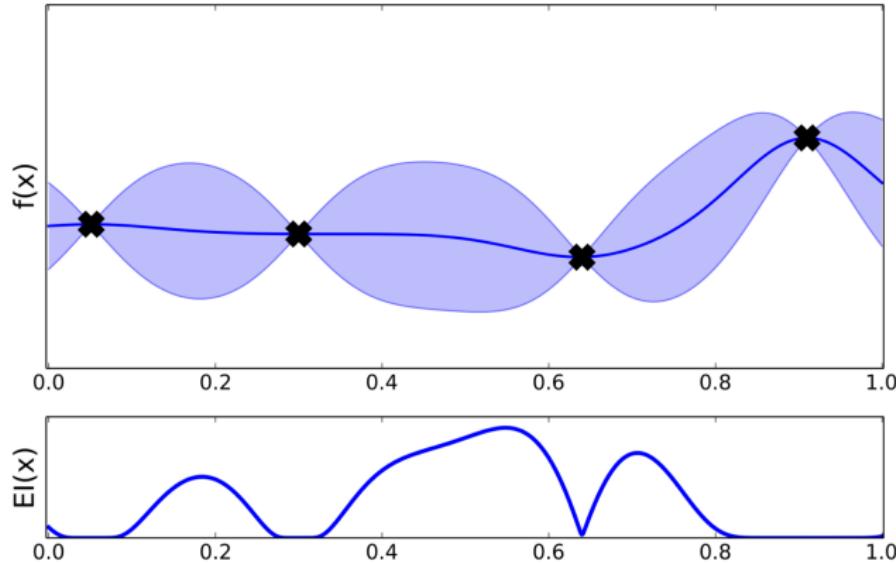


## Expected Improvement

---

Let us denote by  $\Delta(x) = y_{\text{best}} - \mu_*(x)$ , then

$$\begin{aligned}\alpha_{\text{EI}}(x) &= \int_y \max(0, y_{\text{best}} - y_*) p(y_*|x) dy_* = \\ &= \max(0, \Delta(x)) - \sigma_*(x) \varphi \left( \frac{\Delta(x)}{\sigma_*(x)} \right) + |\Delta(x)| \Phi \left( -\frac{|\Delta(x)|}{\sigma_*(x)} \right)\end{aligned}$$

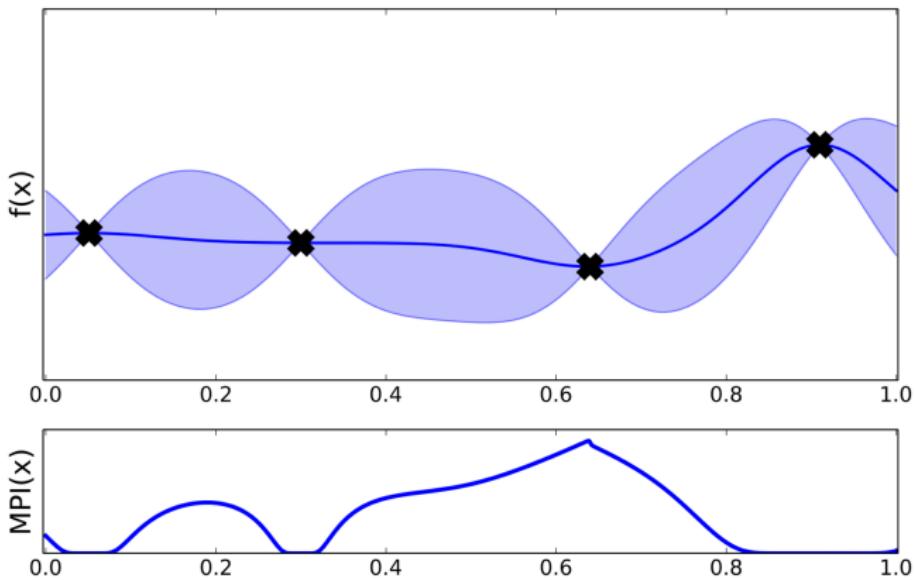


# Maximum Probability of Improvement

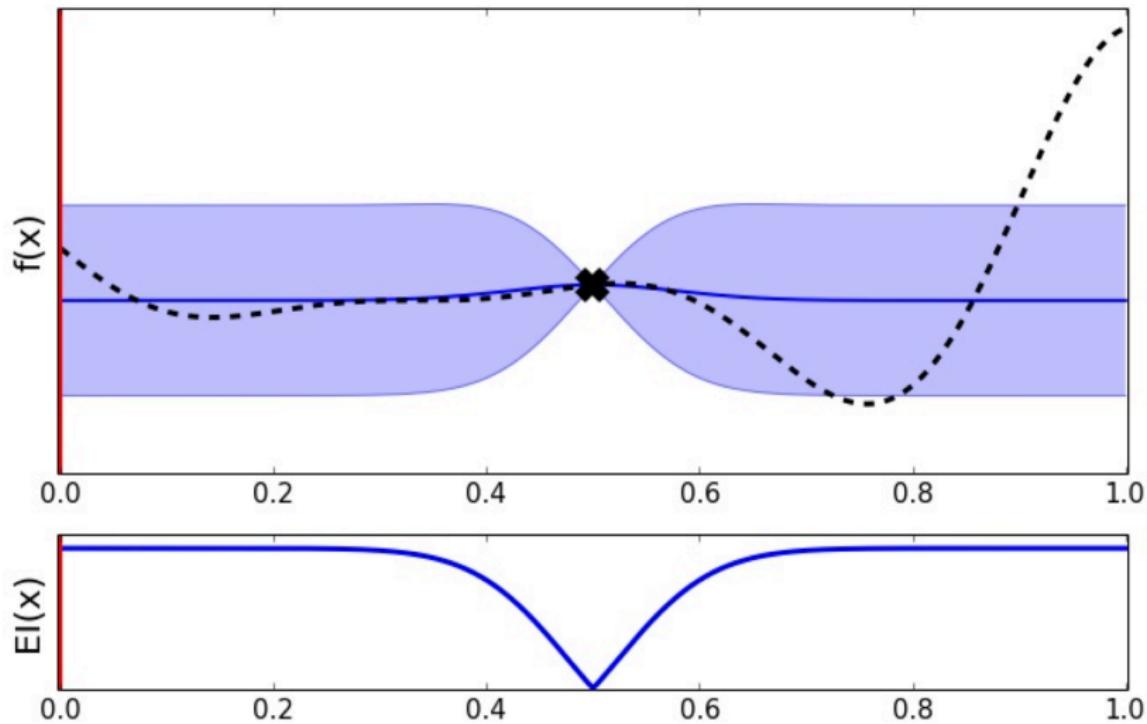
---

$$\gamma(x) = \frac{\mu(x) - y_{\text{best}}}{\sigma(x)}$$

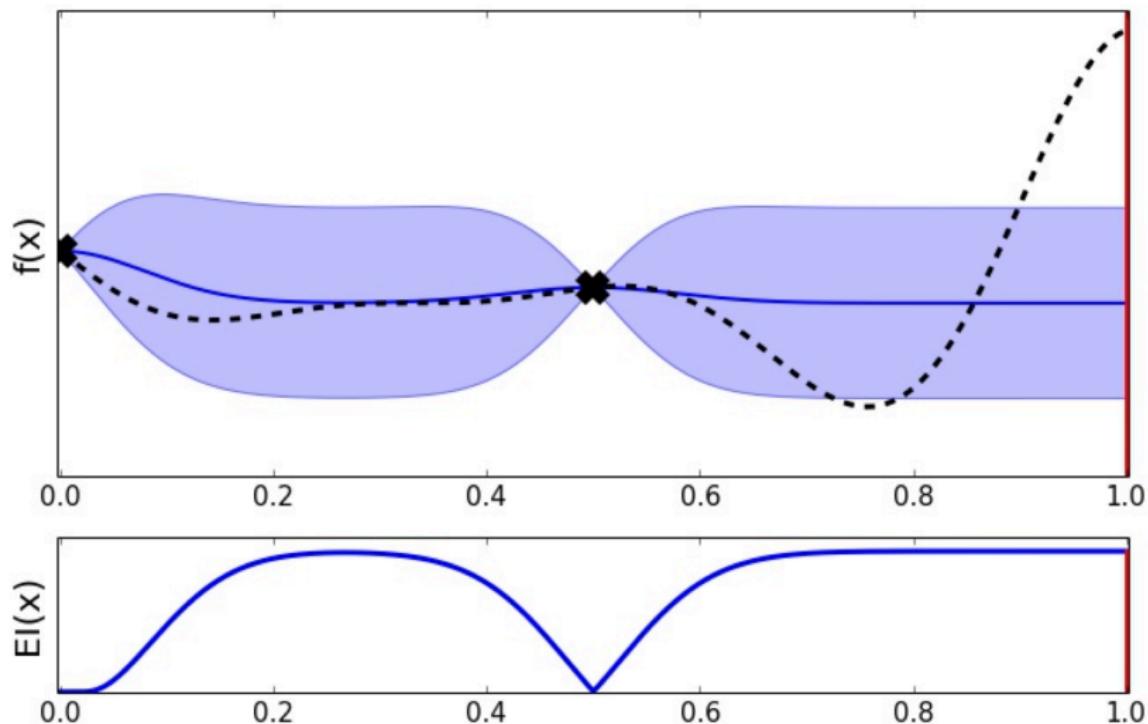
$$\alpha_{\text{MPI}}(x) = \mathbb{P}(f(x) < y_{\text{best}}) = \Phi(\gamma(x))$$



