

A Hierarchical Expected Improvement method for effective Bayesian optimization

C. F. Jeff Wu

School of Data Science
Chinese University of Hong Kong, Shenzhen

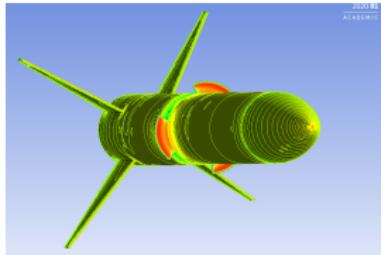
Joint work with **Zhehui Chen** (Google) and **Simon Mak** (Duke University)

¹The paper appeared in JASA 2024.

Outline

- Why **Bayesian optimization?**
- **Kriging** and **Expected Improvement (EI)**
- **Hierarchical Kriging** and **Hierarchical EI (HEI)**
- **Hyperprior specification** for HEI
- Global **convergence rates**
- Two **applications**: process optimization for **semiconductor manufacturing**, hyperparameter tuning for **deep learning**

Motivation



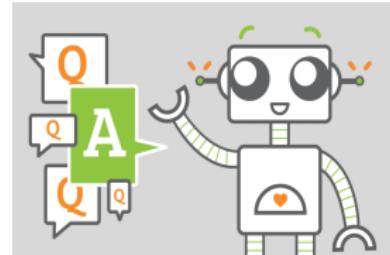
Engineering design



Autonomous vehicles



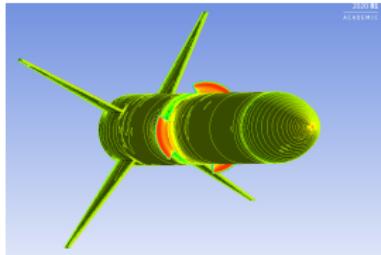
Portfolio management



Q&A bot

Want to find the **best** setting(s) under a **black-box** objective f'n

Motivation



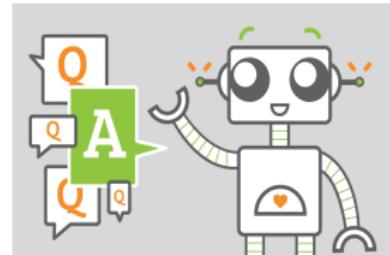
Engineering design



Autonomous vehicles



Portfolio management

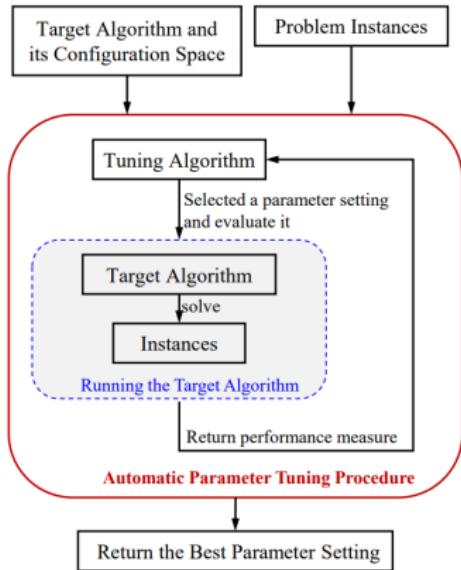


Q&A bot

Want to find the **best** setting(s) under a **black-box** objective f'n

Application in AI

Black-box optimization can be used to tune **hyper-parameters** automatically.



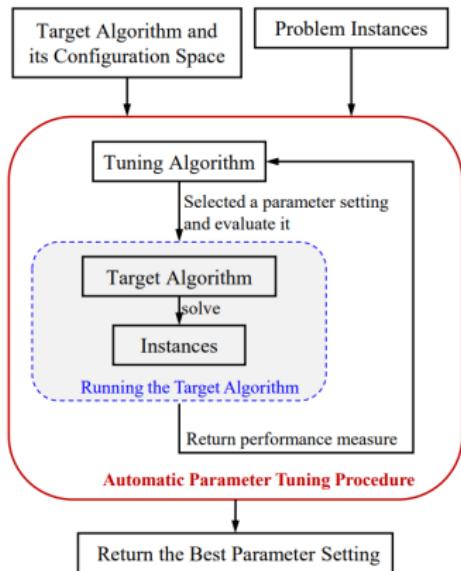
Examples of hyper-parameters in neural networks

- learning rate: $\log U(0.001, 10)$;
- n.layers: 1 to 3;
- batch size: 20 or 100;
- n.hidden units: $\log U(128, 4096)$;

...

Application in AI

Black-box optimization can be used to tune **hyper-parameters** automatically.



Examples of hyper-parameters in neural networks

- learning rate: $\log U(0.001, 10)$;
- n.layers: 1 to 3;
- batch size: 20 or 100;
- n.hidden units: $\log U(128, 4096)$;
- ...

Black-box optimization

$$\boldsymbol{x}^* \in \arg \min_{\boldsymbol{x} \in \Omega} f(\boldsymbol{x})$$

- Ω : **feasible** domain for parameters \boldsymbol{x}
- f : an unknown **black-box objective** function
 - Functional form **unknown** prior to **data**
 - **Queries** on f require **experimentation** and/or **simulations**, can be **expensive**



Goal: find a **good** solution to the **optimization problem** using as **few queries** on f as possible

Black-box optimization

$$\boldsymbol{x}^* \in \arg \min_{\boldsymbol{x} \in \Omega} f(\boldsymbol{x})$$

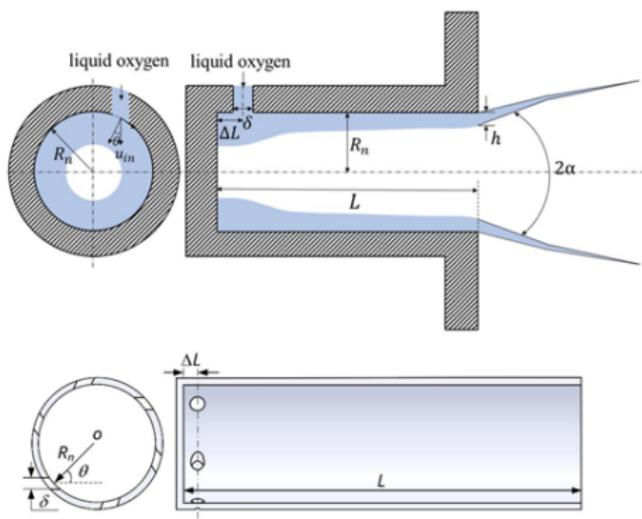
- Ω : **feasible** domain for parameters \boldsymbol{x}
- f : an unknown **black-box objective** function
 - Functional form **unknown** prior to **data**
 - **Queries** on f require **experimentation** and/or **simulations**, can be **expensive**



Goal: find a **good** solution to the **optimization problem** using as **few** **queries** on f as possible

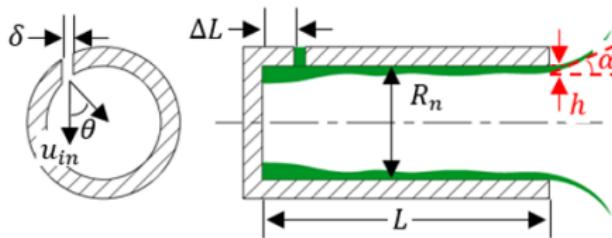
Rocket design

Collaboration project with GT Aerospace (Air Force / SpaceX)

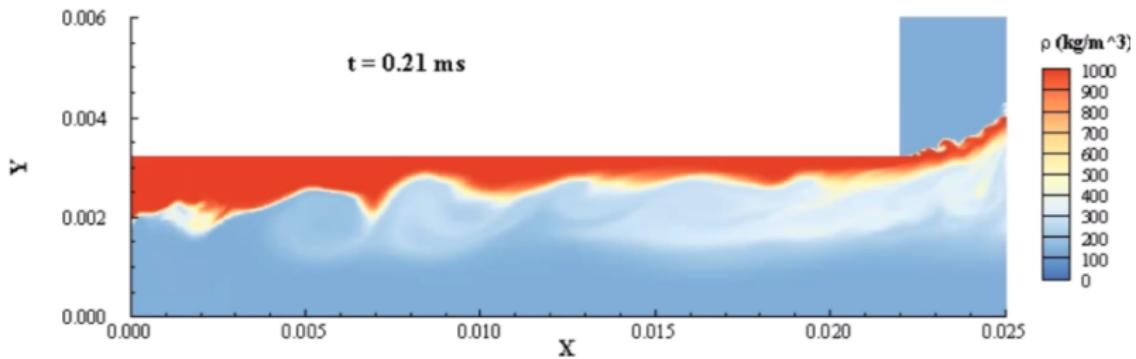


Rocket design

Optimizing **rocket injector** designs over a geometric **design space** (Mak et al. 2018; Yeh et al. 2018; Chang et al. 2019)



L	20-100 mm
R_n	2-5 mm
θ	45°-75°
δ	0.5-2 mm
ΔL	1-4 mm



Rocket design

Simulations are very **expensive**: **14 days** per **injector design** (each evaluation of f), **parallelized** over millions of cores

- **Exploring** the 5D domain Ω by simulations **alone** will take 1000's of simulations \Rightarrow **years** of **computation**

One **solution**: build a **surrogate model** on f

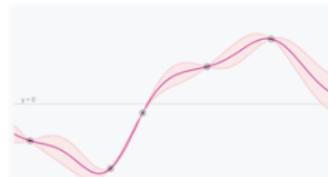
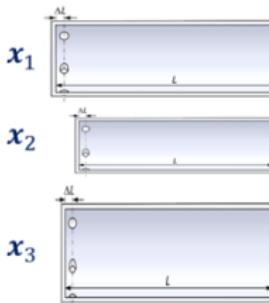
This is **Bayesian optimization!**

Rocket design

Simulations are very **expensive**: **14 days** per **injector design** (each evaluation of f), **parallelized** over millions of cores

- **Exploring** the 5D domain Ω by simulations **alone** will take 1000's of simulations \Rightarrow **years** of **computation**

One **solution**: build a **surrogate model** on f



Run a few experiments at different injector geometries

Train a predictive model using simulation data

Use this model to choose a new geometry for simulation

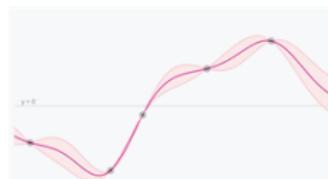
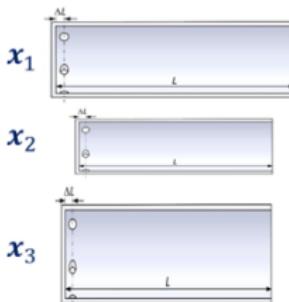
This is **Bayesian optimization!**

Rocket design

Simulations are very **expensive**: **14 days** per **injector design** (each evaluation of f), **parallelized** over millions of cores

- **Exploring** the 5D domain Ω by simulations **alone** will take 1000's of simulations \Rightarrow **years** of **computation**

One **solution**: build a **surrogate model** on f



Run a few experiments at different injector geometries

Train a predictive model using simulation data

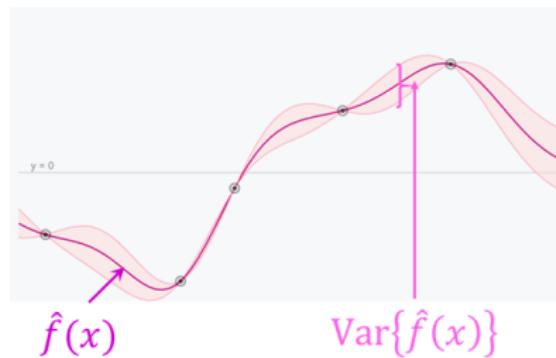
Use this model to choose a new geometry for simulation

This is **Bayesian optimization!**

Bayesian optimization

Two **ingredients** for Bayesian optimization:

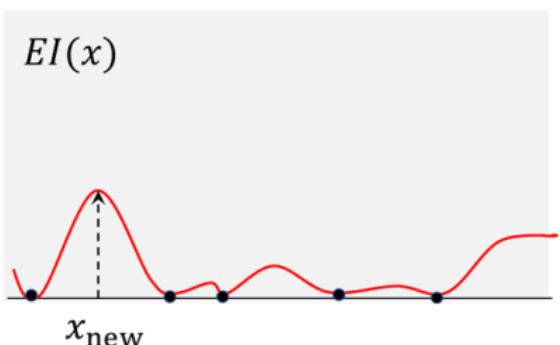
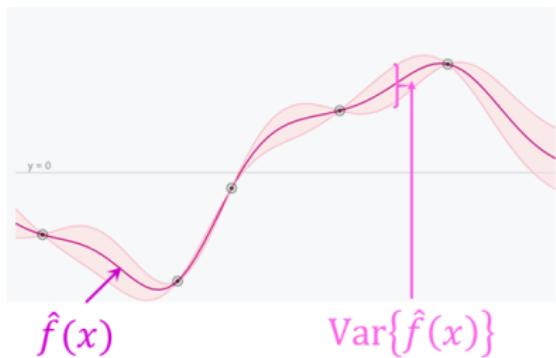
- Bayesian (probabilistic) model on **objective f'n f**:
 - Provide **predictor** of f : $\hat{f} = \mathbb{E}\{\hat{f}| \text{data}\}$
 - Quantify **uncertainty** on f : $\text{Var}\{\hat{f}| \text{data}\}$
- Acquisition function for choosing next point x_{new}
 - **Exploitation**: want x 's with low **predictions** \hat{f}
 - **Exploration**: want x 's with high **uncertainty** $\text{Var}\{\hat{f}\}$



Bayesian optimization

Two **ingredients** for Bayesian optimization:

- **Bayesian** (probabilistic) model on **objective f'n f**:
 - Provide **predictor** of f : $\hat{f} = \mathbb{E}\{\hat{f}| \text{data}\}$
 - Quantify **uncertainty** on f : $\text{Var}\{\hat{f}| \text{data}\}$
- **Acquisition function** for choosing **next** point x_{new}
 - **Exploitation**: want x 's with low **predictions** \hat{f}
 - **Exploration**: want x 's with high **uncertainty** $\text{Var}\{\hat{f}\}$



Bayesian optimization

A **widely-used** setting in **Bayesian optimization**:

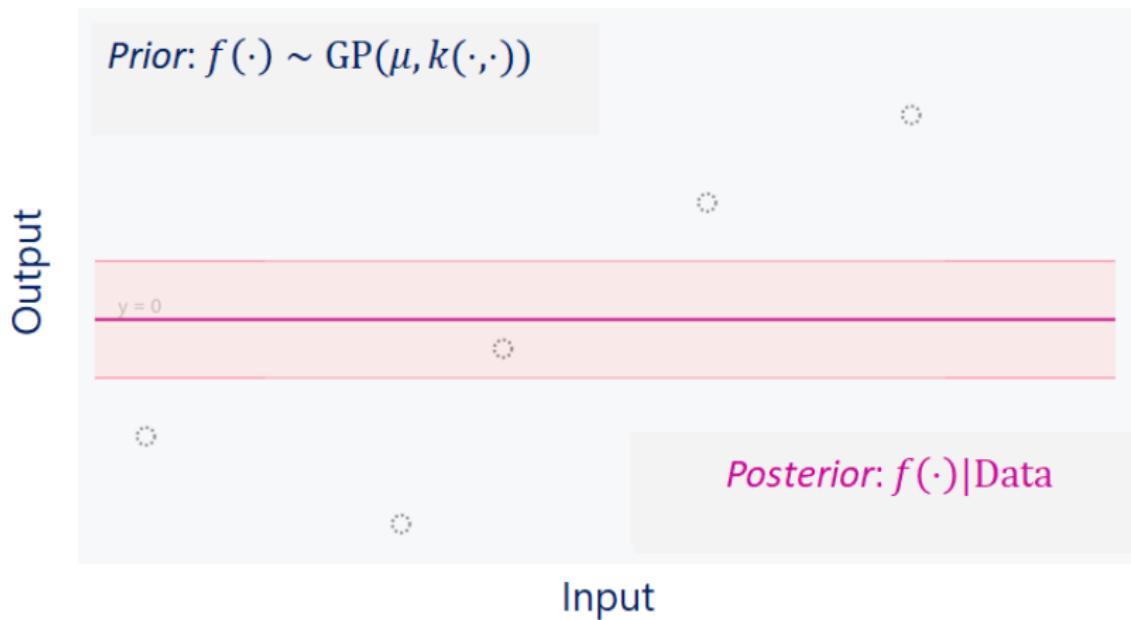
- **Bayesian** model on f : **Gaussian process** (kriging)
 - **Nonparametric** form \Rightarrow **powerful approximation**
 - **Probabilistic** model \Rightarrow **uncertainty quantification** (UQ)
 - **Closed-form** prediction & UQ \Rightarrow **efficient modeling**
- Acquisition function: **Expected Improvement** (Jones et al. 1997)
 - **Closed-form** expression for $\text{EI}(x) \Rightarrow$ **efficient queries**
 - **Interpretable**

