

# Advanced Methods for Bayesian Optimization in Complex Scenarios Course

**Presented by: Eduardo C. Garrido–Merchán.**

Assistant Professor at Universidad Pontificia Comillas

Quantitative Methods Department

Former PhD student of Daniel Hernandez Lobato

June, 2024.

- ▶ Introduction and Motivation.

# Index

- ▶ Introduction and Motivation.
- ▶ Fundamentals of Bayesian Optimization.

# Index

- ▶ Introduction and Motivation.
- ▶ Fundamentals of Bayesian Optimization.
- ▶ Hands on Bayesian optimization with BOTorch.



# Index

- ▶ Introduction and Motivation.
- ▶ Fundamentals of Bayesian Optimization.
- ▶ Hands on Bayesian optimization with BOTorch.
- ▶ Parallel Constrained Multi-objective Bayesian optimization.

# Index

- ▶ Introduction and Motivation.
- ▶ Fundamentals of Bayesian Optimization.
- ▶ Hands on Bayesian optimization with BOTorch.
- ▶ Parallel Constrained Multi-objective Bayesian optimization.
- ▶ Advanced scenarios: Mixed, multi-fidelity, high-dim, topologies.

# Index

- ▶ Introduction and Motivation.
- ▶ Fundamentals of Bayesian Optimization.
- ▶ Hands on Bayesian optimization with BOTorch.
- ▶ Parallel Constrained Multi-objective Bayesian optimization.
- ▶ Advanced scenarios: Mixed, multi-fidelity, high-dim, topologies.
- ▶ Advanced hands on Bayesian optimization with BOTorch.

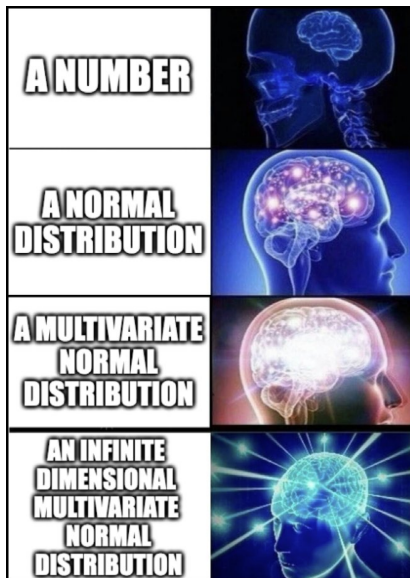
# Index

- ▶ Introduction and Motivation.
- ▶ Fundamentals of Bayesian Optimization.
- ▶ Hands on Bayesian optimization with BOTorch.
- ▶ Parallel Constrained Multi-objective Bayesian optimization.
- ▶ Advanced scenarios: Mixed, multi-fidelity, high-dim, topologies.
- ▶ Advanced hands on Bayesian optimization with BOTorch.
- ▶ References.

# First Session: Bayesian optimization fundamentals

# Introduction.

GPs are great... so now... how can they be used for optimization?



# Challenges in Machine Learning, Finance, Robotics, Engineering, Business...

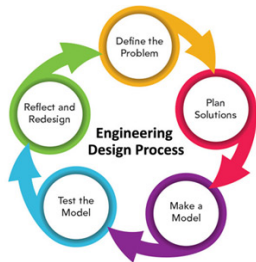
**Companies face complex scenarios dealing with lots of  
scenarios!**





# Challenges in Machine Learning, Finance, Robotics, Engineering, Business...

**Companies face complex scenarios dealing with lots of scenarios!**

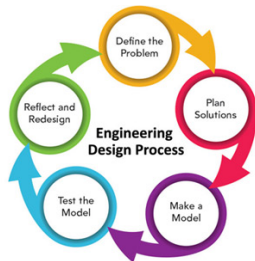


# Challenges in Machine Learning, Finance, Robotics, Engineering, Business...

**Companies face complex scenarios dealing with lots of scenarios!**



► Many choices at each step.

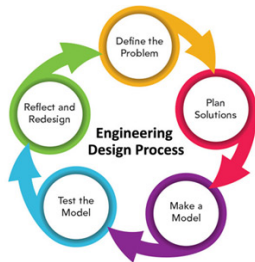


# Challenges in Machine Learning, Finance, Robotics, Engineering, Business...

**Companies face complex scenarios dealing with lots of scenarios!**

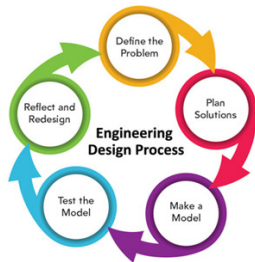


- ▶ Many choices at each step.
- ▶ Complicated and high dimensional.



# Challenges in Machine Learning, Finance, Robotics, Engineering, Business...

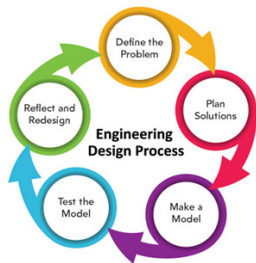
**Companies face complex scenarios dealing with lots of scenarios!**



- ▶ Many choices at each step.
- ▶ Complicated and high dimensional.
- ▶ Difficult for individuals to reason about.

# Challenges in Machine Learning, Finance, Robotics, Engineering, Business...

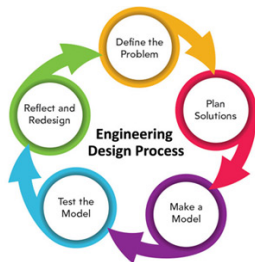
**Companies face complex scenarios dealing with lots of scenarios!**



- ▶ Many choices at each step.
- ▶ Complicated and high dimensional.
- ▶ Difficult for individuals to reason about.
- ▶ Prone to human bias.

# Challenges in Machine Learning, Finance, Robotics, Engineering, Business...

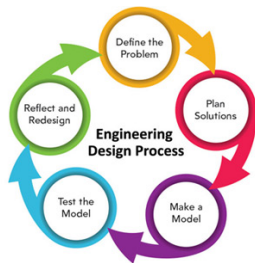
**Companies face complex scenarios dealing with lots of scenarios!**



- ▶ Many choices at each step.
- ▶ Complicated and high dimensional.
- ▶ Difficult for individuals to reason about.
- ▶ Prone to human bias.
- ▶ It may be possible to test various models in parallel.

# Challenges in Machine Learning, Finance, Robotics, Engineering, Business...

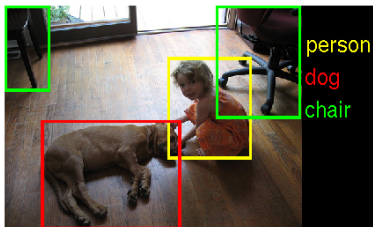
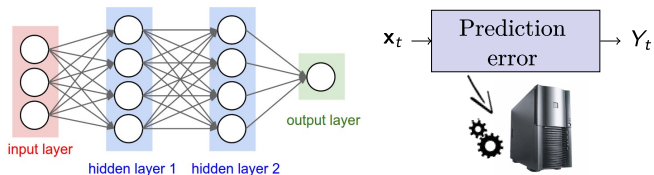
**Companies face complex scenarios dealing with lots of scenarios!**



- ▶ Many choices at each step.
- ▶ Complicated and high dimensional.
- ▶ Difficult for individuals to reason about.
- ▶ Prone to human bias.
- ▶ It may be possible to test various models in parallel.

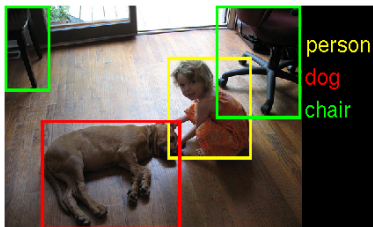
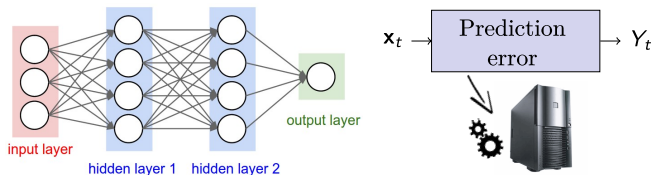
**Optimization is a challenging task in real-life choices!**

## Example: **Deep Neural Network** for object recognition.



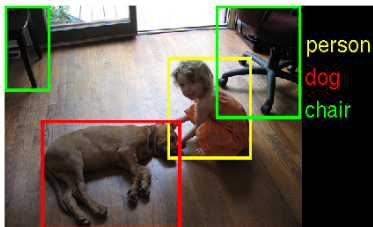
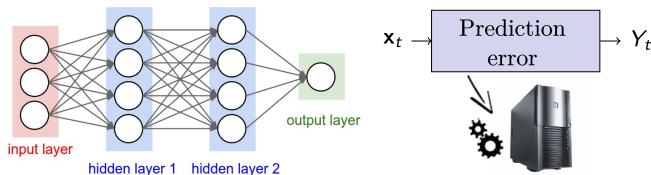


Example: **Deep Neural Network** for object recognition.



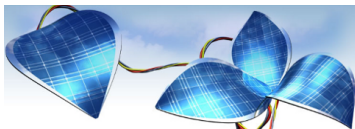
**Parameters to tune:** Number of neurons, number of layers, learning-rate, level of regularization, momentum, etc.

Example: **Deep Neural Network** for object recognition.



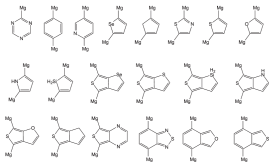
**Parameters to tune:** Number of neurons, number of layers, learning-rate, level of regularization, momentum, etc. If **multiple processors** were available, we could test **various configurations in parallel**, in order to **gain more information**.

Example: new **plastic solar cells** for transforming light into electricity.



### Library generation

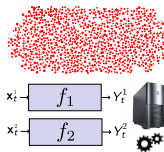
#### Fragments



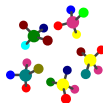
#### Bonding rules



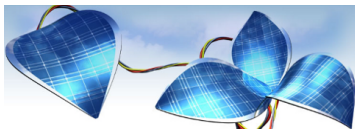
### Performance evaluation



### Interesting molecules

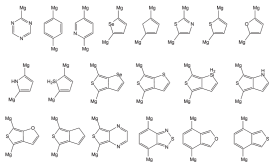


Example: new **plastic solar cells** for transforming light into electricity.



### Library generation

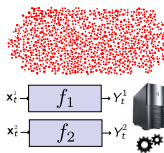
#### Fragments



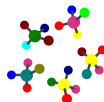
#### Bonding rules



### Performance evaluation



### Interesting molecules

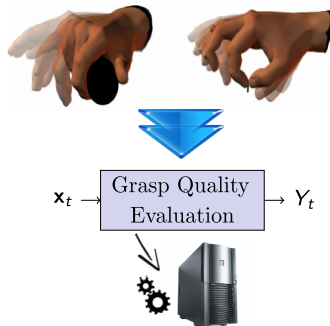


Explore **millions of candidate molecule structures** to identify the compounds with the best properties.

Example: **control system** for a robot that is able to grasp objects.



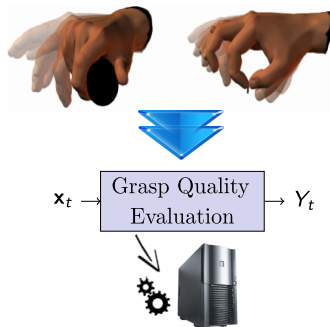
## Finger Joint Trajectories



Example: **control system** for a robot that is able to grasp objects.



### Finger Joint Trajectories



**Parameters to tune:** initial pose for the robot's hand and finger joint trajectories.

# Bayesian optimization, the secret sauce of AlphaGo (and Alphazero)

---

## Bayesian Optimization in AlphaGo

---

Yutian Chen, Aja Huang, Ziyu Wang, Ioannis Antonoglou, Julian Schrittwieser,  
David Silver & Nando de Freitas

DeepMind, London, UK  
yutianc@google.com

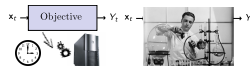
### Abstract

During the development of AlphaGo, its many hyper-parameters were tuned with Bayesian optimization multiple times. This automatic tuning process resulted in substantial improvements in playing strength. For example, prior to the match with Lee Sedol, we tuned the latest AlphaGo agent and this improved its win-rate from 50% to 66.5% in self-play games. This tuned version was deployed in the final match. Of course, since we tuned AlphaGo many times during its development cycle, the compounded contribution was even higher than this percentage. It is our hope that this brief case study will be of interest to Go fans, and also provide Bayesian optimization practitioners with some insights and inspiration.



# Optimization Problems: Common Features

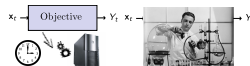
- Very expensive evaluations.





# Optimization Problems: Common Features

- Very expensive evaluations.

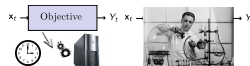


- The objective is a black-box.

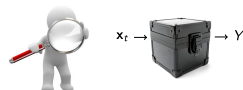


# Optimization Problems: Common Features

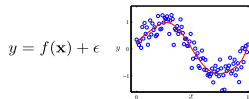
- Very expensive evaluations.



- The objective is a black-box.

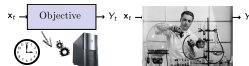


- The evaluation can be noisy.

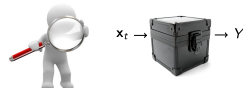


# Optimization Problems: Common Features

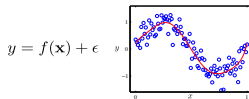
- Very expensive evaluations.



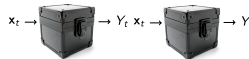
- The objective is a black-box.



- The evaluation can be noisy.



- Evaluations may be done in parallel.



**Bayesian optimization methods can be used to solve these problems!**

# Activity

Think about a real problem that can be solved with Bayesian optimization. 15 mins individual. 10 mins in groups. 3 presentations with slides.

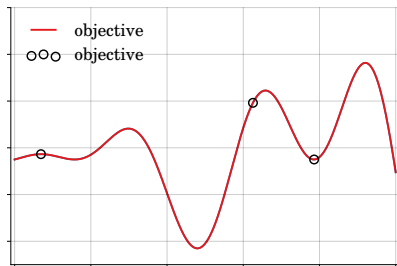
- ▶ Which is the function to optimize?
- ▶ Does it satisfy the BO assumptions?
- ▶ Which are the independent variables that are taken into account?
- ▶ Which are the lower and upper bounds in the optimization?
- ▶ Is  $f(X)$  a sample of a Gaussian process?
- ▶ What is the prior mean of the Gaussian process?
- ▶ Which kernel would you use?

# Bayesian Optimization



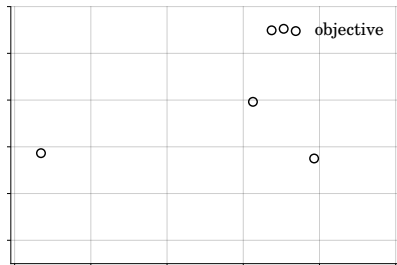
1. Get initial sample.

# Bayesian Optimization



1. Get initial sample.

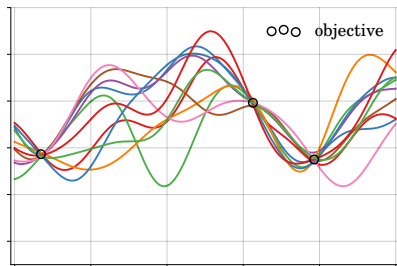
# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

$$p(y|x, \mathcal{D}_n).$$

# Bayesian Optimization

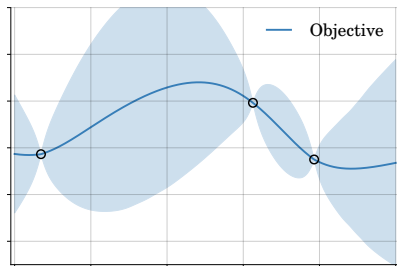


1. Get initial sample.
2. Fit a model to the data:

$$p(y|x, \mathcal{D}_n).$$



# Bayesian Optimization



1. Get initial sample.

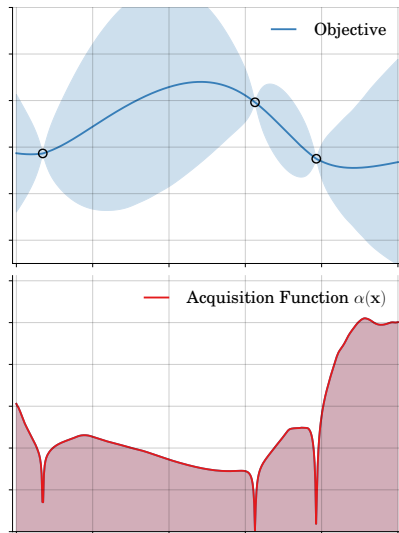
2. Fit a model to the data:

$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = \mathbb{E}_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

# Bayesian Optimization



1. Get initial sample.

2. Fit a model to the data:

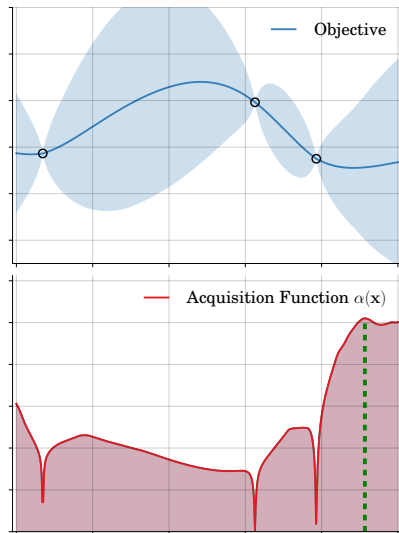
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

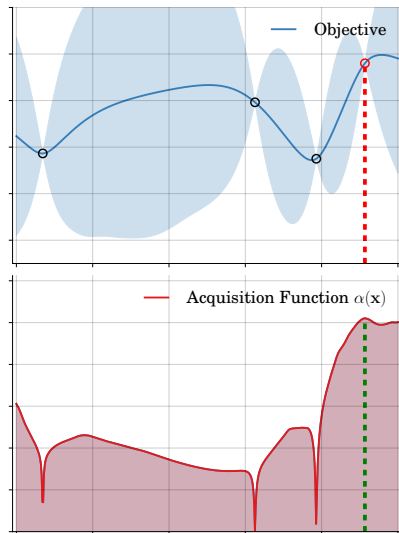
4. Optimize acquisition function  $\alpha(x)$ .

# Bayesian Optimization



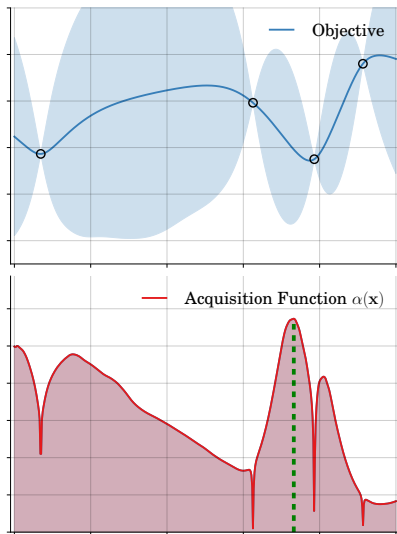
1. Get initial sample.
2. Fit a model to the data:  
$$p(y|x, \mathcal{D}_n).$$
3. Select data collection strategy:  
$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$
4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.

# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:  
$$p(y|x, \mathcal{D}_n).$$
3. Select data collection strategy:  
$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$
4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!

# Bayesian Optimization

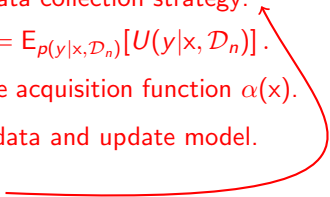


1. Get initial sample.
2. Fit a model to the data:

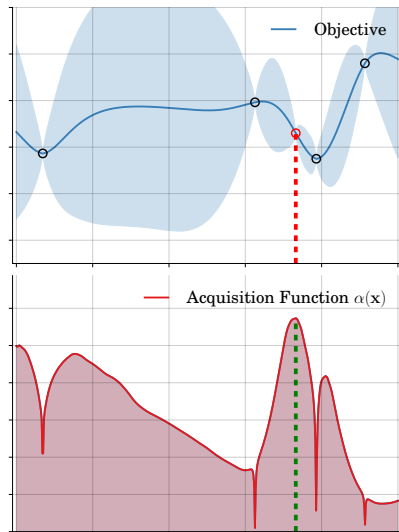
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat! 

# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

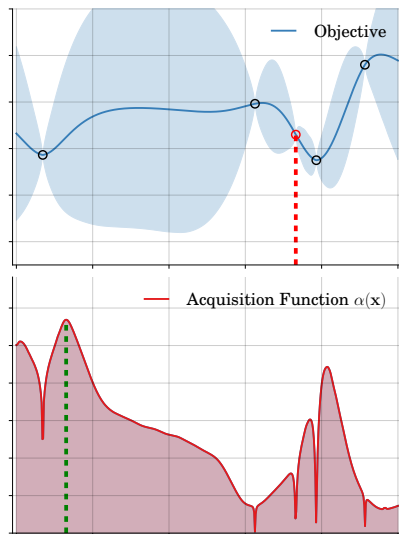
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!

# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

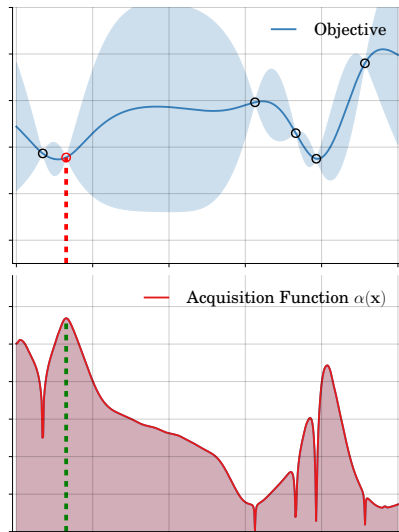
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!

# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

$$p(y|x, \mathcal{D}_n).$$

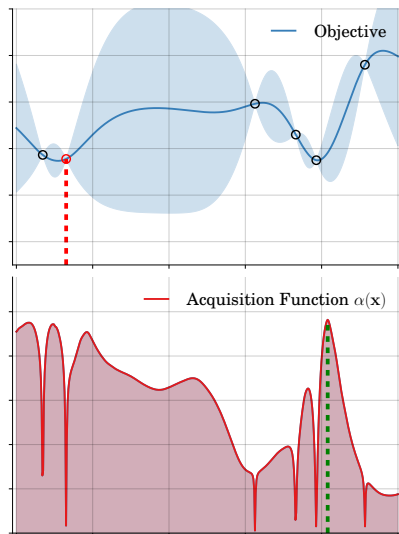
3. Select data collection strategy:

$$\alpha(x) = \mathbb{E}_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!



# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

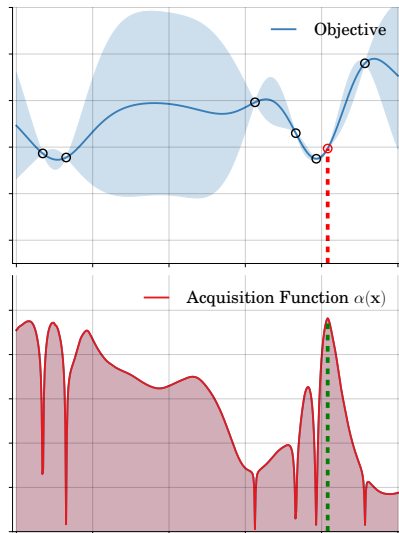
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!

# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

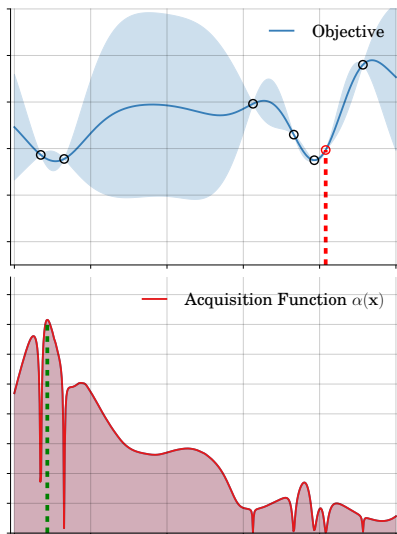
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!

