

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

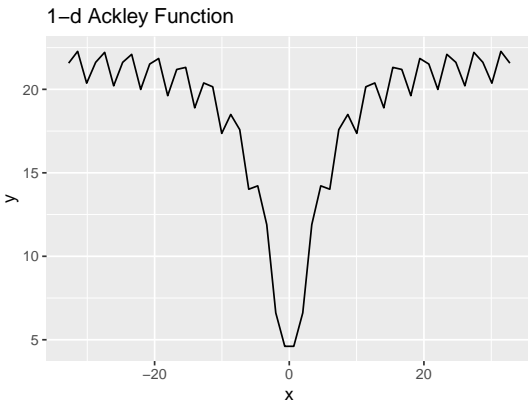
# Agenda

- 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

# Agenda

- 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

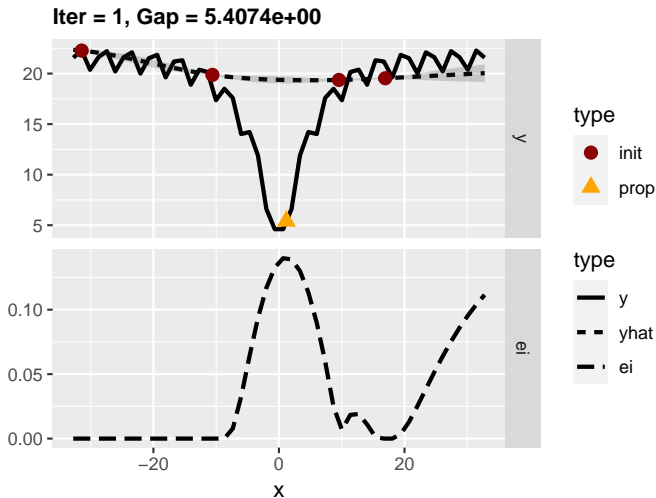
# Bayesian Optimization



---

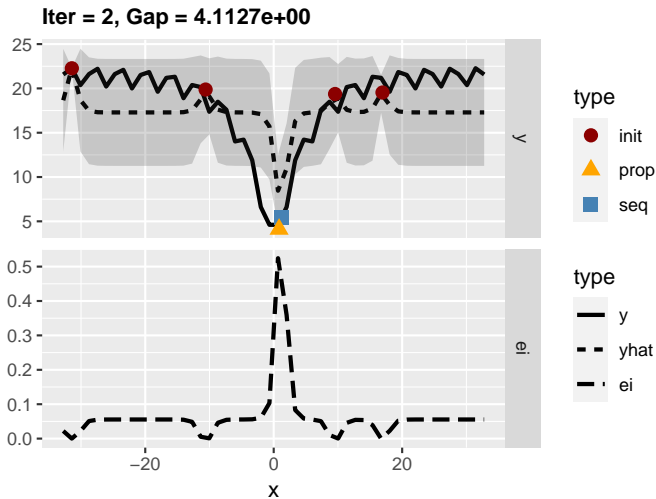
Note: If not otherwise stated, all figures are based on own computations using `ggplot2` [Wickham, 2016], `smoof` [Bossek, 2017] and `mlr(3)MBO` [Bischl et al., 2017]

# Bayesian Optimization



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

# Bayesian Optimization



# Bayesian Optimization

Iteration 3

# Bayesian Optimization

Iteration 4



# Bayesian Optimization

Iteration 5

# Bayesian Optimization

Iteration 6

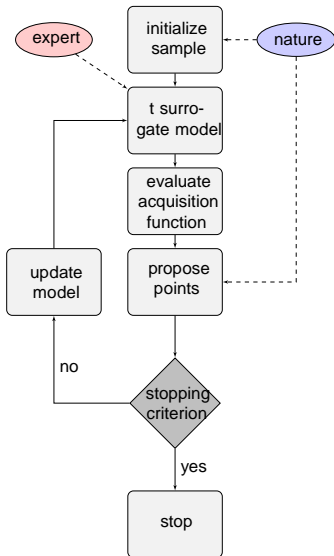
# Bayesian Optimization

Iteration 7

# Bayesian Optimization

Iteration 10

# Bayesian Optimization



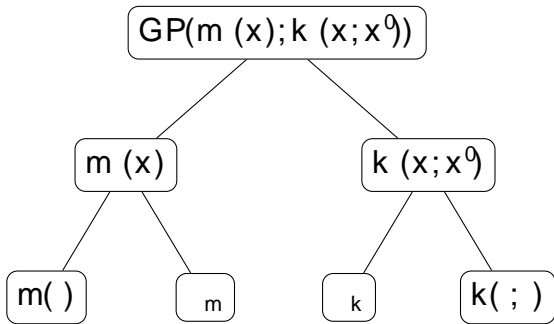
# Agenda

- 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

## Gaussian Processes - Intuition

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

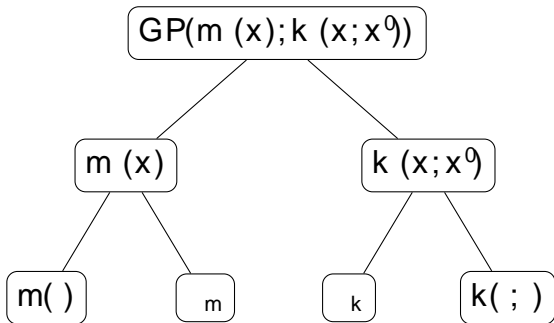
# Gaussian Processes { Prior Components



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



## Gaussian Processes { Prior Components



And: Do they even matter?

# Agenda

- 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

# Agenda

- 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

## Setup

We randomly select 50 synthetic test functions from the packages `moof` [Bossek, 2017], stratified across the covariate space dimensions  $d \in \{2, 3, 4, 7\}$ .

