Guest Lecture for the Automatic Machine Learning Course

# Bayesian Optimization

Elena Raponi
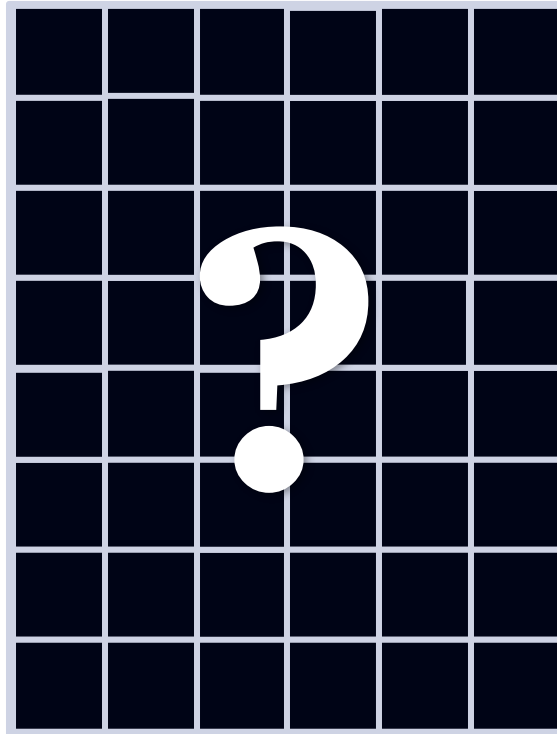
Universiteit Leiden
The Netherlands
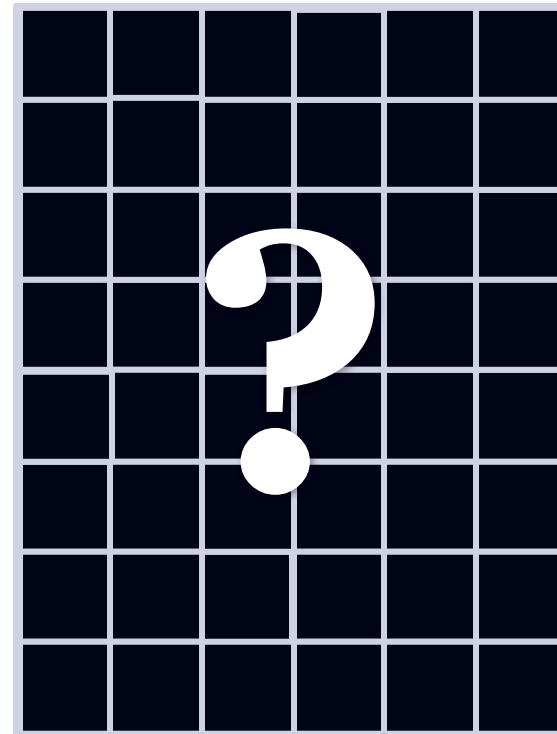
# Model-based optimization– Visual motivation

**Model-based**
Algorithm 1

**Direct search**
Algorithm 2

# Model-based optimization – Visual motivation

**Model-based**
Algorithm 1

€ 5000

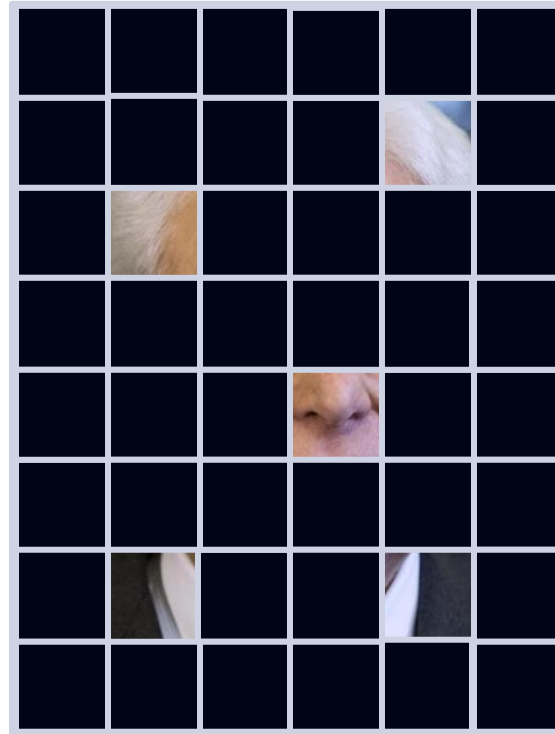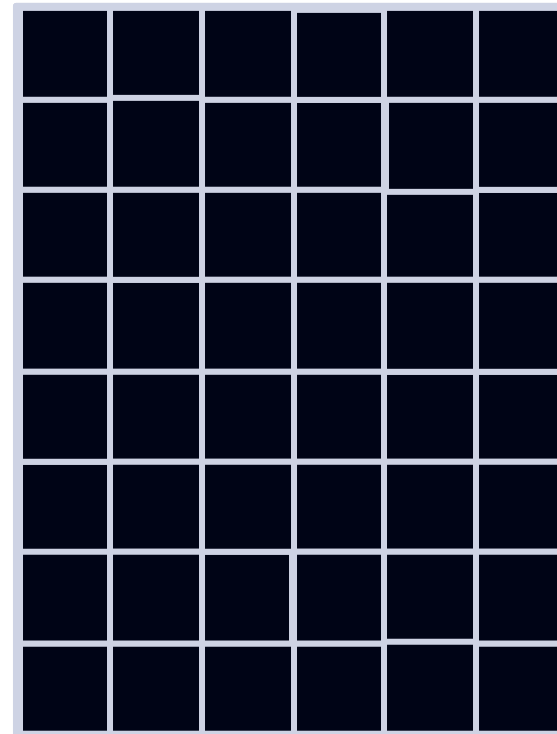**Direct search**
Algorithm 2

# Model-based optimization – Visual motivation



**Model-based**
Algorithm 1

€ 5000

**Direct search**
Algorithm 2

# Model-based optimization – Visual motivation

# Model-based optimization – Visual motivation

# Model-based optimization – Visual motivation

# Model-based optimization – Visual motivation



**Model-based**
Algorithm 1

€ 5000

**Direct search**
Algorithm 2

# Model-based optimization – Visual motivation

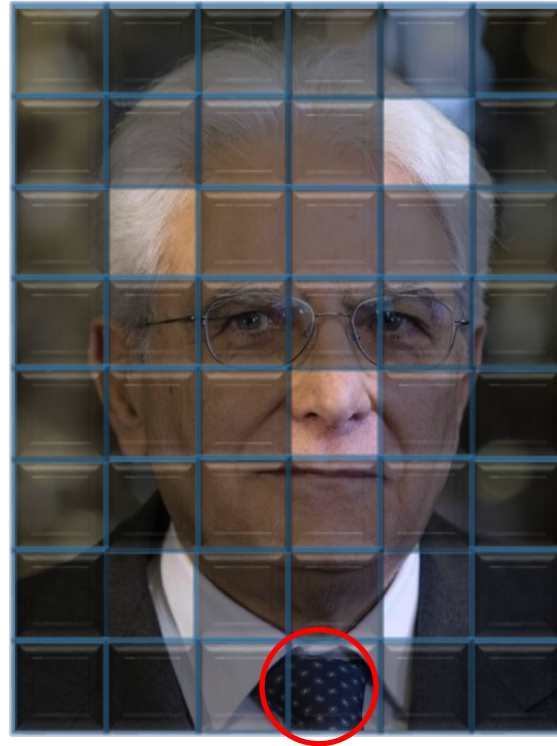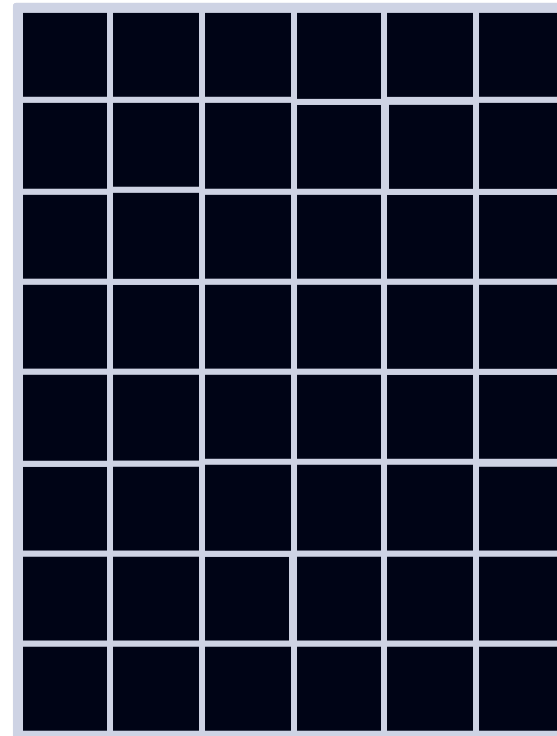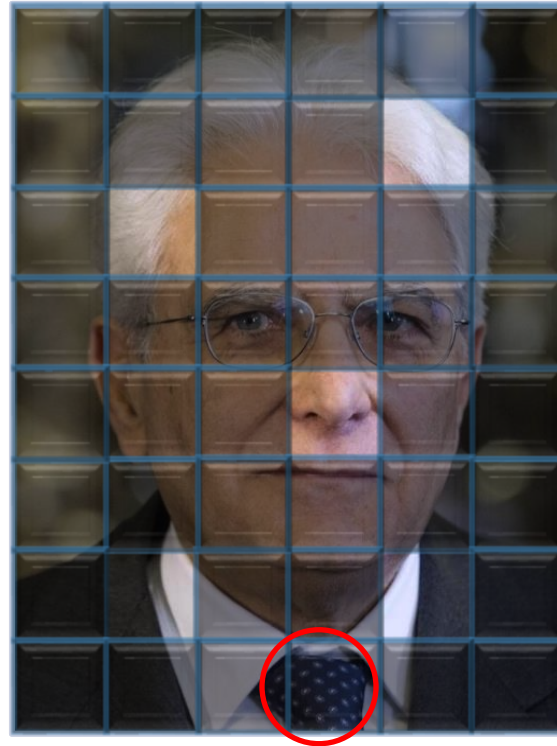# Model-based optimization – Visual motivation



