

Gaussian Processes and Bayesian Optimization

Evgeny Burnaev
Head of ADASE group
Associate Professor, Skoltech



Motivation

Gaussian Process model

GP regression

Learning Gaussian Process model

Bayesian Optimization

Take-home messages

Afterword

Bibliography

Motivation

Gaussian Process model

GP regression

Learning Gaussian Process model

Bayesian Optimization

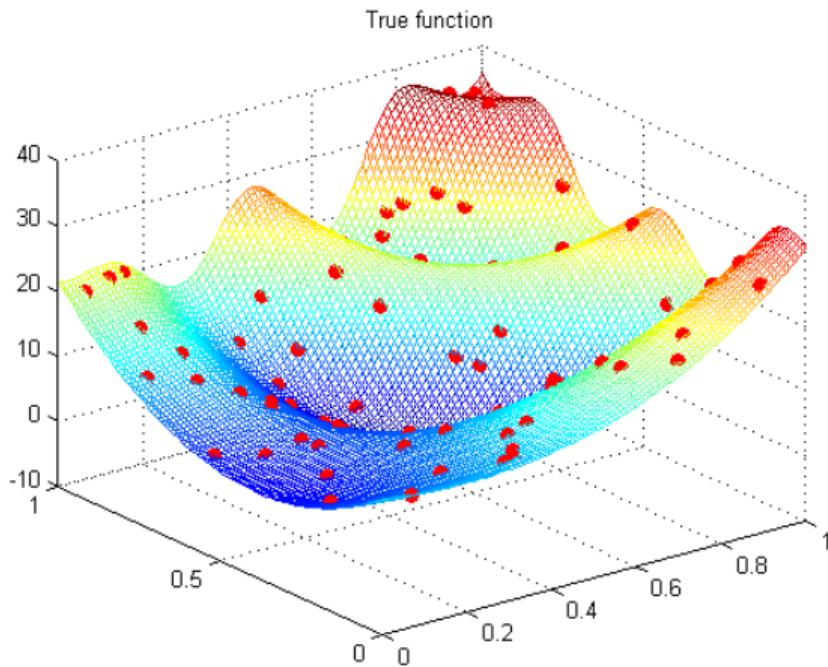
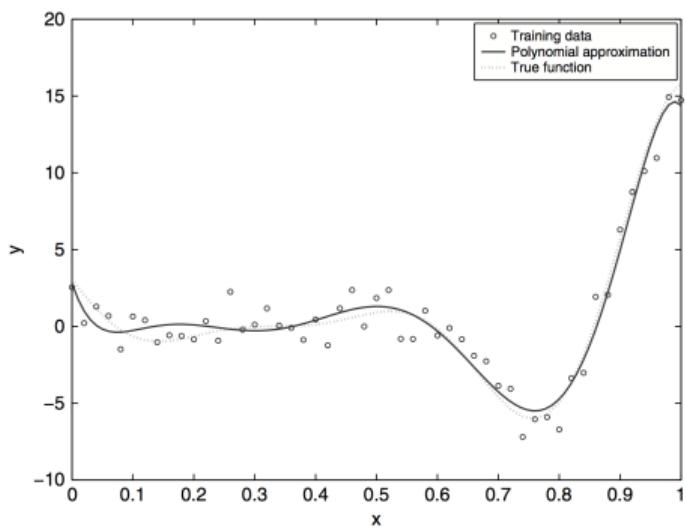
Take-home messages

Afterword

Bibliography

The problem

- ▶ Learn scalar function $f(\mathbf{x})$ of vector \mathbf{x}
- ▶ We have (possibly noisy) sample $S_m = \{\mathbf{X}, \mathbf{y}\} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$



Applications

► Real-valued regression:

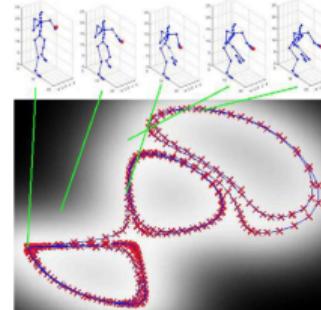
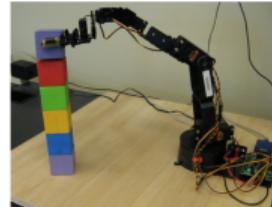
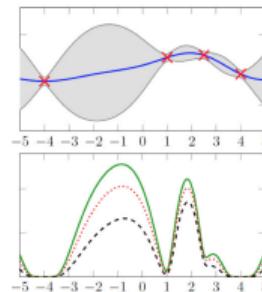
- Model-based predictive control: predict yield, quality, losses, etc.
- Surrogate surfaces for optimization or simulation
- Robotics: target state → required torques

► Classification:

- Recognition: handwritten digits on payment docs, etc.
- Detection: fraud, screening in chemoinformatics

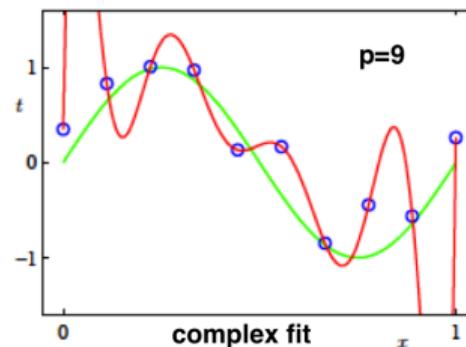
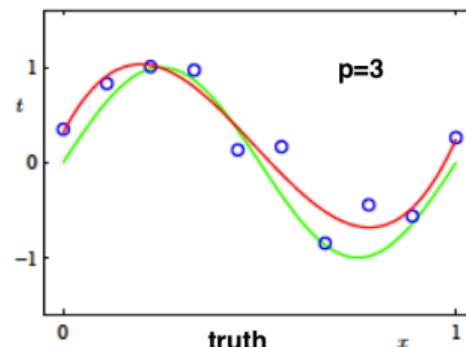
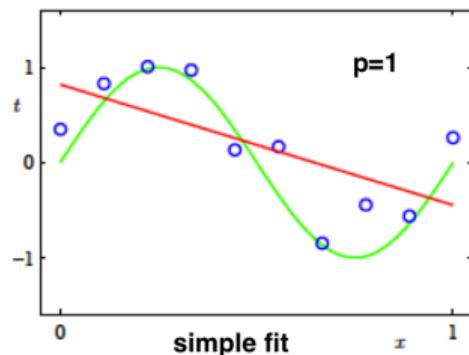
► Ordinal Regression:

- Disease harmfulness prediction
- User ratings (movies, shops, restaurants)



Model Complexity

- Data generation process is often unknown and can be complex:



- Problems:
 - Fitting complicated models can be hard
 - How to find an appropriate model?
 - How to avoid over-fitting?
 - Uncertainty of predictions?

Model-based predictive control

► Problem:

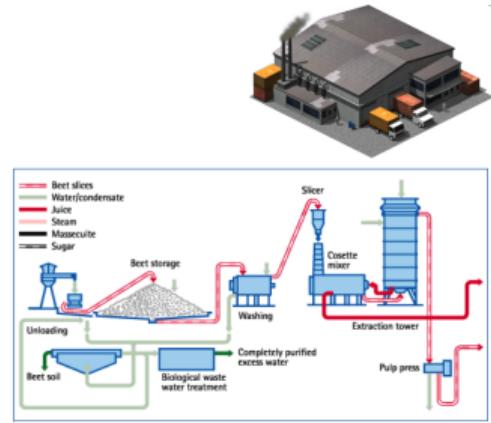
- Object: production of one batch of a product
- Input x : beet chips shape, temperature; wash water temperature, pH and flux, etc.
- Output y : costs, losses, efficiency

► Challenges:

- Heterogeneous data and noise,
- Large volumes of high-dimensional stream data,
- Missing values, outliers, etc.

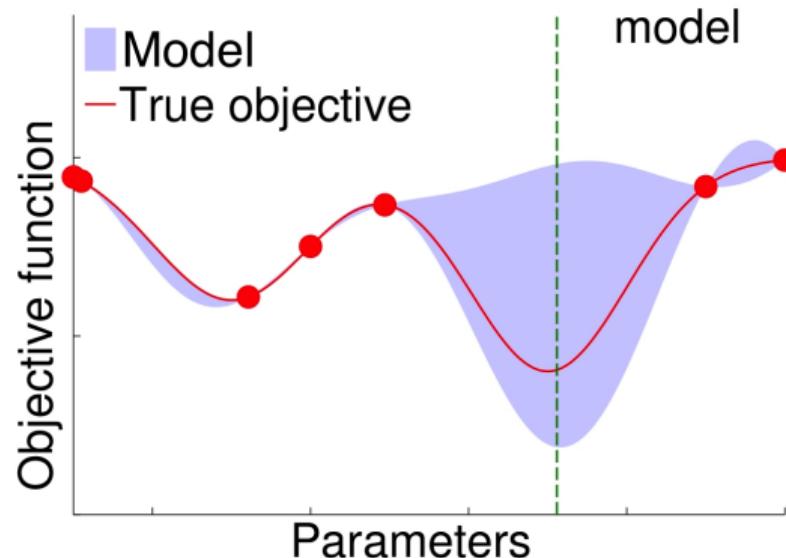
► Knowing error bars is very important:

- Setting $x_1 \rightarrow$ profit of 40 ± 10 units
- Setting $x_2 \rightarrow$ profit of 60 ± 40 units
- Which are the best setting, x_1 or x_2 ?



Optimization

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

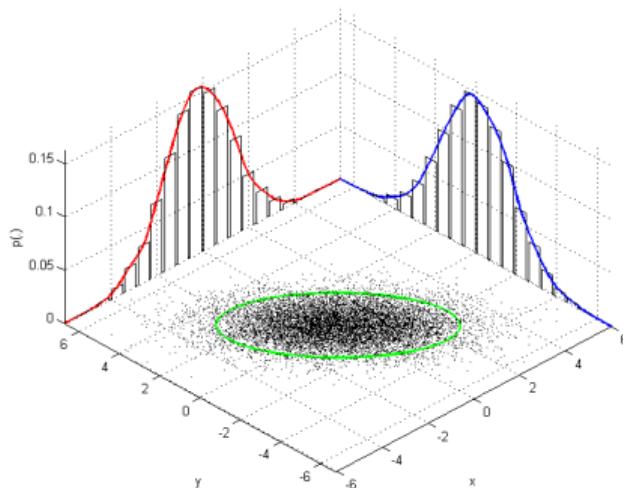


- ▶ Error bars are needed to see if a region is still promising

Bayesian modeling

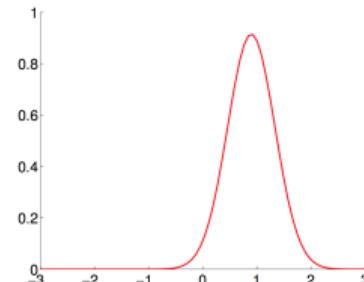
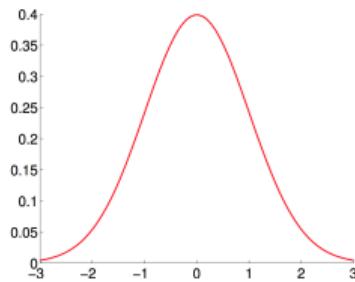
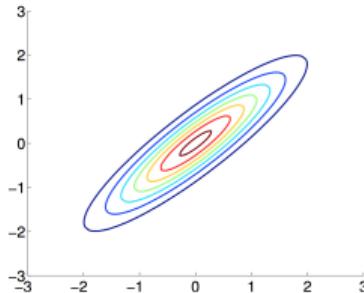
- ▶ Parametric family of functions $f(x; \theta)$
- ▶ Define a prior over θ
- ▶ Perform predictions and estimate uncertainty
- ▶ For flexible models we a.s. get intractable integrals over θ

For Gaussian case usually everything is explicit



Solution of complex ML problems with simple Gaussian models?

The Gaussian Distribution



$$p(\mathbf{f}_1, \mathbf{f}_2) \sim \mathcal{N}(\mathbf{f}_1, \mathbf{f}_2 | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Joint

$$p(\mathbf{f}_1)$$

Marginal

$$p(\mathbf{f}_1 | \mathbf{f}_2)$$

Conditional

The **marginal** and **conditional** distributions are also Gaussians:

$$\begin{pmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \end{pmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{12}^T & \boldsymbol{\Sigma}_{22} \end{bmatrix} \right), \quad p(\mathbf{f}_1) = \int p(\mathbf{f}_1, \mathbf{f}_2) d\mathbf{f}_2 = \mathcal{N}(\mathbf{f}_1 | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$$

$$p(\mathbf{f}_1 | \mathbf{f}_2) = \mathcal{N}(\mathbf{f}_1 | \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{f}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}^T)$$

Motivation

Gaussian Process model

GP regression

Learning Gaussian Process model

Bayesian Optimization

Take-home messages

Afterword

Bibliography

Functional classes

- ▶ **Hypothesis set** $\mathcal{F} \subset Y^X$ is a subset of functions out of which the learner selects his/her hypothesis
 - represents a prior knowledge about the task at hand
 - depends on available features
- ▶ Typical examples (from statistics)
 - Sobolev-type classes:

$$\mathcal{F}_L^k = \left\{ f : \int \left\| \frac{\partial^k f(\mathbf{x})}{\partial \mathbf{x}^k} \right\|^2 d\mathbf{x} \leq L \right\}$$

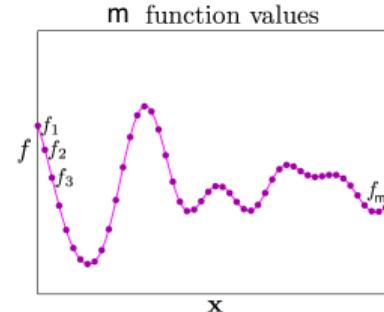
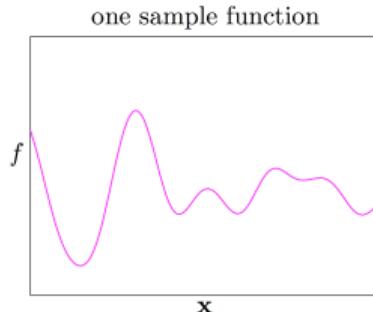
- Lipschitz classes:

$$\mathcal{F}_L = \{f : |f(\mathbf{x}) - f(\mathbf{x}')| \leq L\rho(\mathbf{x}, \mathbf{x}')\}$$

- ▶ How to impose regularity on a function in Bayesian case?

What is a Gaussian process?

- Continuous stochastic process — random function — a set of random variables $f(x)$ indexed by a continuous variable x



- Let us denote by
 - $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ a set of d -dimensional inputs,
 - $\mathbf{f} = \{f_1, \dots, f_m\}$ a set of random function values $f_i = f(\mathbf{x}_i)$
- GP: Any set of function variables $\{f_i\}_{i=1}^m$ has joint Gaussian distribution
$$p(\mathbf{f}|\mathbf{X}) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$$
- Conditional model — density of inputs is not modeled

Mean values and Covariances

- Mean value is constructed from a priori given deterministic function

$$\mu = \{\mu_i\} = \{\mu(\mathbf{x}_i)\}$$

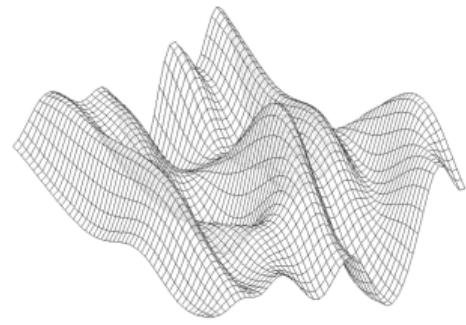
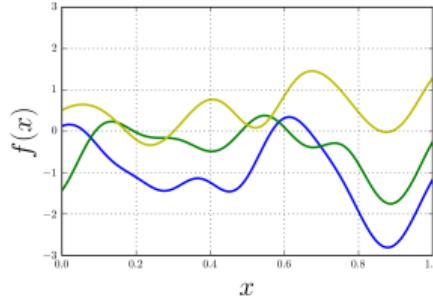
- Covariance matrix is constructed from covariance function

$$\mathbf{K} = \{\mathbf{K}_{ij}\} = \{K(\mathbf{x}_i, \mathbf{x}_j)\}$$

- Covariance function characterizes covariance between points in the process

$$K(\mathbf{x}, \mathbf{x}') = \text{cov}(f(\mathbf{x}), f(\mathbf{x}')) = \mathbb{E} (f(\mathbf{x}) - \mu(\mathbf{x})) (f(\mathbf{x}') - \mu(\mathbf{x}'))$$

$$f(\mathbf{x}) \sim \mathcal{GP}(\cdot | \mu(\mathbf{x}), K(\mathbf{x}, \mathbf{x}'))$$

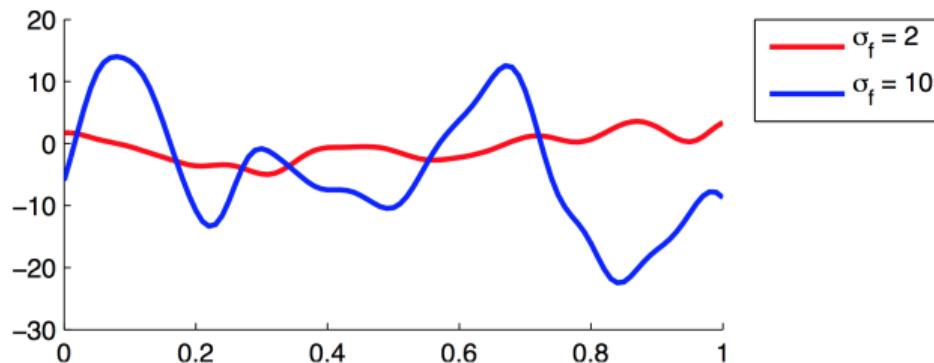


Example: squared-exponential kernel

- The most commonly-used kernel in machine learning

$$K(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp \left\{ -\sum_{i=1}^d \frac{(x_i - x'_i)^2}{2r_i^2} \right\}$$

