

DCE Reading Group @ ATI

Numerical issues in maximum likelihood parameter estimation for  
**Gaussian process interpolation**

Subhasish Basak<sup>1</sup>, Sébastien Petit<sup>1,2</sup>, Julien Bect<sup>1</sup> & Emmanuel Vazquez<sup>1</sup>

March 10, 2021

1. Laboratoire des Signaux et Systèmes, CNRS, CentraleSupélec, Univ. Paris-Saclay
2. Safran Aircraft Engines, France

This work is licensed under a [Creative Commons BY-NC-ND 4.0](https://creativecommons.org/licenses/by-nc-nd/4.0/) license.



<https://creativecommons.org/licenses/by-nc-nd/4.0/legalcode>

DOI: 10.5281/zenodo.4653846

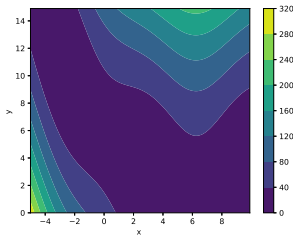
## Motivation & Scope

- **Gaussian processes (GP)**: Popular tool for regression/interpolation, widely used in the Statistics and ML community
  - **Geostatistics** (Stein, 1999),
  - **Design & analysis of computer experiments** (Santner et al., 2003),
  - **Machine Learning** (Rasmussen & Williams, 2006),
  - **Bayesian optimization** (Mockus, 1975; Jones, 1998; Emmerich et al., 2006; ...).
- Applications rely highly on off-the-shelf GP implementations.
- **Problem**: Lack of consistency and robustness (see Erickson et al., 2018) among available software packages (Python, R, Matlab).

## GP modelling with Python packages

- Consider the **Branin** function;  $(x_1, x_2) \in [-5, 10] \times [0, 15]$

$$f(x_1, x_2) = (x_2 - \frac{5.1}{4\pi^2}x_1^2 + \frac{5}{\pi}x_1 - 6)^2 + 10(1 - \frac{1}{8\pi})\cos(x_1) + 10$$



- $\xi \sim \text{GP}(0, k)$  with  $k$  a Matérn kernel ( $\nu = 5/2$ ).
- Training (testing) on dataset of size 50 (500), sampled from a uniform grid.

Optimized negative log likelihood (NLL) & prediction error (ERMSPE).

Estimates	scikit-learn	OpenTURNS	GPy	GPflow	GPy improved
NLL	132.421	163.125	113.707	113.223	112.050
ERMSPE	1.482	3.301	0.259	0.236	0.175

- Efficient **optimization of NLL** is crucial for robust and reliable applications.
- The objective of our article is two-fold
  - Investigate the origin of these inconsistencies.
  - Propose effective strategies for improvement.

# Contents

- 1 Background : GP & MLE**
- 2 Numerical noise**
- 3 Improvement strategies**
- 4 Numerical study**
- 5 Concluding remarks**

# 1 GP & Maximum likelihood estimation

- Consider a data set  $D = \{(x_i, z_i) \in \mathbb{R}^d \times \mathbb{R}, 1 \leq i \leq n\}$  and an additive-noise model

$$Z_i = \xi(x_i) + \varepsilon_i,$$

- where
  - $\xi$  is a **Gaussian process**  $\text{GP}(m, k)$
  - mean function  $m : \mathbb{R}^d \rightarrow \mathbb{R}$ , kernel  $k : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$
  - $\varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma_\varepsilon^2)$ , independent of  $\xi$ .

- Model parameters (to be estimated):
  - $c \rightsquigarrow$  scalar for constant  $m$  ( $c \in \mathbb{R}$ )
  - $\rho_i \rightsquigarrow$  lengthscales of  $k$  ( $1 \leq i \leq d$ )
  - $\sigma^2 \rightsquigarrow$  variance of  $k$  ( $> 0$ )
  - $\sigma_\varepsilon^2 \rightsquigarrow$  noise variance ( $\geq 0$ )
- The predictive **posterior distribution** of  $\xi$  is then obtained as

$$\xi \mid \underline{Z}_n, m, k \sim GP(\hat{\xi}_n, k_n),$$

The posterior mean  $\hat{\xi}_n$  and covariance  $k_n$  is computed by solving a system of linear equations (see Rasmussen and Williams, 2006).