**Model-based**
Algorithm 1

€ 5000

**Direct search**
Algorithm 2

# Model-based optimization – Visual motivation
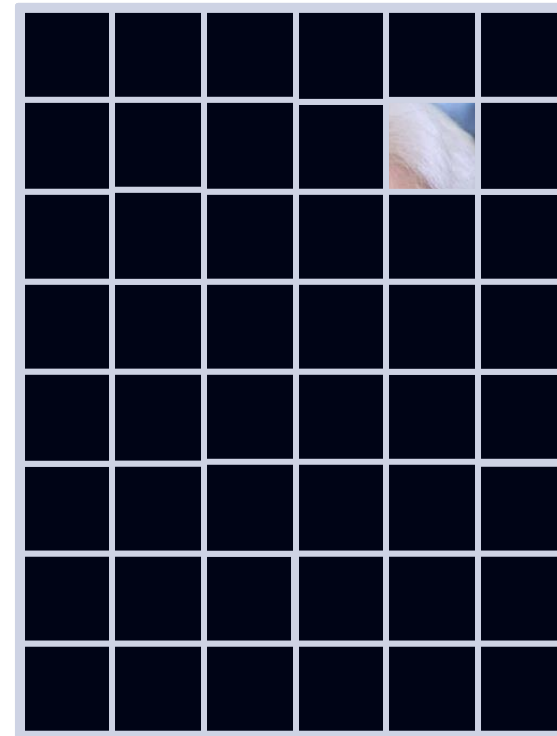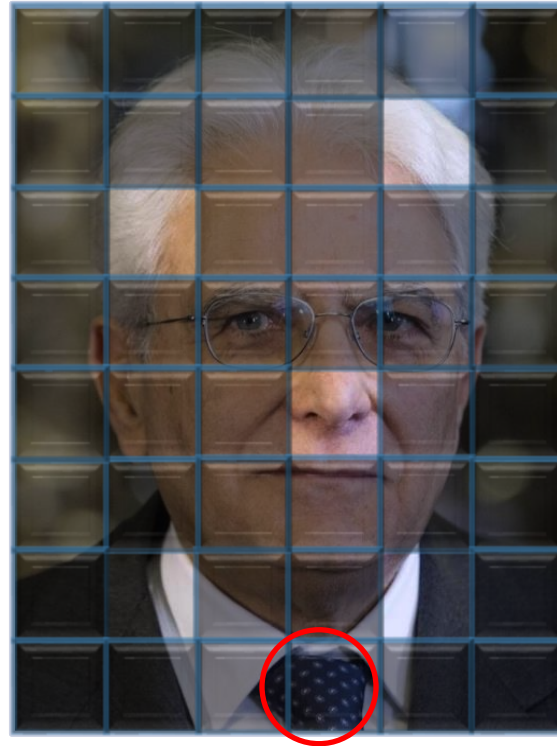
# Model-based optimization – Visual motivation



separable functions

functions with low or moderate conditioning

functions with high conditioning and unimodal structure

multi-modal functions with adequate global structure

multi-modal functions with weak global structure

Remember that not all problems (function landscapes) are friendly!

# The Optimization Challenge

In the context of Bayesian optimization, Gaussian Processes (GPs) are used to model a real-valued function that we want to optimize.

Let $f: X \to \mathbb{R}$ be a real-valued function defined on a domain $X$ (let's assume $X \subset \mathbb{R}^d$ for now)

We want to find a point minimizing $f$:

$$x^* = \operatorname*{argmin}_{x \in X} f(x);$$

$$f^* = \min_{x \in X} f(x) = f(x^*)$$

**Note:** we do not require that the objective function have a known functional form

# The Optimization Challenge

Let $f: X \to \mathbb{R}$ be a real-valued function defined on a domain $X$
(let's assume $X \subset \mathbb{R}^d$ for now)

- Vanilla BO usually performs well for d<20

- f's feasible set is simple, e.g., box constraints

- f is continuous but lacks special structure , e.g., concavity, that would make it easy to optimize

- f is derivative-free: evaluations do not give gradient information

- f is expensive to evaluate: the # of times we can evaluate it is severely limited

- f may be noisy. If noise is present, we'll assume it is independent and normally distributed, with common but unknown variance

$f^*$

$x^*$

# Standard logical flow of BO

**Target**: Find $x^* = \operatorname*{argmin}_{x \in X} f(x)$

Assume a Bayesian **prior** on f (usually a Gaussian process prior)

Build the **posterior** distribution

while (budget is not exhausted) {

- Find x that maximizes **acquisition**(x,**posterior**)
- **Sample** x & **observe** f(x)
- **Update** the **posterior** distribution on f

}

# Optimizing a 1-dimensional problem

Let $f:X\rightarrow\mathbb{R}$ be a real-valued function defined on a domain $X$ (let's assume $X\subset\mathbb{R}$ here)

# Optimizing a 1-dimensional problem

Let $f:X \rightarrow \mathbb{R}$ be a real-valued function defined on a domain $X$ (let's assume $X \subset \mathbb{R}$ here)

# Optimizing a 1-dimensional problem

Let $f: X \to \mathbb{R}$ be a real-valued function defined on a domain $X$ (let's assume $X \subset \mathbb{R}$ here)

# Optimizing a 1-dimensional problem

Let $f : X \to \mathbb{R}$ be a real-valued function defined on a domain $X$ (let's assume $X \subset \mathbb{R}$ here)

# Preliminary step: Initial Sample

Preliminarily to the iterative procedure defining every Bayesian Optimization approach, it is common practice to sample some points in parallel, to have a first insight about the problem.

The initial sample set can be generated through many design schemes.

This initial step is called **Design of Experiments**.



[https://scikit-optimize.github.io/stable/auto_examples/sampler/initial-sampling-method.html]

# Standard logical flow of BO

**Target**: Find $x^* = \underset{x \in X}{\operatorname{argmin}} f(x)$

<mark>Assume a Bayesian **prior** on f (usually a Gaussian process prior)</mark>

<mark>Build the **posterior** distribution</mark>
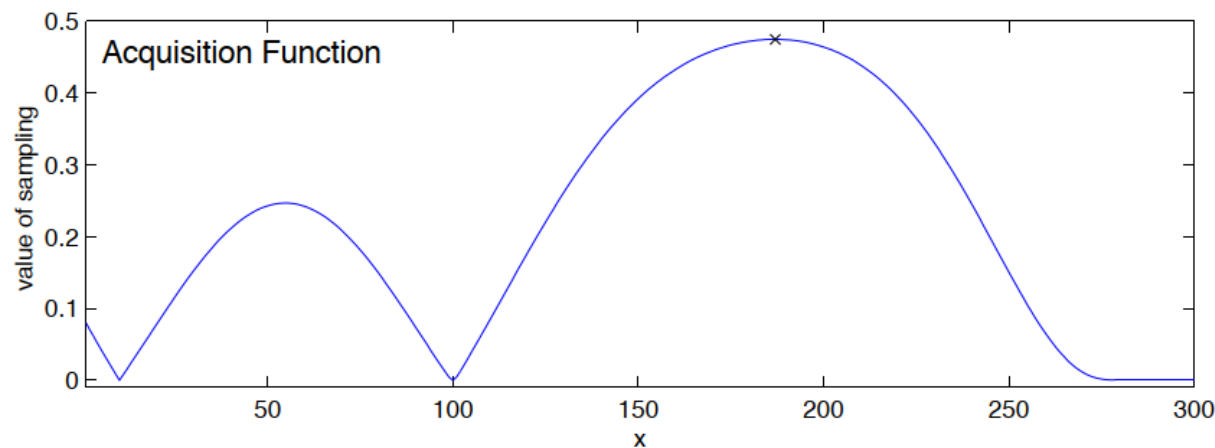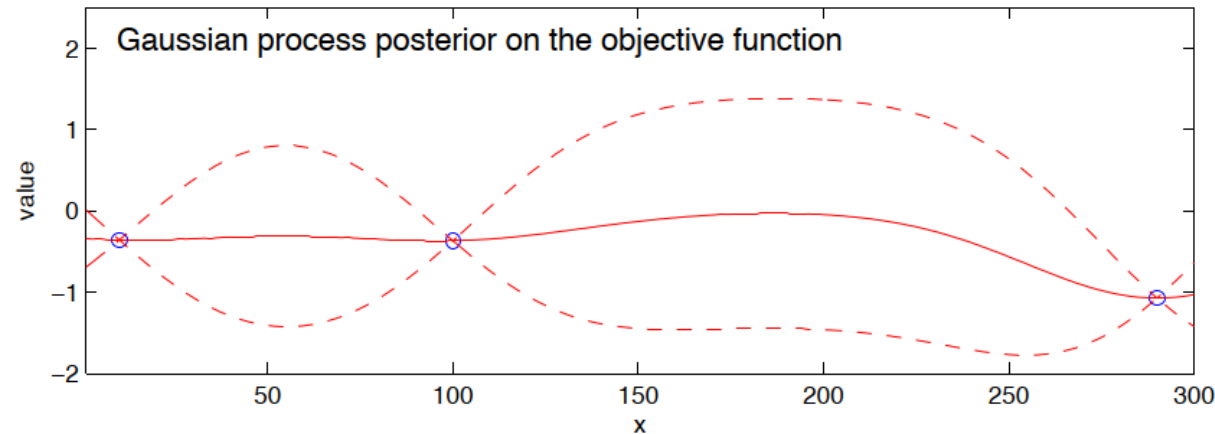
while (budget is not exhausted) {

- Find x that maximizes **acquisition**(x,**posterior**)
- **Sample** x & **observe** f(x)
- **Update** the **posterior** distribution on f

}

# Gaussian Process Regression
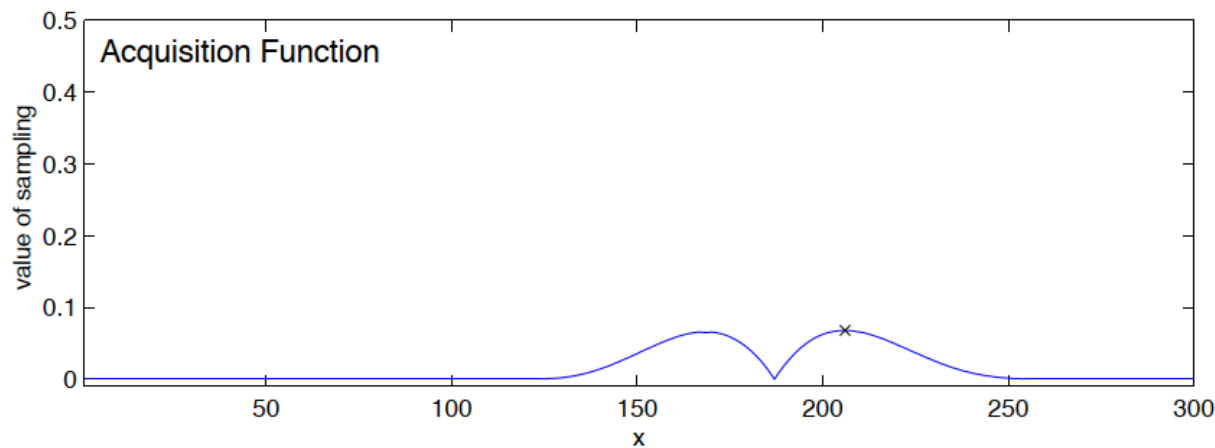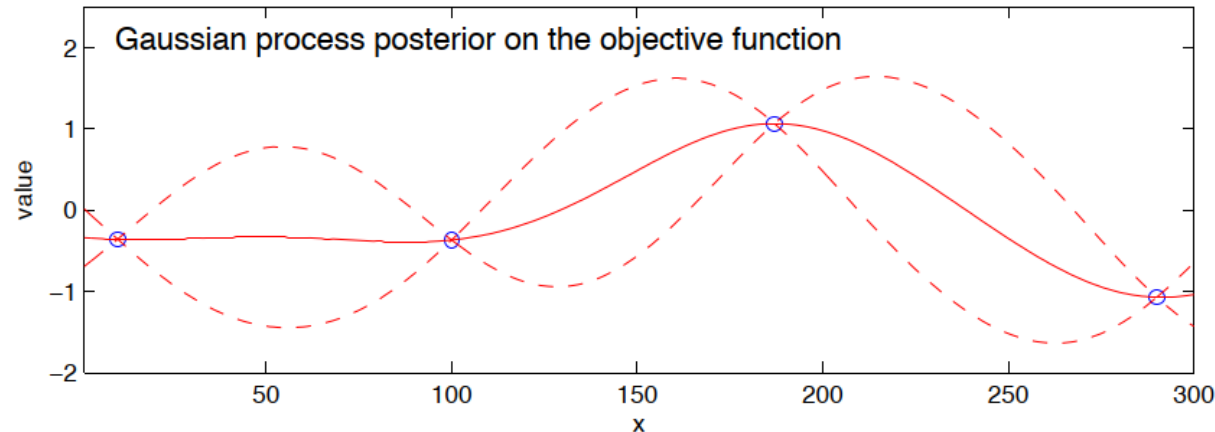
Bayesian optimization usually uses **Gaussian process regression (GPR)**:

- $f$ takes a vector $x$ as input

- We observe $f$ at some $x$ and want to predict its values at other points

- GPR is a Bayesian machine learning method for doing this

- $f$ is treated as random

- $f$'s probability distribution is called "prior distribution" before observing data and "posterior distribution" after observing data (data = $f$'s values at some points.



