

Julian Rodemann<sup>1</sup>, Thomas Augustin<sup>1</sup>

# Towards Prior-Mean Robust Bayesian Optimization

Young Statisticians Session (YSS)

DAGStat 2022

March 30, 2022, Hamburg



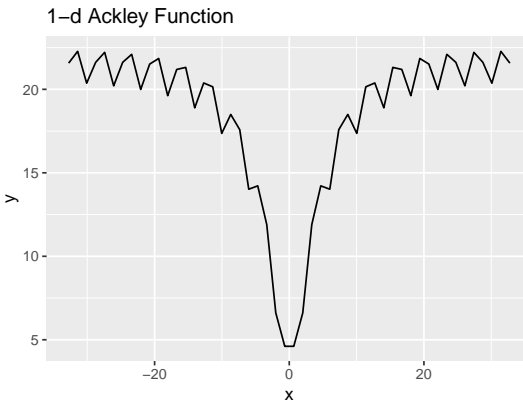
---

<sup>1</sup>Department of Statistics, LMU Munich

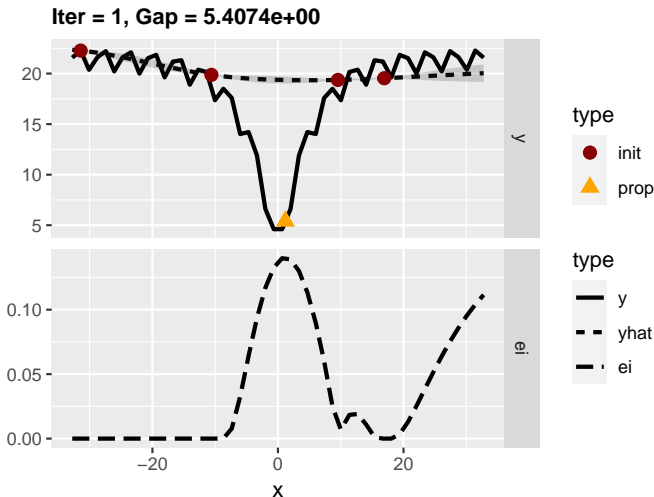


- ① Bayesian Optimization
- ② Gaussian Processes
- ③ Sensitivity Analysis
  - Setup
  - Results
- ④ Prior-Mean-Robust BO (PROBO)
  - Prior near-ignorance models
  - GLCB
- ⑤ Application in Material Science
- ⑥ Discussion
- ⑦ Literature
- ⑧ Appendix

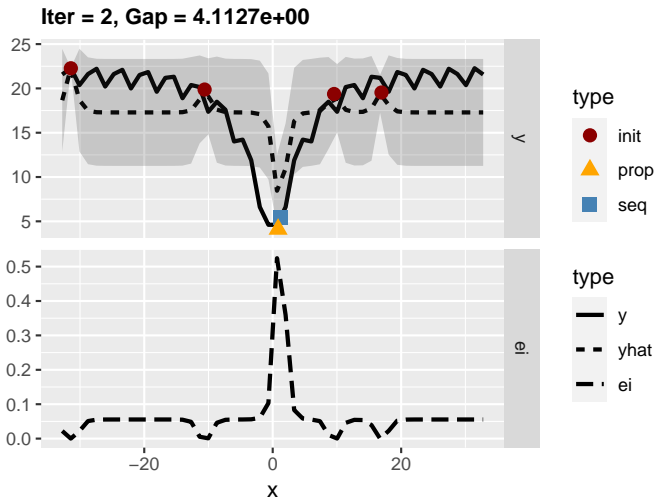
- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Sensitivity Analysis
- 4 Prior-Mean-Robust BO (PROBO)
- 5 Application in Material Science
- 6 Discussion
- 7 Literature
- 8 Appendix