Bayesian optimization

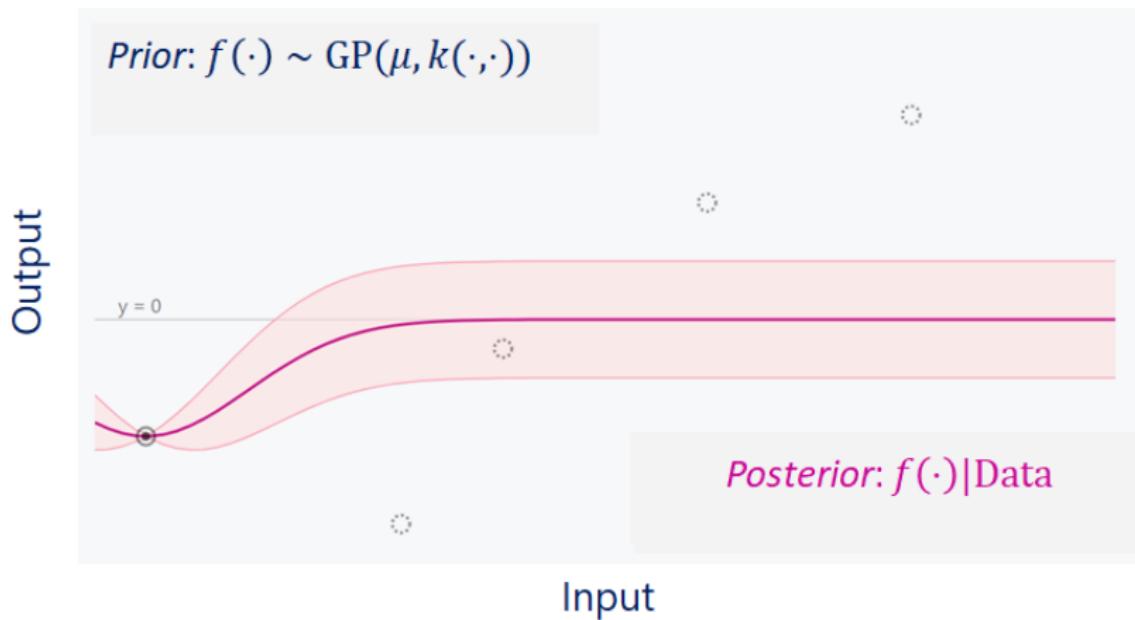
A **widely-used** setting in Bayesian optimization:

- Bayesian model on f : **Gaussian process** (kriging)
 - Nonparametric form \Rightarrow **powerful approximation**
 - Probabilistic model \Rightarrow **uncertainty quantification** (UQ)
 - Closed-form prediction & UQ \Rightarrow **efficient modeling**
- Acquisition function: **Expected Improvement** (Jones et al. 1997)
 - Closed-form expression for $\text{EI}(x)$ \Rightarrow **efficient queries**
 - **Interpretable**

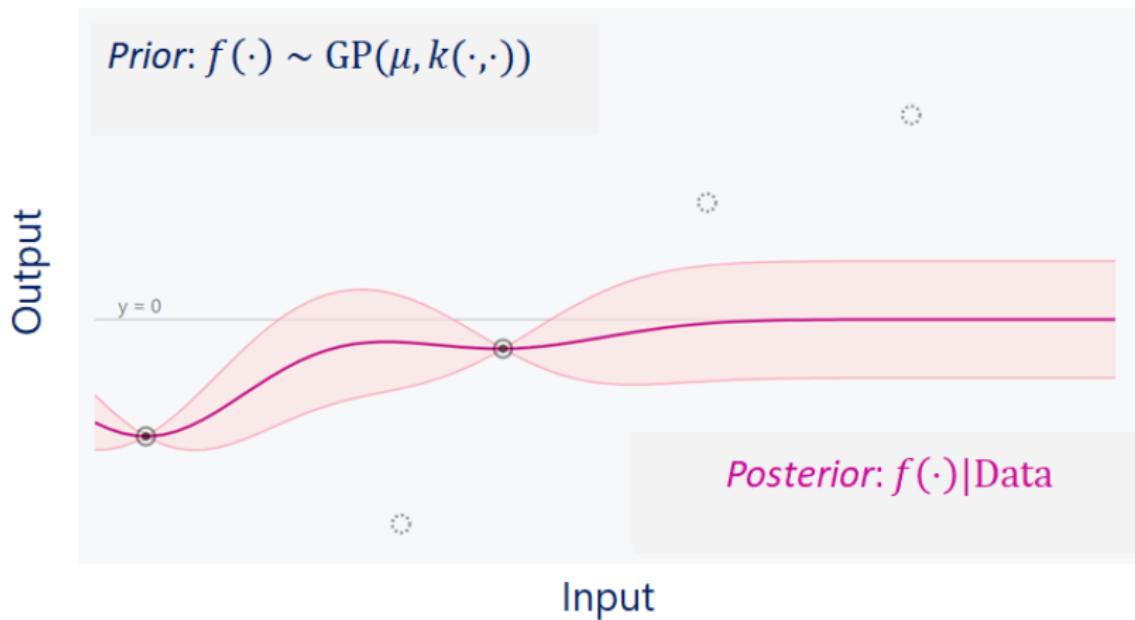
Gaussian process (kriging) model



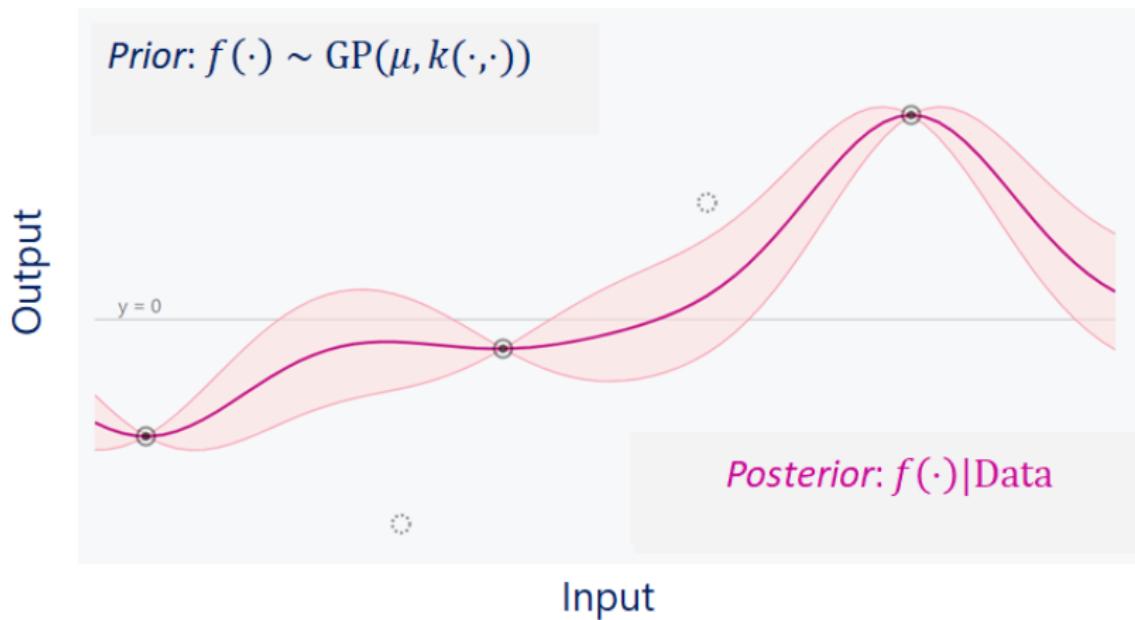
Gaussian process (kriging) model



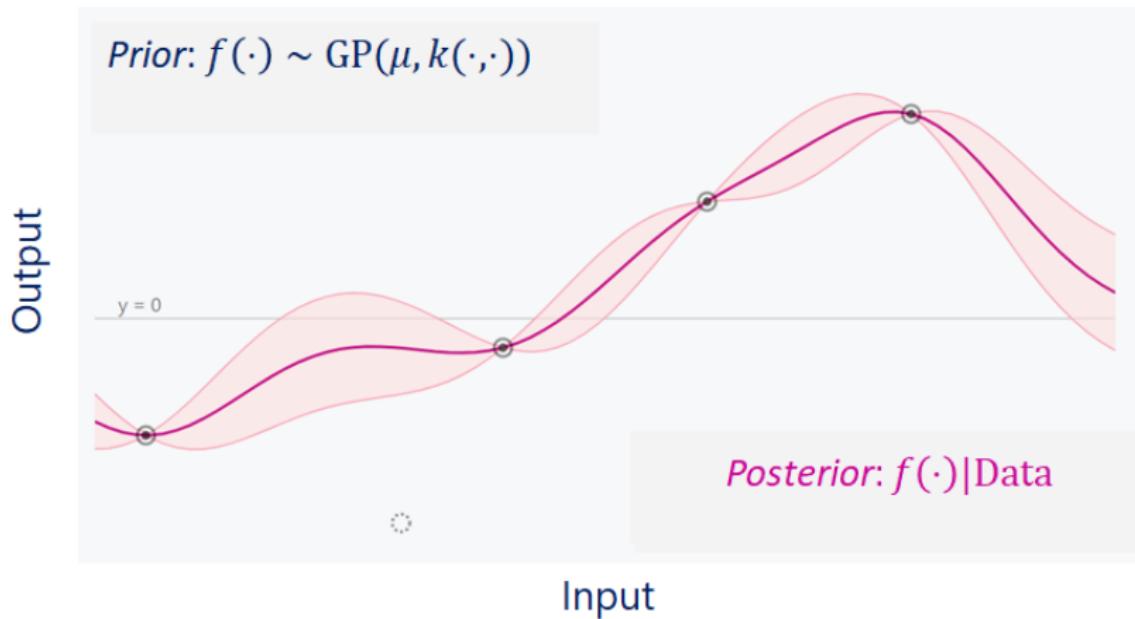
Gaussian process (kriging) model



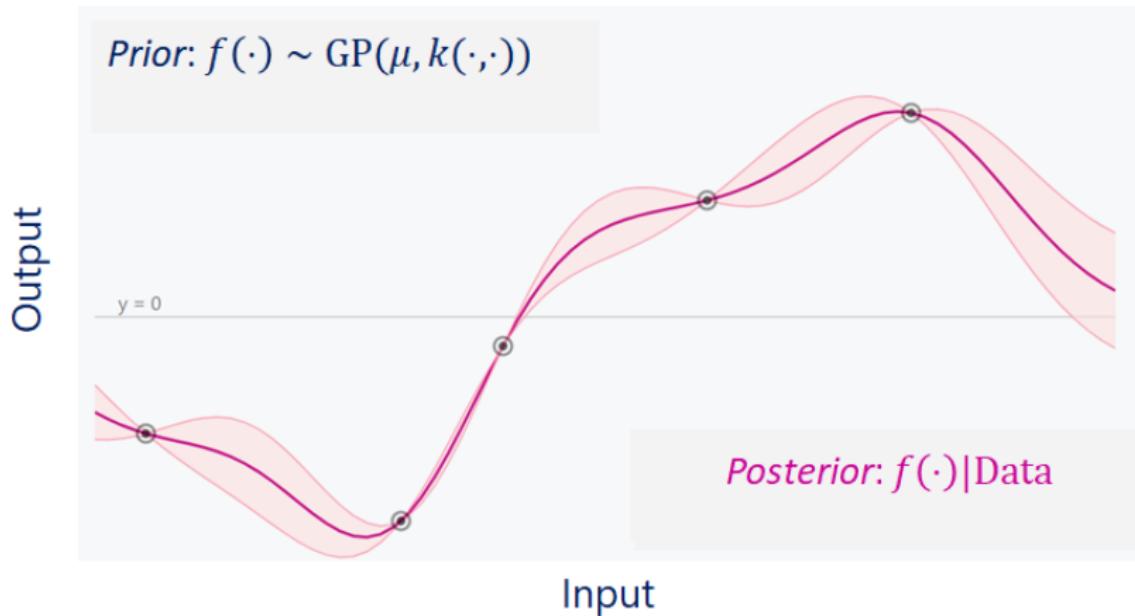
Gaussian process (kriging) model



Gaussian process (kriging) model



Gaussian process (kriging) model



Gaussian process (kriging) model

Universal kriging model (see, e.g., Santner et al. 2013):

$$f(\cdot) \sim GP\{\mu(\cdot), \sigma^2 r(\cdot, \cdot)\}, \quad \boldsymbol{x} \in \mathbb{R}^d.$$

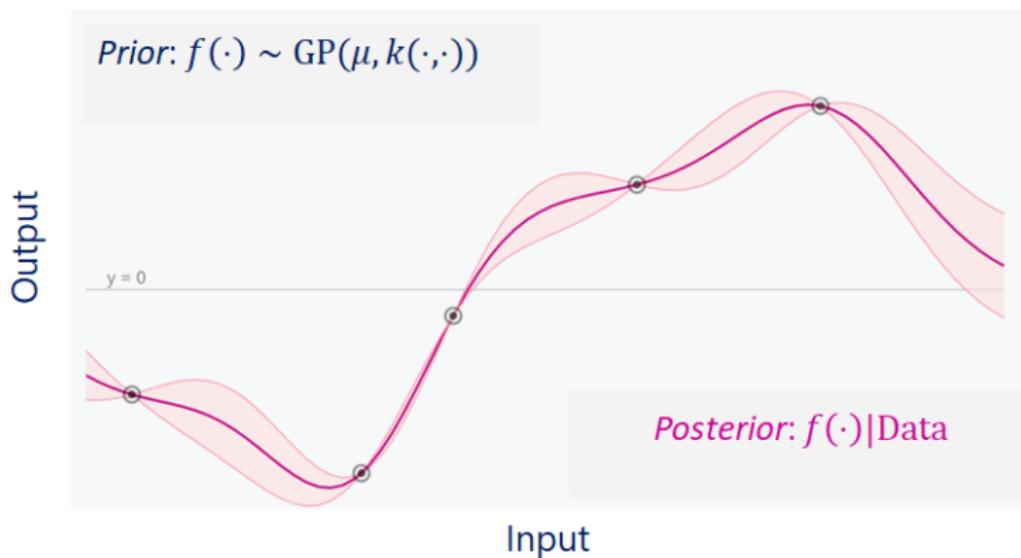
- Mean function: $\mu(\boldsymbol{x}) = \boldsymbol{p}^\top(\boldsymbol{x})\boldsymbol{\beta}$
- Basis functions: $\boldsymbol{p}(\boldsymbol{x}) = [p_1(\boldsymbol{x}), \dots, p_q(\boldsymbol{x})]^\top$
- Process variance: $\sigma^2 = \text{Var}\{f(\boldsymbol{x})\}$
- Correlation function: $r(\boldsymbol{x}, \boldsymbol{x}') = r(\boldsymbol{x} - \boldsymbol{x}')$

When no basis functions are used (i.e., $p(\boldsymbol{x}) = 1$), UK reduces to ordinary kriging (OK)

Gaussian process (kriging) model

Given **data** $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, the posterior of $f(\mathbf{x})$ is **Gaussian**:

$$[f(\mathbf{x})|\mathcal{D}] \sim \mathcal{N}\left(\hat{f}_n(\mathbf{x}), \sigma^2 s^2(\mathbf{x})\right)$$



Gaussian process (kriging) model

Given **data** $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, the posterior of $f(x)$ is **Gaussian**:

$$[f(x)|\mathcal{D}] \sim \mathcal{N}\left(\hat{f}_n(x), \sigma^2 s^2(x)\right)$$

- Posterior **mean**: $\hat{f}_n(x) = \mathbf{p}^\top(x)\hat{\boldsymbol{\beta}} + \mathbf{r}^\top(x)\mathbf{R}^{-1}(\mathbf{y} - \mathbf{P}\hat{\boldsymbol{\beta}})$
 - $\hat{\boldsymbol{\beta}} = (\mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P})^{-1} \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{y}$
 - $\mathbf{y} = [y_1, \dots, y_n]^\top$, $\mathbf{P} = [p(x_1), \dots, p(x_n)]^\top$
 - $\mathbf{R} = (r(x_i, x_j))_{n \times n}$, $\mathbf{r}(x) = [r(x, x_1), \dots, r(x, x_n)]^\top$
- Posterior **variance**: $\sigma^2 s^2(x)$
 - σ^2 : **prior variance**
 - $s^2(x) = r(x, x) - \mathbf{r}^\top(x)\mathbf{R}^{-1}\mathbf{r}(x) + \mathbf{p}^\top(x)\mathbf{G}^{-1}\mathbf{p}(x)$,
where $\mathbf{G} = \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P}$

Gaussian process (kriging) model

Given **data** $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, the posterior of $f(x)$ is **Gaussian**:

$$[f(x)|\mathcal{D}] \sim \mathcal{N}\left(\hat{f}_n(x), \sigma^2 s^2(x)\right)$$

- Posterior **mean**: $\hat{f}_n(x) = \mathbf{p}^\top(x)\hat{\boldsymbol{\beta}} + \mathbf{r}^\top(x)\mathbf{R}^{-1}(\mathbf{y} - \mathbf{P}\hat{\boldsymbol{\beta}})$
 - $\hat{\boldsymbol{\beta}} = (\mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P})^{-1} \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{y}$
 - $\mathbf{y} = [y_1, \dots, y_n]^\top$, $\mathbf{P} = [p(x_1), \dots, p(x_n)]^\top$
 - $\mathbf{R} = (r(x_i, x_j))_{n \times n}$, $\mathbf{r}(x) = [r(x, x_1), \dots, r(x, x_n)]^\top$
- Posterior **variance**: $\sigma^2 s^2(x)$
 - σ^2 : **prior variance**
 - $s^2(x) = r(x, x) - \mathbf{r}^\top(x)\mathbf{R}^{-1}\mathbf{r}(x) + \mathbf{p}^\top(x)\mathbf{G}^{-1}\mathbf{p}(x)$,
where $\mathbf{G} = \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P}$

Gaussian process (kriging) model

Given **data** $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, the posterior of $f(x)$ is **Gaussian**:

$$[f(x)|\mathcal{D}] \sim \mathcal{N}\left(\hat{f}_n(x), \sigma^2 s^2(x)\right)$$