- A GP need not use the Gaussian kernel. There are a lot of other choices
- Consider $\mathbf{x} = \mathbf{x}' \rightarrow$ marginal function variance is σ_f^2

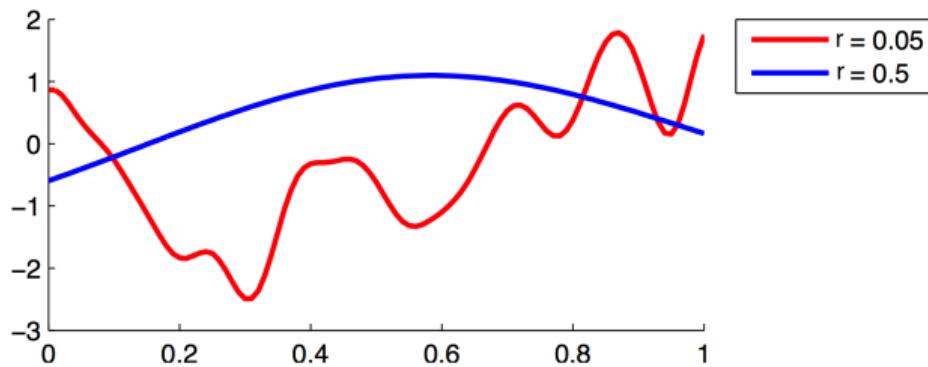


Example: squared-exponential kernel

- The r_i parameters give the overall lengthscale in dimension i

$$K(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp \left\{ - \sum_{i=1}^d \frac{(x_i - x'_i)^2}{2r_i^2} \right\}$$

- Typical distance between peaks $\approx r$



Example: Gaussian process covariance functions

- More general example of covariance function

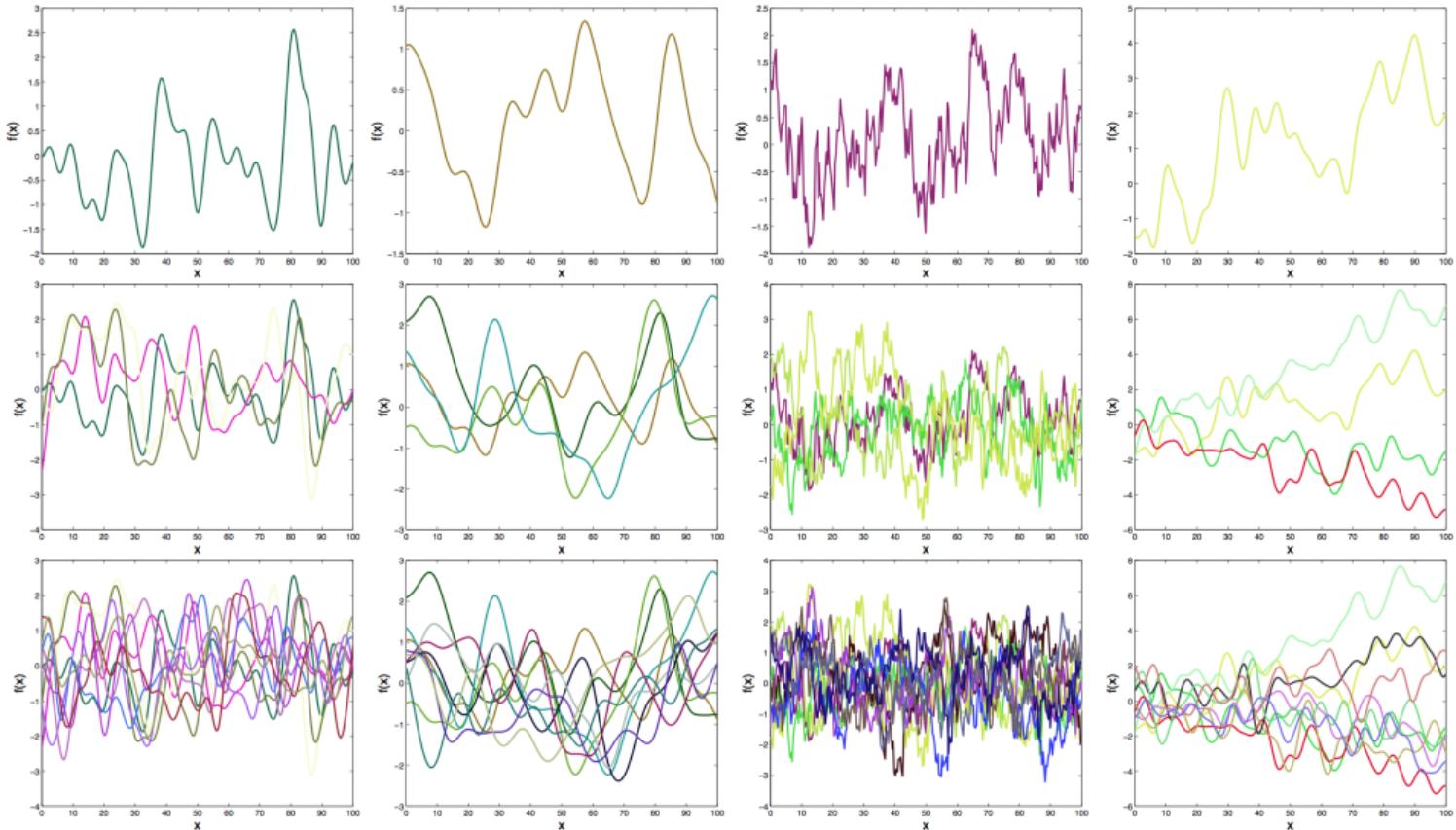
$$K(x, x') = \sigma_f^2 \exp \left\{ - \left(\frac{|x - x'|}{r} \right)^\alpha \right\} + \sigma_1^2 + \sigma_2^2 \delta(x - x')$$

- As usual covariance function parameters are **interpretable**

- σ_f^2 — marginal function variance
- σ_1^2 — variance of bias
- σ_2^2 — noise variance
- r — lengthscale
- $\alpha \geq 1$ — roughness

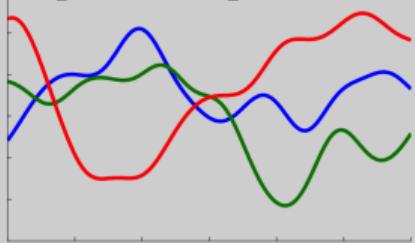
- Once the mean and covariance functions are defined, everything else about GPs follows from the basic rules of probability applied to multivariate Gaussians

Samples from GPs with different $K(x, x')$



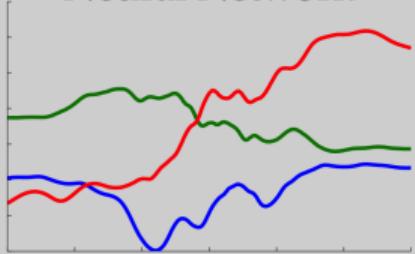
Comparison of GP realizations vs. different covariance functions

Squared-Exponential



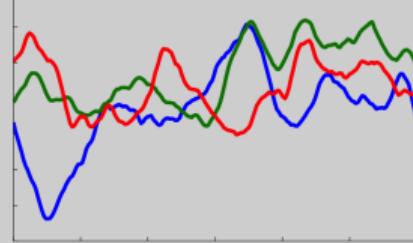
$$K(x, x') = \alpha \exp \left\{ -\frac{1}{2} \sum_{d=1}^D \left(\frac{x_d - x'_d}{\ell_d} \right)^2 \right\}$$

“Neural Network”



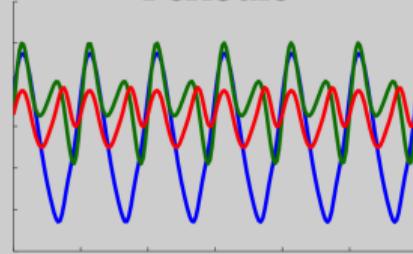
$$K(x, x') = \frac{2}{\pi} \sin^{-1} \left\{ \frac{2x^\top \Sigma x'}{\sqrt{(1 + 2x^\top \Sigma x)(1 + 2x'^\top \Sigma x')}} \right\}$$

Matérn



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

Periodic



$$K(x, x') = \exp \left\{ -\frac{2 \sin^2 \left(\frac{1}{2}(x - x') \right)}{\ell^2} \right\}$$

Example: Constructing new covariances from old

There are several ways to combine covariances (**seminar**)

- ▶ **Sum:** $K(x, x') = K_1(x, x') + K_2(x, x')$
- ▶ **Product:** $K(x, x') = K_1(x, x') \cdot K_2(x, x')$
- ▶ **Convolution:** $K(x, x') = \int dz dz' h(x, z) K(z, z') h(x', z')$ (blurring of process with kernel h)

Motivation

Gaussian Process model

GP regression

Learning Gaussian Process model

Bayesian Optimization

Take-home messages

Afterword

Bibliography

Using GP for nonlinear regression

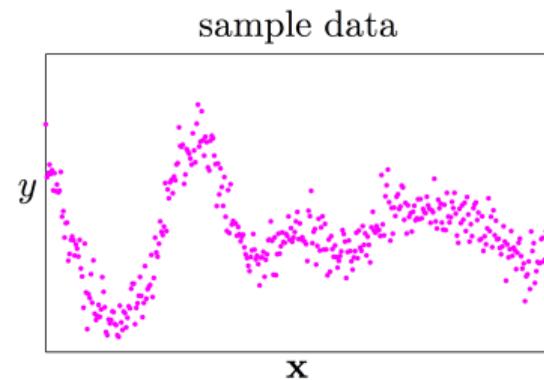
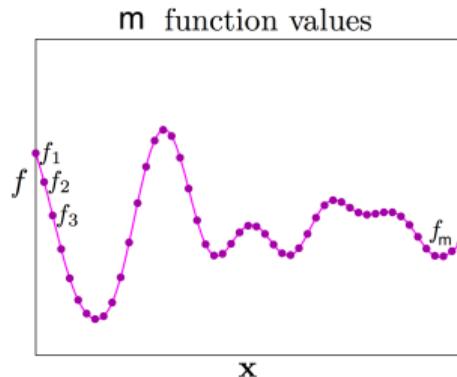
- ▶ Training data set $S_m = \{\mathbf{X}, \mathbf{y}\} = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$
- ▶ Model:

$$y_i = f(\mathbf{x}_i) + \varepsilon_i,$$

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

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

- ▶ We can use heteroscedastic noise and/or non-Gaussian noise models



Using GP for nonlinear regression

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

▶ **Model:**

$$y_i = f(\mathbf{x}_i) + \varepsilon_i,$$

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

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

- ▶ The prior is

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

- ▶ The noise model, or **likelihood** is

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

- ▶ Integrating over the function variables \mathbf{f} we get the **marginal likelihood**

$$p(\mathbf{y}) = \int d\mathbf{f} p(\mathbf{y} | \mathbf{f}) p(\mathbf{f}) d\mathbf{f} = \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I}_m)$$

Prediction

- ▶ Let us denote input test point as \mathbf{x}_* , and output

$$y_* = f_* + \varepsilon_*, \quad f_* = f(\mathbf{x}_*)$$

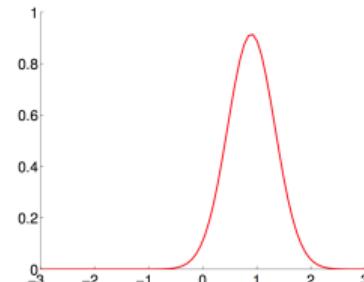
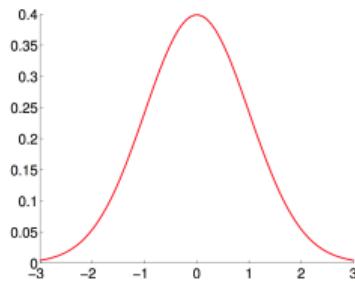
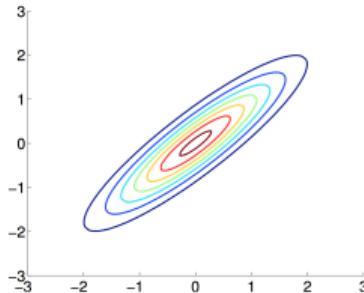
- ▶ Consider joint training and test marginal likelihood

$$p(\mathbf{y}, f_*) = \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma^2 \mathbf{I}_m & \mathbf{k}_* \\ \mathbf{k}_*^T & K_{**} \end{bmatrix} \right),$$

where $\mathbf{k}_* = \{K(\mathbf{x}_*, \mathbf{x}_i)\}_{i=1}^m$ and $K_{**} = K(\mathbf{x}_*, \mathbf{x}_*)$

- ▶ What we know about noiseless value $f(\mathbf{x}_*)$?

The Gaussian Distribution



$$p(\mathbf{f}_1, \mathbf{f}_2) \sim \mathcal{N}(\mathbf{f}_1, \mathbf{f}_2 | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Joint

$$p(\mathbf{f}_1)$$

Marginal

$$p(\mathbf{f}_1 | \mathbf{f}_2)$$

Conditional

The **marginal** and **conditional** distributions are also Gaussians:

$$\begin{pmatrix} \mathbf{f}_1 \\ \mathbf{f}_2 \end{pmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{12}^T & \boldsymbol{\Sigma}_{22} \end{bmatrix}\right), \quad p(\mathbf{f}_1) = \int p(\mathbf{f}_1, \mathbf{f}_2) d\mathbf{f}_2 = \mathcal{N}(\mathbf{f}_1 | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$$

$$p(\mathbf{f}_1 | \mathbf{f}_2) = \mathcal{N}(\mathbf{f}_1 | \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{f}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{12}^T)$$

Prediction

- ▶ Joint training and test marginal likelihood

$$p(\mathbf{y}, f_*) = \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma^2 \mathbf{I}_m & \mathbf{k}_* \\ \mathbf{k}_*^\top & K_{**} \end{bmatrix} \right),$$

where $\mathbf{k}_* = \{K(\mathbf{x}_*, \mathbf{x}_i)\}_{i=1}^m$ and $K_{**} = K(\mathbf{x}_*, \mathbf{x}_*)$

- ▶ Condition on training outputs we get

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

$$\mu_* = \mathbf{k}_*^\top [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{y}, \quad \sigma_*^2 = K_{**} - \mathbf{k}_*^\top [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{k}_*$$

- ▶ In fact μ_* has the form

$$\mu_* = \sum_{i=1}^m \alpha_i K(\mathbf{x}_*, \mathbf{x}_i), \quad \boldsymbol{\alpha} = [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{y}$$

- ▶ Computational complexity: $O(m^3)$ for inversion in $\boldsymbol{\alpha}$, $O(m)$ for μ_* and $O(m^2)$ for σ_*^2

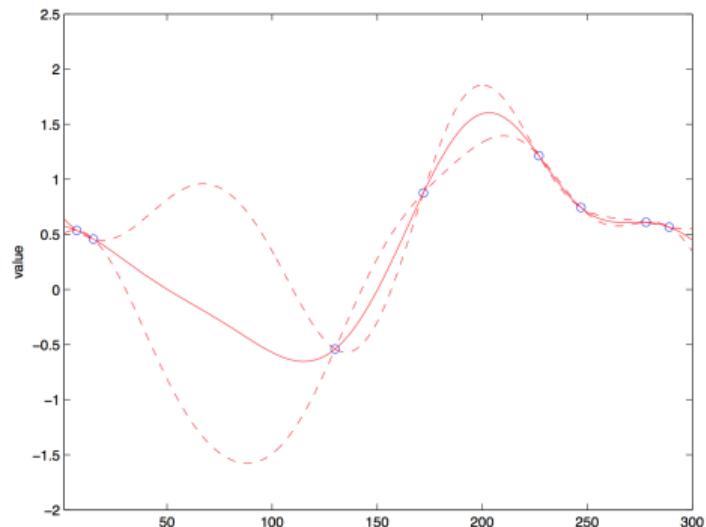
Interpolation

- ▶ Conditional distribution

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

$$\mu_* = \mathbf{k}_*^T [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{y}, \quad \sigma_*^2 = K_{**} - \mathbf{k}_*^T [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{k}_*$$

- ▶ What we know about noiseless value $f_* = f(\mathbf{x}_*)$? Let us assume that $\sigma = 0$ (**interpolation**)



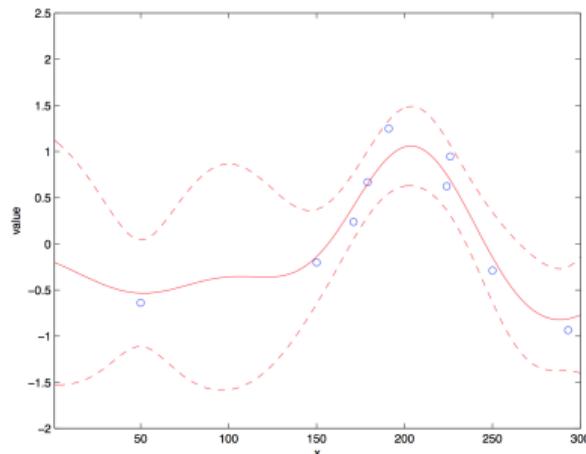
Smoothing

- ▶ Conditional distribution

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

$$\mu_* = \mathbf{k}_*^T [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{y}, \quad \sigma_*^2 = K_{**} - \mathbf{k}_*^T [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{k}_*$$

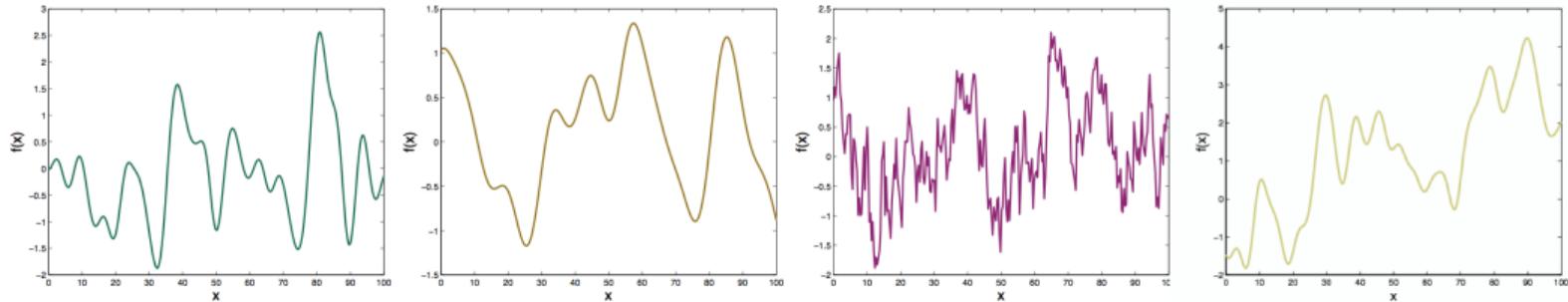
- ▶ What we know about noisy value $f_* = f(\mathbf{x}_*)$?



