

AI and the Uncertainty Challenge in Fundamental Physics
CNRS AISSAI and CNRS IN2P3

An introduction to Bayesian optimization

Emmanuel Vazquez
Laboratoire des Signaux et Systèmes – Paris-Saclay

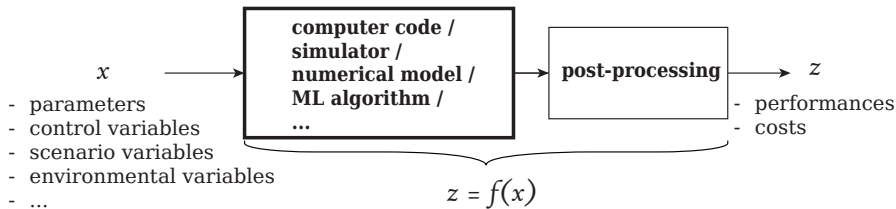
Nov. 27, 2023

Contents

- 1 Computer experiments and optimization
- 2 Optimization of expensive-to-evaluate functions
- 3 Bayesian black-box modeling
- 4 Bayesian optimization
- 5 Concluding remarks

1 Computer experiments and optimization

- Consider a **computer procedure** with inputs and outputs

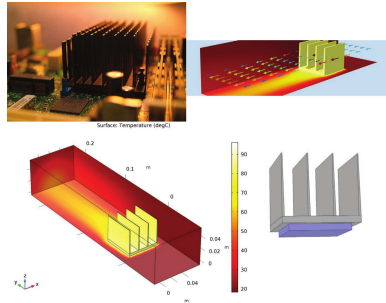


- Selecting a value for $x \in \mathbb{X}$ and observing the resulting output $z \in \mathbb{Z}$ is a **computer experiment**.
- In many cases, we aim to **minimize** or **maximize** the value of z .

Example (1/2): optimization of a system

Find the best values for the design parameters of a system

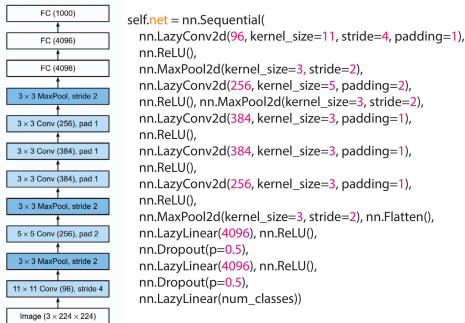
- **Heat sink** → shape controls airflow characteristics and radiation, which have direct impact on **cooling performances**
- $f : \mathbb{X} \subset \mathbb{R}^d \rightarrow \mathbb{R}$ performance as a function of design parameters.
- Computing $f(x)$: **time-consuming!** (solve PDEs using a finite element method)
- Optimization must be done within a **limited budget of simulations**



Simulation of a heat sink in COMSOL Multiphysics

Example (2/2): optimization of a DNN

- $f : \mathbb{X} \rightarrow \mathbb{R} \rightsquigarrow$ **validation loss** of a DNN as a function of its parameters (size of layers, dropout. . .).
- Computing $f(x)$ on large datasets is **resource- and time-consuming**



Old(!) AlexNet (Krizhevsky et al., 2012) is a 8-layer CNN. 20–30 parameters
 Won the ImageNet Large Scale Visual Recognition Challenge 2012.

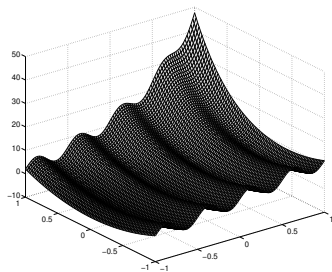
2 Local and grid search optimization for expensive-to-evaluate functions

Illustrative example

- Consider

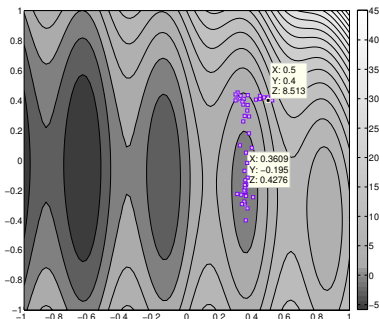
$$f: \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$x \mapsto f(x) = \exp(1.8(x_{[1]} + x_{[2]})) + 5x_{[1]} + 6x_{[2]}^2 + 3\sin(4\pi x_{[1]})$$



- Objective: find an approximation of $x^* = \arg \min_{x \in [-1, 1]^2} f(x)$ with a budget of $N = 60$ experiments

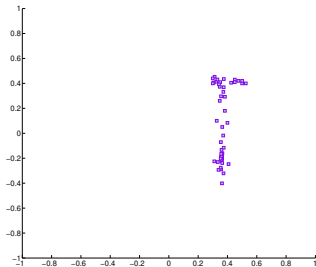
Evaluations points using a Nelder-Mead algorithm



→ the algorithm converges to a **local minimum** (≈ 0.427 , global minimum is ≈ -5.845))

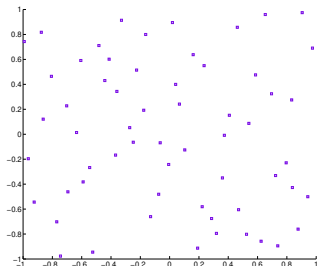
This comes as no surprise (local search algorithm). But above all...

- ~> after having spent the budget of (possibly expensive) evaluations, the behavior of the function is only known in a small region of the search domain



- ~> the global behavior of the function is unknown
- ~> potentially interesting regions have not been explored

How about sampling f uniformly on the search domain?



→ minimum of evaluation results is ≈ -5.823 (global minimum ≈ -5.845)

- Uniformly sampling tends to **minimize the fill-distance**

$$h_N = \max_{x \in \mathbb{X}} \min_i \|x - x_i\|$$

- From a theoretical perspective, minimizing h_N on a search domain is **minimax optimal** for approximating the function's optimum.

- However, **random search** or **grid search** strategies are generally inefficient because they **do not focus on** the more **promising regions** of the search space
- In situations where functions are expensive to evaluate and there is a limited budget for evaluations → strike a **balance** between **local search** and **exploration** of the search domain.