# Bayesian Optimization

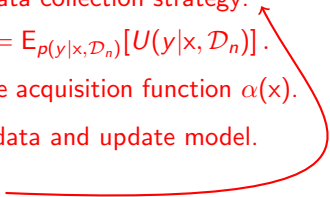


1. Get initial sample.
2. Fit a model to the data:

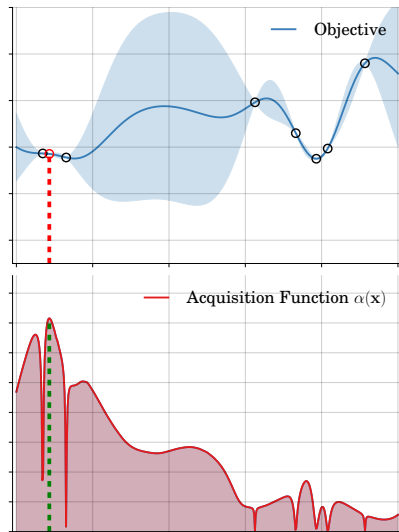
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat! 

# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

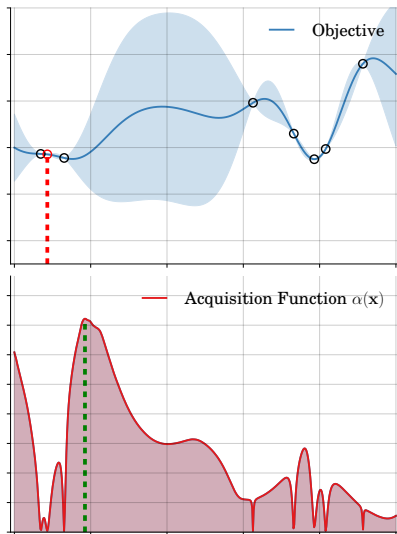
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = \mathbb{E}_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!

# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

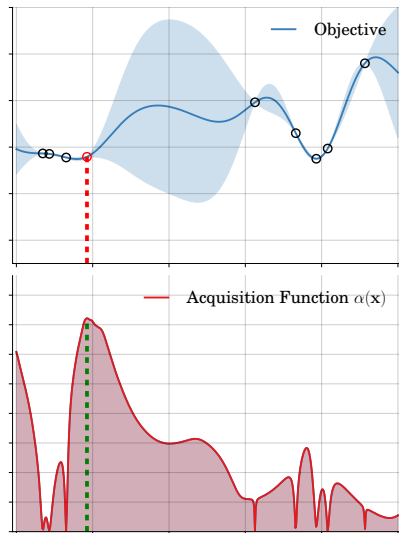
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!

# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

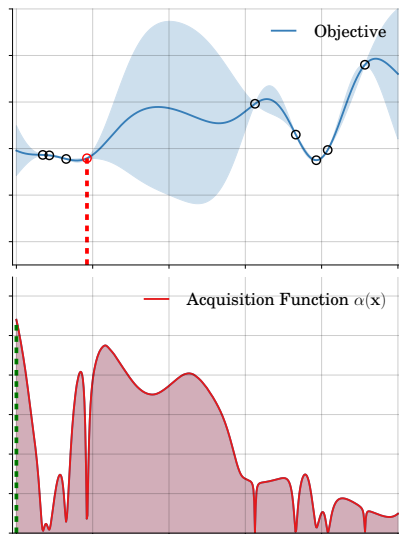
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!

# Bayesian Optimization

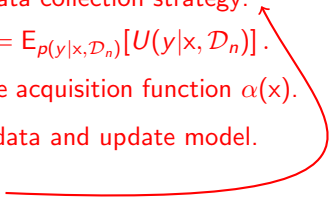


1. Get initial sample.
2. Fit a model to the data:

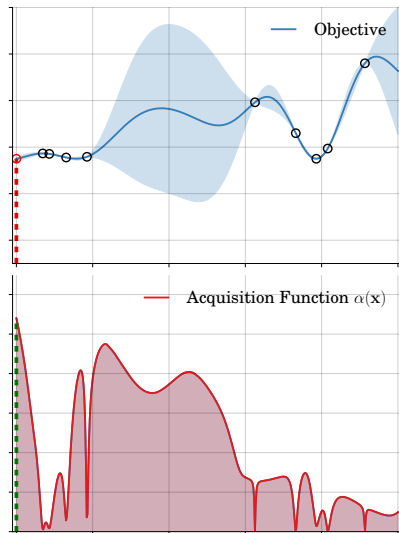
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat! 

# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

$$p(y|x, \mathcal{D}_n).$$

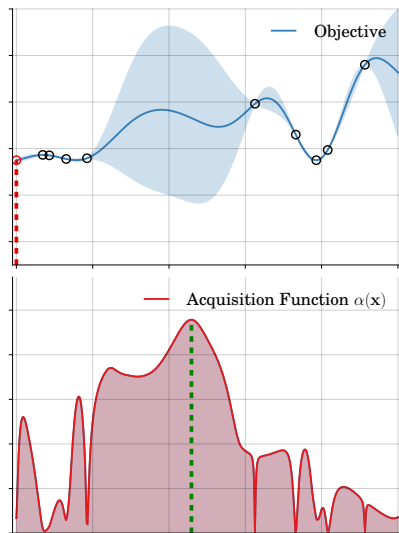
3. Select data collection strategy:

$$\alpha(x) = \mathbb{E}_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!



# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

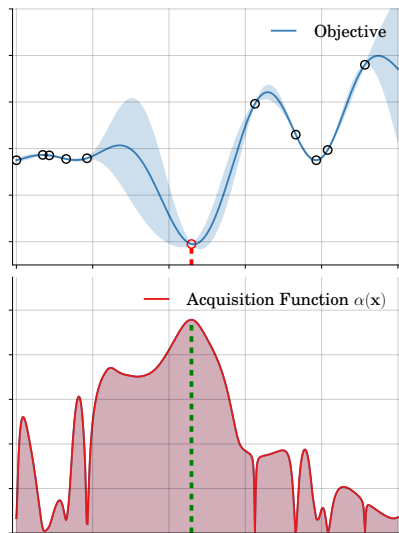
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!

# Bayesian Optimization



1. Get initial sample.
2. Fit a model to the data:

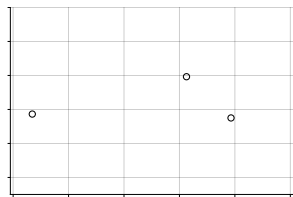
$$p(y|x, \mathcal{D}_n).$$

3. Select data collection strategy:

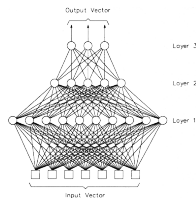
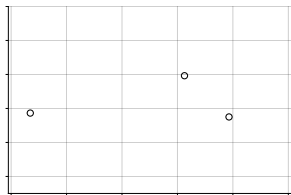
$$\alpha(x) = E_{p(y|x, \mathcal{D}_n)}[U(y|x, \mathcal{D}_n)].$$

4. Optimize acquisition function  $\alpha(x)$ .
5. Collect data and update model.
6. Repeat!

# Fitting a Model to the Data



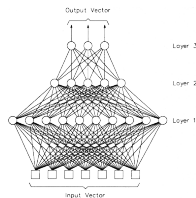
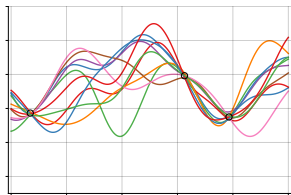
# Fitting a Model to the Data



$$h_j(\mathbf{x}) = \tanh\left(\sum_{i=1}^I x_i w_{ji}\right)$$

$$f(\mathbf{x}) = \sum_{j=1}^H v_j h_j(\mathbf{x})$$

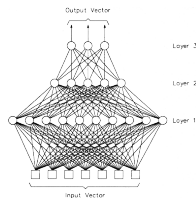
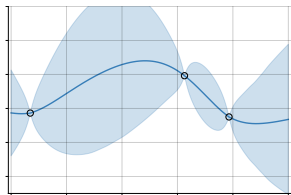
# Fitting a Model to the Data



$$h_j(\mathbf{x}) = \tanh\left(\sum_{i=1}^I x_i w_{ji}\right)$$

$$f(\mathbf{x}) = \sum_{j=1}^H v_j h_j(\mathbf{x})$$

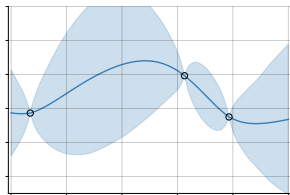
# Fitting a Model to the Data



$$h_j(\mathbf{x}) = \tanh\left(\sum_{i=1}^I x_i w_{ji}\right)$$

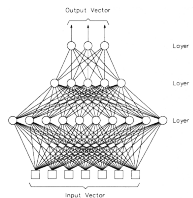
$$f(\mathbf{x}) = \sum_{j=1}^H v_j h_j(\mathbf{x})$$

# Fitting a Model to the Data



**Posterior Dist.**

**Predictive Dist.**



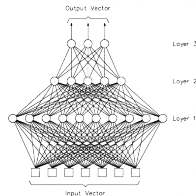
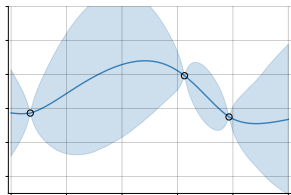
$$h_j(\mathbf{x}) = \tanh\left(\sum_{i=1}^I x_i w_{ji}\right)$$

$$f(\mathbf{x}) = \sum_{j=1}^H v_j h_j(\mathbf{x})$$

$$p(\mathbf{W}|\text{Data}) = p(\mathbf{W})p(\text{Data}|\mathbf{W})/p(\text{Data})$$

$$p(y|\text{Data}, x) = \int p(y|\mathbf{W}, x) p(\mathbf{W}|\text{Data}) d\mathbf{W}$$

# Fitting a Model to the Data



$$h_j(\mathbf{x}) = \tanh\left(\sum_{i=1}^I x_i w_{ji}\right)$$

$$f(\mathbf{x}) = \sum_{j=1}^H v_j h_j(\mathbf{x})$$

**Posterior Dist.**

$$p(\mathbf{W}|\text{Data}) = p(\mathbf{W})p(\text{Data}|\mathbf{W})/p(\text{Data})$$

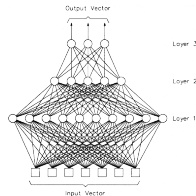
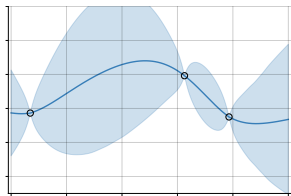
**Predictive Dist.**

$$p(y|\text{Data}, x) = \int p(y|\mathbf{W}, x) p(\mathbf{W}|\text{Data}) d\mathbf{W}$$

**Challenges: The model should be non-parametric (the world is complicated) and computing  $p(\text{Data})$  is intractable!**



# Fitting a Model to the Data



$$h_j(\mathbf{x}) = \tanh\left(\sum_{i=1}^I x_i w_{ji}\right)$$

$$f(\mathbf{x}) = \sum_{j=1}^H v_j h_j(\mathbf{x})$$

Posterior Dist.

$$p(\mathbf{W}|\text{Data}) = p(\mathbf{W})p(\text{Data}|\mathbf{W})/p(\text{Data})$$

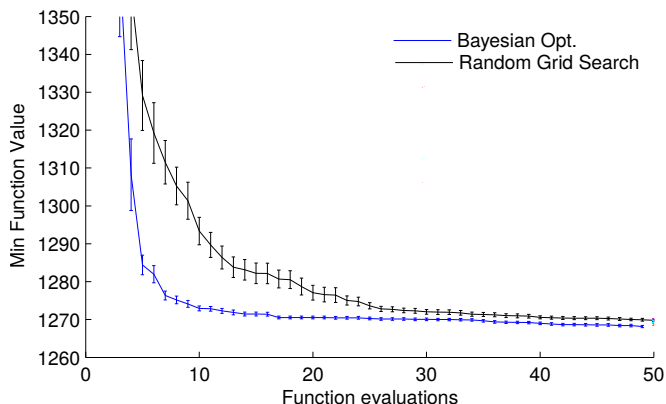
Predictive Dist.

$$p(y|\text{Data}, x) = \int p(y|\mathbf{W}, x) p(\mathbf{W}|\text{Data}) d\mathbf{W}$$

**Challenges: The model should be non-parametric (the world is complicated) and computing  $p(\text{Data})$  is intractable!**

**Solved by setting  $p(\mathbf{W}) = \prod_{ij} \mathcal{N}(w_{ji}|0, \sigma^2 H^{-1})$  and letting  $H \rightarrow \infty$ !**

# Bayesian Optimization vs. Uniform Exploration



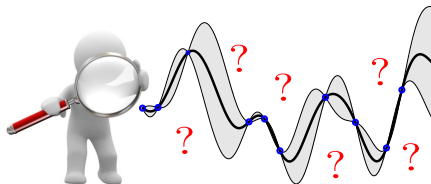
Tuning LDA on a collection of Wikipedia articles (Snoek *et al.*, 2012).

# Using the GP Uncertainty in Optimization

Where to evaluate **next**?

# Using the GP Uncertainty in Optimization

Where to evaluate **next**?



# Using the GP Uncertainty in Optimization

Where to evaluate **next**?



- **Exploration:** seek places with high variance.

# Using the GP Uncertainty in Optimization

Where to evaluate **next**?



- ▶ **Exploration:** seek places with high variance.
- ▶ **Exploitation:** seek places with low mean.

# Using the GP Uncertainty in Optimization

Where to evaluate **next**?



- ▶ **Exploration:** seek places with high variance.
- ▶ **Exploitation:** seek places with low mean.

**The acquisition function balances these two, to choose in an intelligent way the next evaluation point!**

# Using the GP Uncertainty in Optimization

Where to evaluate **next**?



- **Exploration:** seek places with high variance.
- **Exploitation:** seek places with low mean.

**The acquisition function balances these two, to choose in an intelligent way the next evaluation point!**

$$\alpha(x) = \mathbb{E}_{p(y^*|\mathcal{D}_N, x)} [U(y^*|x, \mathcal{D}_N)]$$



## Some Acquisition Functions

Let  $\nu = \min\{y_1, \dots, y_N\}$  and  $\gamma(\mathbf{x}) = \frac{\nu - \mu(\mathbf{x})}{\sigma(\mathbf{x})}$ .

## Some Acquisition Functions

Let  $\nu = \min\{y_1, \dots, y_N\}$  and  $\gamma(x) = \frac{\nu - \mu(x)}{\sigma(x)}$ .

► **Probability of Improvement:**

$$U(y^* | \mathcal{D}_N, x) = \mathbb{I}(y_* < \nu), \quad \alpha(x) = \Phi(\gamma(x))$$

## Some Acquisition Functions

Let  $\nu = \min\{y_1, \dots, y_N\}$  and  $\gamma(x) = \frac{\nu - \mu(x)}{\sigma(x)}$ .

► **Probability of Improvement:**

$$U(y^* | \mathcal{D}_N, x) = \mathbb{I}(y_* < \nu), \quad \alpha(x) = \Phi(\gamma(x))$$

► **Expected Improvement:**

$$U(y^* | \mathcal{D}_N, x) = \max(0, \nu - y^*), \quad \alpha(x) = \sigma(x) (\gamma(x)\Phi(\gamma(x)) + \phi(\gamma(x)))$$

## Some Acquisition Functions

Let  $\nu = \min\{y_1, \dots, y_N\}$  and  $\gamma(x) = \frac{\nu - \mu(x)}{\sigma(x)}$ .

► **Probability of Improvement:**

$$U(y^* | \mathcal{D}_N, x) = \mathbb{I}(y_* < \nu), \quad \alpha(x) = \Phi(\gamma(x))$$

► **Expected Improvement:**

$$U(y^* | \mathcal{D}_N, x) = \max(0, \nu - y^*), \quad \alpha(x) = \sigma(x) (\gamma(x) \Phi(\gamma(x)) + \phi(\gamma(x)))$$

► **Lower Confidence Bound:**

$$\alpha(x) = -(\mu(x) - \kappa\sigma(x))$$

## Some Acquisition Functions

Let  $\nu = \min\{y_1, \dots, y_N\}$  and  $\gamma(x) = \frac{\nu - \mu(x)}{\sigma(x)}$ .

► **Probability of Improvement:**

$$U(y^* | \mathcal{D}_N, x) = \mathbb{I}(y_* < \nu), \quad \alpha(x) = \Phi(\gamma(x))$$

► **Expected Improvement:**

$$U(y^* | \mathcal{D}_N, x) = \max(0, \nu - y^*), \quad \alpha(x) = \sigma(x) (\gamma(x)\Phi(\gamma(x)) + \phi(\gamma(x)))$$

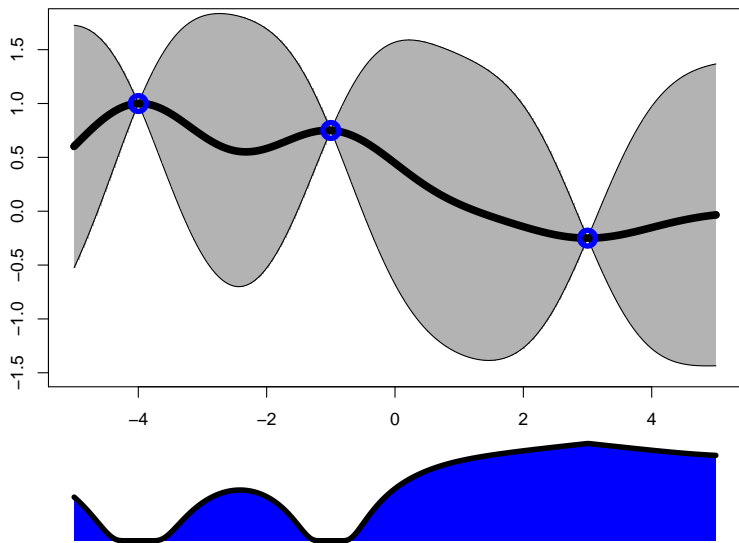
► **Lower Confidence Bound:**

$$\alpha(x) = -(\mu(x) - \kappa\sigma(x))$$

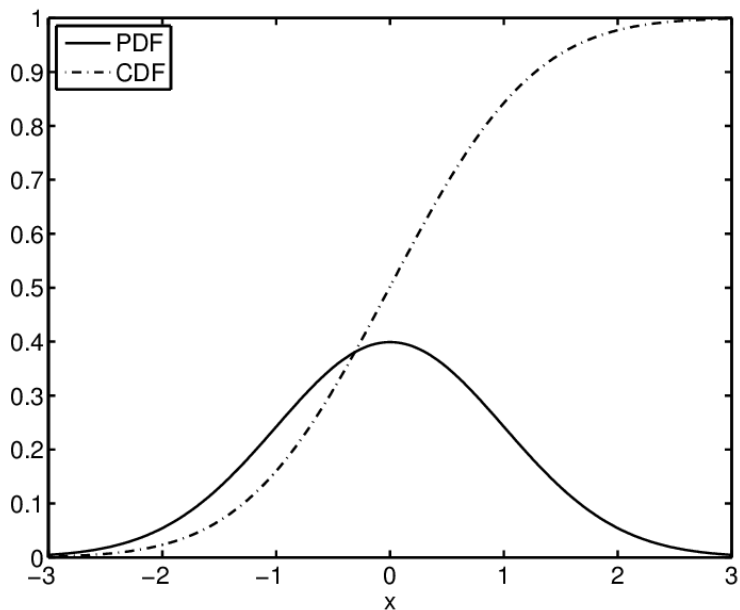
► **Entropy Search:**

$$U(y^* | \mathcal{D}_N, x) = H[p(x_{\min} | \mathcal{D}_N)] - H[p(x_{\min} | \mathcal{D}_N \cup \{x, y^*\})]$$

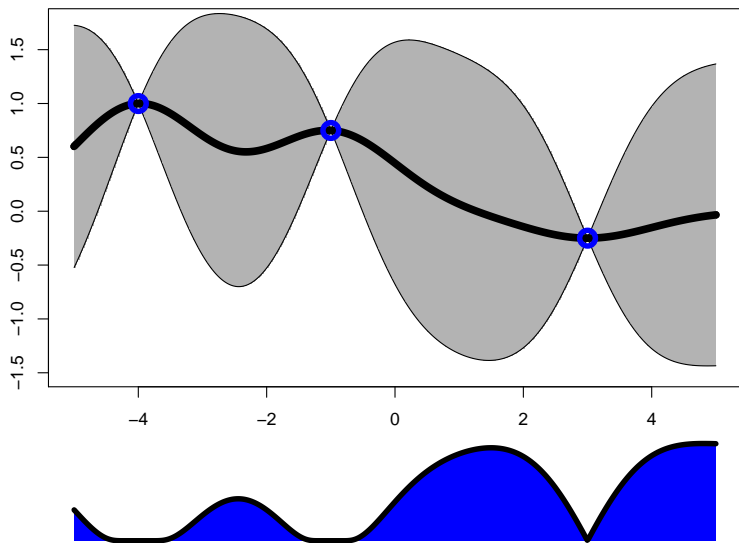
## Some Acquisition Functions: Prob. Improvement