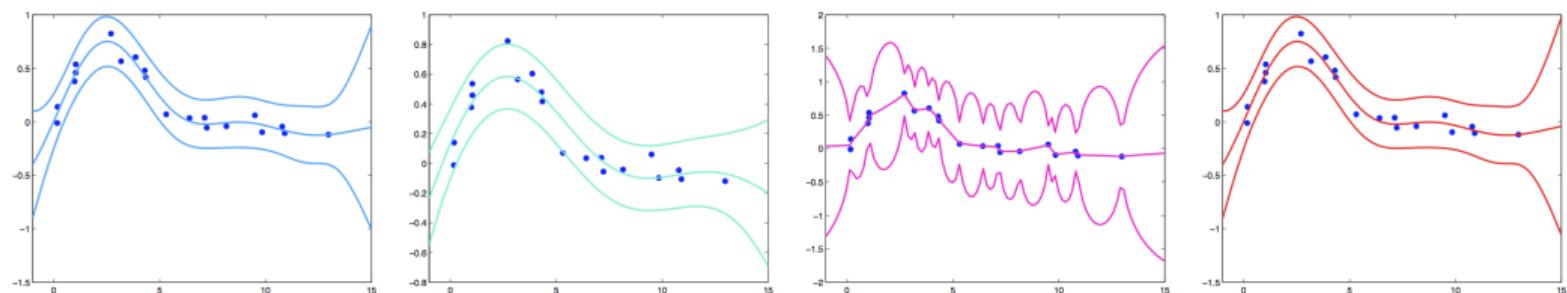
- ▶ $p(y_* | \mathbf{y}) = \mathcal{N}(y_* | \mu_*, \sigma_*^2 + \sigma^2)$ predicts what we'll see next

Prediction using GPs with different $K(x, x')$

A sample from the prior for each covariance function:



Corresponding predictions, mean with two standard deviations:



Motivation

Gaussian Process model

GP regression

Learning Gaussian Process model

Bayesian Optimization

Take-home messages

Afterword

Bibliography

Determination of covariance function parameters

- ▶ Advantage of the probabilistic GP framework — ability to choose hyperparameters and covariances directly from the training data
- ▶ Other models, e.g. SVMs, require cross-validation
- ▶ **GP**: minimize negative log marginal likelihood $\mathcal{L}(\theta)$ wrt covariance function parameters and noise level θ . Since

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I}_m),$$

then

$$\mathcal{L} = -\log p(\mathbf{y}|\theta) = \underbrace{\frac{1}{2} \log \det \mathbf{C}(\theta)}_{\text{regularization}} + \underbrace{\frac{1}{2} \mathbf{y}^T \mathbf{C}^{-1}(\theta) \mathbf{y}}_{\text{data-fit}} + \frac{m}{2} \log(2\pi),$$

where $\mathbf{C} = \mathbf{K} + \sigma^2 \mathbf{I}_m$

- ▶ Uncertainty in function variables \mathbf{f} is taken into account

Gradient based optimization

- ▶ Minimization of $\mathcal{L}(\theta)$ is a **non-convex optimization** problem
- ▶ Standard gradient based techniques, such as CG or quasi-Newton
- ▶ **Gradients:**

$$\frac{\partial \mathcal{L}}{\partial \theta_i} = \frac{1}{2} \text{tr} \left(\mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_i} \right) - \frac{1}{2} \mathbf{y}^T \mathbf{C}^{-1} \frac{\partial \mathbf{C}}{\partial \theta_i} \mathbf{C}^{-1} \mathbf{y}$$

- ▶ Local minima, but usually not much of a problem with few parameters, and/or using restarts
- ▶ Use weighted sums of covariances and let ML choose
- ▶ Each iteration of optimization costs $O(m^3)$

Motivation

Gaussian Process model

GP regression

Learning Gaussian Process model

Bayesian Optimization

Take-home messages

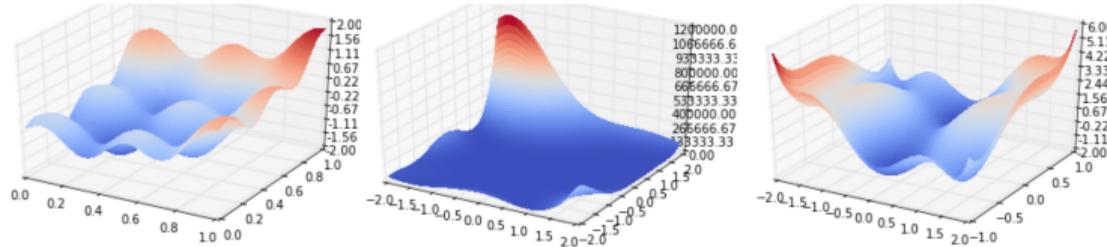
Afterword

Bibliography

Global optimization

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

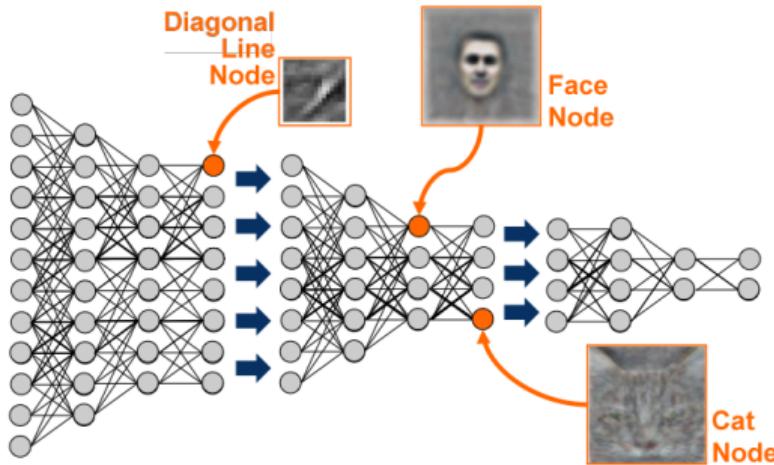
$$\mathbf{x}_{\min} = \arg \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$$



- ▶ 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

Auto ML

Parameter tuning in ML algorithms



x contains parameters

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

$y = f(x)$ is a prediction accuracy, estimated via comp. expensive cross-validation

Other expensive functions

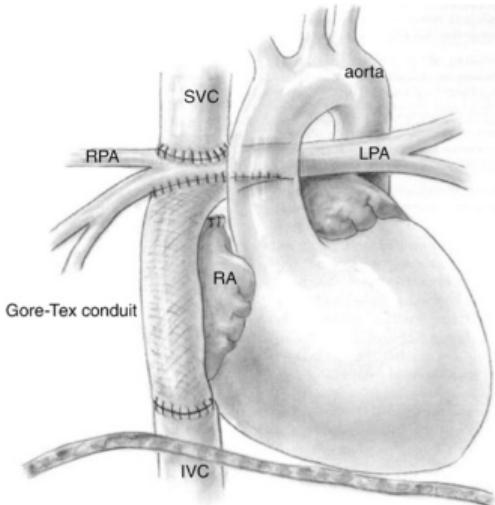
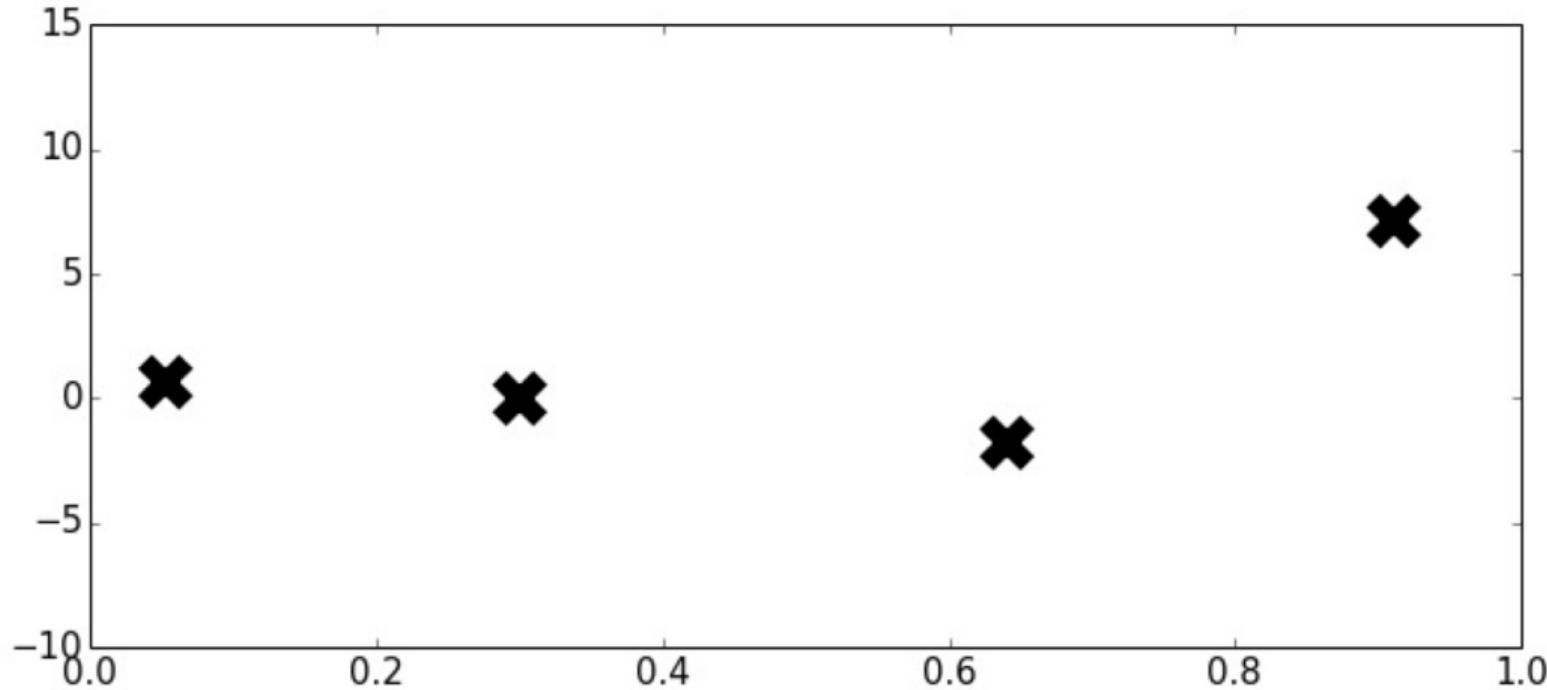


Fig. 1. Extracardiac total cavopulmonary connection. The IVC is disconnected from the right atrium (RA) and connected to the PAs via a Gore-Tex conduit. Figure taken from Reddy et al. [13].

- ▶ 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.

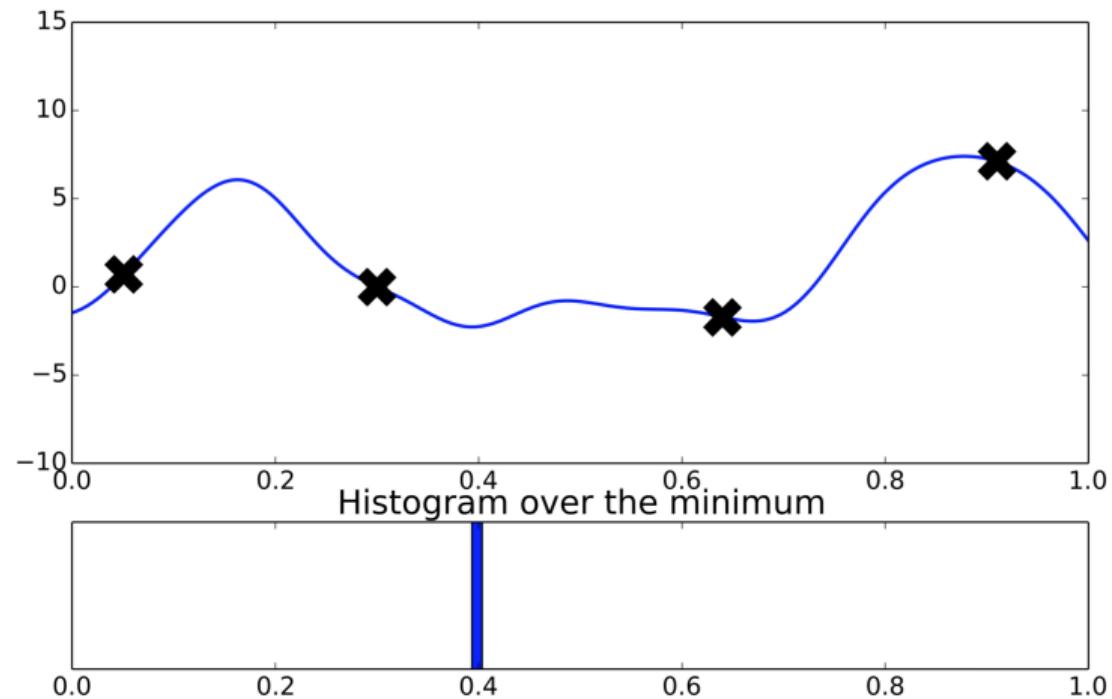
NB! There exists commercial services for optimizing black-box functions:
SIGOPT, Google Vizier, etc.

Typical situation

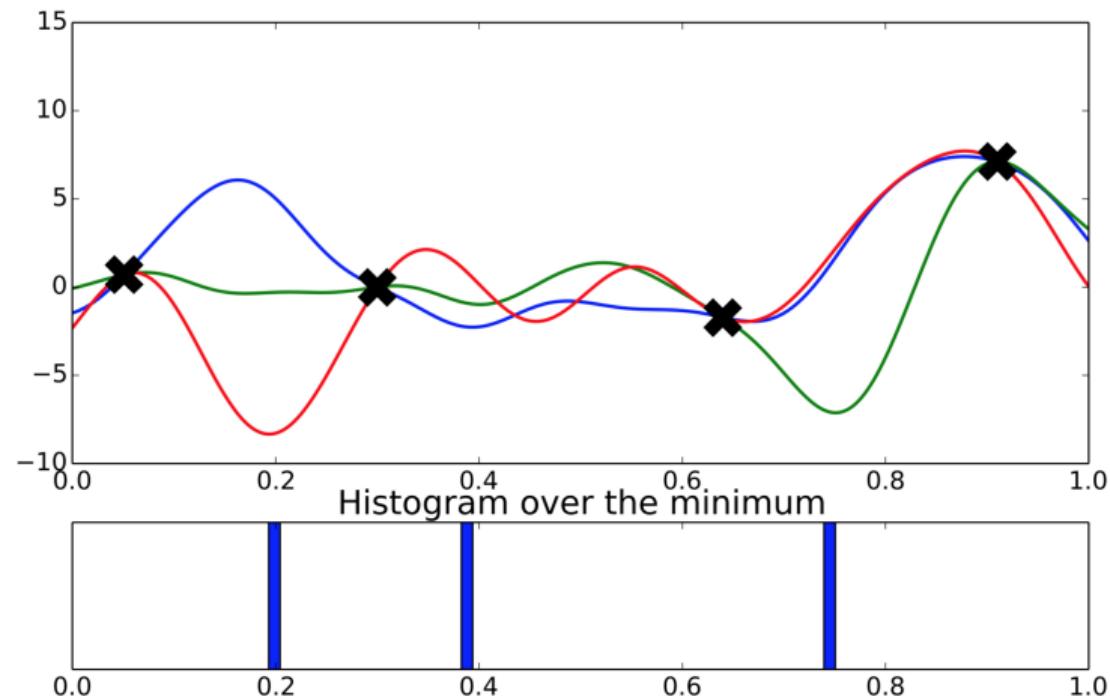


Where is the minimum of f ?
Where should we evaluate the function next?