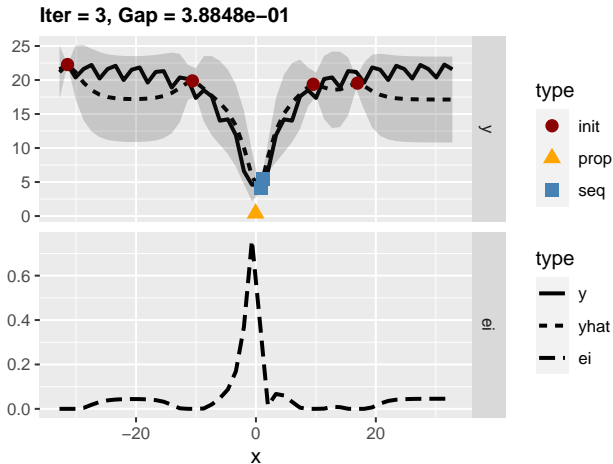
Note: If not otherwise stated, all figures are based on own computations using `ggplot2` [Wickham, 2016], `smoof` [Bossek, 2017] and `mlr(3)MB0` [Bischl et al., 2017]



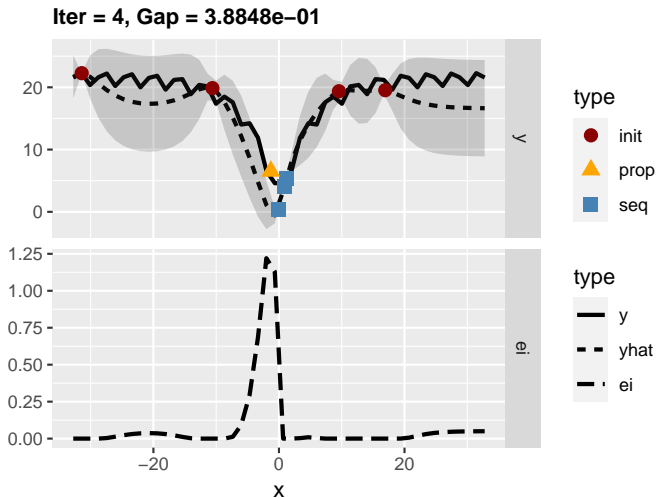
Iteration 1: Surrogate Model (top) and Acquisition Function (bottom)



