

Parallel predictive entropy search for multi-objective Bayesian optimization with constraints

Daniel Hernández-Lobato

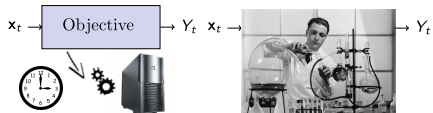
Computer Science Department
Universidad Autónoma de Madrid

<http://dhnz1.org>, daniel.hernandez@uam.es

Joint work with
Eduardo C. Garrido-Merchán and Daniel Fernández Sánchez

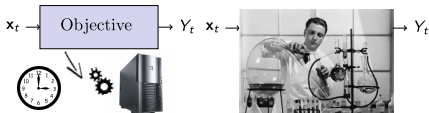
Bayesian Optimization: Common Features

- Very expensive evaluations.



Bayesian Optimization: Common Features

- Very expensive evaluations.

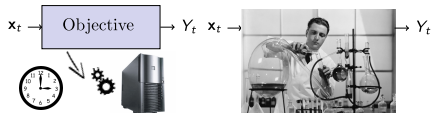


- The objective is a black-box.



Bayesian Optimization: Common Features

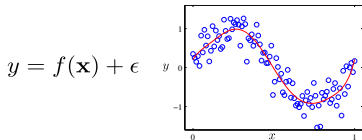
- Very expensive evaluations.



- The objective is a black-box.

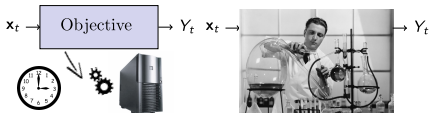


- The evaluation can be noisy.



Bayesian Optimization: Common Features

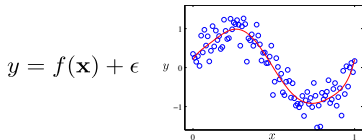
- Very expensive evaluations.



- The objective is a black-box.



- The evaluation can be noisy.



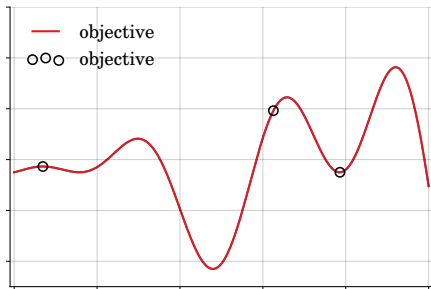
Bayesian optimization methods can be used to solve these problems!

Bayesian Optimization in Practice



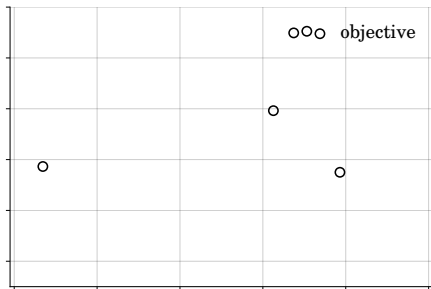
① Get initial sample.

Bayesian Optimization in Practice



① Get initial sample.

Bayesian Optimization in Practice

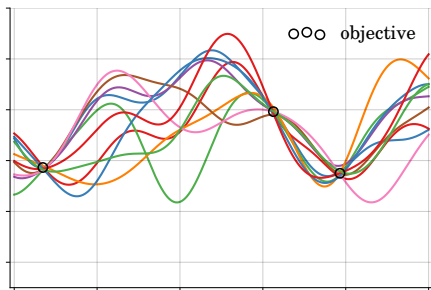


① Get initial sample.

② Fit a model to the data:

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

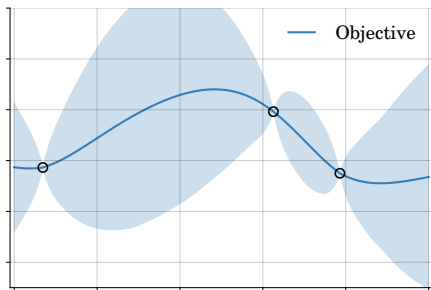
Bayesian Optimization in Practice



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

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

Bayesian Optimization in Practice



① Get initial sample.

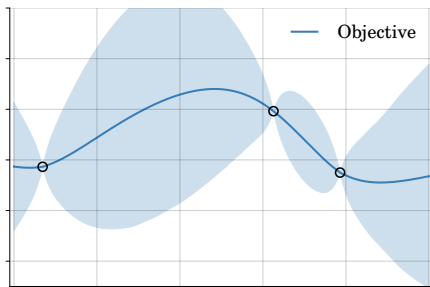
② **Fit a model to the data:**

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

③ Select data collection strategy:

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

Bayesian Optimization in Practice



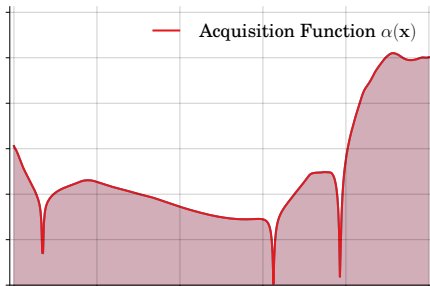
- 1 Get initial sample.
- 2 Fit a model to the data:

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

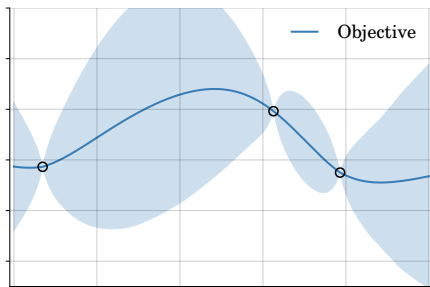
- 3 **Select data collection strategy:**

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

- 4 Optimize acquisition function $\alpha(\mathbf{x})$.



Bayesian Optimization in Practice



1 Get initial sample.

2 Fit a model to the data:

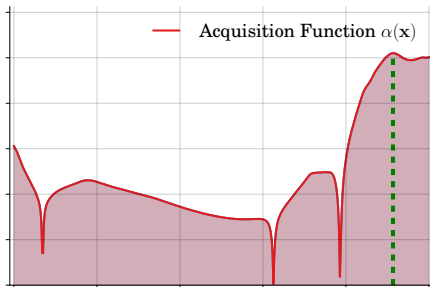
$$p(y|\mathbf{x}, \mathcal{D}_n).$$

3 Select data collection strategy:

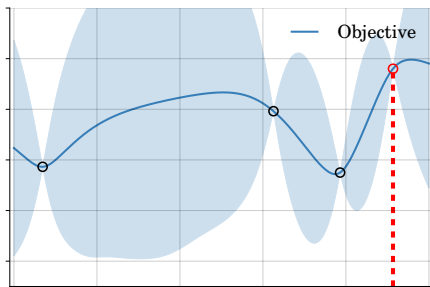
$$\alpha(\mathbf{x}) = \mathbf{E}_{p(y|\mathbf{x}, \mathcal{D}_n)}[U(y|\mathbf{x}, \mathcal{D}_n)].$$

4 Optimize acquisition function $\alpha(\mathbf{x})$.

5 Collect data and update model.



Bayesian Optimization in Practice



1 Get initial sample.

2 Fit a model to the data:

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

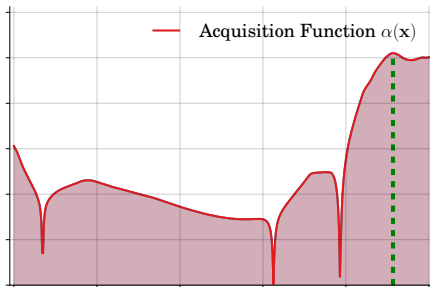
3 Select data collection strategy:

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

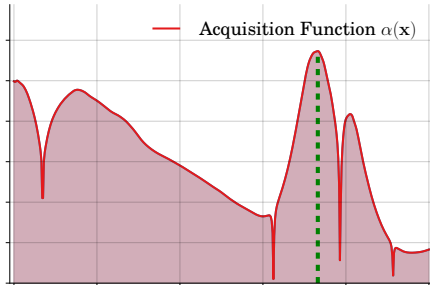
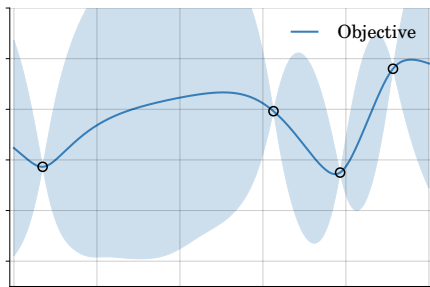
4 Optimize acquisition function $\alpha(\mathbf{x})$.

5 **Collect data and update model.**

6 Repeat!



Bayesian Optimization in Practice



1 Get initial sample.

2 Fit a model to the data:

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

3 Select data collection strategy:

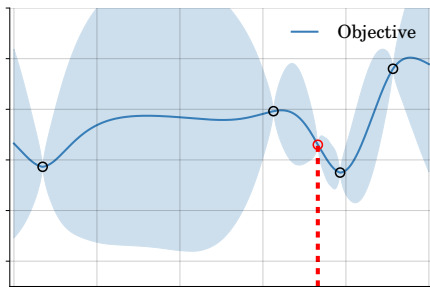
$$\alpha(\mathbf{x}) = \mathbf{E}_{p(y|\mathbf{x}, \mathcal{D}_n)}[U(y|\mathbf{x}, \mathcal{D}_n)].$$

4 Optimize acquisition function $\alpha(\mathbf{x})$.

5 Collect data and update model.

6 Repeat!

Bayesian Optimization in Practice



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

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

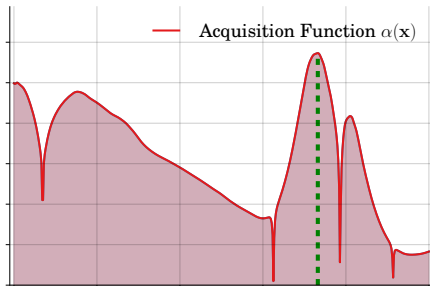
- 3 Select data collection strategy:

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

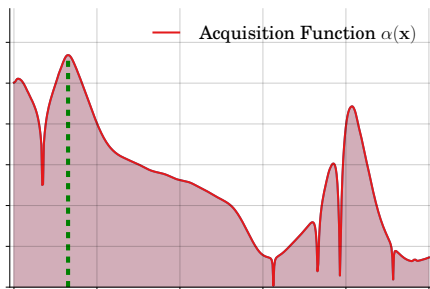
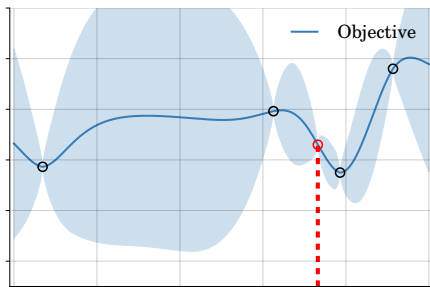
- 4 Optimize acquisition function $\alpha(\mathbf{x})$.

- 5 Collect data and update model.

- 6 Repeat!