## Prob. Improvement: CDF

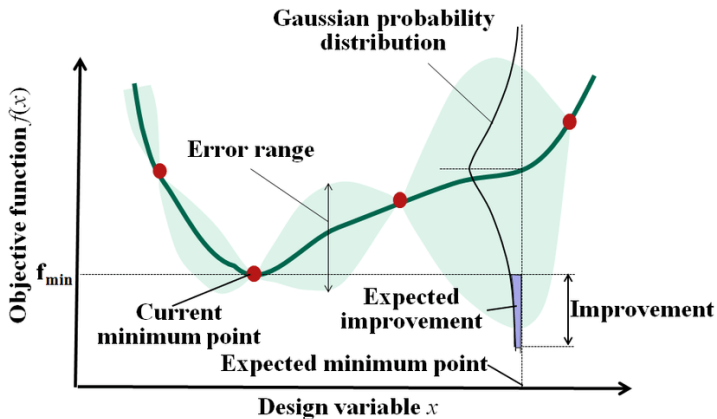


## Some Acquisition Functions: Exp. Improvement

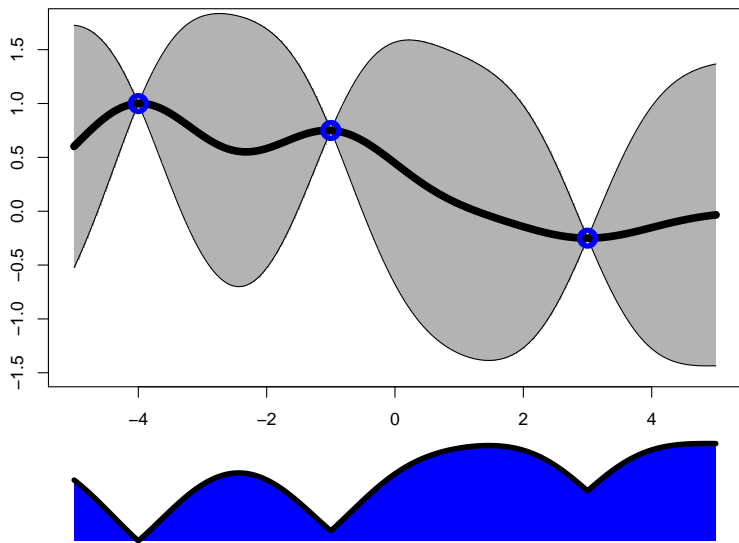




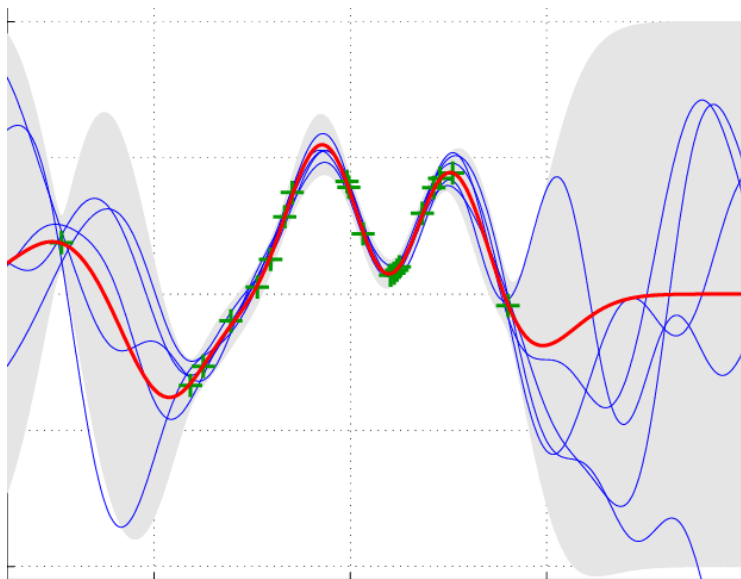
## Expected Improvement: intuition



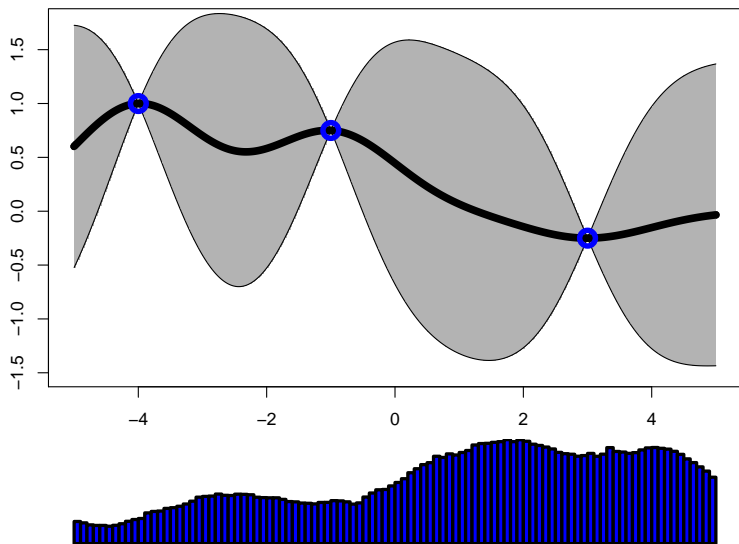
## Some Acquisition Functions: Lower Conf. Bound

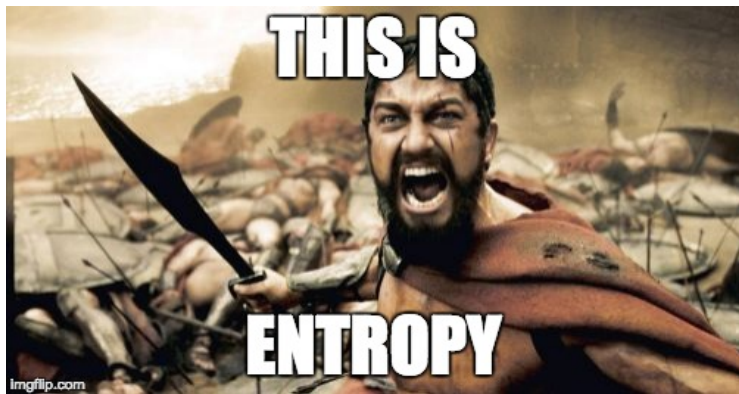


$\kappa=0$ . Use only GP samples and optimize to suggest!



## Some Acquisition Functions: Entropy Search





# Information-based Approach

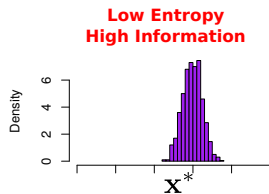
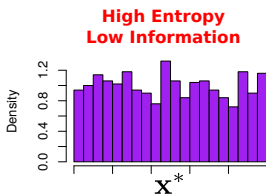
The minimizer,  $x^*$ , can be modelled as a **random variable**!

**Information** is measured by the **entropy** of  $p(x^*|\mathcal{D})$ .

# Information-based Approach

The minimizer,  $x^*$ , can be modelled as a **random variable**!

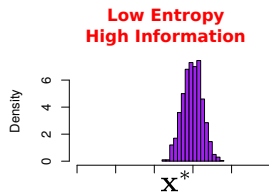
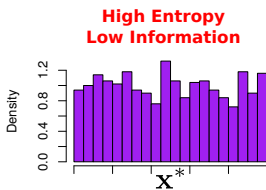
**Information** is measured by the **entropy** of  $p(x^*|\mathcal{D})$ .



# Information-based Approach

The minimizer,  $x^*$ , can be modelled as a **random variable**!

**Information** is measured by the **entropy** of  $p(x^*|\mathcal{D})$ .



The acquisition function is

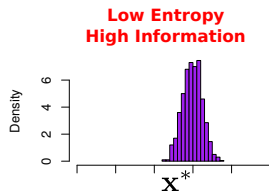
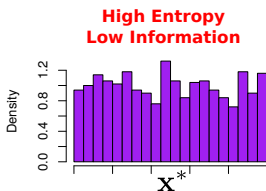
$$\alpha(\mathbf{x}) = H[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] | \mathcal{D}_t, \mathbf{x}] \quad (1)$$



# Information-based Approach

The minimizer,  $x^*$ , can be modelled as a **random variable**!

**Information** is measured by the **entropy** of  $p(x^*|\mathcal{D})$ .



The acquisition function is

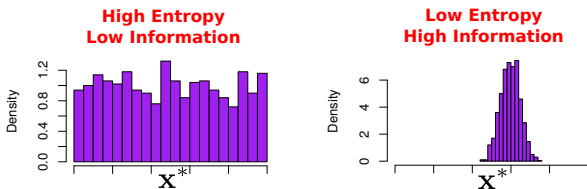
$$\alpha(\mathbf{x}) = \mathbb{E}_{\mathbf{x}^*}[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[\mathbb{E}_{\mathbf{x}^*}[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] | \mathcal{D}_t, \mathbf{x}] \quad (1)$$

How much we know  
about  $\mathbf{x}^*$  now.

# Information-based Approach

The minimizer,  $x^*$ , can be modelled as a **random variable**!

**Information** is measured by the **entropy** of  $p(x^*|\mathcal{D})$ .



The acquisition function is

$$\alpha(\mathbf{x}) = H[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] | \mathcal{D}_t, \mathbf{x}] \quad (1)$$

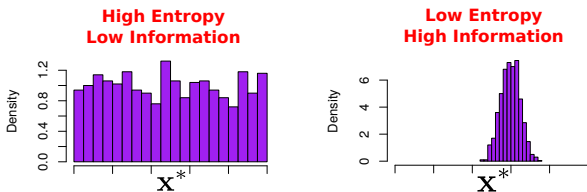
How much we know about  $\mathbf{X}^*$  now.

How much we will know about  $\mathbf{X}^*$  after collecting  $\mathbf{y}$  at  $\mathbf{x}$ .

# Information-based Approach

The minimizer,  $x^*$ , can be modelled as a **random variable**!

**Information** is measured by the **entropy** of  $p(x^*|\mathcal{D})$ .



The acquisition function is

$$\alpha(\mathbf{x}) = H[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] | \mathcal{D}_t, \mathbf{x}] \quad (1)$$

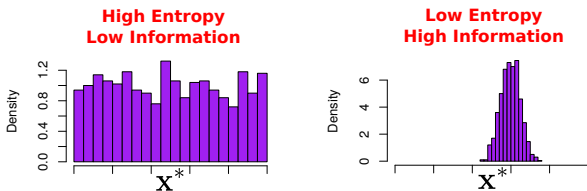
How much we know about  $\mathbf{x}^*$  now.

How much we will know about  $\mathbf{x}^*$  after collecting  $\mathbf{y}$  at  $\mathbf{x}$ .

# Information-based Approach

The minimizer,  $x^*$ , can be modelled as a **random variable**!

**Information** is measured by the **entropy** of  $p(x^*|\mathcal{D})$ .



The acquisition function is

$$\alpha(\mathbf{x}) = H[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] | \mathcal{D}_t, \mathbf{x}] \quad (1)$$

How much we know about  $\mathbf{x}^*$  now.

How much we will know about  $\mathbf{x}^*$  after collecting  $\mathbf{y}$  at  $\mathbf{x}$ .

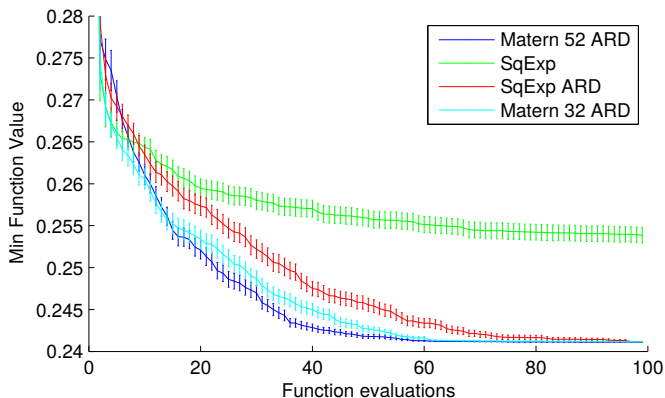
Computing (1) is **very difficult in practice!**

# Bayesian Optimization and Model Selection

- **Covariance function selection:** **critical** to achieve good performance. The default choice for regression (squared exponential) is too smooth. Matérn  $\nu = 5/2$  kernel works better.

# Bayesian Optimization and Model Selection

- **Covariance function selection:** **critical** to achieve good performance. The default choice for regression (squared exponential) is too smooth. Matérn  $\nu = 5/2$  kernel works better.



Structured SVM for protein motif finding (Snoek *et al.*, 2012).

# Bayesian Optimization and Model Selection

- ▶ **Hyper-parameter selection:** with a small number of observations maximizing  $p(y|\theta)$  can give **too confident** uncertainty estimates.

# Bayesian Optimization and Model Selection

- ▶ **Hyper-parameter selection:** with a small number of observations maximizing  $p(y|\theta)$  can give **too confident** uncertainty estimates.
- ▶ **Sampling the hyper-parameters:** computing  $p(\theta|y)$  is **intractable!** Alternative: generate a few samples from  $p(\theta|y)$  using MCMC.



# Bayesian Optimization and Model Selection

- ▶ **Hyper-parameter selection:** with a small number of observations maximizing  $p(y|\theta)$  can give **too confident** uncertainty estimates.
- ▶ **Sampling the hyper-parameters:** computing  $p(\theta|y)$  is **intractable!** Alternative: generate a few samples from  $p(\theta|y)$  using MCMC.

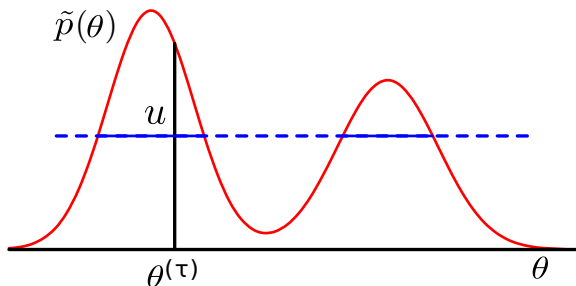
**Slice sampling means no additional hyper-parameters!**

(Neal, 2003)

# Bayesian Optimization and Model Selection

- ▶ **Hyper-parameter selection:** with a small number of observations maximizing  $p(y|\theta)$  can give **too confident** uncertainty estimates.
- ▶ **Sampling the hyper-parameters:** computing  $p(\theta|y)$  is **intractable!** Alternative: generate a few samples from  $p(\theta|y)$  using MCMC.

**Slice sampling means no additional hyper-parameters!**

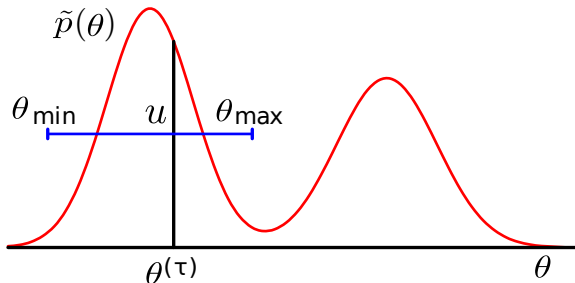


(Neal, 2003)

# Bayesian Optimization and Model Selection

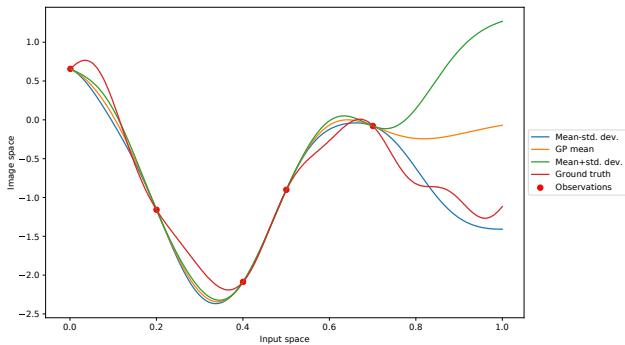
- ▶ **Hyper-parameter selection:** with a small number of observations maximizing  $p(y|\theta)$  can give **too confident** uncertainty estimates.
- ▶ **Sampling the hyper-parameters:** computing  $p(\theta|y)$  is **intractable!** Alternative: generate a few samples from  $p(\theta|y)$  using MCMC.

**Slice sampling means no additional hyper-parameters!**

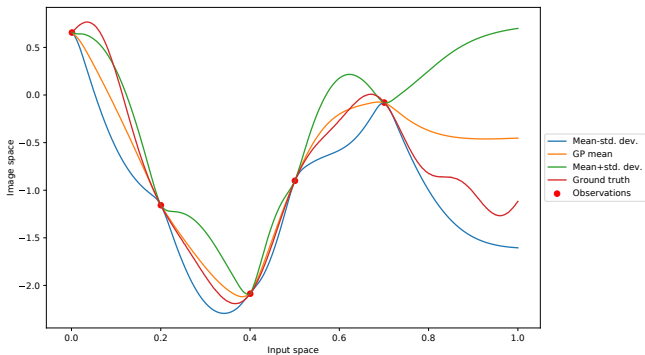


(Neal, 2003)

### GP fitting an unknown function using Maximum Likelihood



### GP fitting an unknown function sampling hyperparameters



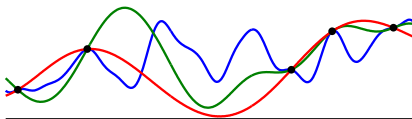
# Integrated Acquisition Function

$$\hat{\alpha}(\mathbf{x}) = \int \alpha(\mathbf{x}; \theta) p(\theta|y) d\theta \approx \frac{1}{K} \sum_{k=1}^K \alpha(\mathbf{x}; \theta^{(k)}) \quad \theta^{(k)} \sim p(\theta|y),$$

# Integrated Acquisition Function

$$\hat{\alpha}(x) = \int \alpha(x; \theta) p(\theta|y) d\theta \approx \frac{1}{K} \sum_{k=1}^K \alpha(x; \theta^{(k)}) \quad \theta^{(k)} \sim p(\theta|y),$$

Posterior samples  
with three different  
length-scales

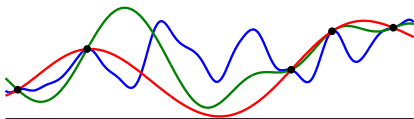


(Snoek *et al.*, 2012)

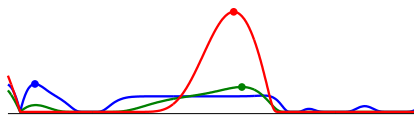
# Integrated Acquisition Function

$$\hat{\alpha}(x) = \int \alpha(x; \theta) p(\theta|y) d\theta \approx \frac{1}{K} \sum_{k=1}^K \alpha(x; \theta^{(k)}) \quad \theta^{(k)} \sim p(\theta|y),$$

Posterior samples  
with three different  
length-scales



Length-scale specific  
expected improvement



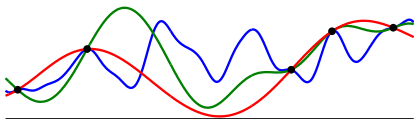
(Snoek *et al.*, 2012)



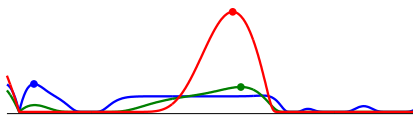
# Integrated Acquisition Function

$$\hat{\alpha}(x) = \int \alpha(x; \theta) p(\theta|y) d\theta \approx \frac{1}{K} \sum_{k=1}^K \alpha(x; \theta^{(k)}) \quad \theta^{(k)} \sim p(\theta|y),$$

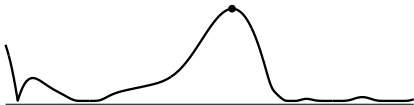
Posterior samples  
with three different  
length-scales



Length-scale specific  
expected improvement

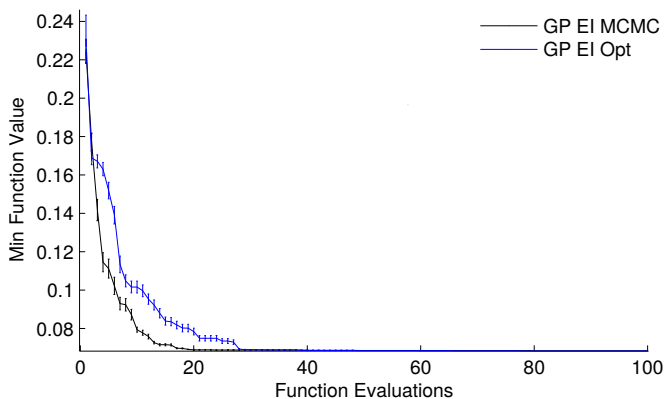


Integrated expected  
improvement



(Snoek *et al.*, 2012)

# MCMC estimation vs. Maximization



Logistic regression on the MNIST (Snoek *et al.*, 2012).

# Cost-sensitive Bayesian Optimization

- ▶ Different inputs may have **different computational costs**, e.g., training a neural network of increasing hidden layers and units.

# Cost-sensitive Bayesian Optimization

- ▶ Different inputs may have **different computational costs**, e.g., training a neural network of increasing hidden layers and units.
- ▶ Better to do **cheap evaluations** before expensive ones!

# Cost-sensitive Bayesian Optimization

- ▶ Different inputs may have **different computational costs**, e.g., training a neural network of increasing hidden layers and units.
- ▶ Better to do **cheap evaluations** before expensive ones!
- ▶ The evaluation costs are **unknown** but they can be **recorded** and then **modeled** with an additional **Gaussian process**.

# Cost-sensitive Bayesian Optimization

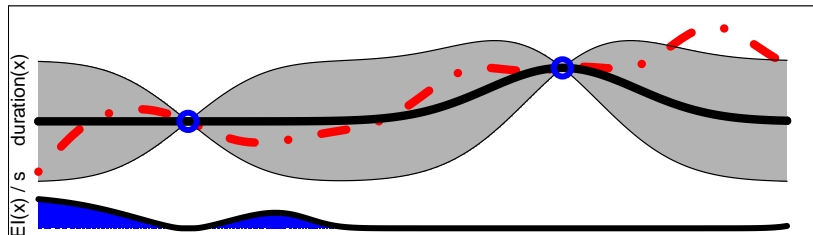
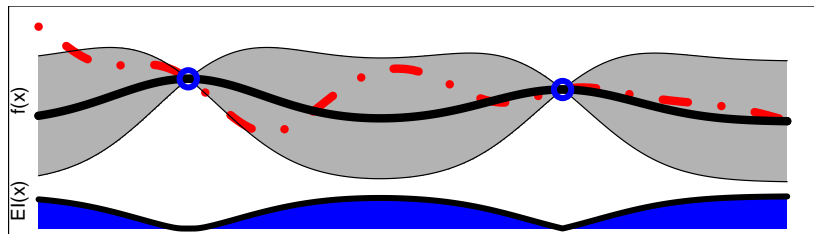
- ▶ Different inputs may have **different computational costs**, e.g., training a neural network of increasing hidden layers and units.
- ▶ Better to do **cheap evaluations** before expensive ones!
- ▶ The evaluation costs are **unknown** but they can be **recorded** and then **modeled** with an additional **Gaussian process**.