## Expected Improvement: Toy Problem

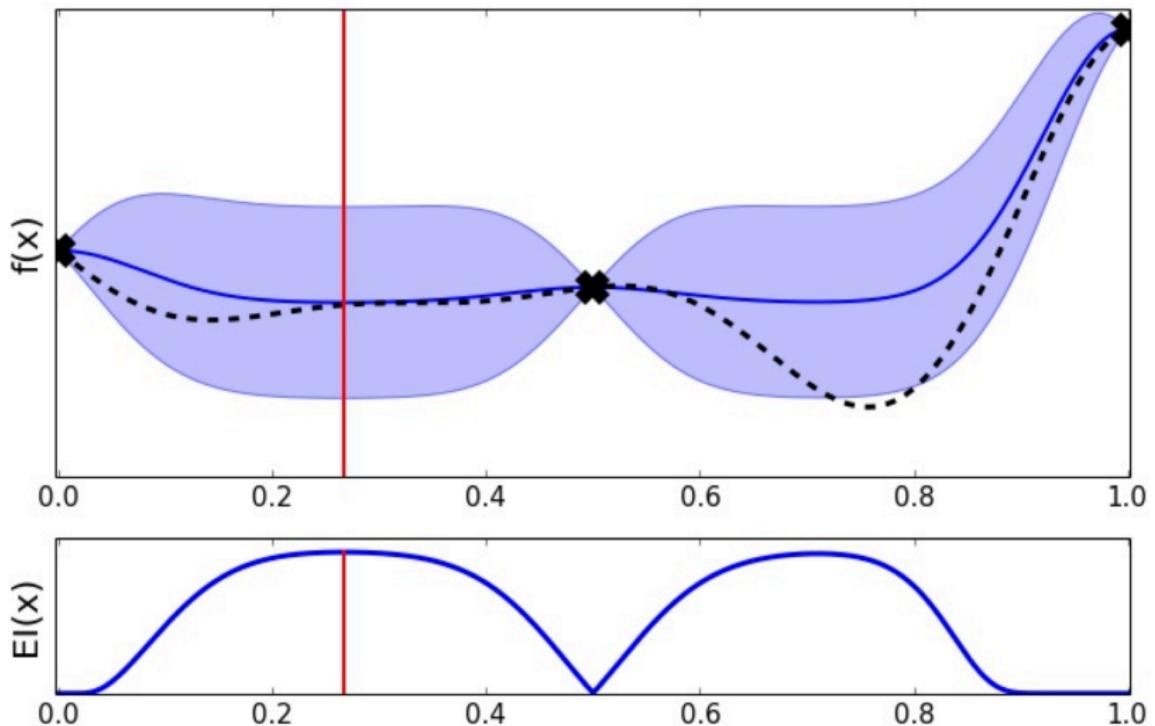


## Expected Improvement: Toy Problem



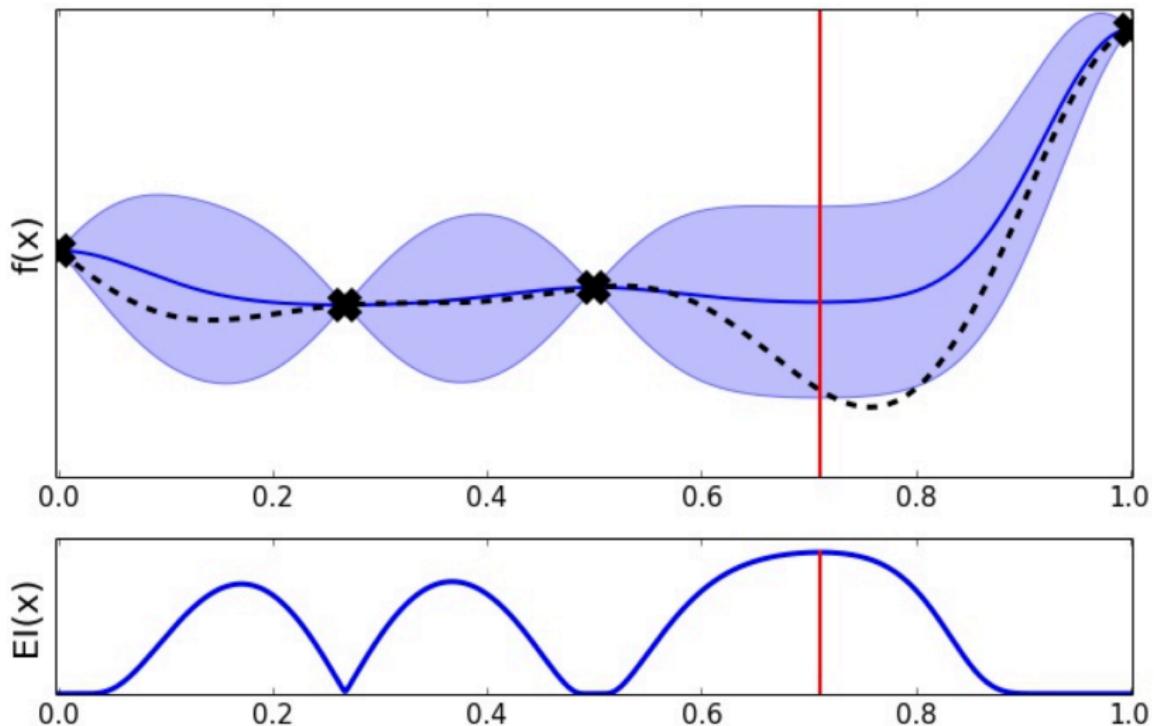
## Expected Improvement: Toy Problem

---



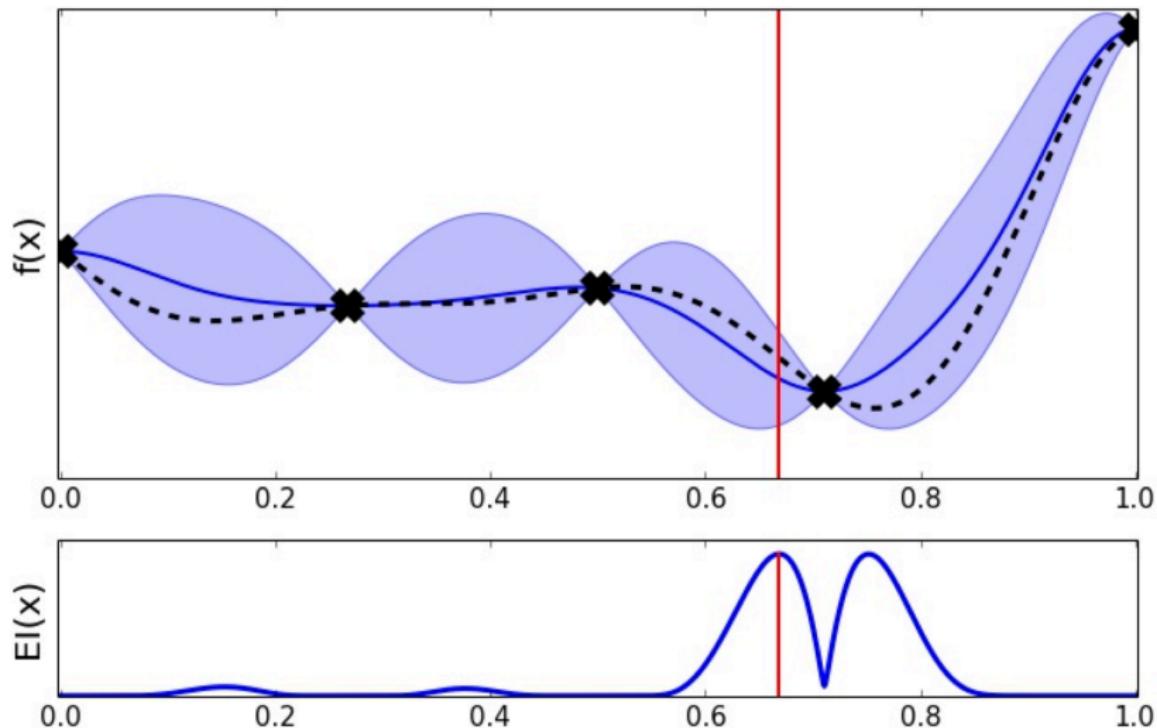
## Expected Improvement: Toy Problem

---



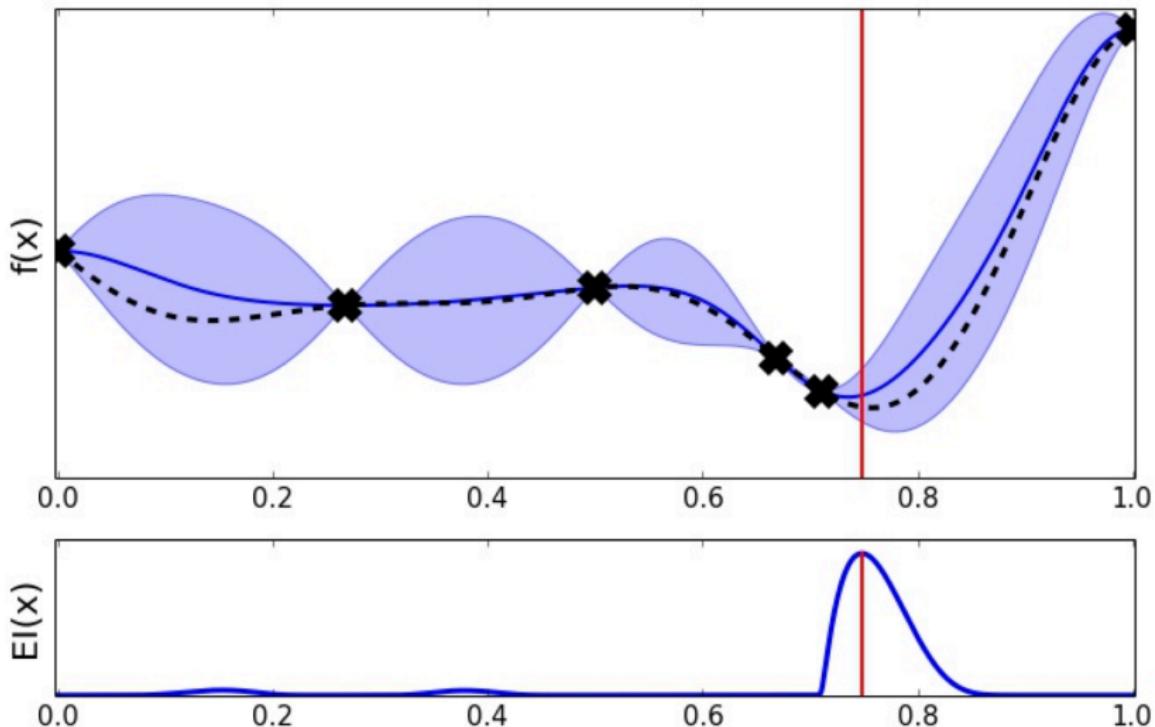
## Expected Improvement: Toy Problem

---



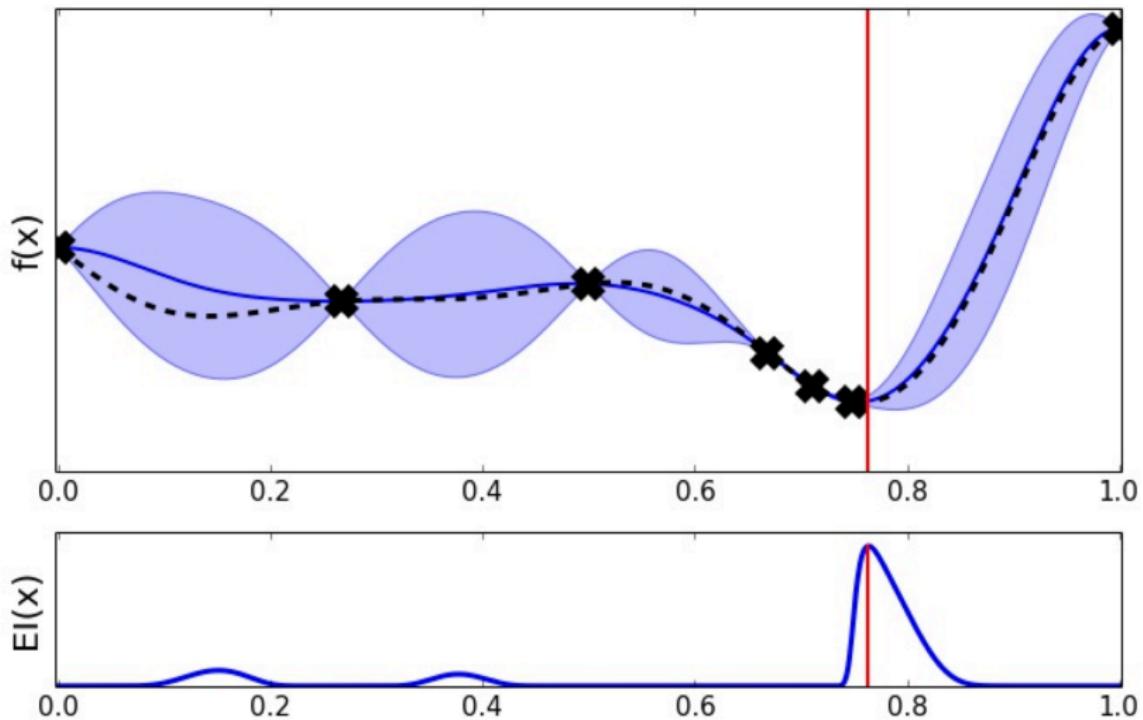