(a), prior   (b), posterior

Figure 1.1: Panel (a) shows four samples drawn from the prior distribution. Panel (b) shows the situation after two datapoints have been observed. The mean prediction is shown as the solid line and four samples from the posterior are shown as dashed lines. In both plots the shaded region denotes twice the standard deviation at each input value $x$.

[C. E. Rasmussen & C. K. I. Williams, Gaussian Processes for Machine Learning, the MIT Press, 2006]

# GPR model – Definition

Let $f: X \to \mathbb{R}$ be a real-valued function defined on a domain $X$,

A Gaussian process is a collection of random variables, any finite number of which have a joint Gaussian distribution.

A Gaussian process is completely specified by its mean function and a covariance function. We define mean function $\mu(x)$ and the covariance function $k(x, x')$ of a real process $f(x)$ as

$$\mu(x) = \mathbb{E}[f(x)]$$

$$k(x, x') = \mathbb{E}\big[\big(f(x) - \mu(x)\big)\big(f(x') - \mu(x')\big)\big],$$

and will write the Gaussian process as

$$f(x) \sim \mathcal{GP}\big(\mu(x), k(x, x')\big).$$

In our case, the random variables represent the value of the function f(x) at location x.

# GPR model – Definition

A Gaussian process is the generalization of multivariate distribution to infinite variables. It is probability distribution over functions.

# GPR model – Definition

A Gaussian process is a generalization of the Gaussian probability distribution. Whereas a probability distribution describes **random variables** which are scalars or vectors (for multivariate distributions), a **stochastic process** governs the properties of functions.

- Random variable Y → get an instance y.

Example:

random event $\omega \in \Omega$ (e.g., throw a dice) $\Longrightarrow$

- if dice $\leq 3$ , $y = 1$
- if $4 \geq$ dice $\leq 5$ , $y = 2$
- if dice $= 6$ , $y = 3$

- Random process Y(x) → get a function y(x).

Example:
a set of RVs indexed
by x random event $\omega \in \Omega$.

We can generate several repetitions.

[A. I. J. Forrester, A. Sóbester, and A. J. Keane, Engineering Design via Surrogate Modelling - A Practical Guide. John Wiley & Sons Ltd., 2008]

# How do we build a GPR model?

1. We start with a set of sample data, $X = \{x^{(1)}, \ldots, x^{(n)}\}^T$, with observed responses $y = \{y^{(1)}, \ldots, y^{(n)}\}^T$, and we want to find an expression for a predicted value at a new point $x$.

2. We are going to view our observed responses as if they are from a stochastic process (even though they may in fact be from a deterministic computer code). We denote this using the set of random vectors:

$$\mathbf{Y} = \begin{pmatrix} Y(\mathbf{x}^{(1)}) \\ \vdots \\ Y(\mathbf{x}^{(n)}) \end{pmatrix}.$$

3. It is a random field that has <span style="color:red">mean $\mathbf{1}\mu$</span>, and the random variables are correlated with each other using a basis/kernel function, like for example an *exponential* basis function:

$$\mathrm{cor}[Y(\mathbf{x}^{(i)}), Y(\mathbf{x}^{(l)})] = \exp\left(-\sum_{j=1}^{k} \theta_j \mid x_j^{(i)} - x_j^{(l)} \mid^{p_j}\right)$$

# How do we build a GPR model?

$$\mathrm{cor}[Y(\mathbf{x}^{(i)}), Y(\mathbf{x}^{(l)})] = \exp\left(-\sum_{j=1}^{k} \theta_j \mid x_j^{(i)} - x_j^{(l)} \mid^{p_j}\right)$$

From the correlation basis function, we can construct an $n \times n$ correlation matrix of all the observed data:

$$\boldsymbol{\Psi} = \begin{pmatrix} \mathrm{cor}[Y(\mathbf{x}^{(1)}), Y(\mathbf{x}^{(1)})] & \cdots & \mathrm{cor}[Y(\mathbf{x}^{(1)}), Y(\mathbf{x}^{(n)})] \\ \vdots & \ddots & \vdots \\ \mathrm{cor}[Y(\mathbf{x}^{(n)}), Y(\mathbf{x}^{(1)})] & \cdots & \mathrm{cor}[Y(\mathbf{x}^{(n)}), Y(\mathbf{x}^{(n)})] \end{pmatrix}$$

and a covariance matrix $\mathrm{Cov}(\mathbf{Y}, \mathbf{Y}) = \sigma^2 \boldsymbol{\Psi}$.

[Note that correlation is a normalized form of covariance, and some references use directly the covariance to derive the closed forms for mean and uncertainty of the model]

# Influence of the kernel hyperparameters

$$\text{cor}[Y(\mathbf{x}^{(i)}), Y(\mathbf{x}^{(l)})] = \exp\left(-\sum_{j=1}^{k} \theta_j \mid x_j^{(i)} - x_j^{(l)} \mid^{p_j}\right)$$

Correlations with varying $\boldsymbol{\theta}$ and $\mathbf{p}$:

- $\mathbf{p}$ changes the smoothness for the correlation function

- $\boldsymbol{\theta}$ is a parameter that affects how far a sample point's influence extends

# Kernel hyperparameters

$$\text{cor}[Y(\mathbf{x}^{(i)}), Y(\mathbf{x}^{(l)})] = \exp\left(-\sum_{j=1}^{k} \theta_j \mid x_j^{(i)} - x_j^{(l)} \mid^{p_j}\right)$$

We now know what our correlations mean, but

1) how do we estimate the values of $\boldsymbol{\theta}$ and $\mathbf{p}$ and

2) where does our observed data $\mathbf{y}$ come in?

**One** answer is to choose $\boldsymbol{\theta}$ and p to maximize the likelihood of $\mathbf{y}$:

$$L(\mathbf{Y}^{(1)}, \ldots, \mathbf{Y}^{(n)} \mid \mu, \sigma) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left[-\frac{\sum(\mathbf{Y}^{(i)} - \mu)^2}{2\sigma^2}\right]$$

Which can be expressed in terms of the sample data as

$$L = \frac{1}{(2\pi\sigma^2)^{n/2}|\mathbf{\Psi}|^{1/2}} \exp\left[-\frac{(\mathbf{y} - \mathbf{1}\mu)^{\mathrm{T}}\mathbf{\Psi}^{-1}(\mathbf{y} - \mathbf{1}\mu)}{2\sigma^2}\right]$$

# Kernel hyperparameters

$$L = \frac{1}{(2\pi\sigma^2)^{n/2}|\mathbf{\Psi}|^{1/2}} \exp\left[-\frac{(\mathbf{y}-\mathbf{1}\mu)^{\mathrm{T}}\mathbf{\Psi}^{-1}(\mathbf{y}-\mathbf{1}\mu)}{2\sigma^2}\right]$$

We need to maximize the likelihood, so we can consider its natural logarithm:

$$\ln(L) = -\frac{n}{2}\ln(2\pi) - \frac{n}{2}\ln(\sigma^2) - \frac{1}{2}\ln|\mathbf{\Psi}| - \frac{(\mathbf{y}-\mathbf{1}\mu)^{\mathrm{T}}\mathbf{\Psi}^{-1}(\mathbf{y}-\mathbf{1}\mu)}{2\sigma^2}$$

By taking the derivatives of this equation and setting them to zero, we obtain the maximum likelihood estimates (MLEs) for $\mu$ and $\sigma^2$:

$$\widehat{\mu} = \frac{\mathbf{1}^{\mathrm{T}}\mathbf{\Psi}^{-1}\mathbf{y}}{\mathbf{1}^{\mathrm{T}}\mathbf{\Psi}^{-1}\mathbf{1}},$$

$$\widehat{\sigma}^2 = \frac{(\mathbf{y}-\mathbf{1}\mu)^{\mathrm{T}}\mathbf{\Psi}^{-1}(\mathbf{y}-\mathbf{1}\mu)}{n}.$$

Substitute back, remove constant terms, and optimize w.r.t. $\theta$ and p!

# GPR Prediction

To predict values at unobserved locations, we need to maximize the likelihood of the sample data **AND** the prediction.

To do this, we augment the observed data with the new prediction ($\tilde{\boldsymbol{y}} = \{\boldsymbol{y}^T, \hat{y}\}^T$) and define a vector of correlations between the observed data and the new prediction:

$$\boldsymbol{\psi} = \begin{pmatrix} \text{cor}[Y(\mathbf{x}^{(1)}), Y(\mathbf{x})] \\ \vdots \\ \text{cor}[Y(\mathbf{x}^{(n)}), Y(\mathbf{x})] \end{pmatrix} = \begin{pmatrix} \psi^{(1)} \\ \vdots \\ \psi^{(n)} \end{pmatrix}$$

Now, we construct an augmented correlation matrix:

$$\tilde{\boldsymbol{\Psi}} = \begin{pmatrix} \boldsymbol{\Psi} & \boldsymbol{\psi} \\ \boldsymbol{\psi}^T & 1 \end{pmatrix}$$

# GPR Prediction

The ln-likelihood of the augmented data is

$$\ln(L) = -\frac{n}{2}\ln(2\pi) - \frac{n}{2}\ln(\widehat{\sigma}^2) - \frac{1}{2}\ln|\widetilde{\boldsymbol{\Psi}}| - \frac{(\widetilde{\mathbf{y}} - \mathbf{1}\widehat{\mu})^{\mathrm{T}}\widetilde{\boldsymbol{\Psi}}^{-1}(\widetilde{\mathbf{y}} - \mathbf{1}\widehat{\mu})}{2\widehat{\sigma}^2}$$