## Expected Improvement per-second:

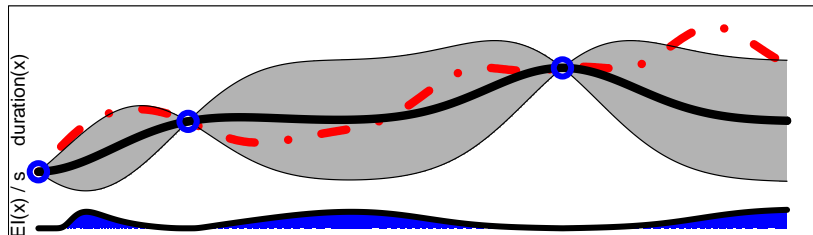
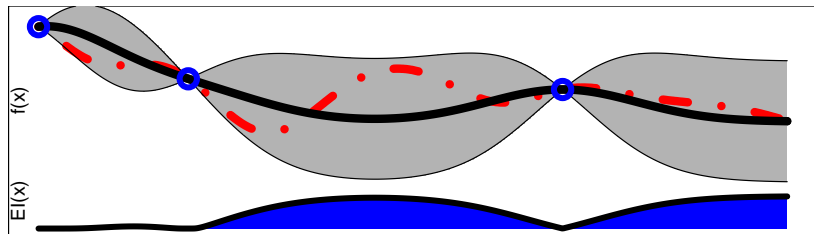
$$\alpha(x) = \frac{\sigma(x) (\gamma(x)\Phi(\gamma(x)) + \phi(\gamma(x)))}{\exp\{\mu_{\log\text{-time}}(x)\}}$$

(Snoek *et al.*, 2012)

# Cost-sensitive Bayesian Optimization

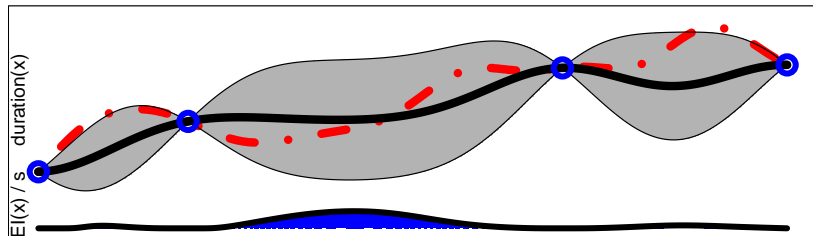
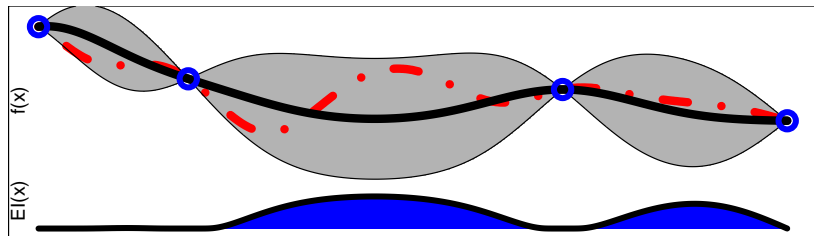


# Cost-sensitive Bayesian Optimization

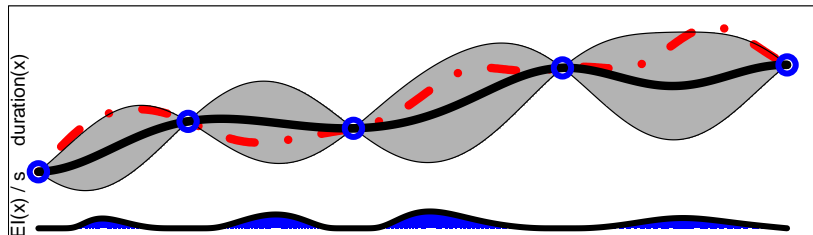
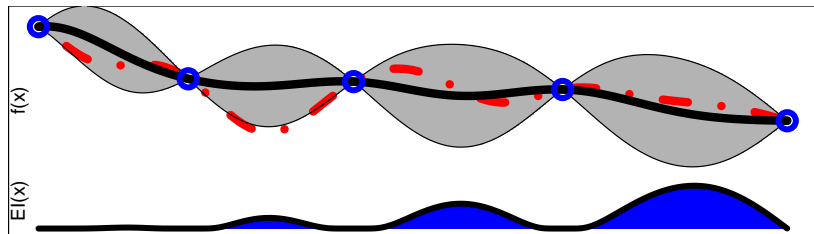




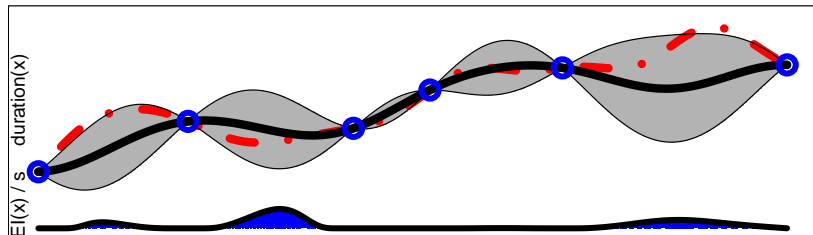
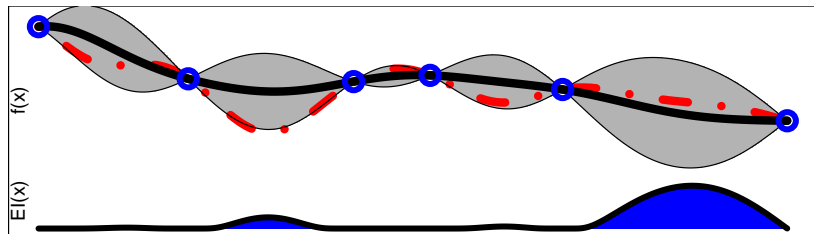
# Cost-sensitive Bayesian Optimization



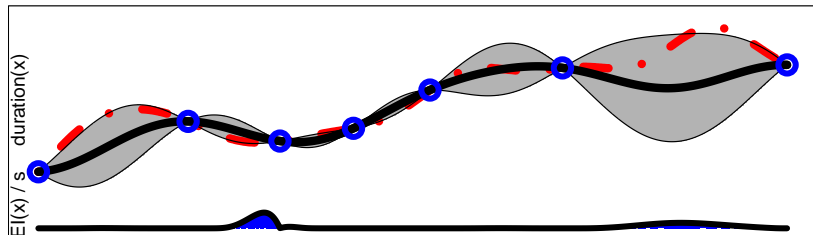
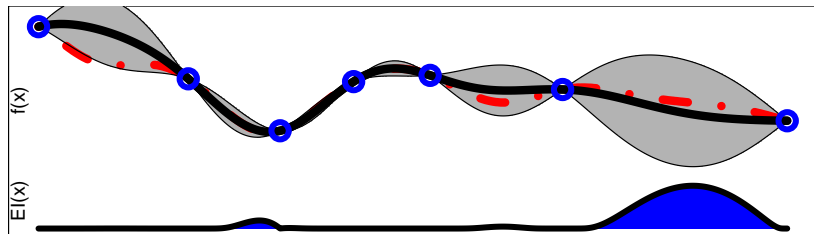
# Cost-sensitive Bayesian Optimization



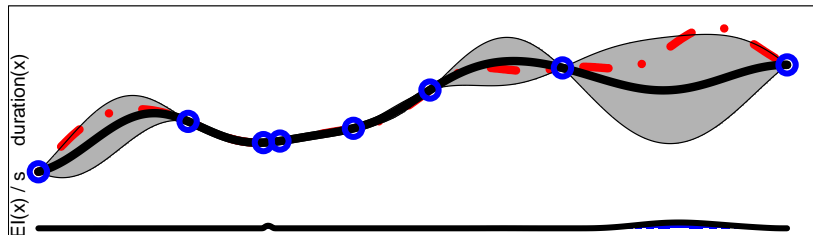
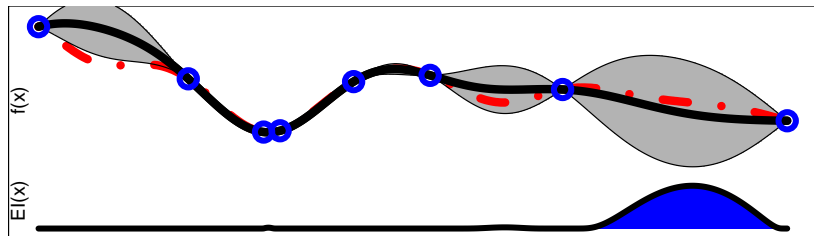
# Cost-sensitive Bayesian Optimization



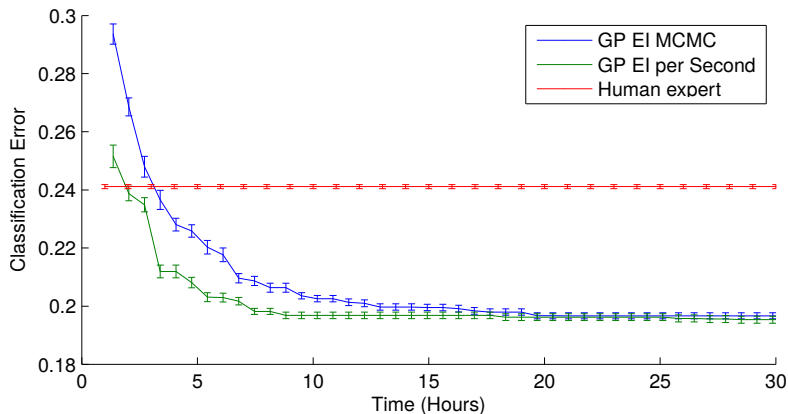
# Cost-sensitive Bayesian Optimization



# Cost-sensitive Bayesian Optimization



# Cost-sensitive Bayesian Optimization



Deep neural network on the CIFAR dataset (Snoek *et al.*, 2012)

# Gaussian processes are not the only probabilistic surrogate model!

1. **Computationally expensive:** GPs scale poorly due to cubic observation complexity.

# Gaussian processes are not the only probabilistic surrogate model!

1. **Computationally expensive:** GPs scale poorly due to cubic observation complexity.
2. **Fixed covariance structure:** GPs assume a fixed structure in data covariance (stationarity), limiting their ability to model complex data patterns.



# Gaussian processes are not the only probabilistic surrogate model!

1. **Computationally expensive:** GPs scale poorly due to cubic observation complexity.
2. **Fixed covariance structure:** GPs assume a fixed structure in data covariance (stationarity), limiting their ability to model complex data patterns.
3. **Difficulties with high dimensionality:** GPs struggle to perform well in high-dimensional search spaces.

# Gaussian processes are not the only probabilistic surrogate model!

1. **Computationally expensive:** GPs scale poorly due to cubic observation complexity.
2. **Fixed covariance structure:** GPs assume a fixed structure in data covariance (stationarity), limiting their ability to model complex data patterns.
3. **Difficulties with high dimensionality:** GPs struggle to perform well in high-dimensional search spaces.
4. **Random Forests:** We use its empirical predictive distribution. Perform well in hierarchical spaces, are robust to outliers and scale well.

# Gaussian processes are not the only probabilistic surrogate model!

1. **Computationally expensive:** GPs scale poorly due to cubic observation complexity.
2. **Fixed covariance structure:** GPs assume a fixed structure in data covariance (stationarity), limiting their ability to model complex data patterns.
3. **Difficulties with high dimensionality:** GPs struggle to perform well in high-dimensional search spaces.
4. **Random Forests:** We use its empirical predictive distribution. Perform well in hierarchical spaces, are robust to outliers and scale well.
5. **Bayesian neural networks:** Able to model complex patterns and scalable.

# Gaussian processes are not the only probabilistic surrogate model!

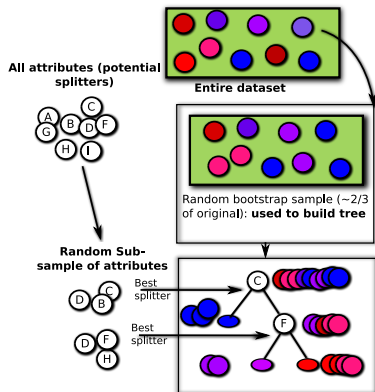
1. **Computationally expensive:** GPs scale poorly due to cubic observation complexity.
2. **Fixed covariance structure:** GPs assume a fixed structure in data covariance (stationarity), limiting their ability to model complex data patterns.
3. **Difficulties with high dimensionality:** GPs struggle to perform well in high-dimensional search spaces.
4. **Random Forests:** We use its empirical predictive distribution. Perform well in hierarchical spaces, are robust to outliers and scale well.
5. **Bayesian neural networks:** Able to model complex patterns and scalable.
6. **Deep Gaussian Process:** Increased expressivity, advantages of GPs.

## Other Models: Random Forest

Ensemble method where the predictors are random regression trees trained on random subsamples of the data.

# Other Models: Random Forest

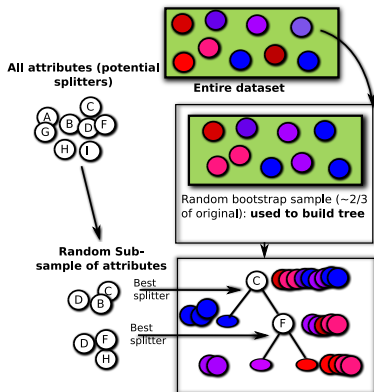
Ensemble method where the predictors are random regression trees trained on random subsamples of the data.



- Trees are grown on different bootstrap samples of the data.

# Other Models: Random Forest

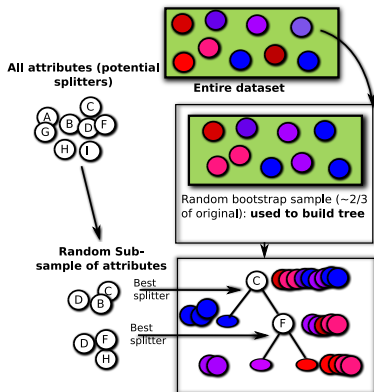
Ensemble method where the predictors are random regression trees trained on random subsamples of the data.



- ▶ Trees are grown on different bootstrap samples of the data.
- ▶ At each node the best splitter is chosen randomly.

# Other Models: Random Forest

Ensemble method where the predictors are random regression trees trained on random subsamples of the data.

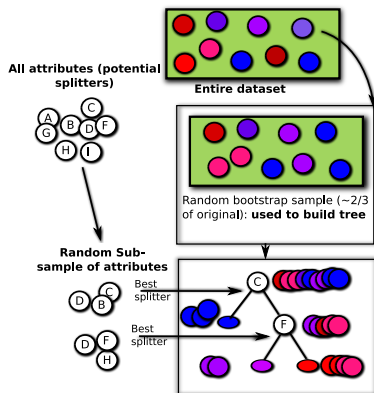


- ▶ Trees are grown on different bootstrap samples of the data.
- ▶ At each node the best splitter is chosen randomly.
- ▶ Leaf nodes predict the average value of the points reaching that node.



# Other Models: Random Forest

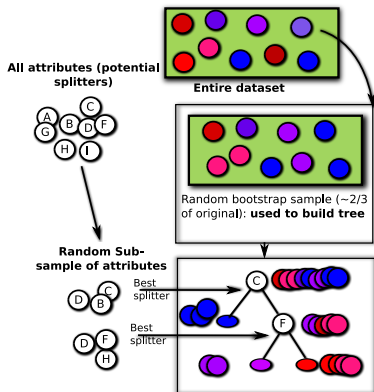
Ensemble method where the predictors are random regression trees trained on random subsamples of the data.



- ▶ Trees are grown on different bootstrap samples of the data.
- ▶ At each node the best splitter is chosen randomly.
- ▶ Leaf nodes predict the average value of the points reaching that node.
- ▶ This guarantees that each tree is slightly different.

# Other Models: Random Forest

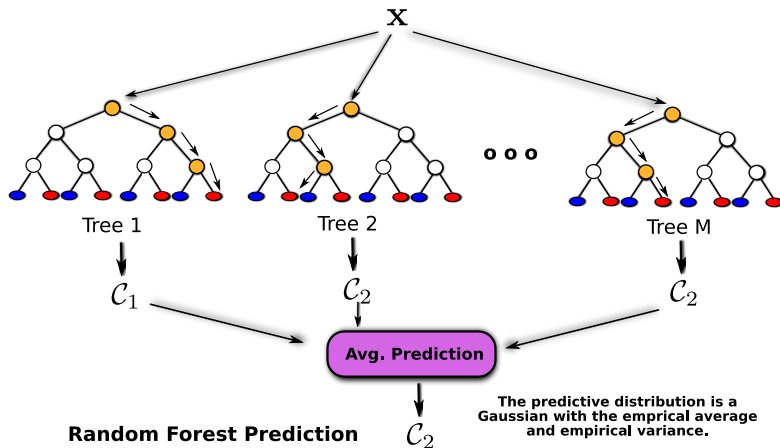
Ensemble method where the predictors are random regression trees trained on random subsamples of the data.



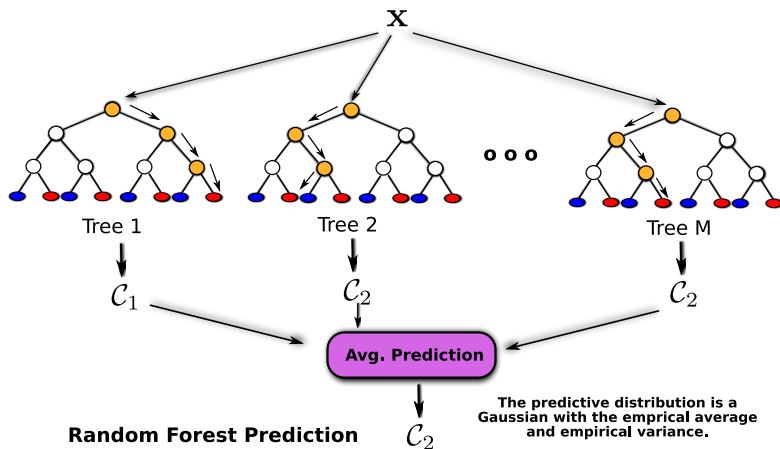
- ▶ Trees are grown on different bootstrap samples of the data.
- ▶ At each node the best splitter is chosen randomly.
- ▶ Leaf nodes predict the average value of the points reaching that node.
- ▶ This guarantees that each tree is slightly different.

**Very cheap to compute and massively paralelizable!**

# Random Forest: Predictive Distribution



# Random Forest: Predictive Distribution

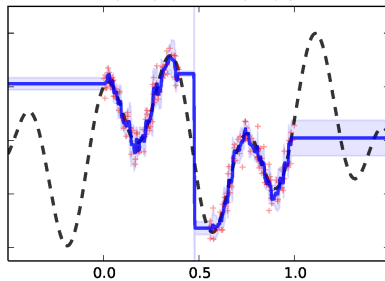


$$p(f^*|\mathcal{D}_n) = \mathcal{N}(f^*|\bar{\mu}, \bar{\nu}^2)$$

(Hutter et al., 2011)

# Random Forest in Practice

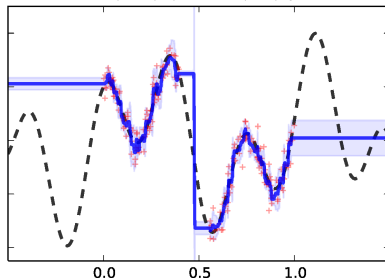
Random Forest



(Shahriari et al., 2016)

# Random Forest in Practice

Random Forest

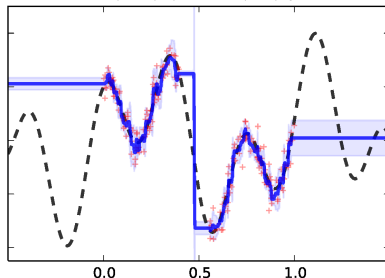


(Shahriari et al., 2016)

- Allows for a lot of evaluations (good when the objective is cheap).

# Random Forest in Practice

Random Forest

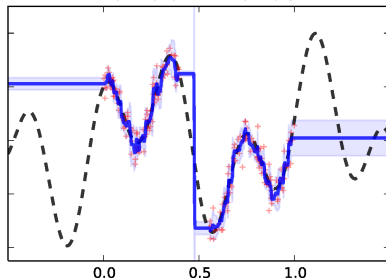


(Shahriari et al., 2016)

- ▶ Allows for a lot of evaluations (good when the objective is cheap).
- ▶ Too confident intervals in far away from the data regions.

# Random Forest in Practice

Random Forest



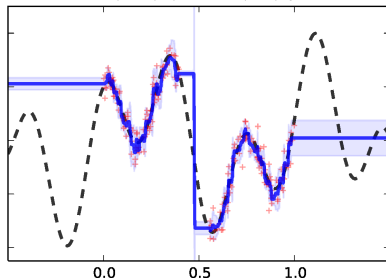
(Shahriari et al., 2016)

- ▶ Allows for a lot of evaluations (good when the objective is cheap).
- ▶ Too confident intervals in far away from the data regions.
- ▶ Conflicting predictions can cause the variance to be too high.



# Random Forest in Practice

Random Forest

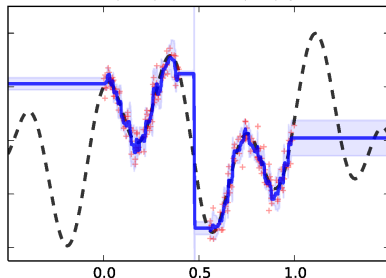


(Shahriari et al., 2016)

- ▶ Allows for a lot of evaluations (good when the objective is cheap).
- ▶ Too confident intervals in far away from the data regions.
- ▶ Conflicting predictions can cause the variance to be too high.
- ▶ Discontinuous: Difficult to optimize the acquisition function.

# Random Forest in Practice

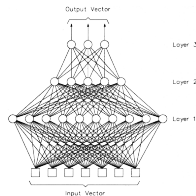
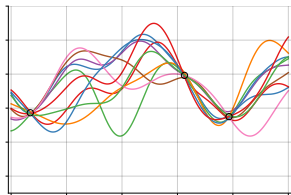
Random Forest



(Shahriari et al., 2016)

- ▶ Allows for a lot of evaluations (good when the objective is cheap).
- ▶ Too confident intervals in far away from the data regions.
- ▶ Conflicting predictions can cause the variance to be too high.
- ▶ Discontinuous: Difficult to optimize the acquisition function.
- ▶ No parameters to tune.

## Other Models: Bayesian Neural Networks

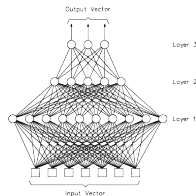
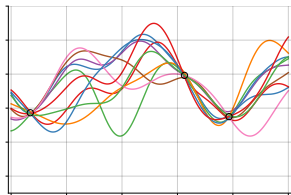


$$h_j(\mathbf{x}) = \tanh \left( \sum_{i=1}^I x_i w_{ji} \right)$$

$$f(\mathbf{x}) = \sum_{j=1}^H v_j h_j(\mathbf{x})$$

- Neural networks scale well to the training data (linear cost).

# Other Models: Bayesian Neural Networks

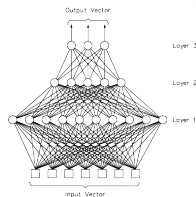
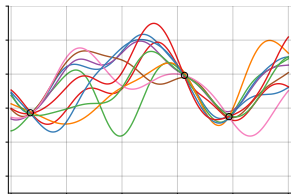


$$h_j(\mathbf{x}) = \tanh \left( \sum_{i=1}^I x_i w_{ji} \right)$$

$$f(\mathbf{x}) = \sum_{j=1}^H v_j h_j(\mathbf{x})$$

- ▶ Neural networks scale well to the training data (linear cost).
- ▶ Trained very fast on GPUs.

# Other Models: Bayesian Neural Networks

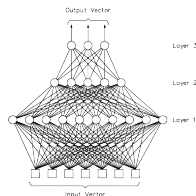
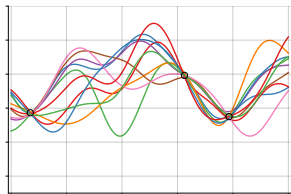


$$h_j(\mathbf{x}) = \tanh\left(\sum_{i=1}^I x_i w_{ji}\right)$$

$$f(\mathbf{x}) = \sum_{j=1}^H v_j h_j(\mathbf{x})$$

- ▶ Neural networks scale well to the training data (linear cost).
- ▶ Trained very fast on GPUs.
- ▶ State of the art prediction results.

# Other Models: Bayesian Neural Networks



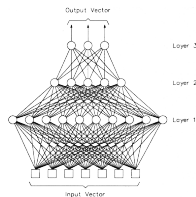
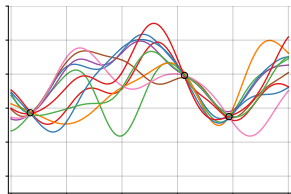
$$h_j(\mathbf{x}) = \tanh \left( \sum_{i=1}^I x_i w_{ji} \right)$$

$$f(\mathbf{x}) = \sum_{j=1}^H v_j h_j(\mathbf{x})$$

- ▶ Neural networks scale well to the training data (linear cost).
- ▶ Trained very fast on GPUs.
- ▶ State of the art prediction results.

They are an alternative to GPs to allow for a large number observations!

# Other Models: Bayesian Neural Networks



$$h_j(\mathbf{x}) = \tanh\left(\sum_{i=1}^I x_i w_{ji}\right)$$

$$f(\mathbf{x}) = \sum_{j=1}^H v_j h_j(\mathbf{x})$$

- ▶ Neural networks scale well to the training data (linear cost).
- ▶ Trained very fast on GPUs.
- ▶ State of the art prediction results.

They are an alternative to GPs to allow for a large number observations!

**The posterior distribution of the networks weights  $\mathbf{W}$  is intractable!**

# Bayesian Neural Networks: Predictive Distribution

Several techniques considered to approximate the predictive distribution:



# Bayesian Neural Networks: Predictive Distribution

Several techniques considered to approximate the predictive distribution:

- ▶ Markov Chain Monte Carlo methods.

# Bayesian Neural Networks: Predictive Distribution

Several techniques considered to approximate the predictive distribution:

- ▶ Markov Chain Monte Carlo methods.
- ▶ Variational Inference.

# Bayesian Neural Networks: Predictive Distribution

Several techniques considered to approximate the predictive distribution:

- ▶ Markov Chain Monte Carlo methods.
- ▶ Variational Inference.
- ▶ Expectation Propagation.

# Bayesian Neural Networks: Predictive Distribution

Several techniques considered to approximate the predictive distribution:

- ▶ Markov Chain Monte Carlo methods.
- ▶ Variational Inference.
- ▶ Expectation Propagation.
- ▶ Reinterpretations of dropout.

# Bayesian Neural Networks: Predictive Distribution

Several techniques considered to approximate the predictive distribution:

- ▶ Markov Chain Monte Carlo methods.
- ▶ Variational Inference.
- ▶ Expectation Propagation.
- ▶ Reinterpretations of dropout.
- ▶ Point estimates and Bayesian linear-models in the last layer.

# Bayesian Neural Networks: Predictive Distribution

Several techniques considered to approximate the predictive distribution:

- ▶ Markov Chain Monte Carlo methods.
- ▶ Variational Inference.
- ▶ Expectation Propagation.
- ▶ Reinterpretations of dropout.
- ▶ Point estimates and Bayesian linear-models in the last layer.