## Expected Improvement: Toy Problem

---



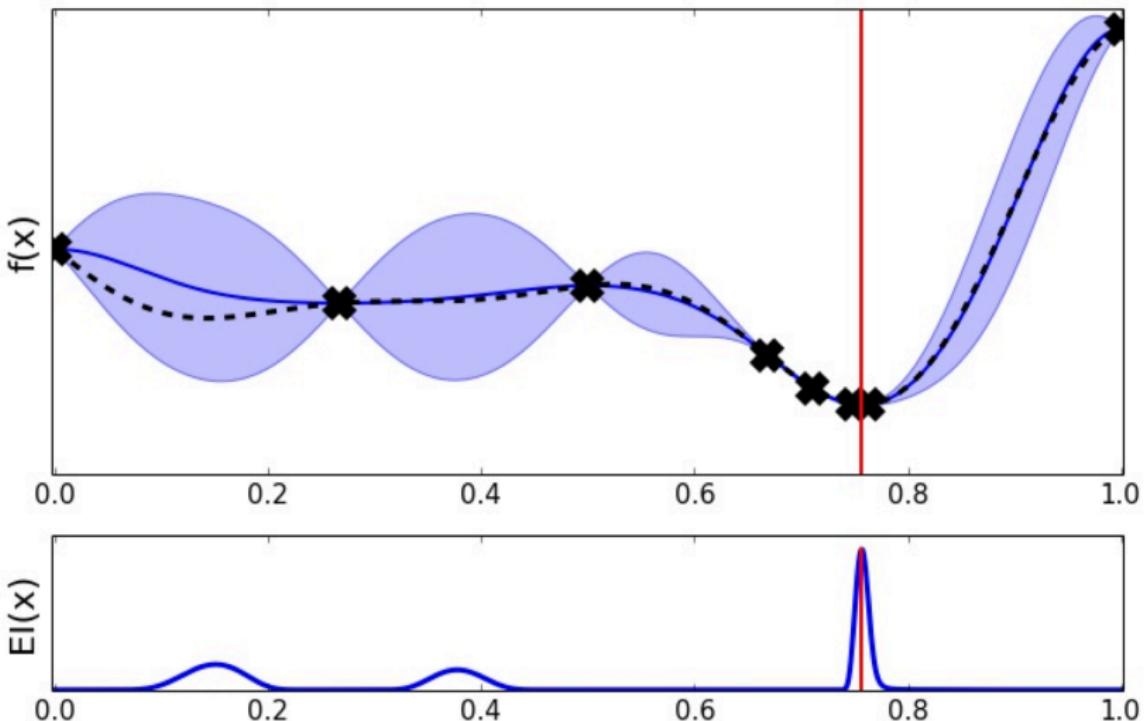
## Expected Improvement: Toy Problem

---



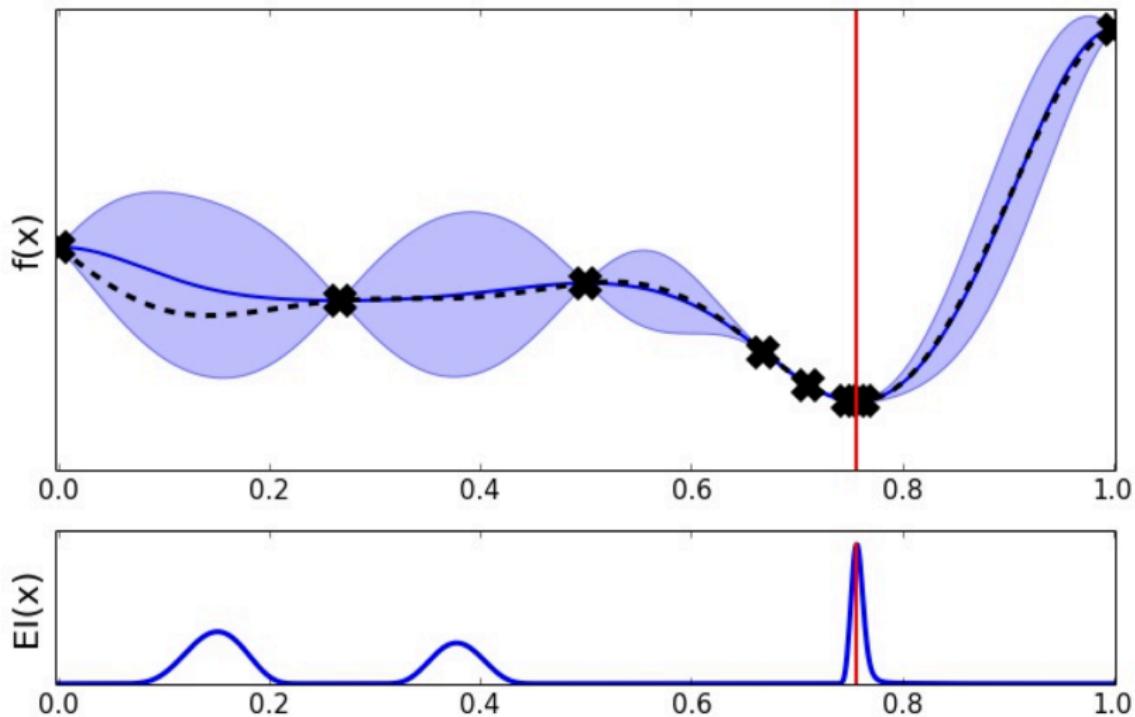
## Expected Improvement: Toy Problem

---



## Expected Improvement: Toy Problem

---



- Scikit-learn dummy data for testing classifiers:  
 $n\_samples=2500$ ,  $n\_features=45$ ,  $n\_informative=15$ ,  