Here, only the last term depends on $\hat{y}$, so we consider only that term and optimize (substituting $\tilde{y}$ and $\widetilde{\Psi}$)

$$\ln(L) \approx \frac{-\begin{pmatrix}\mathbf{y} - \mathbf{1}\widehat{\mu} \\ \widehat{y} - \widehat{\mu}\end{pmatrix}^{\mathrm{T}}\begin{pmatrix}\boldsymbol{\Psi} & \boldsymbol{\psi} \\ \boldsymbol{\psi}^{\mathrm{T}} & 1\end{pmatrix}^{-1}\begin{pmatrix}\mathbf{y} - \mathbf{1}\widehat{\mu} \\ \widehat{y} - \widehat{\mu}\end{pmatrix}}{2\widehat{\sigma}^2}$$

# GPR Prediction

With some matrix calculations (inversions and multiplications) we can obtain:

$$\ln(L) \approx \left( \frac{-1}{2\widehat{\sigma}^2 (1 - \boldsymbol{\psi}^{\mathrm{T}} \boldsymbol{\Psi}^{-1} \boldsymbol{\psi})} \right) (\widehat{y} - \widehat{\mu})^2 + \left( \frac{\boldsymbol{\psi}^{\mathrm{T}} \boldsymbol{\Psi}^{-1} (\mathbf{y} - \mathbf{1}\widehat{\mu})}{\widehat{\sigma}^2 (1 - \boldsymbol{\psi}^{\mathrm{T}} \boldsymbol{\Psi}^{-1} \boldsymbol{\psi})} \right) (\widehat{y} - \widehat{\mu})$$

Which we differentiate with respect to $\hat{y}$ and set to zero:

$$\left( \frac{-1}{\widehat{\sigma}^2 (1 - \boldsymbol{\psi}^{\mathrm{T}} \boldsymbol{\Psi}^{-1} \boldsymbol{\psi})} \right) (\widehat{y} - \widehat{\mu}) + \left( \frac{\boldsymbol{\psi}^{\mathrm{T}} \boldsymbol{\Psi}^{-1} (\mathbf{y} - \mathbf{1}\widehat{\mu})}{\widehat{\sigma}^2 (1 - \boldsymbol{\psi}^{\mathrm{T}} \boldsymbol{\Psi}^{-1} \boldsymbol{\psi})} \right) = 0$$

Thus, our MLE for $\hat{y}$ is:
$$\boxed{\widehat{y}(\mathbf{x}) = \widehat{\mu} + \boldsymbol{\psi}^{\mathrm{T}} \boldsymbol{\Psi}^{-1} (\mathbf{y} - \mathbf{1}\widehat{\mu})}$$

**POSTERIOR MEAN**

**Note that:** the model is constructed in such a way that the prediction goes through all the data points (it interpolates the data):

If we make a prediction at a sampled point $x^{(i)}$, $\psi$ is the $i^{\text{th}}$ column of $\Psi$, meaning that $\psi\Psi^{-1}$ it the $i^{\text{th}}$ unit vector. Hence, $\hat{y}(x) = \hat{\mu} + y^{(i)} - \hat{\mu} = y^{(i)}$, which is exactly the observed objective function value at $x^{(i)}$.

# GPR Prediction

The GPR model also permits the calculation of an estimated error in the model at the queried location $x$.

The committed error in the prediction is denoted as $\hat{s}^2$ and, after some other calculations, it is derived to be defined as:
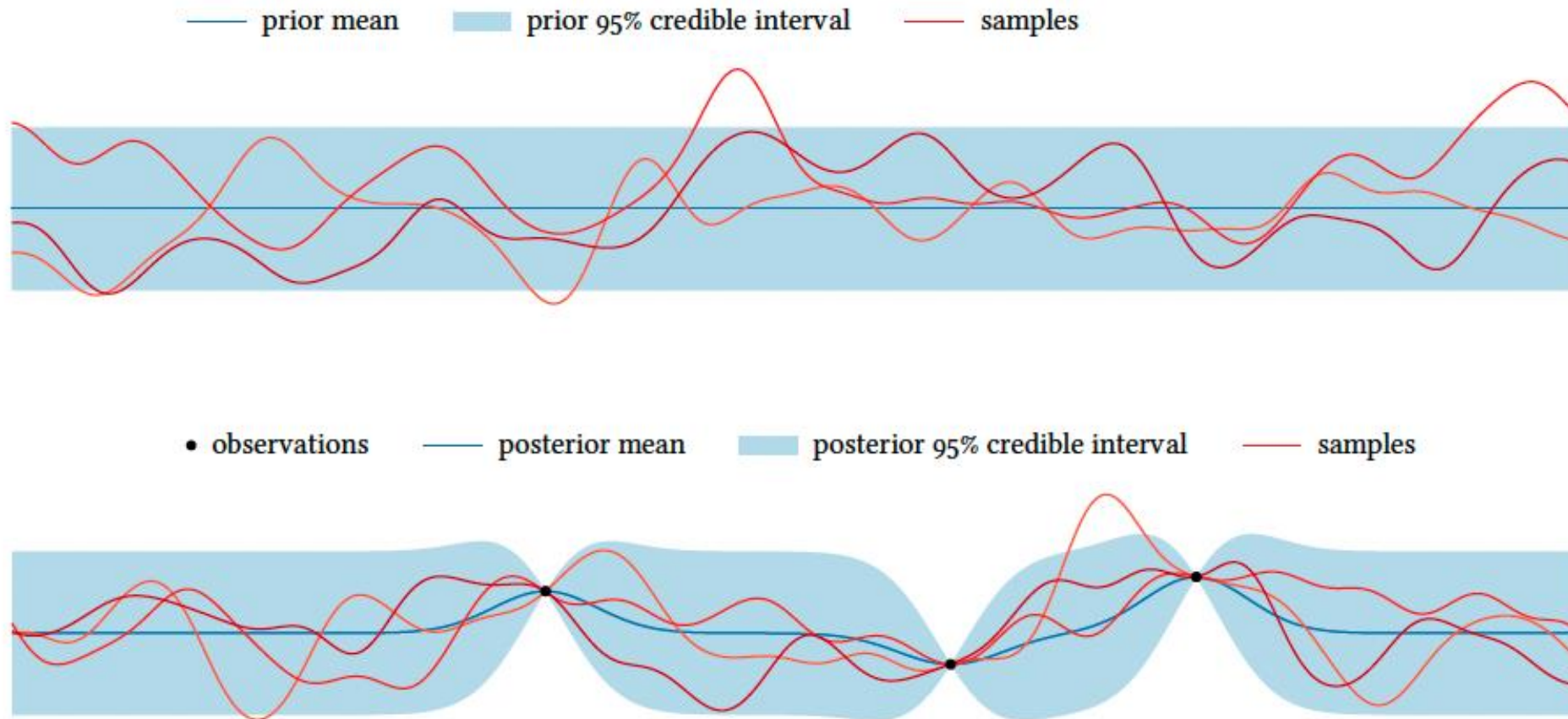
$$\widehat{s}^2(\mathbf{x}) = \sigma^2 \left[ 1 - \boldsymbol{\psi}^{\mathrm{T}} \boldsymbol{\Psi}^{-1} \boldsymbol{\psi} + \frac{1 - \mathbf{1}^{\mathrm{T}} \boldsymbol{\Psi}^{-1} \boldsymbol{\psi}}{\mathbf{1}^{\mathrm{T}} \boldsymbol{\Psi}^{-1} \mathbf{1}} \right]$$

**POSTERIOR VARIANCE**

Based on the posterior mean and posterior variance, we have our GPR which is constructed allover the design domain and can start looking for new points to query (infill points) in an iterative fashion.

[Details can be found in: Sacks, J., Welch, W. J., Mitchell, T. J. and Wynn, H. (1989) Design and analysis of computer experiments. Statistical Science, 4(4), 409−423 ]

# GPR Prediction



[R. Garnett, Bayesian Optimization Book. Accessed: Feb. 08, 2023. [Online]. Available: https://bayesoptbook.com//]

# GPR Summary

- We look for an unknown response value $\widehat{y}$ at the new location $\widehat{\boldsymbol{x}}$
- For this, we assume that known responses and the new response have a joint Gaussian distribution with specified covariance structure

$$\begin{bmatrix} \boldsymbol{y} \\ \widehat{y} \end{bmatrix} \sim \mathcal{N}\left( \begin{bmatrix} \boldsymbol{\mu}(\boldsymbol{x}) \\ \mu(\widehat{\boldsymbol{x}}) \end{bmatrix}, \begin{bmatrix} \mathbf{\Sigma_0} & \boldsymbol{\gamma}(\widehat{\boldsymbol{x}}) \\ \boldsymbol{\gamma}(\widehat{\boldsymbol{x}})^T & \Sigma_0(\widehat{\boldsymbol{x}}, \widehat{\boldsymbol{x}}) \end{bmatrix} \right)$$

- Conditional distribution of the unknown output

$$\widehat{y}|\boldsymbol{y} \sim \mathcal{N}\left( \mu(\widehat{x}) + \boldsymbol{\gamma}(\widehat{\boldsymbol{x}})^T\mathbf{\Sigma_0^{-1}}(\boldsymbol{y} - \mathbb{I}_n\boldsymbol{\mu}(\boldsymbol{x})), \Sigma_0(\widehat{\boldsymbol{x}}, \widehat{\boldsymbol{x}}) - \boldsymbol{\gamma}(\widehat{\boldsymbol{x}})^T\mathbf{\Sigma_0^{-1}}\boldsymbol{\gamma}(\widehat{\boldsymbol{x}}) \right)$$

$$\widehat{\boldsymbol{y}} = \mu(\widehat{x}) + \boldsymbol{\gamma}(\widehat{\boldsymbol{x}})^T\mathbf{\Sigma_0^{-1}}(\boldsymbol{y} - \mathbb{I}_n\boldsymbol{\mu}(\boldsymbol{x}))$$

Posterior mean

$$C(\boldsymbol{x}, \boldsymbol{x}') = \Sigma_0(\boldsymbol{x}, \boldsymbol{x}') - \boldsymbol{\gamma}(\boldsymbol{x})^T\mathbf{\Sigma_0^{-1}}\boldsymbol{\gamma}(\boldsymbol{x}')$$

Posterior covariance

# Alternative Kernels

Other common choices as kernels for GPs are:

### Radial basis function kernel

$$k_{RBF}(x, x') = \exp\left(-\frac{(x - x')^2}{2\ell^2}\right)$$

### Matérn kernel

$$k_M(x, x') = \frac{2^{1-\nu}}{\Gamma(\nu)}\left(\frac{\sqrt{2\nu}|x - x'|}{\ell}\right)^{\nu} K_\nu\left(\frac{\sqrt{2\nu}|x - x'|}{\ell}\right)$$

Where

- $l$ is a length-scale parameter

- $\nu$ is a positive parameter controlling smoothness. Three of the most popular choices are $\nu \in \{\frac{1}{2}, \frac{3}{2}, \frac{5}{2}\}$

$$k_M^{1/2}(x, x') = \exp\left(-\frac{|x - x'|}{\ell}\right)$$

$$k_M^{3/2}(x, x') = \left(1 + \frac{\sqrt{3}|x - x'|}{\ell}\right)\exp\left(-\frac{\sqrt{3}|x - x'|}{\ell}\right)$$

$$k_M^{5/2}(x, x') = \left(1 + \frac{\sqrt{5}|x - x'|}{\ell} + \frac{5(x - x')^2}{3\ell^2}\right)\exp\left(-\frac{\sqrt{5}|x - x'|}{\ell}\right)$$