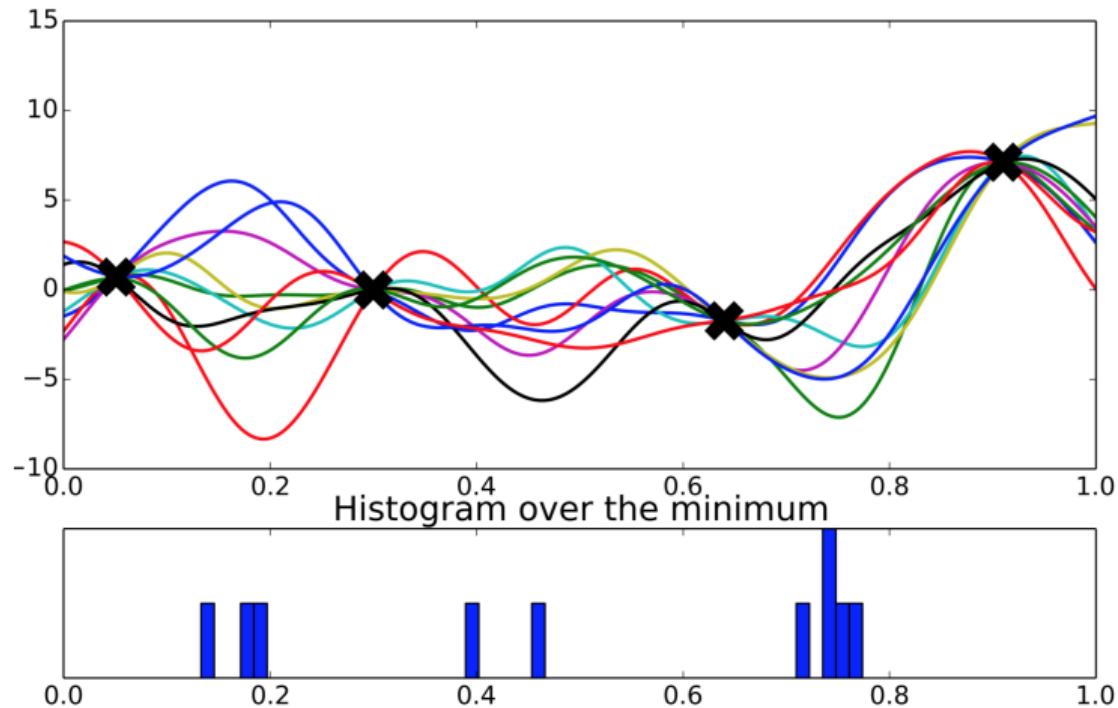
Intuition: one curve



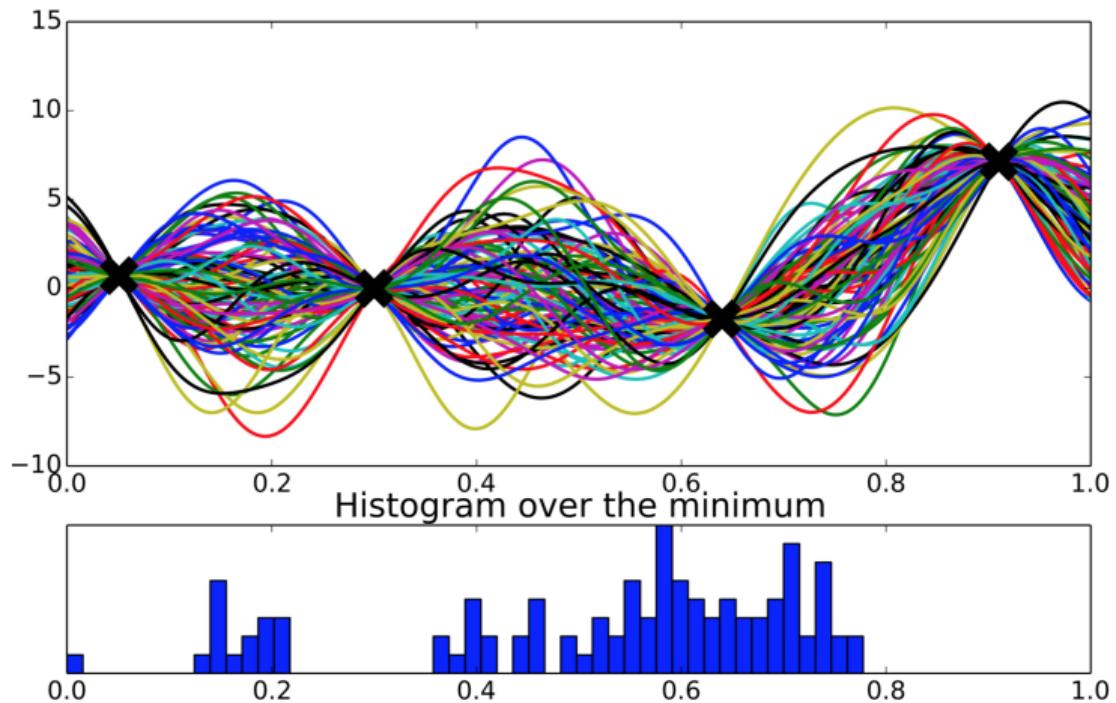
Intuition: three curves



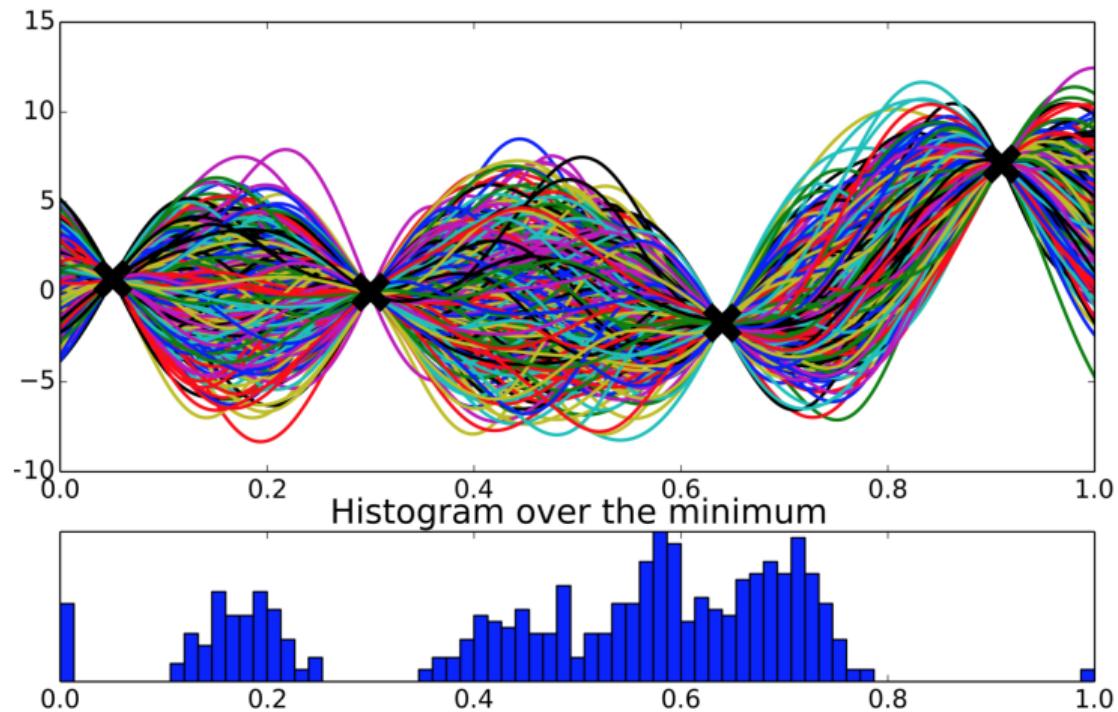
Intuition: ten curves



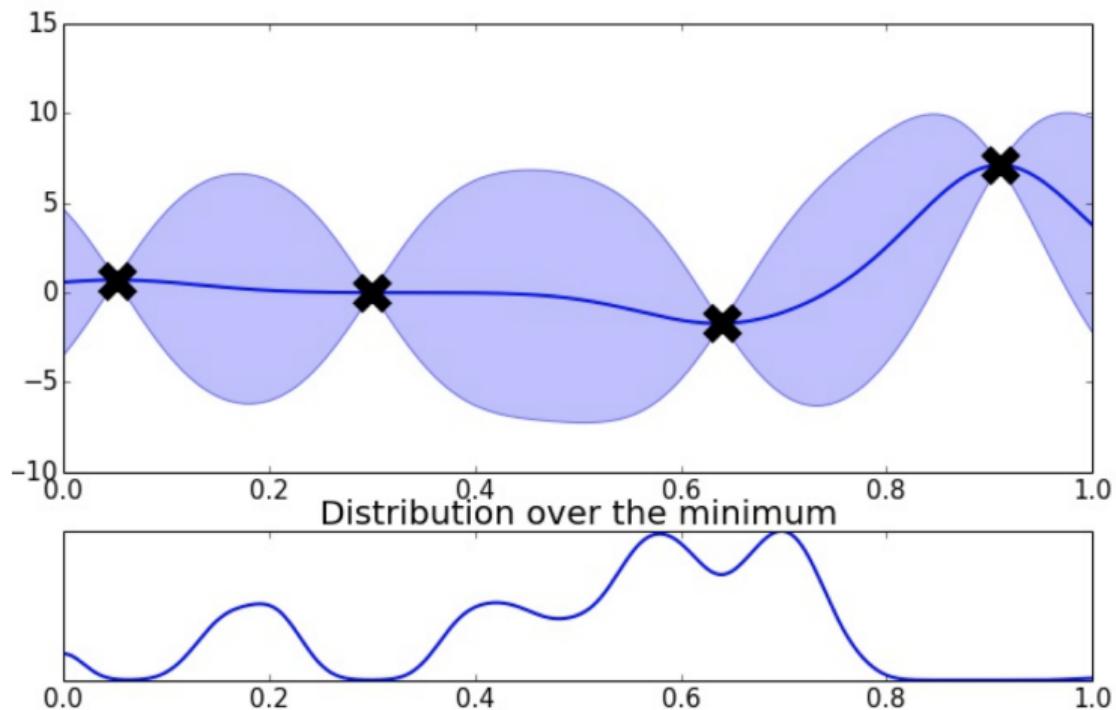
Intuition: hundred curves



Intuition: many curves



Intuition: infinite number of curves



Formalization

- ▶ 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.

Bayesian Optimization

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

Iterate between 2 and 4 until the evaluation budget is over

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

Acquisition functions

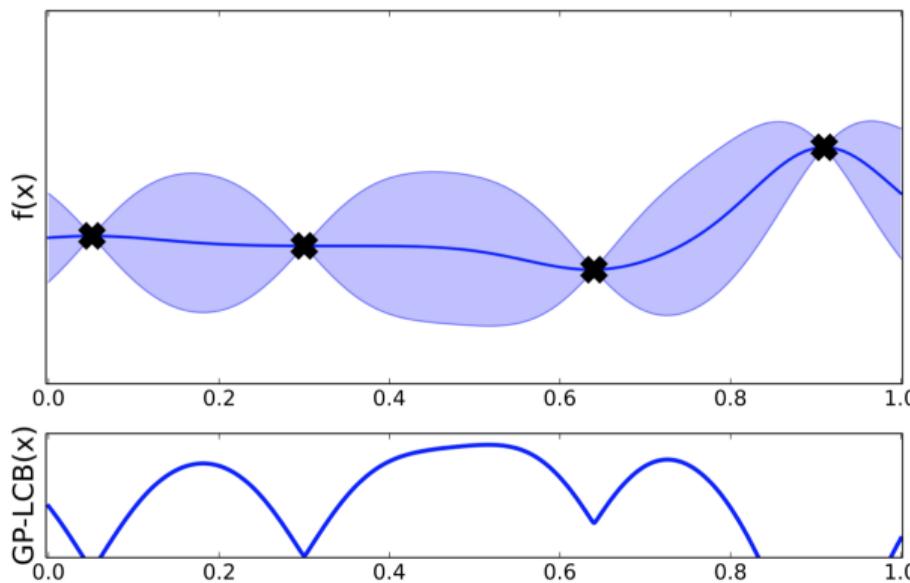
- ▶ Use GP $\mathcal{GP}(\cdot | \mu(\mathbf{x}), K(\mathbf{x}, \mathbf{x}'))$ as a prior for $f(\cdot)$
- ▶ GP has marginal closed-form for the posterior mean $\mu_*(\mathbf{x})$ and variance $\sigma_*^2(\mathbf{x}) \Rightarrow$ efficient calculation of acquisition function
 - ▶ **Exploration:** Evaluate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

GP Upper (lower) Confidence Band

Direct balance between exploration and exploitation (ζ is a user-defined parameter):

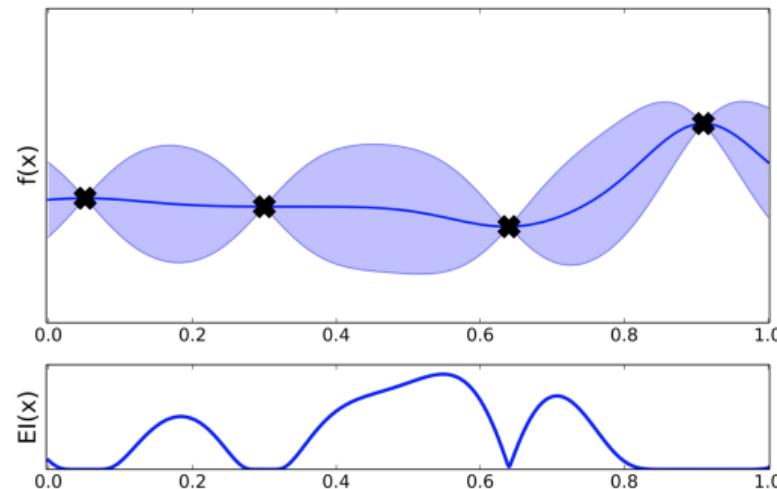
$$\alpha_{LCB}(\mathbf{x}) = -\mu_*(\mathbf{x}) + \zeta \cdot \sigma_*(\mathbf{x})$$



Expected Improvement

Let us denote by $\Delta(\mathbf{x}) = y_{\text{best}} - \mu_*(\mathbf{x})$, $y_{\text{best}} = \min_{i=1,\dots,m} y_i$, then

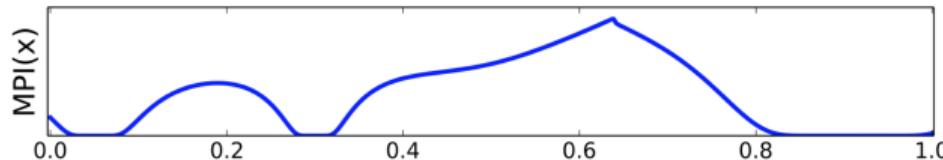
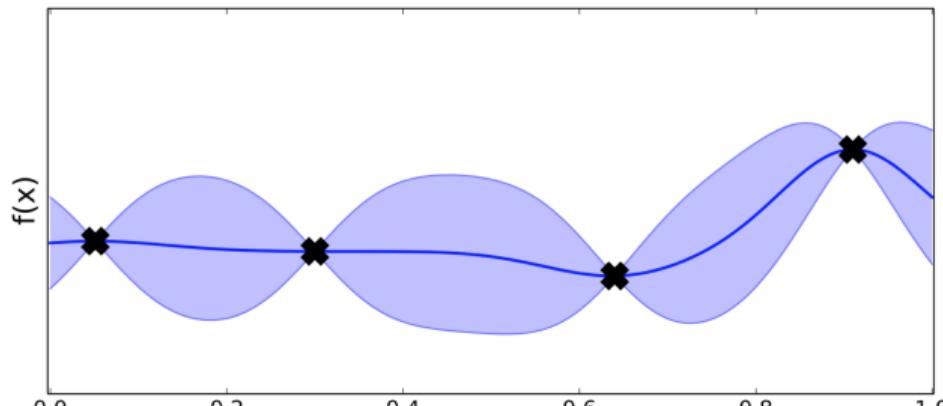
$$\begin{aligned}\alpha_{EI}(\mathbf{x}) &= \int \max(0, y_{\text{best}} - y_*) p(y_* | \mathbf{x}) dy_* = \\ &= \Delta(\mathbf{x}) \Phi(-\Delta(\mathbf{x})/\sigma_*(\mathbf{x})) + \sigma_*(\mathbf{x}) \varphi(\Delta(\mathbf{x})/\sigma_*(\mathbf{x}))\end{aligned}$$



Maximum Probability of Improvement

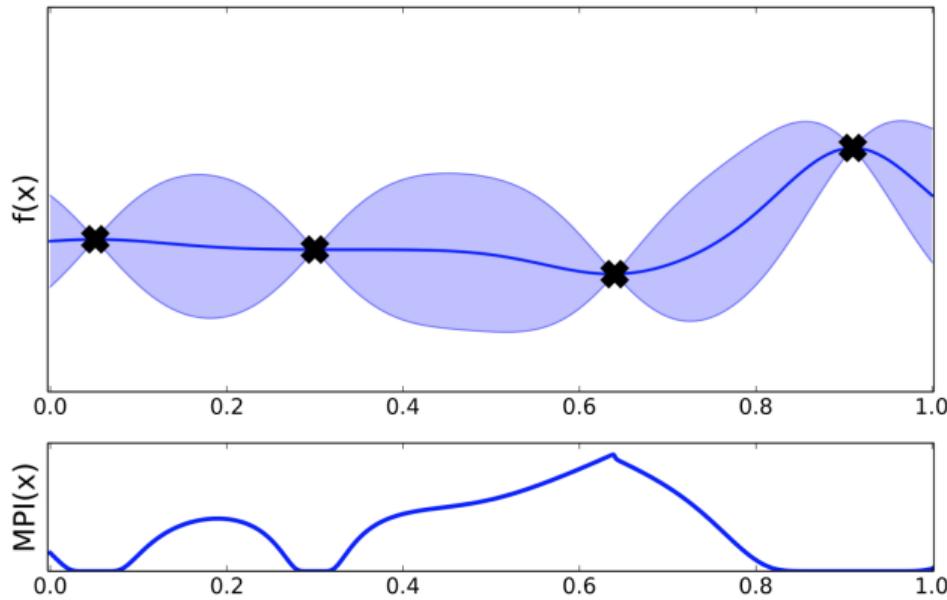
$$\gamma(\mathbf{x}) = \frac{\mu(\mathbf{x}) - y_{\text{best}}}{\sigma(\mathbf{x})}$$