## Parameter estimation with MLE

- Let  $K_\theta = (k(x_i, x_j))_{n \times n} + \sigma_\varepsilon^2 \mathbf{I}_n$ ,  $x_i \in \underline{x}_n$ , be the covariance matrix with parameters  $\theta = (\sigma^2, \rho_1, \dots, \rho_d, \sigma_\varepsilon^2)^\top$ .
- The **likelihood** of  $\underline{Z}_n$  is given by

$$\mathcal{L}(\underline{Z}_n | \theta, c) = \frac{1}{(2\pi)^{\frac{n}{2}} \sqrt{|K_\theta|}} \exp\left(-\frac{1}{2} (\underline{Z}_n - c\mathbf{1}_n)^\top K_\theta^{-1} (\underline{Z}_n - c\mathbf{1}_n)\right). \quad (1)$$

- The negative log likelihood (NLL) is

$$-\log(\mathcal{L}(\underline{Z}_n | \theta, c)) = \frac{1}{2} (\underline{Z}_n - c\mathbf{1}_n)^\top K_\theta^{-1} (\underline{Z}_n - c\mathbf{1}_n) + \frac{1}{2} \log|K_\theta| + \text{constant}. \quad (2)$$

- Most GP packages estimate the parameters by minimizing the NLL.

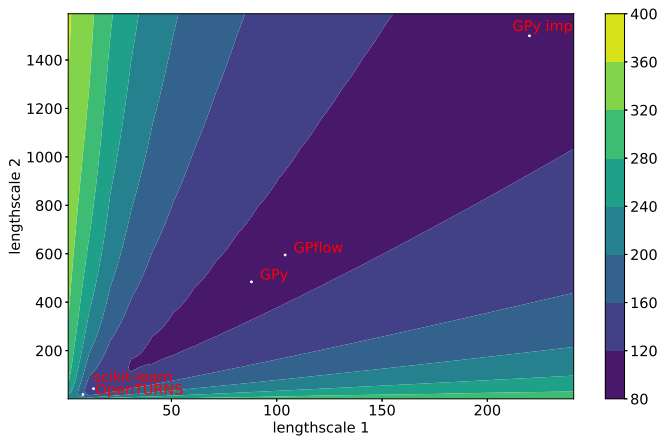


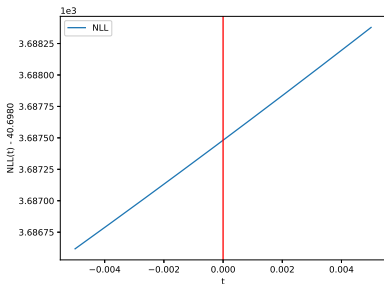
Figure 1: Contour of NLL along with estimated lengthscales.

## Optimizing the NLL

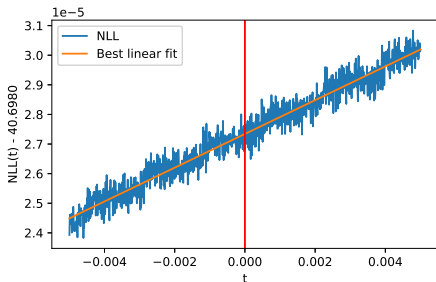
- The optimization is typically performed by **gradient-based** algorithms.
  - GPyTorch uses ADAM (Kingma and Ba, 2015)
  - OpenTURNS uses TNC (Nash, 1984)
  - others generally use L-BFGS-B (Byrd et al., 1995).
- Associated computational cost  $O(n^3 + dn^2)$  (Rasmussen and Williams, 2006; Petit et al., 2020).
- In this approach the **non-convexity** of the likelihood can be of significant concern.

## 2 Numerical noise

- Evaluation of the likelihood is susceptible to numerical noise
- Prevents proper convergence (early stopping) of the optimizer



→ here the optimizer stalls at 0.



- Numerical noise stems from both,  $\frac{1}{2}\underline{Z}_n^\top \mathbf{K}_\theta^{-1} \underline{Z}_n$  and  $\frac{1}{2} \log|\mathbf{K}_\theta|$ .
- Our paper details that this noise is directly linked to the **condition number**  $\kappa(\mathbf{K}_\theta)$ .
- Even when  $\kappa(\mathbf{K}_\theta)$  is standard, the noise can prevent proper convergence of the optimizer.

## A standard solution : using jitter

- A **high condition number** (equivalently an ill-conditioned matrix) can occur if
  - $\sigma_\varepsilon^2 = 0$  (for an interpolating model)
  - very smooth covariance function (e.g. squared exponential)
  - high lengthscale values
  - very close datapoints
- **Jitter**: small positive quantity, added to the diagonal of  $K_\theta$ .
- Lowers  $\kappa(K_\theta)$ , but produces **non interpolating** model.
- Example: GPy uses iterative jitter ranging from  $10^{-6}\sigma^2$  to  $10^{-1}\sigma^2$ .
- The literature includes other standard ways to choose and implement jitter (see Ranjan et al.).