- $K_\nu$ is a modified Bessel function

[Technical blog post by Andy Jones: https://andrewcharlesjones.github.io/journal/matern-kernels.html]

# Alternative Kernels

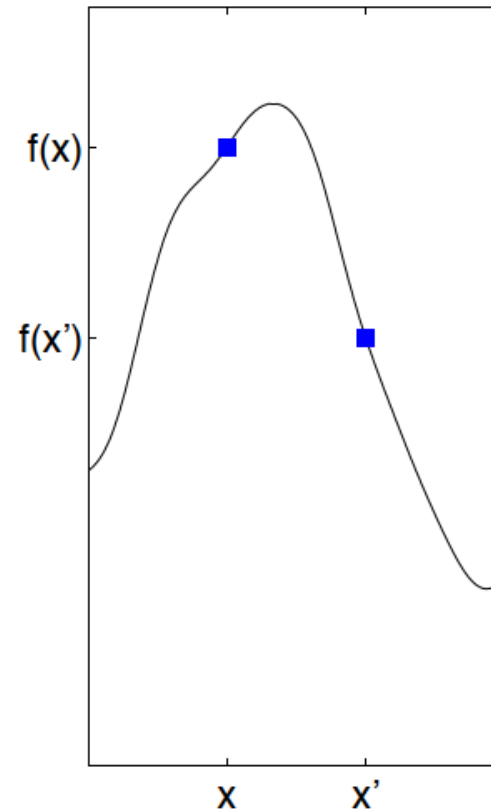In case of a discontinuous function, Matérn karnel is able to better capture the sharp jumps than RBF

(the plot below is for $\nu=3/2$)



$$k_M^{1/2}(x,x') = \exp\left(-\frac{|x-x'|}{\ell}\right)$$

$$k_M^{3/2}(x,x') = \left(1+\frac{\sqrt{3}|x-x'|}{\ell}\right)\exp\left(-\frac{\sqrt{3}|x-x'|}{\ell}\right)$$

$$k_M^{5/2}(x,x') = \left(1+\frac{\sqrt{5}|x-x'|}{\ell}+\frac{5(x-x')^2}{3\ell^2}\right)\exp\left(-\frac{\sqrt{5}|x-x'|}{\ell}\right)$$

[Technical blog post by Andy Jones: https://andrewcharlesjones.github.io/journal/matern-kernels.html]

# "Extra" visual example



- Fix 2 points x and x'

- Consider the vector [f(x),f(x')]

- We will understand how observing f(x') leads to a posterior distribution on f(x)

# "Extra" visual example



x    x'

f(x')

f(x)

# "Extra" visual example

# "Extra" visual example

# "Extra" visual example

# "Extra" visual example



- This procedure can give us the posterior on f(x) for any x

- This posterior is normal. The blue dot is its mean. Error bars give the mean +/- 2 standard deviations.

- Repeating this procedure for every point on the x-axis gives the red solid line (posterior mean) and red dashed lines (posterior credible region)

# Standard logical flow of BO

**Target**: Find $x^* = \underset{x \in X}{\arg\min} f(x)$

Assume a Bayesian **prior** on f (usually a Gaussian process prior)

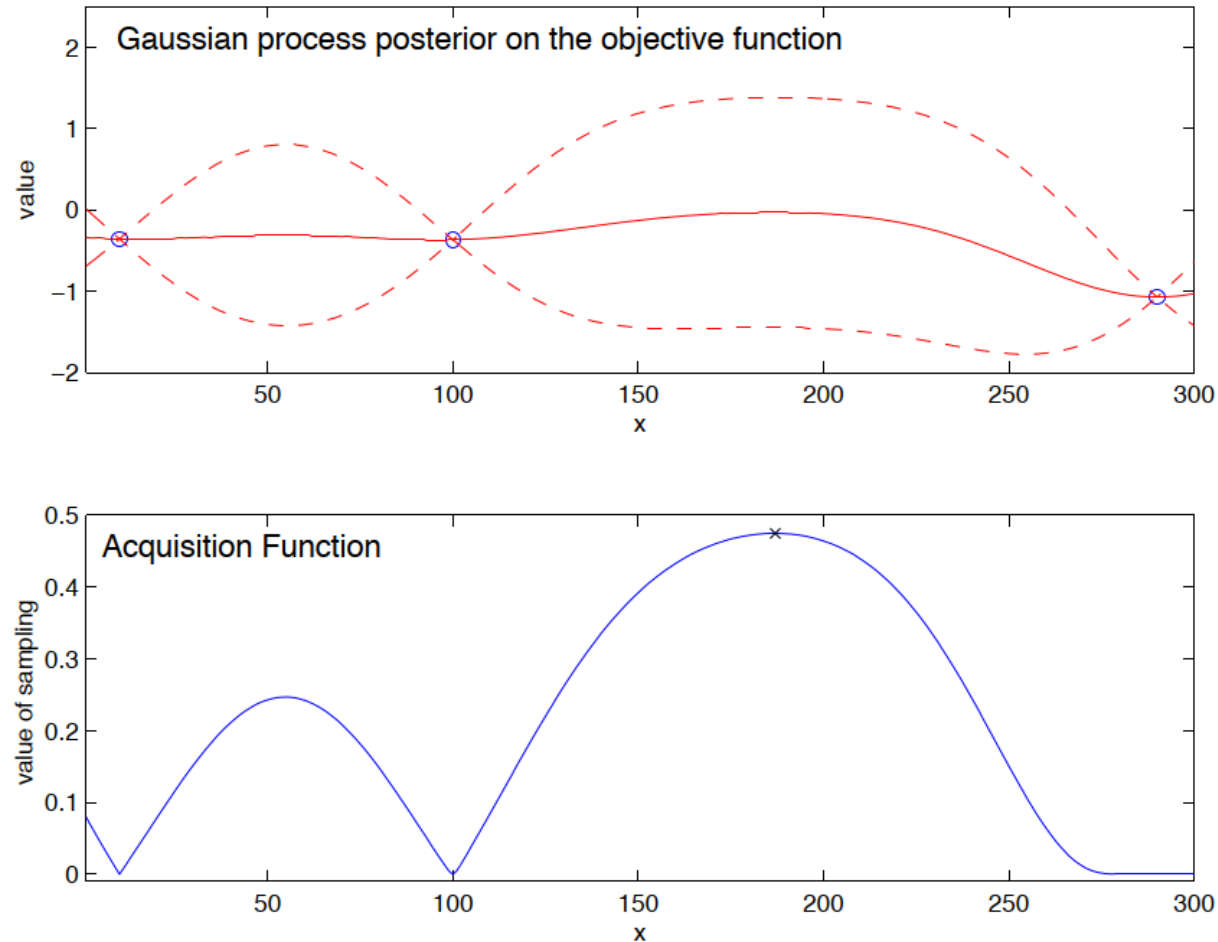Build the **posterior** distribution

while (budget is not exhausted) {

- Find x that maximizes **acquisition**(x,**posterior**)
- **Sample** x & **observe** f(x)
- **Update** the **posterior** distribution on f
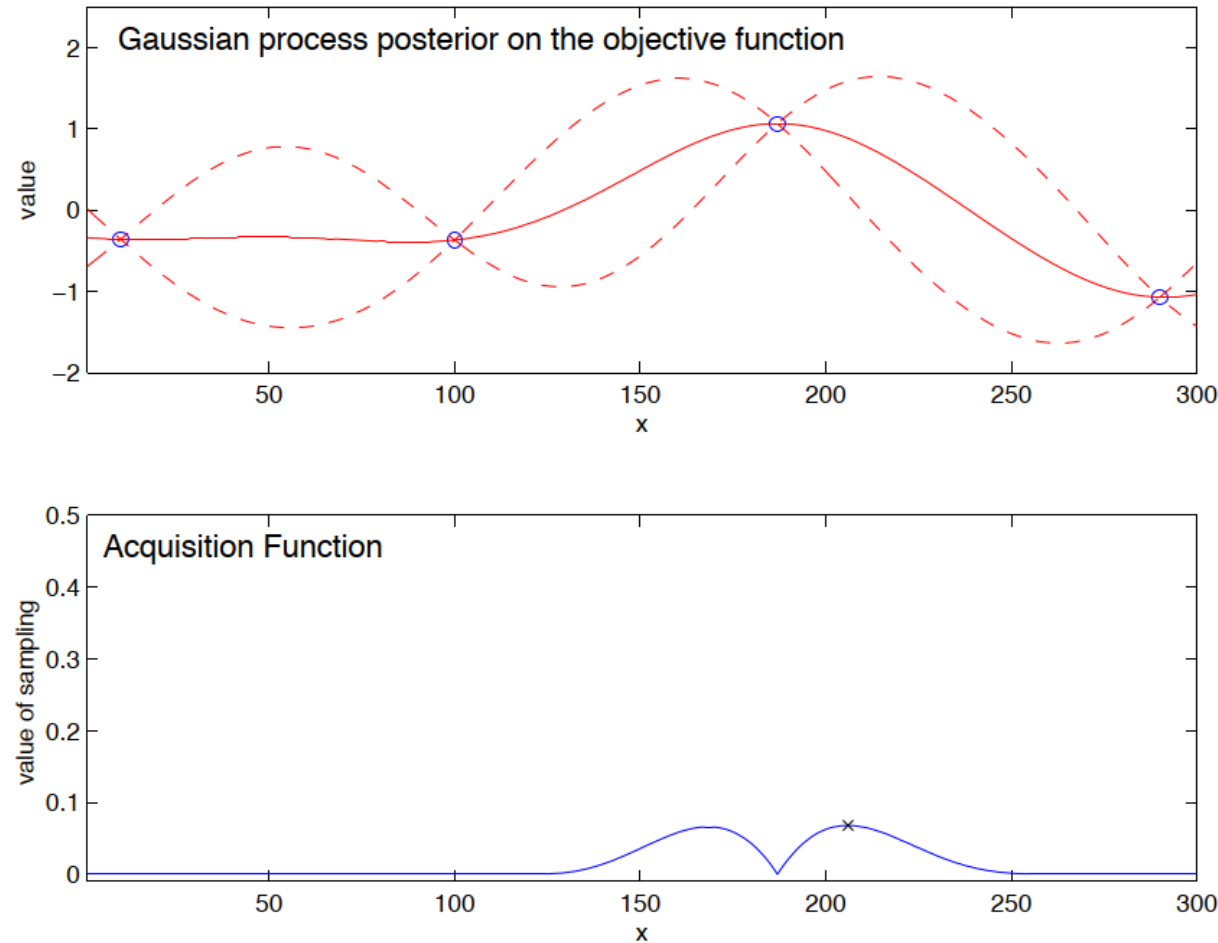
}

# Optimizing a 1-dimensional problem

Let's remember how the iterative and update procedure looks like:

# Optimizing a 1-dimensional problem

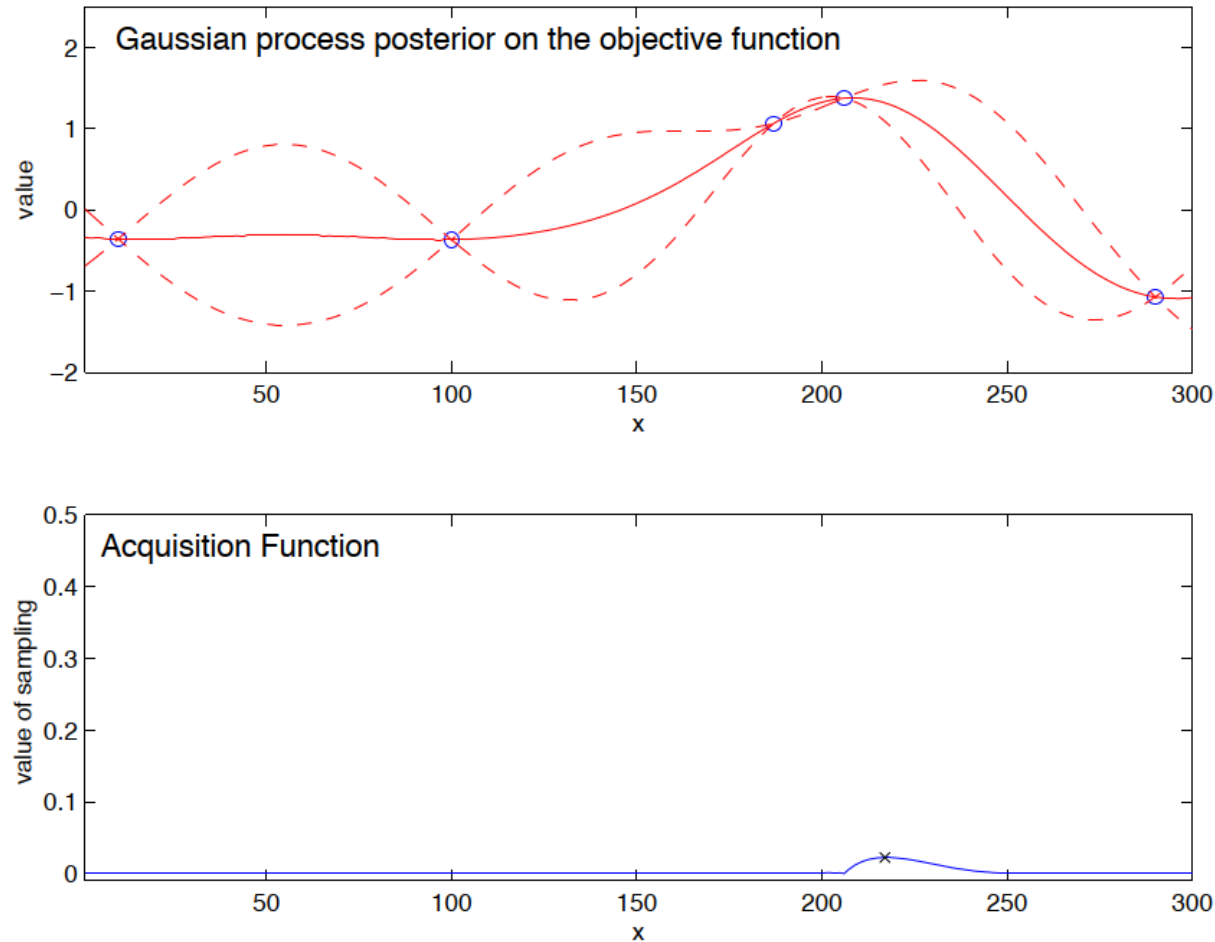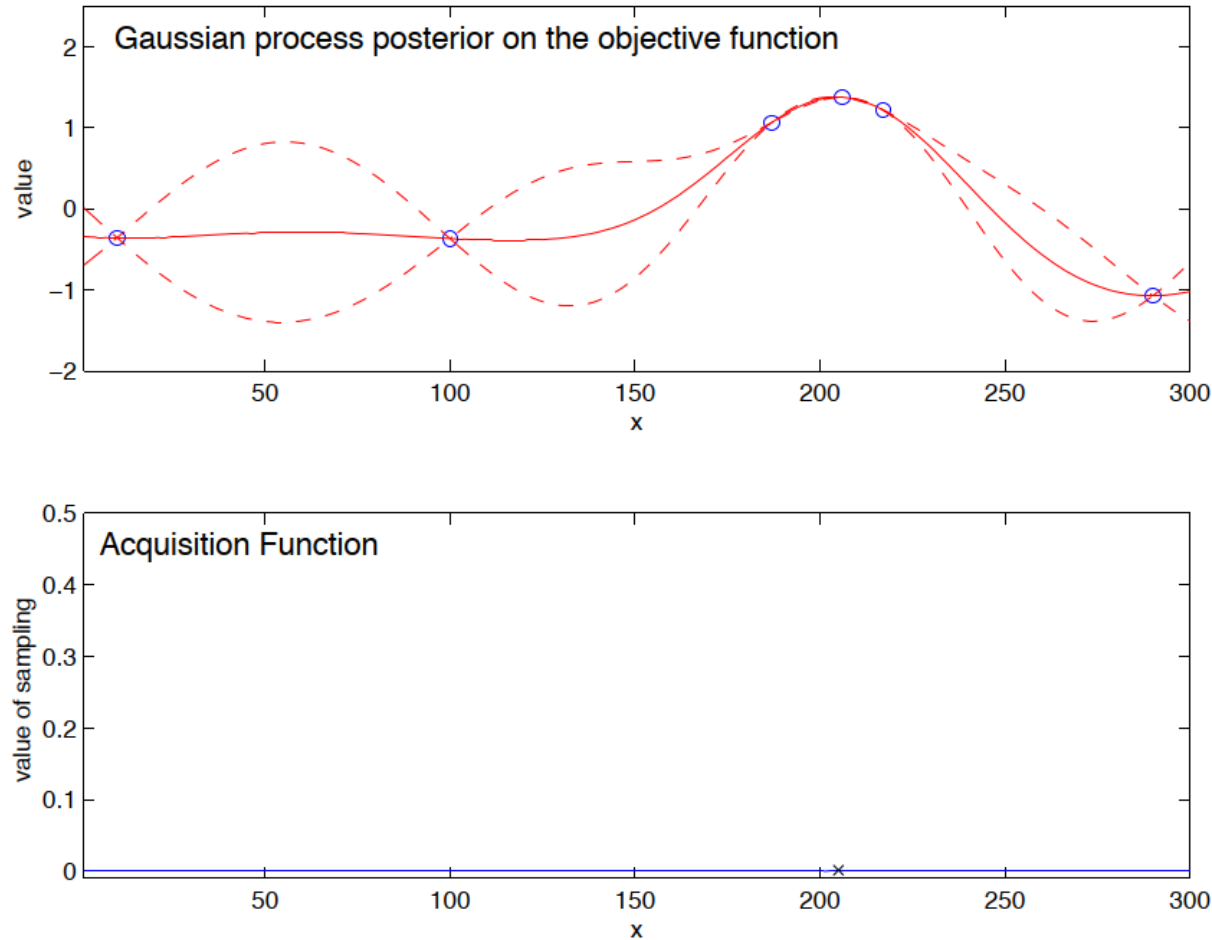Let's remember how the iterative and update procedure looks like:

# Optimizing a 1-dimensional problem

Let's remember how the iterative and update procedure looks like:
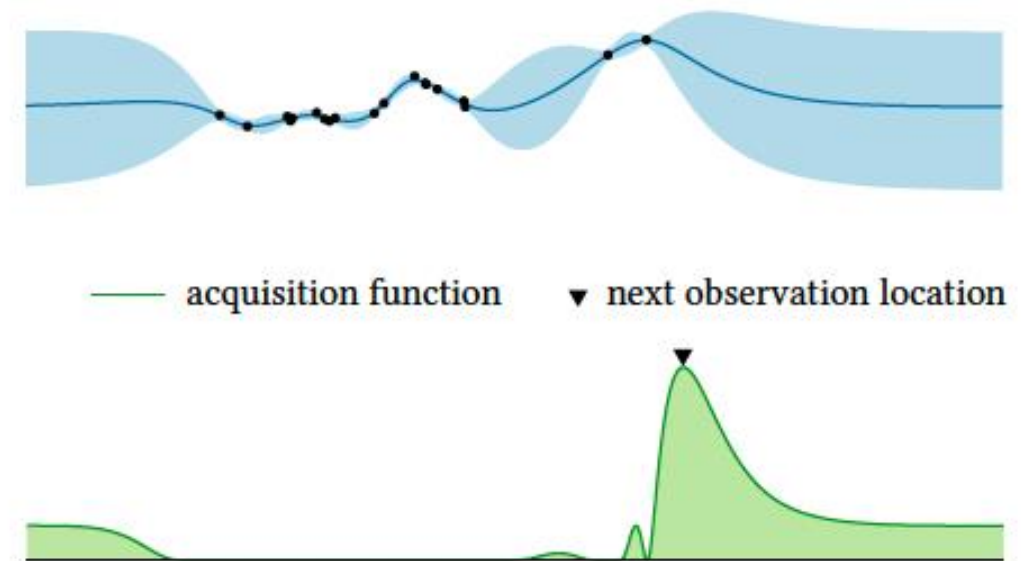
# Optimizing a 1-dimensional problem

Let's remember how the iterative and update procedure looks like:

# Acquisition Function

- The iterative procedure is based on the maximization of an acquisition function (blue curve in previous slides)

- The acquisition function measures the value that would be generated by evaluation of the objective function at a new point x, based on the current posterior distribution over f

- Each time we observe f at a new point, this GPR posterior distribution is updated

- Well-known acquisition functions are:

  1. Lower/Upper confidence bound

  2. Probability of improvement

  3. Expected improvement
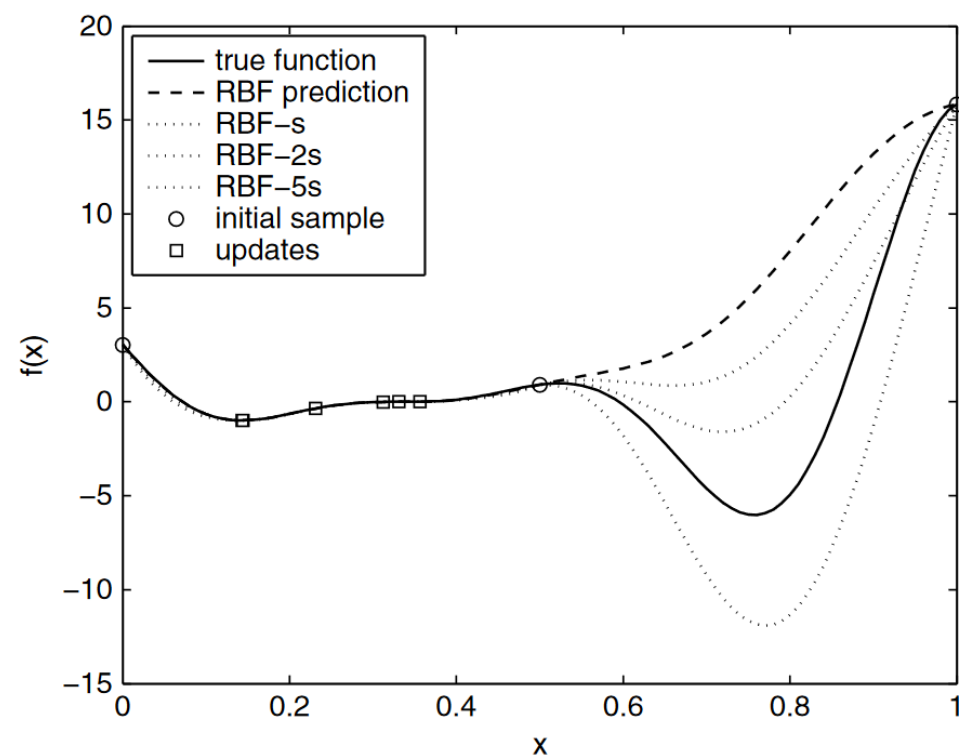
  4. Thompson Sampling

  5. Etc.



—— acquisition function ▼ next observation location

# AF: Lower Confidence Bound (LCB)

The simplest way of balancing exploitation of the prediction $\hat{y}(\boldsymbol{x})$ and exploration using $\hat{s}^2(\boldsymbol{x})$ is to minimize a lower confidence bound:

$$\text{LB}(\mathbf{x}) = \widehat{y}(\mathbf{x}) - A\widehat{s}(\mathbf{x}),$$

where A is a constant that controls the exploitation/exploration balance.