- Posterior **mean**: $\hat{f}_n(x) = \mathbf{p}^\top(x)\hat{\boldsymbol{\beta}} + \mathbf{r}^\top(x)\mathbf{R}^{-1}(\mathbf{y} - \mathbf{P}\hat{\boldsymbol{\beta}})$
 - $\hat{\boldsymbol{\beta}} = (\mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P})^{-1} \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{y}$
 - $\mathbf{y} = [y_1, \dots, y_n]^\top$, $\mathbf{P} = [p(x_1), \dots, p(x_n)]^\top$
 - $\mathbf{R} = (r(x_i, x_j))_{n \times n}$, $\mathbf{r}(x) = [r(x, x_1), \dots, r(x, x_n)]^\top$
- Posterior **variance**: $\sigma^2 s^2(x)$
 - σ^2 : **prior variance**
 - $s^2(x) = r(x, x) - \mathbf{r}^\top(x)\mathbf{R}^{-1}\mathbf{r}(x) + \mathbf{p}^\top(x)\mathbf{G}^{-1}\mathbf{p}(x)$,
where $\mathbf{G} = \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P}$

Gaussian process (kriging) model

Given **data** $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, the posterior of $f(x)$ is **Gaussian**:

$$[f(x)|\mathcal{D}] \sim \mathcal{N}\left(\hat{f}_n(x), \sigma^2 s^2(x)\right)$$

- Posterior **mean**: $\hat{f}_n(x) = \mathbf{p}^\top(x)\hat{\boldsymbol{\beta}} + \mathbf{r}^\top(x)\mathbf{R}^{-1}(\mathbf{y} - \mathbf{P}\hat{\boldsymbol{\beta}})$
 - $\hat{\boldsymbol{\beta}} = (\mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P})^{-1} \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{y}$
 - $\mathbf{y} = [y_1, \dots, y_n]^\top$, $\mathbf{P} = [\mathbf{p}(x_1), \dots, \mathbf{p}(x_n)]^\top$
 - $\mathbf{R} = (r(x_i, x_j))_{n \times n}$, $\mathbf{r}(x) = [r(x, x_1), \dots, r(x, x_n)]^\top$
- Posterior **variance**: $\sigma^2 s^2(x)$
 - σ^2 : **prior** variance
 - $s^2(x) = r(x, x) - \mathbf{r}^\top(x)\mathbf{R}^{-1}\mathbf{r}(x) + \mathbf{p}^\top(x)\mathbf{G}^{-1}\mathbf{p}(x)$,
where $\mathbf{G} = \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P}$

Gaussian process (kriging) model

Given **data** $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$, the posterior of $f(\mathbf{x})$ is **Gaussian**:

$$[f(\mathbf{x}) | \mathcal{D}] \sim \mathcal{N}\left(\hat{f}_n(\mathbf{x}), \sigma^2 s^2(\mathbf{x})\right)$$

- Posterior **mean**: $\hat{f}_n(\mathbf{x}) = \mathbf{p}^\top(\mathbf{x})\hat{\boldsymbol{\beta}} + \mathbf{r}^\top(\mathbf{x})\mathbf{R}^{-1}(\mathbf{y} - \mathbf{P}\hat{\boldsymbol{\beta}})$
 - $\hat{\boldsymbol{\beta}} = (\mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P})^{-1} \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{y}$
 - $\mathbf{y} = [y_1, \dots, y_n]^\top$, $\mathbf{P} = [\mathbf{p}(\mathbf{x}_1), \dots, \mathbf{p}(\mathbf{x}_n)]^\top$
 - $\mathbf{R} = (r(\mathbf{x}_i, \mathbf{x}_j))_{n \times n}$, $\mathbf{r}(\mathbf{x}) = [r(\mathbf{x}, \mathbf{x}_1), \dots, r(\mathbf{x}, \mathbf{x}_n)]^\top$
- Posterior **variance**: $\sigma^2 s^2(\mathbf{x})$
 - σ^2 : **prior** variance
 - $s^2(\mathbf{x}) = r(\mathbf{x}, \mathbf{x}) - \mathbf{r}^\top(\mathbf{x})\mathbf{R}^{-1}\mathbf{r}(\mathbf{x}) + \mathbf{p}^\top(\mathbf{x})\mathbf{G}^{-1}\mathbf{p}(\mathbf{x})$,
where $\mathbf{G} = \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P}$

Gaussian process (kriging) model

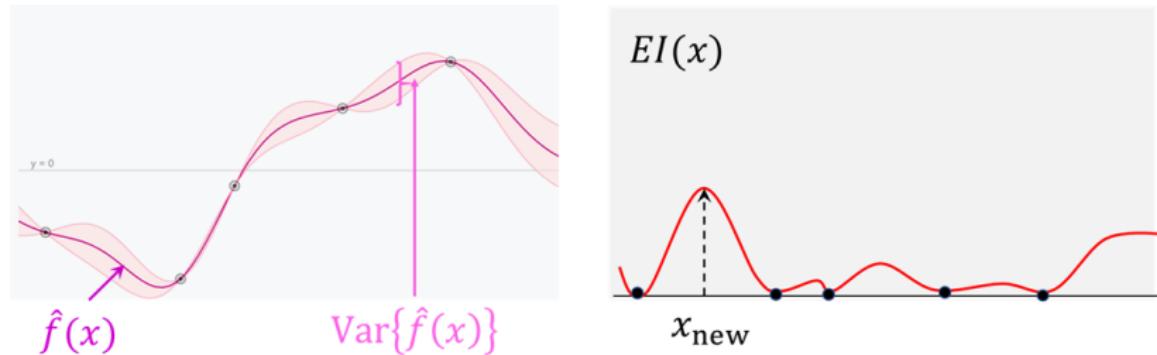
Given **data** $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, the posterior of $f(x)$ is **Gaussian**:

$$[f(x)|\mathcal{D}] \sim \mathcal{N}\left(\hat{f}_n(x), \sigma^2 s^2(x)\right)$$

- Posterior **mean**: $\hat{f}_n(x) = \mathbf{p}^\top(x)\hat{\beta} + \mathbf{r}^\top(x)\mathbf{R}^{-1}(\mathbf{y} - \mathbf{P}\hat{\beta})$
 - $\hat{\beta} = (\mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P})^{-1} \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{y}$
 - $\mathbf{y} = [y_1, \dots, y_n]^\top$, $\mathbf{P} = [\mathbf{p}(x_1), \dots, \mathbf{p}(x_n)]^\top$
 - $\mathbf{R} = (r(x_i, x_j))_{n \times n}$, $\mathbf{r}(x) = [r(x, x_1), \dots, r(x, x_n)]^\top$
- Posterior **variance**: $\sigma^2 s^2(x)$
 - σ^2 : **prior** variance
 - $s^2(x) = r(x, x) - \mathbf{r}^\top(x)\mathbf{R}^{-1}\mathbf{r}(x) + \mathbf{p}^\top(x)\mathbf{G}^{-1}\mathbf{p}(x)$,
where $\mathbf{G} = \mathbf{P}^\top \mathbf{R}^{-1} \mathbf{P}$

Expected Improvement

How to incorporate this **probabilistic** predictor within the **acquisition function?**

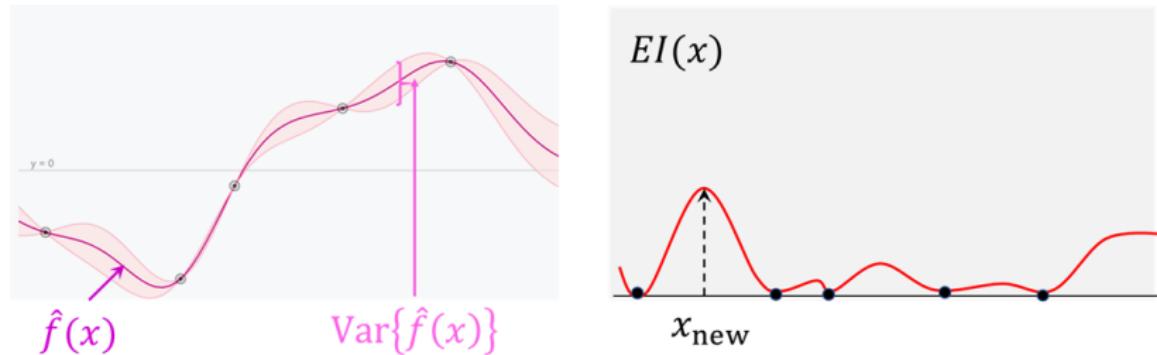


Idea: want to choose a new point x which yields greatest improvement over the current minimum $y^* = \min_{i=1}^n y_i$:

$$\max(y^* - f(x), 0)$$

Expected Improvement

How to incorporate this **probabilistic** predictor within the **acquisition function?**



Idea: want to choose a new point x which yields greatest **improvement** over the current **minimum** $y^* = \min_{i=1}^n y_i$:

$$\max(y^* - f(x), 0)$$

Expected Improvement

Expected improvement (Jones et al. 1997)

$$\begin{aligned}\arg \max_{x \in \Omega} \text{EI}(x) &:= \mathbb{E}[\max(y^* - f(x), 0) | \mathcal{D}] \\ &= I(x)\Phi\left(\frac{I(x)}{\sigma s_n(x)}\right) + \sigma s_n(x)\phi\left(\frac{I(x)}{\sigma s_n(x)}\right).\end{aligned}$$

- $y^* = \min_{i=1}^n y_i$: current **best** value
- $I(x) = y^* - \hat{f}_n(x)$: difference between y^* and predictor $\hat{f}_n(x)$
- ϕ, Φ : p.d.f. and c.d.f. of a standard **normal** distribution
- $\sigma s_n(x)$: posterior standard deviation of $f(x)$

Expected Improvement

Expected improvement (Jones et al. 1997)

$$\begin{aligned}\arg \max_{\boldsymbol{x} \in \Omega} \text{EI}(\boldsymbol{x}) &:= \mathbb{E}[\max(y^* - f(\boldsymbol{x}), 0) | \mathcal{D}] \\ &= I(\boldsymbol{x})\Phi\left(\frac{I(\boldsymbol{x})}{\sigma s_n(\boldsymbol{x})}\right) + \sigma s_n(\boldsymbol{x})\phi\left(\frac{I(\boldsymbol{x})}{\sigma s_n(\boldsymbol{x})}\right).\end{aligned}$$

- $y^* = \min_{i=1}^n y_i$: current **best** value
- $I(\boldsymbol{x}) = y^* - \hat{f}_n(\boldsymbol{x})$: **difference** between y^* and predictor $\hat{f}_n(\boldsymbol{x})$
- ϕ, Φ : p.d.f. and c.d.f. of a standard **normal** distribution
- $\sigma s_n(\boldsymbol{x})$: **posterior standard deviation** of $f(\boldsymbol{x})$

Expected Improvement

Expected improvement (Jones et al. 1997)

$$\begin{aligned}\arg \max_{\boldsymbol{x} \in \Omega} \text{EI}(\boldsymbol{x}) &:= \mathbb{E}[\max(y^* - f(\boldsymbol{x}), 0) | \mathcal{D}] \\ &= I(\boldsymbol{x})\Phi\left(\frac{I(\boldsymbol{x})}{\sigma s_n(\boldsymbol{x})}\right) + \sigma s_n(\boldsymbol{x})\phi\left(\frac{I(\boldsymbol{x})}{\sigma s_n(\boldsymbol{x})}\right).\end{aligned}$$

- $y^* = \min_{i=1}^n y_i$: current **best** value
- $I(\boldsymbol{x}) = y^* - \hat{f}_n(\boldsymbol{x})$: **difference** between y^* and predictor $\hat{f}_n(\boldsymbol{x})$
- ϕ, Φ : p.d.f. and c.d.f. of a standard **normal** distribution
- $\sigma s_n(\boldsymbol{x})$: **posterior standard deviation** of $f(\boldsymbol{x})$

Expected Improvement

EI has an intuitive **interpretation**:

$$\arg \max_{\boldsymbol{x} \in \Omega} \text{EI}(\boldsymbol{x}) := \underbrace{I(\boldsymbol{x}) \Phi \left(\frac{I(\boldsymbol{x})}{\sigma s_n(\boldsymbol{x})} \right)}_{\text{Exploitation}} + \underbrace{\sigma s_n(\boldsymbol{x}) \phi \left(\frac{I(\boldsymbol{x})}{\sigma s_n(\boldsymbol{x})} \right)}_{\text{Exploration}}.$$

- $I(\boldsymbol{x}) = y^* - \hat{f}_n(\boldsymbol{x})$: **difference** between y^* and predictor $\hat{f}_n(\boldsymbol{x})$
- $\sigma s_n(\boldsymbol{x})$: **posterior standard deviation** of $f(\boldsymbol{x})$
- **Exploitation**: favor new points \boldsymbol{x} with **small** predicted objective $\hat{f}_n(\boldsymbol{x})$ – “exploit” the fitted model to find **better** solutions
- **Exploration**: favor new points \boldsymbol{x} with **large** variance – “explore” uncertain regions to find **better** solutions
- EI balances the **exploitation-exploration trade-off** (Kearns & Singh 2002)

Expected Improvement

EI has an intuitive **interpretation**:

$$\arg \max_{\mathbf{x} \in \Omega} \text{EI}(\mathbf{x}) := \underbrace{I(\mathbf{x}) \Phi \left(\frac{I(\mathbf{x})}{\sigma s_n(\mathbf{x})} \right)}_{\text{Exploitation}} + \underbrace{\sigma s_n(\mathbf{x}) \phi \left(\frac{I(\mathbf{x})}{\sigma s_n(\mathbf{x})} \right)}_{\text{Exploration}}.$$

- $I(\mathbf{x}) = y^* - \hat{f}_n(\mathbf{x})$: **difference** between y^* and predictor $\hat{f}_n(\mathbf{x})$
- $\sigma s_n(\mathbf{x})$: **posterior standard deviation** of $f(\mathbf{x})$
- **Exploitation**: favor new points \mathbf{x} with **small** predicted objective $\hat{f}_n(\mathbf{x})$ – “exploit” the fitted model to find **better** solutions
- **Exploration**: favor new points \mathbf{x} with **large** variance – “explore” uncertain regions to find **better** solutions
- EI balances the **exploitation-exploration trade-off** (Kearns & Singh 2002)

Expected Improvement

EI has an intuitive **interpretation**:

$$\arg \max_{\mathbf{x} \in \Omega} EI(\mathbf{x}) := \underbrace{I(\mathbf{x})\Phi\left(\frac{I(\mathbf{x})}{\sigma s_n(\mathbf{x})}\right)}_{\text{Exploitation}} + \underbrace{\sigma s_n(\mathbf{x})\phi\left(\frac{I(\mathbf{x})}{\sigma s_n(\mathbf{x})}\right)}_{\text{Exploration}}.$$

- $I(\mathbf{x}) = y^* - \hat{f}_n(\mathbf{x})$: **difference** between y^* and predictor $\hat{f}_n(\mathbf{x})$
- $\sigma s_n(\mathbf{x})$: **posterior standard deviation** of $f(\mathbf{x})$
- **Exploitation**: favor new points \mathbf{x} with **small** predicted objective $\hat{f}_n(\mathbf{x})$ – “**exploit**” the fitted model to find **better** solutions
- **Exploration**: favor new points \mathbf{x} with **large** variance – “**explore**” uncertain regions to find **better** solutions
- EI balances the **exploitation-exploration trade-off** (Kearns & Singh 2002)

Expected Improvement

EI has an intuitive **interpretation**:

$$\arg \max_{\mathbf{x} \in \Omega} EI(\mathbf{x}) := I(\mathbf{x}) \underbrace{\Phi\left(\frac{I(\mathbf{x})}{\sigma s_n(\mathbf{x})}\right)}_{\text{Exploitation}} + \sigma s_n(\mathbf{x}) \underbrace{\phi\left(\frac{I(\mathbf{x})}{\sigma s_n(\mathbf{x})}\right)}_{\text{Exploration}}.$$

- $I(\mathbf{x}) = y^* - \hat{f}_n(\mathbf{x})$: **difference** between y^* and predictor $\hat{f}_n(\mathbf{x})$
- $\sigma s_n(\mathbf{x})$: **posterior standard deviation** of $f(\mathbf{x})$
- **Exploitation**: favor new points \mathbf{x} with **small** predicted objective $\hat{f}_n(\mathbf{x})$ – “exploit” the fitted model to find **better** solutions
- **Exploration**: favor new points \mathbf{x} with **large** variance – “explore” uncertain regions to find **better** solutions
- EI balances the **exploitation-exploration trade-off** (Kearns & Singh 2002)

Plug-in estimation

In practice, the GP **variance** parameter σ^2 is **unknown** and needs to be **estimated**:

- Estimate σ^2 via MLE

$$\hat{\sigma}^2 = \frac{1}{n}(y - P^\top \hat{\beta})^\top R^{-1}(y - P^\top \hat{\beta})$$

- Plug-in $\hat{\sigma}^2$ into acquisition function $El(x)$
- Optimize $El(x)$ for new query point x_{new}

Sounds reasonable... what can go wrong?

Plug-in estimation

In practice, the GP **variance** parameter σ^2 is **unknown** and needs to be **estimated**:

- **Estimate** σ^2 via MLE

$$\hat{\sigma}^2 = \frac{1}{n}(\mathbf{y} - \mathbf{P}^\top \hat{\boldsymbol{\beta}})^\top \mathbf{R}^{-1}(\mathbf{y} - \mathbf{P}^\top \hat{\boldsymbol{\beta}})$$

- **Plug-in** $\hat{\sigma}^2$ into acquisition function $\text{EI}(\mathbf{x})$
- **Optimize** $\text{EI}(\mathbf{x})$ for new query point \mathbf{x}_{new}

Sounds **reasonable**... what can go **wrong**?