Bayesian Optimization in Practice



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

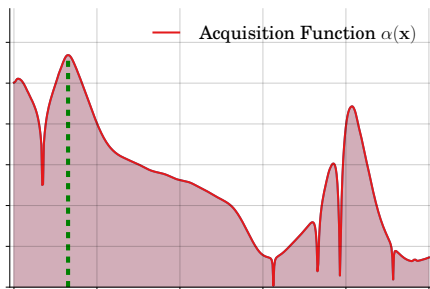
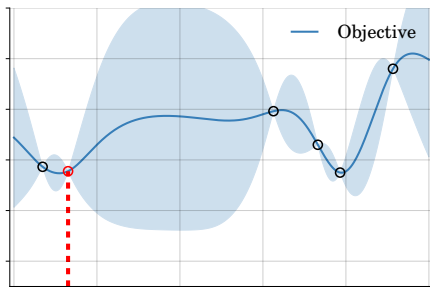
$$p(y|\mathbf{x}, \mathcal{D}_n).$$

- 3 Select data collection strategy:

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

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

Bayesian Optimization in Practice



1 Get initial sample.

2 Fit a model to the data:

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

3 Select data collection strategy:

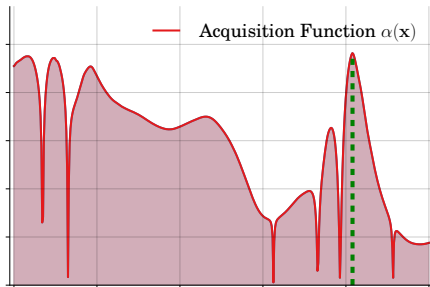
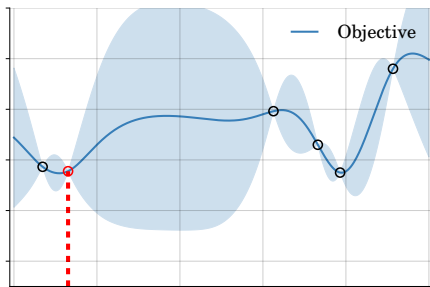
$$\alpha(\mathbf{x}) = \mathbf{E}_{p(y|\mathbf{x}, \mathcal{D}_n)}[U(y|\mathbf{x}, \mathcal{D}_n)].$$

4 Optimize acquisition function $\alpha(\mathbf{x})$.

5 Collect data and update model.

6 Repeat!

Bayesian Optimization in Practice



1 Get initial sample.

2 Fit a model to the data:

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

3 Select data collection strategy:

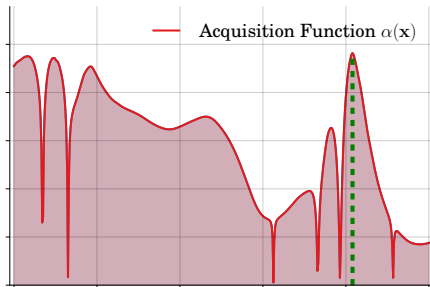
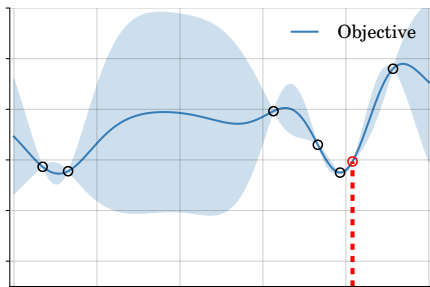
$$\alpha(\mathbf{x}) = \mathbf{E}_{p(y|\mathbf{x}, \mathcal{D}_n)}[U(y|\mathbf{x}, \mathcal{D}_n)].$$

4 Optimize acquisition function $\alpha(\mathbf{x})$.

5 Collect data and update model.

6 Repeat!

Bayesian Optimization in Practice



1 Get initial sample.

2 Fit a model to the data:

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

3 Select data collection strategy:

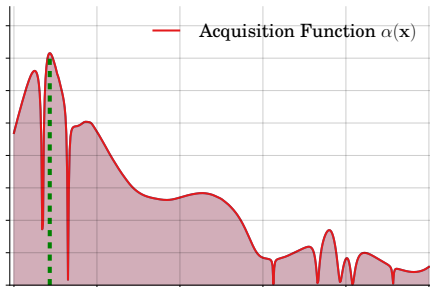
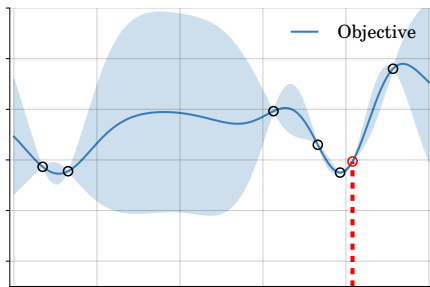
$$\alpha(\mathbf{x}) = \mathbf{E}_{p(y|\mathbf{x}, \mathcal{D}_n)}[U(y|\mathbf{x}, \mathcal{D}_n)].$$

4 Optimize acquisition function $\alpha(\mathbf{x})$.

5 Collect data and update model.

6 Repeat!

Bayesian Optimization in Practice



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

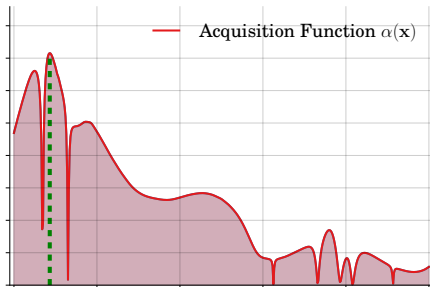
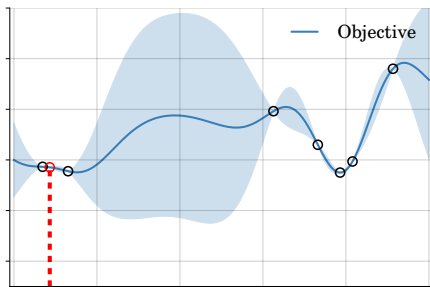
$$p(y|\mathbf{x}, \mathcal{D}_n).$$

- 3 Select data collection strategy:

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

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

Bayesian Optimization in Practice



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

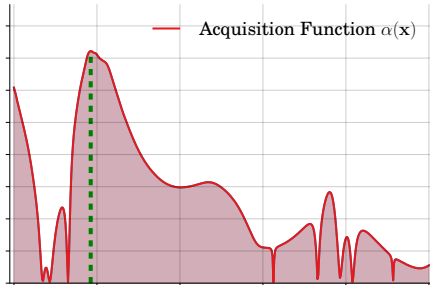
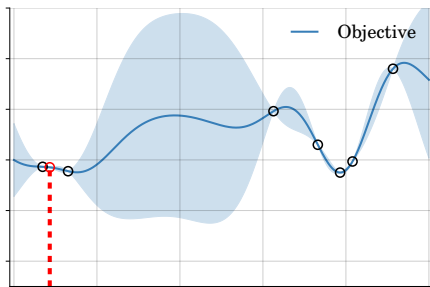
$$p(y|\mathbf{x}, \mathcal{D}_n).$$

- 3 Select data collection strategy:

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

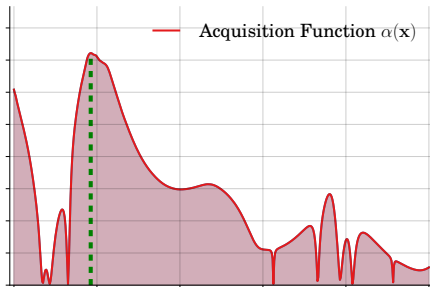
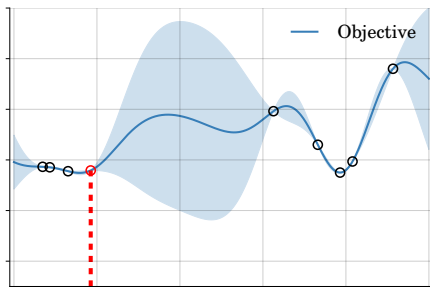
- 4 Optimize acquisition function $\alpha(\mathbf{x})$.
- 5 Collect data and update model.
- 6 Repeat!

Bayesian Optimization in Practice



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

Bayesian Optimization in Practice



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

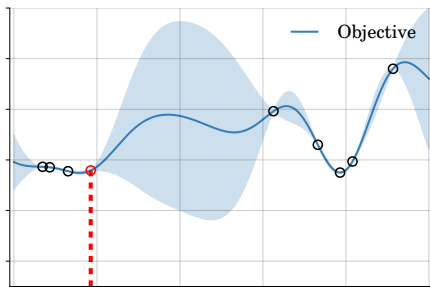
$$p(y|\mathbf{x}, \mathcal{D}_n).$$

- 3 Select data collection strategy:

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

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

Bayesian Optimization in Practice



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

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

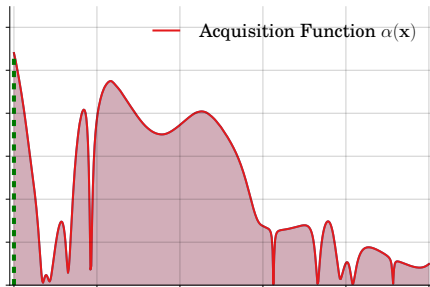
- 3 Select data collection strategy:

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

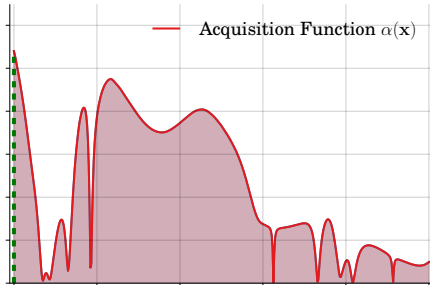
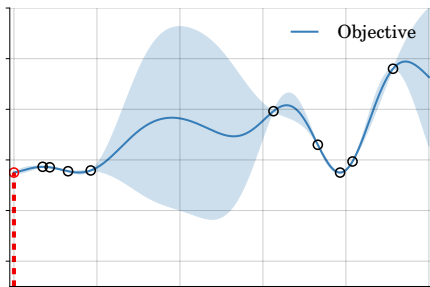
- 4 Optimize acquisition function $\alpha(\mathbf{x})$.

- 5 Collect data and update model.

- 6 Repeat!



Bayesian Optimization in Practice



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

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

- 3 Select data collection strategy:

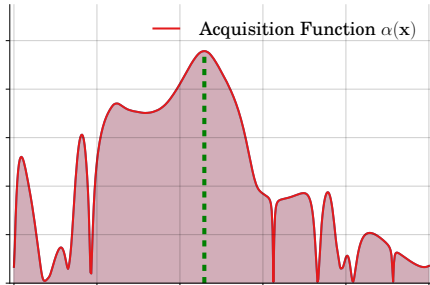
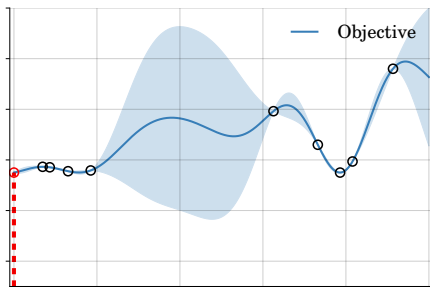
$$\alpha(\mathbf{x}) = \mathbf{E}_{p(y|\mathbf{x}, \mathcal{D}_n)}[U(y|\mathbf{x}, \mathcal{D}_n)].$$

- 4 Optimize acquisition function $\alpha(\mathbf{x})$.

- 5 Collect data and update model.

- 6 Repeat!