↔ **exploration vs exploitation** trade-off

3 Bayesian black-box modeling

- Let $f : \mathbb{X} \rightarrow \mathbb{R}$ be a real function defined on $\mathbb{X} \subseteq \mathbb{R}^d$, where \mathbb{X} is the input/parameter domain of the computer code under study
- f is a **black-box**, only known through evaluation/simulation results: **query an evaluation at x , observe the result**
- Given n simulations points $x_1, \dots, x_n \in \mathbb{X}$, denote by

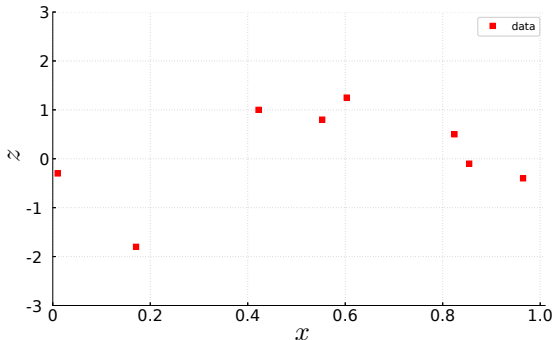
$$z_1 = f(x_1), \dots, z_n = f(x_n)$$

the corresp. simulation results (observations/evaluations of f)

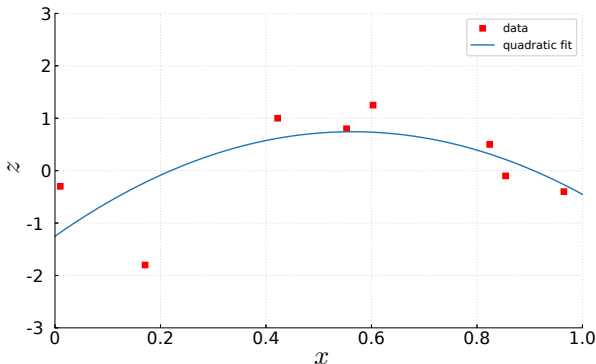
- **Our objective:** use the data $\mathcal{D}_n = (x_i, z_i)_{i=1\dots n}$ to **infer properties** about f
 - **Examples:**
 - given a new $x \in \mathbb{R}^d$, predict the value $f(x)$,
 - predict $x^* = \arg \max_x f(x)$, or $M = \max_x f(x)$
 - Predict the value of f at a given x ?
- the problem is that of constructing an **approximation** / an **estimator** \hat{f}_n of f from \mathcal{D}_n

A “curve fitting” problem

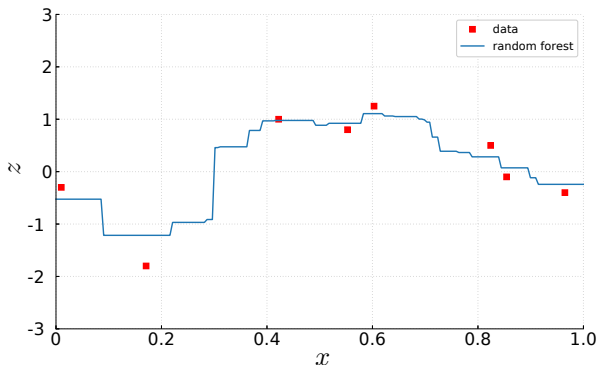
- We are given a data set of n simulation results, i.e., evaluations results of an unknown function $f : [0, 1] \rightarrow \mathbb{R}$, at points x_1, \dots, x_n .
- An example with $n = 8$:



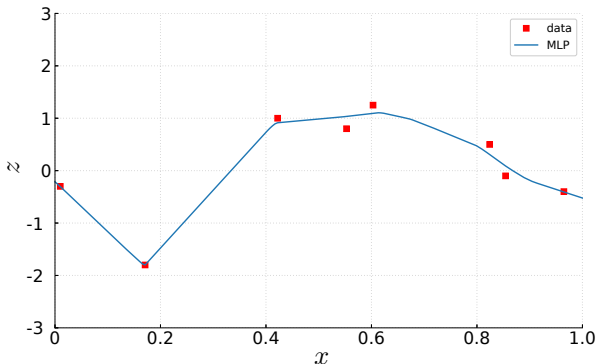
- Any approximation procedure of f consists in building a function $\hat{f}_n = h(\cdot; \theta)$, where θ is a vector of parameters, to be chosen using \mathcal{D}_n and available prior information
- e.g., quadratic fit using least-squares regression (try it in [colab](#))



- Any approximation procedure of f consists in building a function $\hat{f}_n = h(\cdot; \theta)$, where θ is a vector of parameters, to be chosen using \mathcal{D}_n and available prior information
- e.g., random forest regression (try it in [colab](#))



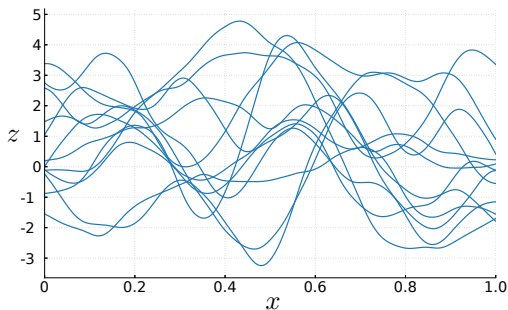
- Any approximation procedure of f consists in building a function $\hat{f}_n = h(\cdot; \theta)$, where θ is a vector of parameters, to be chosen using \mathcal{D}_n and available prior information
- e.g., neural network (multilayer perceptron) (try it in [colab](#))



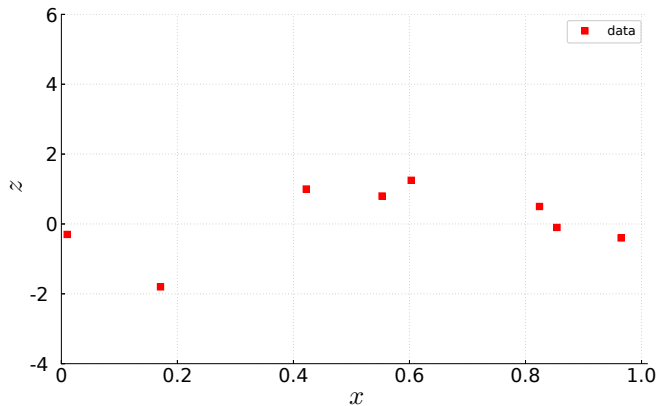
- Why choose one rather the other?
- Moreover, **no uncertainty quantification** in these approaches!

Bayesian approach: Gaussian processes

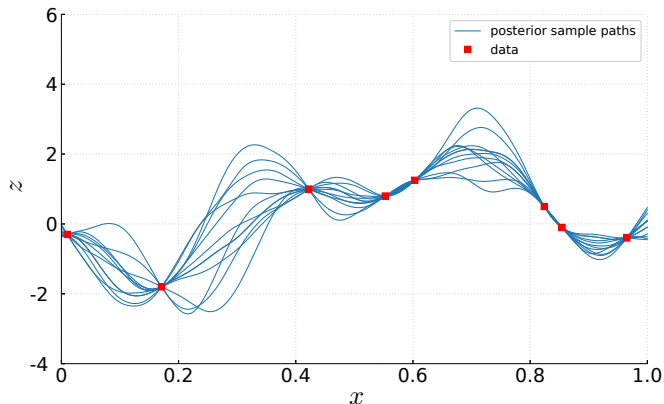
- A **non-parametric** route (for higher model capacity)
- f is modeled by a **random Gaussian process**, hereafter denoted by ξ , which encodes our prior knowledge about f