Plug-in estimation

In practice, the GP **variance** parameter σ^2 is **unknown** and needs to be **estimated**:

- **Estimate** σ^2 via MLE

$$\hat{\sigma}^2 = \frac{1}{n}(\mathbf{y} - \mathbf{P}^\top \hat{\boldsymbol{\beta}})^\top \mathbf{R}^{-1}(\mathbf{y} - \mathbf{P}^\top \hat{\boldsymbol{\beta}})$$

- **Plug-in** $\hat{\sigma}^2$ into acquisition function $\text{EI}(\mathbf{x})$
- **Optimize** $\text{EI}(\mathbf{x})$ for new query point \mathbf{x}_{new}

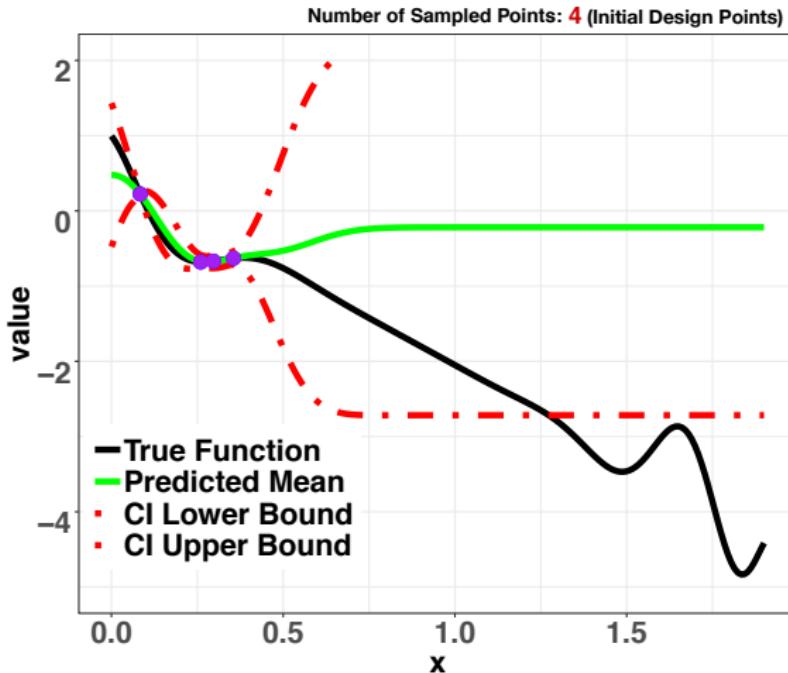
Sounds **reasonable**... what can go **wrong**?

Murphy's Law

Anything that can go wrong...
will go wrong.

Over-exploitation

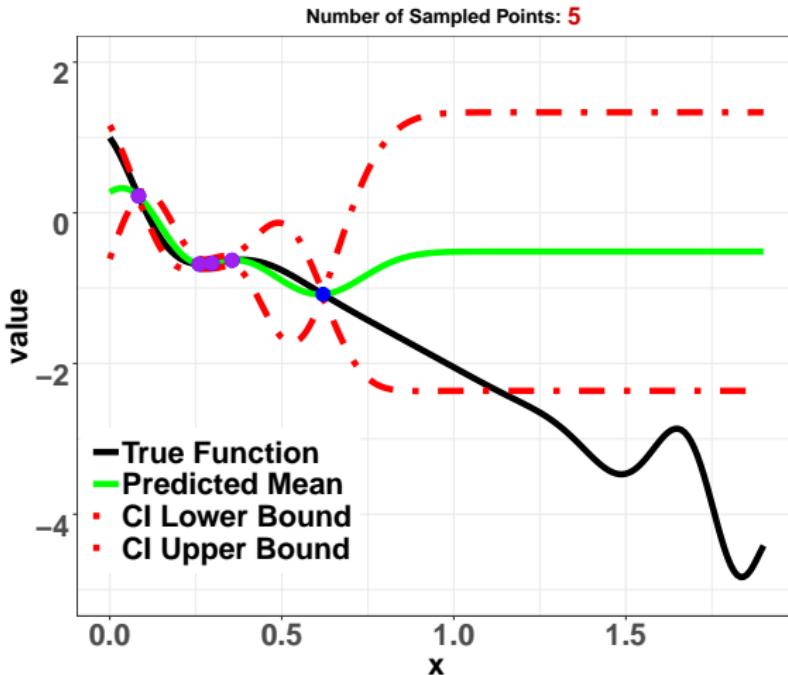
Let's try this **plug-in EI** approach:



... plug-in EI gets **stuck** in local minima!

Over-exploitation

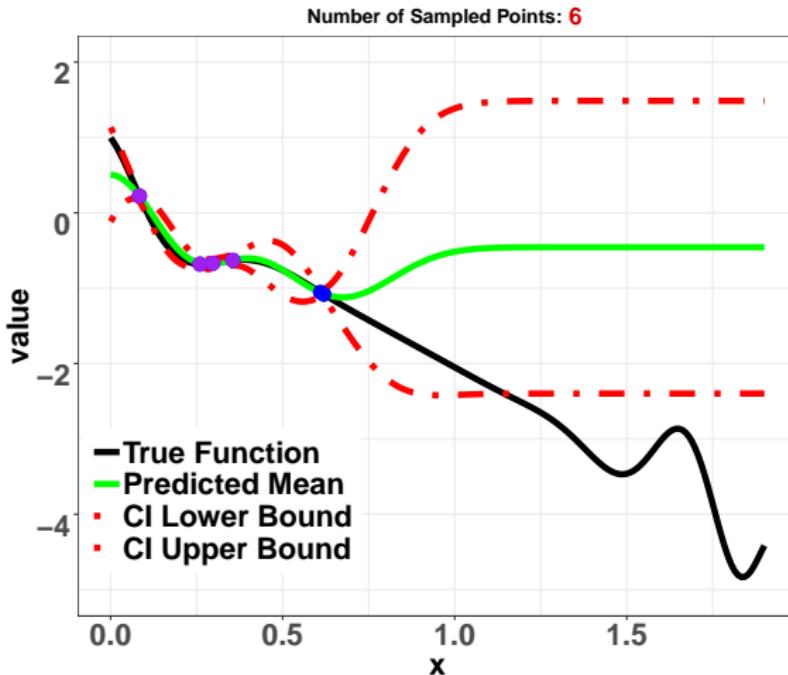
Let's try this **plug-in EI** approach:



... plug-in EI gets **stuck** in local minima!

Over-exploitation

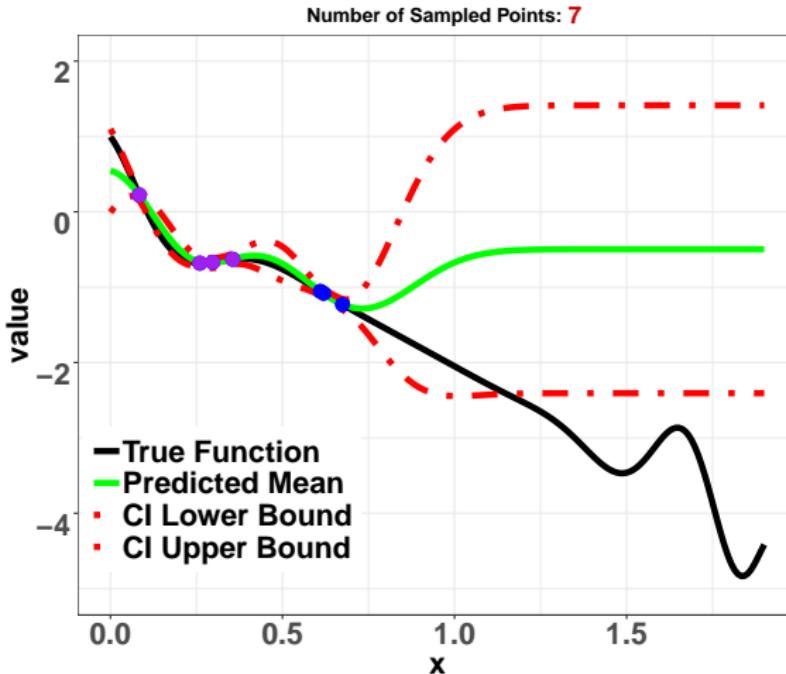
Let's try this **plug-in EI** approach:



... plug-in EI gets **stuck** in local minima!

Over-exploitation

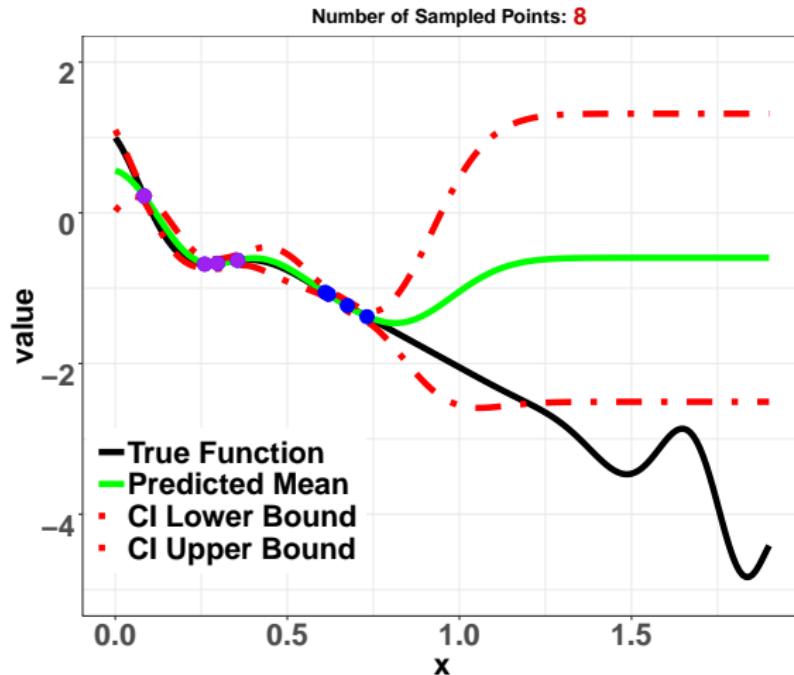
Let's try this **plug-in EI** approach:



... plug-in EI gets **stuck** in local minima!

Over-exploitation

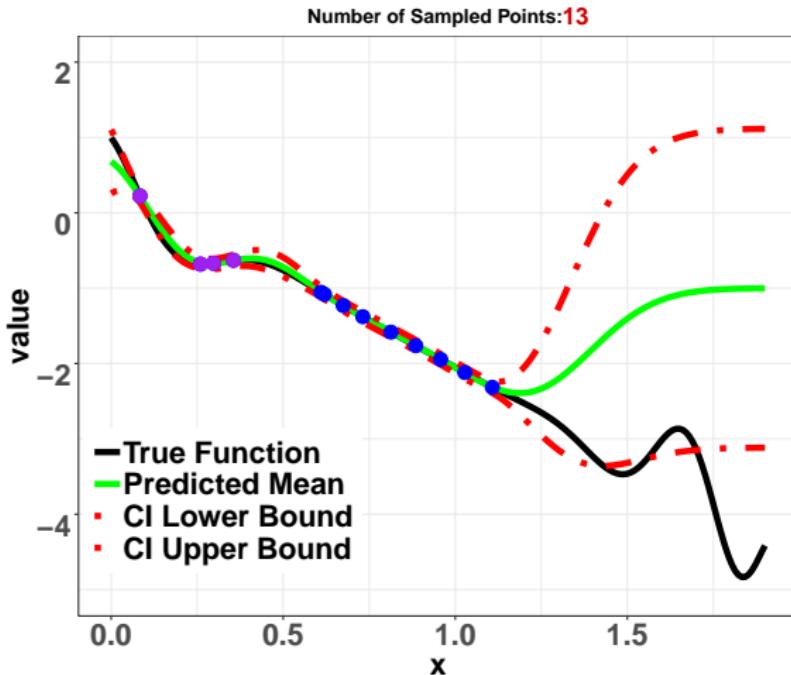
Let's try this **plug-in EI** approach:



... plug-in EI gets **stuck** in local minima!

Over-exploitation

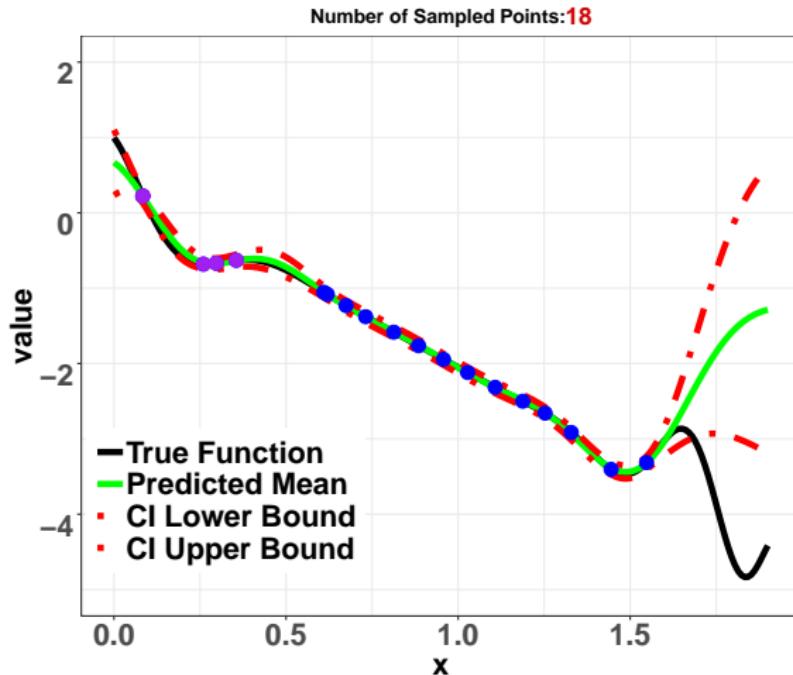
Let's try this **plug-in EI** approach:



... plug-in EI gets **stuck** in local minima!

Over-exploitation

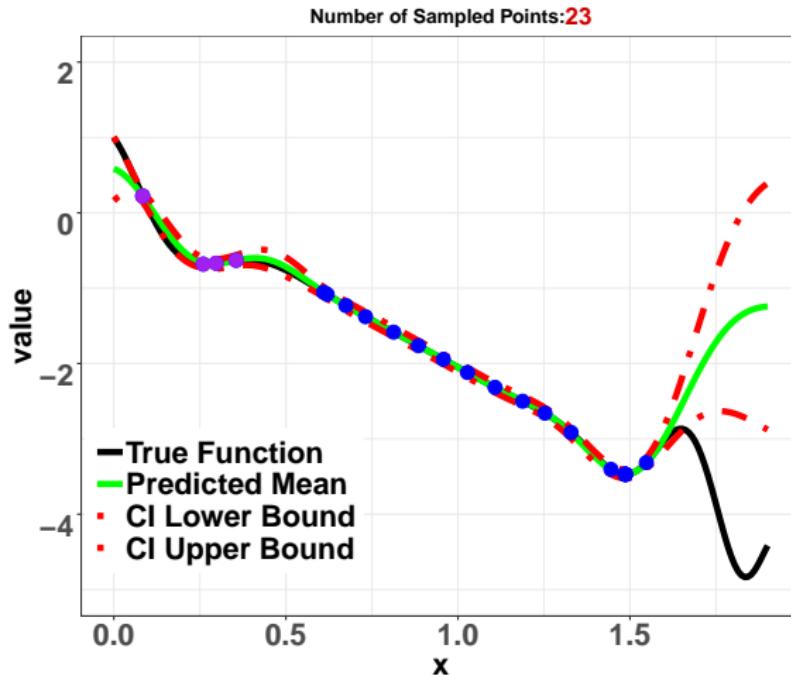
Let's try this **plug-in EI** approach:



... plug-in EI gets **stuck** in local minima!

Over-exploitation

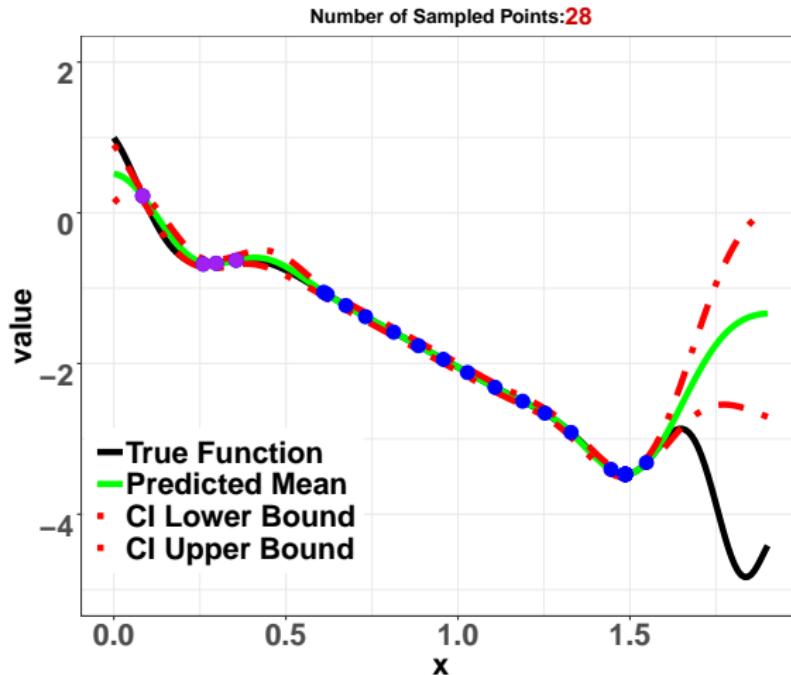
Let's try this **plug-in EI** approach:



... plug-in EI gets **stuck** in local minima!