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



Information-theoretic approaches

$$\alpha_{ES}(x) = \mathcal{H}[p(x_{\min} | y)] - \mathbb{E}_{p(y|y,x)}[\mathcal{H}[p(x_{\min} | y \cup \{x, y\})]]$$



Bayesian Optimization

- BO is an strategy to transform the problem

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

unsolvable!

into a series of problems

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

solvable!

where

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

AutoML example: SVM hyperparameters tuning

- ▶ Scikit-learn dummy data for testing classifiers:
`n_samples=2500, n_features=45, n_informative=15, n_redundant=5`
- ▶ Optimize w.r.t. $\mathbf{x} = (C, \gamma)$, where C — penalization, and γ — kernel width
- ▶ The target function $\mathcal{L}(\mathbf{x})$ is the AUC based on three fold cross-validation

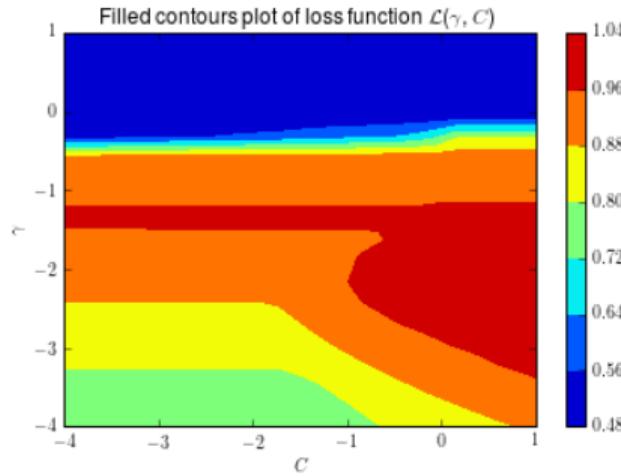
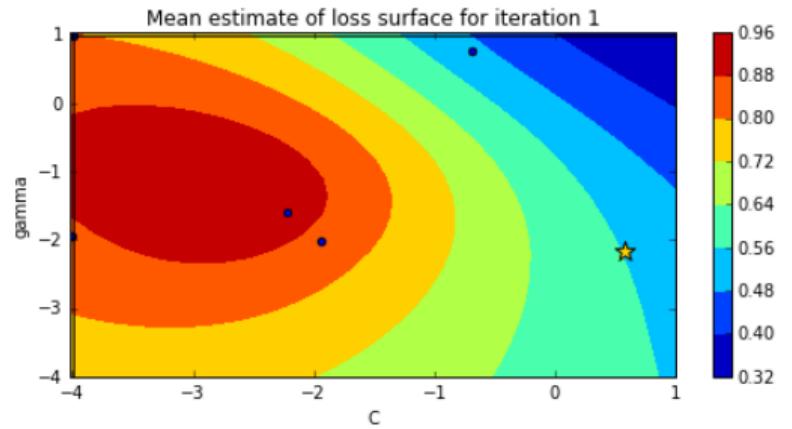
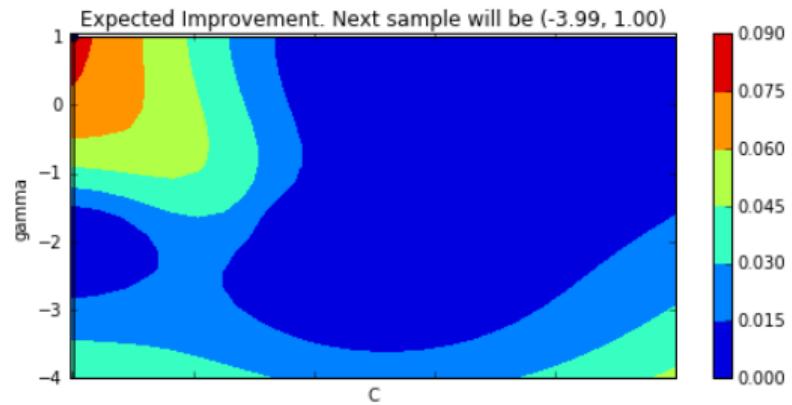
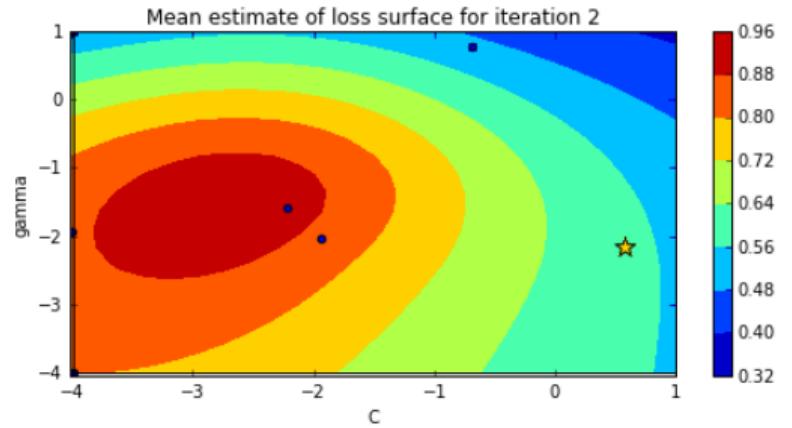
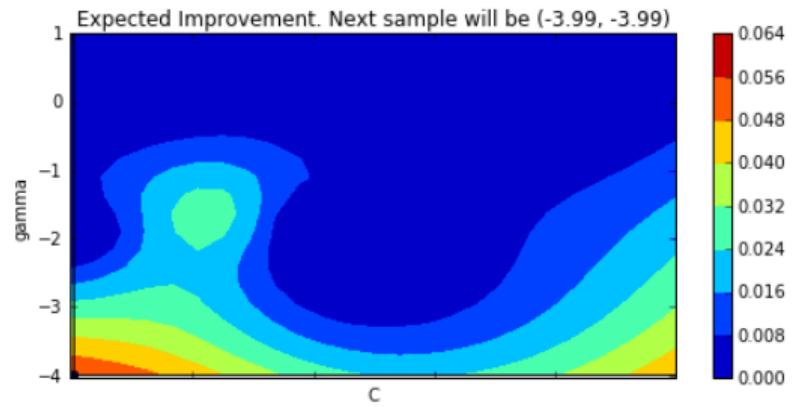


Figure: The target surface w.r.t. \mathbf{x} to see where the true optimum is

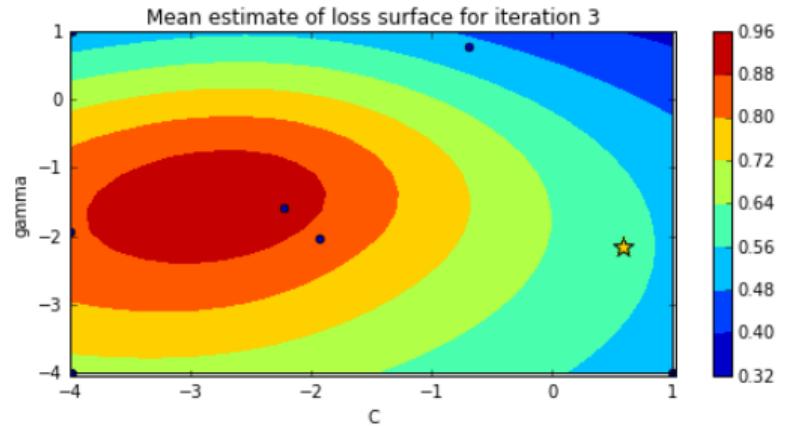
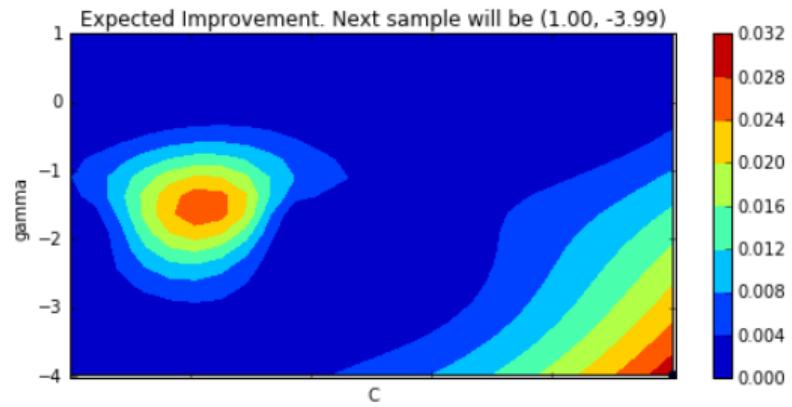
Iterations



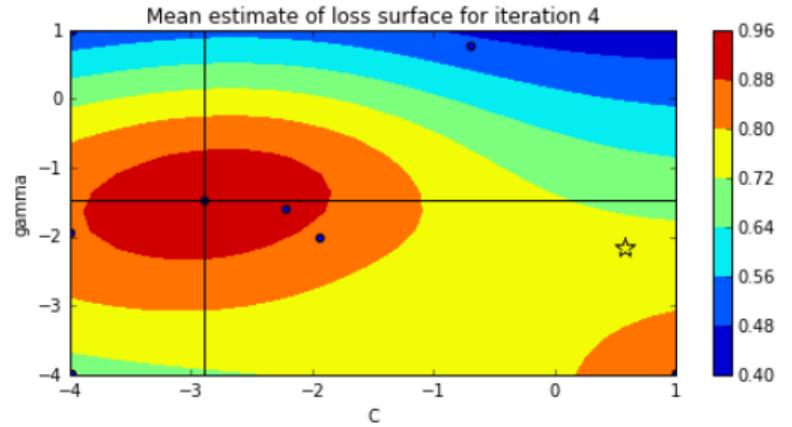
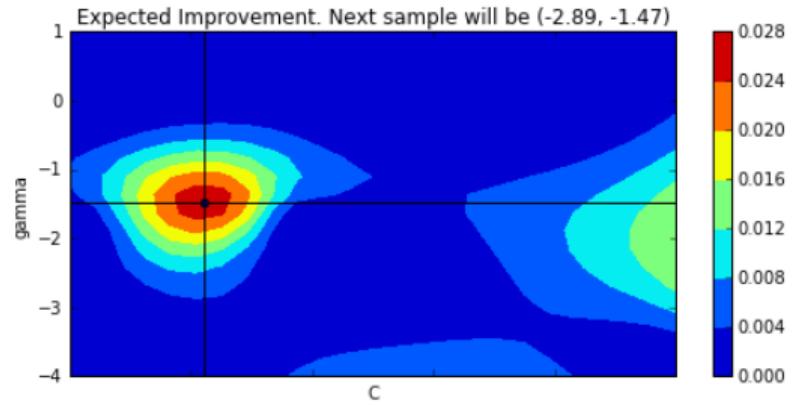
Iterations



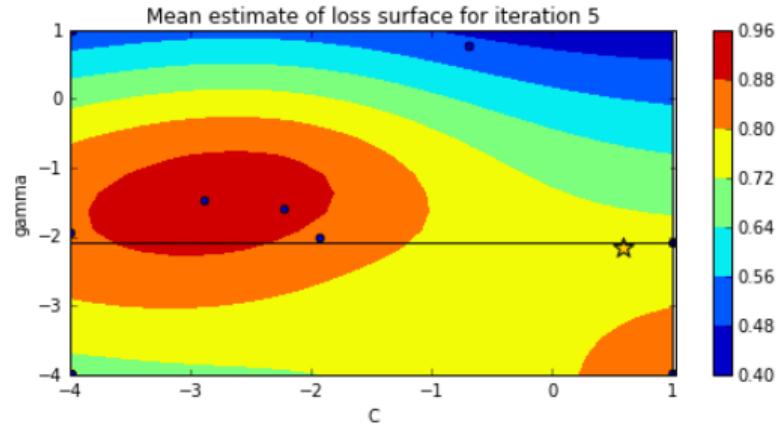
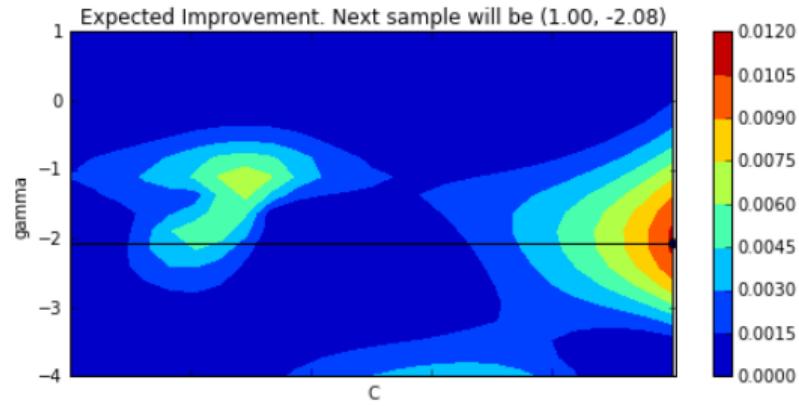
Iterations



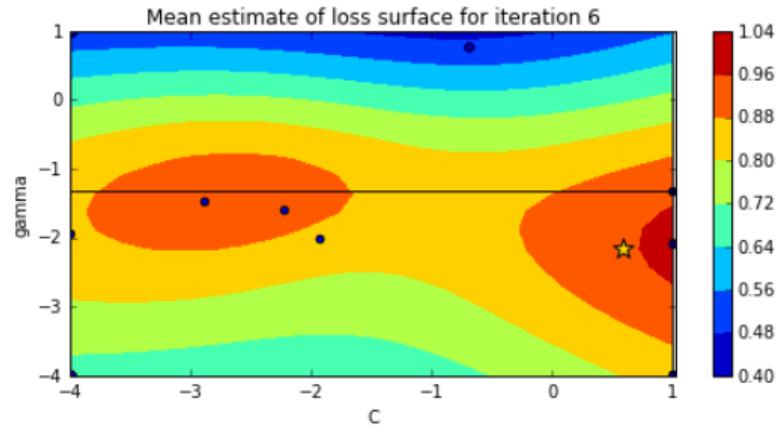
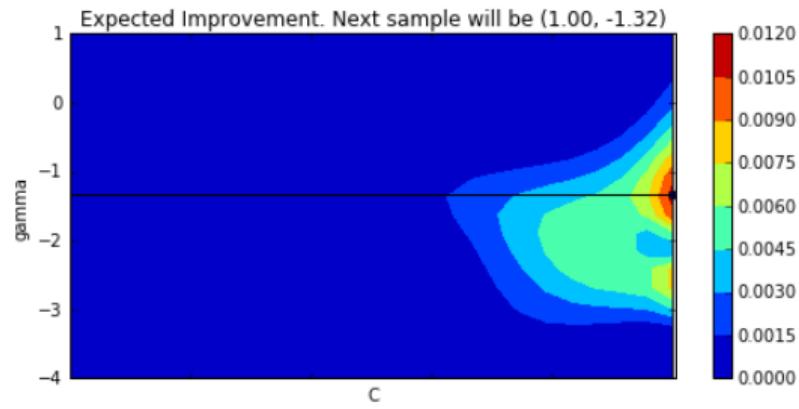
Iterations



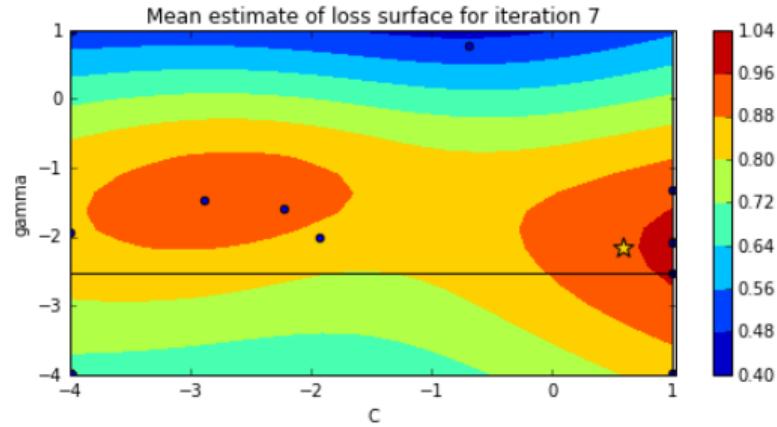
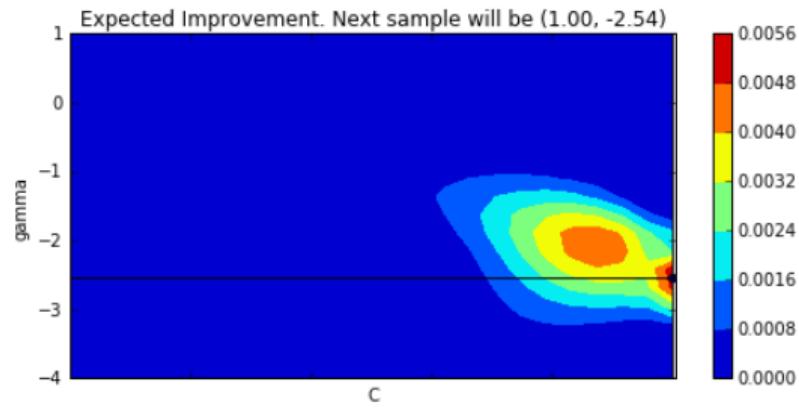
Iterations