Bayesian Optimization in Practice



1 Get initial sample.

2 Fit a model to the data:

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

3 Select data collection strategy:

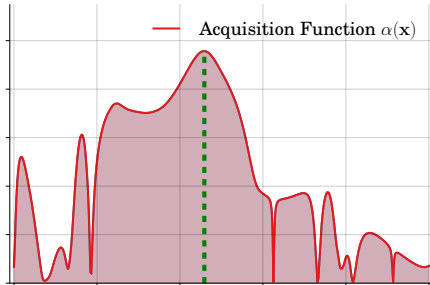
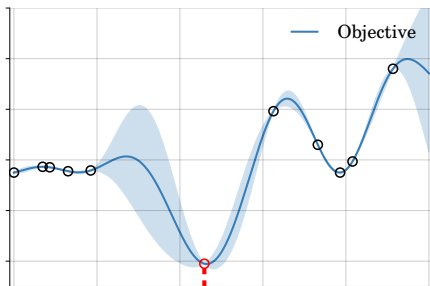
$$\alpha(\mathbf{x}) = \mathbf{E}_{p(y|\mathbf{x}, \mathcal{D}_n)}[U(y|\mathbf{x}, \mathcal{D}_n)].$$

4 Optimize acquisition function $\alpha(\mathbf{x})$.

5 Collect data and update model.

6 Repeat!

Bayesian Optimization in Practice



1 Get initial sample.

2 Fit a model to the data:

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

3 Select data collection strategy:

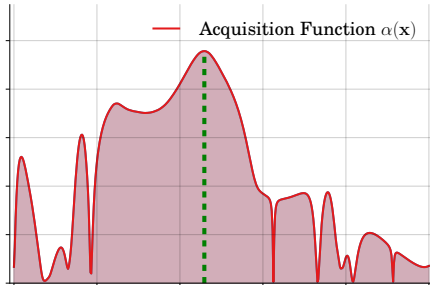
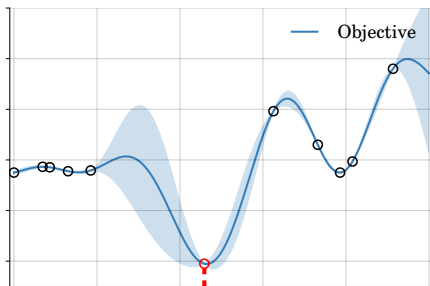
$$\alpha(\mathbf{x}) = \mathbf{E}_{p(y|\mathbf{x}, \mathcal{D}_n)}[U(y|\mathbf{x}, \mathcal{D}_n)].$$

4 Optimize acquisition function $\alpha(\mathbf{x})$.

5 Collect data and update model.

6 Repeat!

Bayesian Optimization in Practice



1 Get initial sample.

2 Fit a model to the data:

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

3 Select data collection strategy:

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

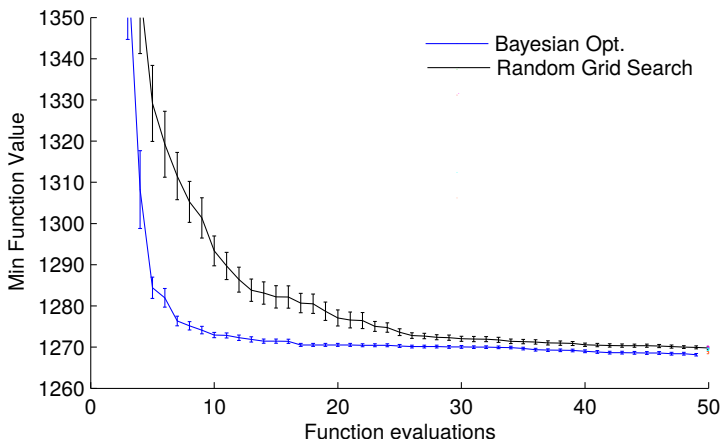
4 Optimize acquisition function $\alpha(\mathbf{x})$.

5 Collect data and update model.

6 Repeat!

The model guides the search focusing on the most-promising regions of the input space!

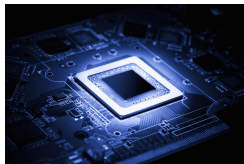
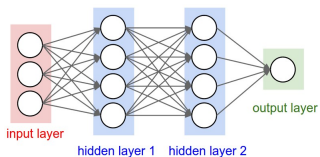
Bayesian Optimization vs. Uniform Exploration



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

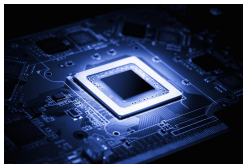
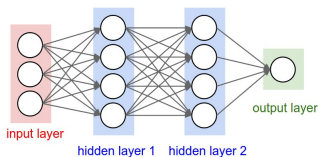
Several Objectives and Constraints

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



Several Objectives and Constraints

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

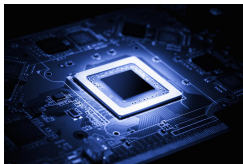
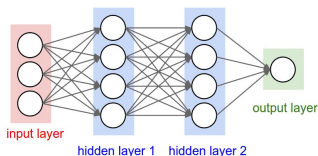


Goals:

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

Several Objectives and Constraints

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



Goals:

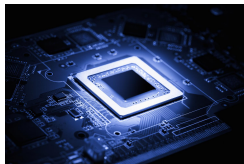
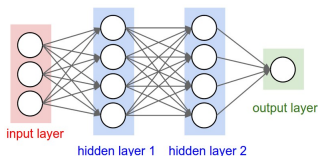
- Minimize **prediction error**.
- Minimize **prediction time**.

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.

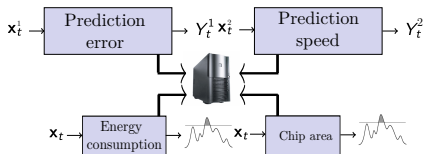


Goals:

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

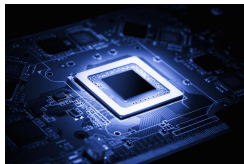
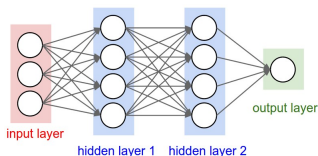
Constrained to:

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



Several Objectives and Constraints

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

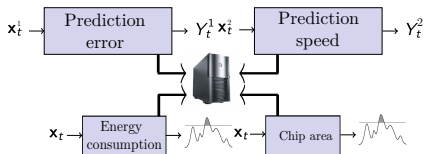


Goals:

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

Constrained to:

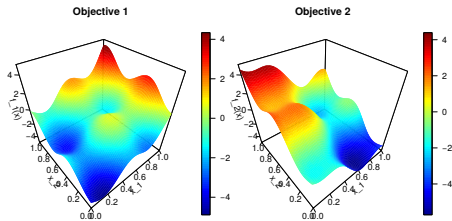
- **Chip area** below a value.
- **Power consumption** below a level.



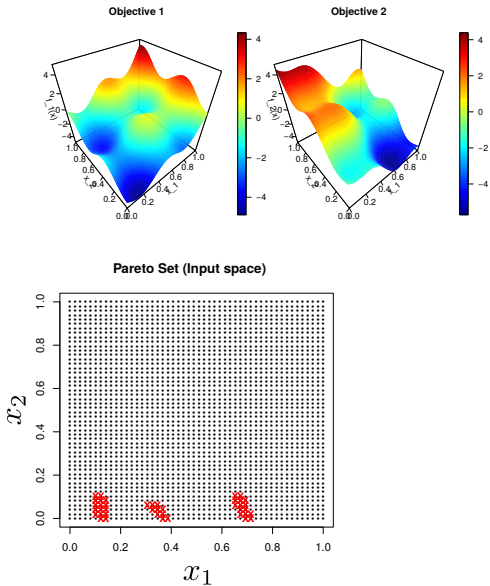
Challenges:

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

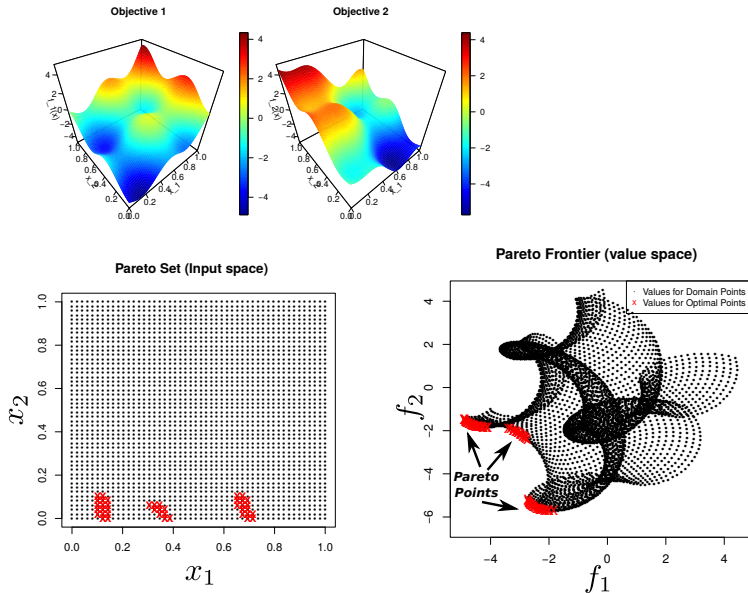
Constrained Multi-Objective Optimization



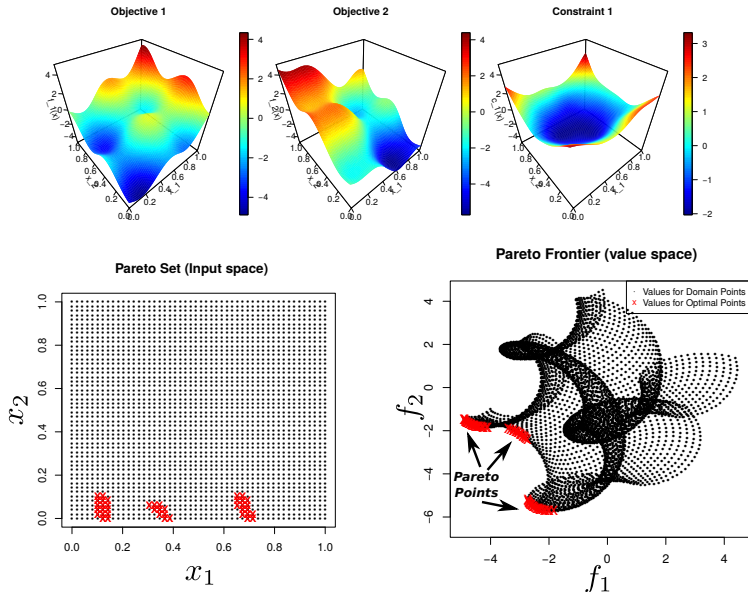
Constrained Multi-Objective Optimization



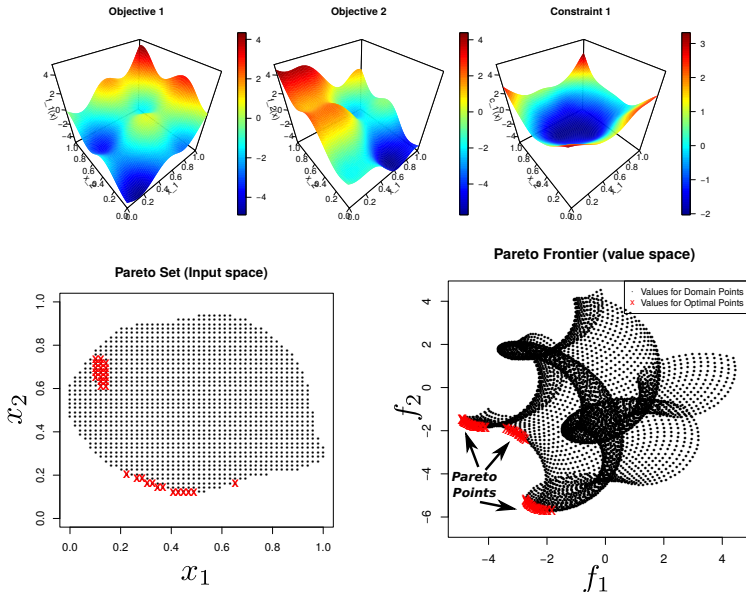
Constrained Multi-Objective Optimization