- As $A \to 0,\ LB(\boldsymbol{x}) \to \hat{y}(\boldsymbol{x})$ (pure exploitation)

- As $A \to \infty$, the effect of $\hat{y}(\boldsymbol{x})$ becomes negligible (pure exploration)



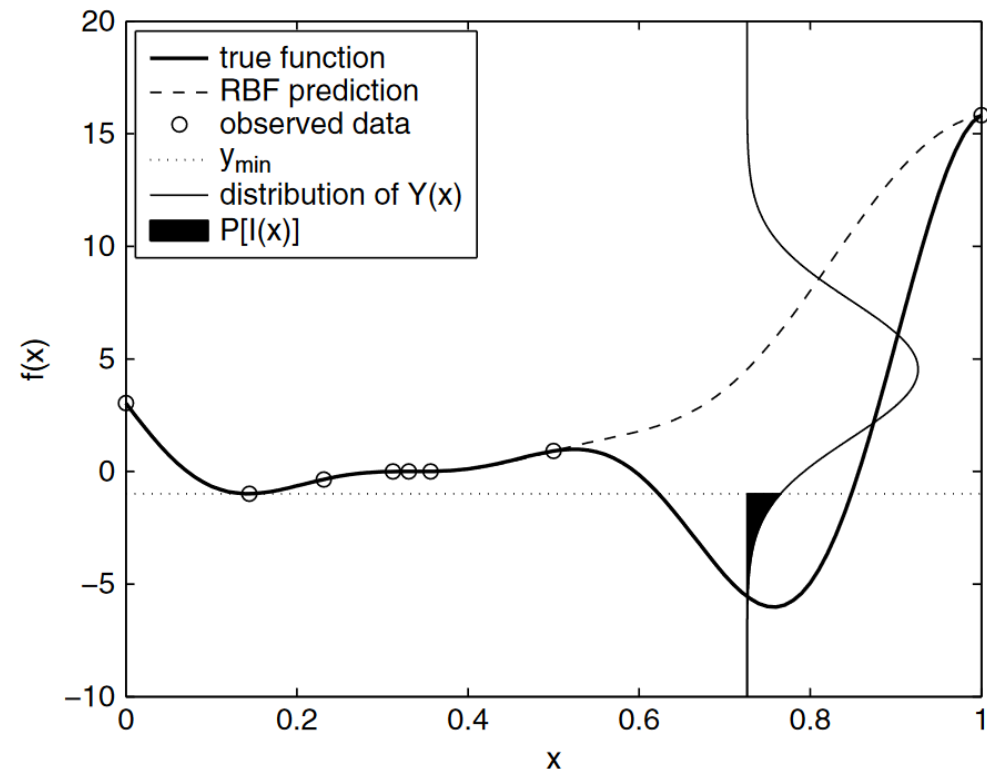The lower confidence bound for varying A.

# AF: Probability of Improvement (PI)

The aim is to position the infill point at a value of **x** that will lead to an improvement over the best observed so far, $y_{\min}$.

We considered the approximated value $\hat{y}(\boldsymbol{x})$ as the realization of a random variable. Hence, we can compute the probability of improvement upon $y_{\min}$,
$I = y_{\min} - Y(\boldsymbol{x})$:

$$P[I(\mathbf{x})] = \frac{1}{\widehat{s}\sqrt{2\pi}} \int_{-\infty}^{0} e^{-[I-\widehat{y}(\mathbf{x})]^2/(2s^2)} dI$$

In the figure, we see a Gaussian distribution with variance $s^2(\boldsymbol{x})$ centered around the predicted mean value at location $\boldsymbol{x}$. The distribution represents the uncertainty in the prediction and the shaded area is the probability of improving with respect to $y_{\min}$ (dotted line).



A graphical interpretation of the probability of improvement.

# AF: Expected Improvement (EI)

Instead of simply finding the probability that there will be some improvement, we can calculate the amount of improvement we expect, given the mean $\hat{\mu}(x)$ and the variance $\hat{s}^2(x)$.

$$I(\mathbf{x}) = \max(y_{min} - Y(\mathbf{x}), 0)$$

$$P[I(\mathbf{x})] = P[Y(\mathbf{x}) < y_{min}]$$

$$P[I(\mathbf{x})] = \Phi\left(\frac{y_{min} - \hat{y}(\mathbf{x})}{\hat{s}(\mathbf{x})}\right) = \frac{1}{\hat{s}(\mathbf{x})\sqrt{2\pi}}\int_{-\infty}^{0} e^{-\frac{[I(\mathbf{x}) - \hat{y}(\mathbf{x})]^2}{2\hat{s}^2(\mathbf{x})}} dI \qquad \longrightarrow \qquad \mathbf{x}_{infill} = \operatorname*{argmax}_{\mathbf{x}} E[I(\mathbf{x})].$$

$$E[I(\mathbf{x})] = \begin{cases} (y_{min} - \hat{y}(\mathbf{x}))\Phi\left(\frac{y_{min} - \hat{y}(\mathbf{x})}{\hat{s}(\mathbf{x})}\right) + \hat{s}(\mathbf{x})\phi\left(\frac{y_{min} - \hat{y}(\mathbf{x})}{\hat{s}(\mathbf{x})}\right) & \text{if } \hat{s}(\mathbf{x}) > 0 \\ 0 & \text{if } \hat{s}(\mathbf{x}) = 0 \end{cases}$$
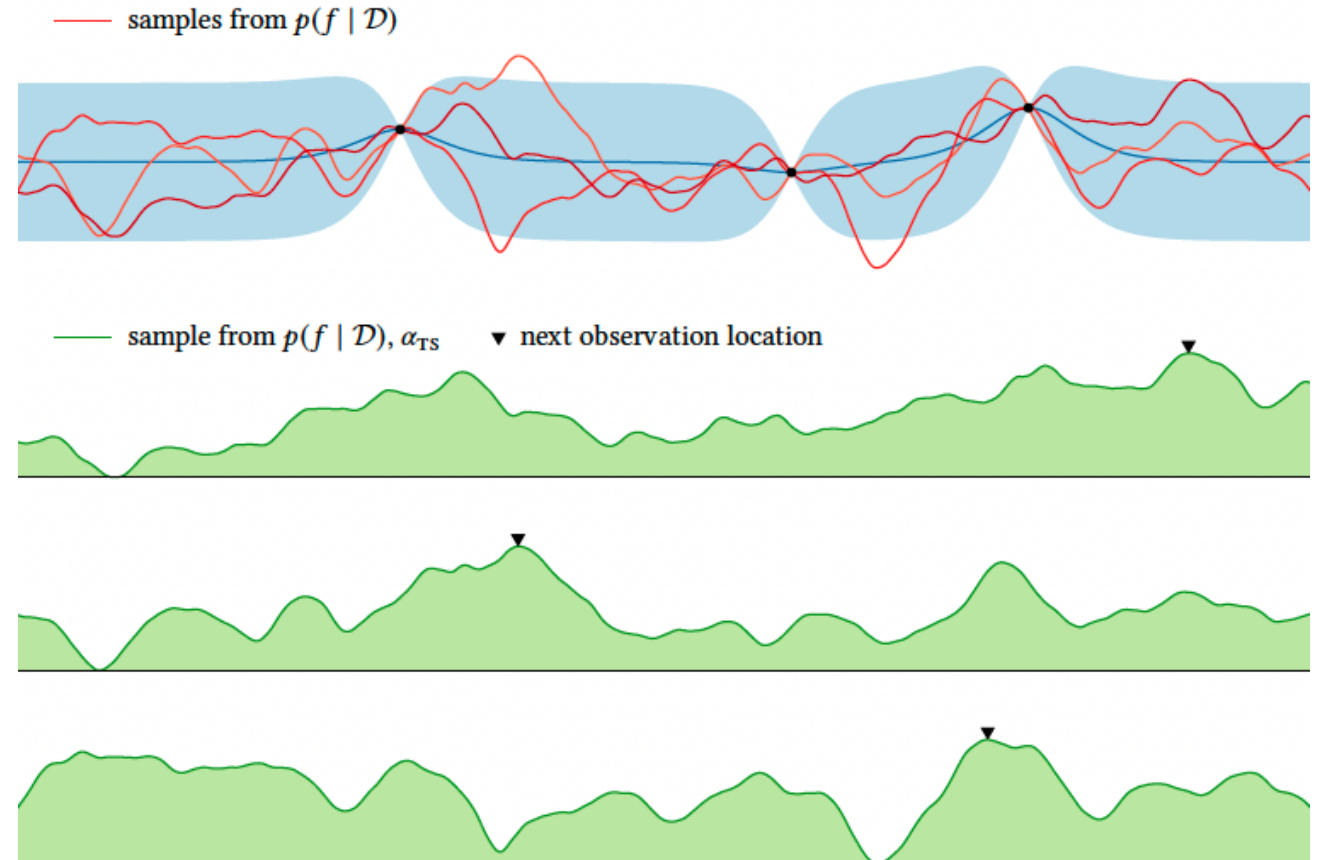
Mean-based term promoting exploitation of low function value regions

Variance-based term to explore sparse regions

where $\Phi(.)$ and $\phi(.)$ are the cumulative distribution function and probability density function, respectively.

# AF: Thompson Sampling (TS)

1. Thompson Sampling selects the next evaluation point by **sampling from the posterior** of the surrogate model.

2. It uses the **sampled function** to identify the point with the highest expected value and chooses it for the next evaluation.

