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

Daniel Hernández–Lobato

Computer Science Department Universidad Autónoma de Madrid

http://dhnzl.org, daniel.hernandez@uam.es

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

• Very expensive evaluations.



• Very expensive evaluations.



• The objective is a black-box.





• The objective is a black-box.

• The evaluation can be noisy.



Bayesian optimization methods can be used to solve these problems!















































Bayesian Optimization vs. Uniform Exploration



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

Optimal design of hardware accelerator for neural network predictions.





Optimal design of hardware accelerator for neural network predictions.





Goals:

- Minimize prediction error.
- Minimize prediction time.

Optimal design of hardware accelerator for neural network predictions.





Goals:

Constrained to:

- Minimize **prediction error**. Chip area below a value.
 - Minimize **prediction time**. **Power consumption** below a level.

Optimal design of hardware accelerator for neural network predictions.





Goals:

Constrained to:

- Minimize **prediction error**.
- Chip area below a value.
 - Minimize **prediction time**. **Power consumption** below a level.



Optimal design of hardware accelerator for neural network predictions.





Goals:

Constrained to:

- Minimize **prediction error**.
- Chip area below a value.
- Minimize prediction time. Power consumption below a level.



Challenges:

- Complicated constraints.
- Conflictive objectives.

Constrained Multi-Objective Optimization



Constrained Multi-Objective Optimization

2

-2

-4










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





The Pareto set \mathcal{X}^* in the feasible space is a **random variable**! **Information** is measured by the **entropy** of $p(\mathcal{X}^*|\mathcal{D}_N)$.





Optimized Samples Drawn from the Posterior











The Pareto set \mathcal{X}^* in the feasible space is a **random variable**! Information is measured by the **entropy** of $p(\mathcal{X}^*|\mathcal{D}_N)$.





Optimized Samples Drawn from the Posterior































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



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

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



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

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

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



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

$$\alpha \left(\mathbf{x} \right) = \mathbf{H} \left[\mathcal{X}^{\star} | \mathcal{D}_{t} \right] - \mathbb{E}_{\mathbf{y}} \left[\mathbf{H} \left[\mathcal{X}^{\star} | \mathcal{D}_{t} \cup \{ \mathbf{x}, \mathbf{y} \} \right] \Big| \mathcal{D}_{t}, \mathbf{x} \right] \quad (1)$$
How much we know about \mathcal{X}^{\star} now.

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





The acquisition function is



$$\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{x}, \mathbf{y}\}] \Big| \mathcal{D}_{t}, \mathbf{x} \Big] \equiv \mathsf{MI}(\mathbf{y}, \mathcal{X}^{\star}) \quad (\mathsf{ESMOC})$$

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

$$\begin{aligned} & \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{x},\mathbf{y}\}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathbf{y},\mathcal{X}^{\star}) & (\mathsf{ESMOC}) \\ & \bullet & \bullet \\ & \mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x},\mathcal{X}^{\star}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{y}) & (\mathsf{PESMOC}) \end{aligned}$$

$$\begin{aligned} & \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{x},\mathbf{y}\}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathbf{y},\mathcal{X}^{\star}) & (\mathsf{ESMOC}) \\ & \mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x},\mathcal{X}^{\star}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{y}) & (\mathsf{PESMOC}) \end{aligned}$$

$$\begin{array}{c} \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{y}} \left[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{x},\mathbf{y}\}] \middle| \mathcal{D}_{t},\mathbf{x} \right] \equiv \mathsf{MI}(\mathbf{y},\mathcal{X}^{\star}) \quad (\mathsf{ESMOC}) \\ \\ \mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x}] - \mathbb{E}_{\mathcal{X}^{\star}} \left[\mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x},\mathcal{X}^{\star}] \middle| \mathcal{D}_{t},\mathbf{x} \right] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{y}) \quad (\mathsf{PESMOC}) \\ \\ \end{array}$$