**Trade-off between accuracy of the predictive distribution and scalability! Still a lot of research going on!**

# Activity: Bayesian optimization of a ML algorithm

Use a generative AI to generate a BO code of a ML hyperparameter problem.

- ▶ Try to execute the code in your machine and fix it.
- ▶ Try to interpret everything.
- ▶ Tune it and improve the quality of the error estimator.
- ▶ Try to improve it with additional acquisition functions.
- ▶ Try to obtain the best score in the titanic problem with BO-ML!

# Software for Bayesian Optimization

Many of the methods described are implemented into **BOTorch** using Python.

<https://botorch.org/>





# Software for Bayesian Optimization

Many of the methods described are implemented into **BOTorch** using Python.

<https://botorch.org/>



BOTorch's super-nice features:

1. **Modularity**: Plugin new models, acquisition functions and optimizers.

# Software for Bayesian Optimization

Many of the methods described are implemented into **BOTorch** using Python.

<https://botorch.org/>



BOTorch's super-nice features:

1. **Modularity**: Plugin new models, acquisition functions and optimizers.
2. **PyTorch**: Easily integrate neural network modules. Native GPU and Autograd support.

# Software for Bayesian Optimization

Many of the methods described are implemented into **BOTorch** using Python.

<https://botorch.org/>



BOTorch's super-nice features:

1. **Modularity**: Plugin new models, acquisition functions and optimizers.
2. **PyTorch**: Easily integrate neural network modules. Native GPU and Autograd support.
3. **Scalability**: Support for scalable GPs via GPyTorch.

# Software for Bayesian Optimization

Many of the methods described are implemented into **BOTorch** using Python.

<https://botorch.org/>



BOTorch's super-nice features:

1. **Modularity**: Plugin new models, acquisition functions and optimizers.
2. **PyTorch**: Easily integrate neural network modules. Native GPU and Autograd support.
3. **Scalability**: Support for scalable GPs via GPyTorch.
4. **Advanced scenarios**: Advanced research coded into tutorials about complex BO scenarios.

# Software for Bayesian Optimization

Many of the methods described are implemented into **BOTorch** using Python.

<https://botorch.org/>



BOTorch's super-nice features:

1. **Modularity**: Plugin new models, acquisition functions and optimizers.
2. **PyTorch**: Easily integrate neural network modules. Native GPU and Autograd support.
3. **Scalability**: Support for scalable GPs via GPyTorch.
4. **Advanced scenarios**: Advanced research coded into tutorials about complex BO scenarios.

**Other tools**: SMAC3 (Python-RFs), GPyOpt (Python3), Spearmint (Python2.7-sampling), mlrMBO (R).

## Time to practice!

1. Bayesian optimization of a benchmark optimization function.
2. Bayesian optimization of the hyper-parameters of a machine learning model.
3. Bayesian optimization of the hyper-parameters of a deep reinforcement learning algorithm.



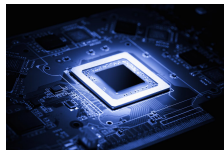
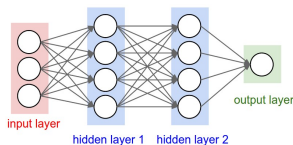
## Second Session: Advanced Bayesian optimization.

# **(Parallel) Multi-objective Bayesian optimization with constraints.**



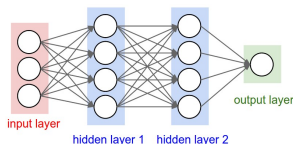
# Several Objectives and Constraints

Optimal design of **hardware accelerator** for neural network predictions.



# Several Objectives and Constraints

Optimal design of **hardware accelerator** for neural network predictions.

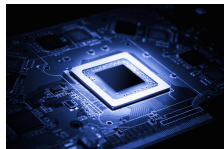
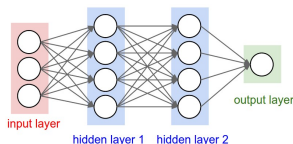


## Goals:

- ▶ Minimize **prediction error**.
- ▶ Minimize **prediction time**.

# Several Objectives and Constraints

Optimal design of **hardware accelerator** for neural network predictions.



## Goals:

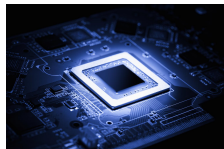
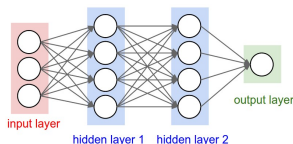
- ▶ Minimize **prediction error**.
- ▶ Minimize **prediction time**.

## Constrained to:

- ▶ **Chip area** below a value.
- ▶ **Power consumption** below a level.

# Several Objectives and Constraints

Optimal design of **hardware accelerator** for neural network predictions.

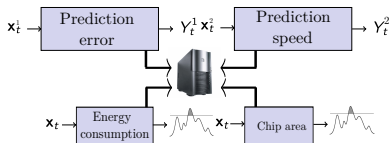


## Goals:

- ▶ Minimize **prediction error**.
- ▶ Minimize **prediction time**.

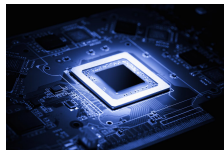
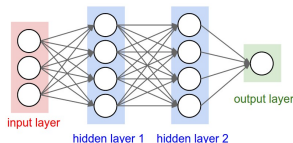
## Constrained to:

- ▶ **Chip area** below a value.
- ▶ **Power consumption** below a level.



# Several Objectives and Constraints

Optimal design of **hardware accelerator** for neural network predictions.

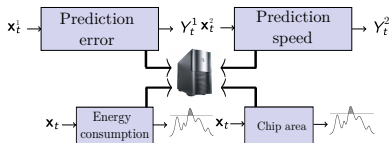


## Goals:

- ▶ Minimize **prediction error**.
- ▶ Minimize **prediction time**.

## Constrained to:

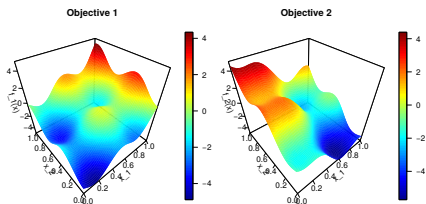
- ▶ **Chip area** below a value.
- ▶ **Power consumption** below a level.



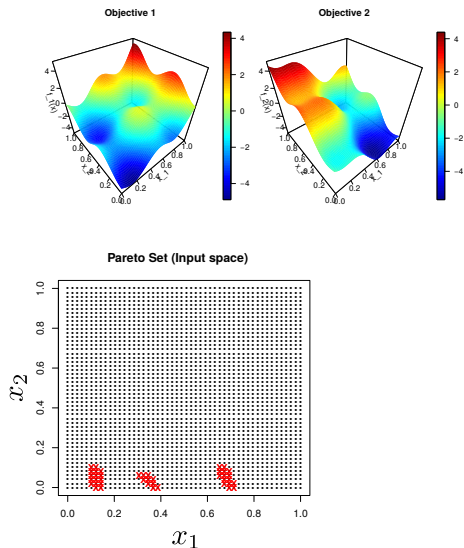
## Challenges:

- ▶ **Complicated** constraints.
- ▶ **Conflicting** objectives.

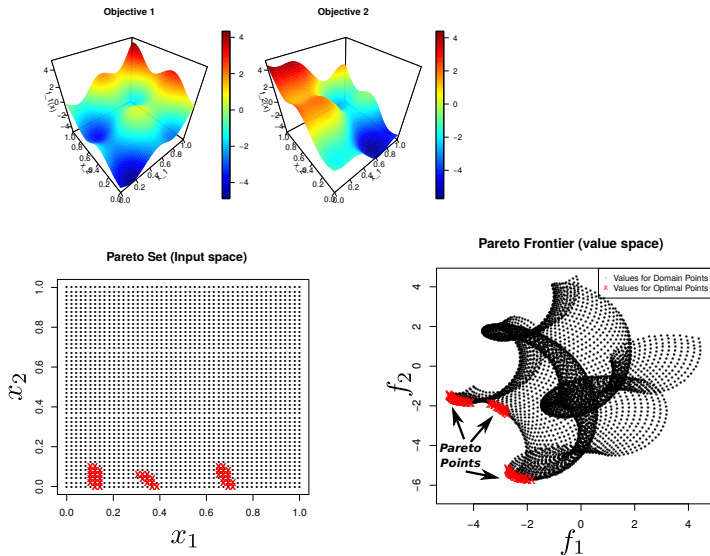
# Constrained Multi-Objective Optimization



# Constrained Multi-Objective Optimization

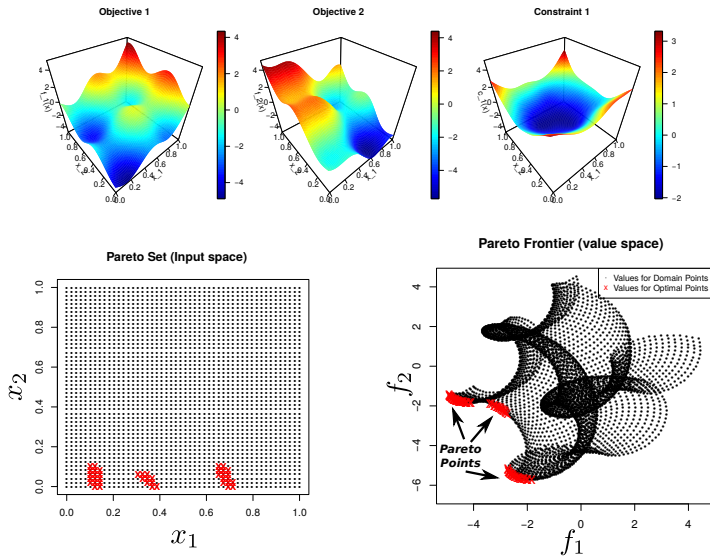


# Constrained Multi-Objective Optimization

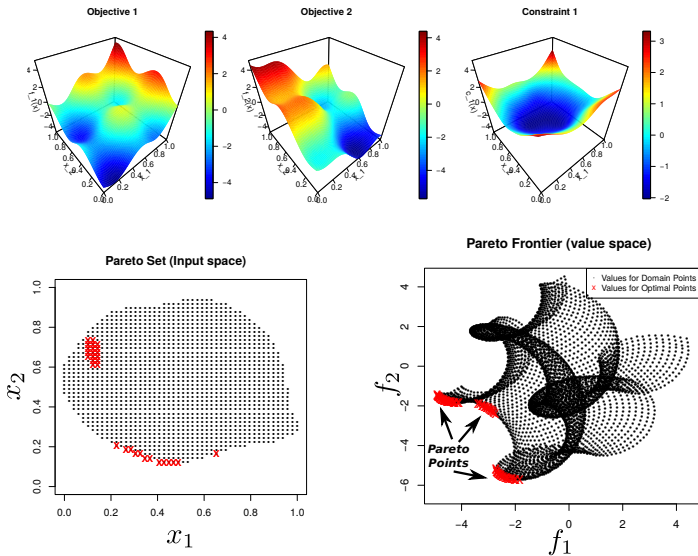




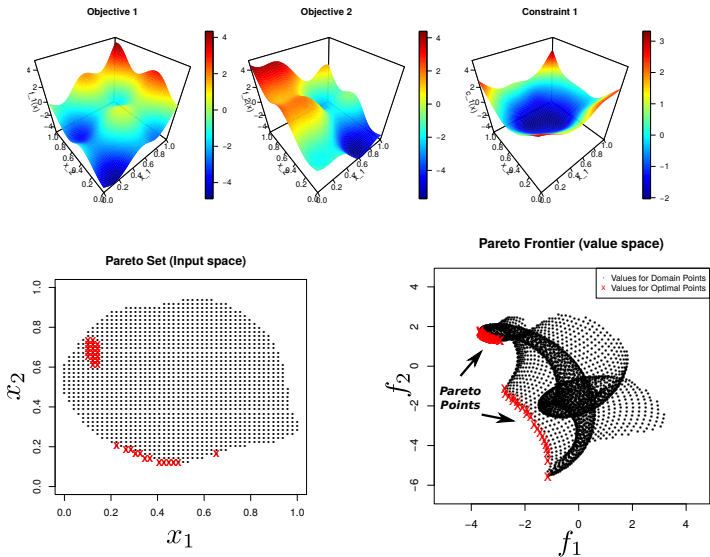
# Constrained Multi-Objective Optimization



# Constrained Multi-Objective Optimization



# Constrained Multi-Objective Optimization



# Bayesian Optimization Methods

**Additional challenges when dealing with several black-boxes.**

# Bayesian Optimization Methods

## Additional challenges when dealing with several black-boxes.

- ▶ Simple approach: evaluate **all** the objectives and constraints at the **same input location**. Expected to be sub-optimal.

# Bayesian Optimization Methods

## Additional challenges when dealing with several black-boxes.

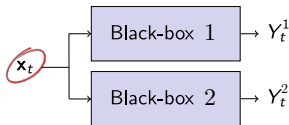
- ▶ Simple approach: evaluate **all** the objectives and constraints at the **same input location**. Expected to be sub-optimal.
- ▶ Advanced approach: make **intelligent** decisions about what **black-box** to **evaluate next** and on **which location**.

# Bayesian Optimization Methods

## Additional challenges when dealing with several black-boxes.

- ▶ Simple approach: evaluate **all** the objectives and constraints at the **same input location**. Expected to be sub-optimal.
- ▶ Advanced approach: make **intelligent** decisions about what **black-box** to **evaluate next** and on **which location**.

## Coupled evaluations

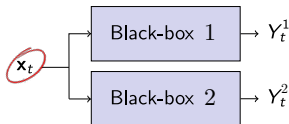


# Bayesian Optimization Methods

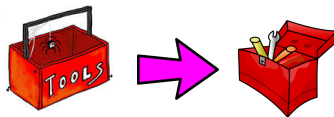
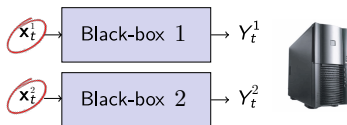
## Additional challenges when dealing with several black-boxes.

- ▶ Simple approach: evaluate **all** the objectives and constraints at the **same input location**. Expected to be sub-optimal.
- ▶ Advanced approach: make **intelligent** decisions about what **black-box** to **evaluate next** and on **which location**.

### Coupled evaluations



### Decoupled evaluations





# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

# Information-based Approach

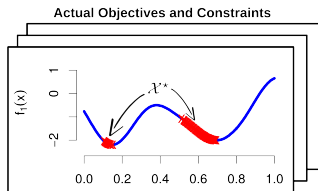
The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .

# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

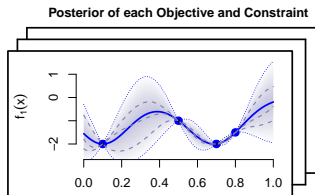
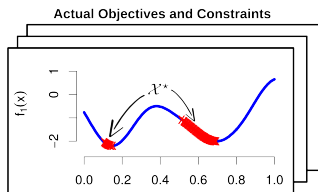
**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

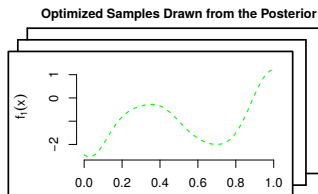
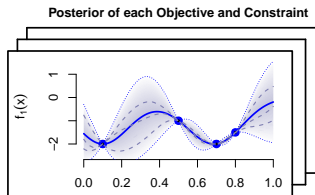
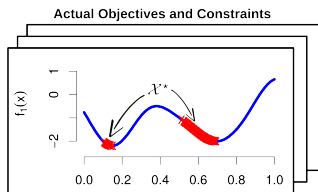
**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

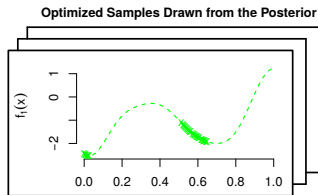
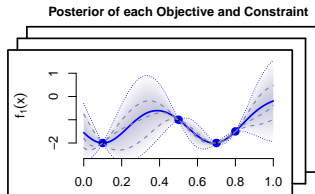
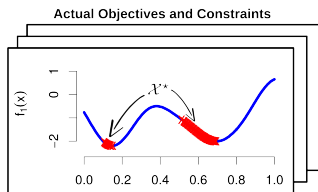
**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

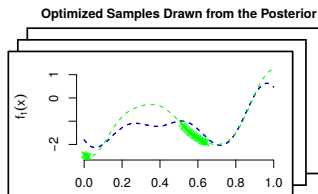
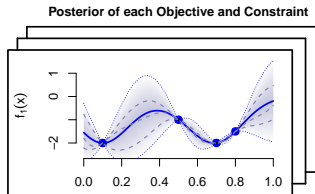
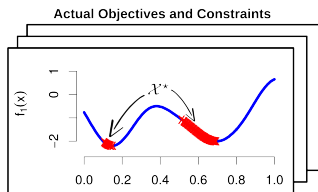
**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

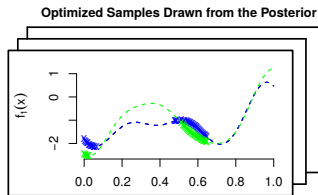
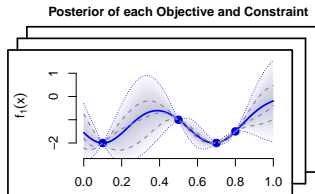
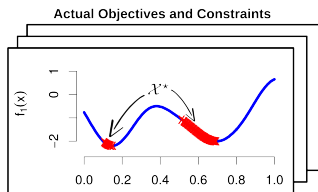
**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .

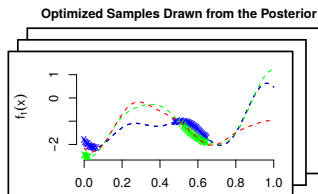
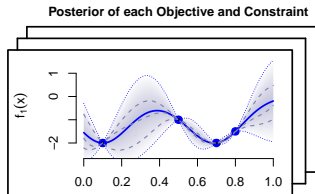
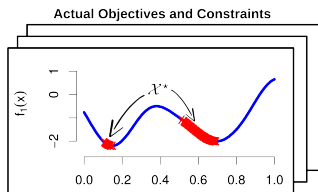




# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

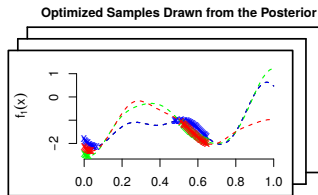
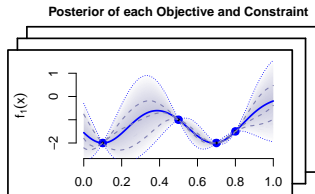
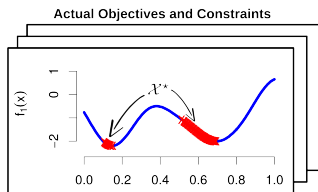
**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

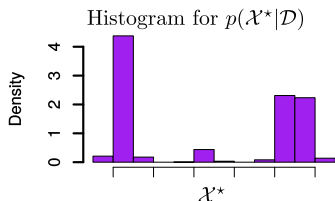
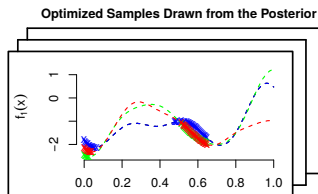
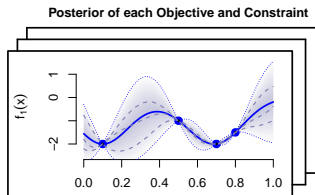
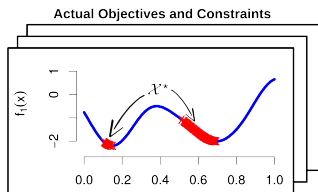
**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



# Information-based Approach

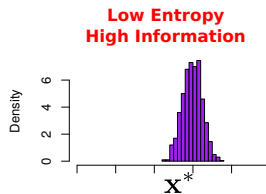
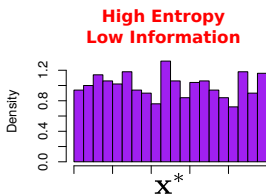
The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .

# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

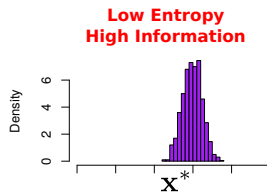
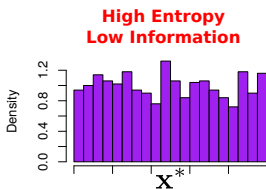
**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



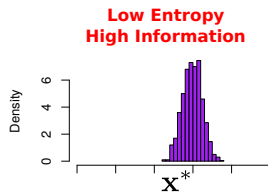
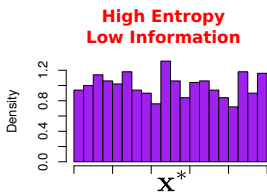
The acquisition function is

$$\alpha(\mathbf{x}) = H[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] | \mathcal{D}_t, \mathbf{x}] \quad (1)$$

# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



The acquisition function is

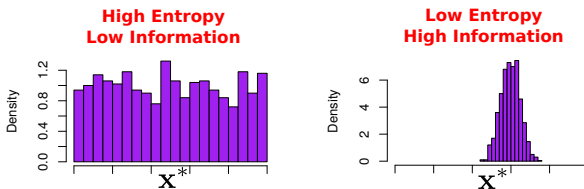
$$\alpha(\mathbf{x}) = \mathbb{H}[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}} \left[ \mathbb{H}[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] \middle| \mathcal{D}_t, \mathbf{x} \right] \quad (1)$$

How much we know  
about  $\mathbf{X}^*$  now.

# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



The acquisition function is

$$\alpha(\mathbf{x}) = H[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] | \mathcal{D}_t, \mathbf{x}] \quad (1)$$

How much we know about  $\mathbf{X}^*$  now.

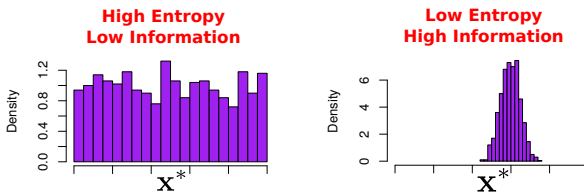
How much we will know about  $\mathbf{X}^*$  after collecting  $\mathbf{y}$  at  $\mathbf{x}$ .



# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



The acquisition function is

$$\alpha(\mathbf{x}) = \mathbb{H}[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[\mathbb{H}[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] | \mathcal{D}_t, \mathbf{x}] \quad (1)$$

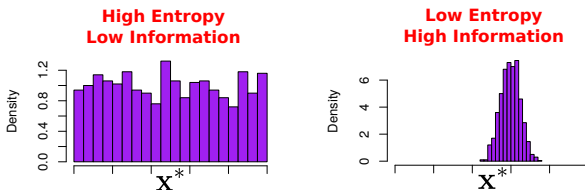
How much we know about  $\mathbf{x}^*$  now.

How much we will know about  $\mathbf{x}^*$  after collecting  $\mathbf{y}$  at  $\mathbf{x}$ .

# Information-based Approach

The Pareto set  $\mathcal{X}^*$  in the feasible space is a **random variable**!

**Information** is measured by the **entropy** of  $p(\mathcal{X}^*|\mathcal{D}_N)$ .



The acquisition function is

$$\alpha(\mathbf{x}) = H[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] | \mathcal{D}_t, \mathbf{x}] \quad (1)$$

How much we know about  $\mathbf{X}^*$  now.

How much we will know about  $\mathbf{X}^*$  after collecting  $\mathbf{y}$  at  $\mathbf{x}$ .

Computing (1) is **very difficult in practice!**

# Predictive Entropy Search (PES)

We **swap**  $y$  and  $\mathcal{X}^*$  to obtain a reformulation of the acquisition function.

