

# Robust Surrogate Models in Bayesian Optimization

Julian Rodemann

Summer Retreat  
Department of Statistics

23.10.2021





- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO
  - Prior near-ignorance models
  - GLCB
- 4 Weighted ML Estimation of  $\theta_m$ 
  - Problem
  - Idea
  - Results
- 5 Questions?
- 6 Questions!
- 7 References

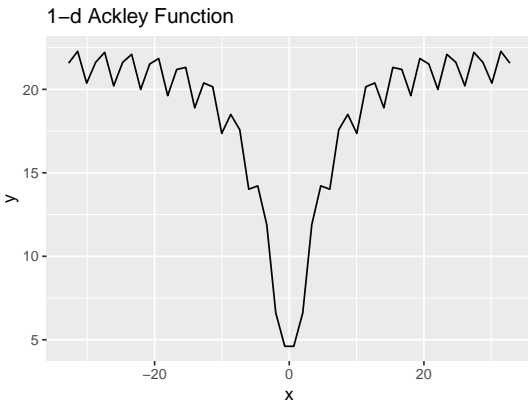


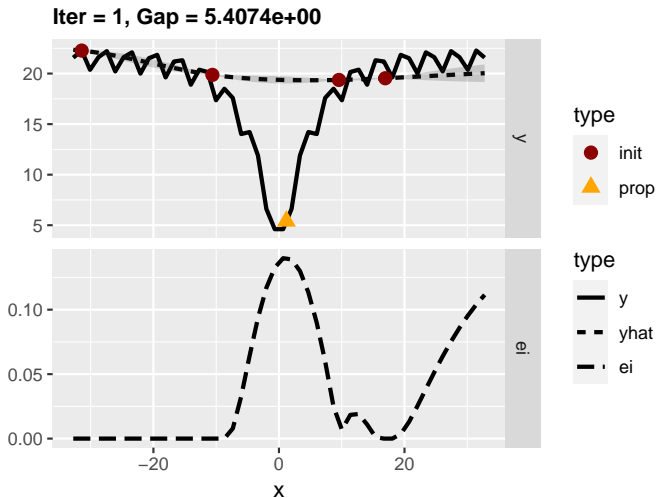
- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO
- 4 Weighted ML Estimation of  $\theta_m$
- 5 Questions?
- 6 Questions!
- 7 References



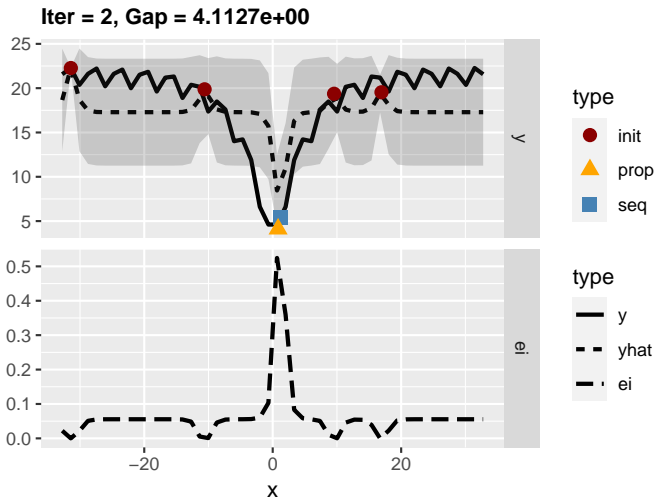
Everyone heard of Bayesian/model-based optimization before?

[skip intro](#)

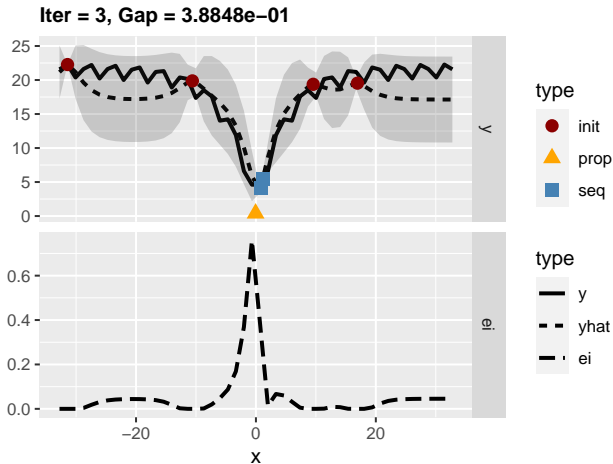




Iteration 1: Surrogate Model (top) and Acquisition Function (bottom)