Iteration 2

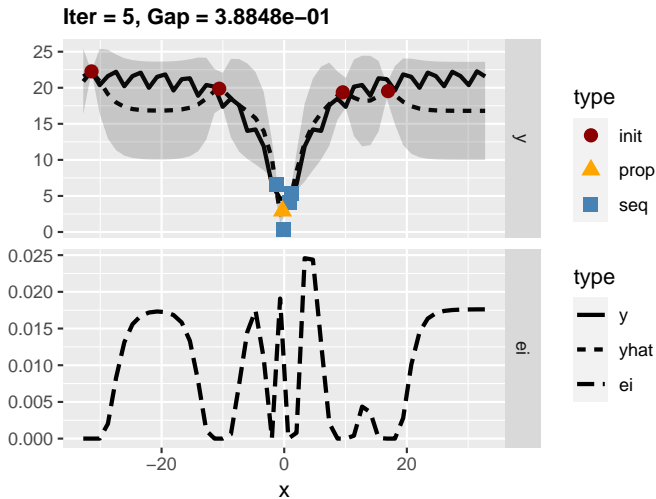


Iteration 3

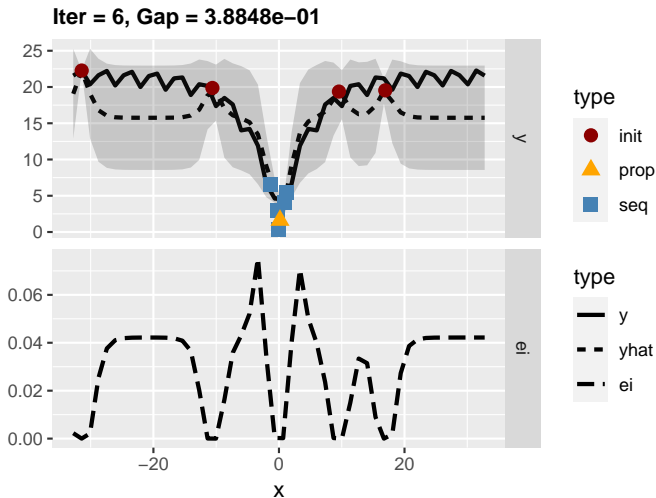


Iteration 4

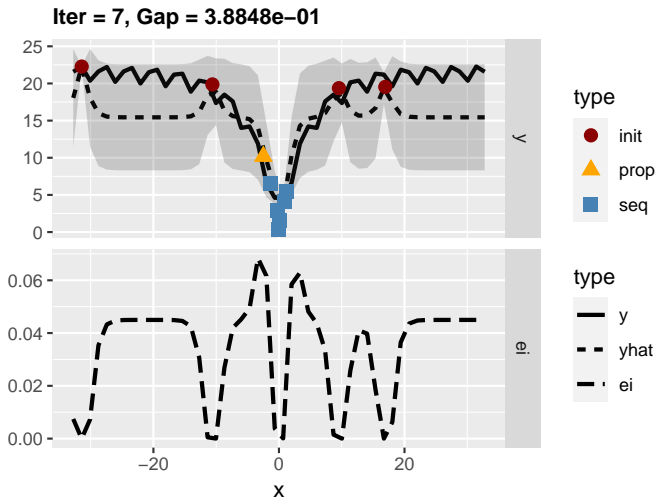




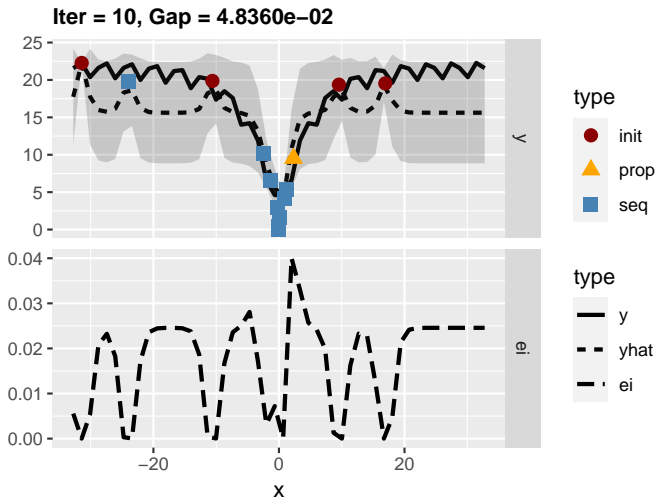
Iteration 5



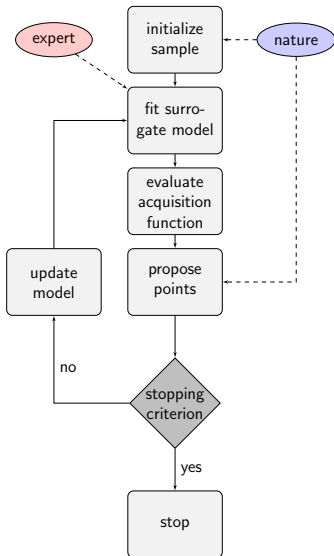
Iteration 6



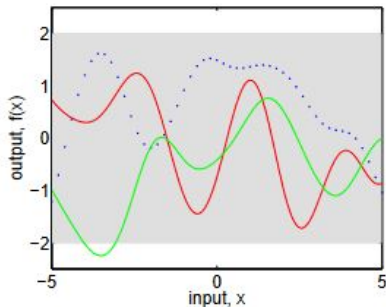
Iteration 7



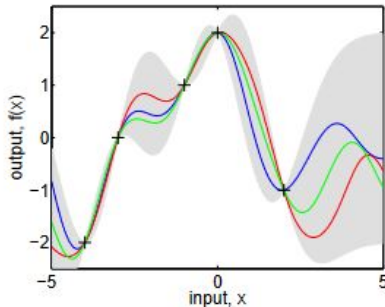
Iteration 10



- 1 Bayesian Optimization
- 2 Gaussian Processes**
- 3 Sensitivity Analysis
- 4 Prior-Mean-Robust BO (PROBO)
- 5 Application in Material Science
- 6 Discussion
- 7 Literature
- 8 Appendix