$$\begin{array}{c} \mathsf{Gaussian} \\ \mathsf{distribution} \end{array}$$

$$\begin{array}{c} \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{y}} \left[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{x},\mathbf{y}\}] \middle| \mathcal{D}_{t},\mathbf{x} \right] \equiv \mathsf{MI}(\mathbf{y},\mathcal{X}^{\star}) \quad (\mathsf{ESMOC}) \\ \\ \mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x}] - \mathbb{E}_{\mathcal{X}^{\star}} \left[\mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x},\mathcal{X}^{\star}] \middle| \mathcal{D}_{t},\mathbf{x} \right] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{y}) \quad (\mathsf{PESMOC}) \\ \\ \\ \\ \mathbf{G}_{aussian} \\ \\ \mathbf{distribution} \end{array}$$

$$\begin{array}{c} \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{x},\mathbf{y}\}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathbf{y},\mathcal{X}^{\star}) \quad (\mathsf{ESMOC}) \\ \\ \mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x},\mathcal{X}^{\star}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{y}) \quad (\mathsf{PESMOC}) \\ \\ \\ \hline \\ \mathbf{G}_{aussian} \\ \text{distribution} \\ \end{array}$$

$$\begin{array}{c} \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{x},\mathbf{y}\}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathbf{y},\mathcal{X}^{\star}) \quad (\mathsf{ESMOC}) \\ \\ \mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x},\mathcal{X}^{\star}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{y}) \quad (\mathsf{PESMOC}) \\ \\ \\ \hline \\ \mathbf{G}_{aussian} \\ \text{distribution} \\ \end{array}$$

$$\begin{array}{c} \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{x},\mathbf{y}\}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathbf{y},\mathcal{X}^{\star}) \quad (\mathsf{ESMOC}) \\ \\ \mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x},\mathcal{X}^{\star}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{y}) \quad (\mathsf{PESMOC}) \\ \\ \end{array}$$

$$\begin{array}{c} \mathsf{Gaussian} \\ \mathsf{distribution} \end{array} \quad \mathsf{Approximated by} \\ \mathsf{sampling from } p(\mathcal{X}^{\star}|\mathcal{D}_{t}) \end{array} \quad \mathsf{Factorized Gaussian approximation} \\ \mathcal{X}^{\star} \operatorname{dominates any other point in } \mathcal{X} \\ \end{array} \right. .$$

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

(Minka, 2001)

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

$$\begin{aligned} & \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{x},\mathbf{y}\}] \, \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathbf{y},\mathcal{X}^{\star}) \quad (\mathsf{ESMOC}) \\ & \mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x},\mathcal{X}^{\star}] \, \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{y}) \quad (\mathsf{PESMOC}) \\ & \mathsf{Gaussian} \\ & \mathsf{Approximated by} \\ & \mathsf{sampling from } p(\mathcal{X}^{\star}|\mathcal{D}_{t}) \Big] \end{aligned}$$

$$\begin{aligned} & \mathsf{Factorized Gaussian approximation} \\ & \mathsf{with expectation propagation.} \\ & \mathcal{X}^{\star} \operatorname{dominates any other point in } \mathcal{X} \Big] . \end{aligned}$$

$$\begin{aligned} & \mathsf{A}(\mathbf{x}) \approx \sum_{c=1}^{C} \log v_{c}^{PD}(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^{M} \left(\sum_{k=1}^{C} \log v_{k}^{CPD}(\mathbf{x}|\mathcal{X}_{(m)}^{\star}) \right) + \\ & \sum_{k=1}^{K} \log v_{k}^{PD}(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^{M} \left(\sum_{k=1}^{K} \log v_{k}^{CPD}(\mathbf{x}|\mathcal{X}_{(m)}^{\star}) \right) = \sum_{i=1}^{C+K} \alpha_{i}(\mathbf{x}) \end{aligned}$$

(Minka, 2001)

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

$$\begin{aligned} & \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{x},\mathbf{y}\}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathbf{y},\mathcal{X}^{\star}) \quad (\mathsf{ESMOC}) \\ & \mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{y}|\mathcal{D}_{t},\mathbf{x},\mathcal{X}^{\star}] \Big| \mathcal{D}_{t},\mathbf{x} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{y}) \quad (\mathsf{PESMOC}) \\ & \mathsf{Gaussian} \\ & \mathsf{distribution} \\ & \mathsf{Approximated by} \\ & \mathsf{sampling from } p(\mathcal{X}^{\star}|\mathcal{D}_{t}) \\ & \mathsf{Factorized Gaussian approximation} \\ & \mathsf{with expectation} \\ & \mathsf{One acquisition} \\ & \mathsf{Probal} \\ & \mathsf{N}^{\star} \operatorname{dominates any} \\ & \mathsf{One acquisition} \\ & \mathsf{Probal} \\ & \mathsf{C}_{c=1}^{C} \log v_{c}^{CPD}(\mathbf{x}|\mathcal{X}_{(m)}) \Big) + \\ & \mathsf{C}_{k=1}^{K} \log v_{k}^{PD}(\mathbf{x}) - \frac{1}{M} \sum_{m=1}^{M} \left(\sum_{k=1}^{K} \log v_{k}^{CPD}(\mathbf{x}|\mathcal{X}_{(m)}^{\star}) \right) = \sum_{i=1}^{C+K} \alpha_{i}(\mathbf{x}) \end{aligned}$$