Iteration 2

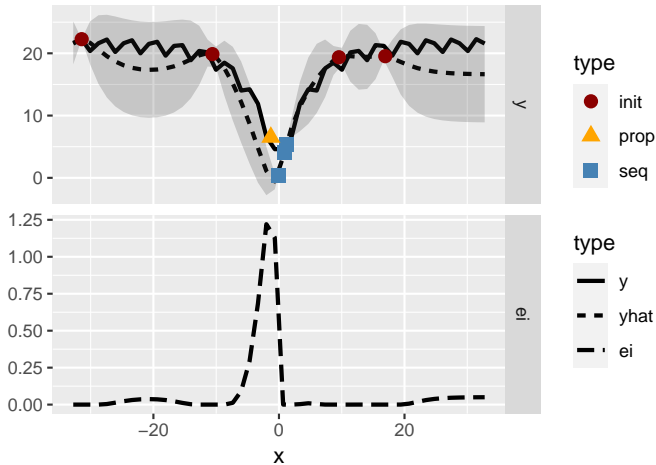


Iteration 3





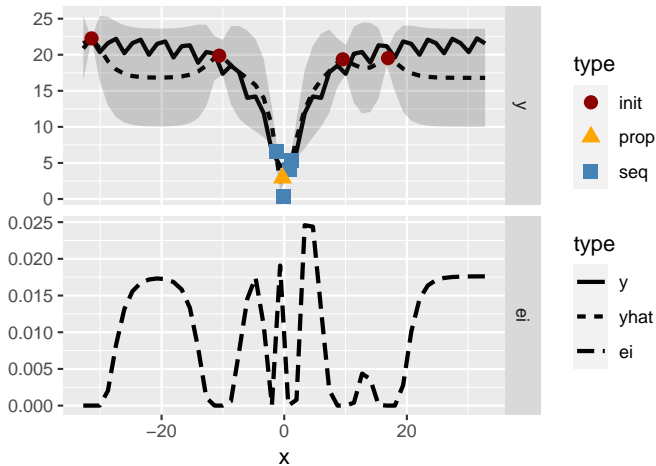
Iter = 4, Gap = 3.8848e-01



Iteration 4



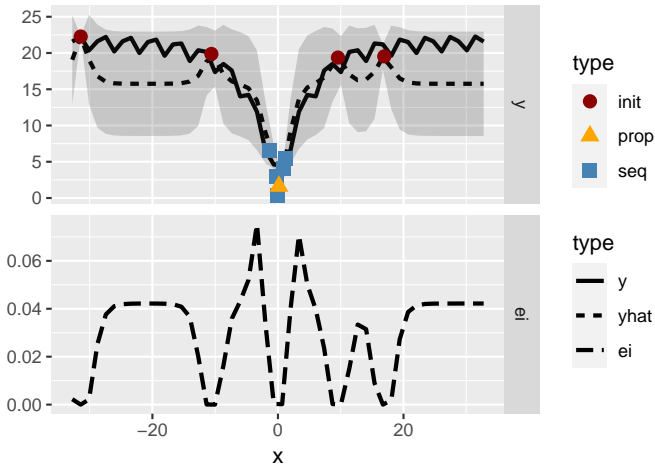
Iter = 5, Gap = 3.8848e-01



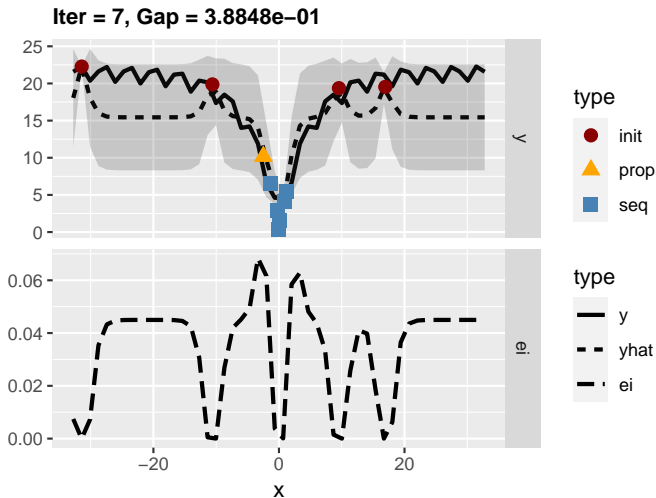
Iteration 5



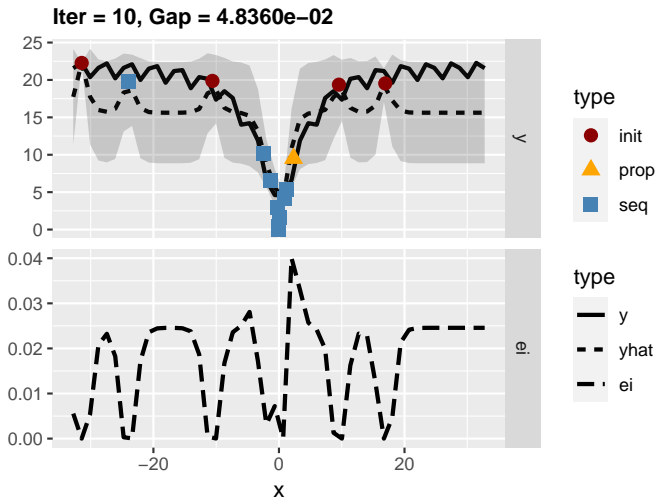
Iter = 6, Gap = 3.8848e-01

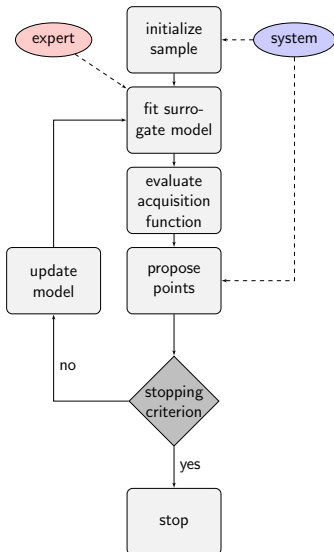


Iteration 6



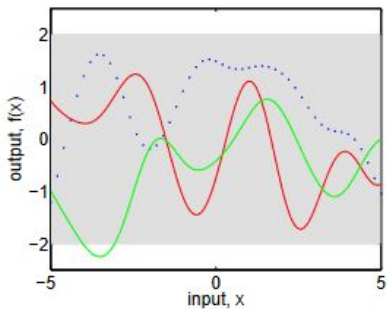
Iteration 7



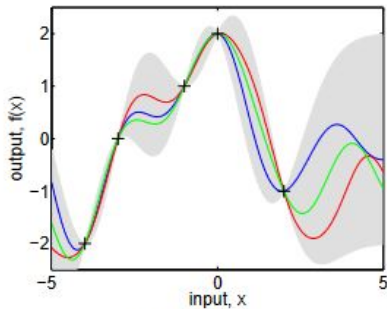




- 1 Bayesian Optimization
- 2 Gaussian Processes**
- 3 Prior-Mean-Robust BO
- 4 Weighted ML Estimation of  $\theta_m$
- 5 Questions?
- 6 Questions!
- 7 References



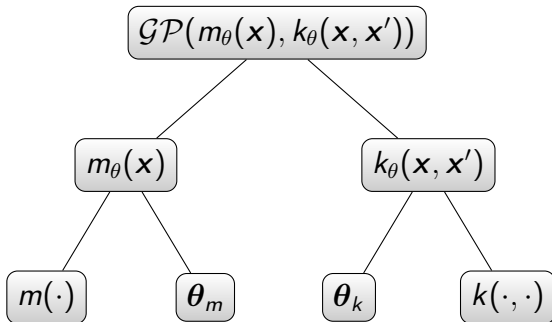
(a), prior



(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?



- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO**
  - Prior near-ignorance models
  - GLCB
- 4 Weighted ML Estimation of  $\theta_m$
- 5 Questions?
- 6 Questions!
- 7 References



- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO**
  - Prior near-ignorance models
  - GLCB
- 4 Weighted ML Estimation of  $\theta_m$
- 5 Questions?