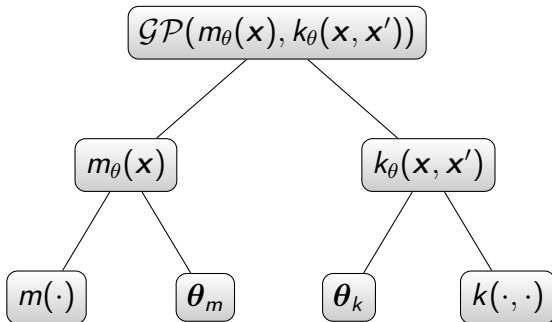


(a), prior



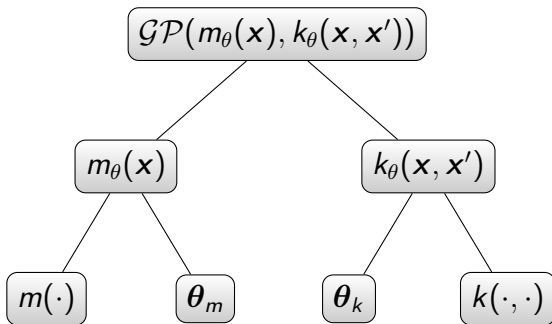
(b), posterior

Functional GP regression: Three functions drawn from prior (a) and posterior (b) GP. Image credits: [Rasmussen, 2003].



How to specify  $m(\cdot)$ ,  $\theta_m$ ,  $\theta_k$  and  $k(\cdot, \cdot)$   
in absence of prior knowledge?





**And:** Do they even matter?

- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Sensitivity Analysis**
  - Setup
  - Results
- 4 Prior-Mean-Robust BO (PROBO)
- 5 Application in Material Science
- 6 Discussion
- 7 Literature

- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Sensitivity Analysis**
  - Setup
  - Results
- 4 Prior-Mean-Robust BO (PROBO)
- 5 Application in Material Science
- 6 Discussion
- 7 Literature



- We randomly select 50 synthetic test functions from the R package `smoof` [Bossek, 2017], stratified across the covariate space dimensions 1, 2, 3, 4 and 7.
- For each of them, a sensitivity analysis is conducted with regard to each of the four prior components.
  - 5 functional forms
  - 5 mean and kernel parameter specifications (relative deviation from global mean)
  - we control for interaction effects
- The initial design of size  $n_{init} = 10$  is randomly sampled anew for each of the  $R = 40$  BO repetitions with  $T = 20$  iterations each.

- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Sensitivity Analysis**
  - Setup
  - Results
- 4 Prior-Mean-Robust BO (PROBO)
- 5 Application in Material Science
- 6 Discussion
- 7 Literature



- **Mean parameters** influence convergence the most, followed by the **kernel's functional form**.
- **Mean functional form** and **Kernel parameters** play a (relatively) negligible role.

- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Sensitivity Analysis
- 4 Prior-Mean-Robust BO (PROBO)**
  - Prior near-ignorance models
  - GLCB
- 5 Application in Material Science
- 6 Discussion
- 7 Literature

- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Sensitivity Analysis
- 4 **Prior-Mean-Robust BO (PROBO)**  
Prior near-ignorance models  
GLCB
- 5 Application in Material Science
- 6 Discussion
- 7 Literature





- Idea: Use set of  $\theta_m$  instead of precise  $\theta_m$ . Fully specify the other components.
- [Mangili, 2015] proposes imprecise Gaussian processes

$$\left\{ \mathcal{GP} \left( Mh, k_{\theta}(x, x') + \frac{1+M}{c} \right) : h = \pm 1, M \geq 0 \right\},$$

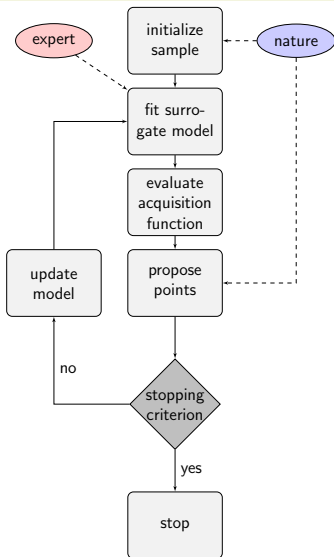
given a base kernel  $k_{\theta}(x, x')$  and a degree of imprecision  $c > 0$ .