 $n\_redundant=5$
- Optimize w.r.t.  $x = (C, \gamma)$ , where  $C$  — penalization, and  $\gamma$  — kernel width
- The target function  $\mathcal{L}(x)$  is the AUC based on three fold cross-validation

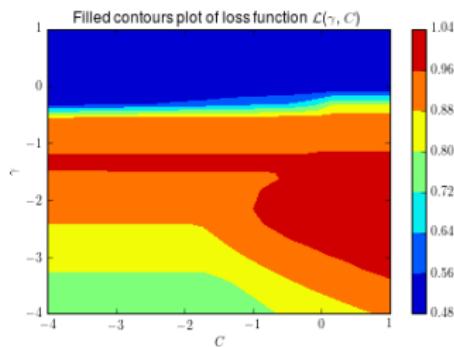
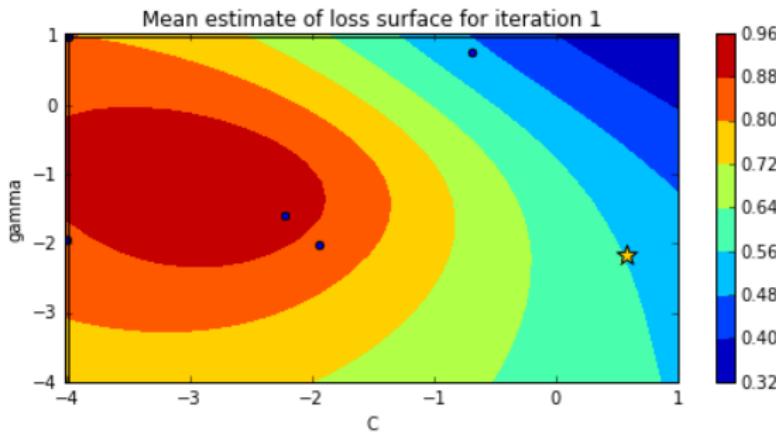
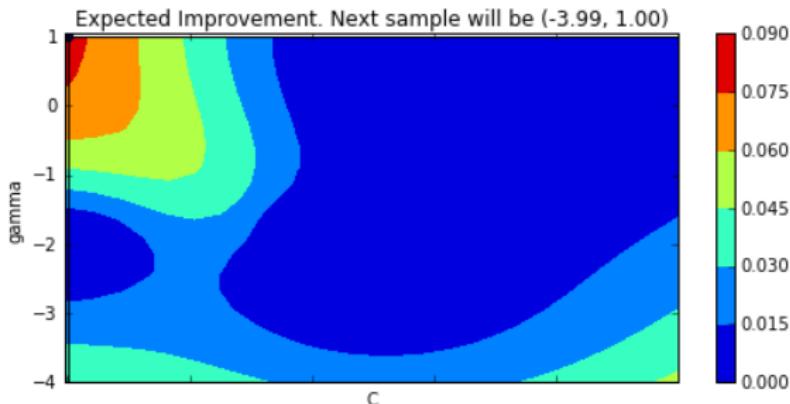
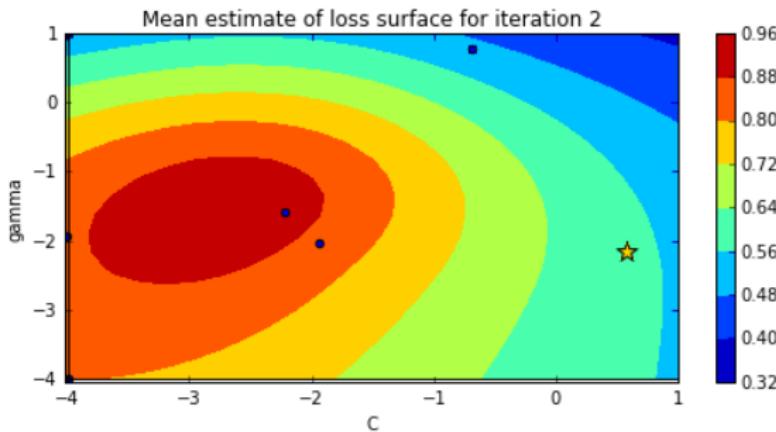
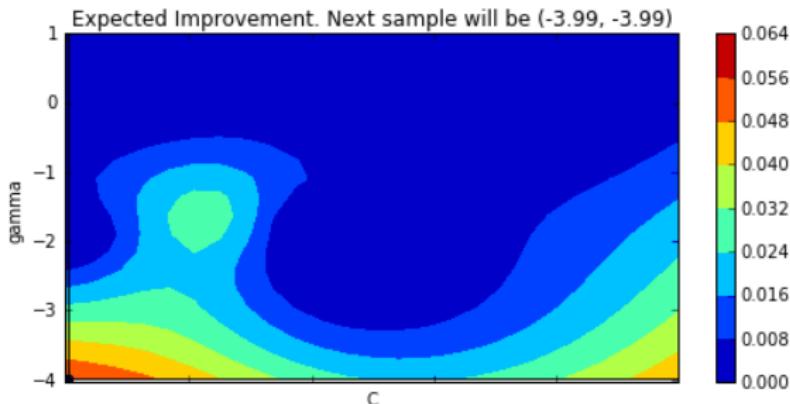