- 6 Questions!
- 7 References



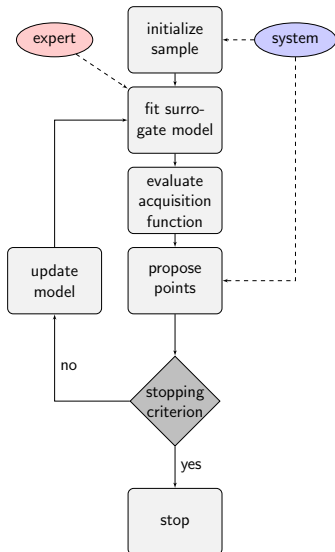
- Idea: Use set of  $\theta_m$  instead of precise  $\theta_m$ . Fully specify the other components.
- [Mangili, 2015] proposes imprecise Gaussian processes  $\{\mathcal{GP}(Mh, k_\theta(x, x') + \frac{1+M}{c}) : h = \pm 1, M \geq 0\}$ , given a base kernel  $k_\theta(x, x')$  and a degree of imprecision  $c > 0$ .
  - set of posteriors with upper and lower mean estimates  $\underline{\mu}(x)_c$ ,  $\bar{\mu}(x)_c$

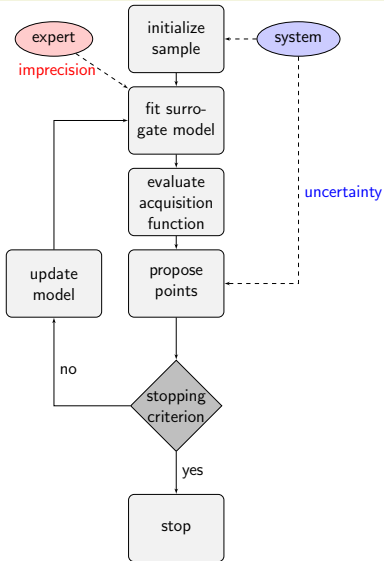


- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO**
  - Prior near-ignorance models
  - GLCB
- 4 Weighted ML Estimation of  $\theta_m$
- 5 Questions?
- 6 Questions!
- 7 References

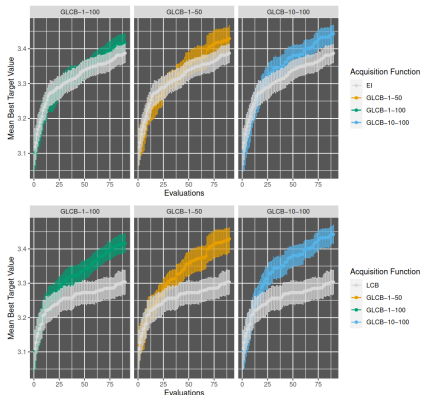


- $LCB(x) = -\hat{\mu}(x) + \tau \cdot \sqrt{\widehat{\text{Var}}(\mu(x))}$
- $GLCB(x) = -\hat{\mu}(x) + \tau \cdot \underbrace{\sqrt{\widehat{\text{Var}}(\mu(x))}}_{\text{"classical" uncertainty}} + \rho \cdot \underbrace{(\bar{\mu}(x)_c - \underline{\mu}(x)_c)}_{\text{prior-induced imprecision}}$ 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