Over-exploitation

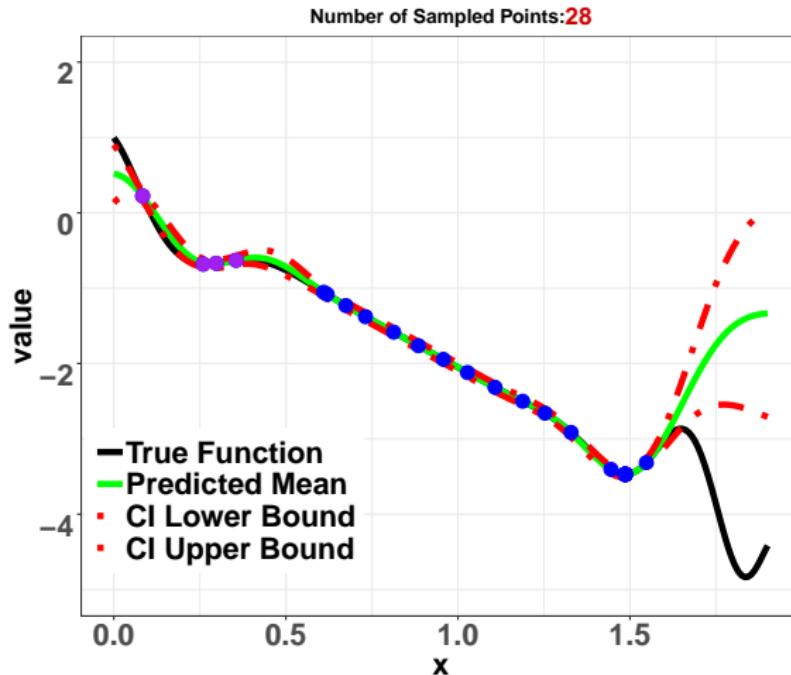
Let's try this **plug-in EI** approach:



... plug-in EI gets **stuck** in local minima!

Over-exploitation

Let's try this **plug-in EI** approach:



... plug-in EI gets **stuck** in local minima!

Over-exploitation

This **over-exploitation** of the **fitted GP model** for (plug-in) EI has been noted in recent literature:

- Qin et al. (2017): EI is too **greedy**, too much **exploitation** near the current best point
- Bull (2011): EI can **fail** to find the global optima!

Theorem (Bull, 2011)

Given any $\epsilon > 0$ and any random initial design F , there exist some $\delta > 0$ and some f in a certain space,

$$\mathbb{P}_F(\min_n f(x_n) - \min_{x \in \Omega} f(x) > \delta) \geq 1 - \epsilon,$$

where $(x_n)_{n=1}^{\infty}$ is a solution sequence maximizing EI with plug-in MLEs $(\hat{\beta}, \hat{\sigma}^2)$.

Why? We're not factoring in uncertainty on the plug-in MLEs!

Over-exploitation

This **over-exploitation** of the **fitted GP model** for (plug-in) EI has been noted in recent literature:

- Qin et al. (2017): EI is too **greedy**, too much **exploitation** near the current best point
- Bull (2011): EI can **fail** to find the global optima!

Theorem (Bull, 2011)

Given any $\epsilon > 0$ and any random initial design F , there exist some $\delta > 0$ and some f in a certain space,

$$\mathbb{P}_F(\min_n f(\mathbf{x}_n) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) > \delta) \geq 1 - \epsilon,$$

where $(\mathbf{x}_n)_{n=1}^{\infty}$ is a solution sequence maximizing EI with plug-in MLEs $(\hat{\beta}, \hat{\sigma}^2)$.

Why? We're not factoring in **uncertainty** on the **plug-in MLEs**!

Solutions?

- Bull (2011): EI with **enlarged variance** and ϵ -**greedy** EI (ϵ -EI):
 - **Enlarged variance:** use $n\hat{\sigma}_n^2$ instead of MLE $\hat{\sigma}_n^2$
 - **ϵ -greedy:** with probability ϵ , **uniformly** select next point; with probability $1 - \epsilon$, use **EI** with enlarged variance
 - Corrects **over-exploitation** by forcing EI to **explore** more
- Srinivas et al. (2010): Upper confidence bound (UCB) on GPs

$$x_{\text{new}} = \arg \max_x \left\{ \hat{f}_n(x) + \rho \hat{\sigma} s_n(x) \right\}, \quad \rho: \text{exploration parameter}$$

These methods encourage **more exploration**, but in *ad hoc* ways!
Can we provide a **principled** solution to **over-exploitation** via a **hierarchical Bayesian** view on EI?

Solutions?

- Bull (2011): EI with **enlarged variance** and ϵ -**greedy** EI (ϵ -EI):
 - **Enlarged variance:** use $n\hat{\sigma}_n^2$ instead of MLE $\hat{\sigma}_n^2$
 - **ϵ -greedy:** with probability ϵ , **uniformly** select next point;
with probability $1 - \epsilon$, use **EI** with enlarged variance
 - Corrects **over-exploitation** by forcing EI to **explore** more
- Srinivas et al. (2010): Upper confidence bound (UCB) on GPs

$$\mathbf{x}_{\text{new}} = \arg \max_{\mathbf{x}} \left\{ \hat{f}_n(\mathbf{x}) + \rho \hat{\sigma} s_n(\mathbf{x}) \right\}, \quad \rho: \text{exploration parameter}$$

These methods encourage **more exploration**, but in *ad hoc* ways!
Can we provide a **principled** solution to **over-exploitation** via a
hierarchical Bayesian view on EI?

Hierarchical GP Model

Hierarchical GP (kriging) model ([Handcock & Stein 1993](#)):

$$\textbf{UK} : f(\cdot) = \mu(\cdot) + Z(\cdot), \quad Z(\cdot) \sim GP\{0, \sigma^2 r(\cdot, \cdot)\}$$

$$\textbf{Mean} : \mu(x) = p^\top(x)\beta, \quad \textbf{Prior: } [\beta] \propto \mathbf{1}$$

$$\textbf{Covariance} : \sigma^2 r(\cdot, \cdot), \quad \textbf{Prior: } [\sigma^2] \sim \text{Inv-Gamma}(a, b)$$

Still a **closed-form posterior** distribution ([Santner et al., 2018](#)):

$$[f(x)|\mathcal{D}] \sim T(2\tilde{a}, \hat{f}_n(x), \tilde{\sigma}^2 s^2(x)).$$

- $T(2\tilde{a}, \cdot, \cdot)$: **Student** t -distribution with d.f. $2\tilde{a}$
- $\tilde{a} = a + (n - q)/2$, $\tilde{b} = b + n\hat{\sigma}^2/2$, $\tilde{\sigma}^2 = \tilde{b}/\tilde{a}$

Hierarchical GP Model

Hierarchical GP (kriging) model (Handcock & Stein 1993):

$$\text{UK} : f(\cdot) = \mu(\cdot) + Z(\cdot), \quad Z(\cdot) \sim GP\{0, \sigma^2 r(\cdot, \cdot)\}$$

$$\text{Mean} : \mu(\mathbf{x}) = \mathbf{p}^\top(\mathbf{x})\boldsymbol{\beta}, \quad \text{Prior: } [\boldsymbol{\beta}] \propto \mathbf{1}$$

$$\text{Covariance} : \sigma^2 r(\cdot, \cdot), \quad \text{Prior: } [\sigma^2] \sim \text{Inv-Gamma}(a, b)$$

Still a **closed-form posterior** distribution (Santner et al., 2018):

$$[f(\mathbf{x}) | \mathcal{D}] \sim T(2\tilde{a}, \hat{f}_n(\mathbf{x}), \tilde{\sigma}^2 s^2(\mathbf{x})).$$

- $T(2\tilde{a}, \cdot, \cdot)$: **Student** t -distribution with d.f. $2\tilde{a}$
- $\tilde{a} = a + (n - q)/2$, $\tilde{b} = b + n\hat{\sigma}^2/2$, $\tilde{\sigma}^2 = \tilde{b}/\tilde{a}$

Hierarchical Expected Improvement

Do we get a **closed-form** expression for **expected improvement** under this **hierarchical** GP model?

$$\text{HEI}(x) = \mathbb{E}[\max(y^* - f(x), 0) | \mathcal{D}]$$

$$= I(x)\Phi_{2\tilde{a}}\left(\frac{I(x)}{\tilde{\sigma}s_n(x)}\right) + \kappa\tilde{\sigma}s_n(x)\phi_{2\tilde{a}-2}\left(\frac{I(x)}{\kappa\tilde{\sigma}s_n(x)}\right)$$

- ϕ_ν , Φ_ν : p.d.f. and c.d.f. of a ***t-distribution*** with d.f. ν

Why is this **important**?

- Function **queries** are obtained via **maximizing** $\text{HEI}(x)$

$$x_{\text{new}} = \arg \max_{x \in \Omega} \text{HEI}(x)$$

- **Closed-form** $\text{HEI}(x) \Rightarrow$ **efficient** queries, **improved** optimization

Hierarchical Expected Improvement

Do we get a **closed-form** expression for **expected improvement** under this **hierarchical** GP model?

$$\text{HEI}(\mathbf{x}) = \mathbb{E}[\max(y^* - f(\mathbf{x}), 0) | \mathcal{D}]$$

$$= I(\mathbf{x})\Phi_{2\tilde{a}}\left(\frac{I(\mathbf{x})}{\tilde{\sigma}s_n(\mathbf{x})}\right) + \kappa\tilde{\sigma}s_n(\mathbf{x})\phi_{2\tilde{a}-2}\left(\frac{I(\mathbf{x})}{\kappa\tilde{\sigma}s_n(\mathbf{x})}\right)$$

- ϕ_ν , Φ_ν : p.d.f. and c.d.f. of a ***t-distribution*** with d.f. ν

Why is this **important**?

- Function queries are obtained via **maximizing** $\text{HEI}(\mathbf{x})$

$$x_{\text{new}} = \arg \max_{x \in \Omega} \text{HEI}(x)$$

- **Closed-form** $\text{HEI}(x) \Rightarrow$ **efficient** queries, **improved** optimization

Hierarchical Expected Improvement

Do we get a **closed-form** expression for **expected improvement** under this **hierarchical** GP model?

$$\text{HEI}(\mathbf{x}) = \mathbb{E}[\max(y^* - f(\mathbf{x}), 0) | \mathcal{D}]$$

$$= I(\mathbf{x})\Phi_{2\tilde{a}}\left(\frac{I(\mathbf{x})}{\tilde{\sigma}s_n(\mathbf{x})}\right) + \kappa\tilde{\sigma}s_n(\mathbf{x})\phi_{2\tilde{a}-2}\left(\frac{I(\mathbf{x})}{\kappa\tilde{\sigma}s_n(\mathbf{x})}\right)$$

- ϕ_ν , Φ_ν : p.d.f. and c.d.f. of a ***t-distribution*** with d.f. ν

Why is this **important**?

- Function **queries** are obtained via **maximizing** $\text{HEI}(\mathbf{x})$

$$\mathbf{x}_{\text{new}} = \arg \max_{\mathbf{x} \in \Omega} \text{HEI}(\mathbf{x})$$

- **Closed-form** $\text{HEI}(\mathbf{x}) \Rightarrow$ **efficient** queries, **improved** optimization

HEI vs. EI

How is HEI **different** from EI?

$$\text{HEI}(\mathbf{x}) = I(\mathbf{x})\Phi_{2\tilde{a}}\left(\frac{I(\mathbf{x})}{\tilde{\sigma}s_n(\mathbf{x})}\right) + \kappa\tilde{\sigma}s_n(\mathbf{x})\phi_{2\tilde{a}-2}\left(\frac{I(\mathbf{x})}{\kappa\tilde{\sigma}s_n(\mathbf{x})}\right),$$

vs.

$$\text{EI}(\mathbf{x}) = I(\mathbf{x})\Phi\left(\frac{I(\mathbf{x})}{\hat{\sigma}s_n(\mathbf{x})}\right) + \hat{\sigma}s_n(\mathbf{x})\phi\left(\frac{I(\mathbf{x})}{\hat{\sigma}s_n(\mathbf{x})}\right).$$

HEI encourages **more exploration**:

- **Finite-sample** correction $\kappa = \sqrt{\frac{\tilde{a}}{a-1}} > 1$
- $\tilde{\sigma} > \hat{\sigma}$ for weak hyperpriors on σ^2
- t distribution has **heavier** tails

HEI vs. EI

How is HEI **different** from EI?

$$\text{HEI}(\mathbf{x}) = I(\mathbf{x}) \Phi_{2\tilde{a}} \left(\frac{I(\mathbf{x})}{\tilde{\sigma}s_n(\mathbf{x})} \right) + \kappa \tilde{\sigma} s_n(\mathbf{x}) \phi_{2\tilde{a}-2} \left(\frac{I(\mathbf{x})}{\kappa \tilde{\sigma} s_n(\mathbf{x})} \right),$$

vs.

$$\text{EI}(\mathbf{x}) = I(\mathbf{x}) \Phi \left(\frac{I(\mathbf{x})}{\hat{\sigma}s_n(\mathbf{x})} \right) + \hat{\sigma} s_n(\mathbf{x}) \phi \left(\frac{I(\mathbf{x})}{\hat{\sigma}s_n(\mathbf{x})} \right).$$

HEI encourages **more exploration**:

- Finite-sample correction $\kappa = \sqrt{\frac{\tilde{a}}{a-1}} > 1$
- $\tilde{\sigma} > \hat{\sigma}$ for weak hyperpriors on σ^2
- t distribution has **heavier** tails

HEI vs. EI

How is HEI **different** from EI?

$$\text{HEI}(\mathbf{x}) = I(\mathbf{x})\Phi_{2\tilde{a}}\left(\frac{I(\mathbf{x})}{\tilde{\sigma}s_n(\mathbf{x})}\right) + \kappa\tilde{\sigma}s_n(\mathbf{x})\phi_{2\tilde{a}-2}\left(\frac{I(\mathbf{x})}{\kappa\tilde{\sigma}s_n(\mathbf{x})}\right),$$

vs.

$$\text{EI}(\mathbf{x}) = I(\mathbf{x})\Phi\left(\frac{I(\mathbf{x})}{\hat{\sigma}s_n(\mathbf{x})}\right) + \hat{\sigma}s_n(\mathbf{x})\phi\left(\frac{I(\mathbf{x})}{\hat{\sigma}s_n(\mathbf{x})}\right).$$

HEI encourages **more exploration**:

- **Finite-sample** correction $\kappa = \sqrt{\frac{\tilde{a}}{a-1}} > 1$
- $\tilde{\sigma} > \hat{\sigma}$ for weak hyperpriors on σ^2
- t distribution has **heavier** tails

HEI vs. EI

How is HEI **different** from EI?

$$\text{HEI}(\mathbf{x}) = I(\mathbf{x})\Phi_{2\tilde{a}}\left(\frac{I(\mathbf{x})}{\tilde{\sigma}s_n(\mathbf{x})}\right) + \kappa\tilde{\sigma}s_n(\mathbf{x})\phi_{2\tilde{a}-2}\left(\frac{I(\mathbf{x})}{\kappa\tilde{\sigma}s_n(\mathbf{x})}\right),$$

vs.

$$\text{EI}(\mathbf{x}) = I(\mathbf{x})\Phi\left(\frac{I(\mathbf{x})}{\hat{\sigma}s_n(\mathbf{x})}\right) + \hat{\sigma}s_n(\mathbf{x})\phi\left(\frac{I(\mathbf{x})}{\hat{\sigma}s_n(\mathbf{x})}\right).$$

HEI encourages **more exploration**:

- **Finite-sample** correction $\kappa = \sqrt{\frac{\tilde{a}}{\tilde{a}-1}} > 1$
- $\tilde{\sigma} > \hat{\sigma}$ for weak hyperpriors on σ^2
- t distribution has **heavier** tails

HEI vs. EI

How is HEI **different** from EI?

$$\text{HEI}(\mathbf{x}) = I(\mathbf{x}) \Phi_{2\tilde{a}} \left(\frac{I(\mathbf{x})}{\tilde{\sigma} s_n(\mathbf{x})} \right) + \kappa \tilde{\sigma} s_n(\mathbf{x}) \phi_{2\tilde{a}-2} \left(\frac{I(\mathbf{x})}{\kappa \tilde{\sigma} s_n(\mathbf{x})} \right),$$

vs.

$$\text{EI}(\mathbf{x}) = I(\mathbf{x}) \Phi \left(\frac{I(\mathbf{x})}{\hat{\sigma} s_n(\mathbf{x})} \right) + \hat{\sigma} s_n(\mathbf{x}) \phi \left(\frac{I(\mathbf{x})}{\hat{\sigma} s_n(\mathbf{x})} \right).$$

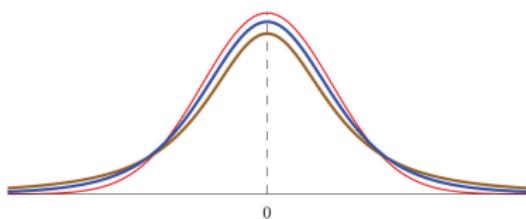
HEI encourages **more exploration**:

Standard normal

t-distribution with $df = 5$

t-distribution with $df = 2$

- **Finite-sample** correction $\kappa = \sqrt{\frac{\tilde{a}}{a-1}} > 1$
- $\tilde{\sigma} > \hat{\sigma}$ for weak hyperpriors on σ^2
- *t* distribution has **heavier** tails



HEI vs. EI

How is HEI **different** from EI?

$$\text{HEI}(\mathbf{x}) = I(\mathbf{x}) \Phi_{2\tilde{a}} \left(\frac{I(\mathbf{x})}{\tilde{\sigma} s_n(\mathbf{x})} \right) + \kappa \tilde{\sigma} s_n(\mathbf{x}) \phi_{2\tilde{a}-2} \left(\frac{I(\mathbf{x})}{\kappa \tilde{\sigma} s_n(\mathbf{x})} \right),$$

vs.

$$\text{EI}(\mathbf{x}) = I(\mathbf{x}) \Phi \left(\frac{I(\mathbf{x})}{\hat{\sigma} s_n(\mathbf{x})} \right) + \hat{\sigma} s_n(\mathbf{x}) \phi \left(\frac{I(\mathbf{x})}{\hat{\sigma} s_n(\mathbf{x})} \right).$$

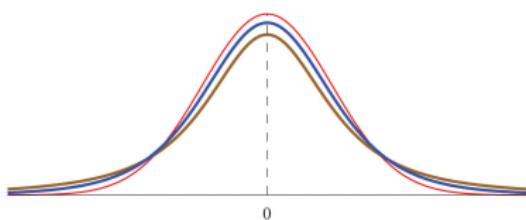
HEI encourages **more exploration**:

Standard normal

t-distribution with $df = 5$

t-distribution with $df = 2$

- **Finite-sample** correction $\kappa = \sqrt{\frac{\tilde{a}}{\tilde{a}-1}} > 1$
- $\tilde{\sigma} > \hat{\sigma}$ for weak hyperpriors on σ^2
- *t* distribution has **heavier** tails



Hyperparameter specification

How to set **hyperparameters** a and b in $[\sigma^2] \sim \text{Inv-Gamma}(a, b)$?