Figure – The target surface w.r.t.  $x$  to see where the true optimum is

# Iterations

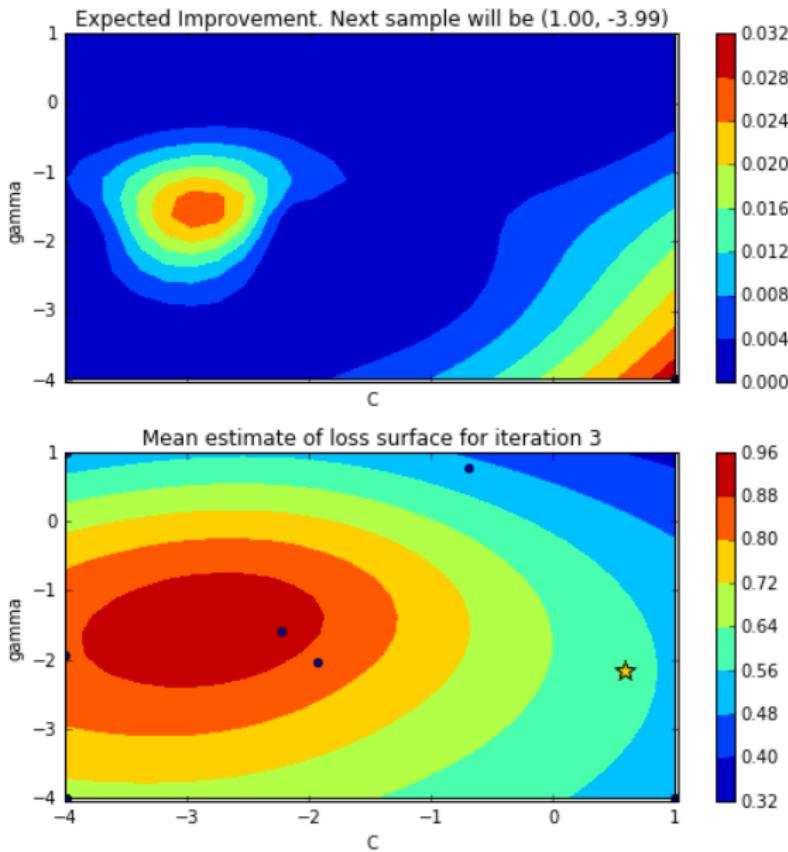
---





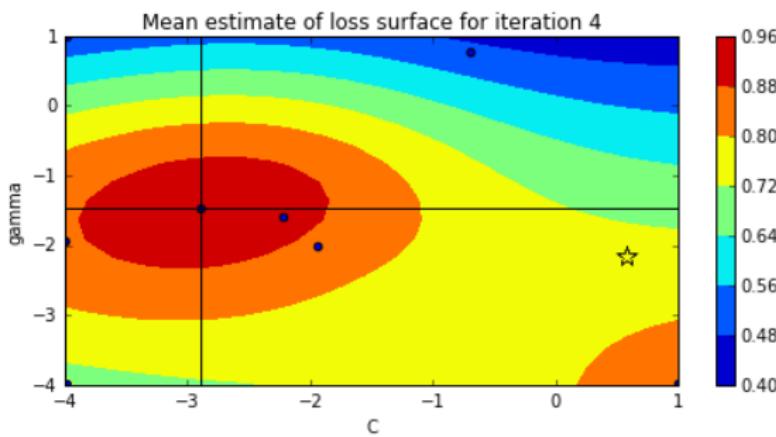
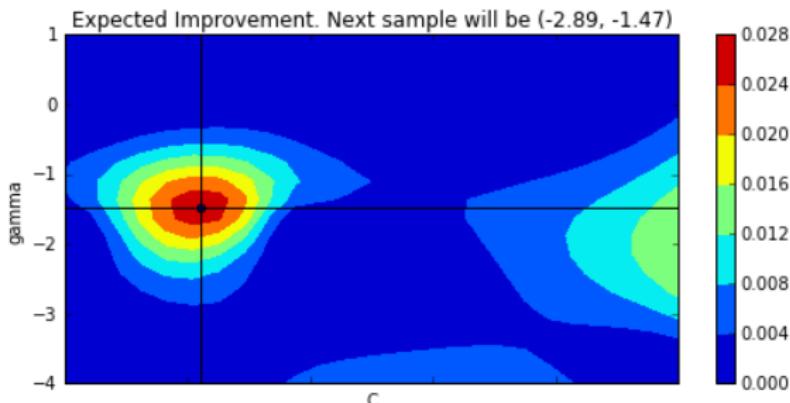
# Iterations

---



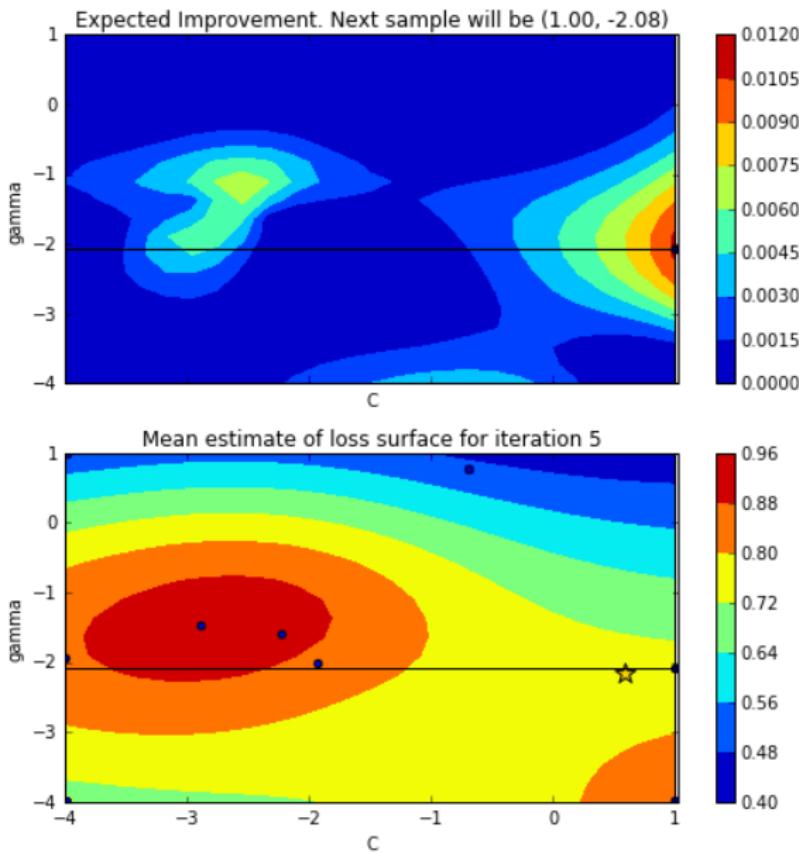
# Iterations

---



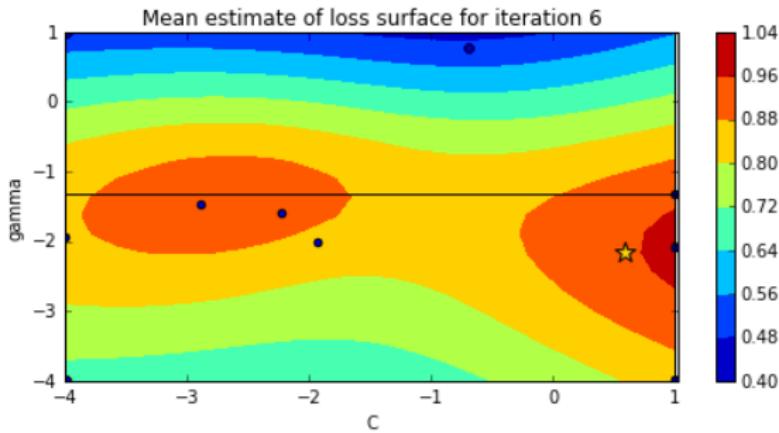
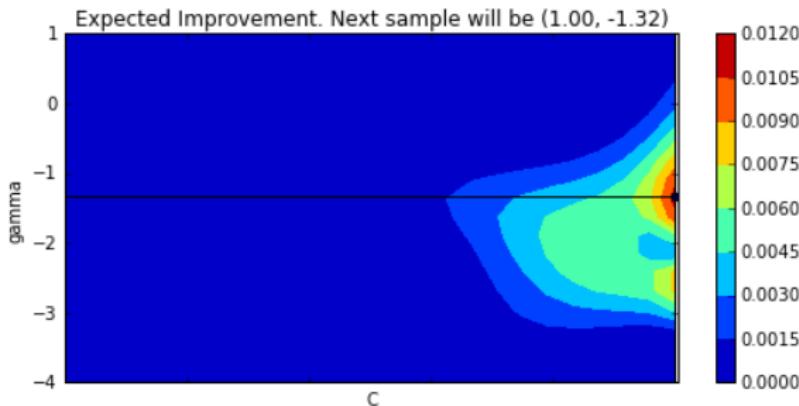
# Iterations