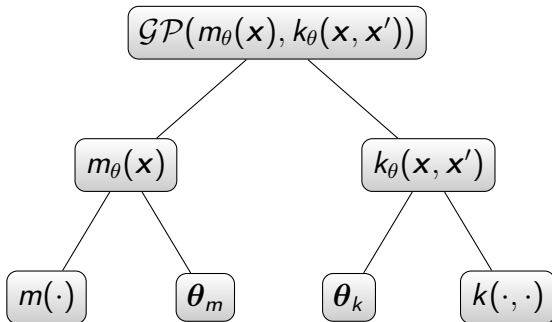
BO with GLCB on Graphene function. GLCB-1-50 means GLCB with  $\rho = 1$ ,  $c = 50$ . Data source: [Wahab et al., 2020].



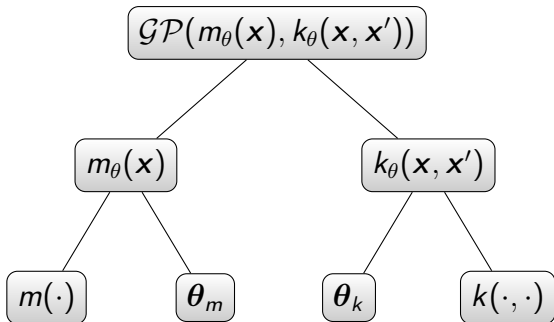
- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO
- 4 Weighted ML Estimation of  $\theta_m$**   
Problem  
Idea  
Results
- 5 Questions?
- 6 Questions!
- 7 References



- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO
- 4 Weighted ML Estimation of  $\theta_m$**   
Problem  
Idea  
Results
- 5 Questions?
- 6 Questions!
- 7 References



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



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

**Answer**<sup>1</sup>: “Empirical Bayes”, i.e.  $\hat{\theta}_m = \arg \max_{\theta_m} \mathcal{L}(\theta_m | \mathbf{X}_t)$ , where  $\mathbf{X}_t$  is the incumbent design of iteration  $t$ .

<sup>1</sup>given by most software libraries



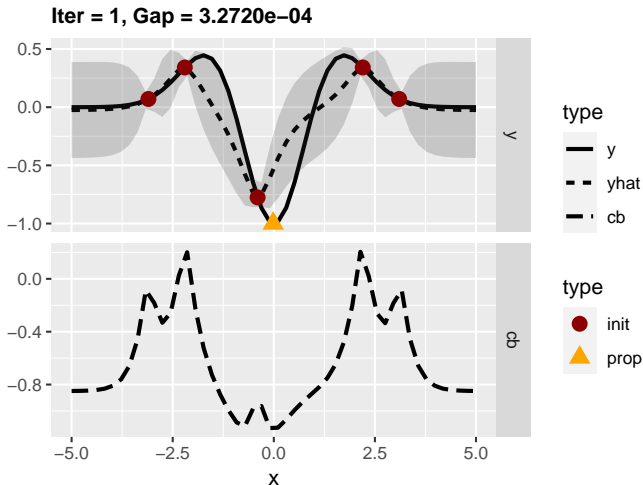
- $\hat{\theta}_m = \arg \max_{\theta_m} \mathcal{L}(\theta_m | \mathbf{X}_t)$  requires an iid. sample

$$\mathcal{L}(\theta | \mathbf{X}_t) \stackrel{\text{ind.}}{=} \prod_{i=1}^n \ell(\theta_m | \mathbf{x}_i) \stackrel{\text{iden.}}{=} \prod_{i=1}^n f(\theta_m | \mathbf{x}_i),$$

- iid. assumption is fulfilled for initial sample
- ⚡ in the course of the optimization, however,  $\mathbf{X}_t$  becomes biased
- ⚡ can slow down BO performance

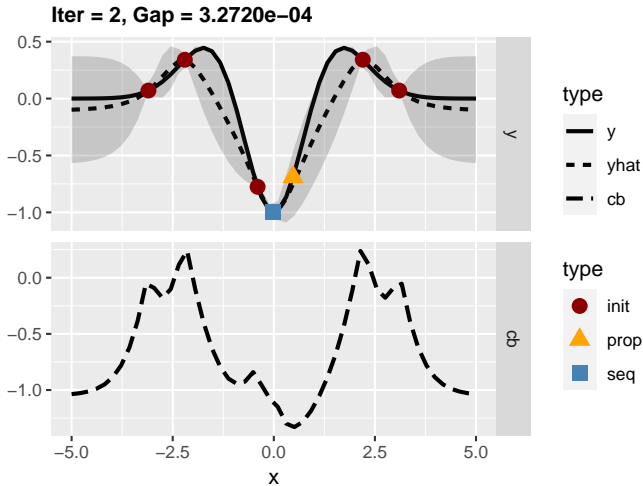


Image credits: Pixabay (cc license)