(Minka, 2001)
The predictions must be compatible with \mathcal{X}^{\star} !

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

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

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

Noiseless conditional predictive distribution:

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

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

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

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

Noiseless conditional predictive distribution:

$$p(\mathbf{f}|\mathcal{D}, \mathbf{x}, \mathcal{X}^{\star}) \propto \int \underbrace{p(\mathbf{f}|\mathbf{x}, \mathcal{F})}_{\text{Black-box values at } \mathbf{x}} \times \underbrace{p(\mathcal{F}|\mathcal{D})}_{\text{Post. dist.}} \times \underbrace{p(\mathcal{X}^{\star}|\mathcal{F})}_{\text{Guarantees } \mathcal{X}^{\star} \text{ optimal.}} d\mathcal{F}$$
where \mathcal{F} informally representes all potential black-box functions.
• Unconditional posterior.

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

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

Noiseless conditional predictive distribution:

$$p(\mathbf{f}|\mathcal{D}, \mathbf{x}, \mathcal{X}^{\star}) \propto \int \underbrace{p(\mathbf{f}|\mathbf{x}, \mathcal{F})}_{\text{Black-box values at } \mathbf{x}} \times \underbrace{p(\mathcal{F}|\mathcal{D})}_{\text{Post. dist.}} \times \underbrace{p(\mathcal{X}^{\star}|\mathcal{F})}_{\text{Guarantees } \mathcal{X}^{\star} \text{ optimal.}} d\mathcal{F}$$
where \mathcal{F} informally representes all potential black-box functions.
• Unconditional posterior.

• Takes value 1 if \mathcal{X}^{\star} is optimal given \mathcal{F} and zero otherwise. /

The factor that guarantees optimality is:

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

The factor that guarantees optimality is:

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

Takes value 0 if x^{*} is infeasible and zero othersie!

The factor that guarantees optimality is:

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

Takes value 0 if \mathbf{x}^{\star} is infeasible and zero othersie!

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

The factor that guarantees optimality is:

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

Takes value 0 if \mathbf{x}^{\star} is infeasible and zero othersie!

• 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!

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

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

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

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

$$\mathsf{KL}[\hat{p}_j||q] \quad ext{for}\, j=1,\ldots, \mathsf{N}\,, \quad ext{where} \quad egin{array}{cc} \hat{p}_j(m{z}) & \propto & f_j(m{z}) \prod_{i
eq j} ilde{f}_i(m{z}) \ q(m{z}) & \propto & ilde{f}_j(m{z}) \prod_{i
eq j} ilde{f}_i(m{z}) \,. \end{cases}$$

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

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

$$\mathsf{KL}[\hat{p}_j||q] \quad ext{for}\, j=1,\ldots, \mathsf{N}\,, \quad ext{where} \quad egin{array}{cc} \hat{p}_j(m{z}) & \propto & f_j(m{z}) \prod_{i
eq j} ilde{f}_i(m{z}) \ q(m{z}) & \propto & ilde{f}_j(m{z}) \prod_{i
eq j} ilde{f}_i(m{z}) \,. \end{cases}$$

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













Traditional Bayesian optimization is sequential!

Traditional Bayesian optimization is sequential!



Traditional Bayesian optimization is sequential!



Computing clusters let us do many things at once!

Traditional Bayesian optimization is sequential!



Computing clusters let us do many things at once!



Traditional Bayesian optimization is sequential!



Computing clusters let us do many things at once!



Traditional Bayesian optimization is sequential!



Computing clusters let us do many things at once!