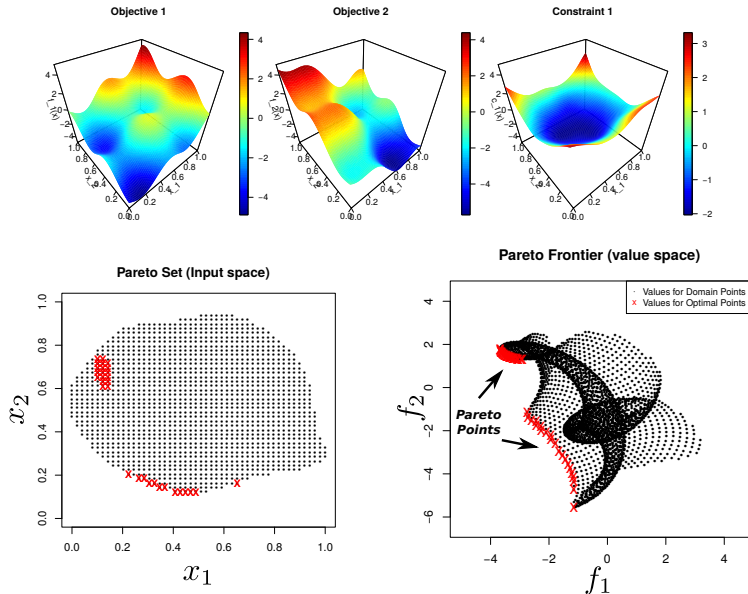
Constrained Multi-Objective Optimization



Constrained Multi-Objective Optimization



Constrained Multi-Objective Optimization



Information-based Approach

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

Information-based Approach

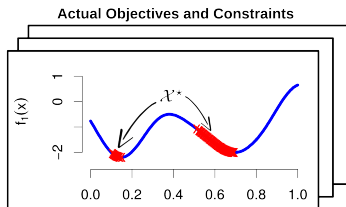
The Pareto set \mathcal{X}^* in the feasible space is a **random variable!**

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**!

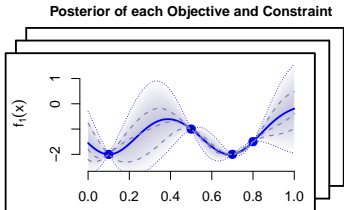
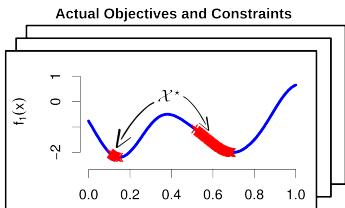
Information is measured by the **entropy** of $p(\mathcal{X}^*|\mathcal{D}_N)$.



Information-based Approach

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

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

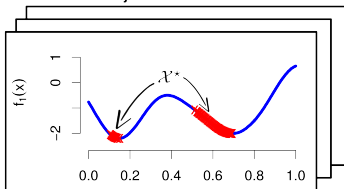


Information-based Approach

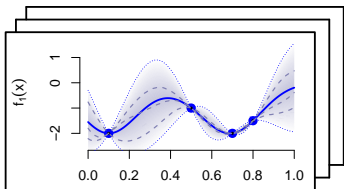
The Pareto set \mathcal{X}^* in the feasible space is a **random variable**!

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

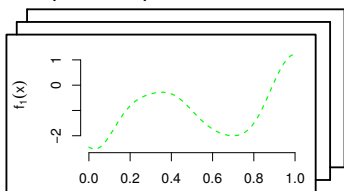
Actual Objectives and Constraints



Posterior of each Objective and Constraint



Optimized Samples Drawn from the Posterior

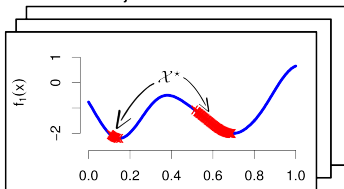


Information-based Approach

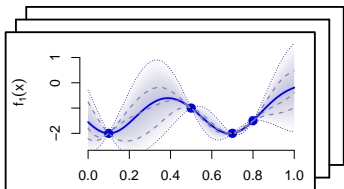
The Pareto set \mathcal{X}^* in the feasible space is a **random variable**!

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

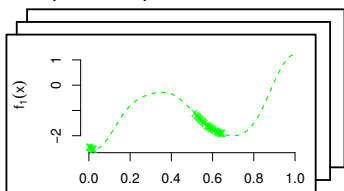
Actual Objectives and Constraints



Posterior of each Objective and Constraint



Optimized Samples Drawn from the Posterior

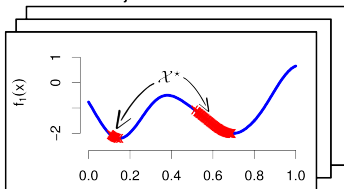


Information-based Approach

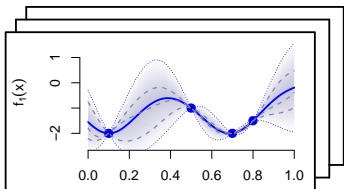
The Pareto set \mathcal{X}^* in the feasible space is a **random variable**!

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

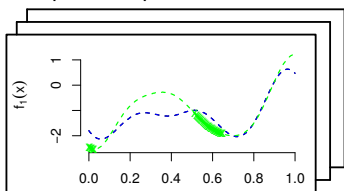
Actual Objectives and Constraints



Posterior of each Objective and Constraint



Optimized Samples Drawn from the Posterior

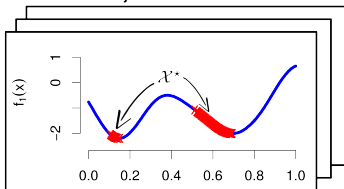


Information-based Approach

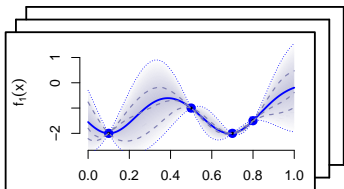
The Pareto set \mathcal{X}^* in the feasible space is a **random variable**!

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

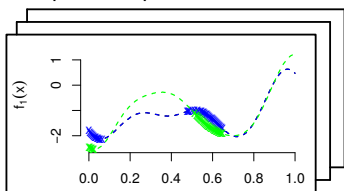
Actual Objectives and Constraints



Posterior of each Objective and Constraint



Optimized Samples Drawn from the Posterior

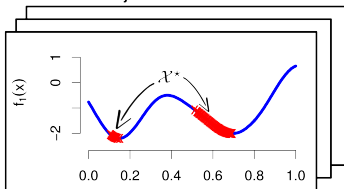


Information-based Approach

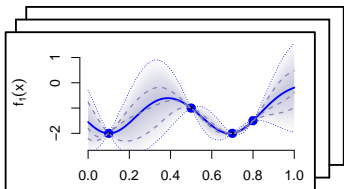
The Pareto set \mathcal{X}^* in the feasible space is a **random variable**!

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

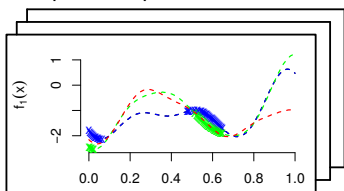
Actual Objectives and Constraints



Posterior of each Objective and Constraint



Optimized Samples Drawn from the Posterior

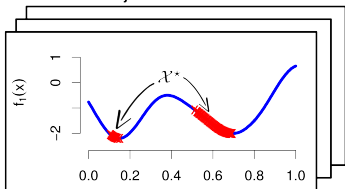


Information-based Approach

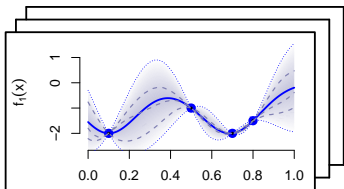
The Pareto set \mathcal{X}^* in the feasible space is a **random variable**!

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

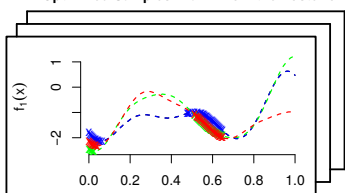
Actual Objectives and Constraints



Posterior of each Objective and Constraint



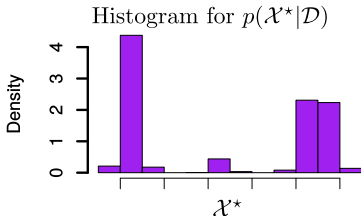
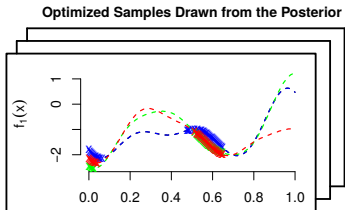
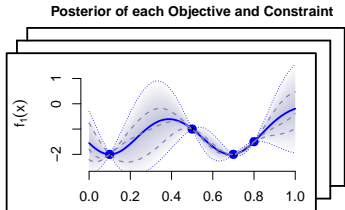
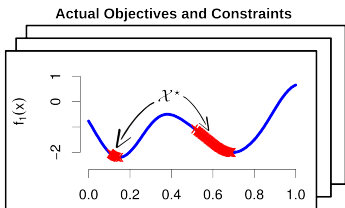
Optimized Samples Drawn from the Posterior



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

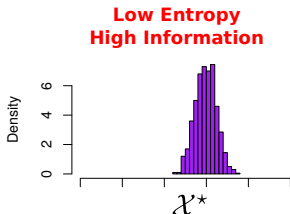
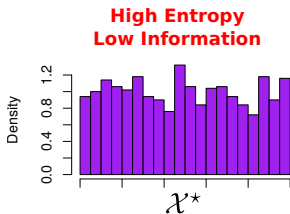
The Pareto set \mathcal{X}^* in the feasible space is a **random variable!**

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!**

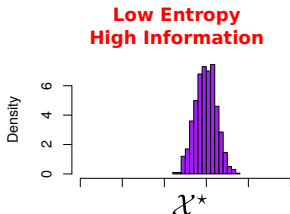
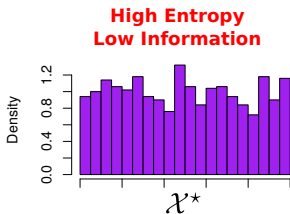
Information is measured by the **entropy** of $p(\mathcal{X}^*|\mathcal{D}_N)$.



Information-based Approach

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

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



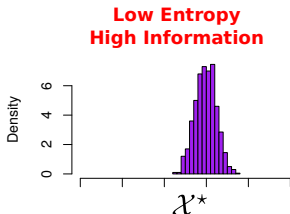
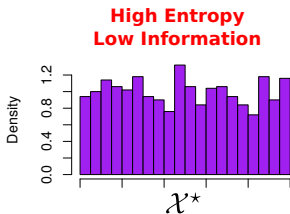
The acquisition function is

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

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

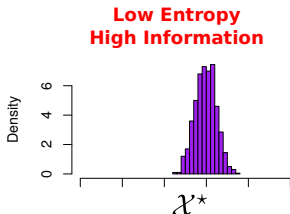
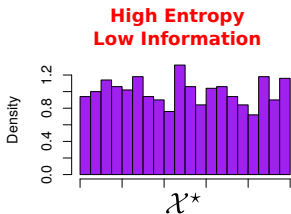
$$\alpha(\mathbf{x}) = \mathbb{H}[\mathcal{X}^*|\mathcal{D}_t] - \mathbb{E}_{\mathbf{y}} \left[\mathbb{H}[\mathcal{X}^*|\mathcal{D}_t \cup \{\mathbf{x}, \mathbf{y}\}] \mid \mathcal{D}_t, \mathbf{x} \right] \quad (1)$$

How much we know
about \mathcal{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

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

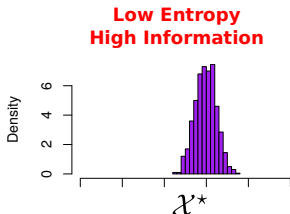
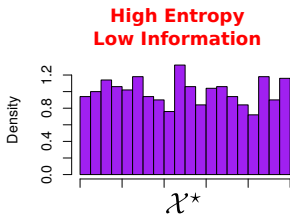
How much we know about \mathcal{X}^* now.