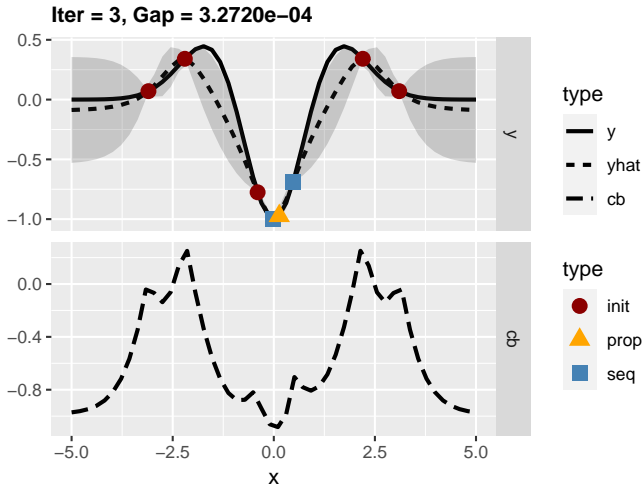


Iteration 1 of BO on Mexican Hat Function

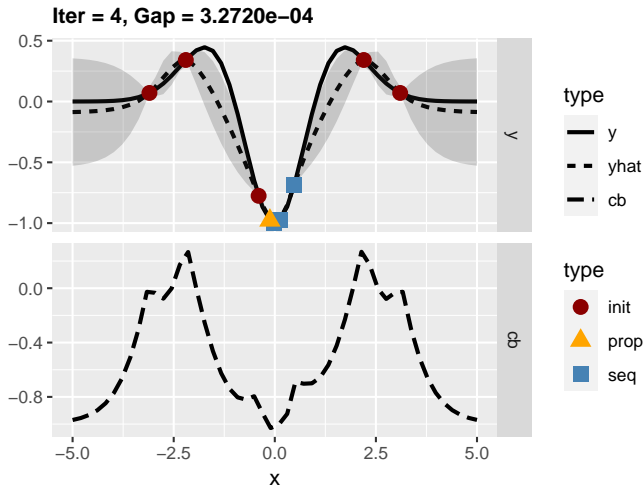




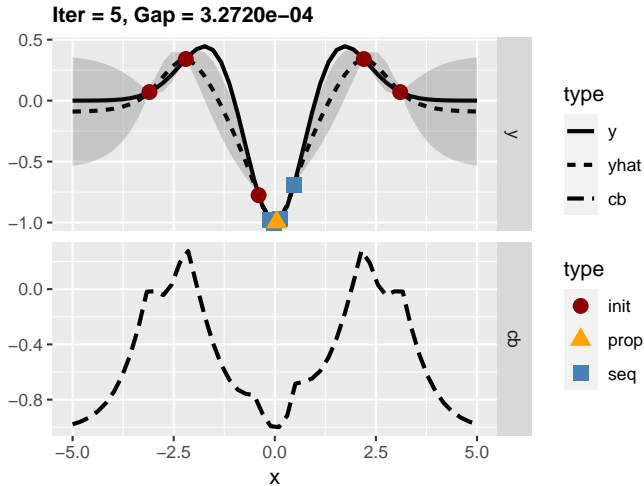
Iteration 2



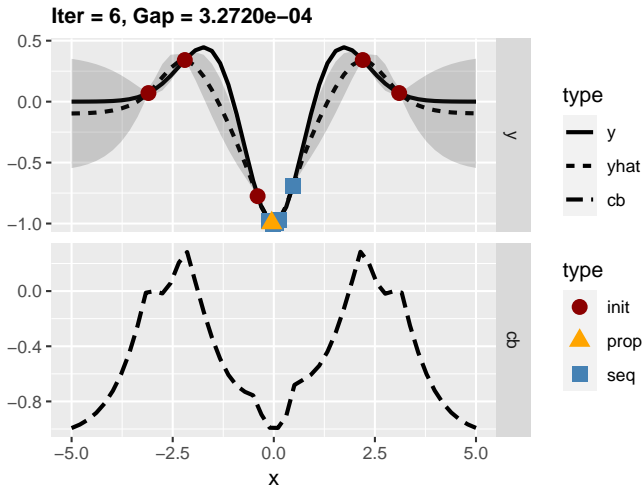
Iteration 3



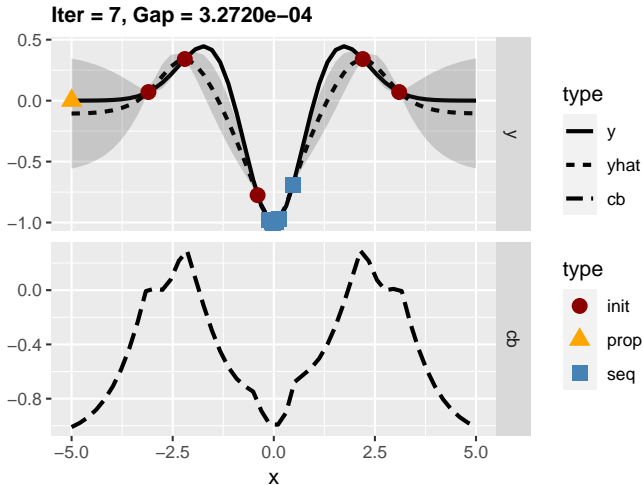
Iteration 4



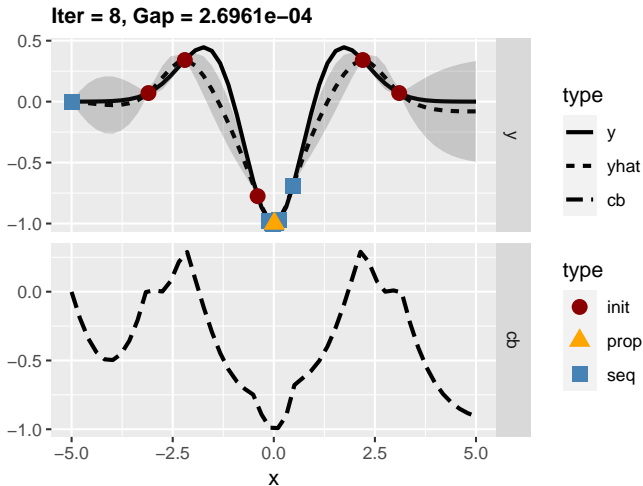
Iteration 5



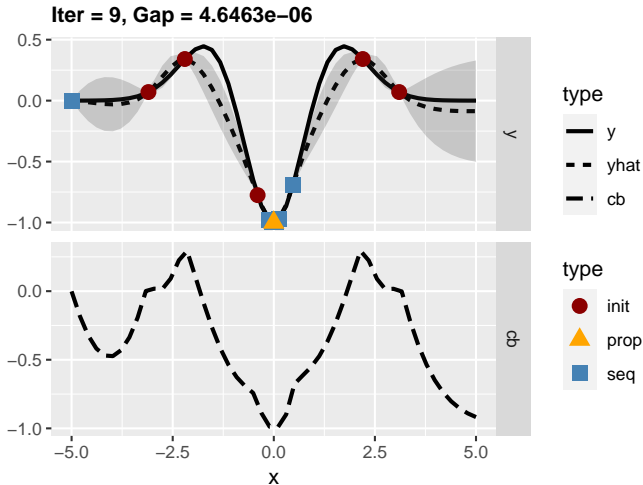
Iteration 6



Iteration 7

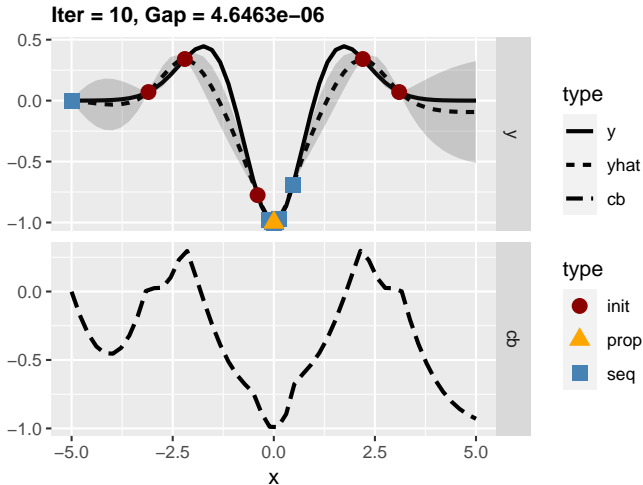


Iteration 8

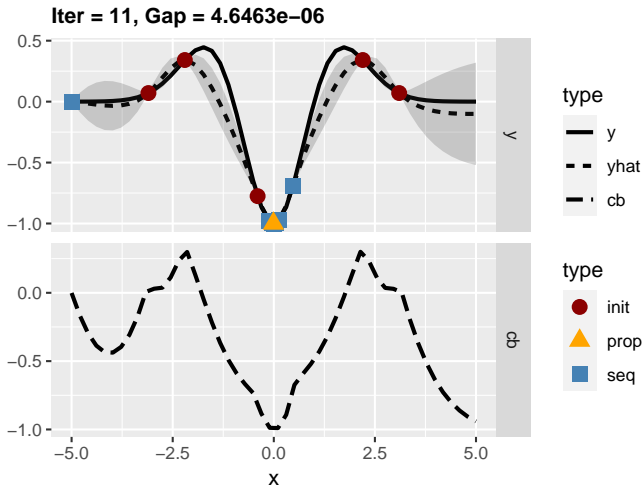


Iteration 9

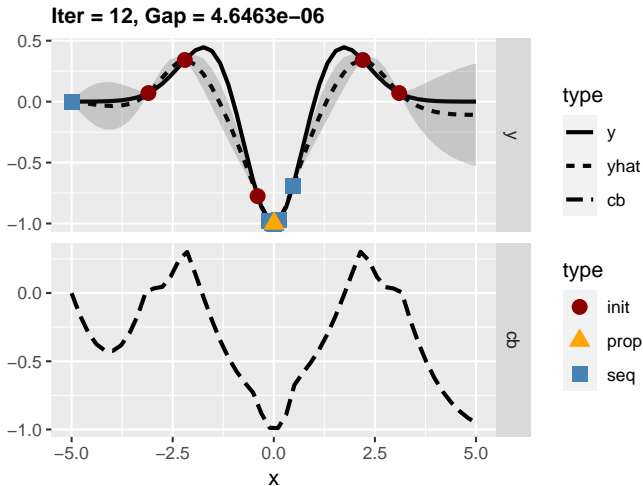




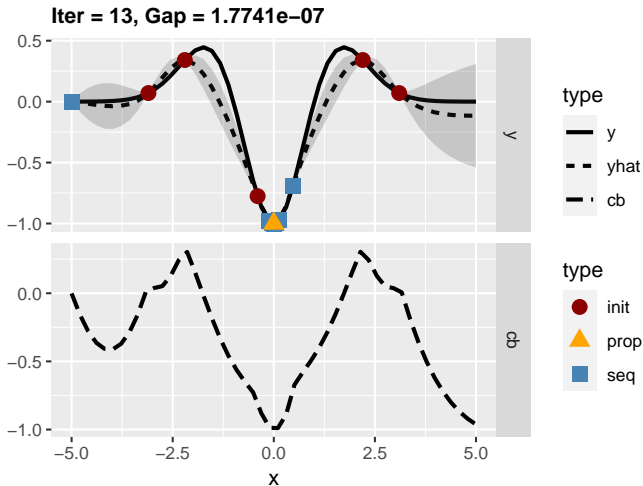
Iteration 10



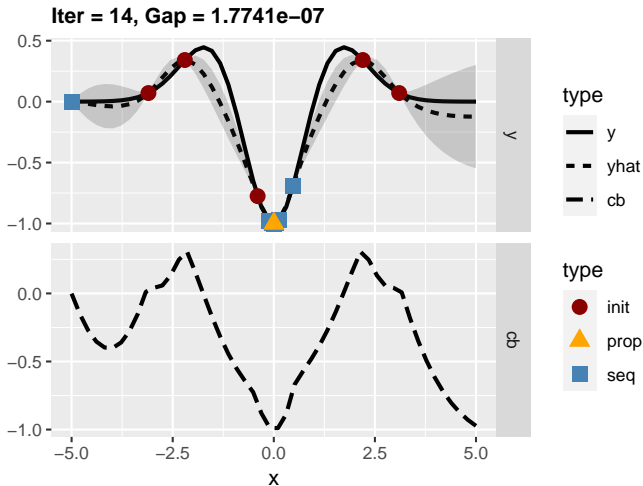
Iteration 11



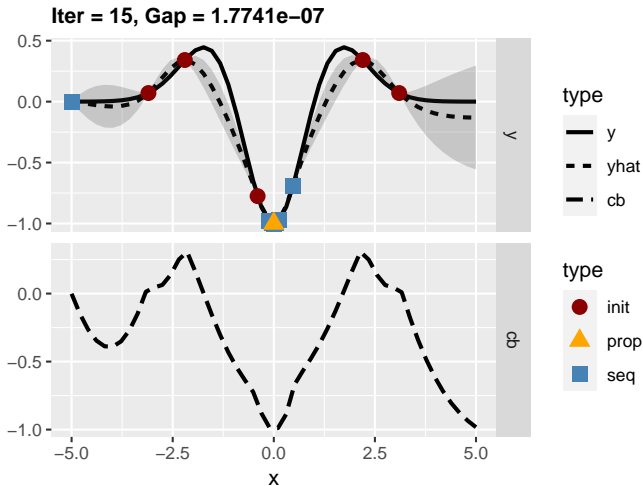
Iteration 12



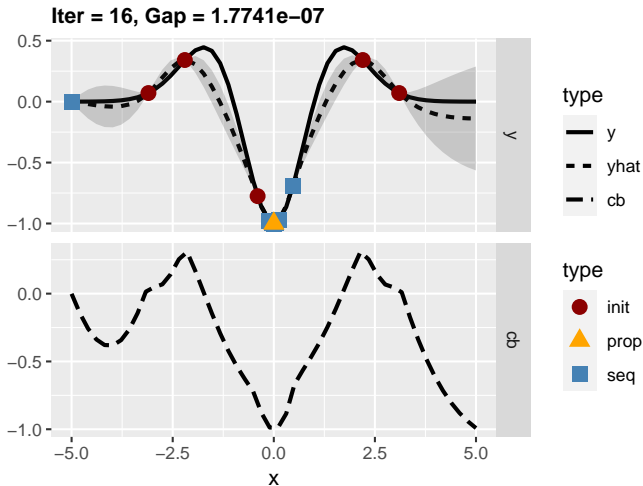
Iteration 13



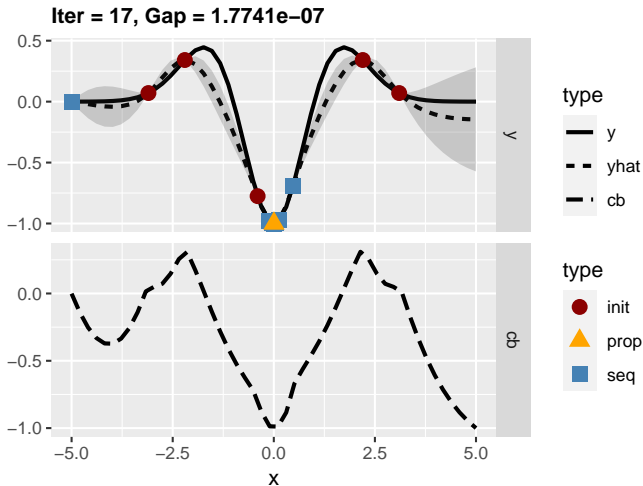
Iteration 14



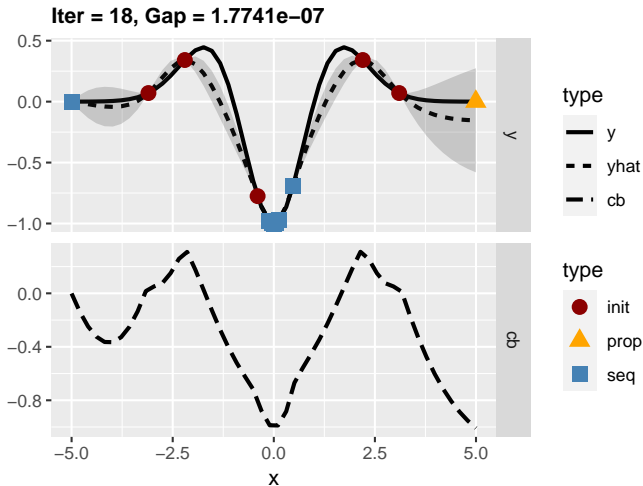