Parallel experiments should be highly informative but different!

$$\begin{split} \mathsf{H}\big[\mathcal{X}^{\star}\big|\mathcal{D}_{t}\big] - \mathbb{E}_{\mathbf{Y}}\!\Big[\mathsf{H}\big[\mathcal{X}^{\star}\big|\mathcal{D}_{t}\cup\{\mathbf{Y},\mathbf{X}\}\big]\Big|\mathcal{D}_{t},\mathbf{X}\Big] \equiv \mathsf{MI}(\mathbf{Y},\mathcal{X}^{\star}) & \text{(Parallel ESMOC)} \end{split}$$

$$\begin{split} \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] &- \mathbb{E}_{\mathbf{Y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{Y},\mathbf{X}\}] \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathbf{Y},\mathcal{X}^{\star}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ \mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X}] &- \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X},\mathcal{X}^{\star}] \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{Y}) & (\mathsf{Parallel} \\ \mathsf{PESMOC}) \\ \end{split}$$

$$\begin{split} &\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{Y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{Y},\mathbf{X}\}] \, \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathbf{Y},\mathcal{X}^{\star}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ &\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X},\mathcal{X}^{\star}] \, \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{Y}) & (\mathsf{Parallel} \\ \mathsf{PESMOC}) \\ \end{split}$$

$$\begin{split} & \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{Y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{Y},\mathbf{X}\}] \, \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathbf{Y},\mathcal{X}^{\star}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ & \mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X},\mathcal{X}^{\star}] \, \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{Y}) & (\mathsf{Parallel} \\ \mathsf{PESMOC}) \\ & \mathsf{PESMOC}) \end{split}$$

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

$$\begin{array}{c} \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{Y}} \left[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{Y},\mathbf{X}\}] \middle| \mathcal{D}_{t},\mathbf{X} \right] \equiv \mathsf{MI}(\mathbf{Y},\mathcal{X}^{\star}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ \mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X}] - \mathbb{E}_{\mathcal{X}^{\star}} \left[\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X},\mathcal{X}^{\star}] \middle| \mathcal{D}_{t},\mathbf{X} \right] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{Y}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ \mathbf{Gaussian} \\ \mathsf{distribution} \end{array}$$

$$\begin{array}{c} \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{Y}} \left[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{Y},\mathbf{X}\}] \middle| \mathcal{D}_{t},\mathbf{X} \right] \equiv \mathsf{MI}(\mathbf{Y},\mathcal{X}^{\star}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ \mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X}] - \mathbb{E}_{\mathcal{X}^{\star}} \left[\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X},\mathcal{X}^{\star}] \middle| \mathcal{D}_{t},\mathbf{X} \right] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{Y}) & (\mathsf{Parallel} \\ \mathsf{PESMOC}) \\ \end{array}$$

$$\begin{array}{c} \mathsf{Gaussian} \\ \mathsf{distribution} \end{array}$$

$$\begin{array}{c} \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{Y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{Y},\mathbf{X}\}] \middle| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathbf{Y},\mathcal{X}^{\star}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ \mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X},\mathcal{X}^{\star}] \middle| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{Y}) & (\mathsf{Parallel} \\ \mathsf{PESMOC}) \\ \end{array}$$

$$\begin{array}{c} \mathsf{Gaussian} \\ \mathsf{distribution} \end{array} \quad \begin{array}{c} \mathsf{Approximated} \ \mathsf{by} \\ \mathsf{sampling} \ \mathsf{from} \ \rho(\mathcal{X}^{\star}|\mathcal{D}_{t}) \end{array}$$

$$\begin{array}{c} \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{Y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{Y},\mathbf{X}\}] \, \big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathbf{Y},\mathcal{X}^{\star}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ \mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X},\mathcal{X}^{\star}] \, \big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{Y}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ & \mathsf{Gaussian} \\ \mathsf{distribution} & \mathsf{Approximated} \ \mathsf{by} \\ \mathsf{sampling} \ \mathsf{from} \ p(\mathcal{X}^{\star}|\mathcal{D}_{t}) \end{array}$$

$$\begin{array}{c} \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{Y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{Y},\mathbf{X}\}] \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathbf{Y},\mathcal{X}^{\star}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ \mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X},\mathcal{X}^{\star}] \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{Y}) & (\mathsf{Parallel} \\ \mathsf{PESMOC}) \\ \end{array}$$