How much we will know about \mathcal{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

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

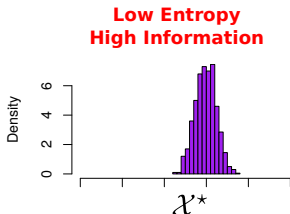
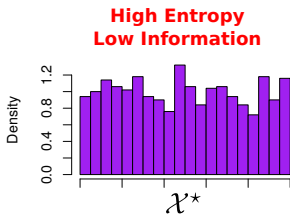
How much we know about \mathcal{X}^* now.

How much we will know about \mathcal{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

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

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

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

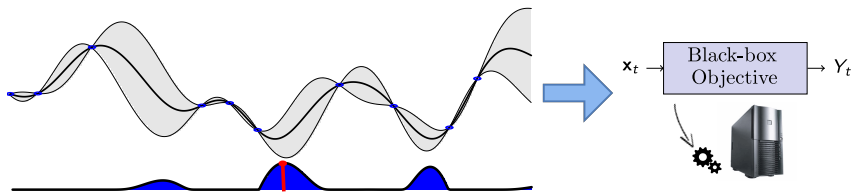
Computing (1) is **very difficult in practice!**

Parallel Bayesian Optimization

Traditional Bayesian optimization is **sequential!**

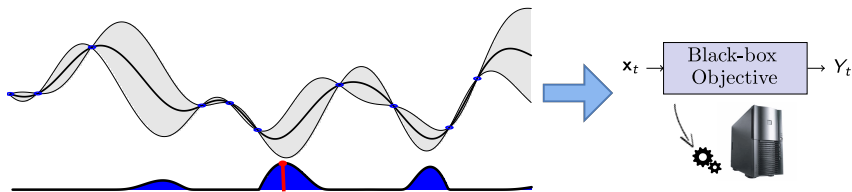
Parallel Bayesian Optimization

Traditional Bayesian optimization is **sequential!**



Parallel Bayesian Optimization

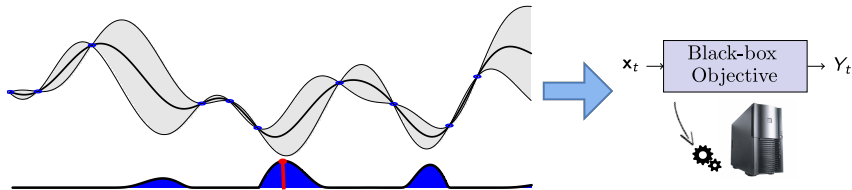
Traditional Bayesian optimization is **sequential**!



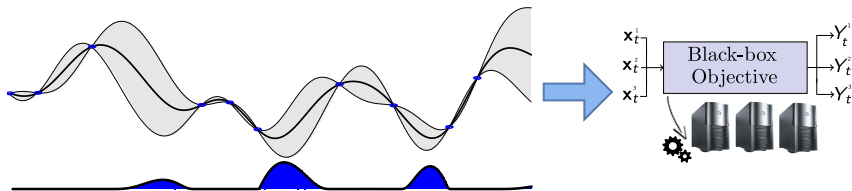
Computing clusters let us do **many things** at once!

Parallel Bayesian Optimization

Traditional Bayesian optimization is **sequential**!

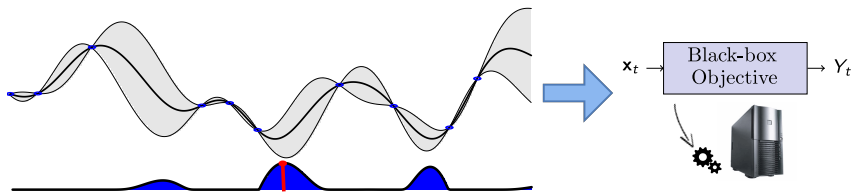


Computing clusters let us do **many things** at once!

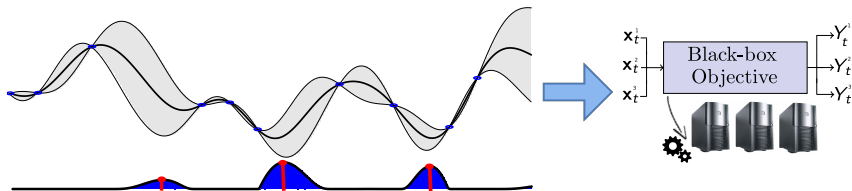


Parallel Bayesian Optimization

Traditional Bayesian optimization is **sequential**!

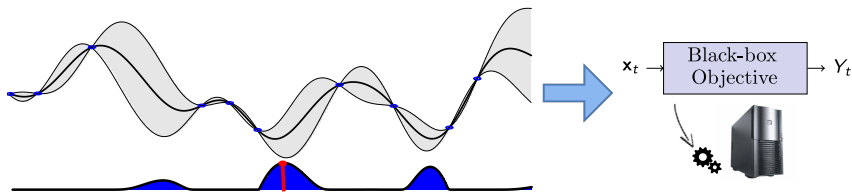


Computing clusters let us do **many things** at once!

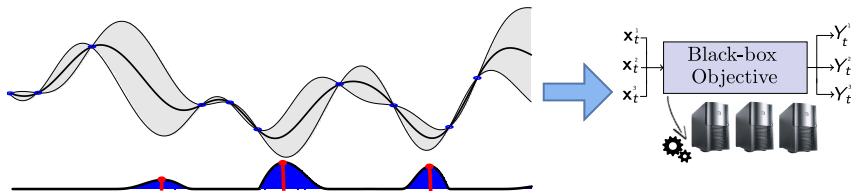


Parallel Bayesian Optimization

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 of Q points $\mathbf{X} = \{\mathbf{x}_q\}_{q=1}^Q$ to minimize the entropy of \mathcal{X}^* .

Parallel Predictive Entropy Search

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

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

Parallel Predictive Entropy Search

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

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

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

Parallel Predictive Entropy Search

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

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

Parallel Predictive Entropy Search

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

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

Parallel Predictive Entropy Search

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

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

Parallel Predictive Entropy Search

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

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

Gaussian distribution

Parallel Predictive Entropy Search

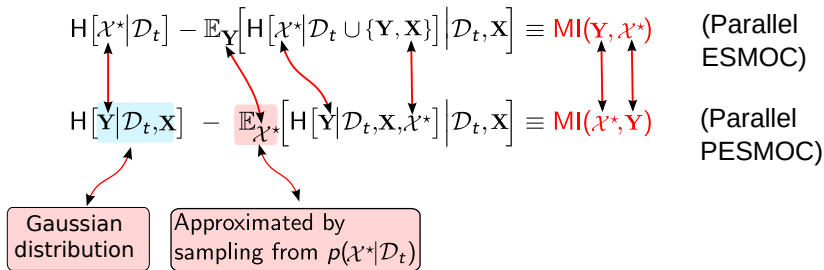
Choose a set of Q points $\mathbf{X} = \{\mathbf{x}_q\}_{q=1}^Q$ to minimize the entropy of \mathcal{X}^* .

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

Gaussian distribution

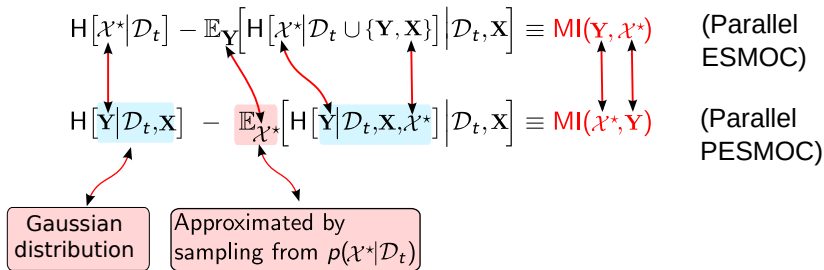
Parallel Predictive Entropy Search

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



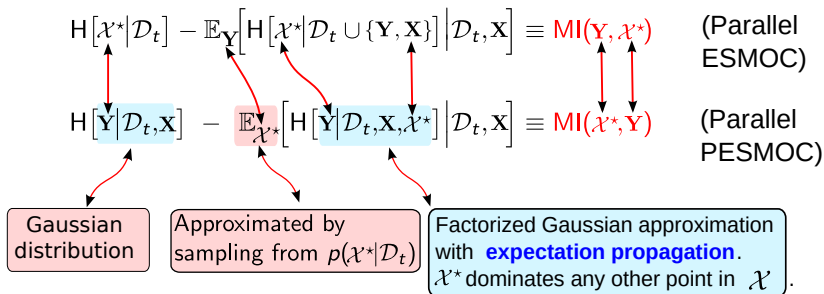
Parallel Predictive Entropy Search

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



Parallel Predictive Entropy Search

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



(Minka, 2001)

Parallel Predictive Entropy Search

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

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

$$H[\mathbf{Y} | \mathcal{D}_t, \mathbf{X}] - \mathbb{E}_{\mathcal{X}^*} \left[H[\mathbf{Y} | \mathcal{D}_t, \mathbf{X}, \mathcal{X}^*] \middle| \mathcal{D}_t, \mathbf{X} \right] \equiv \text{MI}(\mathcal{X}^*, \mathbf{Y}) \quad (\text{Parallel PESMOC})$$

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} .