- HEI-Weak: Weakly informative (Gelman et al 2004)

$a = b$ small, e.g., 0.1

- HEI-MMAP: Empirical Bayes (EB, Robbins 1956)

Estimate hyperprior parameters via marginal likelihood

- HEI-DSD: Data-size dependent

From **global convergence theorem** (introduced later), set a constant and $b = \Theta(n) = cn$, where a and c are estimated via EB

Hyperparameter specification

How to set **hyperparameters** a and b in $[\sigma^2] \sim \text{Inv-Gamma}(a, b)$?

- **HEI-Weak: Weakly informative** (Gelman et al 2004)

$a = b$ small, e.g., 0.1

- **HEI-MMAP: Empirical Bayes** (EB, Robbins 1956)

Estimate hyperprior parameters via marginal likelihood

- **HEI-DSD: Data-size dependent**

From **global convergence theorem** (introduced later), set a constant and $b = \Theta(n) = cn$, where a and c are estimated via EB

Hyperparameter specification

How to set **hyperparameters** a and b in $[\sigma^2] \sim \text{Inv-Gamma}(a, b)$?

- **HEI-Weak: Weakly informative** ([Gelman et al 2004](#))

$a = b$ small, e.g., 0.1

- **HEI-MMAP: Empirical Bayes** ([EB, Robbins 1956](#))

Estimate hyperprior parameters via marginal likelihood

- **HEI-DSD: Data-size dependent**

From **global convergence theorem** (introduced later), set a constant and $b = \Theta(n) = cn$, where a and c are estimated via EB

Hyperparameter specification

How to set **hyperparameters** a and b in $[\sigma^2] \sim \text{Inv-Gamma}(a, b)$?

- **HEI-Weak: Weakly informative** ([Gelman et al 2004](#))

$a = b$ small, e.g., 0.1

- **HEI-MMAP: Empirical Bayes** ([EB, Robbins 1956](#))

Estimate hyperprior parameters via marginal likelihood

- **HEI-DSD: Data-size dependent**

From **global convergence theorem** ([introduced later](#)), set a constant and $b = \Theta(n) = cn$, where a and c are estimated via EB

Empirical Bayes

- EB approximates a **fully Bayesian hierarchical model** by estimating hyperparameters via **marginal likelihood maximization** (Carlin & Louis 2000)
 - Here, find (\hat{a}^*, \hat{b}^*) which maximizes

$$(\hat{a}^*, \hat{b}^*) = \arg \max_{a,b} \int p(y|\beta, \sigma^2) \pi(\beta) \pi(\sigma^2|a, b) d\beta d\sigma^2$$

... unfortunately this problem is unbounded!

Empirical Bayes

- EB approximates a **fully Bayesian hierarchical model** by estimating hyperparameters via **marginal likelihood maximization** ([Carlin & Louis 2000](#))
- **Here**, find (a^*, b^*) which **maximizes**:

$$(a^*, b^*) = \arg \max_{a,b} \int p(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}) \pi(\sigma^2|a, b) d\boldsymbol{\beta} d\sigma^2$$

... unfortunately this problem is **unbounded!**

Empirical Bayes

- EB approximates a **fully Bayesian hierarchical model** by estimating hyperparameters via **marginal likelihood maximization** ([Carlin & Louis 2000](#))
- **Here**, find (a^*, b^*) which **maximizes**:

$$(a^*, b^*) = \arg \max_{a,b} \int p(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}) \pi(\sigma^2|a, b) d\boldsymbol{\beta} d\sigma^2$$

... unfortunately this problem is **unbounded!**

Empirical Bayes

- EB approximates a **fully Bayesian hierarchical model** by estimating hyperparameters via **marginal likelihood maximization** ([Carlin & Louis 2000](#))
- Here, find (a^*, b^*) which **maximizes**:

$$(a^*, b^*) = \arg \max_{a,b} \int p(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}) \pi(\sigma^2|a, b) d\boldsymbol{\beta} d\sigma^2$$

... unfortunately this problem is **unbounded!**

Murphy's Law

Anything that can go wrong...
will go wrong.

Empirical Bayes

- Marginal maximum a posteriori (MMAP, Doucet et al. 2002): add another layer of hyperpriors on (a, b) for regularization

Proposition

Assume the following independent hyperpriors:

$$[a] \sim \text{Gamma}(\zeta, \iota), \quad [b] \propto 1.$$

Then the MMAP estimator

$$(a^*, b^*) = \arg \max_{a,b} \int p(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}) \pi(\sigma^2|a, b) d\boldsymbol{\beta} d\sigma^2 \pi(a, b)$$

exists and is finite.

- We can then use the estimated prior $\sigma^2 \sim IG(a^*, b^*)$ with MMAP estimates (a^*, b^*) within HEI
- Key: HEI with EB approximates a **fully Bayesian** EI procedure, without need for **expensive** MCMC (all **closed**-form!)

Empirical Bayes

- Marginal maximum a posteriori (MMAP, Doucet et al. 2002): add another layer of hyperpriors on (a, b) for regularization

Proposition

Assume the following independent hyperpriors:

$$[a] \sim \text{Gamma}(\zeta, \iota), \quad [b] \propto 1.$$

Then the MMAP estimator

$$(a^*, b^*) = \arg \max_{a,b} \int p(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}) \pi(\sigma^2|a, b) d\boldsymbol{\beta} d\sigma^2 \pi(a, b)$$

exists and is finite.

- We can then use the estimated prior $\sigma^2 \sim IG(a^*, b^*)$ with MMAP estimates (a^*, b^*) within HEI
- Key: HEI with EB approximates a fully Bayesian EI procedure, without need for expensive MCMC (all closed-form!)

Empirical Bayes

- Marginal maximum a posteriori (MMAP, Doucet et al. 2002): add another layer of hyperpriors on (a, b) for regularization

Proposition

Assume the following independent hyperpriors:

$$[a] \sim \text{Gamma}(\zeta, \iota), \quad [b] \propto 1.$$

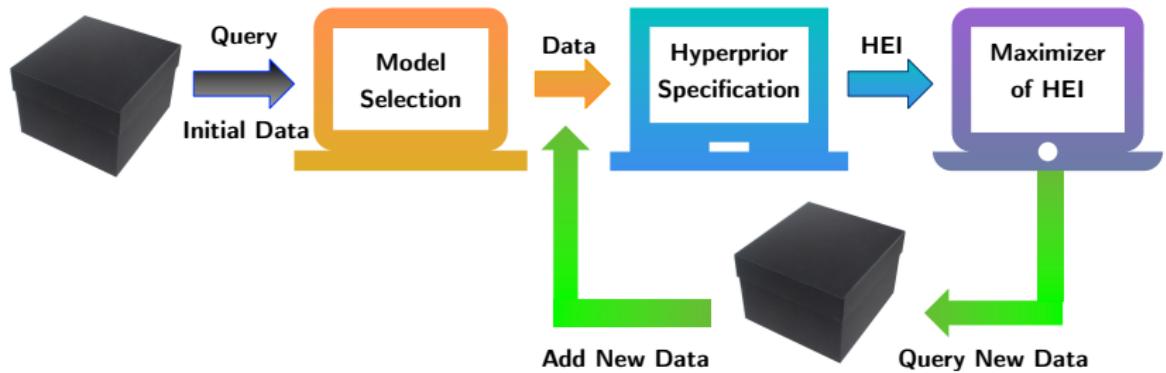
Then the MMAP estimator

$$(a^*, b^*) = \arg \max_{a,b} \int p(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}) \pi(\sigma^2|a, b) d\boldsymbol{\beta} d\sigma^2 \pi(a, b)$$

exists and is finite.

- We can then use the estimated prior $\sigma^2 \sim IG(a^*, b^*)$ with MMAP estimates (a^*, b^*) within HEI
- Key: HEI with EB approximates a **fully Bayesian** EI procedure, without need for **expensive** MCMC (all **closed**-form!)

Algorithm



Algorithm

Algorithm 1 Hierarchical Expected Improvement for Bayesian Optimization

Initialization

- Generate n_{ini} space-filling design points $\{\mathbf{x}_1, \dots, \mathbf{x}_{n_{\text{ini}}}\}$ on Ω .
- Evaluate function points $y_i = f(\mathbf{x}_i)$, yielding the initial dataset $\mathcal{D}_{n_{\text{ini}}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n_{\text{ini}}}$.

Model selection

- Select model order via BIC using (21).
- Estimate hyperparameters (a, b) via MMAP using (18).

Optimization

for $n \leftarrow n_{\text{ini}}$ to $n_{\text{tot}} - 1$ do

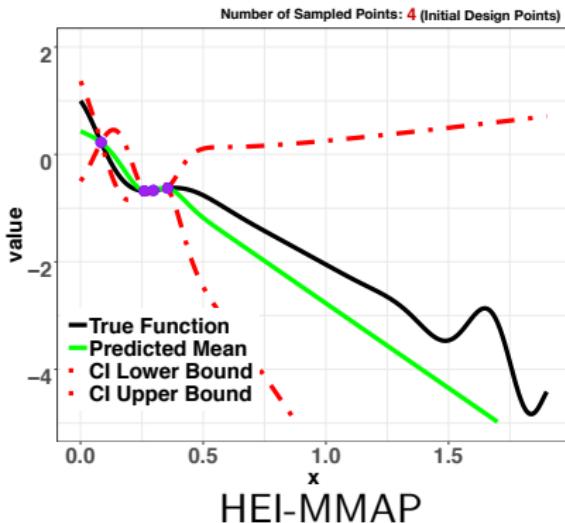
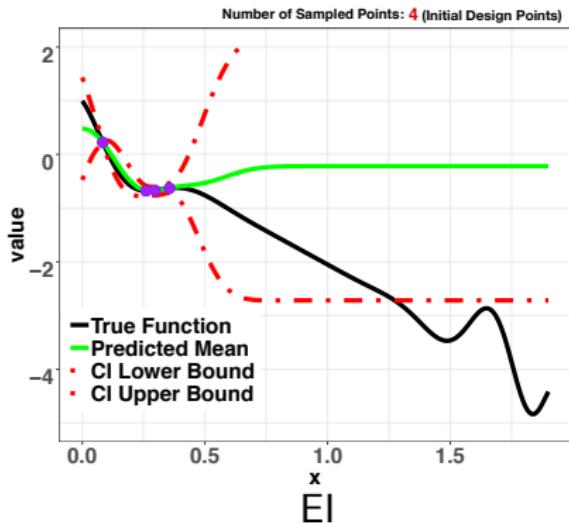
- Given \mathcal{D}_n , estimate length-scale parameters θ via MAP and compute $\text{HEI}_n(\mathbf{x})$.
- Obtain the next evaluation point \mathbf{x}_{n+1} by maximizing $\text{HEI}_n(\mathbf{x})$:

$$\mathbf{x}_{n+1} \leftarrow \underset{\mathbf{x} \in \Omega}{\operatorname{argmax}} \text{HEI}_n(\mathbf{x}).$$

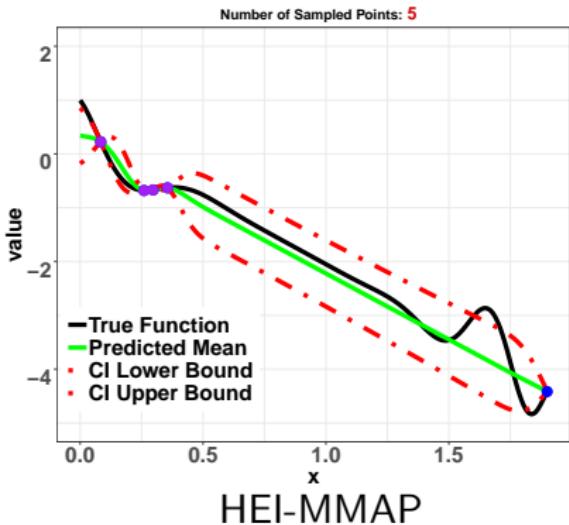
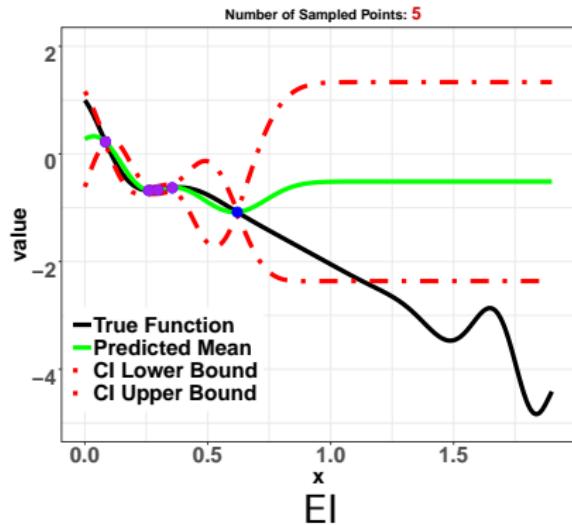
- Evaluate $y_{n+1} = f(\mathbf{x}_{n+1})$, and update data $\mathcal{D}_{n+1} = \mathcal{D}_n \cup \{(\mathbf{x}_{n+1}, y_{n+1})\}$.

Return: The best observed solution \mathbf{x}_{i^*} , where $i^* = \operatorname{argmin}_{i=1}^{n_{\text{tot}}} f(\mathbf{x}_i)$.

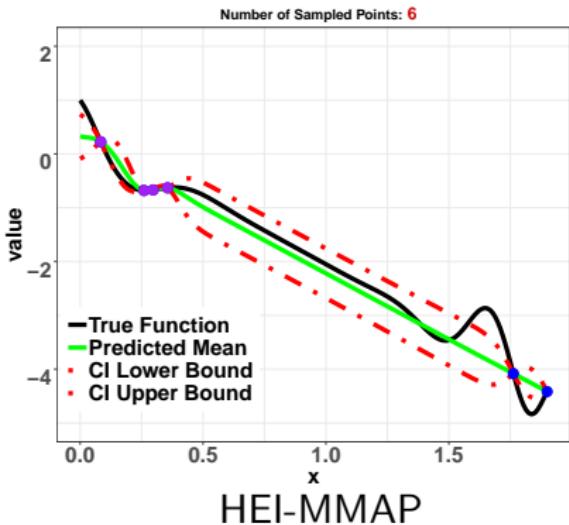
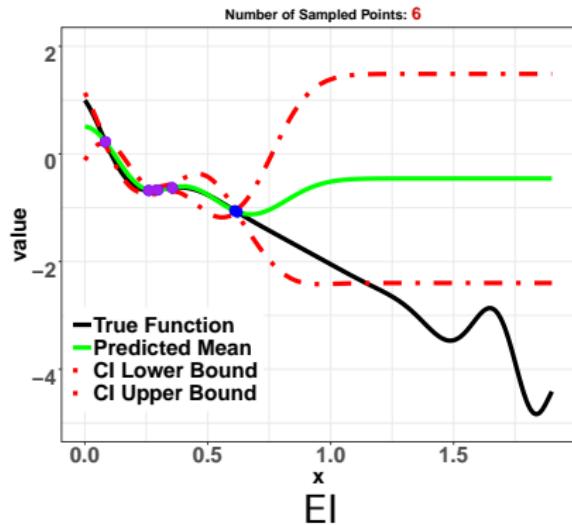
Back to example