$$\begin{array}{c} \mathsf{Gaussian} \\ \mathsf{distribution} \end{array} \quad \begin{array}{c} \mathsf{Approximated} \ \mathsf{by} \\ \mathsf{sampling} \ \mathsf{from} \ p(\mathcal{X}^{\star}|\mathcal{D}_{t}) \end{array} \right] \\ \begin{array}{c} \mathsf{Factorized} \ \mathsf{Gaussian} \ \mathsf{approximation} \\ \mathsf{with} \ \mathsf{expectation} \ \mathsf{propagation}. \\ \mathcal{X}^{\star} \ \mathsf{dominates} \ \mathsf{any} \ \mathsf{other point in} \ \mathcal{X} \end{array} \right].$$

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

$$\begin{split} & \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{Y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{Y},\mathbf{X}\}] \, \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathbf{Y},\mathcal{X}^{\star}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ & \mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X}] - \mathbb{E}_{\mathcal{X}^{\star}}^{\mathsf{L}} \Big[\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X},\mathcal{X}^{\star}] \, \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star},\mathbf{Y}) & (\mathsf{Parallel} \\ \mathsf{PESMOC}) \\ & \mathsf{Gaussian} \\ \mathsf{distribution} & \mathsf{Approximated} \ \mathsf{by} \\ \mathsf{sampling} \ \mathsf{from} \ p(\mathcal{X}^{\star}|\mathcal{D}_{t}) \\ & \mathsf{Factorized} \ \mathsf{Gaussian} \ \mathsf{approximation} \\ \mathsf{with} \ \mathsf{expectation} \ \mathsf{propagation}. \\ & \mathcal{X}^{\star} \ \mathsf{dominates} \ \mathsf{any} \ \mathsf{other} \ \mathsf{point} \ \mathsf{i} \ \mathcal{X} \\ & \mathsf{A}(\mathbf{x}) \approx \sum_{c=1}^{C} \log |\mathbf{V}_{c}^{\mathrm{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^{M} \left(\sum_{c=1}^{C} \log |\mathbf{V}_{c}^{\mathrm{CPD}}(\mathbf{X}|\mathcal{X}_{m}^{\star})| \right) + \\ & \sum_{k=1}^{K} \log |\mathbf{V}_{k}^{\mathrm{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^{M} \left(\sum_{k=1}^{K} \log |\mathbf{V}_{k}^{\mathrm{CPD}}(\mathbf{X}|\mathcal{X}_{m}^{\star})| \right) \\ \end{split}$$

(Minka, 2001)
Choose a set of Q points $\mathbf{X} = {\mathbf{x}_q}_{q=1}^Q$ to minimize the entropy of \mathcal{X}^{\star} .

$$\begin{split} & \mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{Y}} \Big[\mathsf{H}[\mathcal{X}^{\star}|\mathcal{D}_{t} \cup \{\mathbf{Y}, \mathbf{X}\}] \Big| \mathcal{D}_{t}, \mathbf{X} \Big] \equiv \mathsf{MI}(\mathbf{Y}, \mathcal{X}^{\star}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ & \mathsf{H}[\mathbf{Y}|\mathcal{D}_{t}, \mathbf{X}] - \mathbb{E}_{\mathcal{X}^{\star}} \Big[\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t}, \mathbf{X}, \mathcal{X}^{\star}] \Big| \mathcal{D}_{t}, \mathbf{X} \Big] \equiv \mathsf{MI}(\mathcal{X}^{\star}, \mathbf{Y}) & (\mathsf{Parallel} \\ \mathsf{PESMOC}) \\ & \mathsf{Gaussian} \\ \mathsf{distribution} & \mathsf{Approximated by} \\ \mathsf{sampling from } p(\mathcal{X}^{\star}|\mathcal{D}_{t}) \\ & \mathsf{Min}(\mathsf{expectation propagation}, \mathcal{X}^{\star}) \\ & \mathcal{X}^{\star} \mathsf{dominates any other point in } \mathcal{X} \\ & \mathsf{A}(\mathbf{x}) \approx \sum_{c=1}^{C} \log |\mathbf{V}_{c}^{\mathsf{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^{M} \left(\sum_{k=1}^{C} \log |\mathbf{V}_{k}^{\mathsf{CPD}}(\mathbf{X}|\mathcal{X}_{m}^{\star})| \right) + \\ & \sum_{k=1}^{K} \log |\mathbf{V}_{k}^{\mathsf{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^{M} \left(\sum_{k=1}^{K} \log |\mathbf{V}_{k}^{\mathsf{CPD}}(\mathbf{X}|\mathcal{X}_{m}^{\star})| \right) = \sum_{i=1}^{K+C} \alpha_{k}(\mathbf{X}) \\ & \mathsf{A}(\mathbf{X}) \in \mathsf{A}(\mathbf{X}) = \frac{1}{M} \sum_{m=1}^{M} \left(\sum_{k=1}^{K} \log |\mathbf{V}_{k}^{\mathsf{CPD}}(\mathbf{X}|\mathcal{X}_{m}^{\star})| \right) = \mathbb{E}_{i=1}^{K+C} \alpha_{k}(\mathbf{X}) \\ & \mathsf{A}(\mathbf{X}) = \mathbb{E}_{i=1}^{K+C} \alpha_{k}(\mathbf{$$