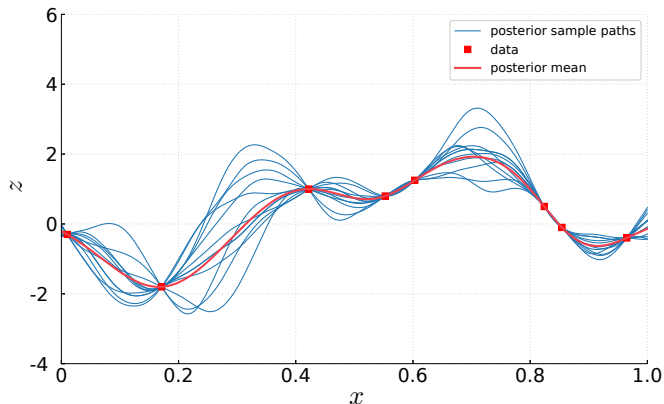
Data points / observations



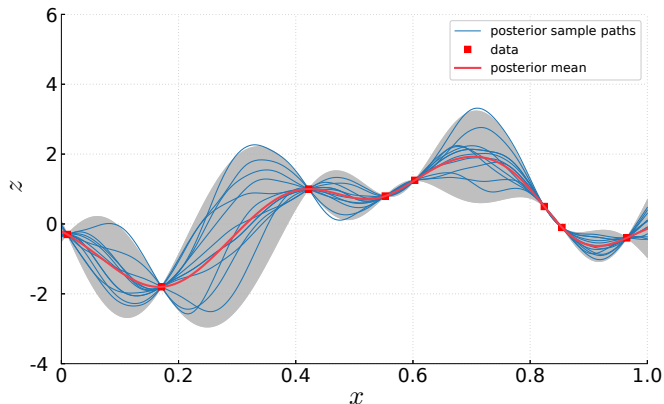
Posterior sample paths of ξ conditioned on observations
(try it in [colab](#))



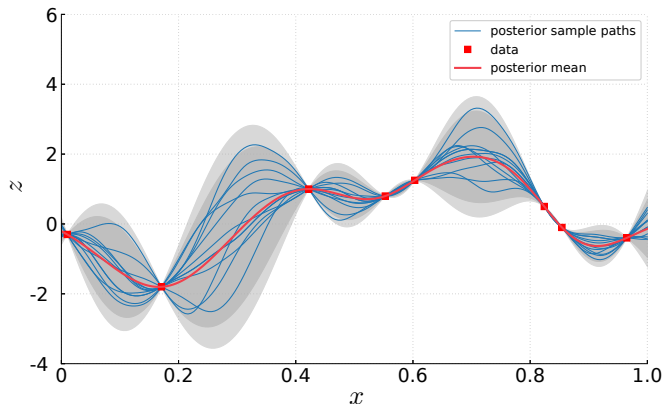
Posterior sample paths of ξ conditioned on observations
(try it in [colab](#))



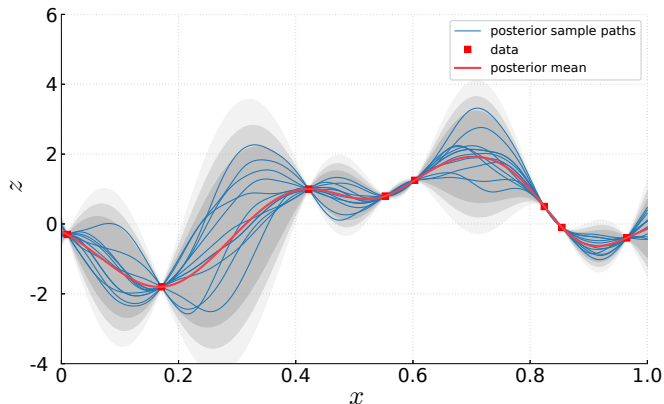
Posterior sample paths of ξ conditioned on observations
(try it in [colab](#))



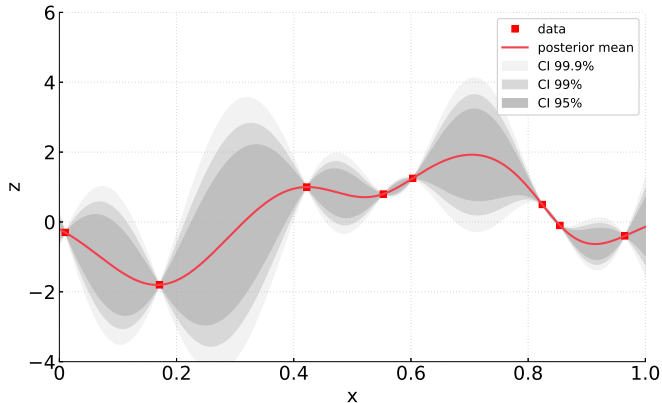
Posterior sample paths of ξ conditioned on observations
(try it in [colab](#))



Posterior sample paths of ξ conditioned on observations
(try it in [colab](#))



Posterior distribution of ξ , obtained by solving a system of linear equations \leftrightarrow
aka kriging prediction, G. Matheron \sim 1970



4 Bayesian optimization

- Objective: find an **approximation** of

$$\begin{cases} M & = \max_{x \in \mathbb{X}} f(x) \\ x^* & \in \arg \max_{x \in \mathbb{X}} f(x) \end{cases}$$

using a sequence of evaluations of f at points $X_1, X_2, \dots \in \mathbb{X}$

- The construction of an **optimization strategy** $\underline{X} : f \mapsto (X_1, X_2, X_3 \dots)$ is a **sequential decision problem**: at iteration n , we must choose a new evaluation point X_{n+1} using evaluation results of f at X_1, \dots, X_n

Bayesian optimization

- Start with a **loss function**: the efficiency of the optimization strategy \underline{X} at iteration n can be measured using the loss

$$\varepsilon_n(\underline{X}, f) = M - M_n$$

with $M_n = f(X_1) \vee \dots \vee f(X_n)$

- Bayesian optimization**: f is considered as a sample path of a **random process** ξ defined on some probability space $(\Omega, \mathcal{B}, P_0)$, with parameter $x \in \mathbb{X}$
- Information at iteration n is a realization of

$$\mathcal{F}_n = \{X_1, \xi(X_1) \dots, X_n, \xi(X_n)\}$$

- Denote by E_n the **conditional expectation** $E(\cdot | \mathcal{F}_n)$

- At each step, given the outcome for \mathcal{F}_n , we want to **minimize the risk** (expected loss)

$$H_n = \mathbb{E}_n[\varepsilon_n(\underline{X}, \xi)] = \mathbb{E}_n(M - M_n)$$

- H_n can be viewed as a measure of **residual uncertainty** about M
- At iteration n , we choose X_{n+1} to minimize the expectation of H_{n+1} :

$$X_{n+1} = \arg \min_{x \in \mathbb{X}} \mathbb{E}_n(H_{n+1} \mid X_{n+1} = x)$$

- This is a **Bayesian decision-theoretic point of view** for optimization, explored by **J. Mockus, A. Žilinskas** and their coauthors (~ 1970 – 1990)
- Extended to other settings under the name of **SUR**:
Sequential Uncertainty Reduction
(V. & Martinez 2006, V. & Bect 2009, Bect et al. 2012, Bect et al. 2019...)

- We have

$$\begin{aligned} X_{n+1} &= \arg \min_{x \in \mathbb{X}} \mathbf{E}_n (M - M_{n+1} \mid X_{n+1} = x) \\ &= \arg \max_{x \in \mathbb{X}} \mathbf{E}_n (M_{n+1} \mid X_{n+1} = x) \\ &= \arg \max_{x \in \mathbb{X}} \mathbf{E}_n (M_n \vee \xi(x)) \\ &= \arg \max_{x \in \mathbb{X}} \rho_n(x) \end{aligned}$$

with $\rho_n(x) := \mathbf{E}_n((\xi(x) - M_n)_+)$, and $z_+ = \max(z, 0)$.