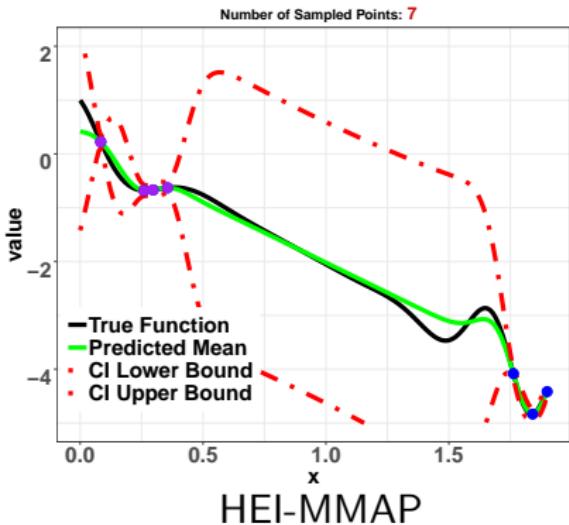
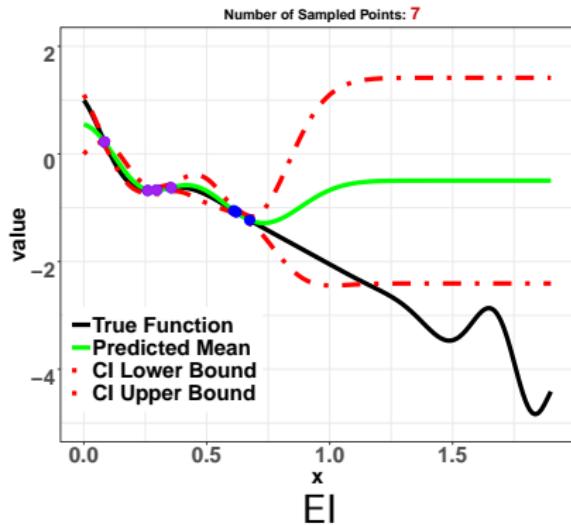
Back to example



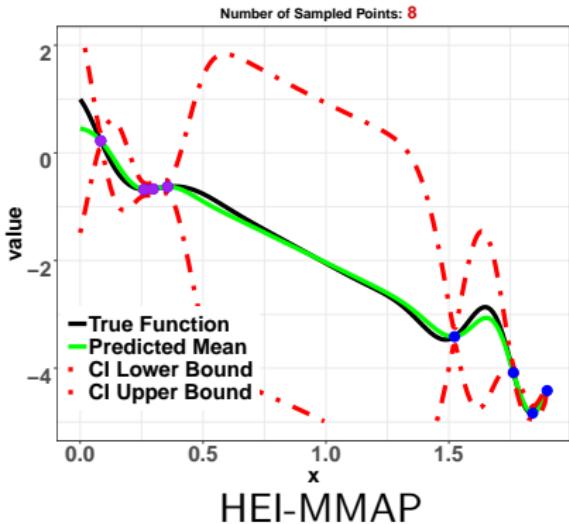
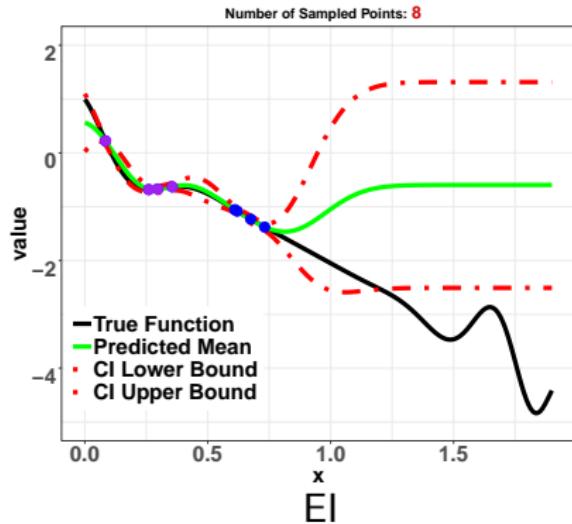
Back to example



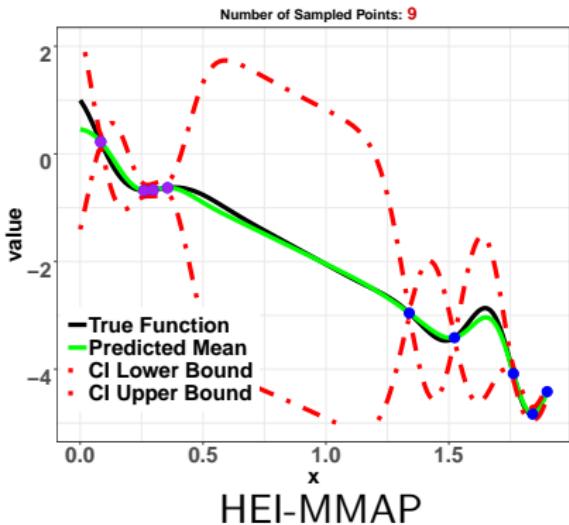
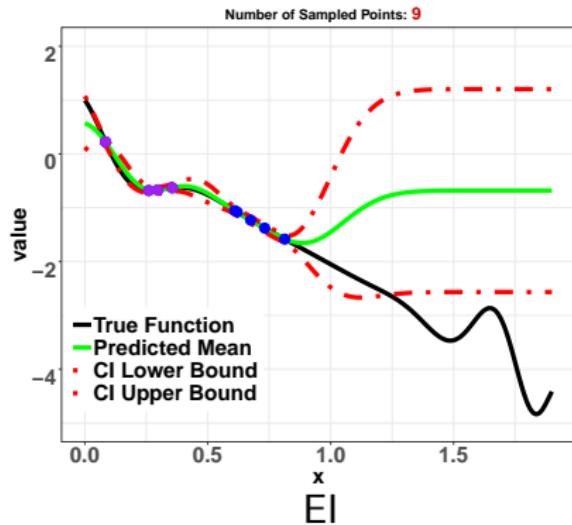
Back to example



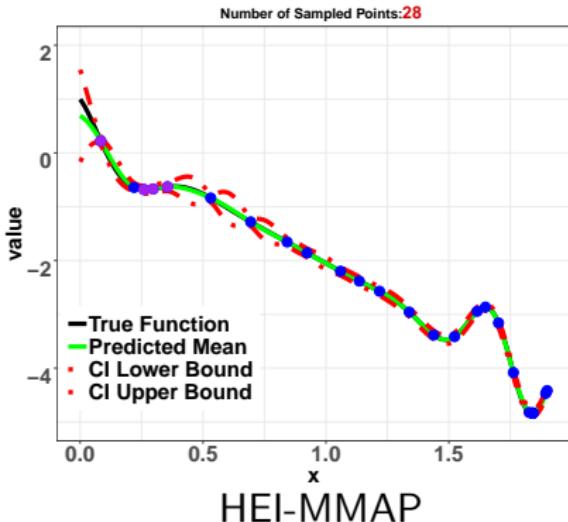
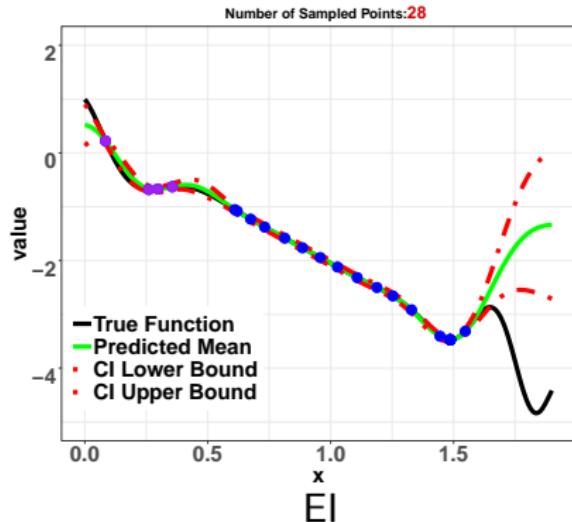
Back to example



Back to example



Back to example



- EI gets **stuck** at **local optimum**, HEI **finds** the **global optimum**

Global convergence

Theorem (Chen, Mak, Wu 2023)

Suppose a is a constant, $b = \Theta(n)$, and correlation $r(\mathbf{x})$ satisfies

$$|r(\mathbf{x}) - Q_{\lfloor 2\nu \rfloor}(\mathbf{x})| = \mathcal{O}(\|\mathbf{x}\|^{2\nu}(-\log \|\mathbf{x}\|)^{2\alpha}), \quad \text{as } \|\mathbf{x}\| \rightarrow 0,$$

where $Q_{\lfloor 2\nu \rfloor}(\mathbf{x})$ is a polynomial of degree $\lfloor 2\nu \rfloor$. Then for any f in the RKHS $\mathcal{H}_r(\Omega)$, we have:

$$\min_n f(\mathbf{x}_n) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) = \begin{cases} \mathcal{O}(n^{-\nu/d} (\log n)^\alpha), & \nu \leq 1 \\ \mathcal{O}(n^{-1/d}), & \nu > 1 \end{cases},$$

where $(\mathbf{x}_n)_{n=1}^\infty$ is the solution sequence by maximizing HEI.

- As sample size $n \rightarrow \infty$, $\min_n f(\mathbf{x}_n) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) \rightarrow 0$
- This shows that, unlike **plug-in EI**, HEI is **guaranteed** to converge to a **global minima**

Global convergence

Theorem (Chen, Mak, Wu 2023)

Suppose a is a constant, $b = \Theta(n)$, and correlation $r(\mathbf{x})$ satisfies

$$|r(\mathbf{x}) - Q_{\lfloor 2\nu \rfloor}(\mathbf{x})| = \mathcal{O}(\|\mathbf{x}\|^{2\nu}(-\log \|\mathbf{x}\|)^{2\alpha}), \quad \text{as } \|\mathbf{x}\| \rightarrow 0,$$

where $Q_{\lfloor 2\nu \rfloor}(\mathbf{x})$ is a polynomial of degree $\lfloor 2\nu \rfloor$. Then for any f in the RKHS $\mathcal{H}_r(\Omega)$, we have:

$$\min_n f(\mathbf{x}_n) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) = \begin{cases} \mathcal{O}(n^{-\nu/d} (\log n)^\alpha), & \nu \leq 1 \\ \mathcal{O}(n^{-1/d}), & \nu > 1 \end{cases},$$

where $(\mathbf{x}_n)_{n=1}^\infty$ is the solution sequence by maximizing HEI.

- As sample size $n \rightarrow \infty$, $\min_n f(\mathbf{x}_n) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) \rightarrow 0$
- This shows that, unlike **plug-in EI**, HEI is **guaranteed** to converge to a **global minima**

Global convergence

Theorem (Chen, Mak, Wu 2023)

Suppose a is a constant, $b = \Theta(n)$, and correlation $r(\mathbf{x})$ satisfies

$$|r(\mathbf{x}) - Q_{\lfloor 2\nu \rfloor}(\mathbf{x})| = \mathcal{O}(\|\mathbf{x}\|^{2\nu}(-\log \|\mathbf{x}\|)^{2\alpha}), \quad \text{as } \|\mathbf{x}\| \rightarrow 0,$$

where $Q_{\lfloor 2\nu \rfloor}(\mathbf{x})$ is a polynomial of degree $\lfloor 2\nu \rfloor$. Then for any f in the RKHS $\mathcal{H}_r(\Omega)$, we have:

$$\min_n f(\mathbf{x}_n) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) = \begin{cases} \mathcal{O}(n^{-\nu/d}(\log n)^\alpha), & \nu \leq 1 \\ \mathcal{O}(n^{-1/d}), & \nu > 1 \end{cases},$$

where $(\mathbf{x}_n)_{n=1}^\infty$ is the solution sequence by maximizing HEI.

- Theory requires a **data size dependent** hyperprior on $\sigma^2 \sim IG(a, b)$, namely $b = \Theta(n)$
- **Matérn correlation** (Stein 2012) satisfies the conditions on $r(\mathbf{x})$

Minimax rate

γ -stability condition (Wynne et al. 2020):

$$s_n(\mathbf{x}_{n+1}) \geq \gamma \|s_n(\mathbf{x})\|_\infty \quad \text{for all } n = 1, 2, \dots, \quad (1)$$

- $s_n(\mathbf{x})$: **posterior variance** term from GP

With this, we can then show the following optimization rate:

Theorem (Chen, Mak, Wu 2023)

Assume the earlier conditions hold. Then for any $f \in \mathcal{H}_r(\Omega)$:

$$\min_n f(\mathbf{x}_n) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) = O(n^{-\nu/d}),$$

where $(\mathbf{x}_n)_{n=1}^\infty$ is the solution sequence maximizing HEI with the γ -stability condition, with $\gamma \in (0, 1)$.

This is precisely the **minimax** rate for $f \in \mathcal{H}_r(\Omega)$ (Bull 2011)!

Minimax rate

γ -**stability** condition ([Wynne et al. 2020](#)):

$$s_n(\mathbf{x}_{n+1}) \geq \gamma \|s_n(\mathbf{x})\|_\infty \quad \text{for all } n = 1, 2, \dots, \quad (1)$$

- $s_n(\mathbf{x})$: **posterior variance** term from GP

With this, we can then show the following optimization rate:

Theorem (Chen, Mak, Wu 2023)

Assume the earlier conditions hold. Then for any $f \in \mathcal{H}_r(\Omega)$:

$$\min_n f(\mathbf{x}_n) - \min_{\mathbf{x} \in \Omega} f(\mathbf{x}) = O(n^{-\nu/d}),$$

where $(\mathbf{x}_n)_{n=1}^\infty$ is the solution sequence maximizing HEI with the γ -stability condition, with $\gamma \in (0, 1)$.

This is precisely the **minimax** rate for $f \in \mathcal{H}_r(\Omega)$ (Bull 2011)!

Minimax rate

Algorithm 1 Hierarchical Expected Improvement for Bayesian Optimization

Initialization

- Generate n_{ini} space-filling design points $\{\mathbf{x}_1, \dots, \mathbf{x}_{n_{\text{ini}}}\}$ on Ω .
- Evaluate function points $y_i = f(\mathbf{x}_i)$, yielding the initial dataset $\mathcal{D}_{n_{\text{ini}}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n_{\text{ini}}}$.

Model selection

- Select model order via BIC using (21).
- Estimate hyperparameters (a, b) via MMAP using (18).

Optimization

for $n \leftarrow n_{\text{ini}}$ to $n_{\text{tot}} - 1$ do

- Given \mathcal{D}_n , estimate length-scale parameters θ via MAP and compute $\text{HEI}_n(\mathbf{x})$.

- Obtain the next evaluation point \mathbf{x}_{n+1} by maximizing $\text{HEI}_n(\mathbf{x})$:

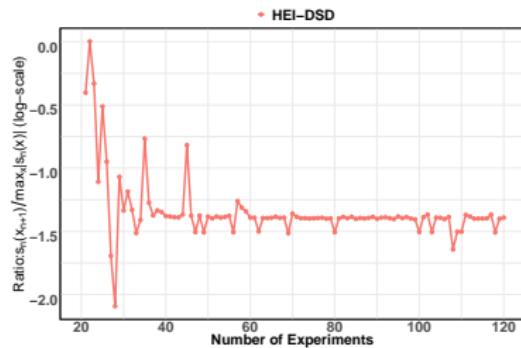
**Add in stability
condition as constraint**
$$\mathbf{x}_{n+1} \leftarrow \underset{\mathbf{x} \in \Omega}{\operatorname{argmax}} \text{HEI}_n(\mathbf{x}).$$

- Evaluate $y_{n+1} = f(\mathbf{x}_{n+1})$, and update data $\mathcal{D}_{n+1} = \mathcal{D}_n \cup \{(\mathbf{x}_{n+1}, y_{n+1})\}$.

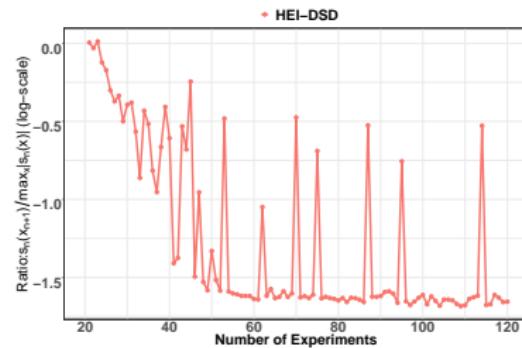
Return: The best observed solution \mathbf{x}_{i^*} , where $i^* = \operatorname{argmin}_{i=1}^{n_{\text{tot}}} f(\mathbf{x}_i)$.

Minimax rate

... but in practice, this γ -stability condition seems to **hold** without an **explicit** constraint on the HEI



Branin



Camel Three-Hump

Minimax rate

Why is our rate $O(n^{-\nu/d})$ **minimax**?

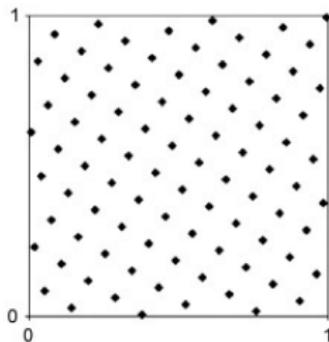
- Bull (2011): the minimax rate is $O(n^{-\nu/d})$ for $f \in \mathcal{H}_r(\Omega)$
- This rate, however, can be achieved by a **non-adaptive** quasi-uniform strategy, which performs **terribly** in practice!

Rates don't tell the whole story!

Minimax rate

Why is our rate $O(n^{-\nu/d})$ **minimax**?

- Bull (2011): the minimax rate is $O(n^{-\nu/d})$ for $f \in \mathcal{H}_r(\Omega)$
- This rate, however, can be achieved by a **non-adaptive quasi-uniform** strategy, which performs **terribly** in practice!

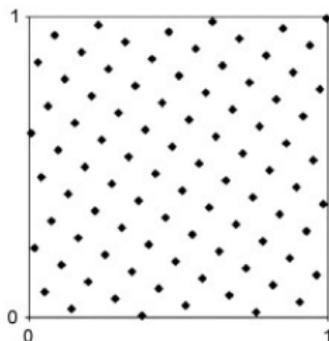


Rates don't tell the whole story!

Minimax rate

Why is our rate $O(n^{-\nu/d})$ **minimax**?

- Bull (2011): the minimax rate is $O(n^{-\nu/d})$ for $f \in \mathcal{H}_r(\Omega)$
- This rate, however, can be achieved by a **non-adaptive quasi-uniform** strategy, which performs **terribly** in practice!