(Minka, 2001)

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

$$\begin{aligned} & \mathsf{H}[\mathcal{X}^{*}|\mathcal{D}_{t}] - \mathbb{E}_{\mathbf{Y}} \Big[\mathsf{H}[\mathcal{X}^{*}|\mathcal{D}_{t} \cup \{\mathbf{Y},\mathbf{X}\}] \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathbf{Y},\mathcal{X}^{*}) & (\mathsf{Parallel} \\ \mathsf{ESMOC}) \\ & \mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X}] - \mathbb{E}_{\mathcal{X}^{*}} \Big[\mathsf{H}[\mathbf{Y}|\mathcal{D}_{t},\mathbf{X},\mathcal{X}^{*}] \Big| \mathcal{D}_{t},\mathbf{X} \Big] \equiv \mathsf{MI}(\mathcal{X}^{*},\mathbf{Y}) & (\mathsf{Parallel} \\ \mathsf{PESMOC}) \\ & \mathsf{Gaussian} \\ \mathsf{distribution} & \mathsf{Approximated by} \\ \mathsf{sampling from } p(\mathcal{X}^{*}|\mathcal{D}_{t}) \\ & \mathsf{Min} \text{ expectation production} \\ & \mathsf{A}(\mathbf{x}) \approx \sum_{c=1}^{C} \log |\mathbf{V}_{c}^{\mathsf{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^{M} \left(\sum_{c=1}^{C} \log |\mathbf{V}_{c}^{\mathsf{CPD}}(\mathbf{X}|\mathcal{X}_{m}^{*})| \right) + \\ & \sum_{k=1}^{K} \log |\mathbf{V}_{k}^{\mathsf{PD}}(\mathbf{X})| - \frac{1}{M} \sum_{m=1}^{M} \left(\sum_{k=1}^{K} \log |\mathbf{V}_{k}^{\mathsf{CPD}}(\mathbf{X}|\mathcal{X}_{m}^{*})| \right) = \sum_{i=1}^{K+C} \alpha_{k}(\mathbf{X}) \\ & \mathsf{Min} (\mathbf{X}) = \sum_{i=1}^{K+C} \alpha_{k}(\mathbf{X})$$

(Minka, 2001)

Consdierations:

• The cost is linear in the number of objectives K and constraints C.

- The cost is linear in the number of objectives K and constraints C.
- The cost is cubic in the batch size due to the determinants.

- The cost is linear in the number of objectives K and constraints C.
- The cost is cubic in the batch size due to the determinants.
- It easily allows for batch decoupled evaluations.

- The cost is linear in the number of objectives K and constraints C.
- The cost is cubic in the batch size due to the determinants.
- It easily allows for batch decoupled evaluations.
- Optimizing the acquisition requires gradient computations.

Consdierations:

- The cost is linear in the number of objectives K and constraints C.
- The cost is cubic in the batch size due to the determinants.
- It easily allows for batch decoupled evaluations.
- Optimizing the acquisition requires gradient computations.

Automatic gradient computation by keeping fixed the approximate factors and using automatic differentiation tools (Autograd)!