# Predictive Entropy Search (PES)

We **swap**  $y$  and  $\mathcal{X}^*$  to obtain a reformulation of the acquisition function.

$$H[\mathcal{X}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathcal{X}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC})$$

# Predictive Entropy Search (PES)

We **swap  $y$  and  $\mathcal{X}^*$**  to obtain a reformulation of the acquisition function.

$$H[\mathcal{X}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathcal{X}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC})$$

$$H[\mathbf{y}|\mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathcal{X}^*}[H[\mathbf{y}|\mathcal{D}_t, \mathbf{x}, \mathcal{X}^*] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathcal{X}^*, \mathbf{y}) \quad (\text{PESMOC})$$

# Predictive Entropy Search (PES)

We **swap  $y$  and  $\mathcal{X}^*$**  to obtain a reformulation of the acquisition function.

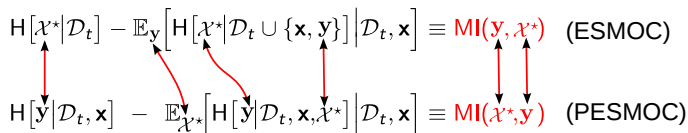
$$H[\mathcal{X}^*|\mathcal{D}_t] - \mathbb{E}_y[H[\mathcal{X}^*|\mathcal{D}_t \cup \{\mathbf{x}, y\}]|\mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC})$$

$$H[y|\mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathcal{X}^*}[H[y|\mathcal{D}_t, \mathbf{x}, \mathcal{X}^*]|\mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathcal{X}^*, y) \quad (\text{PESMOC})$$



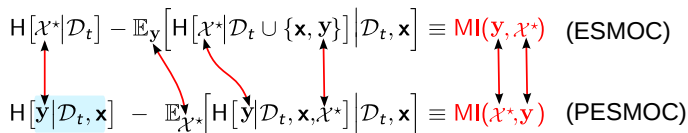
# Predictive Entropy Search (PES)

We **swap  $y$  and  $\mathcal{X}^*$**  to obtain a reformulation of the acquisition function.

$$\begin{aligned} H[\mathcal{X}^* | \mathcal{D}_t] - \mathbb{E}_y [H[\mathcal{X}^* | \mathcal{D}_t \cup \{\mathbf{x}, y\}] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC}) \\ H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathcal{X}^*} [H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathcal{X}^*] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathcal{X}^*, \mathbf{y}) \quad (\text{PESMOC}) \end{aligned}$$


# Predictive Entropy Search (PES)

We **swap  $y$  and  $\mathcal{X}^*$**  to obtain a reformulation of the acquisition function.

$$\begin{aligned} H[\mathcal{X}^* | \mathcal{D}_t] - \mathbb{E}_y [H[\mathcal{X}^* | \mathcal{D}_t \cup \{\mathbf{x}, y\}] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC}) \\ H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathcal{X}^*} [H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathcal{X}^*] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathcal{X}^*, \mathbf{y}) \quad (\text{PESMOC}) \end{aligned}$$




# Predictive Entropy Search (PES)

We **swap  $y$  and  $\mathcal{X}^*$**  to obtain a reformulation of the acquisition function.

$$H[\mathcal{X}^* | \mathcal{D}_t] - \mathbb{E}_y[H[\mathcal{X}^* | \mathcal{D}_t \cup \{\mathbf{x}, y\}] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC})$$

$$H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathcal{X}^*}[H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathcal{X}^*] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathcal{X}^*, \mathbf{y}) \quad (\text{PESMOC})$$

Gaussian  
distribution

# Predictive Entropy Search (PES)

We **swap  $y$  and  $\mathcal{X}^*$**  to obtain a reformulation of the acquisition function.

$$\begin{aligned} H[\mathcal{X}^* | \mathcal{D}_t] - \mathbb{E}_y [H[\mathcal{X}^* | \mathcal{D}_t \cup \{\mathbf{x}, y\}] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC}) \\ H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathcal{X}^*} [H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathcal{X}^*] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathcal{X}^*, \mathbf{y}) \quad (\text{PESMOC}) \end{aligned}$$

Gaussian distribution

# Predictive Entropy Search (PES)

We **swap**  $y$  and  $\mathcal{X}^*$  to obtain a reformulation of the acquisition function.

$$\begin{aligned} H[\mathcal{X}^* | \mathcal{D}_t] - \mathbb{E}_y [H[\mathcal{X}^* | \mathcal{D}_t \cup \{\mathbf{x}, y\}] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC}) \\ H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathcal{X}^*} [H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathcal{X}^*] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathcal{X}^*, \mathbf{y}) \quad (\text{PESMOC}) \end{aligned}$$

Gaussian distribution

Approximated by sampling from  $p(\mathcal{X}^* | \mathcal{D}_t)$

# Predictive Entropy Search (PES)

We **swap**  $y$  and  $\mathcal{X}^*$  to obtain a reformulation of the acquisition function.

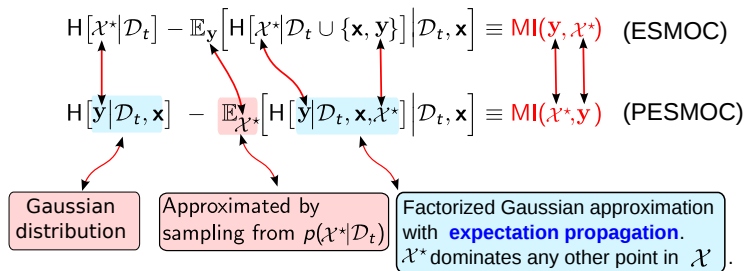
$$\begin{aligned} H[\mathcal{X}^* | \mathcal{D}_t] - \mathbb{E}_y[H[\mathcal{X}^* | \mathcal{D}_t \cup \{\mathbf{x}, y\}] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC}) \\ H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathcal{X}^*}[H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathcal{X}^*] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathcal{X}^*, \mathbf{y}) \quad (\text{PESMOC}) \end{aligned}$$

Gaussian distribution

Approximated by sampling from  $p(\mathcal{X}^* | \mathcal{D}_t)$

# Predictive Entropy Search (PES)

We **swap**  $y$  and  $\mathcal{X}^*$  to obtain a reformulation of the acquisition function.



(Minka, 2001)

# Predictive Entropy Search (PES)

We **swap**  $y$  and  $\mathcal{X}^*$  to obtain a reformulation of the acquisition function.

$$H[\mathcal{X}^* | \mathcal{D}_t] - \mathbb{E}_y[H[\mathcal{X}^* | \mathcal{D}_t \cup \{\mathbf{x}, y\}] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC})$$

$$H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathcal{X}^*}[H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathcal{X}^*] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathcal{X}^*, \mathbf{y}) \quad (\text{PESMOC})$$

Gaussian distribution

Approximated by sampling from  $p(\mathcal{X}^* | \mathcal{D}_t)$

Factorized Gaussian approximation with **expectation propagation**.  
 $\mathcal{X}^*$  dominates any other point in  $\mathcal{X}$ .

$$\alpha(\mathbf{x}) \approx \sum_{c=1}^C \log v_c^{PD}(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^M \left( \sum_{c=1}^C \log v_c^{CPD}(\mathbf{x} | \mathcal{X}_{(m)}^*) \right) + \sum_{k=1}^K \log v_k^{PD}(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^M \left( \sum_{k=1}^K \log v_k^{CPD}(\mathbf{x} | \mathcal{X}_{(m)}^*) \right)$$

(Minka, 2001)

# Predictive Entropy Search (PES)

We **swap**  $y$  and  $\mathcal{X}^*$  to obtain a reformulation of the acquisition function.

$$H[\mathcal{X}^* | \mathcal{D}_t] - \mathbb{E}_y[H[\mathcal{X}^* | \mathcal{D}_t \cup \{\mathbf{x}, y\}] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC})$$

$$H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathcal{X}^*}[H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathcal{X}^*] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathcal{X}^*, \mathbf{y}) \quad (\text{PESMOC})$$

Gaussian distribution

Approximated by sampling from  $p(\mathcal{X}^* | \mathcal{D}_t)$

Factorized Gaussian approximation with **expectation propagation**.  
 $\mathcal{X}^*$  dominates any other point in  $\mathcal{X}$ .

$$\alpha(\mathbf{x}) \approx \sum_{c=1}^C \log v_c^{PD}(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^M \left( \sum_{c=1}^C \log v_c^{CPD}(\mathbf{x} | \mathcal{X}_{(m)}^*) \right) + \sum_{k=1}^K \log v_k^{PD}(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^M \left( \sum_{k=1}^K \log v_k^{CPD}(\mathbf{x} | \mathcal{X}_{(m)}^*) \right) = \sum_{i=1}^{C+K} \alpha_i(\mathbf{x})$$

(Minka, 2001)

# Predictive Entropy Search (PES)

We **swap**  $y$  and  $\mathcal{X}^*$  to obtain a reformulation of the acquisition function.

$$H[\mathcal{X}^* | \mathcal{D}_t] - \mathbb{E}_y[H[\mathcal{X}^* | \mathcal{D}_t \cup \{\mathbf{x}, y\}] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{y}, \mathcal{X}^*) \quad (\text{ESMOC})$$

$$H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathcal{X}^*}[H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathcal{X}^*] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathcal{X}^*, \mathbf{y}) \quad (\text{PESMOC})$$

Gaussian distribution

Approximated by sampling from  $p(\mathcal{X}^* | \mathcal{D}_t)$

Factorized Gaussian approximation with **expectation**  $\mathcal{X}^*$  dominates any

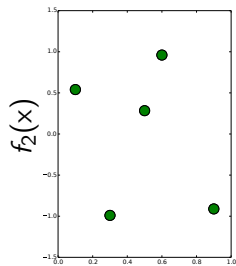
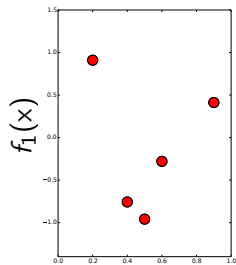
One acquisition per black-box

$$\alpha(\mathbf{x}) \approx \sum_{c=1}^C \log v_c^{PD}(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^M \left( \sum_{c=1}^C \log v_c^{CPD}(\mathbf{x} | \mathcal{X}_{(m)}^*) \right) + \sum_{k=1}^K \log v_k^{PD}(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^M \left( \sum_{k=1}^K \log v_k^{CPD}(\mathbf{x} | \mathcal{X}_{(m)}^*) \right) = \sum_{i=1}^{C+K} \alpha_i(\mathbf{x})$$

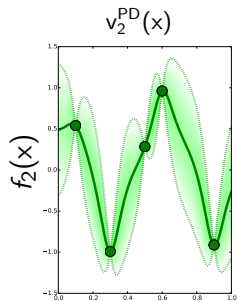
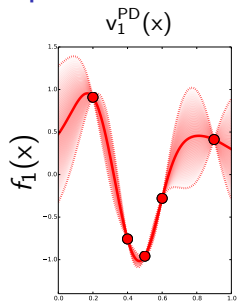
(Minka, 2001)



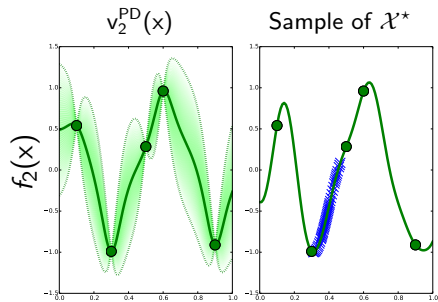
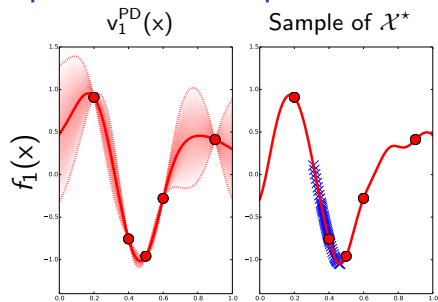
# Example of PES' acquisition



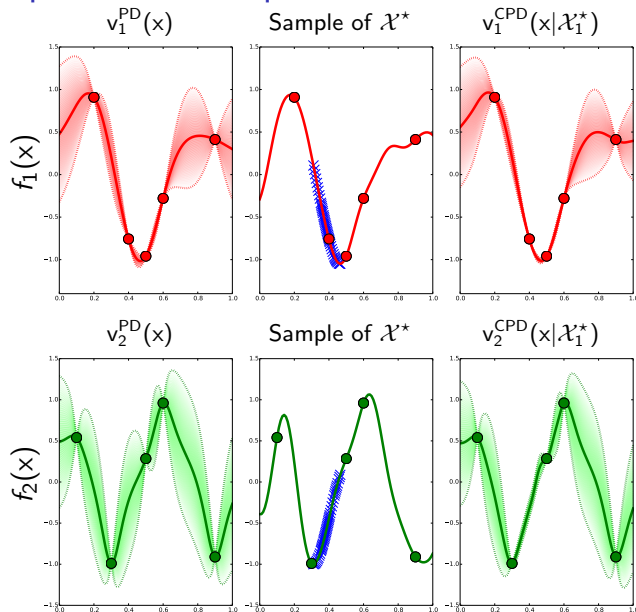
# Example of PES' acquisition



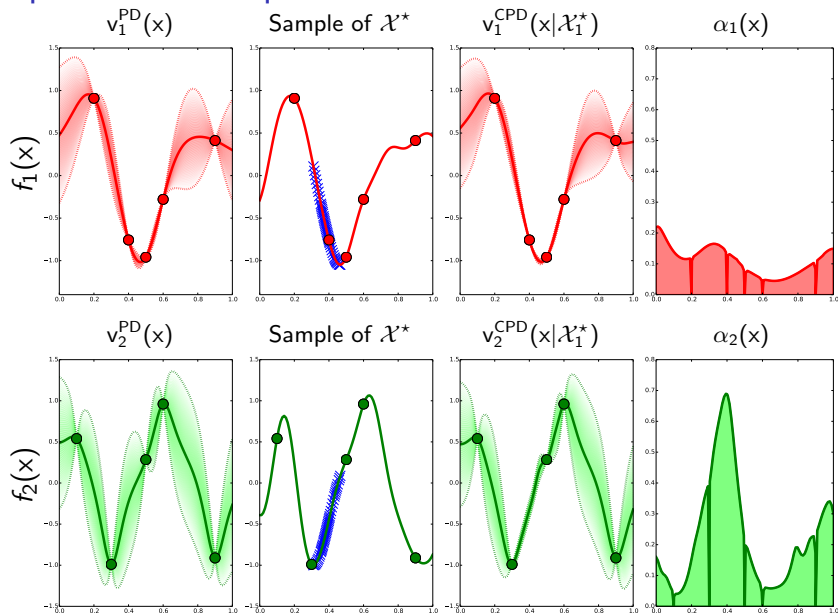
# Example of PES' acquisition



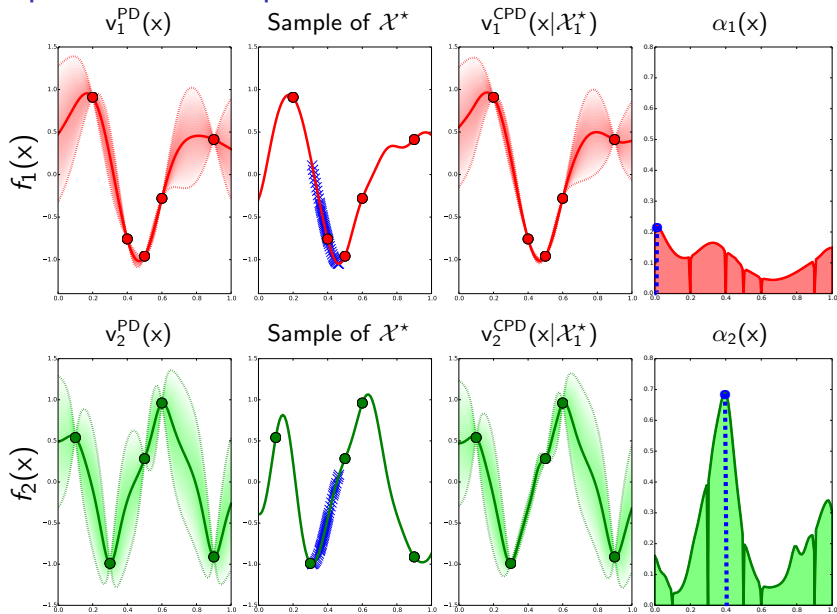
# Example of PES' acquisition



# Example of PES' acquisition

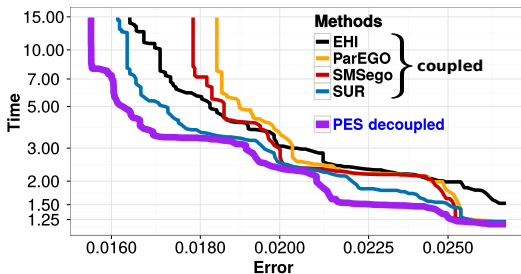


# Example of PES' acquisition



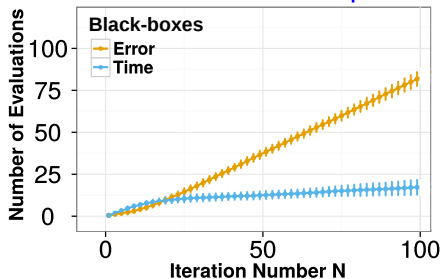
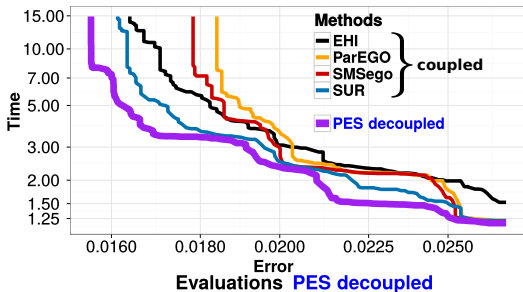
# Finding a Fast and Accurate Neural Network

Average Pareto Front 100 Function Evaluations



# Finding a Fast and Accurate Neural Network

Average Pareto Front 100 Function Evaluations

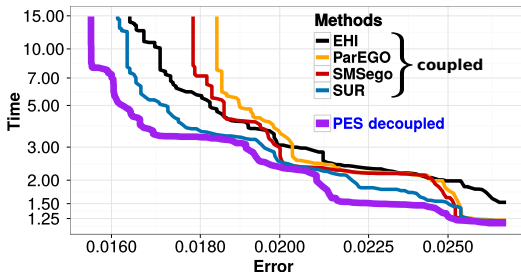


(Hernández-Lobato *et al.*, 2016)

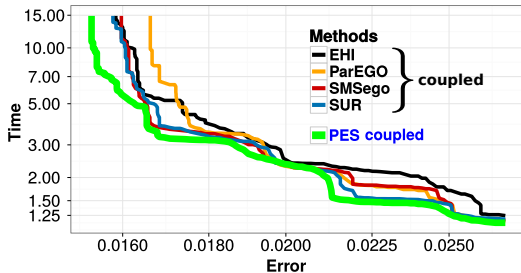


# Finding a Fast and Accurate Neural Network

Average Pareto Front 100 Function Evaluations

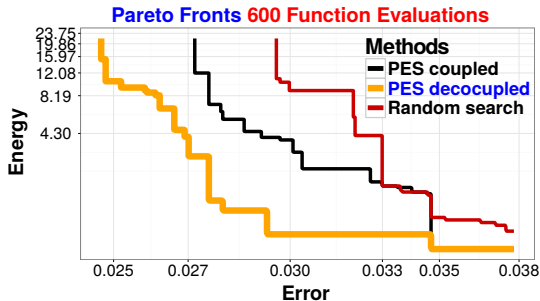


Average Pareto Front 200 Function Evaluations

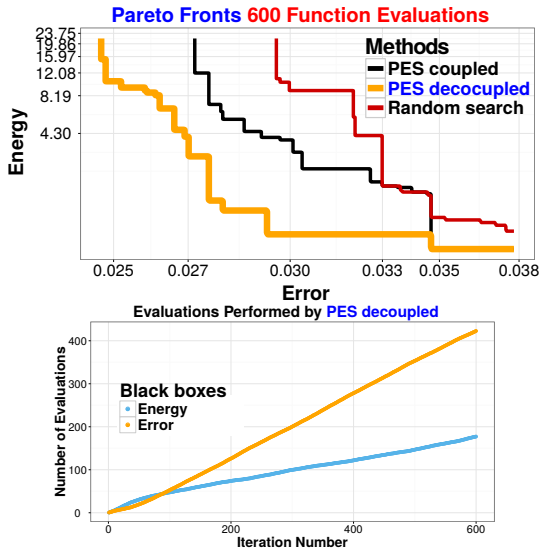


(Hernández-Lobato *et al.*, 2016)

# Low energy hardware accelerator



# Low energy hardware accelerator



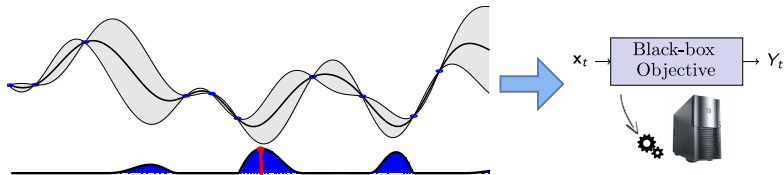
(Hernández-Lobato *et al.*, 2016)

# Parallel Bayesian Optimization

Traditional Bayesian optimization is **sequential**!

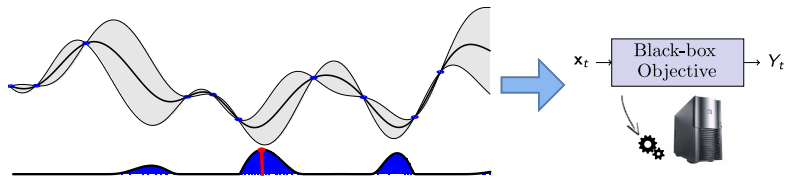
# Parallel Bayesian Optimization

Traditional Bayesian optimization is **sequential**!



# Parallel Bayesian Optimization

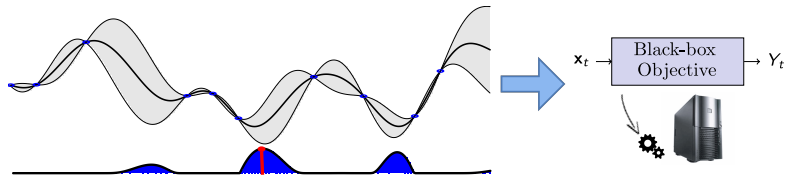
Traditional Bayesian optimization is **sequential**!



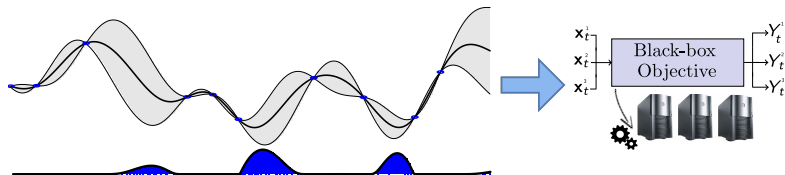
**Computing clusters** let us do **many things** at once!

# Parallel Bayesian Optimization

Traditional Bayesian optimization is **sequential**!

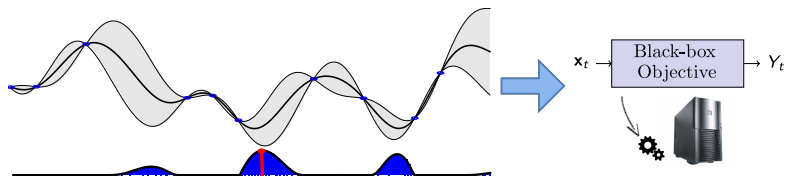


Computing clusters let us do **many things** at once!

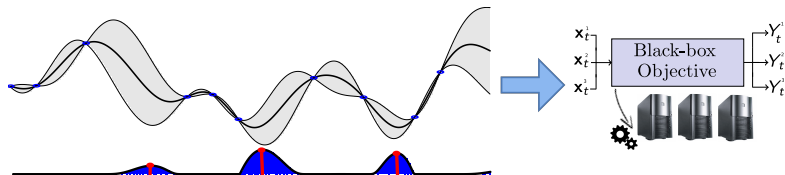


# Parallel Bayesian Optimization

Traditional Bayesian optimization is **sequential**!