Rates don't tell the whole story!

Test functions

HEI methods:

- HEI-Weak, HEI-MMAP, HEI-DSD

Plug-in EI methods ([Jones et al 1997](#)):

- EI-OK, EI-UK

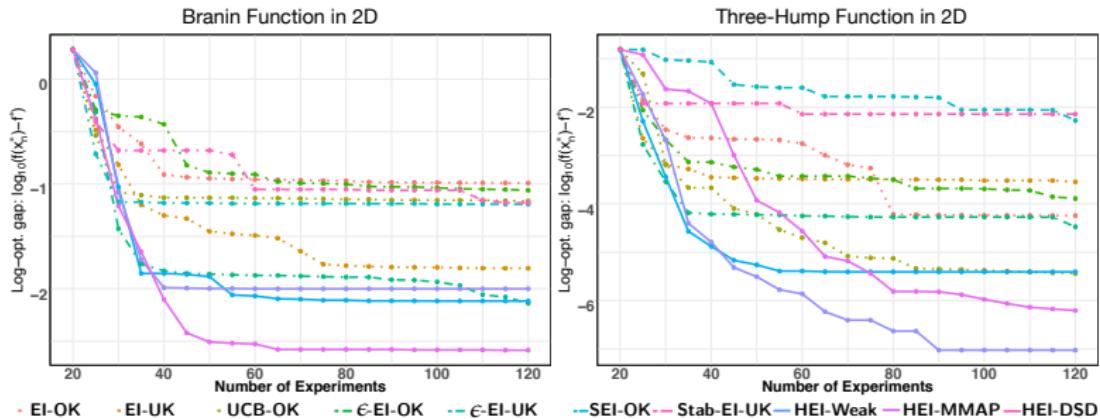
ϵ -**greedy EI** methods ([Bull 2011](#)):

- ϵ -EI-OK, ϵ -EI-UK, with $\epsilon = 0.1$ ([Sutton & Barto 2018](#))

UCB-EI methods ([Srinivas et al 2010](#)):

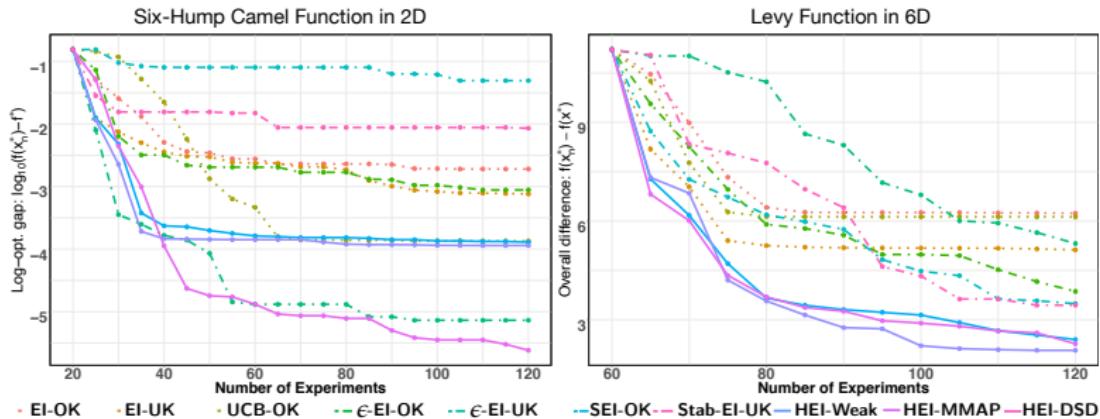
- Exploration parameter $\rho = 2.96$

Test functions



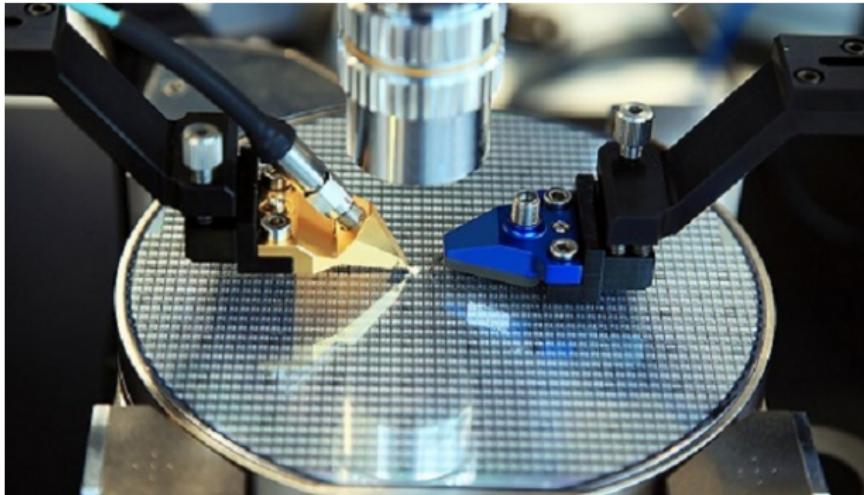
- HEI-MMAP & HEI-DSD perform **best**, EI **stuck in local optima**
- ϵ -greedy approaches **not bad**, but not as **effective** as HEI
(ad-hoc fix to **over-exploitation**)
- HEI provides a **principled** correction via **hierarchical modeling**

Test functions



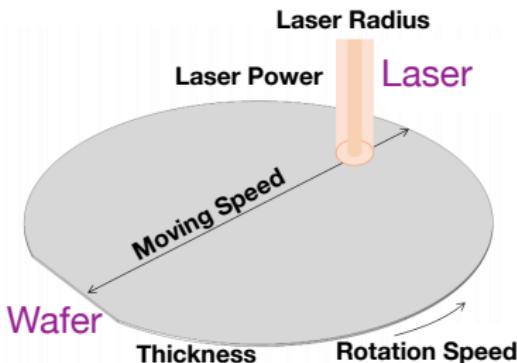
- HEI-MMAP & HEI-DSD perform **best**, EI **stuck** in **local optima**
- ϵ -greedy approaches **not bad**, but not as **effective** as HEI
(**ad-hoc** fix to **over-exploitation**)
- HEI provides a **principled** correction via **hierarchical modeling**

Semiconductor manufacturing optimization



- **Semiconductor manufacturing**: process for manufacturing integrated circuit chips used in **electrical devices**
- **Thermal processing** is a crucial stage, allowing for necessary **chemical reactions** and **surface oxidation** ([Singh et al 2000](#))

Semiconductor manufacturing optimization



- $d = 5$ control variables for optimization
- Goal: heat silicon wafer to a target temperature $T = 600\text{F}$ quickly and smoothly

Semiconductor manufacturing optimization

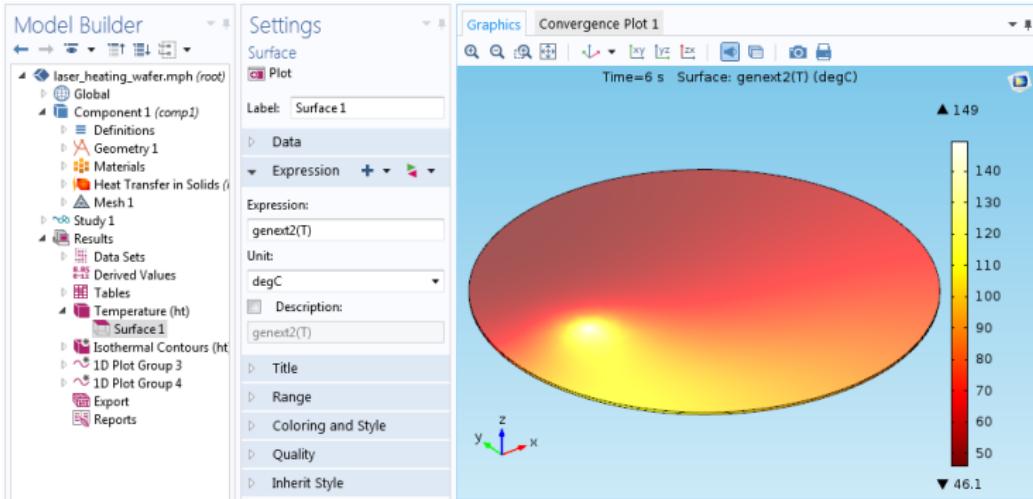
Objective function $f(x)$ to minimize:

$$f(x) := \sum_{t=1}^{60} \max_{\mathbf{s} \in \mathcal{S}} \underbrace{\left| \mathcal{T}_t(\mathbf{s}; x) - 600 \right|}_{\text{Gap to Target}}$$

- \mathbf{s} : **location** on wafer \mathcal{S}
- $\mathcal{T}_t(\mathbf{s}; x)$: **temperature** at time t and location \mathbf{s}

Settings x with **small** $f(x)$ ensure that thermal processing is performed **quickly** and **smoothly**

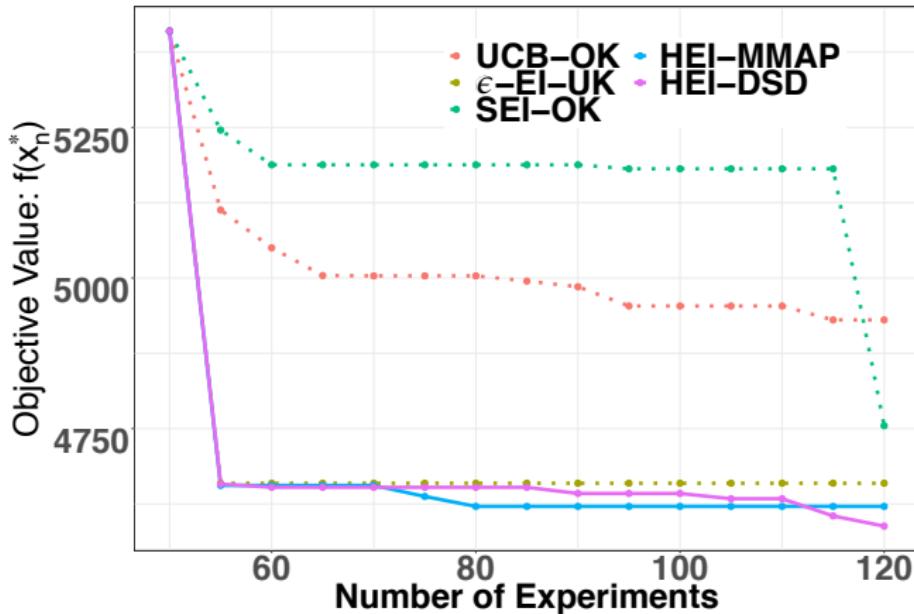
Semiconductor manufacturing optimization



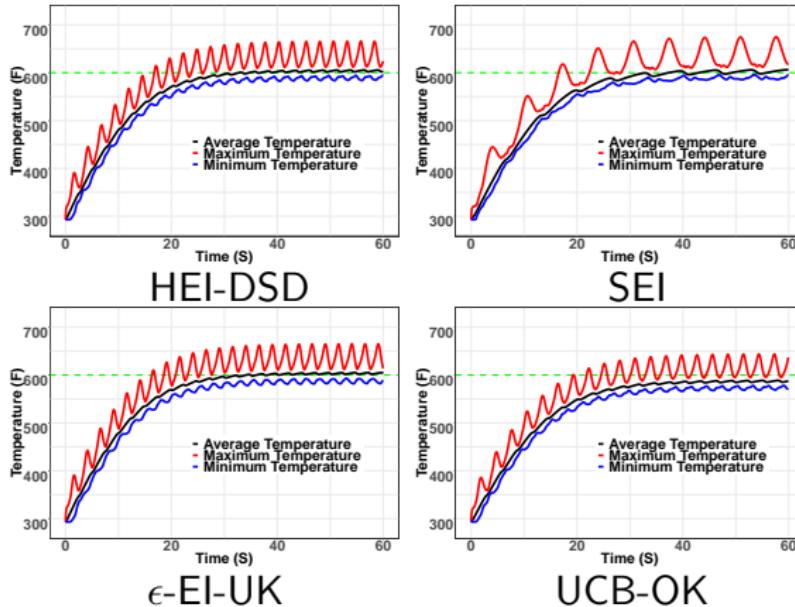
f is **expensive** and **black-box**:

- **Experiments** are **costly** and require **prototyping**
- **Simulations** can be performed on **COMSOL**, but require **30 minutes** for **each** setting

Semiconductor manufacturing optimization

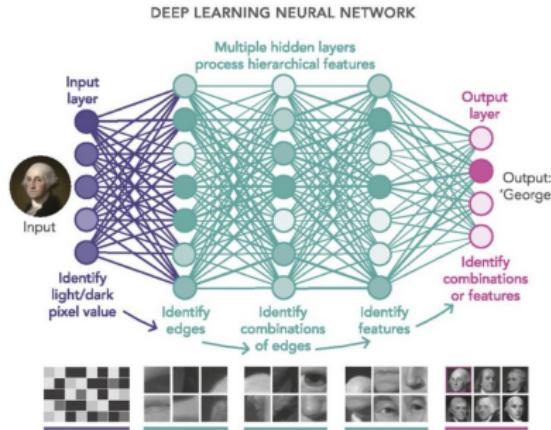


Semiconductor manufacturing optimization



- Visually, HEI-DSD and ϵ -EI-UK give the **best** settings: **average** temperature reaches **target** T **quickly** and **smoothly** (variation over wafer is **small**)

Hyperparameter tuning in deep learning

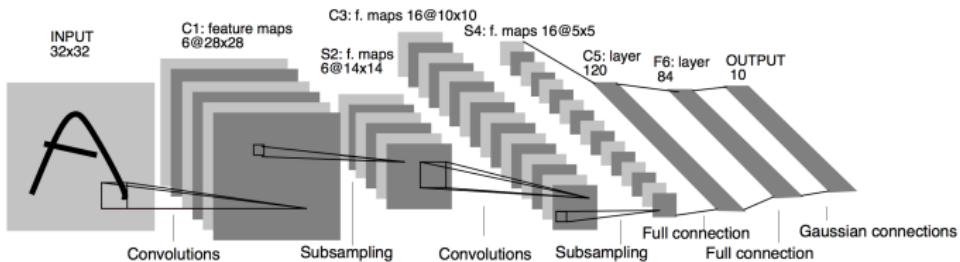


Neural networks (NNs) are widely used in **ML** problems:

- ... but in order to work **well**, many **hyperparameters** need to be carefully **tuned**
- NN training is **expensive** even for a **single** hyperparameter setting (several **hours** for real image databases, such as CIFAR)

Hyperparameter tuning in deep learning

- **Architecture:** LeNet-5 ([LeCun et al. 1998](#))



- **Formulation:** [NN training](#) solves the following problem

$$\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \ell(\mathcal{F}(\mathbf{x}; \mathbf{w}), y)$$

- ℓ : **loss function** (**softmax** used here)
- $\mathcal{F}(\cdot; \mathbf{w})$: **neural network** with weight matrix \mathbf{w}
- $(\mathbf{x}, y) \sim \mathcal{D}$: **image** and corresponding **label**

Hyperparameter tuning in deep learning

- **Training Method:**

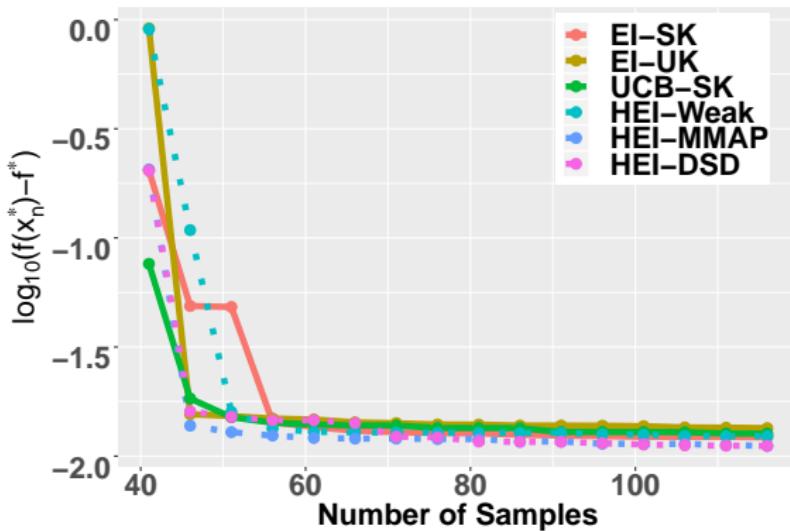
Momentum stochastic gradient descent with weight decay:

$$\boldsymbol{w}^{(t+1)} = \boldsymbol{w}^{(t)} - \underbrace{\eta \cdot \nabla_{\boldsymbol{w}} \ell(\mathcal{F}(\boldsymbol{x}_{i_t}, \boldsymbol{w}^{(t)}), y_{i_t})}_{\text{Gradient Descent Term}} + \underbrace{\mu \cdot (\boldsymbol{w}^{(t)} - \boldsymbol{w}^{(t-1)})}_{\text{Momentum Term}} - \underbrace{\lambda \cdot \boldsymbol{w}^{(t)}}_{\text{Weight Decay Term}}$$

- **Hyperparameters:**

- η : learning rate
- μ : momentum
- τ : weight matrix initialization variance
- λ : weight decay

Hyperparameter tuning in deep learning



- *y-axis:* log-validation error of trained NN
- HEI-DSD & HEI-MMAP perform **best**: lowest errors
- Existing tuning method (**gradient-based tuning**, Bengio 2000) gives a **higher log-error** of -1.854

Conclusions

- Exploration is important for learning and optimizing expensive black-box functions
- A closed-form acquisition function allows for efficient optimization
- A hierarchical Bayesian model which considers parameter uncertainties can encourage more exploration and correct the greediness (overexploitation) of EI

Questions?