Comparison of Exact and Approximate Acquisition

1 dimensional problem, Batch size = 2:



Comparison of Exact and Approximate Acquisition

1 dimensional problem, Batch size = 2:

Approximate acquisition function



Comparison of Exact and Approximate Acquisition

1 dimensional problem, Batch size = 2:

Approximate acquisition function



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

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

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

- 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.

- 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.

- 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:

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

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

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

- 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 (niosy 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 Experimnets



Synthetic Experimnets



PPESMOC 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_PESMOC	191.5 ± 7.0	347.2 ± 6.0	$\textbf{405.49} \pm \textbf{5.8}$	801.05 ± 27.8	1957.72 34.1
PS_BMOO	379.4 ± 13.1	551.1 ± 21.7	593.86 ± 18.0	897.4 ± 29.6	1870.42 42.77
qEHVI	65.2 ± 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 ± 26.9	912.74 ± 26.3	957.3 ± 25.7	1045.7 ± 30.53	1269.35 26.62
PS_PESMOC	191.5 ± 7.0	347.2 ± 6.0	$\textbf{405.49} \pm \textbf{5.8}$	801.05 ± 27.8	1957.72 34.1
PS_BMOO	379.4 ± 13.1	551.1 ± 21.7	593.86 ± 18.0	897.4 ± 29.6	1870.42 42.77
qEHVI	65.2 ± 1.8	417.9 ± 21.9	1174.9 ± 54.3		
qNEHVI	89.5 ± 2.3	401.4 ± 23.9	1169.4 ± 56.1		

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

Benchmark Experimnets



Benchmark Experimnets



PPESMOC performs better or similar to the other strategies!

• Dataset: German Credit

- Dataset: German Credit
 - Number of instances: 1000

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

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

- 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.

- 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).

- 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 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
100 200	0.325 ± 0.007 0.334 ± 0.005	$\begin{array}{c} 0.327 \pm 0.007 \\ 0.335 \pm 0.006 \end{array}$	$\frac{0.295}{0.313} \pm \frac{0.014}{0.010}$	$\frac{0.298}{0.310} \pm \frac{0.009}{0.007}$		$\frac{\underline{0.294}}{\underline{0.309}} \pm \frac{\underline{0.013}}{\underline{0.010}}$

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 a = 0.05.

# Eval.	PPESMOC	PS_PESMOC	PS_BMOO	P_RANDOM	qNEHVI	qEHVI
100 200	0.325 ± 0.007 0.334 ± 0.005	$\begin{array}{c} 0.327 \pm 0.007 \\ 0.335 \pm 0.006 \end{array}$	$\frac{0.295}{0.313} \pm \frac{0.014}{0.010}$	$\frac{0.298}{0.310} \pm \frac{0.009}{0.007}$	$\frac{0.299}{0.3154} \pm \frac{0.011}{0.008}$	$\frac{0.294}{0.309} \pm \frac{0.013}{0.010}$

PPESMOC performs better than or similar to the other strategies!

Optimal Neural Network

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

- Dataset: MNIST
 - Number of instances: 60,000
 - Number of features: 28x28=784

- Dataset: MNIST
 - Number of instances: 60,000
 - Number of features: 28x28=784
- Ensemble Parameters:

- Dataset: MNIST
 - Number of instances: 60,000
 - Number of features: 28x28=784
- Ensemble Parameters:
 - Hidden layers, neurons per layer, learning rate, dropout rate, ℓ_1 penalty, ℓ_2 penalty, memory partition, loop unrolling.

- Dataset: MNIST
 - Number of instances: 60,000
 - Number of features: 28x28=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).

- Dataset: MNIST
 - Number of instances: 60,000
 - Number of features: 28x28=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.



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	$1.020\ \pm\ 0.014$	1.014 ± 0.029	0.982 ± 0.095	0.993 ± 0.035	0.999 ± 0.050	$\underline{0.996} \pm \underline{0.041}$



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	$1.020\ \pm\ 0.014$	1.014 ± 0.029	$\underline{0.982} \pm \underline{0.095}$	$\underline{0.993} \pm \underline{0.035}$	$\underline{0.999} \pm \underline{0.050}$	$\underline{0.996} \pm \underline{0.041}$

PPESMOC performs slithgly better than the other strategies!

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

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

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

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

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

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

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

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

- 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.

- 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.

- 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!



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.