---



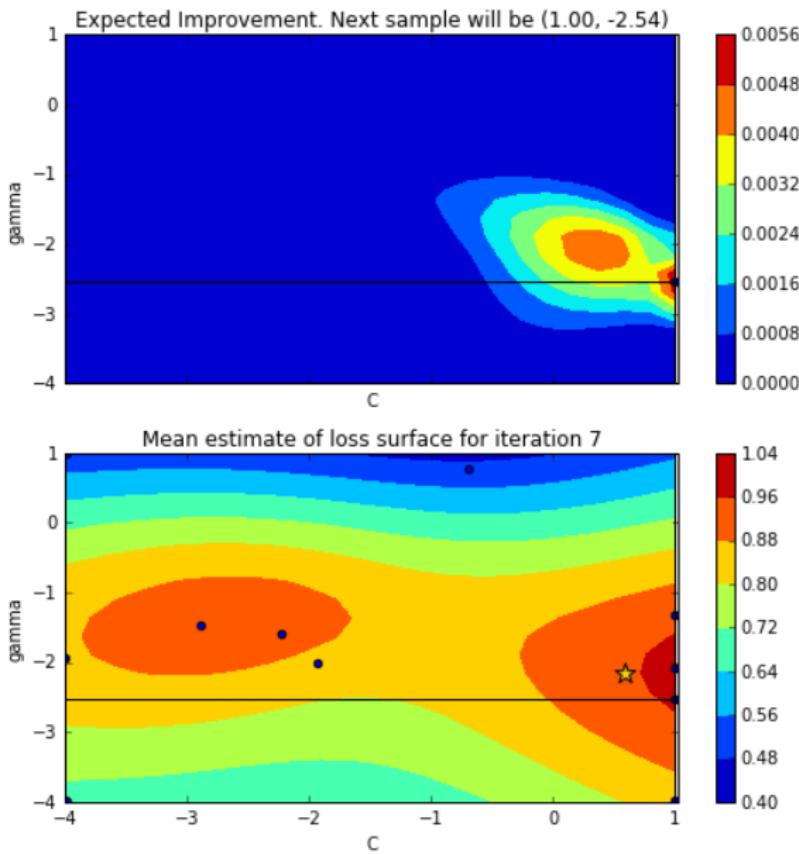
## Iterations

---



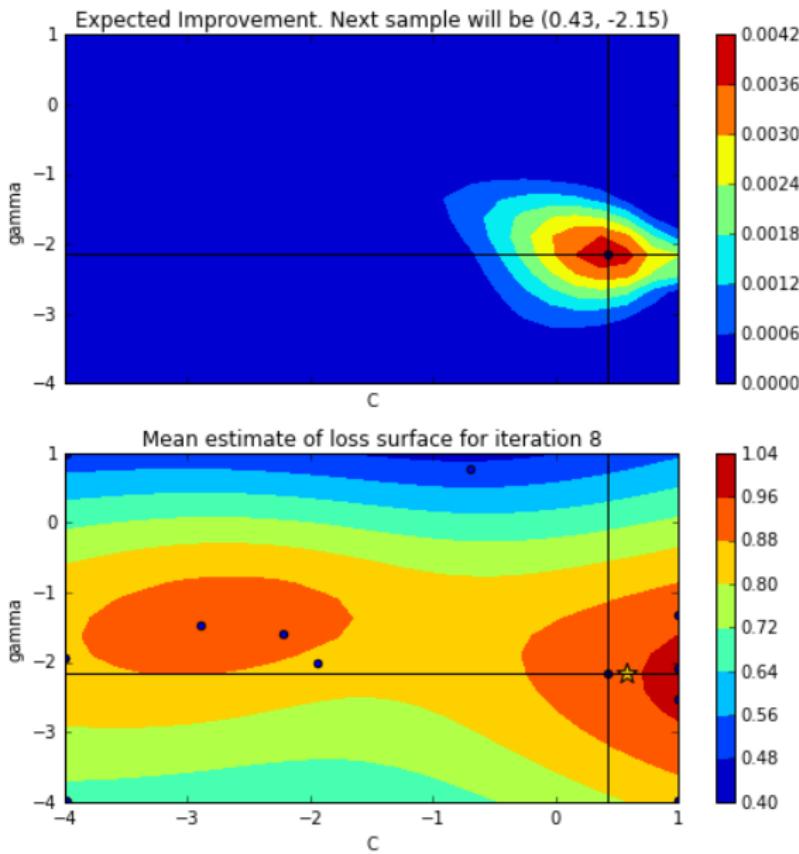
# Iterations

---



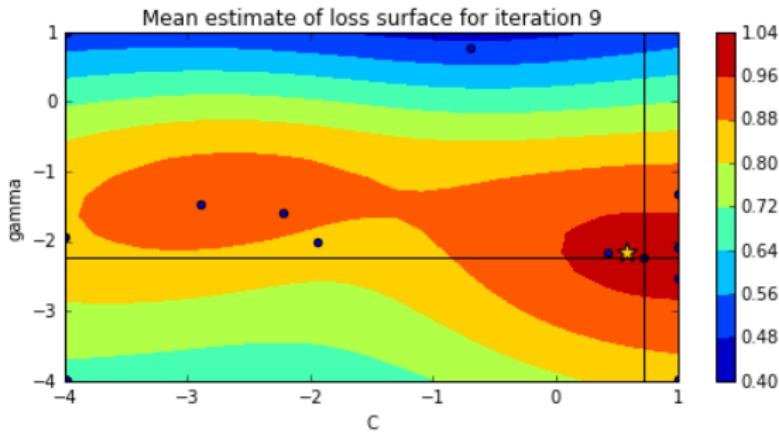
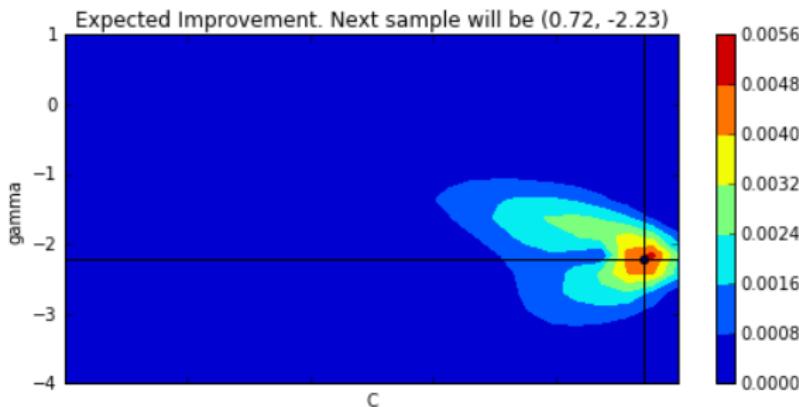
# Iterations