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

(Minka, 2001)

Parallel Predictive Entropy Search

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

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

$$H[\mathbf{Y} | \mathcal{D}_t, \mathbf{X}] - \mathbb{E}_{\mathcal{X}^*} \left[H[\mathbf{Y} | \mathcal{D}_t, \mathbf{X}, \mathcal{X}^*] \middle| \mathcal{D}_t, \mathbf{X} \right] \equiv \text{MI}(\mathcal{X}^*, \mathbf{Y}) \quad (\text{Parallel PESMOC})$$

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} .

$$\alpha(\mathbf{x}) \approx \sum_{c=1}^C \log |\mathbf{V}_c^{\text{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^M \left(\sum_{c=1}^C \log |\mathbf{V}_c^{\text{CPD}}(\mathbf{X} | \mathcal{X}_m^*)| \right) + \sum_{k=1}^K \log |\mathbf{V}_k^{\text{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^M \left(\sum_{k=1}^K \log |\mathbf{V}_k^{\text{CPD}}(\mathbf{X} | \mathcal{X}_m^*)| \right) = \sum_{i=1}^{K+C} \alpha_k(\mathbf{X})$$

(Minka, 2001)

Parallel Predictive Entropy Search

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

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

$$H[\mathbf{Y} | \mathcal{D}_t, \mathbf{X}] - \mathbb{E}_{\mathcal{X}^*} \left[H[\mathbf{Y} | \mathcal{D}_t, \mathbf{X}, \mathcal{X}^*] \middle| \mathcal{D}_t, \mathbf{X} \right] \equiv \text{MI}(\mathcal{X}^*, \mathbf{Y}) \quad (\text{Parallel PESMOC})$$

Gaussian distribution

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

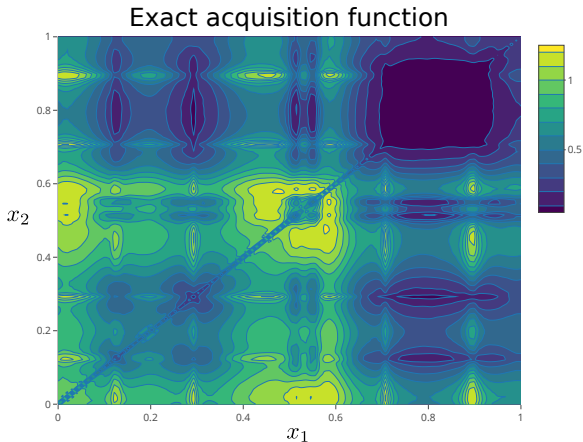
Factorized Gaussian approximation with **expectation propagation**. One acquisition per black-box. \mathcal{X}^* dominates any other.

$$\alpha(\mathbf{x}) \approx \sum_{c=1}^C \log |\mathbf{V}_c^{\text{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^M \left(\sum_{c=1}^C \log |\mathbf{V}_c^{\text{CPD}}(\mathbf{X} | \mathcal{X}_m^*)| \right) + \sum_{k=1}^K \log |\mathbf{V}_k^{\text{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^M \left(\sum_{k=1}^K \log |\mathbf{V}_k^{\text{CPD}}(\mathbf{X} | \mathcal{X}_m^*)| \right) = \sum_{i=1}^{K+C} \alpha_k(\mathbf{X})$$

(Minka, 2001)

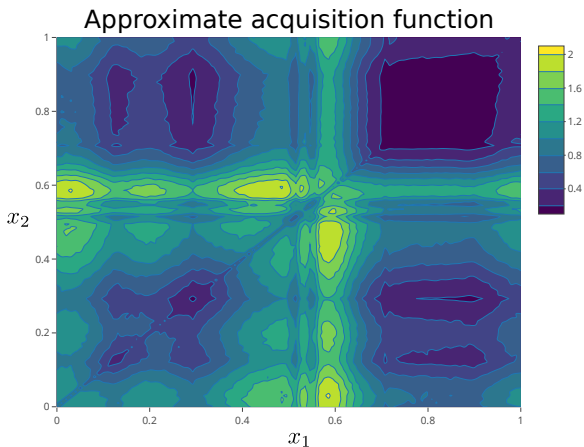
Comparison of Exact and Approximate Acquisition

1 dimensional problem, Batch size = 2:



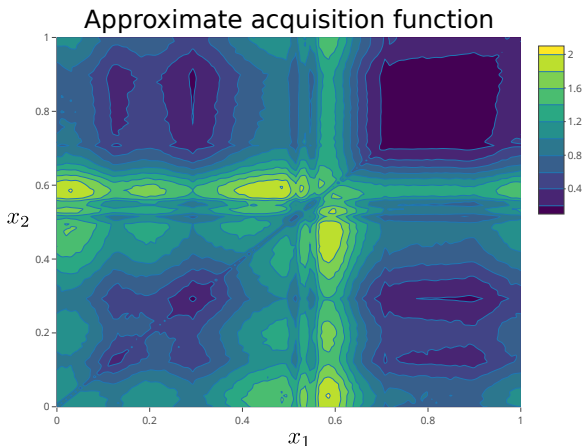
Comparison of Exact and Approximate Acquisition

1 dimensional problem, Batch size = 2:



Comparison of Exact and Approximate Acquisition

1 dimensional problem, Batch size = 2:



The acquisition is symmetric and the approximation is large where the exact acquisition is large, as expected!

Related Methods

- Parallel Sequential: Transforms any sequential BO method into a batch method.

Related Methods

- Parallel Sequential: Transforms any sequential BO method into a batch method.
 - Repeats the optimization of the acquisition B times.

Related Methods

- Parallel Sequential: Transforms any sequential BO method into a batch method.
 - Repeats the optimization of the acquisition B times.
 - Models are updated. Pending points assigned the predictive mean.

Related Methods

- Parallel Sequential: Transforms any sequential BO method into a batch method.
 - Repeats the optimization of the acquisition B times.
 - Models are updated. Pending points assigned the predictive mean.

Expected to be very expensive for large batch sizes B !

Related Methods

- Parallel Sequential: Transforms any sequential BO method into a batch method.
 - Repeats the optimization of the acquisition B times.
 - Models are updated. Pending points assigned the predictive mean.

Expected to be very expensive for large batch sizes B !

- Expected Hyper-volume Improvement Strategies:

Related Methods

- Parallel Sequential: Transforms any sequential BO method into a batch method.
 - Repeats the optimization of the acquisition B times.
 - Models are updated. Pending points assigned the predictive mean.

Expected to be very expensive for large batch sizes B !

- Expected Hyper-volume Improvement Strategies:
 - Two versions: qEHVI (noiseless evals.) and qNEHVI (noisy evals.).

Related Methods

- Parallel Sequential: Transforms any sequential BO method into a batch method.
 - Repeats the optimization of the acquisition B times.
 - Models are updated. Pending points assigned the predictive mean.

Expected to be very expensive for large batch sizes B !

- Expected Hyper-volume Improvement Strategies:
 - Two versions: qEHVI (noiseless evals.) and qNEHVI (noisy evals.).
 - Constraints incorporated by multiplying by the feasibility prob.

Related Methods

- Parallel Sequential: Transforms any sequential BO method into a batch method.
 - Repeats the optimization of the acquisition B times.
 - Models are updated. Pending points assigned the predictive mean.

Expected to be very expensive for large batch sizes B !

- Expected Hyper-volume Improvement Strategies:
 - Two versions: qEHVI (noiseless evals.) and qNEHVI (noisy evals.).
 - Constraints incorporated by multiplying by the feasibility prob.
 - Expectations approximated by Monte Carlo.

Related Methods

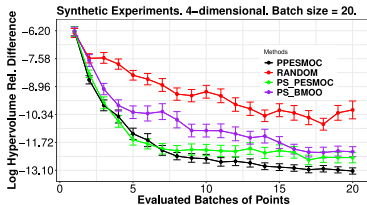
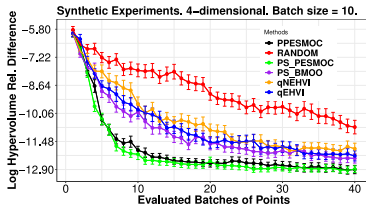
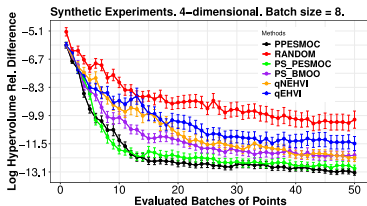
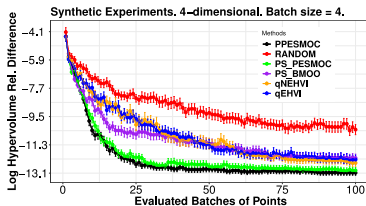
- Parallel Sequential: Transforms any sequential BO method into a batch method.
 - Repeats the optimization of the acquisition B times.
 - Models are updated. Pending points assigned the predictive mean.

Expected to be very expensive for large batch sizes B !

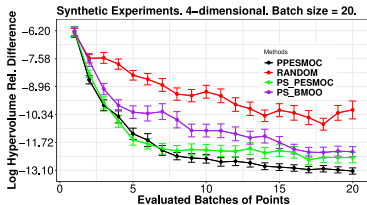
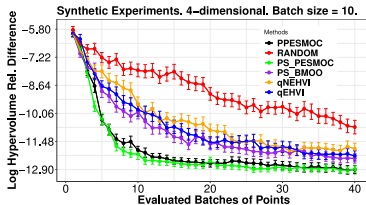
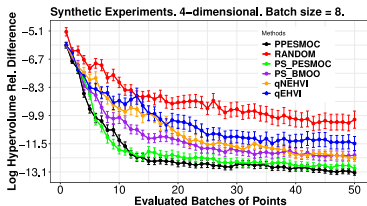
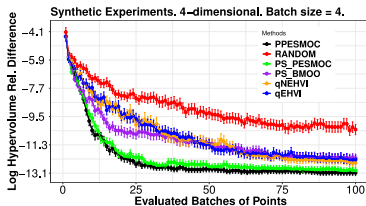
- Expected Hyper-volume Improvement Strategies:
 - Two versions: qEHVI (noiseless evals.) and qNEHVI (noisy evals.).
 - Constraints incorporated by multiplying by the feasibility prob.
 - Expectations approximated by Monte Carlo.

MC approximation is zero after a few evaluations and they have high cost w.r.t. B (even exponential)!

Synthetic Experiments



Synthetic Experiments



PPSMOC performs better than or similar to the other strategies!

Time to Choose the next Batch

Table 1

Mean of the time in seconds to choose the next batch of points by PPESMOC and the parallel sequential approaches. For $B = 50$, underlined results are significantly different with respect to PPESMOC results according to the Wilcoxon test at $\alpha = 0.05$.

Method	$B = 4$	$B = 8$	$B = 10$	$B = 20$	$B = 50$
PPESMOC	696.0 ± 26.9	912.74 ± 26.3	957.3 ± 25.7	1045.7 ± 30.53	1269.35 26.62
PS_PPESMOC	191.5 ± 7.0	347.2 \pm 6.0	405.49 \pm 5.8	801.05 \pm 27.8	<u>1957.72</u> <u>34.1</u>
PS_BMOO	379.4 ± 13.1	551.1 ± 21.7	593.86 ± 18.0	897.4 ± 29.6	<u>1870.42</u> <u>42.77</u>
qEHVI	65.2 \pm 1.8	417.9 ± 21.9	1174.9 ± 54.3		
qNEHVI	89.5 ± 2.3	401.4 ± 23.9	1169.4 ± 56.1		

Time to Choose the next Batch

Table 1

Mean of the time in seconds to choose the next batch of points by PPESMOC and the parallel sequential approaches. For $B = 50$, underlined results are significantly different with respect to PPESMOC results according to the Wilcoxon test at $\alpha = 0.05$.

Method	$B = 4$	$B = 8$	$B = 10$	$B = 20$	$B = 50$
PPESMOC	696.0 \pm 26.9	912.74 \pm 26.3	957.3 \pm 25.7	1045.7 \pm 30.53	1269.35 26.62
PS_PESMOC	191.5 \pm 7.0	347.2 \pm 6.0	405.49 \pm 5.8	801.05 \pm 27.8	<u>1957.72</u> <u>34.1</u>
PS_BMOO	379.4 \pm 13.1	551.1 \pm 21.7	593.86 \pm 18.0	897.4 \pm 29.6	<u>1870.42</u> <u>42.77</u>
qEHVI	65.2 \pm 1.8	417.9 \pm 21.9	1174.9 \pm 54.3		
qNEHVI	89.5 \pm 2.3	401.4 \pm 23.9	1169.4 \pm 56.1		

PPESMOC scales significantly better w.r.t. the batch size B for large values of B !