- ρ_n corresponds to the **average excursion of $\xi(x)$ above the current maximum** of past evaluation results
- ρ_n is called the **expected improvement (EI)** sampling criterion

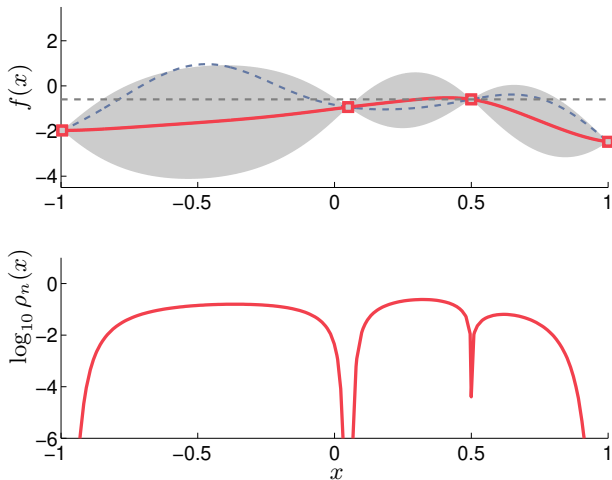
- The strategy $X_{n+1} = \arg \max_{x \in \mathbb{X}} \rho_n(x)$ using a **Gaussian process** for ξ has been popularized by D. Jones et al. 1998 under the name of **EGO (Efficient Global Optimization)**
- When ξ is a GP, with known mean and covariance functions, $\rho_n(x)$ has a closed-form expression:

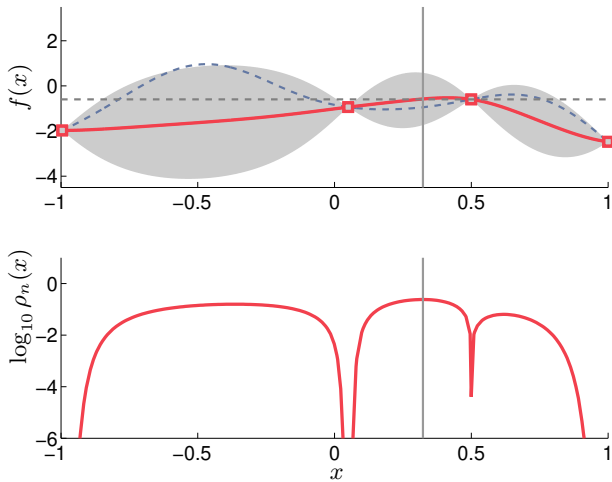
$$\rho_n(x) = \gamma \left(\widehat{\xi}_n(x; \underline{X}_n) - M_n, \sigma_n^2(x) \right),$$

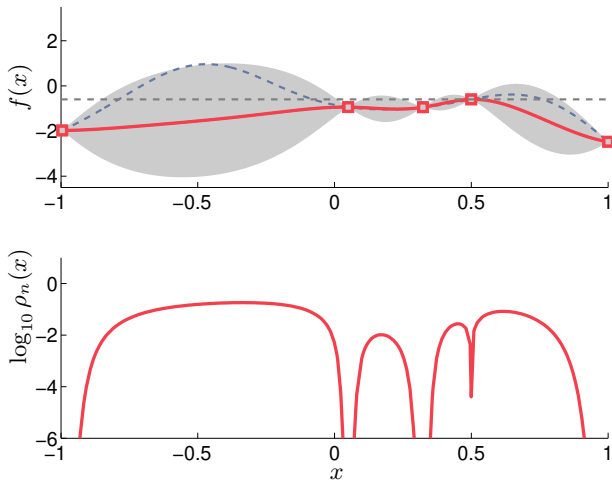
where

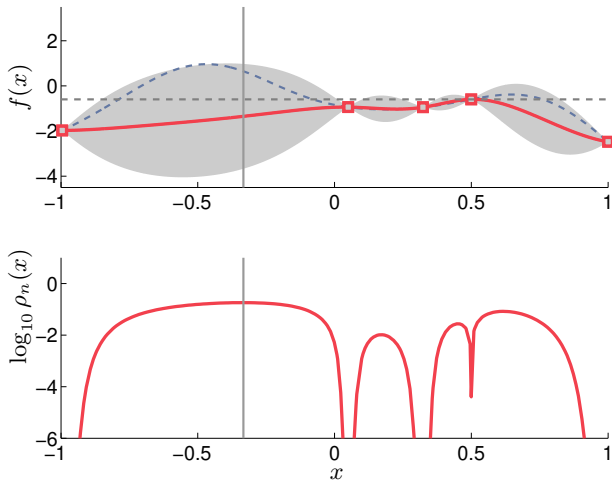
$$\gamma(z, s) = \begin{cases} \sqrt{s} \Phi' \left(\frac{z}{\sqrt{s}} \right) + z \Phi \left(\frac{z}{\sqrt{s}} \right) & \text{if } s > 0, \\ \max(z, 0) & \text{if } s = 0. \end{cases}$$

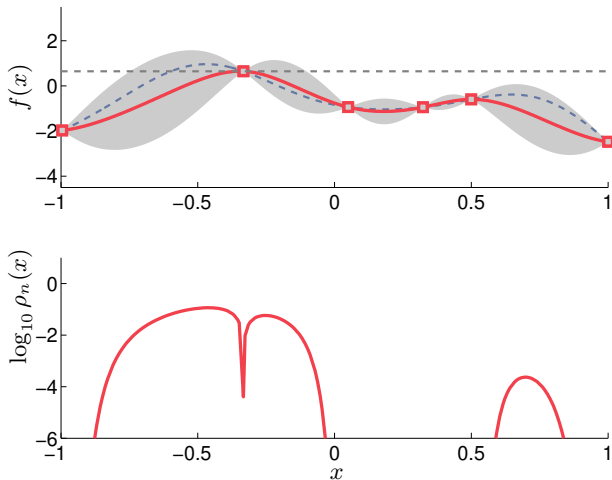
and $\widehat{\xi}_n(x; \underline{X}_n)$ and $\sigma_n^2(x)$ are the **posterior mean** and the **posterior variance** of $\xi(x)$

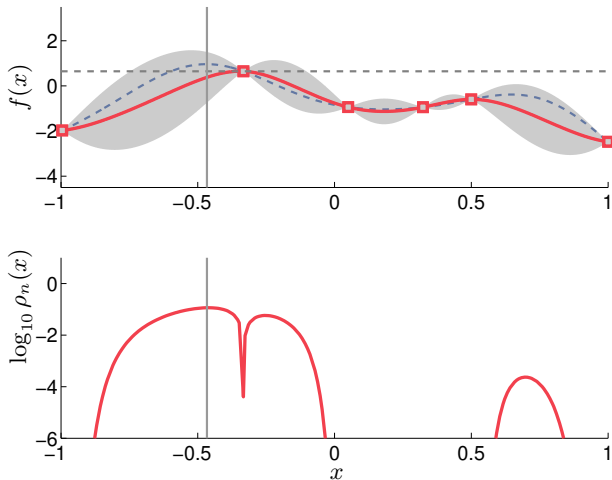
Function to be maximized (dashed blue line) and kriging prediction