---



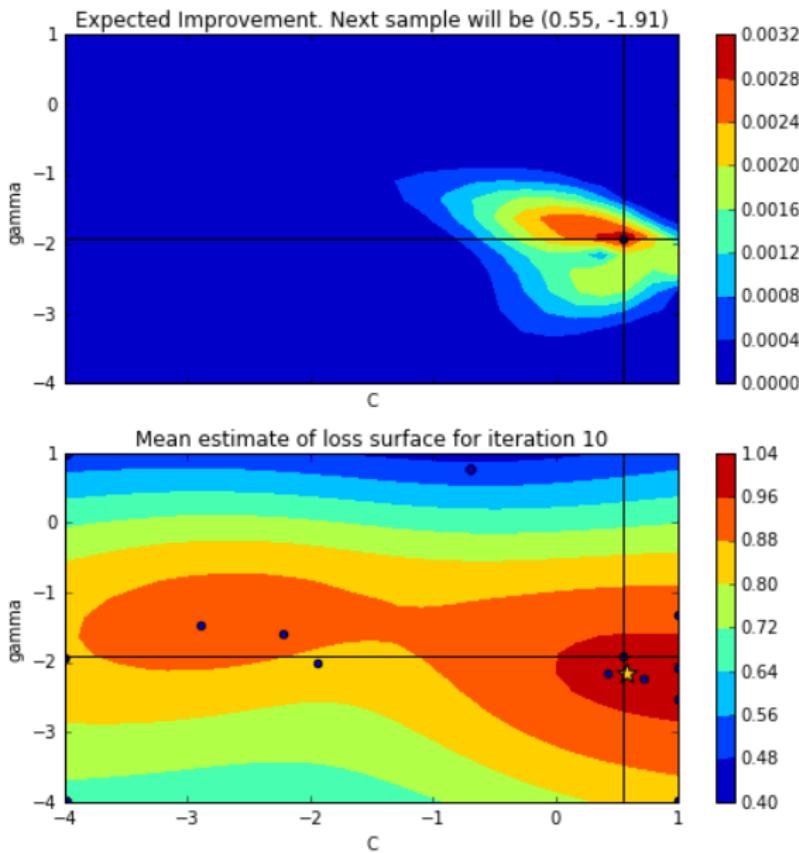
# Iterations

---



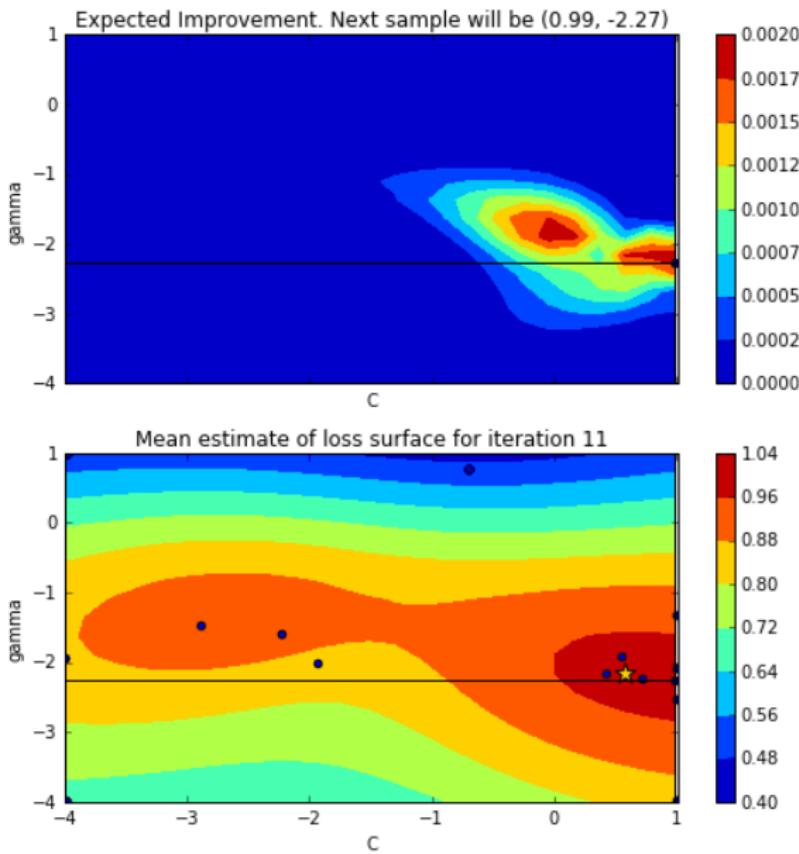
# Iterations

---



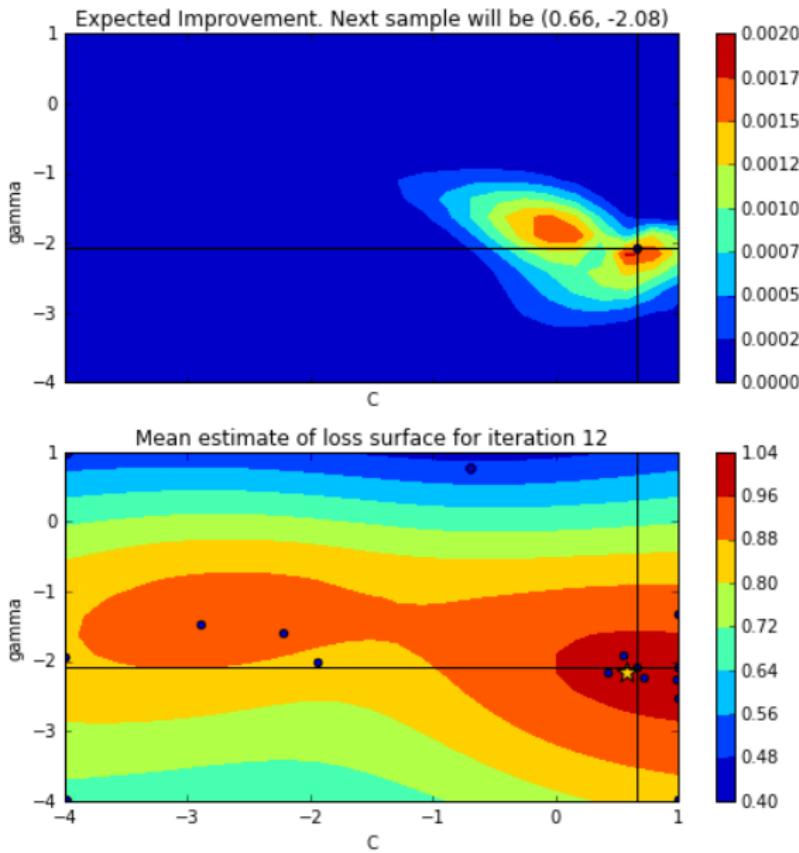
# Iterations

---



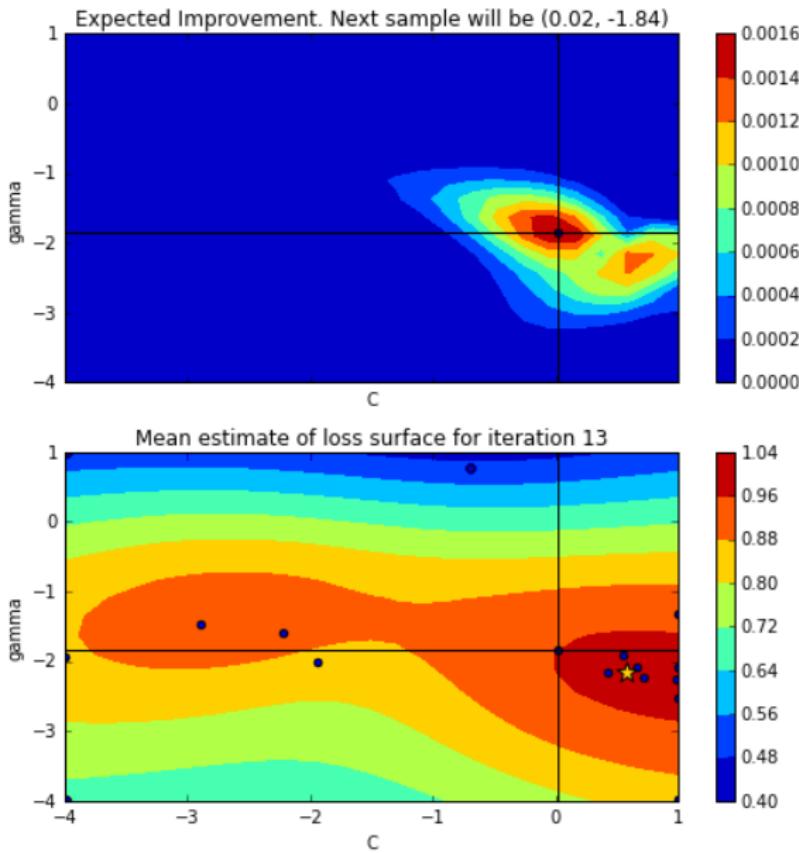
# Iterations

---



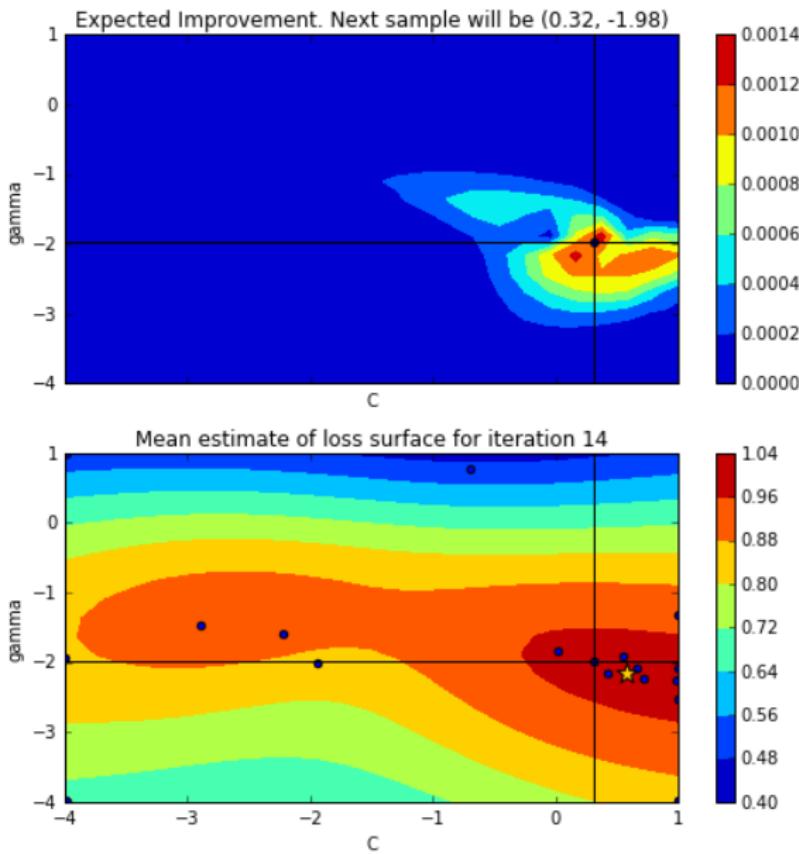
# Iterations

---



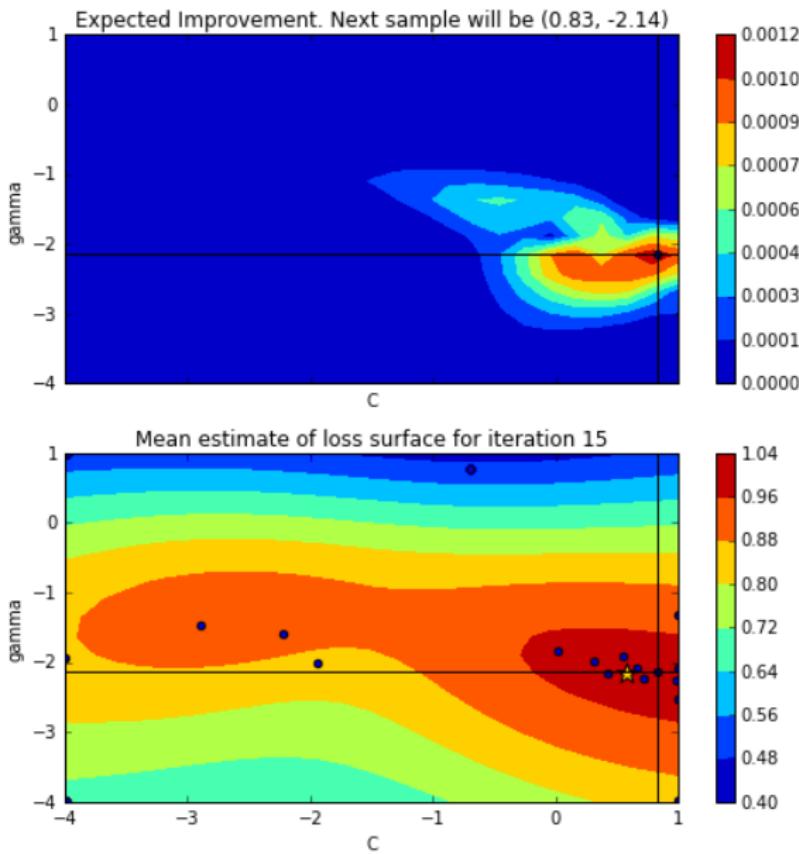