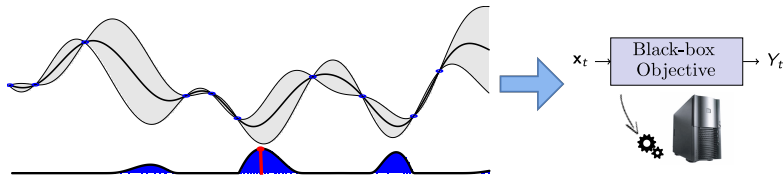
Computing clusters let us do **many things** at once!



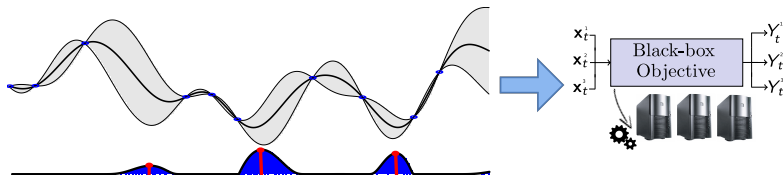


# Parallel Bayesian Optimization

Traditional Bayesian optimization is **sequential**!



Computing clusters let us do **many things** at once!



**Parallel experiments should be highly informative but different!**

# Parallel Predictive Entropy Search

Choose a set  $Q$  points  $\mathcal{S}_t = \{\mathbf{x}_q\}_{q=1}^Q$  to minimize the entropy of  $\mathbf{x}^*$ .

$$H[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}_q, y_q\}_{q=1}^Q]|\mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{y}, \mathbf{x}^*) \quad (\text{Parallel ES})$$

(Shah and Ghahramani, 2015)

# Parallel Predictive Entropy Search

Choose a set  $Q$  points  $\mathcal{S}_t = \{\mathbf{x}_q\}_{q=1}^Q$  to minimize the entropy of  $\mathbf{x}^*$ .

$$H[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}_q, y_q\}_{q=1}^Q} | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{y}, \mathbf{x}^*) \quad (\text{Parallel ES})$$

$$H[y|\mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathbf{x}^*}[H[y|\mathcal{D}_t, \mathbf{x}, \mathbf{x}^*] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{x}^*, y) \quad (\text{Parallel PES})$$

(Shah and Ghahramani, 2015)

# Parallel Predictive Entropy Search

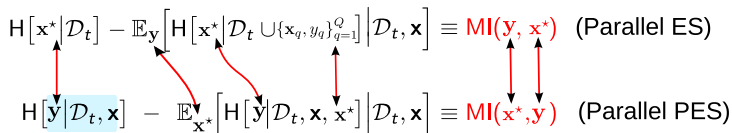
Choose a set  $Q$  points  $\mathcal{S}_t = \{\mathbf{x}_q\}_{q=1}^Q$  to minimize the entropy of  $\mathbf{x}^*$ .

$$\begin{aligned} H[\mathbf{x}^* | \mathcal{D}_t] - \mathbb{E}_{\mathbf{y}} [H[\mathbf{x}^* | \mathcal{D}_t \cup \{\mathbf{x}_q, y_q\}_{q=1}^Q} | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathbf{y}, \mathbf{x}^*) \quad (\text{Parallel ES}) \\ H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathbf{x}^*} [H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathbf{x}^*] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathbf{x}^*, \mathbf{y}) \quad (\text{Parallel PES}) \end{aligned}$$

(Shah and Ghahramani, 2015)

# Parallel Predictive Entropy Search

Choose a set  $Q$  points  $\mathcal{S}_t = \{\mathbf{x}_q\}_{q=1}^Q$  to minimize the entropy of  $\mathbf{x}^*$ .

$$\begin{aligned} H[\mathbf{x}^* | \mathcal{D}_t] - \mathbb{E}_{\mathbf{y}} [H[\mathbf{x}^* | \mathcal{D}_t \cup \{\mathbf{x}_q, \mathbf{y}_q\}_{q=1}^Q} | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathbf{y}, \mathbf{x}^*) \quad (\text{Parallel ES}) \\ H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathbf{x}^*} [H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathbf{x}^*] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathbf{x}^*, \mathbf{y}) \quad (\text{Parallel PES}) \end{aligned}$$


(Shah and Ghahramani, 2015)

# Parallel Predictive Entropy Search

Choose a set  $Q$  points  $\mathcal{S}_t = \{\mathbf{x}_q\}_{q=1}^Q$  to minimize the entropy of  $\mathbf{x}^*$ .

$$H[\mathbf{x}^* | \mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^* | \mathcal{D}_t \cup \{\mathbf{x}_q, y_q\}_{q=1}^Q} | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{y}, \mathbf{x}^*) \quad (\text{Parallel ES})$$

$$H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathbf{x}^*}[H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathbf{x}^*] | \mathcal{D}_t, \mathbf{x}] \equiv \text{MI}(\mathbf{x}^*, \mathbf{y}) \quad (\text{Parallel PES})$$

Multi-variate  
Gaussian  
distribution

(Shah and Ghahramani, 2015)

# Parallel Predictive Entropy Search

Choose a set  $Q$  points  $\mathcal{S}_t = \{\mathbf{x}_q\}_{q=1}^Q$  to minimize the entropy of  $\mathbf{x}^*$ .

$$\begin{aligned} H[\mathbf{x}^* | \mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^* | \mathcal{D}_t \cup \{\mathbf{x}_q, y_q\}_{q=1}^Q} | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathbf{y}, \mathbf{x}^*) \quad (\text{Parallel ES}) \\ H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathbf{x}^*}[H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathbf{x}^*] | \mathcal{D}_t, \mathbf{x}] &\equiv \text{MI}(\mathbf{x}^*, \mathbf{y}) \quad (\text{Parallel PES}) \end{aligned}$$

Diagram illustrating the relationship between the two equations and the underlying distribution:

- A red box labeled "Multi-variate Gaussian distribution" has a red arrow pointing to the  $\mathbf{x}$  in the second equation.
- A red arrow points from the  $\mathbf{y}$  in the second equation to the  $\mathbf{y}$  in the first equation.
- A red arrow points from the  $\mathbf{x}^*$  in the second equation to the  $\mathbf{x}^*$  in the first equation.
- A red arrow points from the  $\mathbf{x}_q$  in the first equation to the  $\mathbf{x}$  in the second equation.
- A red arrow points from the  $y_q$  in the first equation to the  $\mathbf{y}$  in the second equation.

(Shah and Ghahramani, 2015)

# Parallel Predictive Entropy Search

Choose a set  $Q$  points  $\mathcal{S}_t = \{\mathbf{x}_q\}_{q=1}^Q$  to minimize the entropy of  $\mathbf{x}^*$ .

$$\begin{aligned}
 & H[\mathbf{x}^* | \mathcal{D}_t] - \mathbb{E}_{\mathbf{y}} \left[ H[\mathbf{x}^* | \mathcal{D}_t \cup \{\mathbf{x}_q, \mathbf{y}_q\}_{q=1}^Q} \middle| \mathcal{D}_t, \mathbf{x} \right] \equiv \text{MI}(\mathbf{y}, \mathbf{x}^*) \quad (\text{Parallel ES}) \\
 & H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathbf{x}^*} \left[ H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathbf{x}^*] \middle| \mathcal{D}_t, \mathbf{x} \right] \equiv \text{MI}(\mathbf{x}^*, \mathbf{y}) \quad (\text{Parallel PES})
 \end{aligned}$$

Multi-variate Gaussian distribution

Approximated by sampling from  $p(\mathbf{x}^* | \mathcal{D}_t)$

(Shah and Ghahramani, 2015)



# Parallel Predictive Entropy Search

Choose a set  $Q$  points  $\mathcal{S}_t = \{\mathbf{x}_q\}_{q=1}^Q$  to minimize the entropy of  $\mathbf{x}^*$ .

$$\begin{aligned}
 & H[\mathbf{x}^* | \mathcal{D}_t] - \mathbb{E}_{\mathbf{y}} \left[ H[\mathbf{x}^* | \mathcal{D}_t \cup \{\mathbf{x}_q, \mathbf{y}_q\}_{q=1}^Q} \middle| \mathcal{D}_t, \mathbf{x} \right] \equiv \text{MI}(\mathbf{y}, \mathbf{x}^*) \quad (\text{Parallel ES}) \\
 & H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}] - \mathbb{E}_{\mathbf{x}^*} \left[ H[\mathbf{y} | \mathcal{D}_t, \mathbf{x}, \mathbf{x}^*] \middle| \mathcal{D}_t, \mathbf{x} \right] \equiv \text{MI}(\mathbf{x}^*, \mathbf{y}) \quad (\text{Parallel PES})
 \end{aligned}$$

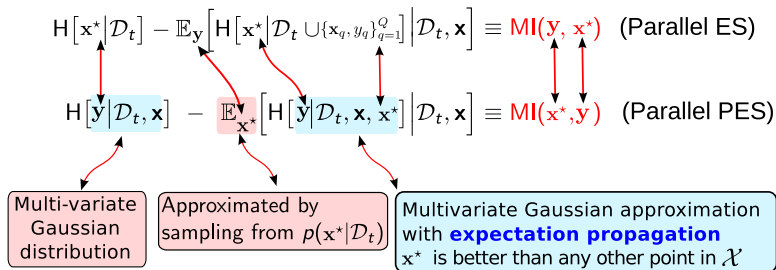
Multi-variate Gaussian distribution

Approximated by sampling from  $p(\mathbf{x}^* | \mathcal{D}_t)$

(Shah and Ghahramani, 2015)

# Parallel Predictive Entropy Search

Choose a set  $Q$  points  $\mathcal{S}_t = \{\mathbf{x}_q\}_{q=1}^Q$  to minimize the entropy of  $\mathbf{x}^*$ .



(Shah and Ghahramani, 2015)

# Parallel Predictive Entropy Search

Choose a set  $Q$  points  $\mathcal{S}_t = \{\mathbf{x}_q\}_{q=1}^Q$  to minimize the entropy of  $\mathbf{x}^*$ .

$$H[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}_q, \mathbf{y}_q\}_{q=1}^Q}|\mathcal{D}_t, \mathbf{X}] \equiv \text{MI}(\mathbf{y}, \mathbf{x}^*) \quad (\text{Parallel ES})$$

$$H[\mathbf{y}|\mathcal{D}_t, \mathbf{X}] - \mathbb{E}_{\mathbf{x}^*}[H[\mathbf{y}|\mathcal{D}_t, \mathbf{X}, \mathbf{x}^*]|\mathcal{D}_t, \mathbf{X}] \equiv \text{MI}(\mathbf{x}^*, \mathbf{y}) \quad (\text{Parallel PES})$$

Multi-variate  
Gaussian  
distribution

Approximated by  
sampling from  $p(\mathbf{x}^*|\mathcal{D}_t)$

Multivariate Gaussian approximation  
with **expectation propagation**  
 $\mathbf{x}^*$  is better than any other point in  $\mathcal{X}$

$$\alpha(\mathbf{X}) = \log |\mathbf{V}^{\text{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^M \log |\mathbf{V}^{\text{CPD}}(\mathbf{X}|\mathbf{x}_{(m)}^*)|$$

(Shah and Ghahramani, 2015)

# Parallel Predictive Entropy Search

Choose a set  $Q$  points  $\mathcal{S}_t = \{\mathbf{x}_q\}_{q=1}^Q$  to minimize the entropy of  $\mathbf{x}^*$ .

$$H[\mathbf{x}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}}[H[\mathbf{x}^*|\mathcal{D}_t \cup \{\mathbf{x}_q, \mathbf{y}_q\}_{q=1}^Q}|\mathcal{D}_t, \mathbf{X}] \equiv \text{MI}(\mathbf{y}, \mathbf{x}^*) \quad (\text{Parallel ES})$$

$$H[\mathbf{y}|\mathcal{D}_t, \mathbf{X}] - \mathbb{E}_{\mathbf{x}^*}[H[\mathbf{y}|\mathcal{D}_t, \mathbf{X}, \mathbf{x}^*]|\mathcal{D}_t, \mathbf{X}] \equiv \text{MI}(\mathbf{x}^*, \mathbf{y}) \quad (\text{Parallel PES})$$

Multi-variate  
Gaussian  
distribution

Approximated by  
sampling from  $p(\mathbf{x}^*|\mathcal{D}_t)$

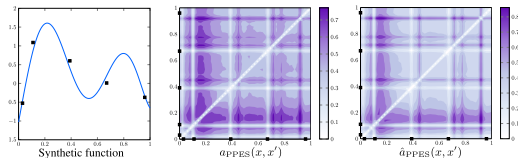
Multivariate Gaussian approximation  
with **expectation propagation**  
 $\mathbf{x}^*$  is better than any other point in  $\mathcal{X}$

$$\alpha(\mathbf{X}) = \log |\mathbf{V}^{\text{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^M \log |\mathbf{V}^{\text{CPD}}(\mathbf{X}|\mathbf{x}_{(m)}^*)|$$

**It is possible to compute the gradient of  $\alpha(\cdot)$  w.r.t. each  $\mathbf{x}_q \in \mathcal{S}_t$ !**

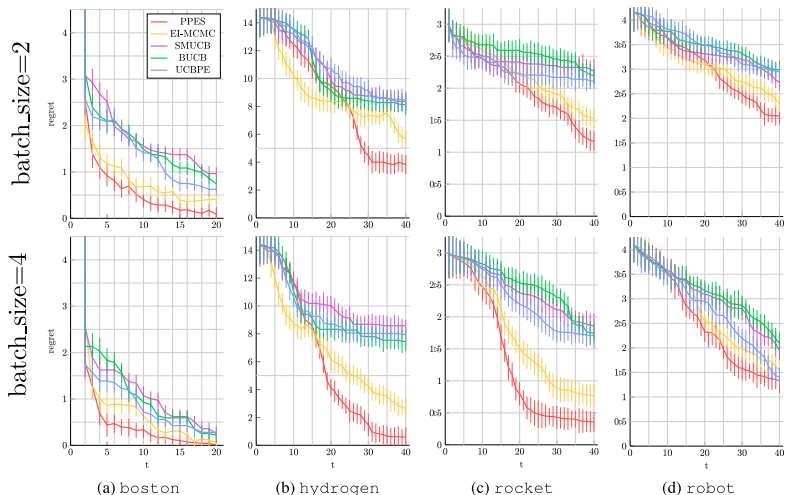
(Shah and Ghahramani, 2015)

# Parallel Predictive Entropy Search: Level Curves



(Shah and Ghahramani, 2015)

# Parallel Predictive Entropy Search: Results



(Shah and Ghahramani, 2015)

# BO with Integer-valued and Categorical Variables

Standard GPs assume continuous input variables which makes BO with integer-valued or categorical challenging.

# BO with Integer-valued and Categorical Variables

Standard GPs assume continuous input variables which makes BO with integer-valued or categorical challenging.

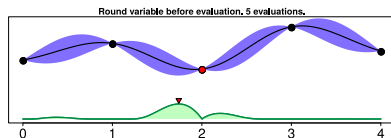
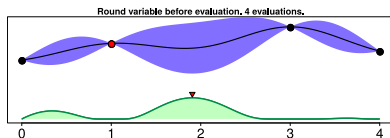
A naive approach is to round the suggested value to the closest integer or to the closest one-hot encoding.



# BO with Integer-valued and Categorical Variables

Standard GPs assume continuous input variables which makes BO with integer-valued or categorical challenging.

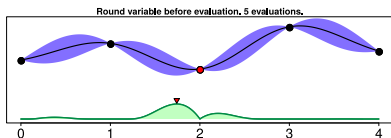
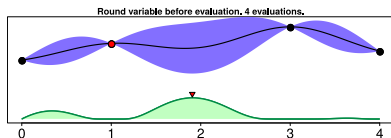
A naive approach is to round the suggested value to the closest integer or to the closest one-hot encoding.



# BO with Integer-valued and Categorical Variables

Standard GPs assume continuous input variables which makes BO with integer-valued or categorical challenging.

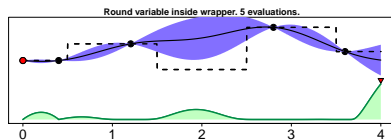
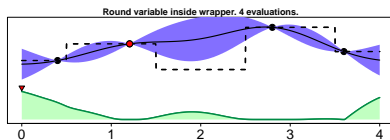
A naive approach is to round the suggested value to the closest integer or to the closest one-hot encoding.



**The BO algorithm may get stuck and may always perform the next evaluation at the same input location!**

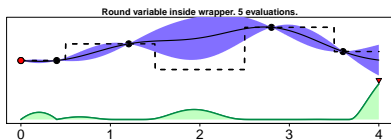
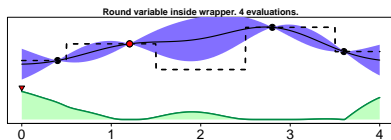
# BO with Integer-valued and Categorical Variables

Rounding inside of the wrapper works but makes the objective flat!



# BO with Integer-valued and Categorical Variables

Rounding inside of the wrapper works but makes the objective flat!



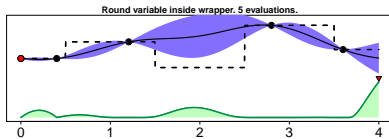
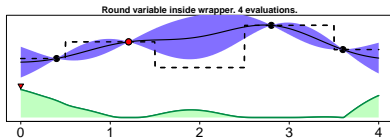
A modified GP covariance function accounts for this:

$$C_{\text{new}}(\mathbf{x}_n, \mathbf{x}_{n'}) = C(T(\mathbf{x}_n), T(\mathbf{x}_{n'}); \boldsymbol{\theta})$$

where  $T(\cdot)$  does the rounding to the closest integer or one-hot encoding.

# BO with Integer-valued and Categorical Variables

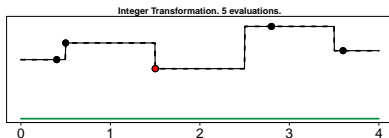
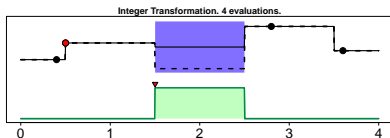
Rounding inside of the wrapper works but makes the objective flat!



A modified GP covariance function accounts for this:

$$C_{\text{new}}(x_n, x_{n'}) = C(T(x_n), T(x_{n'}); \theta)$$

where  $T(\cdot)$  does the rounding to the closest integer or one-hot encoding.

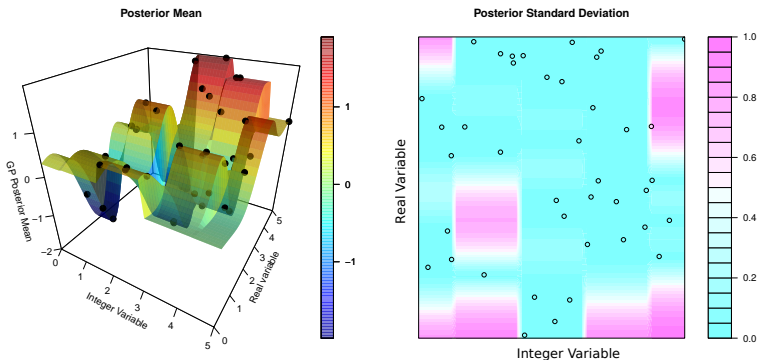


# BO with Integer-valued and Categorical Variables

The GP predictive distribution is constant across all variables that lead to the same integer or one-hot-encoding.

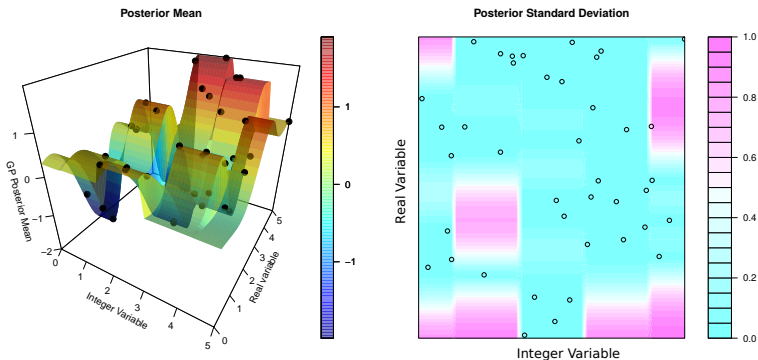
# BO with Integer-valued and Categorical Variables

The GP predictive distribution is constant across all variables that lead to the same integer or one-hot-encoding.



# BO with Integer-valued and Categorical Variables

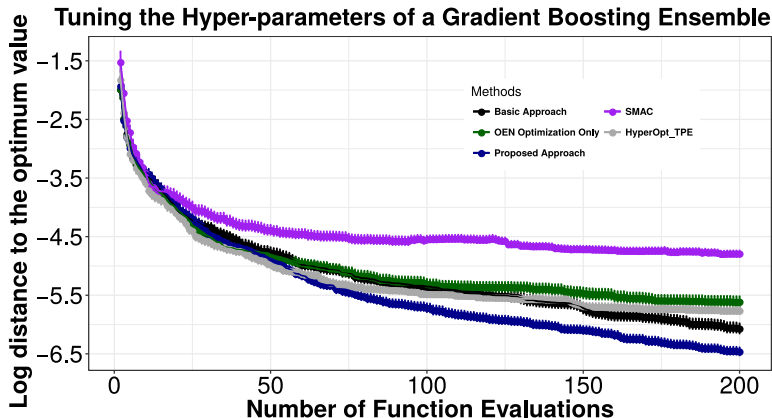
The GP predictive distribution is constant across all variables that lead to the same integer or one-hot-encoding.



Similar results for categorical variables!



# BO with Integer-valued and Categorical Variables



One continuous variable and two integer-valued variables.

# Freeze-Thaw Bayesian Optimization

Common aspects of many machine learning algorithms:

# Freeze-Thaw Bayesian Optimization

Common aspects of many machine learning algorithms:

1. A minimization step must be performed with, e.g., gradient descent.

# Freeze-Thaw Bayesian Optimization

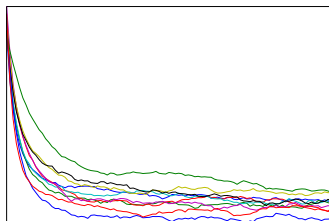
Common aspects of many machine learning algorithms:

1. A minimization step must be performed with, e.g., gradient descent.
2. There are hyper-parameters that impact the final performance.

# Freeze-Thaw Bayesian Optimization

Common aspects of many machine learning algorithms:

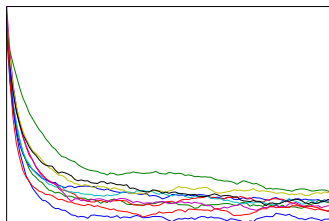
1. A minimization step must be performed with, e.g., gradient descent.
2. There are hyper-parameters that impact the final performance.



# Freeze-Thaw Bayesian Optimization

Common aspects of many machine learning algorithms:

1. A minimization step must be performed with, e.g., gradient descent.
2. There are hyper-parameters that impact the final performance.

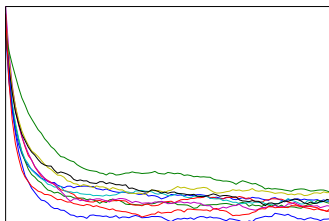


Can we use partial training information and a model to determine which hyper-parameter configuration is going to be optimal?

# Freeze-Thaw Bayesian Optimization

Common aspects of many machine learning algorithms:

1. A minimization step must be performed with, e.g., gradient descent.
2. There are hyper-parameters that impact the final performance.



Can we use partial training information and a model to determine which hyper-parameter configuration is going to be optimal?

**Yes, that is precisely what Freeze-Thaw BO does!**

(Swersky et al., 2014)

# A GP Kernel for Training Curves

We want to specify a kernel that supports exponentially decaying functions of the form  $\exp\{-\lambda t\}$  for  $t, \lambda \geq 0$ .