3. Draws different samples of the GPR model and optimizes over each of them.

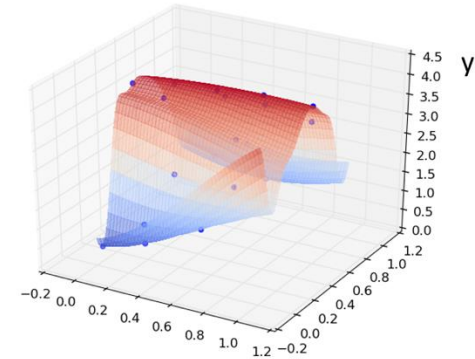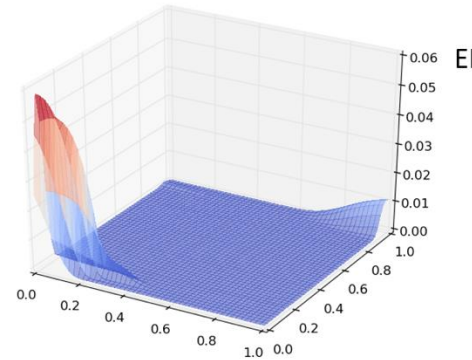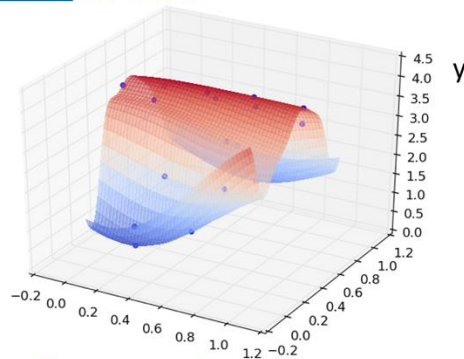4. Variance in the suggestion is given by the repetitive nature of the procedure.



[R. Garnett, Bayesian Optimization Book. Accessed: Feb. 08, 2023. [Online]. Available: https://bayesoptbook.com//]
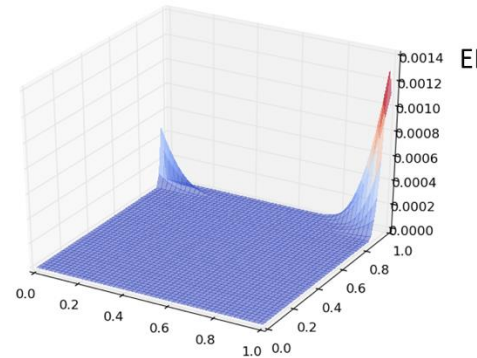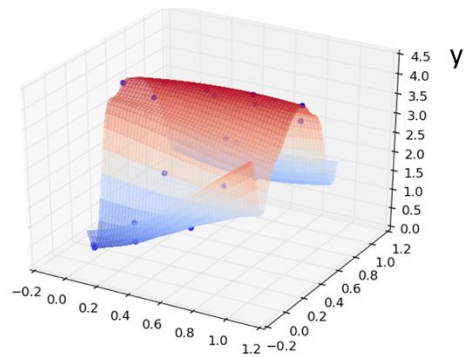
# Seeked behavior in AFs...

**...Good trade-off between exploitation and exploration!**



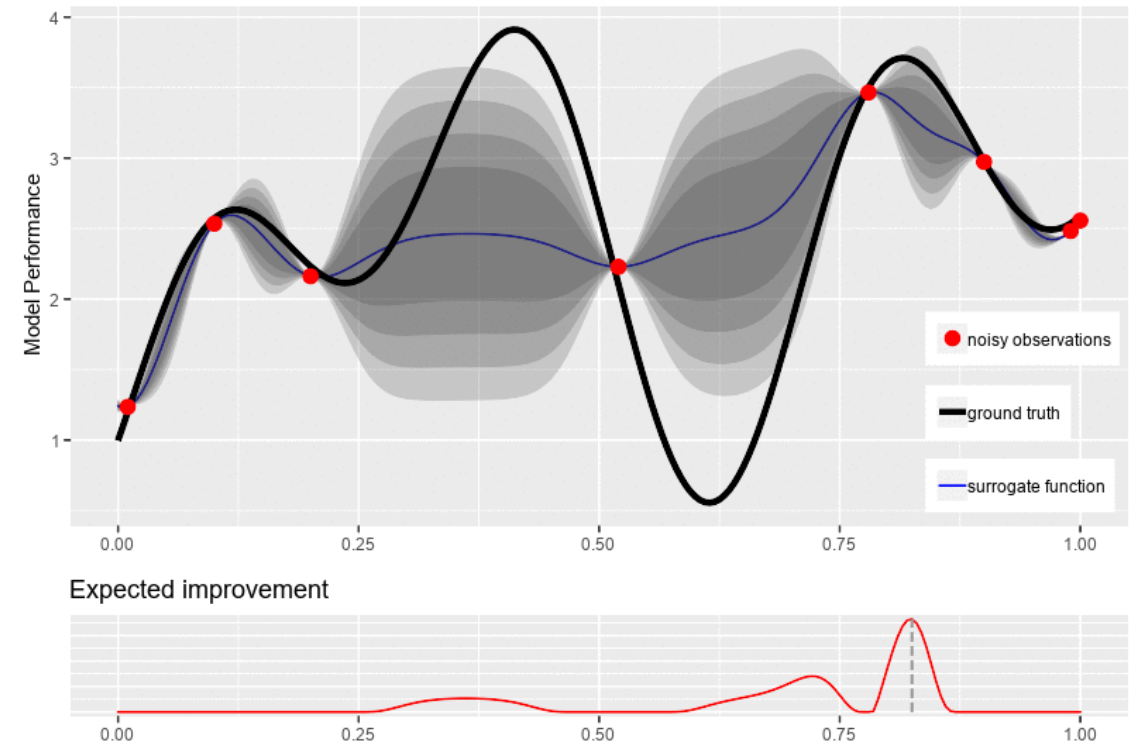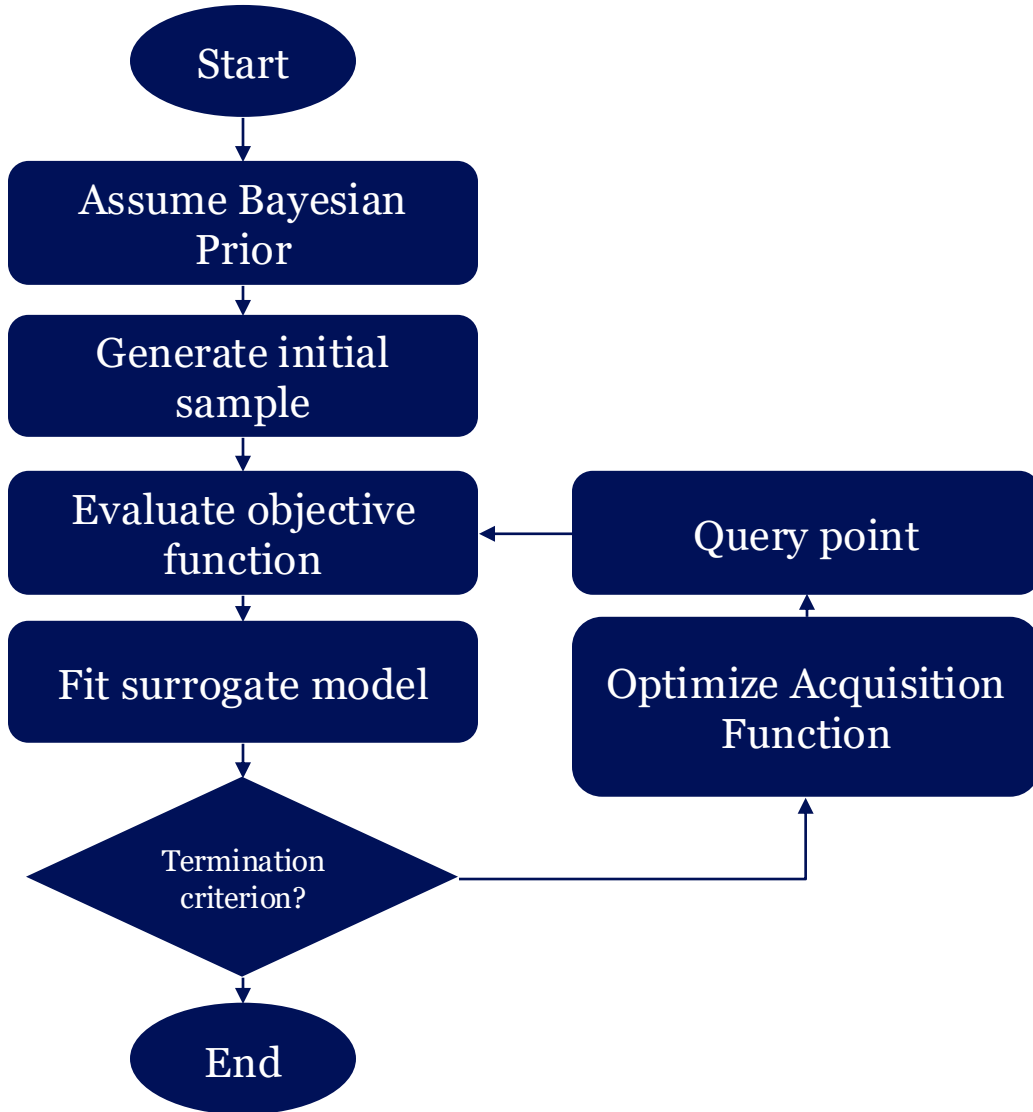Example of Kriging-surrogate and Expected Improvement functions update
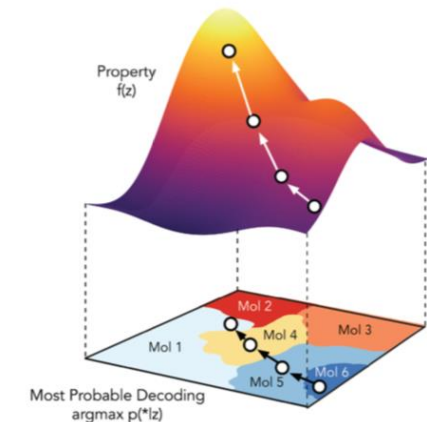2-variables problem

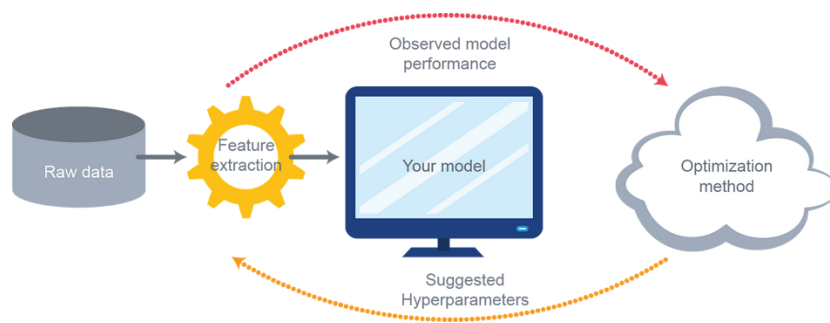Exploitive search

Explorative search

# Altogether – Bayesian Optimization

# Bayesian Optimization is a big field of research!

Today, we haven't touched many subtopics:

- BO for mixed-integer/discrete spaces

- BO for noisy problems

- High-dimensional BO

- Constrained BO

- Applications of BO to real-world problems (HPO, Engineering, Chemistry...)

# Thank you!