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  $n_{\text{init}} = 10$  is randomly sampled anew for each of the  $R = 40$  BO repetitions with  $T = 20$  iterations each.

# Agenda

- 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

## Results

Mean parameters influence convergence the most, followed by the kernel's functional form .

Mean functional form and Kernel parameters play a (relatively) negligible role.

# Agenda

- 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

# Agenda

- 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



## Prior near-ignorance models

Idea: Use set of  $m$  instead of precise  $m$ . Fully specify the other components.

[Mangili, 2015] proposes imprecise Gaussian processes

$$\text{GP } Mh; k(x; x^0) + \frac{1 + M}{c} : h = 1; M \geq 0 ;$$

given a base kernel  $k(x; x^0)$  and a degree of imprecision  $c > 0$ .

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

# Agenda

- 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

# Generalized Lower Confidence Bound (GLCB)

$$\text{LCB}(x) = b(x) + \frac{q}{\sqrt{|\{z\}|}} \text{Var}(x)$$

"classical" uncertainty

$$\text{GLCB}(x) = b(x) + \frac{q}{\sqrt{|\{z\}|}} \text{Var}(x) + \frac{(-x)_c}{\sqrt{|\{z\}|}} (x)_c$$

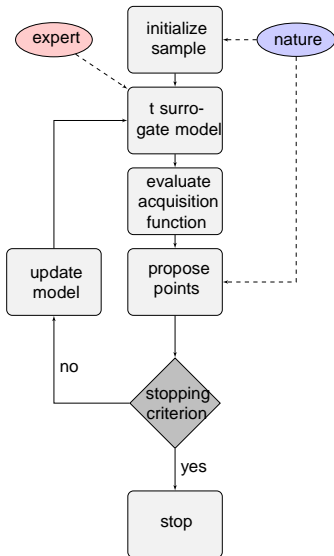
"classical" uncertainty

prior-induced imprecision

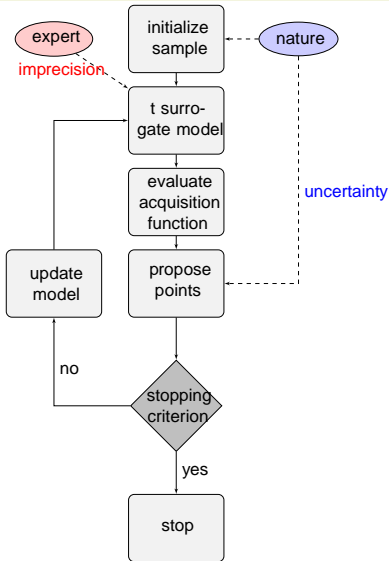
is the degree of **risk**-aversion

is the degree of **ambiguity**-aversion

# Bayesian Optimization



# Bayesian Optimization



## Generalized Lower Confidence Bound (GLCB)

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

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

---

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

## Generalized Lower Confidence Bound (GLCB)

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

$$\bar{\hat{h}}(x) - \underline{\hat{h}}(x) = 2c \frac{\mathbf{k}_x^\top \mathbf{s}_k}{S_k} \quad (2)$$

In this case, GLCB's hyperparameter  $\alpha$  and  $c$  collapse to one.

# Agenda

- 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



# Application in Material Science

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

## GLCB { Results

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

# Agenda

- 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

# Discussion

## Limitations

robust only with regard to possible misspecification of the mean function parameter given a constant trend  
how to specify?

## Venues for future work

locally

multivariate extensions

Can we ensure  $\frac{S_k y_j}{S_k} \approx 1 + \frac{c}{S_k}$  such that hyperparameters  $\alpha$  and  $\beta$  collapse to one?

globally

Imprecise probabilities offer a vivid framework to represent ignorance in surrogate-assisted derivative-free optimization

## Discussion

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!

## PROBO: Literature

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)

# Agenda

- 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

## Literature I

-  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*.



## Literature II



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.

## Literature III

- 📄 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.

# Agenda

- 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**

## Mean Optimization Path

### Definition (Mean Optimization Path)

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

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

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

# Accumulated Difference of Mean Optimization Paths

## 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  $MOP_{t,s} = \frac{1}{R} \sum_{r=1}^R \Psi(x)_{r,t;s}$ . The *accumulated difference (AD)* for this experiment shall then be:

$$AD = \sum_{t=1}^T \max_s MOP_{t,s} - \min_s MOP_{t,s} :$$

## Results

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.

## Upper and lower mean estimates

In order to derive upper and lower bounds for the mean estimate, let  $k(x; x^0)$  be a kernel function as defined in [Rasmussen, 2003]. The finitely positive semi-definite matrix  $K_n$  is then formed by applying  $k(x; x^0)$  on the training data vector  $x \in X$ :

$$K_n = [k(x_i; x_j^0)]_{ij} \quad (3)$$

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

## Upper and lower mean estimates

Then [Mangili, 2015] shows that if  $j \frac{S_k y}{S_k} j = 1 + \frac{c}{S_k}$ :

$$\hat{\bar{}}(x) = k_x^T K_n^{-1} y + (1 - k_x^T S_k) \frac{S_k^T}{S_k} y + c \frac{j 1 - k_x^T S_k j}{S_k} \quad (4)$$

$$\hat{\underline{}}(x) = k_x^T K_n^{-1} y + (1 - k_x^T S_k) \frac{S_k^T}{S_k} y - c \frac{j 1 - k_x^T S_k j}{S_k} \quad (5)$$



## Upper and lower mean estimates

If  $j \frac{S_k y}{S_k} j > 1 + \frac{c}{S_k}$ :

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

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