Iteration 15



Iteration 16







Iteration 18



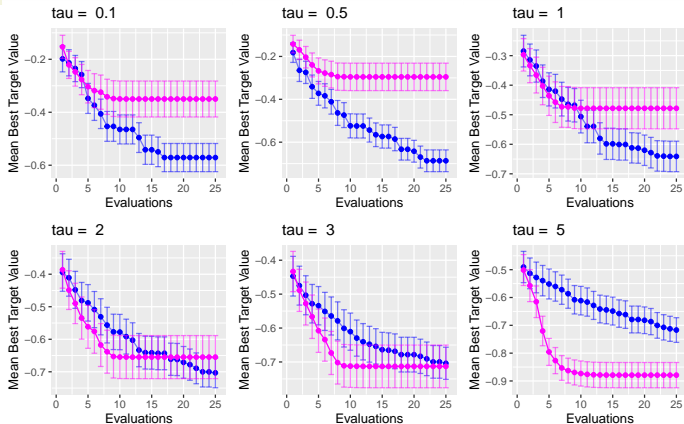
- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO
- 4 Weighted ML Estimation of  $\theta_m$** 
  - Problem
  - Idea
  - Results
- 5 Questions?
- 6 Questions!
- 7 References



- Idea: Weight by potential gain of information at time of proposal
- Weights:
  - Use variance (standard error) estimation at proposed point
  - Compare to variances at  $n$  randomly sampled points
  - Use empirical distribution function  $F(\bullet)$
  - The weight  $w_j$  of  $\mathbf{x}_j$  then is  $w_j = \frac{F(\mathbf{x}_j)}{\sum_i^n F(\mathbf{x}_i)}$ .
- Estimation:
  - Draw  $\mathbf{x}_j$  (with replacement) with probability  $w_j$
  - $\hat{\theta}_m(\mathbf{X}_t) = \arg \max_{\theta_m} \mathcal{L}(\theta_m | \mathbf{X}_t)$ , where  $\mathbf{X}_t$  is the design matrix of the so-generated sample



- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO
- 4 Weighted ML Estimation of  $\theta_m$** 
  - Problem
  - Idea
  - Results
- 5 Questions?
- 6 Questions!
- 7 References



Benchmarking of BO with weighted ML (blue) against classic unweighted ML (magenta) on Mexican hat function with varying  $\tau$  in LCB. 200 BO runs with initial sample size 4. Error bars depict 95%-CIs.



- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO
- 4 Weighted ML Estimation of  $\theta_m$
- 5 Questions?**
- 6 Questions!
- 7 References



Do you have any questions?



- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO
- 4 Weighted ML Estimation of  $\theta_m$
- 5 Questions?
- 6 Questions!**
- 7 References





- On what type of problems is weighted ML superior to standard ML?
  - Distance from local to global optima
  - Edges?
- Why inferior to standard ML on wiggly functions?
- Benchmark GLCB (and weighted ML?) against “integrated acquisition function” [Snoek et al., 2012], i.e. against (improper) uniform hyperpriors<sup>2</sup>




---

<sup>2</sup>This is the native bayesian way to represent uncertainty over the prior. However, note that such a uniform prior reflects *indifference* rather than *ignorance*.



- 1 Bayesian Optimization
- 2 Gaussian Processes
- 3 Prior-Mean-Robust BO
- 4 Weighted ML Estimation of  $\theta_m$
- 5 Questions?
- 6 Questions!
- 7 References



-  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*.



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.





Snoek, J., Larochelle, H., and Adams, R. P. (2012).

Practical Bayesian optimization of machine learning algorithms.

In *Proceedings of the 25th International Conference on Neural Information Processing Systems-Volume 2*, pages 2951–2959.



-  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.