Function to be maximized (dashed blue line) and kriging prediction

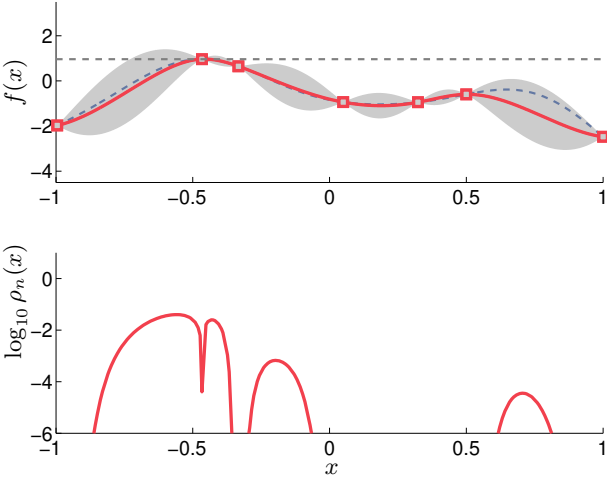
Function to be maximized (dashed blue line) and kriging prediction

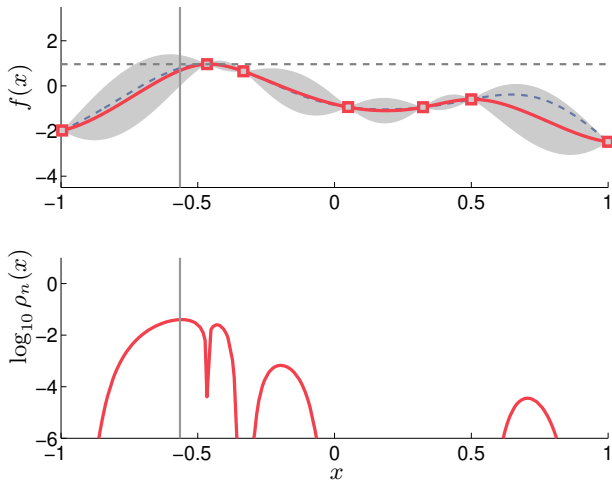
Function to be maximized (dashed blue line) and kriging prediction

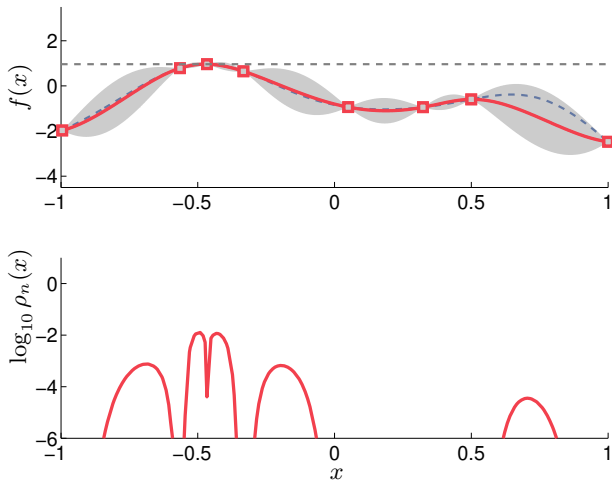
Function to be maximized (dashed blue line) and kriging prediction

Function to be maximized (dashed blue line) and kriging prediction

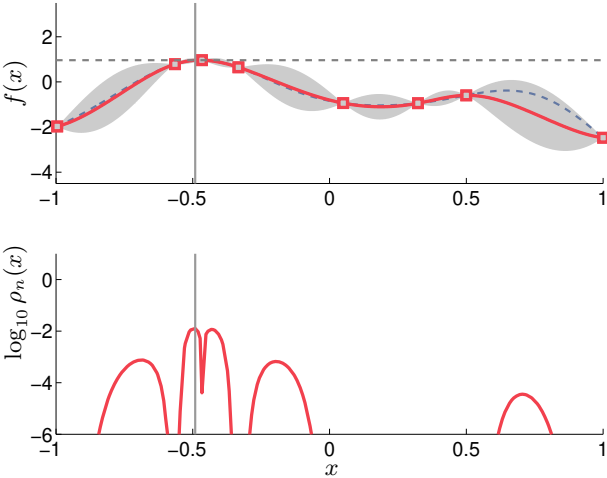
Function to be maximized (dashed blue line) and kriging prediction

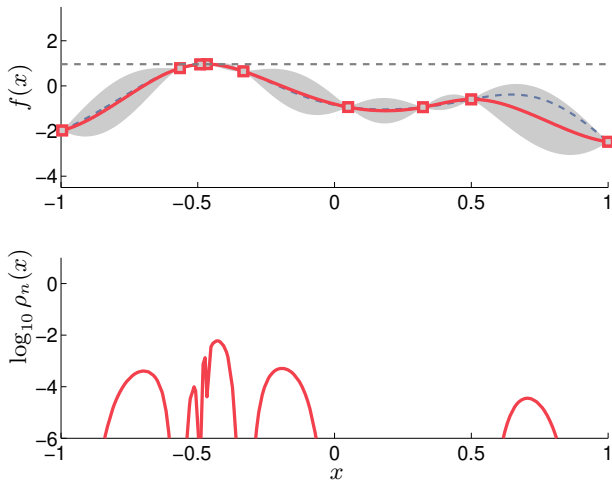


Function to be maximized (dashed blue line) and kriging prediction

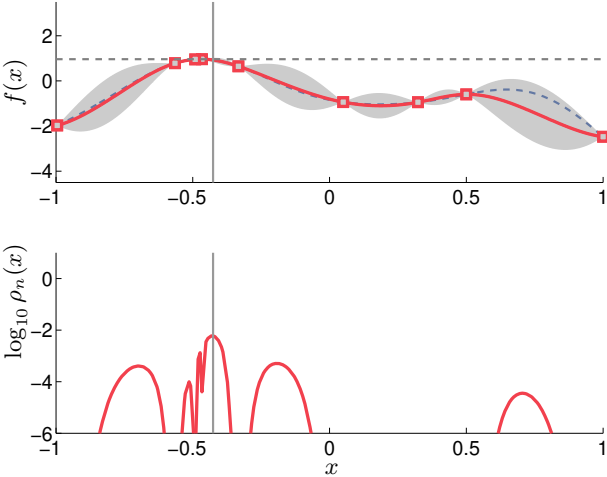
Function to be maximized (dashed blue line) and kriging prediction

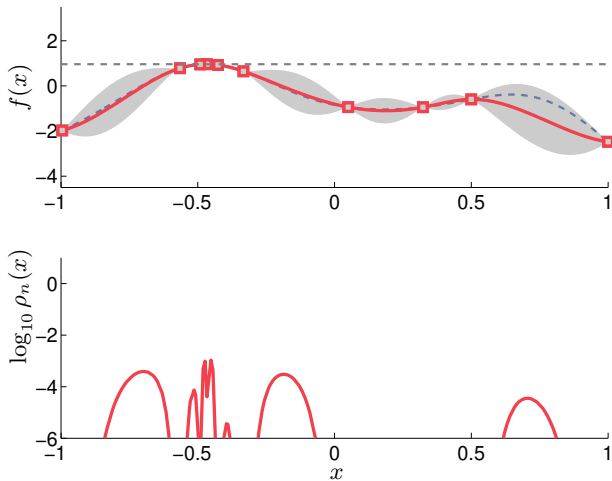
Function to be maximized (dashed blue line) and kriging prediction