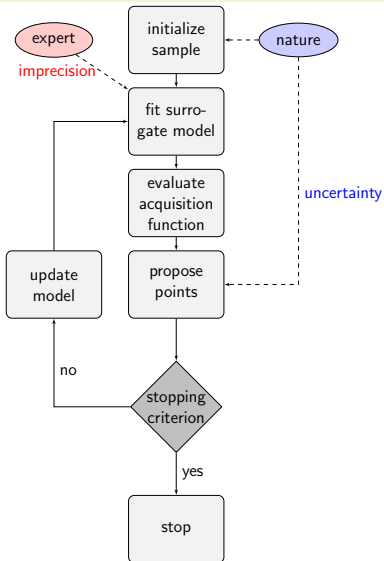
→ results in a set of posteriors whose upper and lower mean estimates  $\underline{\hat{\mu}}(x)_c, \overline{\hat{\mu}}(x)_c$  can be derived

- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Sensitivity Analysis
- 4 Prior-Mean-Robust BO (PROBO)**
  - Prior near-ignorance models
  - GLCB
- 5 Application in Material Science
- 6 Discussion
- 7 Literature



- $LCB(x) = -\hat{\mu}(x) + \tau \cdot \underbrace{\sqrt{\widehat{\text{Var}}(\mu(x))}}_{\text{"classical" uncertainty}}$
- $GLCB(x) = -\hat{\mu}(x) + \tau \cdot \underbrace{\sqrt{\widehat{\text{Var}}(\mu(x))}}_{\text{"classical" uncertainty}} + \rho \cdot \underbrace{(\bar{\mu}(x)_c - \underline{\mu}(x)_c)}_{\text{prior-induced imprecision}}$ 
  - $\tau$  is the degree of **risk**-aversion
  - $\rho$  is the degree of **ambiguity**-aversion







Notably,  $\bar{\hat{\mu}}(\mathbf{x}) - \underline{\hat{\mu}}(\mathbf{x})$  simplifies to an expression only dependent on predictive kernels  $\mathbf{k}_x = [k_\theta(x, x_1), \dots, k_\theta(x, x_n)]^T$ , the base kernel matrix  $\mathbf{K}_n$  (from training) and the degree of imprecision  $c$ . For some<sup>1</sup> values of  $c$  (depending on observations):

$$\bar{\hat{\mu}}(\mathbf{x}) - \underline{\hat{\mu}}(\mathbf{x}) = (1 - \mathbf{k}_x^T \mathbf{s}_k) \left( \frac{\mathbf{s}_k^T \mathbf{y}}{\mathbf{s}_k} + \frac{c}{\mathbf{s}_k} - \frac{\mathbf{s}_k^T \mathbf{y}}{c + \mathbf{s}_k} \right) \quad (1)$$

---

<sup>1</sup>For a thorough case distinction, please refer to the Appendix.



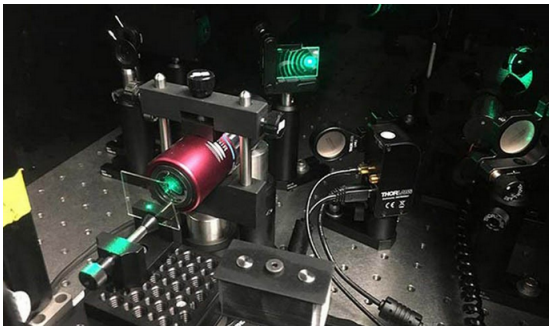
For sufficiently high  $c$ , the model imprecision  $\bar{\hat{\mu}}(\mathbf{x}) - \underline{\hat{\mu}}(\mathbf{x})$  even simplifies further:

$$\bar{\hat{\mu}}(\mathbf{x}) - \underline{\hat{\mu}}(\mathbf{x}) = 2c \frac{|1 - \mathbf{k}_x^T \mathbf{s}_k|}{\mathbf{S}_k} \quad (2)$$

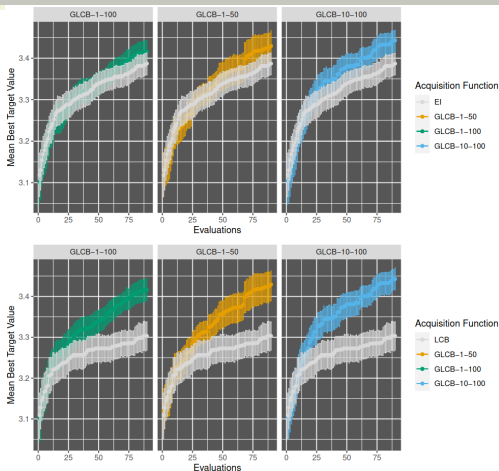
In this case, GLCB's hyperparameters  $\rho$  and  $c$  collapse to one.

- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Sensitivity Analysis
- 4 Prior-Mean-Robust BO (PROBO)
- 5 Application in Material Science**
- 6 Discussion
- 7 Literature
- 8 Appendix





Experimental set-up of graphene production: "The preparation of a sample to be irradiated requires about **one week**." [Kotthoff, 2019]



BO with GLCB on Graphene function. GLCB-1-50 means GLCB with  $\rho = 1$ ,  $c = 50$ . Data source: [Wahab et al., 2020].

- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Sensitivity Analysis
- 4 Prior-Mean-Robust BO (PROBO)
- 5 Application in Material Science
- 6 Discussion**
- 7 Literature
- 8 Appendix



- Limitations
  - robust only with regard to possible misspecification of the mean function parameter given a constant trend
  - how to specify  $c$ ?
- Venues for future work
  - locally
    - multivariate extensions
    - Can we ensure  $|\frac{\mathbf{s}_k \mathbf{y}}{\mathbf{s}_k}| \leq 1 + \frac{c}{\mathbf{s}_k}$  such that hyperparameters  $c$  and  $\rho$  collapse to one?
  - globally
    - Imprecise probabilities offer vivid framework to represent ignorance in surrogate-assisted derivative-free optimization



- Thanks a lot for your attention!
- Feel free to try out PROBO yourself:  
<https://github.com/rodemann/gp-imprecision-in-bo>
- We are looking forward to your feedback and comments of any kind!



- Rodemann, J.: *Robust Generalizations of Stochastic Derivative-Free Optimization*. Master's thesis, LMU Munich (2021) <sup>1</sup>
- Rodemann, J., Augustin, T.: *Accounting for Gaussian Process Imprecision in Bayesian Optimization*. In: Honda, K., Entani, T., Ubukata, S., Huynh, V.N., Inuiguchi, M. (eds.) IUKM. Springer Lecture Notes in Computer Science (LNCS). pp. 92–104. Springer, Cham (2022)

---


<sup>1</sup>[https://epub.ub.uni-muenchen.de/77441/1/MA\\_Rodemann.pdf](https://epub.ub.uni-muenchen.de/77441/1/MA_Rodemann.pdf)

- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Sensitivity Analysis
- 4 Prior-Mean-Robust BO (PROBO)
- 5 Application in Material Science
- 6 Discussion
- 7 Literature**
- 8 Appendix

 Benavoli, A. and Zaffalon, M. (2015).

Prior near ignorance for inferences in the  $k$ -parameter exponential family.

*Statistics*, 49(5):1104–1140.

 Bischl, B., Richter, J., Bossek, J., Horn, D., Thomas, J., and Lang, M. (2017).

mlrmo: A modular framework for model-based optimization of expensive black-box functions.

*arXiv preprint arXiv:1703.03373*.

 Bossek, J. (2017).

smoof: Single- and multi-objective optimization test functions.

*The R Journal*.





Kotthoff, L. (2019).

Ai for materials science: Tuning laser-induced graphene production and beyond.



Mangili, F. (2015).

A prior near-ignorance Gaussian process model for nonparametric regression.

In *ISIPTA '15: Proceedings of the 9th International Symposium on Imprecise Probability: Theories and Applications*, pages 187–196.



Rasmussen, C. E. (2003).

Gaussian processes in machine learning.

In *Summer school on machine learning*, pages 63–71. Springer.



Wahab, H., Jain, V., Tyrrell, A. S., Seas, M. A., Kotthoff, L., and Johnson, P. A. (2020).

Machine-learning-assisted fabrication: Bayesian optimization of laser-induced graphene patterning using in-situ raman analysis. *Carbon*, 167:609–619.



Wickham, H. (2016).

*ggplot2: Elegant Graphics for Data Analysis*.  
Springer-Verlag New York.

- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Sensitivity Analysis
- 4 Prior-Mean-Robust BO (PROBO)
- 5 Application in Material Science
- 6 Discussion
- 7 Literature
- 8 Appendix**

## Definition (Mean Optimization Path)

Given  $R$  repetitions of Bayesian optimization applied on a test function  $\Psi(\mathbf{x})$  with  $T$  iterations each, let  $\Psi(\mathbf{x}^*)_{r,t}$  be the best incumbent target value at iteration  $t \in \{1, \dots, T\}$  from repetition  $r \in \{1, \dots, R\}$ . The elements

$$MOP_t = \frac{1}{R} \sum_{r=1}^R \Psi(\mathbf{x}^*)_{r,t}$$

shall then constitute the  $T$ -dimensional vector  $MOP$ , which we call *mean optimization path (MOP)* henceforth.



## Definition (Accumulated Difference of MOPs)

Consider an experiment comparing  $S$  different prior specifications on a test function with  $R$  repetitions per specification and  $T$  iterations per repetition. Let the results be stored in a  $T \times S$ -matrix of mean optimization paths for iterations  $t \in \{1, \dots, T\}$  and prior specification  $s \in \{1, \dots, S\}$  (e.g. constant, linear, quadratic etc. trend as mean functional form) with entries  $MOP_{t,s} = \frac{1}{R} \sum_{r=1}^R \Psi(\mathbf{x}^*)_{r,t,s}$ . The *accumulated difference (AD)* for this experiment shall then be:

$$AD = \sum_{t=1}^T \left( \max_s MOP_{t,s} - \min_s MOP_{t,s} \right).$$

Mean functional form	Kernel functional form	Mean parameters	Kernel parameters
42.49	68.20	77.91	11.40

**Table:** Sum of relative ADs of all 50 MOPs per prior specification. Comparisons between mean and kernel are more valid than between functional form and parameters.



In order to derive upper and lower bounds for the mean estimate, let  $k_\theta(x, x')$  be a kernel function as defined in [Rasmussen, 2003]. The finitely positive semi-definite matrix  $\mathbf{K}_n$  is then formed by applying  $k_\theta(x, x')$  on the training data vector  $x \in \mathcal{X}$ :

$$\mathbf{K}_n = [k_\theta(x_i, x'_j)]_{ij}. \quad (3)$$

Let  $x$  be a scalar input of test data, whose  $f(x)$  is to be predicted. Then  $\mathbf{k}_x = [k_\theta(x, x_1), \dots, k_\theta(x, x_n)]^T$  is the vector of covariances between  $x$  and the training data. Furthermore, name the training target vector  $\mathbf{y}$  and define  $\mathbf{s}_k = \mathbf{K}_n^{-1} \mathbf{1}_n$  as well as  $\mathbf{S}_k = \mathbf{1}_n^T \mathbf{K}_n^{-1} \mathbf{1}_n$ .

Then [Mangili, 2015] shows that if  $|\frac{\mathbf{s}_k \mathbf{y}}{\mathbf{s}_k}| \leq 1 + \frac{c}{\mathbf{s}_k}$ :

$$\bar{\hat{\mu}}(x) = \mathbf{k}_x^T \mathbf{K}_n^{-1} \mathbf{y} + (1 - \mathbf{k}_x^T \mathbf{s}_k) \frac{\mathbf{s}_k^T}{\mathbf{s}_k} \mathbf{y} + c \frac{|1 - \mathbf{k}_x^T \mathbf{s}_k|}{\mathbf{s}_k} \quad (4)$$

$$\underline{\hat{\mu}}(x) = \mathbf{k}_x^T \mathbf{K}_n^{-1} \mathbf{y} + (1 - \mathbf{k}_x^T \mathbf{s}_k) \frac{\mathbf{s}_k^T}{\mathbf{s}_k} \mathbf{y} - c \frac{|1 - \mathbf{k}_x^T \mathbf{s}_k|}{\mathbf{s}_k} \quad (5)$$





If  $\left| \frac{s_k y}{s_k} \right| > 1 + \frac{c}{s_k}$ :

$$\bar{\hat{\mu}}(x) = k_x^T K_n^{-1} y + (1 - k_x^T s_k) \frac{s_k^T y}{s_k} + c \frac{1 - k_x^T s_k}{s_k} \quad (6)$$

$$\underline{\hat{\mu}}(x) = k_x^T K_n^{-1} y + (1 - k_x^T s_k) \frac{s_k^T y}{c + s_k} \quad (7)$$