Optimal Ensemble on the German Dataset

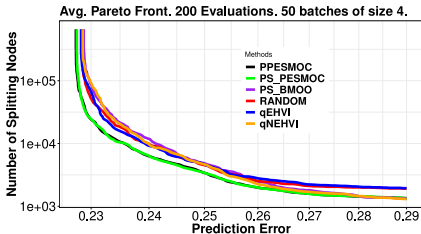
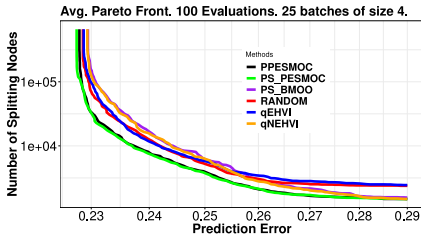


Table 2

Average hyper-volume in the task of finding an optimal ensemble of trees. Underlined results are significantly different with respect to PPESMOC results according to the Wilcoxon test at $\alpha = 0.05$.

# Eval.	PPESMOC	PS_PESMOC	PS_BMOO	P_RANDOM	qNEHVI	qEHVI
100	0.325 ± 0.007	0.327 ± 0.007	<u>0.295 ± 0.014</u>	<u>0.298 ± 0.009</u>	<u>0.299 ± 0.011</u>	<u>0.294 ± 0.013</u>
200	0.334 ± 0.005	0.335 ± 0.006	<u>0.313 ± 0.010</u>	<u>0.310 ± 0.007</u>	<u>0.3154 ± 0.008</u>	<u>0.309 ± 0.010</u>

Optimal Ensemble on the German Dataset

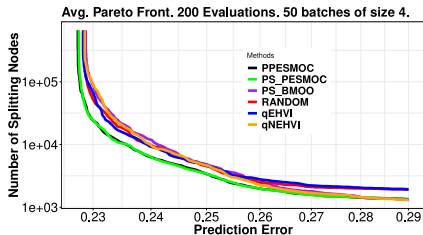
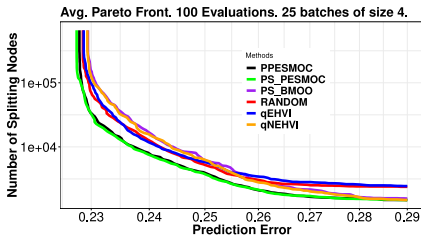


Table 2

Average hyper-volume in the task of finding an optimal ensemble of trees. Underlined results are significantly different with respect to PPESMOC results according to the Wilcoxon test at $\alpha = 0.05$.

# Eval.	PPESMOC	PS_PESMOC	PS_BMOO	P_RANDOM	qNEHVI	qEHVI
100	0.325 \pm 0.007	0.327 \pm 0.007	<u>0.295 \pm 0.014</u>	<u>0.298 \pm 0.009</u>	<u>0.299 \pm 0.011</u>	<u>0.294 \pm 0.013</u>
200	0.334 \pm 0.005	0.335 \pm 0.006	<u>0.313 \pm 0.010</u>	<u>0.310 \pm 0.007</u>	<u>0.3154 \pm 0.008</u>	<u>0.309 \pm 0.010</u>

PPESMOC performs better than or similar to the other strategies!

Optimal Neural Network on MNIST

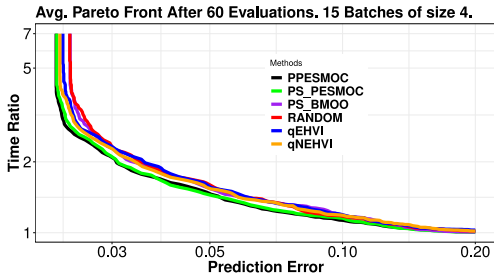


Table 3

Avg. hyper-volume of each method in the neural network experiment. Underlined results are significantly different with respect to PPESMOC results according to the Wilcoxon test at $\alpha = 0.05$.

# Eval.	PPESMOC	PS_PESMOC	PS_BMOO	P_RANDOM	qNEHVI	qEHVI
60	1.020 ± 0.014	1.014 ± 0.029	<u>0.982 ± 0.095</u>	<u>0.993 ± 0.035</u>	<u>0.999 ± 0.050</u>	<u>0.996 ± 0.041</u>

Optimal Neural Network on MNIST

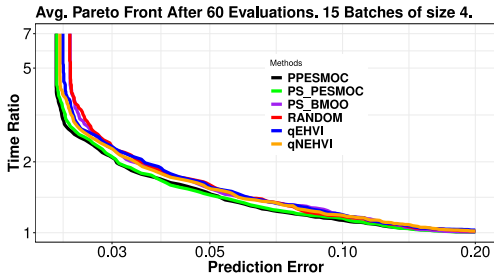


Table 3

Avg. hyper-volume of each method in the neural network experiment. Underlined results are significantly different with respect to PPESMOC results according to the Wilcoxon test at $\alpha = 0.05$.

# Eval.	PPESMOC	PS_PESMOC	PS_BMOO	P_RANDOM	qNEHVI	qEHVI
60	1.020 ± 0.014	1.014 ± 0.029	<u>0.982 ± 0.095</u>	<u>0.993 ± 0.035</u>	<u>0.999 ± 0.050</u>	<u>0.996 ± 0.041</u>

PPESMOC performs slightly better than the other strategies!

Conclusions

- PPESMOC is the first entropy-based BO method for several objectives and constraints and parallel black-box evaluations.

Conclusions

- PPESMOC is the first entropy-based BO method for several objectives and constraints and parallel black-box evaluations.
- PPESMOC performs similar or better than other methods from the literature having a smaller computational cost w.r.t. the batch size.

Conclusions

- PPESMOC is the first entropy-based BO method for several objectives and constraints and parallel black-box evaluations.
- PPESMOC performs similar or better than other methods from the literature having a smaller computational cost w.r.t. the batch size.
- If the batch size B is small, Parallel Sequential methods based on PESMOC may be the better approach.

Conclusions

- PPESMOC is the first entropy-based BO method for several objectives and constraints and parallel black-box evaluations.
- PPESMOC performs similar or better than other methods from the literature having a smaller computational cost w.r.t. the batch size.
- If the batch size B is small, Parallel Sequential methods based on PESMOC may be the better approach.

Thank you for your attention!



e l l i s

UNIT
MADRID



Comunidad
de Madrid

Partially funded by the Autonomous Community of Madrid

References I

- Daulton, S., Balandat, M., & Bakshy, E. (2020). Differentiable expected hypervolume improvement for parallel multi-objective Bayesian optimization. In *Advances in neural information processing systems* (pp. 9851–9864).
- Daulton, S., Balandat, M., & Bakshy, E. (2021). Parallel Bayesian optimization of multiple noisy objectives with expected hypervolume improvement. In *Advances in neural information processing systems* (pp. 2187–2200).
- Feliot, P., Bect, J., & Vazquez, E. (2017). A Bayesian approach to constrained single-and multi-objective optimization. *Journal of Global Optimization*, 67, 97–133.
- Garrido-Merchán, E., & Hernández-Lobato, D. (2019). Predictive entropy search for multi-objective Bayesian optimization with constraints. *Neurocomputing*, 361, 50–68.
- Snoek, J., Larochelle, H., & Adams, R. (2012). Practical Bayesian optimization of machine learning algorithms. In *Advances in neural information processing systems* (pp. 2951–2959).

References II

- Villemonteix, J., Vazquez, E., & Walter, E. (2009). An informational approach to the global optimization of expensive-to-evaluate functions. *Journal of Global Optimization*, 44(509).
- Tu, B., Gandy, A., Kantas, N., & Shafei, B. (2022). Joint entropy search for multi-objective bayesian optimization. *Advances in Neural Information Processing Systems*, 35, 9922-9938.

Conditional Predictive Distribution I

The predictions must be compatible with \mathcal{X}^* !

$$p(\mathbf{y}|\mathcal{D}, \mathbf{x}, \mathcal{X}^*) = \int \underbrace{p(\mathbf{y}|\mathbf{f})}_{\text{Noise}} \underbrace{p(\mathbf{f}|\mathcal{D}, \mathbf{x}, \mathcal{X}^*)}_{\text{Noiseless pred. dist.}} d\mathbf{f}$$

Conditional Predictive Distribution I

The predictions must be compatible with \mathcal{X}^* !

$$p(\mathbf{y}|\mathcal{D}, \mathbf{x}, \mathcal{X}^*) = \int \underbrace{p(\mathbf{y}|\mathbf{f})}_{\text{Noise}} \underbrace{p(\mathbf{f}|\mathcal{D}, \mathbf{x}, \mathcal{X}^*)}_{\text{Noiseless pred. dist.}} d\mathbf{f}$$

Noiseless conditional predictive distribution:

$$p(\mathbf{f}|\mathcal{D}, \mathbf{x}, \mathcal{X}^*) \propto \int \underbrace{p(\mathbf{f}|\mathbf{x}, \mathcal{F})}_{\text{Black-box values at } \mathbf{x}} \times \underbrace{p(\mathcal{F}|\mathcal{D})}_{\text{Post. dist.}} \times \underbrace{p(\mathcal{X}^*|\mathcal{F})}_{\text{Guarantees } \mathcal{X}^* \text{ optimal.}} d\mathcal{F}$$

where \mathcal{F} informally represents all potential black-box functions.

Conditional Predictive Distribution I

The predictions must be compatible with \mathcal{X}^* !

$$p(\mathbf{y}|\mathcal{D}, \mathbf{x}, \mathcal{X}^*) = \int \underbrace{p(\mathbf{y}|\mathbf{f})}_{\text{Noise}} \underbrace{p(\mathbf{f}|\mathcal{D}, \mathbf{x}, \mathcal{X}^*)}_{\text{Noiseless pred. dist.}} d\mathbf{f}$$

Noiseless conditional predictive distribution:

$$p(\mathbf{f}|\mathcal{D}, \mathbf{x}, \mathcal{X}^*) \propto \int \underbrace{p(\mathbf{f}|\mathbf{x}, \mathcal{F})}_{\text{Black-box values at } \mathbf{x}} \times \underbrace{p(\mathcal{F}|\mathcal{D})}_{\text{Post. dist.}} \times \underbrace{p(\mathcal{X}^*|\mathcal{F})}_{\text{Guarantees } \mathcal{X}^* \text{ optimal.}} d\mathcal{F}$$

where \mathcal{F} informally represents all potential black-box functions.

- Unconditional posterior.

Conditional Predictive Distribution I

The predictions must be compatible with \mathcal{X}^* !

$$p(\mathbf{y}|\mathcal{D}, \mathbf{x}, \mathcal{X}^*) = \int \underbrace{p(\mathbf{y}|\mathbf{f})}_{\text{Noise}} \underbrace{p(\mathbf{f}|\mathcal{D}, \mathbf{x}, \mathcal{X}^*)}_{\text{Noiseless pred. dist.}} d\mathbf{f}$$

Noiseless conditional predictive distribution:

$$p(\mathbf{f}|\mathcal{D}, \mathbf{x}, \mathcal{X}^*) \propto \int \underbrace{p(\mathbf{f}|\mathbf{x}, \mathcal{F})}_{\text{Black-box values at } \mathbf{x}} \times \underbrace{p(\mathcal{F}|\mathcal{D})}_{\text{Post. dist.}} \times \underbrace{p(\mathcal{X}^*|\mathcal{F})}_{\text{Guarantees } \mathcal{X}^* \text{ optimal.}} d\mathcal{F}$$

where \mathcal{F} informally represents all potential black-box functions.

- Unconditional posterior.
- Takes value 1 if \mathcal{X}^* is optimal given \mathcal{F} and zero otherwise.