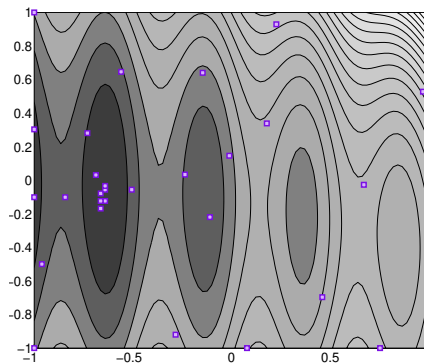
Function to be maximized (dashed blue line) and kriging prediction

Function to be maximized (dashed blue line) and kriging prediction



Function to be maximized (dashed blue line) and kriging prediction

2d illustration



Exploration vs Exploitation (31 evaluations)

	M_n with $N = 60$
LHS	-5.823
DIRECT	-5.839
EI (31 evals)	-5.845
Global minimum	-5.845

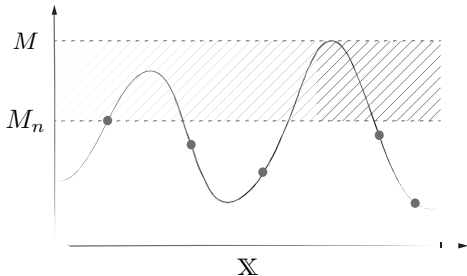
- The EGO/EI algorithm has been used in **countless** of engineering design problems and is widely available (Python, Matlab, R...)
- The **convergence of the algorithm** is well established (V. & Bect 2010)
- Convergence rates are partially known (Bull 2011, Ryzhov 2016...)

Other Bayesian strategies for mono-objective optimization?

- In a **global optimization** problem, it is generally of interest to obtain a good approximation of **both M and x^***
- The loss function $\varepsilon_n(\underline{X}, f) = M - M_n$ does not measure directly the distance of x_n^* to x^*
- Better alternative(s)?

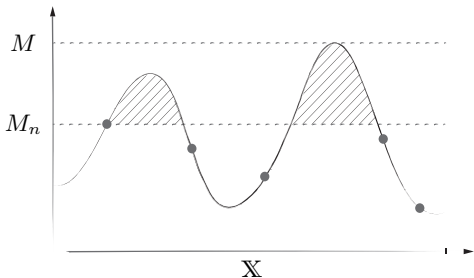
Expected integrated expected improvement

- Note that $\varepsilon_n(\underline{X}, f) = M - M_n \propto \lambda(\mathbb{X}) (M - M_n)$



→ coarse measure of the uncertainty about the pair (M, x^*)

- Idea: use the **integral loss** $\varepsilon_n(\underline{X}, f) = \int_{\mathbf{X}} (f(x) - M_n)_+ \lambda(dx)$



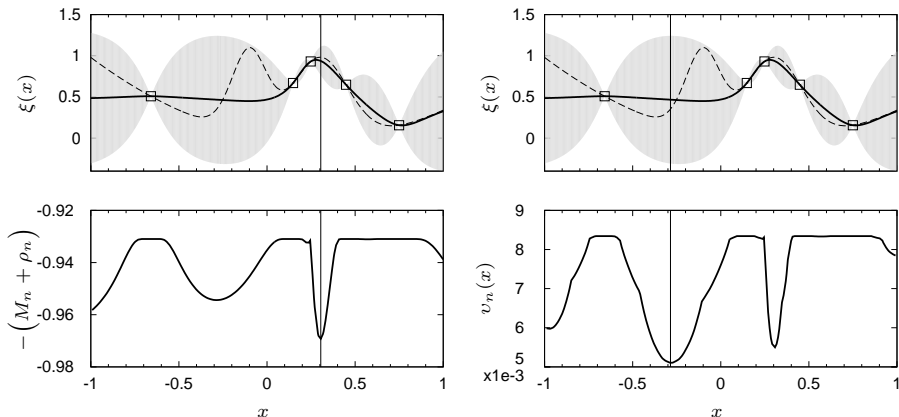
→ ε_n gets smaller when uncertainty about x^* decreases

- SUR strategy:

$$\begin{aligned} X_{n+1} &= \arg \min_{x \in \mathbb{X}} \mathbf{E}_n \left(\int_{\mathbb{X}} (\xi(y) - M_{n+1})_+ dy \mid X_{n+1} = x \right) \\ &= \arg \min_{x \in \mathbb{X}} \mathbf{E}_n \left(\int_{\mathbb{X}} \mathbf{E}_{n+1}((\xi(y) - M_{n+1})_+) dy \mid X_{n+1} = x \right) \\ &= \arg \min_{x \in \mathbb{X}} v_n(x) := \mathbf{E}_n \left(\int_{\mathbb{X}} \rho_{n+1}(y) dy \mid X_{n+1} = x \right) \end{aligned}$$

- We call v_n the **Expected Integrated Expected Improvement (EI²)** (V. & Bect, 2014)

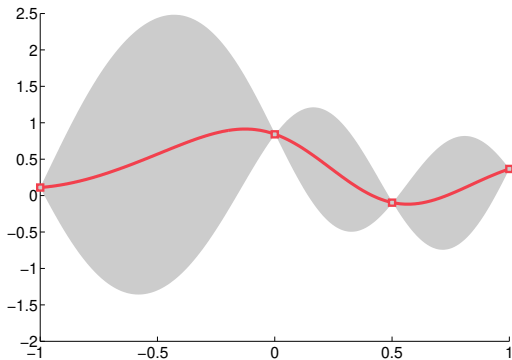
Large expected improvement in a small region, smaller expected improvement over a large region of the search domain \rightarrow here, v_n favors better exploration than ρ_n



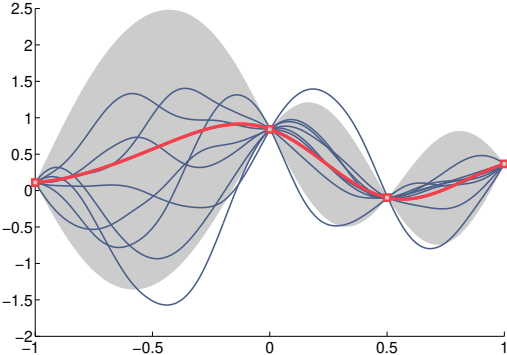
Entropy-based strategies

- **Informational Approach to Global Optimization (IAGO)**
(Villemonaix, Walter & V., 2007–2008)
- **Entropy Search** for Information-Efficient Global Optimization
(Hennig, Schuler, 2012)
- **Predictive Entropy Search** for Efficient Global Optimization of Black-box Functions
(Hernandez-Lobato, Hoffman, Ghahramani, 2014)
- ...