# Iterations

---



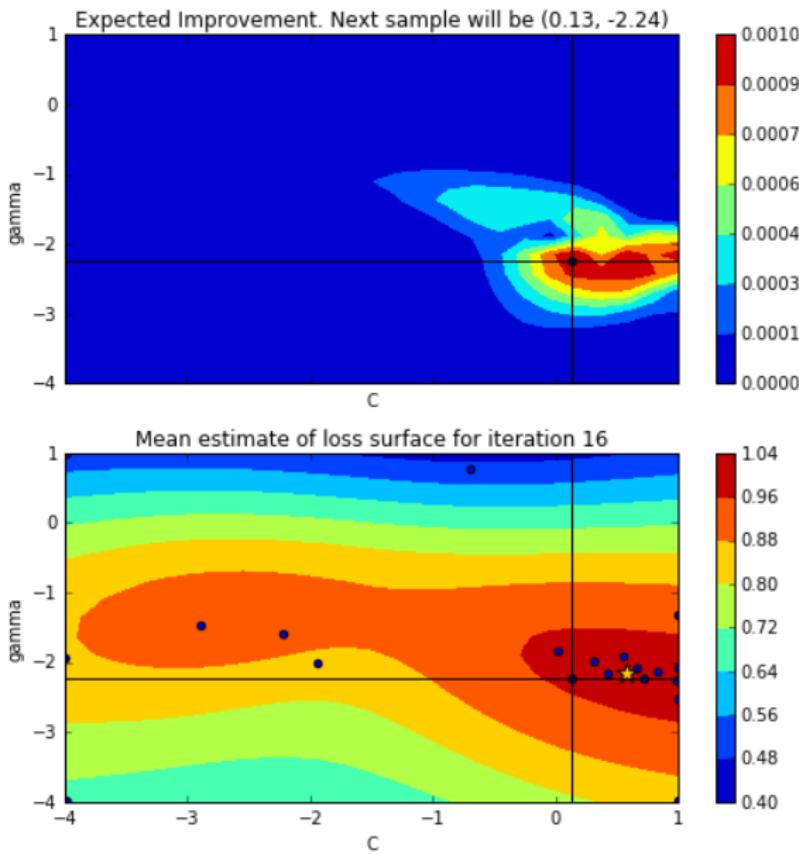
# Iterations

---



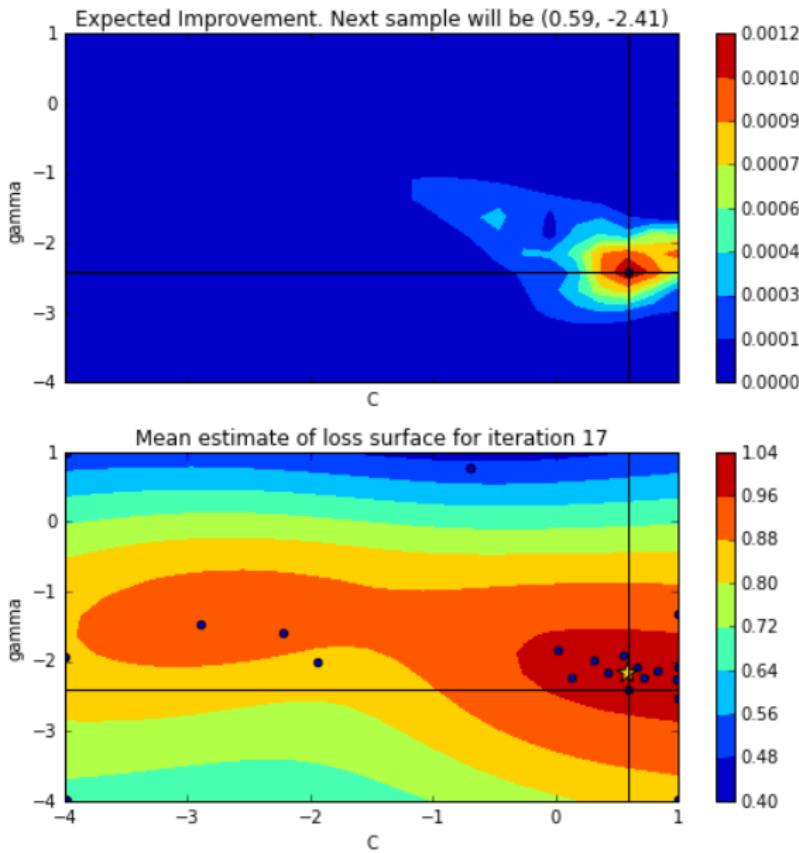
# Iterations

---



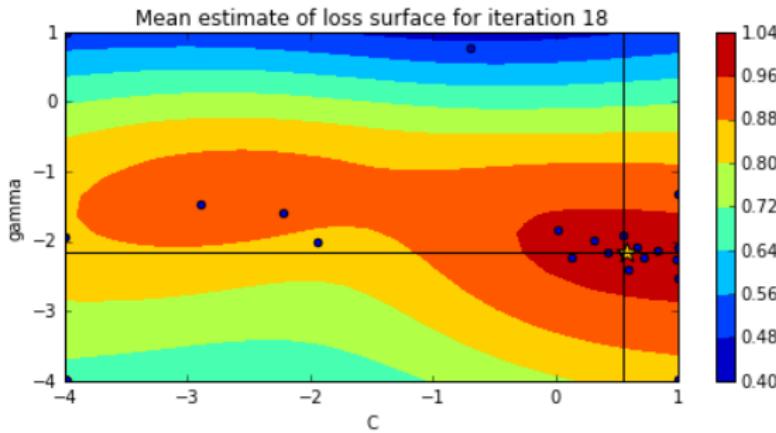
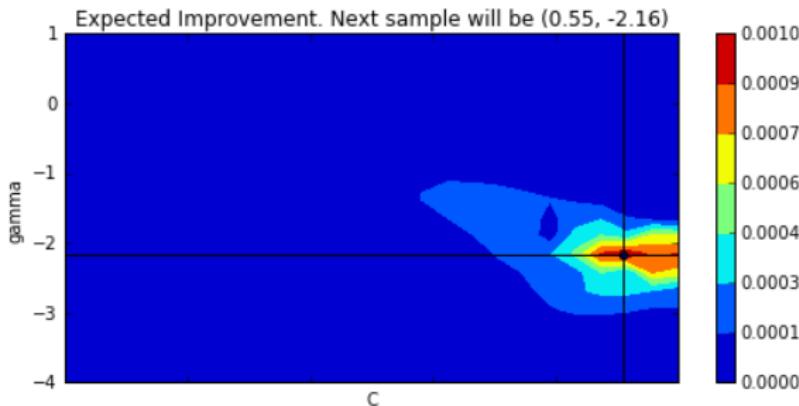
# Iterations

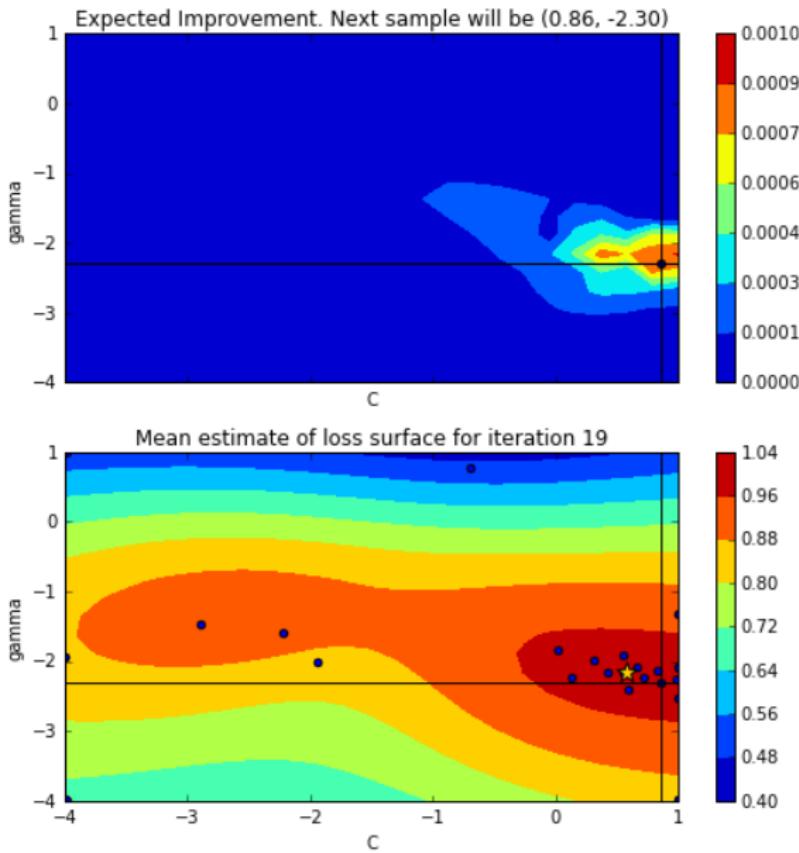
---



# Iterations

---





# Iterations

---