## Effect of jitter

$\sigma_{\varepsilon}^2 / \sigma^2$	0.0	$10^{-8}$	$10^{-6}$	$10^{-4}$	$10^{-2}$
$\kappa(\mathbf{K}_{\theta})$	$10^{11}$	$10^9$	$10^{7.5}$	$10^{5.5}$	$10^{3.5}$
$\sqrt{\text{SSR}/\text{SST}}$	$3.3 \cdot 10^{-10}$	$1.2 \cdot 10^{-3}$	0.028	0.29	0.75
NLL	40.69	45.13	62.32	88.81	124.76

- Model's predictive performance deteriorates, causing possible convergence problems in **Bayesian optimization**.
- **Conclusion**: Jitter is not a satisfactory solution to numerical noise.

### 3 Improvement strategies

- **Goal:** Prevent early stopping of the optimizer and robust estimation of parameters.
- Strategies considered for improving the optimizer
  - parameter initialization methods
  - stopping criterion for optimization
  - “restart” strategies
  - parameterization of the covariance



## Initialization strategies

The choice of good initialization point seems important.

- **Moment-based**: empirical moments of the data used to initialize the parameters

$$c_{\text{init}} = \text{mean}(Z_1, \dots, Z_n), \quad (3)$$

$$\sigma_{\text{init}}^2 = \text{var}(Z_1, \dots, Z_n), \quad (4)$$

$$\rho_{k, \text{init}} = \text{std}(x_{1, [k]}, \dots, x_{n, [k]}), \quad k = 1, \dots, d, \quad (5)$$

→ Available in GPy.

- **Profiled:** Initialize lengthscales using (5) then take the generalized least square (GLS) solutions for  $(c, \sigma^2)$

$$c_{\text{GLS}} = (\mathbf{1}_n^\top \mathbf{K}_\theta^{-1} \mathbf{1}_n)^{-1} \mathbf{1}_n^\top \mathbf{K}_\theta^{-1} \underline{\mathbf{Z}}_n, \quad (6)$$

$$\sigma_{\text{GLS}}^2 = \frac{1}{n} (\underline{\mathbf{Z}}_n - c_{\text{GLS}} \mathbf{1}_n)^\top \mathbf{K}_\theta^{-1} (\underline{\mathbf{Z}}_n - c_{\text{GLS}} \mathbf{1}_n), \quad (7)$$

- **Grid-search:** For each point in a grid of lengthscales  $\{\alpha_1 \rho_0, \dots, \alpha_L \rho_0\}$ , with  $\alpha_i \geq 0$  and

$$\rho_{0,[k]} = \sqrt{d} \left( \max_{1 \leq i \leq n} x_{i,[k]} - \min_{1 \leq i \leq n} x_{i,[k]} \right), \quad 1 \leq k \leq d.$$

the likelihood is optimized w.r.t  $c$  and  $\sigma^2$  using (6) and (7). The lengthscale with the **best likelihood** value is selected.

→ Available in MATLAB/GNU Octave toolbox STK.

## Stopping condition

- Tuning the **stopping criteria** restricts/extends the search of the optimizer.
- Criteria for L-BFGS-B algorithm in GPy
  - **maxiter** : Total no. of iterations
  - **bfgs\_factr** : Threshold for the difference in functional values for two consecutive iteration
  - **gtol** : Threshold for the gradient

## Restarts & multi-starts

- Numerical noise and non-convexity causes early stopping of the optimizer
- Solution: Carrying out several optimization runs with different initialization points.
- **Restart:**
  - Runs the algorithm  $N_{\text{opt}}$  times iteratively
  - Clears the memory (Hessian approximations) each time
  - The last estimated parameters used as the new initial values.
- **Multi-start:**
  - Runs the algorithm  $N_{\text{opt}}$  times with different initialization points
  - Perturbations of the primary initial point are considered
  - The parameter with the best likelihood value over all runs is selected.

## Parameterization of the covariance function

- $\theta \in \Theta \subset \mathbb{R}_+^p$  is optimized on a transformed domain  $\Theta' \subset \mathbb{R}^p$ .
- The aim is to **reshape** the likelihood, facilitating smoother convergence.
- A **monotonic one-to-one mapping**  $\Delta : \Theta \rightarrow \Theta'$ , is applied before optimizing NLL.

$$\theta'_{opt} = \arg \min_{\theta' \in \Theta'} -\log(\mathcal{L}(\underline{Z}_n | \Delta(\theta), c)) \quad (8)$$