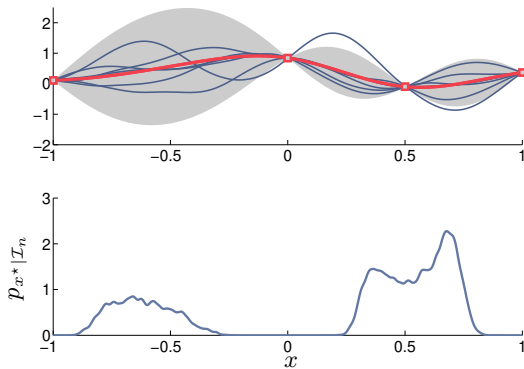
Empirical posterior density of the minimizer



Empirical posterior density of the minimizer



Empirical posterior density of the minimizer



- Assumption: search domain \mathbb{X} is finite
- Define a loss function as the residual uncertainty about x^* measured using the Shannon entropy

$$\varepsilon_n(\underline{X}, \xi) = H(x^*; \mathcal{F}_n) = - \sum_{x \in \mathbb{X}} P_n(x^* = x) \log P_n(x^* = x),$$

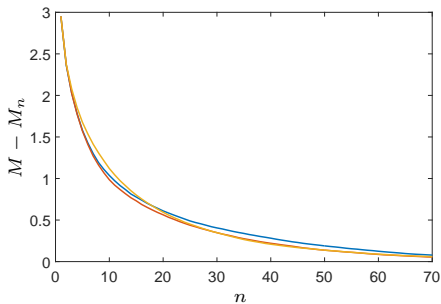
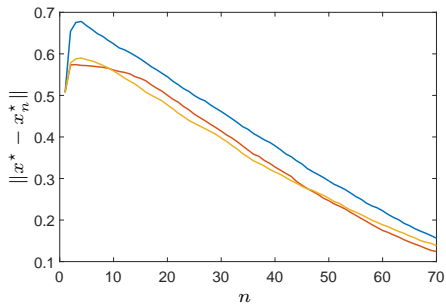
with P_n the conditional distribution $P_0(\cdot | \mathcal{F}_n)$

- SUR strategy:

$$X_{n+1} = \arg \min_{x \in \mathbb{X}} E_n (\varepsilon_{n+1}(\underline{X}, \xi) | X_{n+1} = x)$$

Empirical comparison

Average errors using EI (blue), EI² (red), and IAGO (yellow), from 3000 sample paths of a GP on \mathbb{R}^3 , with zero-mean and isotropic Matérn covariance function, simulated on a set of 1000 points in $[0, 1]^3$.



Multi-objective optimization

- Consider a set of objective functions $f_j : \mathbb{X} \rightarrow \mathbb{R}$, $j = 1, \dots, p$,
to be minimized
- Objective: build an approximation of the **Pareto front**

$$\Gamma = \{x \in \mathbb{X} : \nexists x' \in \mathbb{X} \text{ such that } f(x') \prec f(x)\},$$

where \prec stands for the **Pareto domination rule** defined by

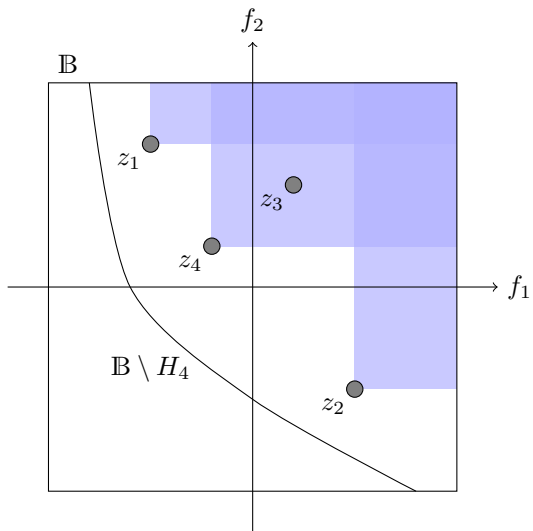
$$z' = (z'_1, \dots, z'_p) \prec z = (z_1, \dots, z_p) \iff \begin{cases} \forall i \leq p, & z'_i \leq z_i, \\ \exists j \leq p, & z'_j < z_j. \end{cases}$$

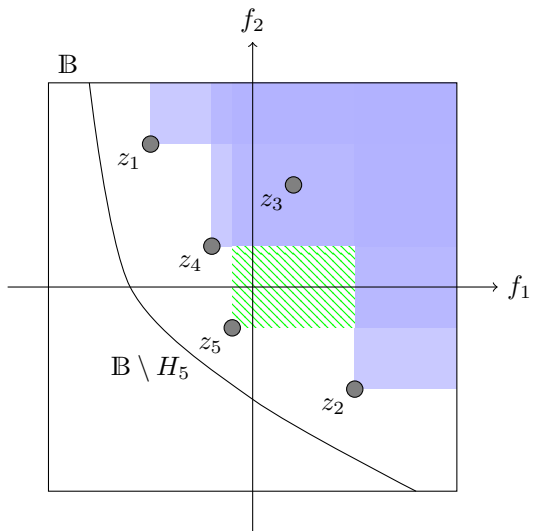
- Define

- $\mathbb{B} = \{z \in \mathbb{R}^p; z \leq z^{\text{upp}}\}, z^{\text{upp}} \in \mathbb{R}^p$

- $H = \{z \in \mathbb{B}; \exists x \in \mathbb{X}, f(x) \prec z\}$

- $H_n = \{z \in \mathbb{B}; \exists i \leq n, f(X_i) \prec z\}$





- Idea: use the **volume of the non-dominated region** as a **loss function**:

$$\varepsilon_n(\underline{X}, f) = |\mathbb{B} \setminus H_n| ,$$

where

$$H_n = \{z \in \mathbb{B}; \exists i \leq n, f(X_i) \prec z\} ,$$

- Define the **improvement** yielded by a new evaluation result

$$f(X_{n+1}) = (f_1(X_{n+1}), \dots, f_p(X_{n+1}))$$

as the **increase of the volume of the dominated region**:

$$I_n(X_{n+1}) = |\mathbb{B} \setminus H_n| - |\mathbb{B} \setminus H_{n+1}| = |H_{n+1} \setminus H_n| = |H_{n+1}| - |H_n| ,$$

since $H_n \subset H_{n+1} \subset H$.