Conditional Predictive Distribution II

The factor that guarantees optimality is:

$$p(\mathcal{X}^*|\mathcal{F}) = \prod_{\mathbf{x}^* \in \mathcal{X}^*} \left(\left[\prod_{c=1}^C \Theta(\text{cons}_c(\mathbf{x}^*)) \right] \left[\prod_{\mathbf{x}' \in \mathcal{X}} \Omega(\mathbf{x}', \mathbf{x}^*) \right] \right)$$

Conditional Predictive Distribution II

The factor that guarantees optimality is:

$$p(\mathcal{X}^*|\mathcal{F}) = \prod_{\mathbf{x}^* \in \mathcal{X}^*} \left(\left[\prod_{c=1}^C \Theta(\text{cons}_c(\mathbf{x}^*)) \right] \left[\prod_{\mathbf{x}' \in \mathcal{X}} \Omega(\mathbf{x}', \mathbf{x}^*) \right] \right)$$

- Takes value 0 if \mathbf{x}^* is infeasible and zero otherwise!

Conditional Predictive Distribution II

The factor that guarantees optimality is:

$$p(\mathcal{X}^*|\mathcal{F}) = \prod_{\mathbf{x}^* \in \mathcal{X}^*} \left(\left[\prod_{c=1}^C \Theta(\text{cons}_c(\mathbf{x}^*)) \right] \left[\prod_{\mathbf{x}' \in \mathcal{X}} \Omega(\mathbf{x}', \mathbf{x}^*) \right] \right)$$

- Takes value 0 if \mathbf{x}^* is infeasible and zero otherwise!
- Takes value 0 if \mathbf{x}^* is dominated by feasible \mathbf{x}' and zero otherwise.

Conditional Predictive Distribution II

The factor that guarantees optimality is:

$$p(\mathcal{X}^*|\mathcal{F}) = \prod_{\mathbf{x}^* \in \mathcal{X}^*} \left(\left[\prod_{c=1}^C \Theta(\text{cons}_c(\mathbf{x}^*)) \right] \left[\prod_{\mathbf{x}' \in \mathcal{X}} \Omega(\mathbf{x}', \mathbf{x}^*) \right] \right)$$

- Takes value 0 if \mathbf{x}^* is infeasible and zero otherwise!
- Takes value 0 if \mathbf{x}^* is dominated by feasible \mathbf{x}' and zero otherwise.

The factors are all step functions! The set \mathcal{X} is approximated using the evaluations!

Expectation Propagation

Approximates $p(\mathbf{z}) \propto f_0(\mathbf{z}) \prod_{j=1}^N f_j(\mathbf{z})$ with $q(\mathbf{z}) \propto f_0(\mathbf{z}) \prod_{j=1}^N \tilde{f}_j(\mathbf{z})$

Expectation Propagation

Approximates $p(\mathbf{z}) \propto f_0(\mathbf{z}) \prod_{j=1}^N f_j(\mathbf{z})$ with $q(\mathbf{z}) \propto f_0(\mathbf{z}) \prod_{j=1}^N \tilde{f}_j(\mathbf{z})$

$$p(\mathbf{z}) \propto f_0(\mathbf{z}) f_1(\mathbf{z}) f_2(\mathbf{z}) f_3(\mathbf{z}) \approx q(\mathbf{z}) \propto f_0(\mathbf{z}) \tilde{f}_1(\mathbf{z}) \tilde{f}_2(\mathbf{z}) \tilde{f}_3(\mathbf{z})$$


Expectation Propagation

Approximates $p(\mathbf{z}) \propto f_0(\mathbf{z}) \prod_{j=1}^N f_j(\mathbf{z})$ with $q(\mathbf{z}) \propto f_0(\mathbf{z}) \prod_{j=1}^N \tilde{f}_j(\mathbf{z})$

$$p(\mathbf{z}) \propto f_0(\mathbf{z}) f_1(\mathbf{z}) f_2(\mathbf{z}) f_3(\mathbf{z}) \approx q(\mathbf{z}) \propto f_0(\mathbf{z}) \tilde{f}_1(\mathbf{z}) \tilde{f}_2(\mathbf{z}) \tilde{f}_3(\mathbf{z})$$

The \tilde{f}_j are tuned by minimizing the KL-divergence

$$\text{KL}[\hat{p}_j || q] \quad \text{for } j = 1, \dots, N, \quad \text{where} \quad \begin{aligned} \hat{p}_j(\mathbf{z}) &\propto f_j(\mathbf{z}) \prod_{i \neq j} \tilde{f}_i(\mathbf{z}) \\ q(\mathbf{z}) &\propto \tilde{f}_j(\mathbf{z}) \prod_{i \neq j} \tilde{f}_i(\mathbf{z}) \end{aligned}$$

Expectation Propagation

Approximates $p(\mathbf{z}) \propto f_0(\mathbf{z}) \prod_{j=1}^N f_j(\mathbf{z})$ with $q(\mathbf{z}) \propto f_0(\mathbf{z}) \prod_{j=1}^N \tilde{f}_j(\mathbf{z})$

$$p(\mathbf{z}) \propto f_0(\mathbf{z}) f_1(\mathbf{z}) f_2(\mathbf{z}) f_3(\mathbf{z}) \approx q(\mathbf{z}) \propto f_0(\mathbf{z}) \tilde{f}_1(\mathbf{z}) \tilde{f}_2(\mathbf{z}) \tilde{f}_3(\mathbf{z})$$

The \tilde{f}_j are tuned by minimizing the KL-divergence

$$\text{KL}[\hat{p}_j \| q] \quad \text{for } j = 1, \dots, N, \quad \text{where} \quad \begin{aligned} \hat{p}_j(\mathbf{z}) &\propto f_j(\mathbf{z}) \prod_{i \neq j} \tilde{f}_i(\mathbf{z}) \\ q(\mathbf{z}) &\propto \tilde{f}_j(\mathbf{z}) \prod_{i \neq j} \tilde{f}_i(\mathbf{z}) \end{aligned}$$

The latent variables \mathbf{z} are in our case the objectives and the constraints values at each \mathbf{x}^* and each \mathbf{x}' !

Optimal Ensemble of Decision Trees

- Dataset: German Credit

Optimal Ensemble of Decision Trees

- Dataset: German Credit
 - Number of instances: 1000

Optimal Ensemble of Decision Trees

- Dataset: German Credit
 - Number of instances: 1000
 - Number of features: 20

Optimal Ensemble of Decision Trees

- Dataset: German Credit
 - Number of instances: 1000
 - Number of features: 20
- Ensemble Parameters:

Optimal Ensemble of Decision Trees

- Dataset: German Credit
 - Number of instances: 1000
 - Number of features: 20
- Ensemble Parameters:
 - Ensemble size, random chosen attributes considered at each split, minimum number of samples required to split a node, sub-sampling probability, fraction of labels changed.

Optimal Ensemble of Decision Trees

- Dataset: German Credit
 - Number of instances: 1000
 - Number of features: 20
- Ensemble Parameters:
 - Ensemble size, random chosen attributes considered at each split, minimum number of samples required to split a node, sub-sampling probability, fraction of labels changed.
- Objectives: Ensemble size in log-number of nodes and prediction error (10-fold-cv).

Optimal Ensemble of Decision Trees

- Dataset: German Credit
 - Number of instances: 1000
 - Number of features: 20
- Ensemble Parameters:
 - Ensemble size, random chosen attributes considered at each split, minimum number of samples required to split a node, sub-sampling probability, fraction of labels changed.
- Objectives: Ensemble size in log-number of nodes and prediction error (10-fold-cv).
- Constraints: time for predictions sped-up at least 25% when using a dynamic pruning technique.

Optimal Neural Network on MNIST

- Dataset: MNIST

Optimal Neural Network on MNIST

- Dataset: MNIST
 - Number of instances: 60,000

Optimal Neural Network on MNIST

- Dataset: MNIST
 - Number of instances: 60,000
 - Number of features: $28 \times 28 = 784$

Optimal Neural Network on MNIST

- Dataset: MNIST
 - Number of instances: 60,000
 - Number of features: $28 \times 28 = 784$
- Ensemble Parameters:

Optimal Neural Network on MNIST

- Dataset: MNIST
 - Number of instances: 60,000
 - Number of features: $28 \times 28 = 784$
- Ensemble Parameters:
 - Hidden layers, neurons per layer, learning rate, dropout rate, ℓ_1 penalty, ℓ_2 penalty, memory partition, loop unrolling.

Optimal Neural Network on MNIST

- Dataset: MNIST
 - Number of instances: 60,000
 - Number of features: $28 \times 28 = 784$
- Ensemble Parameters:
 - Hidden layers, neurons per layer, learning rate, dropout rate, ℓ_1 penalty, ℓ_2 penalty, memory partition, loop unrolling.
- Objectives: network error and prediction time (validation set).

Optimal Neural Network on MNIST

- Dataset: MNIST
 - Number of instances: 60,000
 - Number of features: $28 \times 28 = 784$
- Ensemble Parameters:
 - Hidden layers, neurons per layer, learning rate, dropout rate, ℓ_1 penalty, ℓ_2 penalty, memory partition, loop unrolling.
- Objectives: network error and prediction time (validation set).
- Constraints: chip area below threshold.

Potential Extensions / Improvements

- Incorporate information about the Pareto front \mathcal{Y}^* (JES):

Potential Extensions / Improvements

- Incorporate information about the Pareto front \mathcal{Y}^* (JES):
 - Conditioning to \mathcal{Y}^* , which can be done simply by updating each GP.

Potential Extensions / Improvements

- Incorporate information about the Pareto front \mathcal{Y}^* (JES):
 - Conditioning to \mathcal{Y}^* , which can be done simply by updating each GP.
 - Measuring $I(\{\mathcal{X}^*, \mathcal{Y}^*\}; \mathbf{Y})$ is expected to improve results!

Potential Extensions / Improvements

- Incorporate information about the Pareto front \mathcal{Y}^* (JES):
 - Conditioning to \mathcal{Y}^* , which can be done simply by updating each GP.
 - Measuring $I(\{\mathcal{X}^*, \mathcal{Y}^*\}; \mathbf{Y})$ is expected to improve results!
- Use decoupled information for evaluation:

Potential Extensions / Improvements

- Incorporate information about the Pareto front \mathcal{Y}^* (JES):
 - Conditioning to \mathcal{Y}^* , which can be done simply by updating each GP.
 - Measuring $I(\{\mathcal{X}^*, \mathcal{Y}^*\}; \mathbf{Y})$ is expected to improve results!
- Use decoupled information for evaluation:
 - Easily identifies on which black-box to evaluate each batch.

Potential Extensions / Improvements

- Incorporate information about the Pareto front \mathcal{Y}^* (JES):
 - Conditioning to \mathcal{Y}^* , which can be done simply by updating each GP.
 - Measuring $I(\{\mathcal{X}^*, \mathcal{Y}^*\}; \mathbf{Y})$ is expected to improve results!
- Use decoupled information for evaluation:
 - Easily identifies on which black-box to evaluate each batch.
 - Requires optimizing one acquisition per black-box.

Potential Extensions / Improvements

- Incorporate information about the Pareto front \mathcal{Y}^* (JES):
 - Conditioning to \mathcal{Y}^* , which can be done simply by updating each GP.
 - Measuring $I(\{\mathcal{X}^*, \mathcal{Y}^*\}; \mathbf{Y})$ is expected to improve results!
- Use decoupled information for evaluation:
 - Easily identifies on which black-box to evaluate each batch.
 - Requires optimizing one acquisition per black-box.
 - Expected to give better results if more informative black-boxes.