# A GP Kernel for Training Curves

We want to specify a kernel that supports exponentially decaying functions of the form  $\exp\{-\lambda t\}$  for  $t, \lambda \geq 0$ .

The covariance between inputs  $t$  and  $t'$  is:

$$C(t, t') = \int_0^\infty e^{-\lambda t} e^{-\lambda t'} \psi(\lambda; \alpha, \beta) d\lambda = \frac{\beta^\alpha}{(t + t' + \beta)^\alpha}$$

where  $\psi(\lambda; \alpha, \beta)$  is a gamma distribution with parameters  $\alpha$  and  $\beta$ .

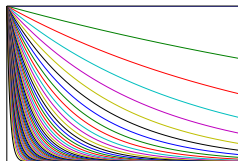
# A GP Kernel for Training Curves

We want to specify a kernel that supports exponentially decaying functions of the form  $\exp\{-\lambda t\}$  for  $t, \lambda \geq 0$ .

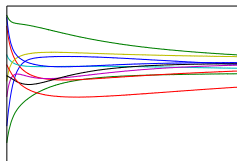
The covariance between inputs  $t$  and  $t'$  is:

$$C(t, t') = \int_0^\infty e^{-\lambda t} e^{-\lambda t'} \psi(\lambda; \alpha, \beta) d\lambda = \frac{\beta^\alpha}{(t + t' + \beta)^\alpha}$$

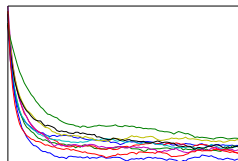
where  $\psi(\lambda; \alpha, \beta)$  is a gamma distribution with parameters  $\alpha$  and  $\beta$ .



(a) Exponential Decay Basis



(b) Samples



(c) Training Curve Samples

# Inference on Asymptotic Values

**A standard GP is used as the prior for the asymptotic values of each training curve.**

# Inference on Asymptotic Values

**A standard GP is used as the prior for the asymptotic values of each training curve.**

Hierarchical generative model:

$$p(\{y_n\}_{n=1}^N | \{x_n\}_{n=1}^N) = \int \left[ \prod_{n=1}^N \mathcal{N}(y_n | f_n, K_{t_n}) \right] \mathcal{N}(f | m, K_x) df$$

where

$x_n \equiv n$  configuration ,

$y_n \equiv n$  observed curve ,

$f_n \equiv n$  asymptotic value ,

$m \equiv$  prior asymptotic mean values ,

$K_{t_n} \equiv$  covariances for curve values ,  $K_x \equiv$  cov. for asymptotic values

# Inference on Asymptotic Values

**A standard GP is used as the prior for the asymptotic values of each training curve.**

Hierarchical generative model:

$$p(\{y_n\}_{n=1}^N | \{x_n\}_{n=1}^N) = \int \left[ \prod_{n=1}^N \mathcal{N}(y_n | f_n, K_{t_n}) \right] \mathcal{N}(f | m, K_x) df$$

where

$x_n \equiv n$  configuration ,

$y_n \equiv n$  observed curve ,

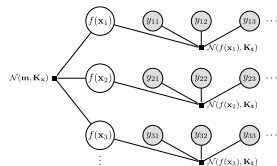
$f_n \equiv n$  asymptotic value ,

$m \equiv$  prior asymptotic mean values ,

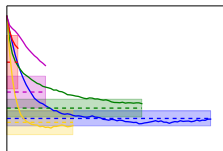
$K_{t_n} \equiv$  covariances for curve values ,  $K_x \equiv$  cov. for asymptotic values

**The joint distribution of  $\{y\}_{n=1}^N$  and  $f$  is Gaussian and hence so it is the predictive distribution  $p(f | \{y\}_{n=1}^N)$ !**

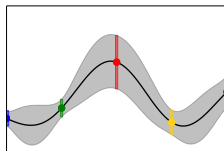
# Inference on Asymptotic Values and BO



(a) Graphical Model

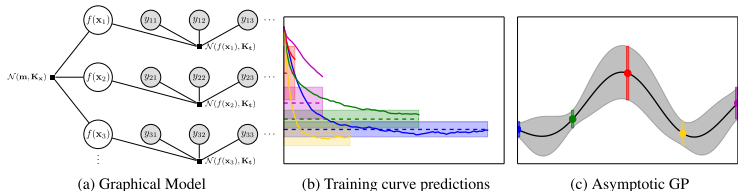


(b) Training curve predictions



(c) Asymptotic GP

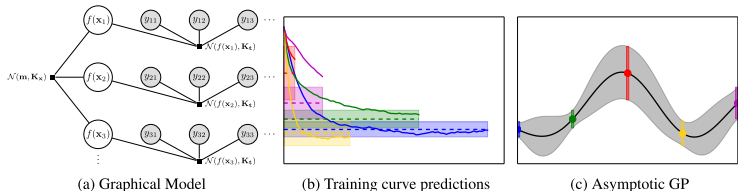
# Inference on Asymptotic Values and BO



## Bayesian Optimization:

- $p(f | \{y_n\}_{n=1}^N, \{\mathbf{x}_n\}_{n=1}^N)$  determines asymptotic values.

# Inference on Asymptotic Values and BO

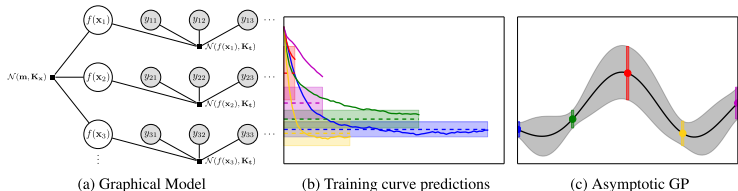


## Bayesian Optimization:

- ▶  $p(f | \{y_n\}_{n=1}^N, \{\mathbf{x}_n\}_{n=1}^N)$  determines asymptotic values.
- ▶ This distribution can be used to make intelligent decisions!



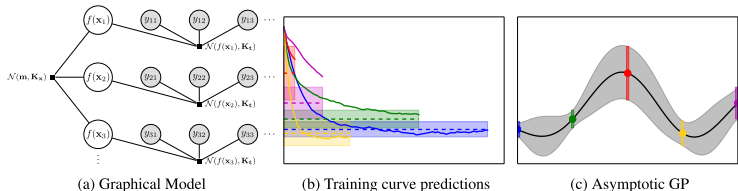
# Inference on Asymptotic Values and BO



## Bayesian Optimization:

- ▶  $p(f | \{y_n\}_{n=1}^N, \{\mathbf{x}_n\}_{n=1}^N)$  determines asymptotic values.
- ▶ This distribution can be used to make intelligent decisions!
- ▶ Shall we train more one configuration or shall we start a new one?

# Inference on Asymptotic Values and BO

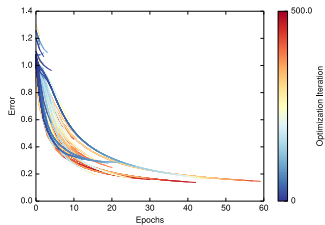
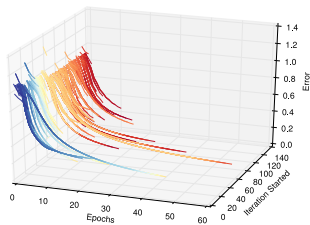


## Bayesian Optimization:

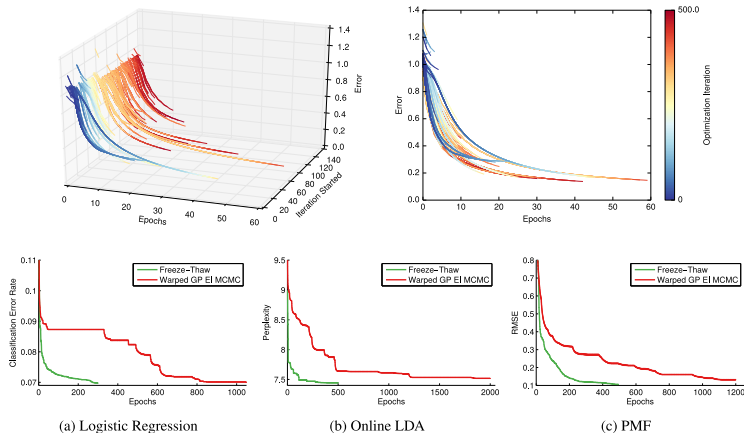
- ▶  $p(f | \{y_n\}_{n=1}^N, \{x_n\}_{n=1}^N)$  determines asymptotic values.
- ▶ This distribution can be used to make intelligent decisions!
- ▶ Shall we train more one configuration or shall we start a new one?
- ▶ A combination of EI and ES is used as the acquisition function.

(Swersky et al., 2014)

# Freeze-Thaw BO in practice

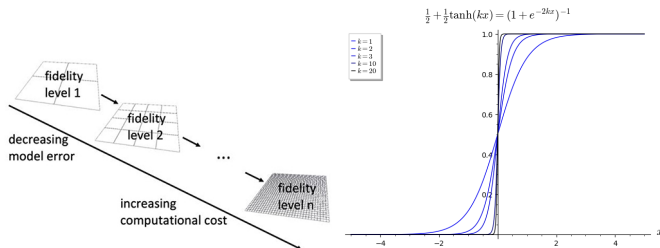


# Freeze-Thaw BO in practice



(Swersky et al., 2014)

# Multi-fidelity Bayesian optimization



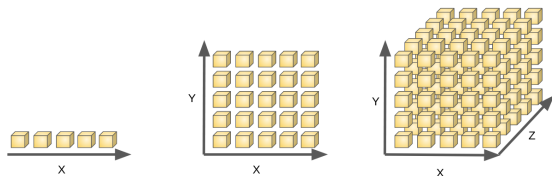
- ▶ **Concept of Fidelity:** Fidelity refers to the accuracy or reliability of the information, we use different levels of accuracy.
- ▶ **Example:** # layers in NNs, # of timesteps for DRL.
- ▶ *We assume that lower fidelities are correlated with higher.*
- ▶ **Cost-Efficiency:** Leveraging computationally cheaper versions of the functions to guide the search process.
- ▶ **The trick:** The acquisition function should balance exploration at cheaper fidelities with exploitation at the highest fidelity level.

# Multi-fidelity Bayesian optimization example

- ▶ Trace-aware knowledge-gradient acquisition function.
- ▶ It values observations of a point ( $x$ ) and a set of fidelities ( $\mathcal{S}$ ) according to the ratio of the reduction in expected loss that it induces, to its computational cost.
- ▶ It measures the value of information per unit cost of sampling.
- ▶ It uses a function  $L()$  to measure the extent to which observing trace information improves the quality of the solution.  $L(0)$  will be the minimum.
- ▶ Its analytical expression basically penalizes the cost wrt the information obtained:

$$takg(x, \mathcal{S}) = \frac{L(0) - L(x, \mathcal{S})}{cost(x, \max(\mathcal{S}))}$$

# High-dimensional Bayesian optimization



- ▶ GPs empirical performance tends to be lower if  $d > 7$ .
- ▶ **The problem:** The search space grows exponentially with the number of dimensions.
- ▶ **The trick:** We can project the high-dimensional problem into a lower-dimensional subspace that explains it well using embeddings.
- ▶ *Find the hidden most explicative manifold for the data, then optimize there!*
- ▶ Approaches differ in the type of embeddings (e.g. random), assumptions about the function (e.g. being a sum of functions), or use of models (Deep GP, GP-LVM.)

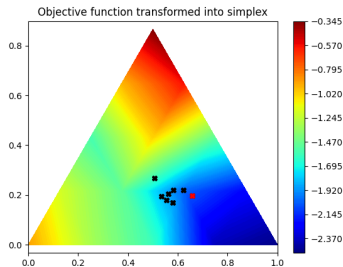
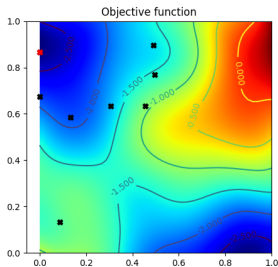
# High-dimensional Bayesian optimization example

- ▶ Sparse Axis-Aligned Subspace Bayesian Optimization (SAASBO)
- ▶ **Goal:** Identify sparse subspaces relevant to modeling the unknown objective function.
- ▶ **Assumption:** High function variability being captured by axis-aligned blocks of input features.
- ▶ **Method:** Use complex GP prior to consider a smaller class of functions.
- ▶ **Effect:** Turn most non-explicative dims off,
- ▶ **Integration:** Perform hyper-parameter sampling with NUTS and EI.

[kernel variance]	$\sigma_k^2 \sim \mathcal{LN}(0, 10^2)$
[global shrinkage]	$\tau \sim \mathcal{HC}(\alpha)$
[length scales]	$\rho_i \sim \mathcal{HC}(\tau)$
[function values]	$\mathbf{f} \sim \mathcal{N}(\mathbf{0}, K_{\mathbf{xx}}^\psi)$
[observations]	$\mathbf{y} \sim \mathcal{N}(\mathbf{f}, \sigma^2 \mathbf{1}_N)$



# Bayesian optimization in a simplex (Portfolio optimization)



- ▶ Common financial metrics as Sharpe or Sortino ratio can be configured for a specific portfolio.
- ▶ For example using real-time ESG values of the assets.
- ▶ Garrido-Merchán, E. C., Piris, G. G., & Vaca, M. C. (2023). Bayesian optimization of ESG (Environmental Social Governance) financial investments. Environmental Research Communications, 5(5), 055003.

# Conclusions and Further Work

- ▶ BO is a state-of-the-art class of methods used to optimize expensive and noisy black-box functions.
- ▶ We can generalize BO to tackle advanced scenarios: parallel constrained multi-obj, high-dim, multi-fidelity...
- ▶ BO can be applied in a wide array of applications: ML, DRL, finance, robotics, materials, business operations...
- ▶ Further work topics: causality, meta-BO, transfer learning, adaptation to specific domains.



# Referencias

- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2019). Predictive Entropy Search for Multi-Objective Bayesian Optimization with Constraints. *Neurocomputing*, 361, 50-68.

# Referencias

- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2019). Predictive Entropy Search for Multi-Objective Bayesian Optimization with Constraints. *Neurocomputing*, 361, 50-68.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2023). Parallel Predictive Entropy Search for Multi-objective Bayesian Optimization with Constraints. *Expert Systems with Applications*, pre-print.

# Referencias

- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2019). Predictive Entropy Search for Multi-Objective Bayesian Optimization with Constraints. *Neurocomputing*, 361, 50-68.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2023). Parallel Predictive Entropy Search for Multi-objective Bayesian Optimization with Constraints. *Expert Systems with Applications*, pre-print.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2020). Dealing with Categorical and Integer-valued Variables in Bayesian Optimization with Gaussian Processes. *Neurocomputing*, 380, 20-35.

# Referencias

- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2019). Predictive Entropy Search for Multi-Objective Bayesian Optimization with Constraints. *Neurocomputing*, 361, 50-68.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2023). Parallel Predictive Entropy Search for Multi-objective Bayesian Optimization with Constraints. *Expert Systems with Applications*, pre-print.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2020). Dealing with Categorical and Integer-valued Variables in Bayesian Optimization with Gaussian Processes. *Neurocomputing*, 380, 20-35.
- ▶ Cornejo-Bueno, L., Garrido-Merchán, E. C., Hernández-Lobato, D., & Salcedo-Sanz, S. (2018). Bayesian Optimization of a Hybrid System for Robust Ocean Wave Features Prediction. *Neurocomputing*, 275, 818-828.

# Referencias

- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2019). Predictive Entropy Search for Multi-Objective Bayesian Optimization with Constraints. *Neurocomputing*, 361, 50-68.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2023). Parallel Predictive Entropy Search for Multi-objective Bayesian Optimization with Constraints. *Expert Systems with Applications*, pre-print.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2020). Dealing with Categorical and Integer-valued Variables in Bayesian Optimization with Gaussian Processes. *Neurocomputing*, 380, 20-35.
- ▶ Cornejo-Bueno, L., Garrido-Merchán, E. C., Hernández-Lobato, D., & Salcedo-Sanz, S. (2018). Bayesian Optimization of a Hybrid System for Robust Ocean Wave Features Prediction. *Neurocomputing*, 275, 818-828.
- ▶ Cornejo-Bueno, L., Garrido-Merchán, E. C., Hernández-Lobato, D., & Salcedo-Sanz, S. (2017, June). Bayesian optimization of a hybrid prediction system for optimal wave energy estimation problems. In *International Work-Conference on Artificial Neural Networks* (pp. 648-660). Springer, Cham.

# Referencias

- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2019). Predictive Entropy Search for Multi-Objective Bayesian Optimization with Constraints. *Neurocomputing*, 361, 50-68.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2023). Parallel Predictive Entropy Search for Multi-objective Bayesian Optimization with Constraints. *Expert Systems with Applications*, pre-print.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2020). Dealing with Categorical and Integer-valued Variables in Bayesian Optimization with Gaussian Processes. *Neurocomputing*, 380, 20-35.
- ▶ Cornejo-Bueno, L., Garrido-Merchán, E. C., Hernández-Lobato, D., & Salcedo-Sanz, S. (2018). Bayesian Optimization of a Hybrid System for Robust Ocean Wave Features Prediction. *Neurocomputing*, 275, 818-828.
- ▶ Cornejo-Bueno, L., Garrido-Merchán, E. C., Hernández-Lobato, D., & Salcedo-Sanz, S. (2017, June). Bayesian optimization of a hybrid prediction system for optimal wave energy estimation problems. In *International Work-Conference on Artificial Neural Networks* (pp. 648-660). Springer, Cham.
- ▶ Balázs, C., van Beekveld, M., Caron, S., Dillon, B. M., Farmer, B., Fowlie, A., Garrido-Merchán, E. C., ... & White, M. (2021). A Comparison of Optimisation Algorithms for High-Dimensional Particle and Astrophysics Applications. *Journal of High Energy Physics*, 2021(5), 1-46.



# Referencias

- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2019). Predictive Entropy Search for Multi-Objective Bayesian Optimization with Constraints. *Neurocomputing*, 361, 50-68.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2023). Parallel Predictive Entropy Search for Multi-objective Bayesian Optimization with Constraints. *Expert Systems with Applications*, pre-print.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2020). Dealing with Categorical and Integer-valued Variables in Bayesian Optimization with Gaussian Processes. *Neurocomputing*, 380, 20-35.
- ▶ Cornejo-Bueno, L., Garrido-Merchán, E. C., Hernández-Lobato, D., & Salcedo-Sanz, S. (2018). Bayesian Optimization of a Hybrid System for Robust Ocean Wave Features Prediction. *Neurocomputing*, 275, 818-828.
- ▶ Cornejo-Bueno, L., Garrido-Merchán, E. C., Hernández-Lobato, D., & Salcedo-Sanz, S. (2017, June). Bayesian optimization of a hybrid prediction system for optimal wave energy estimation problems. In *International Work-Conference on Artificial Neural Networks* (pp. 648-660). Springer, Cham.
- ▶ Balázs, C., van Beekveld, M., Caron, S., Dillon, B. M., Farmer, B., Fowlie, A., Garrido-Merchán, E. C., ... & White, M. (2021). A Comparison of Optimisation Algorithms for High-Dimensional Particle and Astrophysics Applications. *Journal of High Energy Physics*, 2021(5), 1-46.
- ▶ Córdoba, I., Garrido-Merchán, E. C., Hernández-Lobato, D., Bielza, C., & Larranaga, P. (2018, October). Bayesian Optimization of the PC Algorithm for Learning Gaussian Bayesian Networks. In *Conference of the Spanish Association for Artificial Intelligence* (pp. 44-54). Springer, Cham.

# Referencias

- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2019). Predictive Entropy Search for Multi-Objective Bayesian Optimization with Constraints. *Neurocomputing*, 361, 50-68.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2023). Parallel Predictive Entropy Search for Multi-objective Bayesian Optimization with Constraints. *Expert Systems with Applications*, pre-print.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2020). Dealing with Categorical and Integer-valued Variables in Bayesian Optimization with Gaussian Processes. *Neurocomputing*, 380, 20-35.
- ▶ Cornejo-Bueno, L., Garrido-Merchán, E. C., Hernández-Lobato, D., & Salcedo-Sanz, S. (2018). Bayesian Optimization of a Hybrid System for Robust Ocean Wave Features Prediction. *Neurocomputing*, 275, 818-828.
- ▶ Cornejo-Bueno, L., Garrido-Merchán, E. C., Hernández-Lobato, D., & Salcedo-Sanz, S. (2017, June). Bayesian optimization of a hybrid prediction system for optimal wave energy estimation problems. In *International Work-Conference on Artificial Neural Networks* (pp. 648-660). Springer, Cham.
- ▶ Balázs, C., van Beekveld, M., Caron, S., Dillon, B. M., Farmer, B., Fowlie, A., Garrido-Merchán, E. C., ... & White, M. (2021). A Comparison of Optimisation Algorithms for High-Dimensional Particle and Astrophysics Applications. *Journal of High Energy Physics*, 2021(5), 1-46.
- ▶ Córdoba, I., Garrido-Merchán, E. C., Hernández-Lobato, D., Bielza, C., & Larranaga, P. (2018, October). Bayesian Optimization of the PC Algorithm for Learning Gaussian Bayesian Networks. In *Conference of the Spanish Association for Artificial Intelligence* (pp. 44-54). Springer, Cham.
- ▶ Garrido-Merchán, E. C., & Albarca-Molina, A. (2018, November). Suggesting Cooking Recipes through Simulation and Bayesian Optimization. In *International Conference on Intelligent Data Engineering and Automated Learning* (pp. 277-284). Springer, Cham.

# Referencias

- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2019). Predictive Entropy Search for Multi-Objective Bayesian Optimization with Constraints. *Neurocomputing*, 361, 50-68.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2023). Parallel Predictive Entropy Search for Multi-objective Bayesian Optimization with Constraints. *Expert Systems with Applications*, pre-print.
- ▶ Garrido-Merchán, E. C., & Hernández-Lobato, D. (2020). Dealing with Categorical and Integer-valued Variables in Bayesian Optimization with Gaussian Processes. *Neurocomputing*, 380, 20-35.
- ▶ Cornejo-Bueno, L., Garrido-Merchán, E. C., Hernández-Lobato, D., & Salcedo-Sanz, S. (2018). Bayesian Optimization of a Hybrid System for Robust Ocean Wave Features Prediction. *Neurocomputing*, 275, 818-828.
- ▶ Cornejo-Bueno, L., Garrido-Merchán, E. C., Hernández-Lobato, D., & Salcedo-Sanz, S. (2017, June). Bayesian optimization of a hybrid prediction system for optimal wave energy estimation problems. In *International Work-Conference on Artificial Neural Networks* (pp. 648-660). Springer, Cham.
- ▶ Balázs, C., van Beekveld, M., Caron, S., Dillon, B. M., Farmer, B., Fowlie, A., Garrido-Merchán, E. C., ... & White, M. (2021). A Comparison of Optimisation Algorithms for High-Dimensional Particle and Astrophysics Applications. *Journal of High Energy Physics*, 2021(5), 1-46.
- ▶ Córdoba, I., Garrido-Merchán, E. C., Hernández-Lobato, D., Bielza, C., & Larranaga, P. (2018, October). Bayesian Optimization of the PC Algorithm for Learning Gaussian Bayesian Networks. In *Conference of the Spanish Association for Artificial Intelligence* (pp. 44-54). Springer, Cham.
- ▶ Garrido-Merchán, E. C., & Albarca-Molina, A. (2018, November). Suggesting Cooking Recipes through Simulation and Bayesian Optimization. In *International Conference on Intelligent Data Engineering and Automated Learning* (pp. 277-284). Springer, Cham.
- ▶ Jariego-Pérez, L. C., & Garrido-Merchán, E. C. (2020). Towards Automatic Bayesian Optimization: A first step involving acquisition functions. *CAEPIA*. 2021. Accepted.

**Thank you for your attention.**