Expected hyper-volume improvement (EHVI)

- Given a vector-valued Gaussian random process model $\xi = (\xi_1, \dots, \xi_p)$ of $f = (f_1, \dots, f_p)$, define a **multi-objective EI criterion** as

$$\begin{aligned} \rho_n(x) &= \mathbb{E}_n (I_n(x)) \\ &= \mathbb{E}_n \left(\int_{\mathbb{B} \setminus H_n} \mathbb{1}_{\xi(x) \prec z} dz \right) \\ &= \int_{\mathbb{B} \setminus H_n} \mathbb{E}_n (\mathbb{1}_{\xi(x) \prec z}) dz \\ &= \int_{\mathbb{B} \setminus H_n} \mathbb{P}_n (\xi(x) \prec z) dz, \end{aligned}$$

- First proposed by Emmerich and coworkers (2005–2008)
- Extension to the **constrained** multi-objective setting (Feliot, Bect & V. 2017)

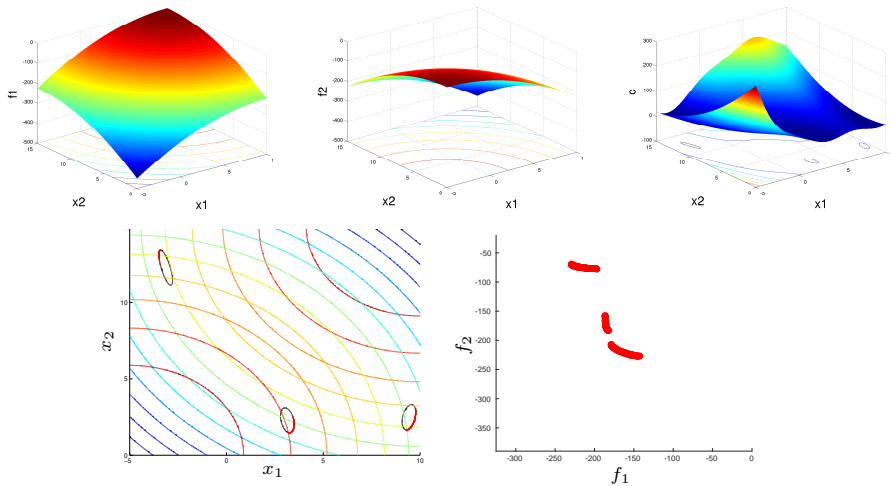
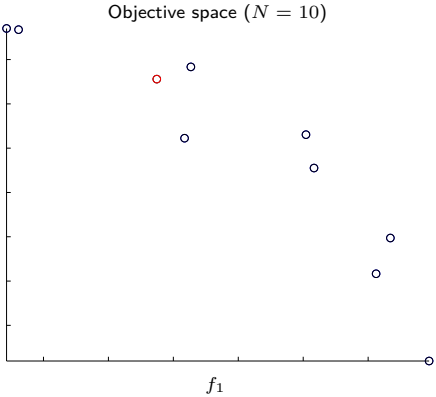
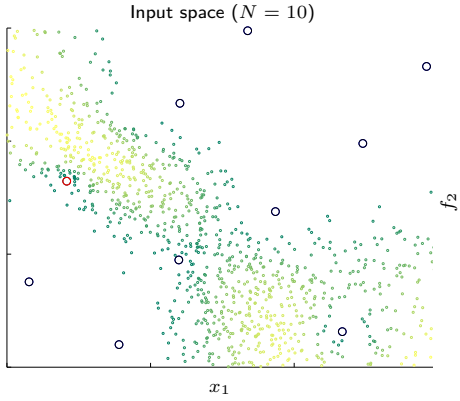
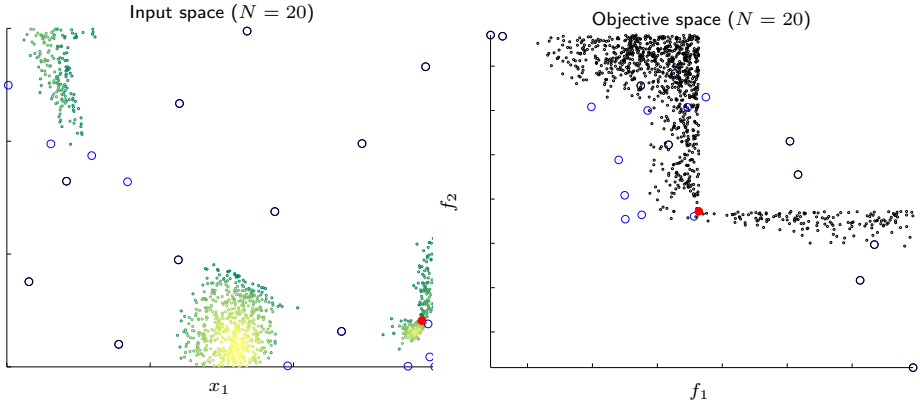


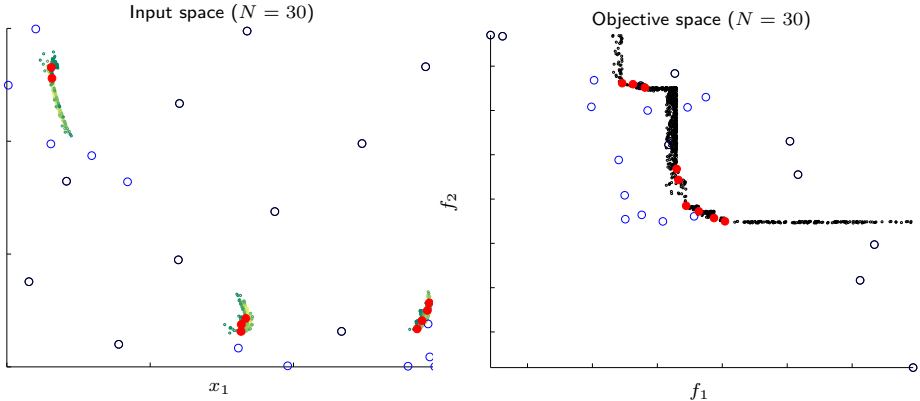
Illustration with $d = 2, p = 2, q = 1$



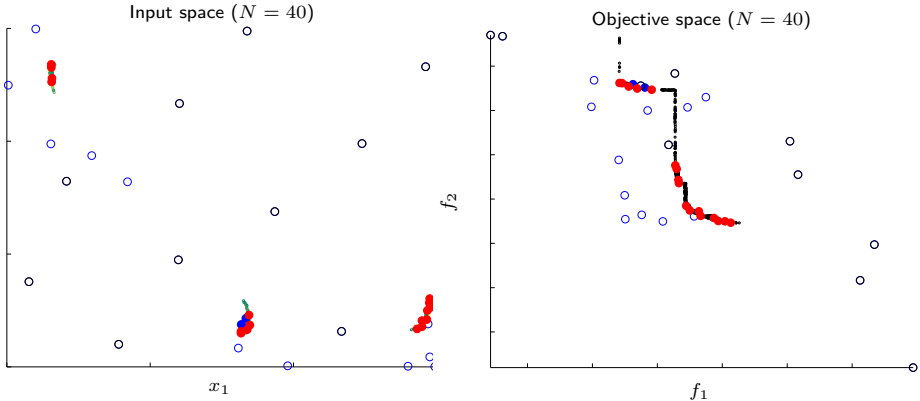
BMOO with Sequential Monte Carlo (Feliot, Bect & Vazquez 2017)



BMOO with Sequential Monte Carlo (Feliot, Bect & Vazquez 2017)



BMOO with Sequential Monte Carlo (Feliot, Bect & Vazquez 2017)



BMOO with Sequential Monte Carlo (Feliot, Bect & Vazquez 2017)

5 Concluding remarks

- **Bayesian optimization**, along with Bayesian techniques for the estimation of **excursion sets** (a closely related problem), have been **prominent topics** in recent years (2015–today)
- **Numerous publications** in the ML and UQ communities on applications, sampling criteria, high dimensionality. . .
- Software is available: scikit-learn, BoTorch. . .
- However, I strongly advise against using these implementations in black-box mode: **performance depends on the GP model!**
- Please feel free to contact me if you want to discuss about BO