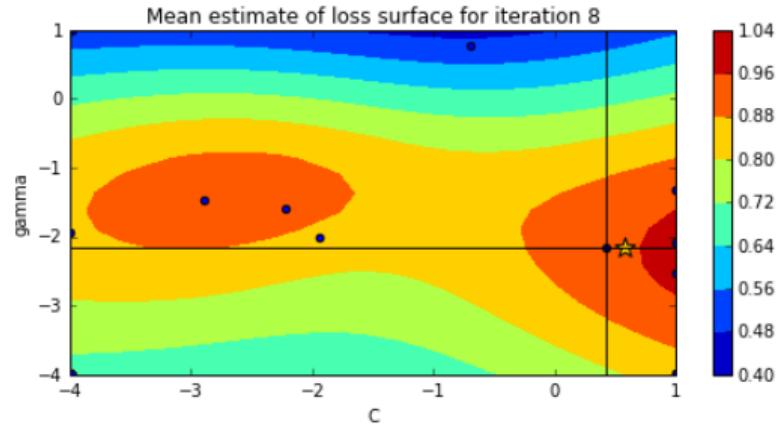
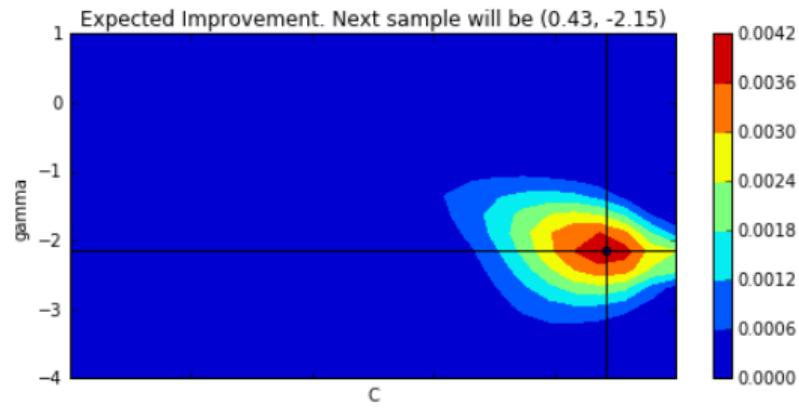
Iterations



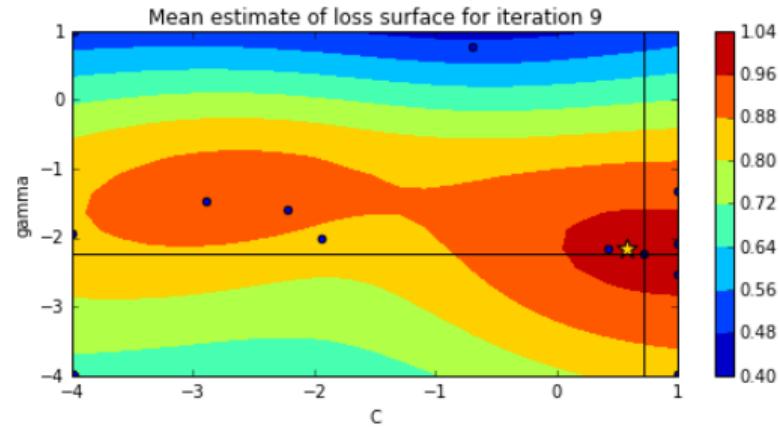
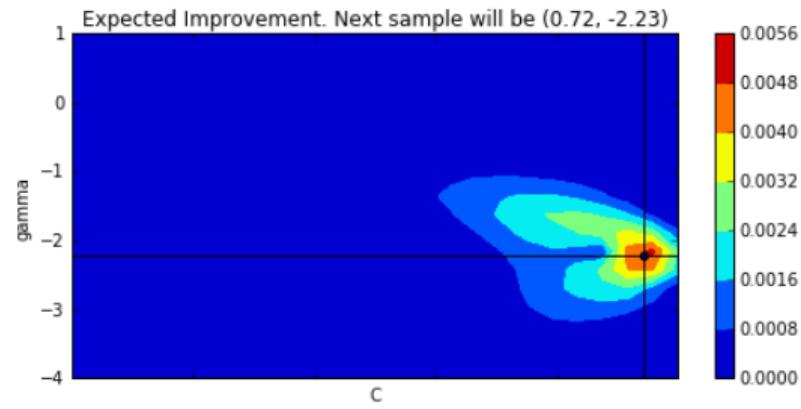
Iterations



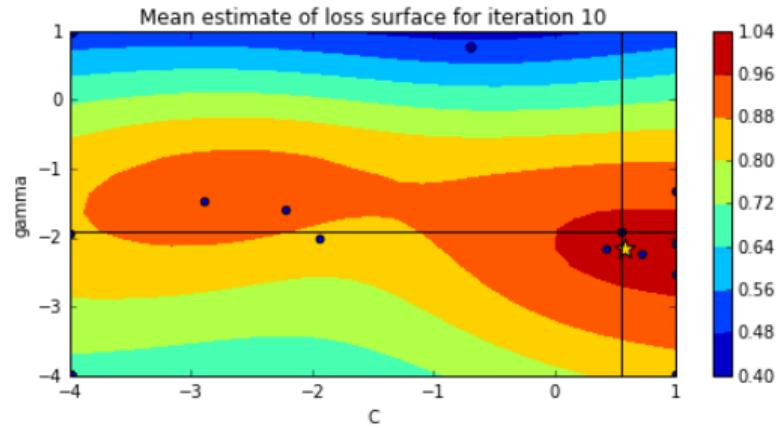
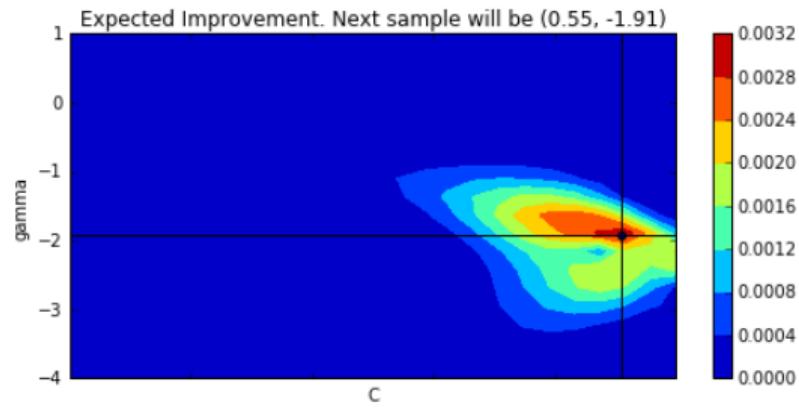
Iterations



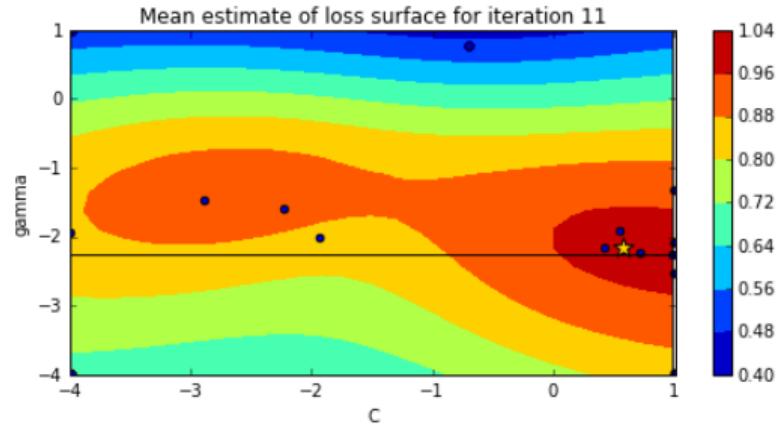
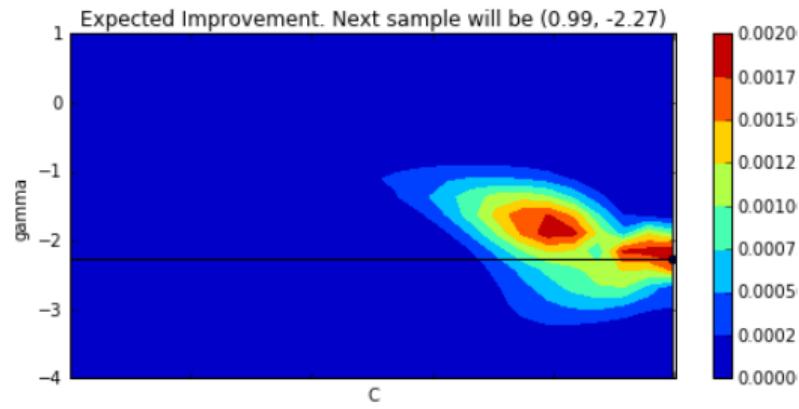
Iterations



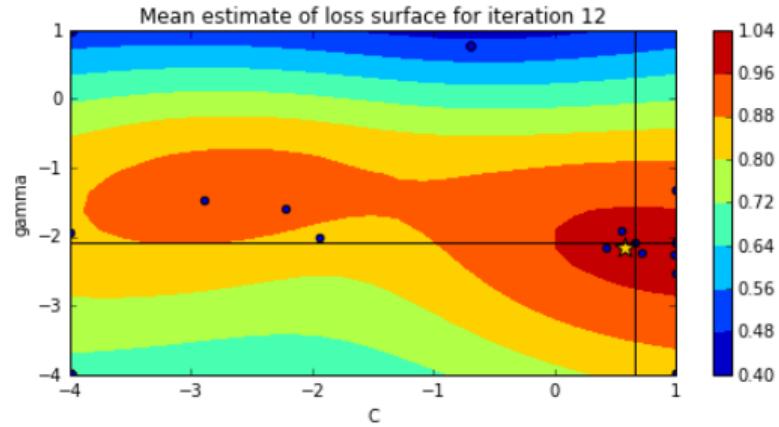
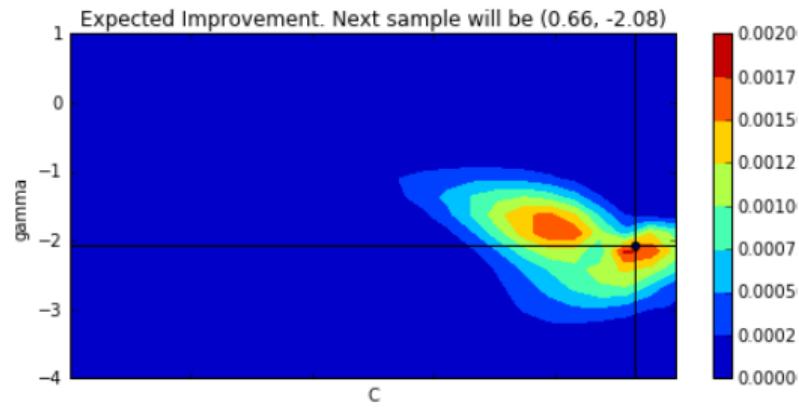
Iterations



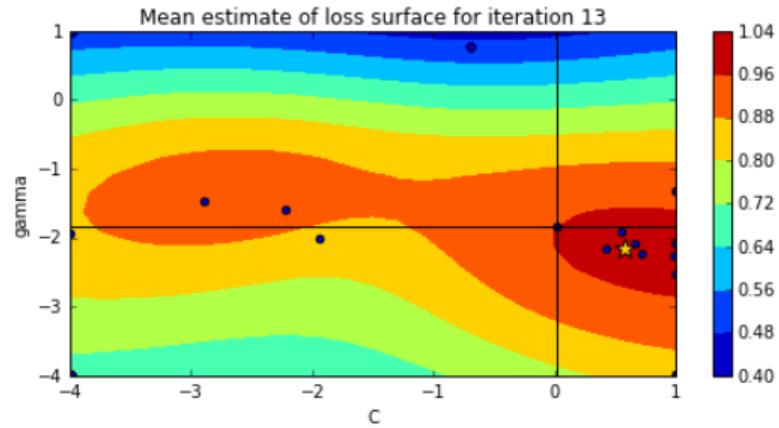
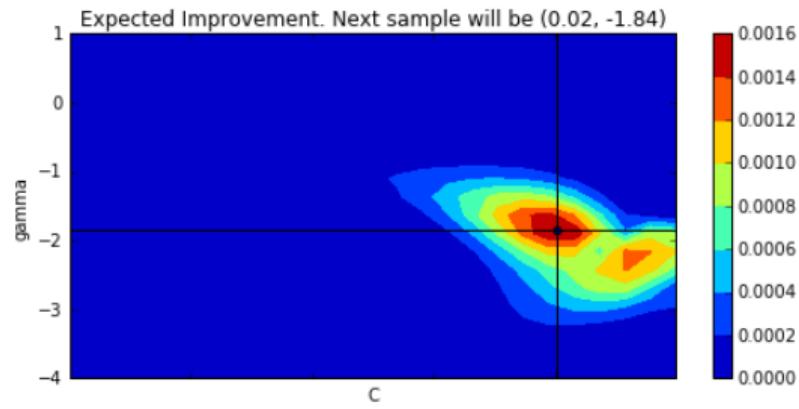
Iterations



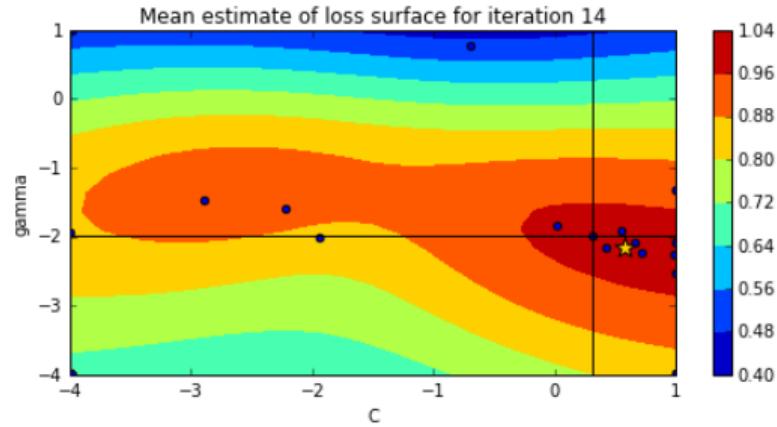
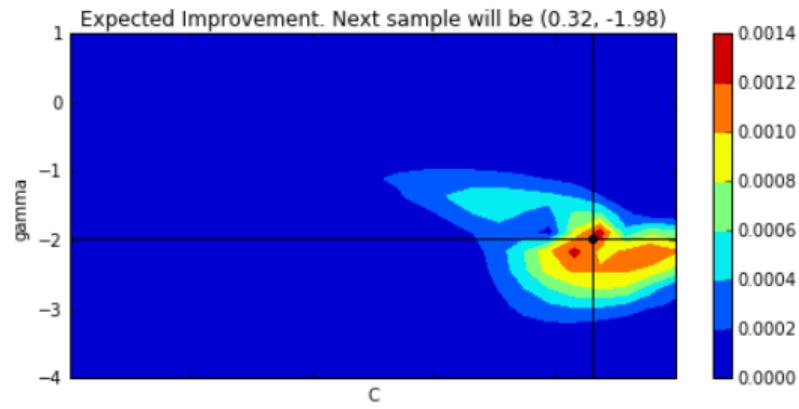
Iterations



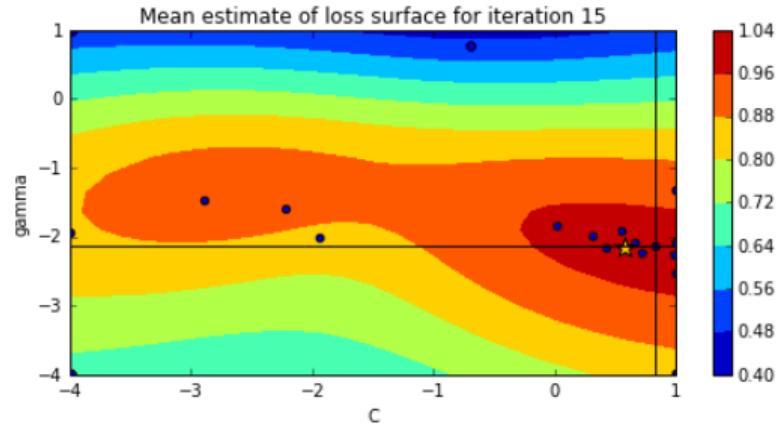
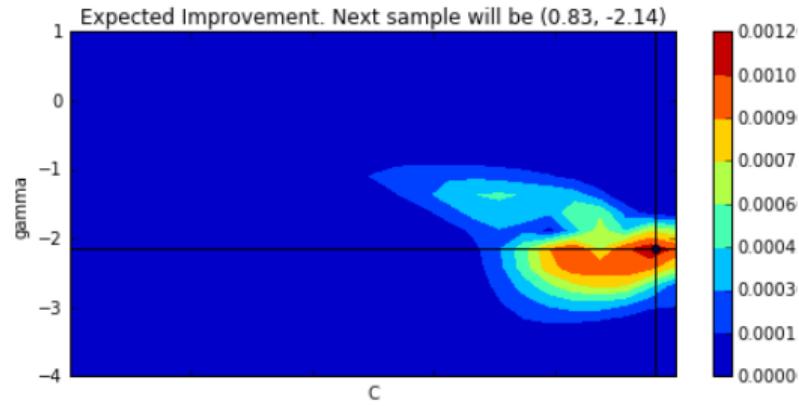
Iterations



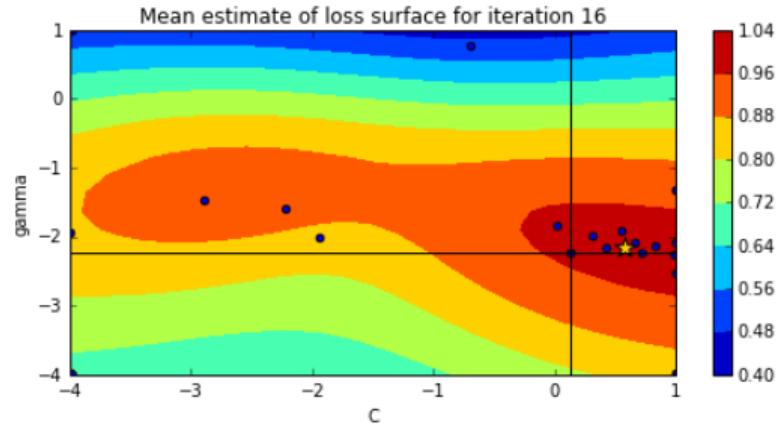
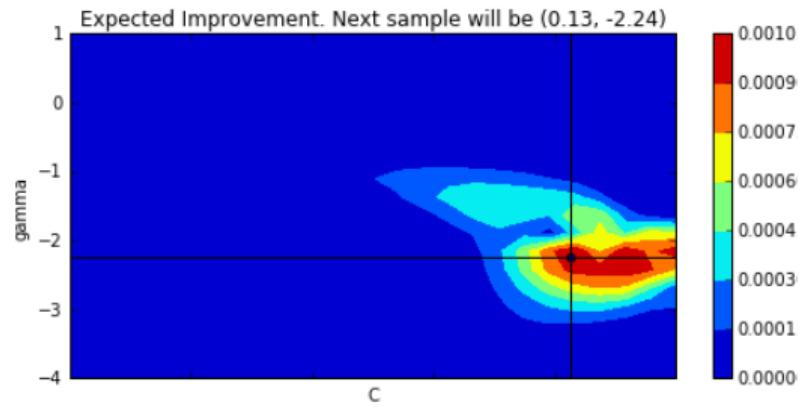
Iterations



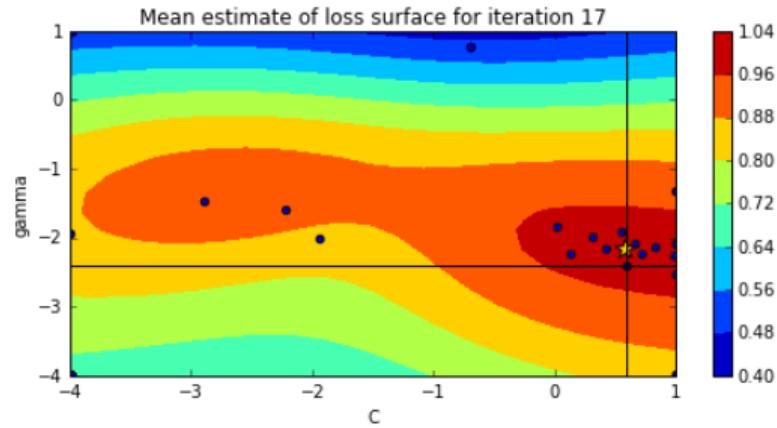
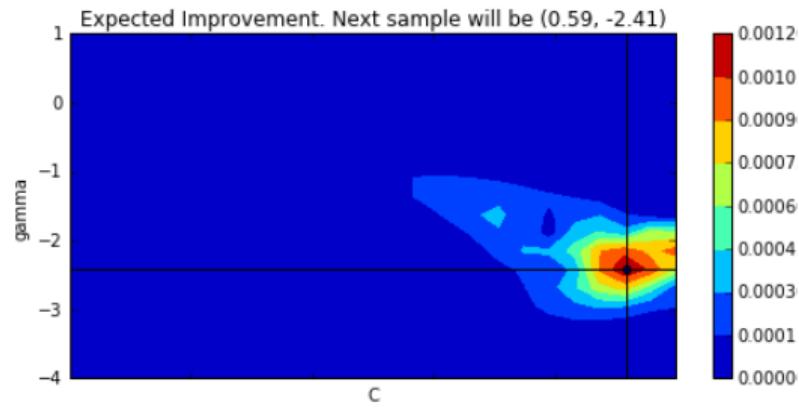
Iterations



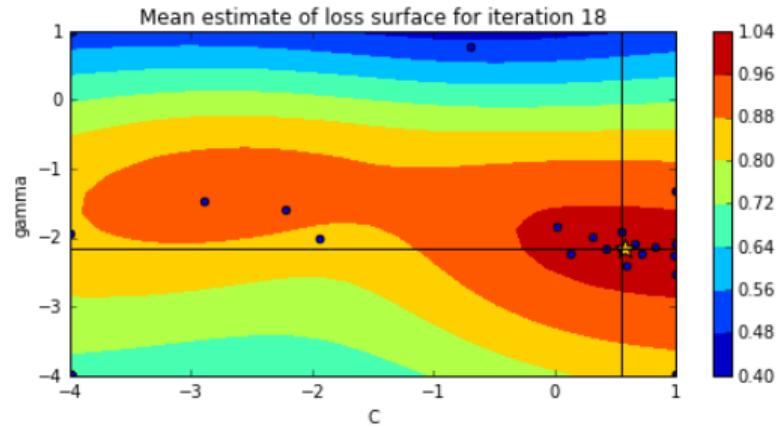
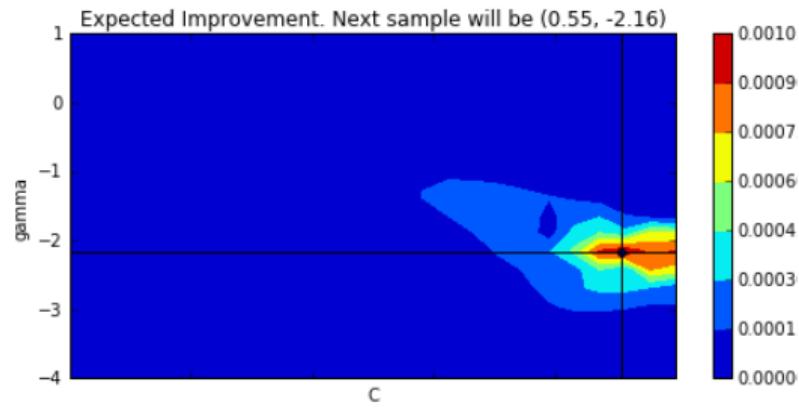
Iterations



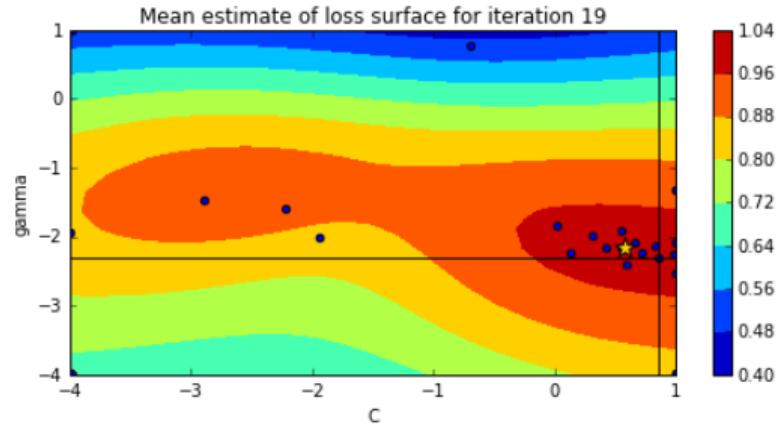
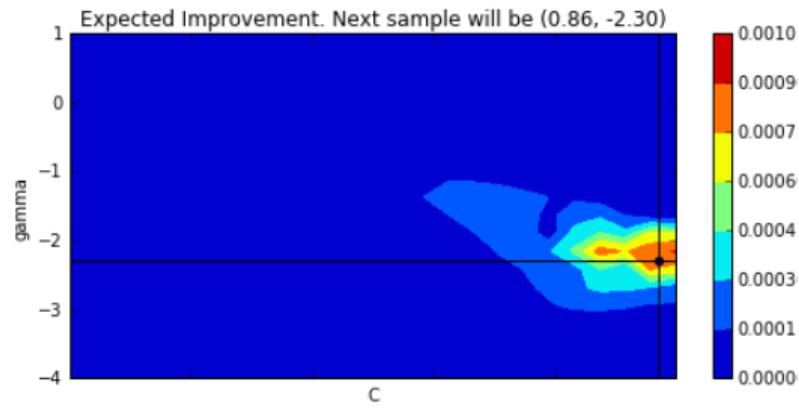
Iterations



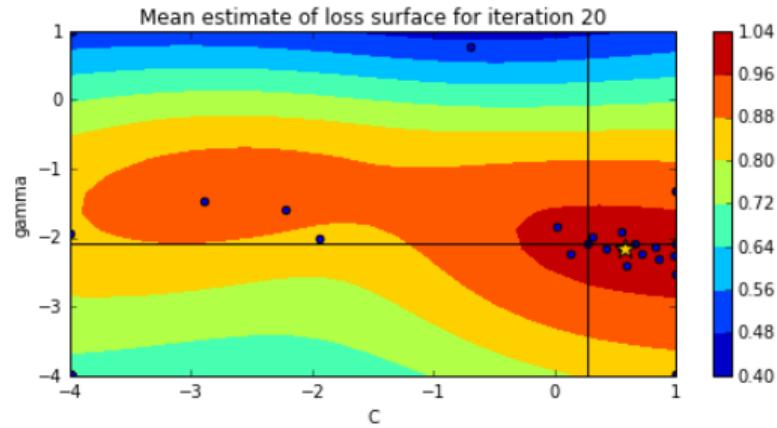
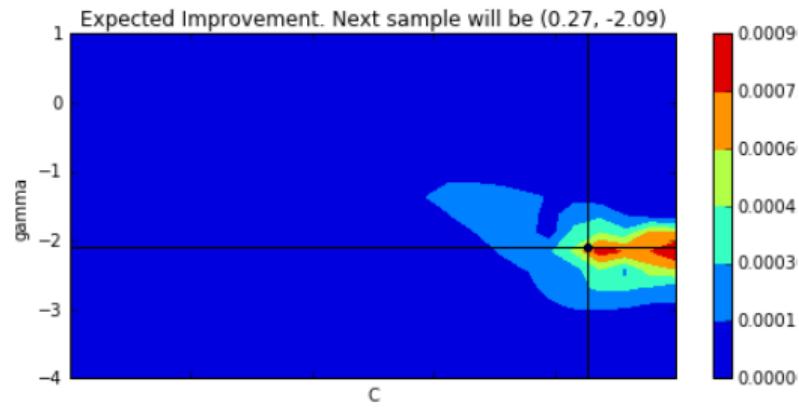
Iterations



Iterations



Iterations



Motivation

Gaussian Process model

GP regression

Learning Gaussian Process model

Bayesian Optimization

Take-home messages

Afterword

Bibliography

My message...



About GP

- ▶ **Simple to use**

- Just matrix operations (if likelihoods are Gaussian)
- Few parameters: relatively easy to set or sample over
- Predictions are often very good
- Handle uncertainty in unknown function f
- Can learn kernel parameters automatically from data
- Can incorporate interpretable noise models and priors over functions
- Can combine automatic feature selection with learning using ARD

- ▶ **No magic bullet**

- ▶ **The need for approximate inference:**

- Sometimes Gaussian likelihoods aren't enough
- $O(m^3)$ and $O(m^2)$ costs are bad news for big problems

What's next?

- ▶ Multi-task and structured GP learning
- ▶ Large scale GP based on Random Fourier features
- ▶ **Deep GP based on random Fourier features**
- ▶ **Deep GP with Deep NN input features**
- ▶ Multi-fidelity and multi-criteria Bayesian optimization
- ▶ Learning to optimize
- ▶ Other probabilistic models, e.g. modeling and optimization based on point processes

Motivation

Gaussian Process model

GP regression

Learning Gaussian Process model

Bayesian Optimization

Take-home messages

Afterword

Bibliography

What else?



Kernel features

- ▶ A lot of kernels can be represented as

$$K(\mathbf{x}, \mathbf{x}') = \int_{\mathbb{R}^d} \phi(\mathbf{x}, \mathbf{w}) \phi(\mathbf{x}', \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$$

- ▶ Thus we can approximate kernels using random sampling

$$K(\mathbf{x}, \mathbf{x}') \approx \frac{1}{k} \sum_{i=1}^k \phi(\mathbf{x}, \mathbf{w}_i) \phi(\mathbf{x}', \mathbf{w}_i) = \langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle, \quad \mathbf{w}_i \sim p(\mathbf{w})$$

with a finite dimensional features $\psi(\mathbf{x})$

- ▶ The problem reduces to linear and we get approximate non-linear kernel models in linear time $O(m)$

Degenerate behaviour

Let us consider the squared exponential covariance function

$$K(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp \left(- \sum_{i=1}^d \theta_i^2 (x_i - x'_i)^2 \right)$$

When optimizing the likelihood wrt covariance function parameters $\{\boldsymbol{\theta} = (\theta_1, \dots, \theta_d), \sigma_f, \sigma\}$, the following weird behavior can occur

- ▶ MLE optimum is located in the area where $\|\boldsymbol{\theta}\| \sim 0$, and the conditional number of \mathbf{K} is big
- ▶ For $\|\boldsymbol{\theta}\| \rightarrow \infty$ matrix $\mathbf{K} \rightarrow \mathbf{I}_m$. Thus, we have a degenerate approximation
- ▶ Dependence of the likelihood $\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta})$ on $\boldsymbol{\theta}$ can be weak

Degenerate approximation

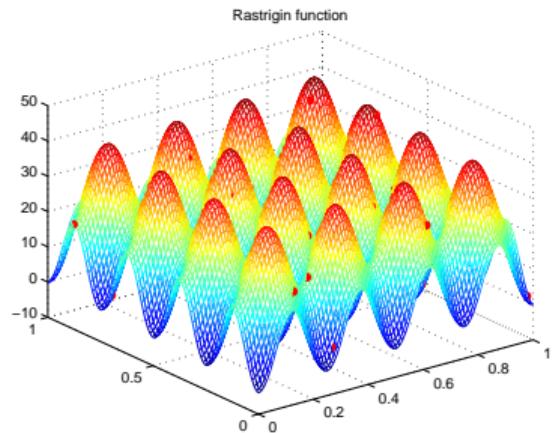


Figure: Rastrigin function

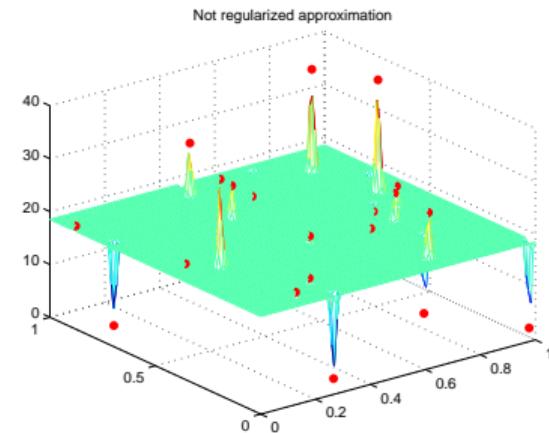


Figure: An approximation of Rastrigin function

Numerical noise

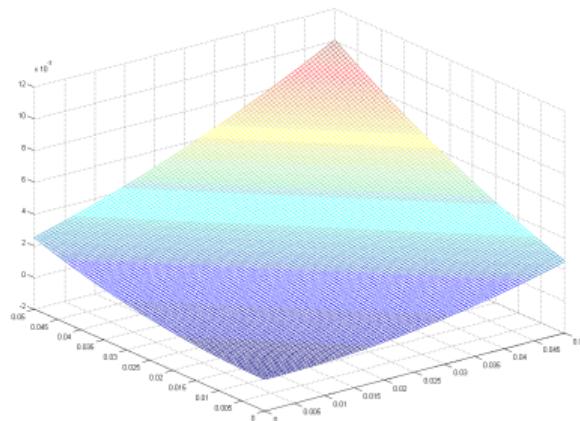
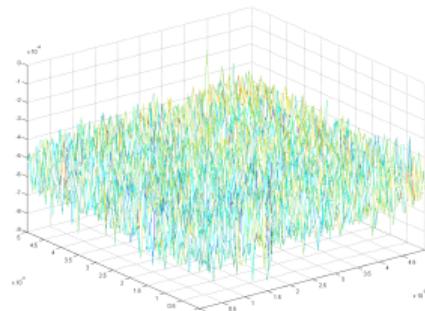
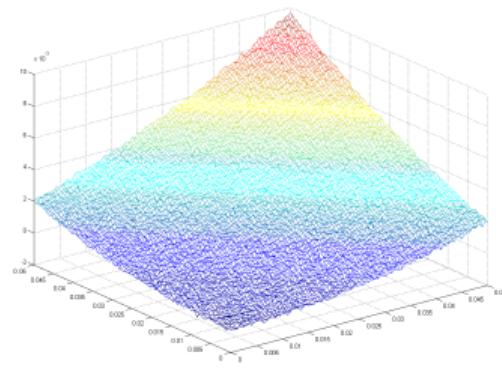


Figure: $y(\mathbf{x}) = (x_1 + x_2)^2$



Bayesian regularization

Let us introduce the following priors on covariance function parameters:

$$\theta_i^2 \in \Gamma(\alpha, \beta), i = 1, \dots, d, \quad \sigma_f^2 \in \Gamma(\alpha_\sigma, \beta_\sigma) \quad (1)$$

Thus log-prior of parameters has the form

$$\begin{aligned} \log p(\{\boldsymbol{\theta}, \sigma_f\}) &= m(\alpha \log \beta - \log \Gamma(\alpha)) + 2(\alpha - 1) \sum_{i=1}^m \log \theta_i + \\ &- \beta \sum_{i=1}^m \theta_i^2 + (\alpha_\sigma \log \beta_\sigma - \log \Gamma(\alpha_\sigma)) + 2(\alpha_\sigma - 1) \log \sigma_f - \beta_\sigma \sigma_f^2 \end{aligned}$$

We obtain MAP estimates for $\{\boldsymbol{\theta}, \sigma_f\}$:

$$\begin{aligned} &\log p(\mathbf{y}|\mathbf{X}, \{\boldsymbol{\theta}, \sigma_f\}) + \log p_{\alpha, \beta, \alpha_\sigma, \beta_\sigma}(\{\boldsymbol{\theta}, \sigma_f\}) \\ &+ \log p(\alpha, \beta, \alpha_\sigma, \beta_\sigma) \rightarrow \max_{\{\boldsymbol{\theta}, \sigma_f\}, \{\alpha, \beta, \alpha_\sigma, \beta_\sigma\}} \end{aligned}$$

Here parameters $\{\alpha, \beta, \alpha_\sigma, \beta_\sigma\}$ are also optimized along with imposed hyper-prior $p(\alpha, \beta, \alpha_\sigma, \beta_\sigma)$, parameters of which are fixed to some constants

Regularized approximation vs. degeneracy

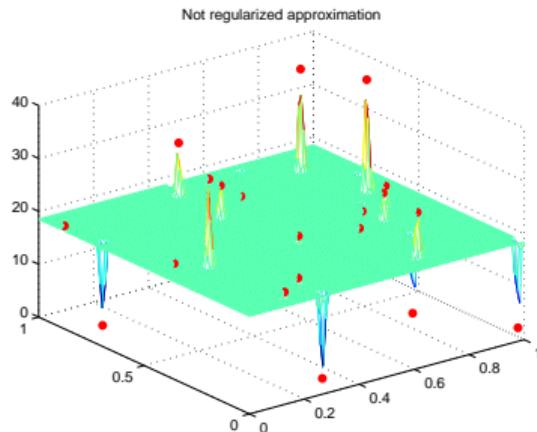


Figure: Degenerate approximation

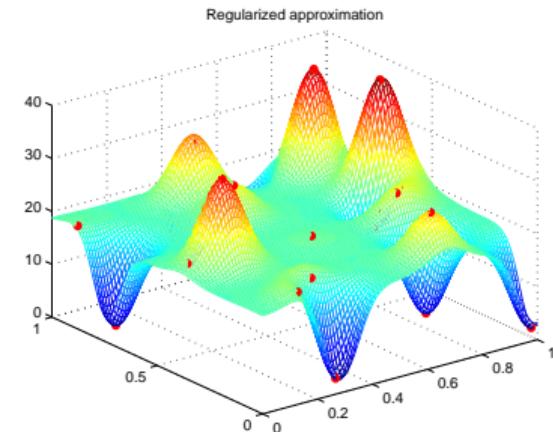
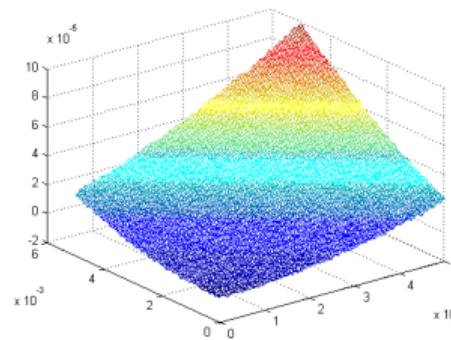
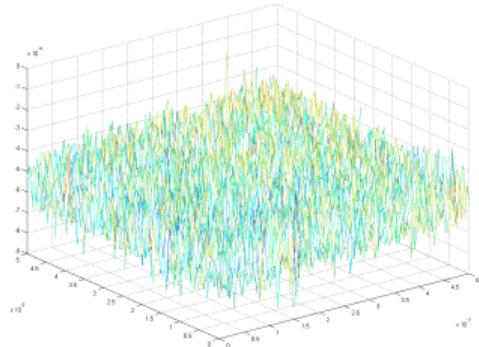
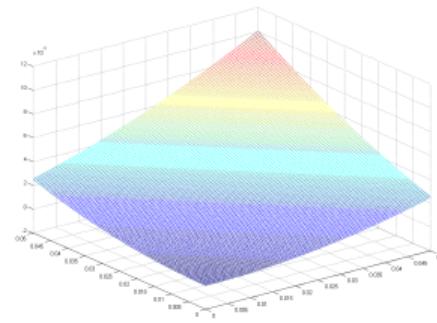
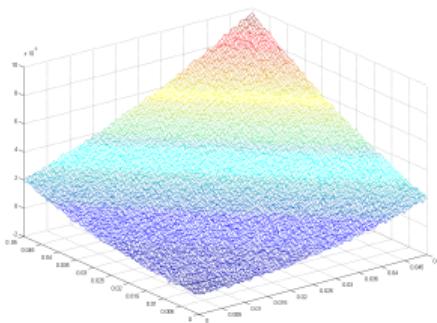


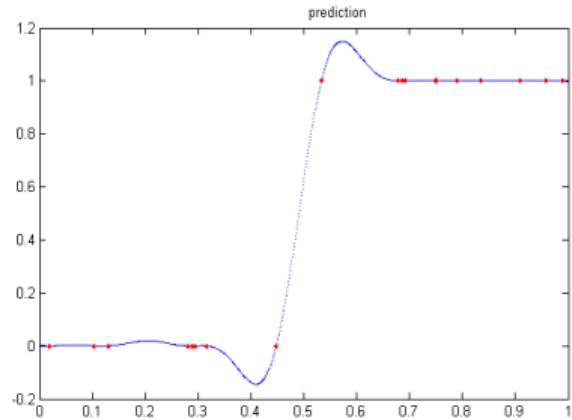
Figure: Approximation with imposed Bayesian regularization

Regularized approximation vs. numerical noise



Approximation of Heaviside function

Approximation of Heaviside function for $x \in [0, 1]$, training sample size $m = 20$



Non-stationary Gaussian process

Let us assume that

$$y(\mathbf{x}) = \sum_{i=1}^Q \alpha_i \psi_i(\mathbf{x}) + f(\mathbf{x}) + \varepsilon(\mathbf{x}),$$

where

- ▶ $f(\mathbf{x})$ is a GP,
- ▶ $\varepsilon(\mathbf{x})$ is a Gaussian white noise,
- ▶ $\{\alpha_i\}_{i=1}^Q$ are i.i.d. r.v. with zero mean and variance $\frac{\sigma_0^2}{Q}$,
- ▶ $\{\psi_i(\mathbf{x})\}_{i=1}^Q$ is a set of functions

Thus the covariance function of $y(\mathbf{x})$ has the form:

$$\bar{K}(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}, \mathbf{x}') + \sigma_0^2 \psi(\mathbf{x}) \psi(\mathbf{x}') + \sigma^2 \delta(\mathbf{x}, \mathbf{x}'),$$

where $\psi(\mathbf{x}) = \{\psi_i(\mathbf{x}), i = 1, \dots, Q\}$

Posterior parameters

Posterior mean of $y(\mathbf{x})$ has the form:

$$\mu_*(\mathbf{x}) = \bar{\mathbf{k}}(\mathbf{x})\bar{\mathbf{K}}^{-1}\mathbf{y},$$

where

- ▶ $\Psi = \{\psi(\mathbf{x}_i), i = 1, \dots, m\}$,
- ▶ $\bar{\mathbf{k}}(\mathbf{x}) = \mathbf{k}(\mathbf{x}) + \sigma_0^2 \psi^T(\mathbf{x})\Psi$,
- ▶ $\bar{\mathbf{K}} = \mathbf{K} + \sigma_0^2 \Psi^T \Psi + \sigma^2 \mathbf{I}_m$

Posterior variance has the form:

$$\sigma_*^2(\mathbf{x}) = K(\mathbf{x}, \mathbf{x}) + \sigma^2 + \sigma_0^2 \psi^T(\mathbf{x})\psi(\mathbf{x}) - \bar{\mathbf{k}}^T(\mathbf{x})\bar{\mathbf{K}}^{-1}\bar{\mathbf{k}}(\mathbf{x})$$

Estimation of parameters

Minus log-likelihood has a similar structure compared to a stationary process, but we should replace \mathbf{K} by $\bar{\mathbf{K}}$:

$$-\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \frac{1}{2}\mathbf{y}^T \bar{\mathbf{K}}^{-1} \mathbf{y} + \frac{1}{2} \log |\bar{\mathbf{K}}| - \frac{m}{2} \log 2\pi,$$

where $\{\boldsymbol{\theta}, \sigma_0, \sigma\}$ are parameters

If the set of functions is fixed, then estimates of parameters can be obtained via usual maximization of the log-likelihood

$$\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) \rightarrow \max_{\boldsymbol{\theta}}$$

Dictionaries of function

We consider three families of functions:

1. Sigmoid functions

$$\psi_j(\mathbf{x}) = \sigma \left(\sum_{i=1}^d \beta_{j,i} x_i \right),$$

where $\sigma(x) = \frac{e^x - 1}{e^x + 1}$, $\beta_{j,i} \in \mathbb{R}$ are parameters

2. Radial basis functions (RBF)

$$\psi_j(\mathbf{x}) = \exp \left(-\frac{\|\mathbf{x} - \mathbf{c}_j\|_2^2}{r_j^2} \right),$$

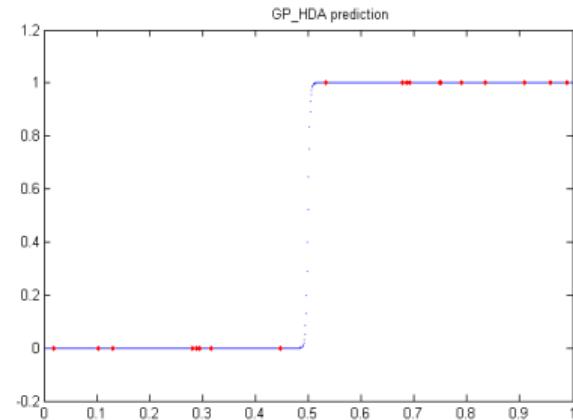
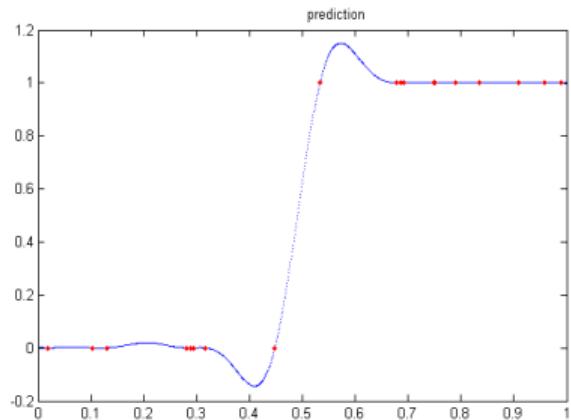
where $\mathbf{c}_j \in \mathbb{R}^d$, $r_j \in \mathbb{R}$ are parameters of the function.

3. Linear basis functions

$$\psi_j(\mathbf{x}) = x_j, j = 1, 2, \dots, d$$

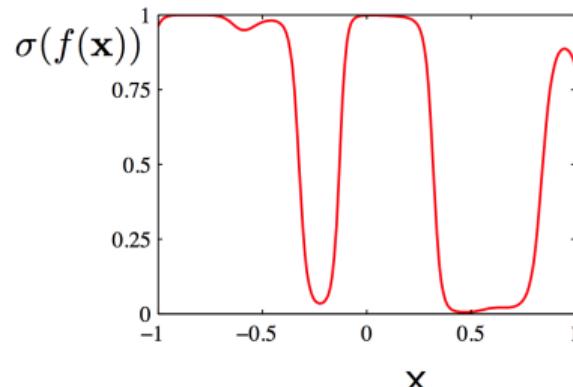
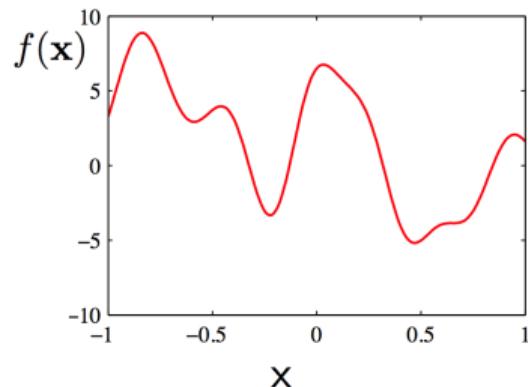
Approximation of Heaviside function

Approximation of Heaviside function for $x \in [0, 1]$, training sample size $m = 20$



Gaussian Process classification

- **Binary classification task:** $y = \pm 1$
- GLM likelihood: $p(y = +1|\mathbf{x}, \mathbf{w}) = \pi(\mathbf{x}) = \sigma(\mathbf{x}^T \mathbf{w})$
- $\sigma(z)$ — **sigmoid** function such as the logistic or cumulative normal
- **Weight space viewpoint:** prior $p(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \Sigma_w)$
- **Function space viewpoint:** let $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$, then likelihood $\pi(\mathbf{x}) = \sigma(f(\mathbf{x}))$, Gaussian prior on \mathbf{f}



GP classification

- ▶ Place a GP prior directly on $f(\mathbf{x})$
- ▶ Use a sigmoidal likelihood: $p(y = +1|f) = \sigma(f)$
- ▶ Non-Gaussian likelihood makes integrating over \mathbf{f} intractable:

$$p(f_*|\mathbf{y}) = \int d\mathbf{f} p(f_*|\mathbf{f})p(\mathbf{f}|\mathbf{y}),$$

where the posterior $p(\mathbf{f}|\mathbf{y}) \sim p(\mathbf{y}|\mathbf{f})p(\mathbf{f})$

- ▶ Make tractable by using a Gaussian approximation to posterior. Then prediction:

$$p(y_* = +1|\mathbf{y}) = \int df_* \sigma(f_*) p(f_*|\mathbf{y})$$

GP classification

Two common ways to make Gaussian approximation to posterior:

1. **Laplace approximation**. Second order Taylor approximation about mode of posterior
 2. **Expectation propagation (EP)**. EP can be thought of as approximately minimizing $\text{KL}[p(\mathbf{f}|\mathbf{y})||q(\mathbf{f}|\mathbf{y})]$ by an iterative procedure
- ▶ Kuss and Rasmussen [2005] evaluate both methods experimentally and find EP to be often better
 - ▶ Classification accuracy on digits data sets comparable to SVMs
 - ▶ **Advantages**: uncertainty quantification, probabilistic predictions, kernel parameters by ML

Motivation

Gaussian Process model

GP regression

Learning Gaussian Process model

Bayesian Optimization

Take-home messages

Afterword

Bibliography

Bibliography: books, papers

- ▶ Forrester, A., Sobester, A., and Keane, A. (2008). Engineering design via surrogate modelling: a practical guide. Wiley, West Sussex, UK.
- ▶ Rasmussen, C. and Williams, C. (2006). Gaussian Processes for Machine Learning. MIT Press, Cambridge, MA.
- ▶ Cressie, N. (1993). Statistics for Spatial Data, revised edition. Wiley Series in Probability and Mathematical Statistics: Applied Probability and Statistics. Wiley Interscience, New York
- ▶ Jones, D., Schonlau, M., and Welch, W. (1998). Efficient Global Optimization of Expensive Black-Box Functions. *Journal of Global Optimization*, 13(4):455-492.
- ▶ Villemonteix, J., Vazquez, E., and Walter, E. (2009). An informational approach to the global optimization of expensive-to-evaluate functions. *Journal of Global Optimization*, 44(4):509534.
- ▶ B. Scholkopf and A.J. Smola. Learning with Kernels, Support Vector Machines, Regularization, Optimization, and Beyond. Adaptive Computation and Machine Learning. The MIT Press, Cambridge, MA, USA, 2002.

Bibliography: tutorials

- ▶ Javier Gonzalez, Gaussian Processes for Global Optimization, 2015
- ▶ Novi Quadrianto, Gaussian Processes, 2017
- ▶ Ryan P. Adams, Tutorial on! Bayesian Optimization for Machine Learning
- ▶ Peter I. Frazier, Tutorial: Bayesian Methods for Global and Simulation Optimization, 2011
- ▶ Marc Deisenroth, Gaussian Processes for Big Data Problems, 2015
- ▶ Ed Snelson, Tutorial: Gaussian process models for machine learning, 2006
- ▶ Zoubin Ghahramani, A Tutorial on Gaussian Processes (or why I don't use SVMs), 2011
- ▶ Iain Murray, Introduction to Gaussian Process, 2008
- ▶ Cedric Archambeau, Bayesian Optimization, 2016

Bibliography: additional reading

- ▶ Burnaev E., Zaytsev A. Minimax approach to variable fidelity data interpolation. Proceedings of Machine Learning Research 54:652-661, Volume 54: AISTATS 2017.
- ▶ Burnaev E., Zaytsev A. Large Scale Variable Fidelity Surrogate Modeling. Ann Math Artif Intell (2017), pp. 1-20. doi:10.1007/s10472-017-9545-y
- ▶ E. Burnaev, M. Belyaev, E. Kapushev. Computationally efficient algorithm for Gaussian Processes based regression in case of structured samples. Comp. Math. and Math. Phys., 2016, Vol. 56, No. 4, pp. 499-513, 2016.
- ▶ E. Burnaev, M. Panov, A. Zaytsev. Regression on the Basis of Nonstationary Gaussian Processes with Bayesian Regularization, Journal of Communications Technology and Electronics, 2016, Vol. 61, No. 6, pp. 661-671.
- ▶ Burnaev E, Nazarov I. Conformalized Kernel Ridge Regression. 15th IEEE International Conference on Machine Learning and Applications (ICMLA), pp. 45 - 52, 2016.
- ▶ E. Burnaev and M. Panov. Adaptive design of experiments based on gaussian processes. LNAI. Proc. of SLDS 2015. Vol. 9047, p. 116-126, 2015. Springer.
- ▶ E. Burnaev, V. Vovk. Efficiency of conformalized ridge regression. JMLR 35:605-622, 2014
- ▶ A. Zaytsev, E. Burnaev, V. Spokoiny. Properties of the Bayesian Parameter Estimation of a Regression Based on Gaussian Processes. Journal of Math. Sci., Vol. 203, Issue 6, pp. 789-798, 15 Nov 2014.