- $\theta'_{opt}$  is inverted using  $\Delta^{-1}$  to obtain the true parameter values.
- Examples :
  - **invsoftplus**( $s$ ):  $\log(\exp(\theta/s) - 1)$ ,  $s$  = parameter for **input standardization**.
  - **log**:  $\log(\theta)$

## 4 Numerical study

- **Optimization schemes** (combination of different improvement strategies) are compared to find the optimal strategies.
- Metric: empirical cumulative distributions (ECDFs) of differences of NLL values corresponding to a **brute-force** scheme with robust MLE.
- Data: simulated with a multi-dimensional uniformity criterion (**LHS-MDU**; Deutsch and Deutsch, 2012) from test functions in optimization literature.
- Common setup
  - GPy version 1.9.9 with L-BFGS-B optimizer
  - Constant mean-function
  - Anisotropic Matérn( $\nu = 5/2$ ) kernel
  - $\sigma_\varepsilon^2 = 0$

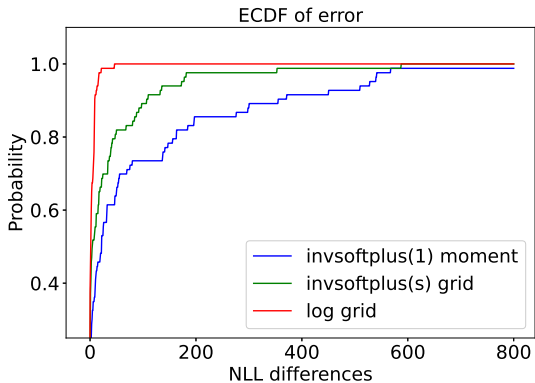


Figure 2: Best **initialization** method for each of the parameterizations.

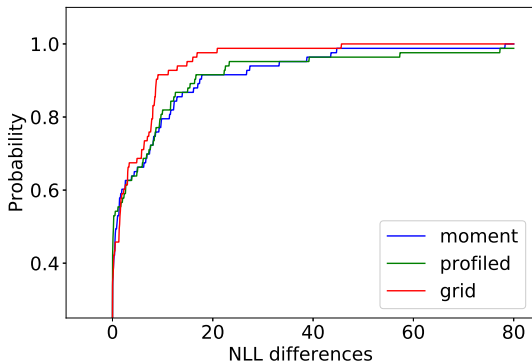


Figure 3: Different initializations for the log **parameterization**.



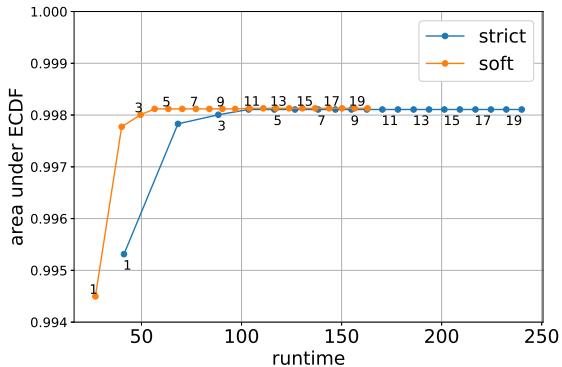


Figure 4: Area under the ECDF against run time for **restart** strategy with  $N_{\text{opt}} = 1, \dots, 20$ .

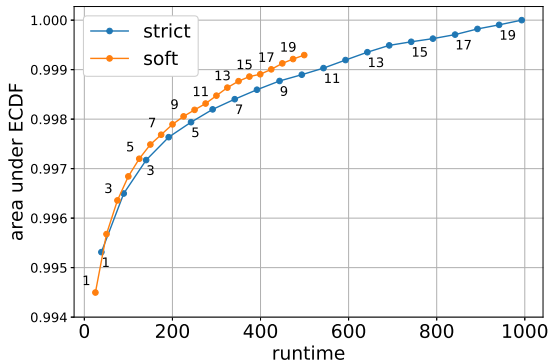


Figure 5: Area under the ECDF against run time for **multi-restart** strategy with  $N_{\text{opt}} = 1, \dots, 20$ .

## The optimal choice

- The recommended configuration consists of :
  - **log** parameterization
  - **grid-search** initialization
  - **soft** (GPY's default) stopping condition
  - small number  $N_{\text{opt}} = 5$  of **restarts**
- This **improved** setup is compared with GPY's default setup, with Leave One Out metrics (LOO) on the Borehole function (160 data points).

$$f(x) = \frac{2\pi T_u (H_u - H_l)}{\ln\left(\frac{r}{r_w}\right) \left[ 1 + \frac{2LT_u}{\ln\left(\frac{r}{r_w}\right) r_w^2 K_w} + \frac{T_u}{T_l} \right]}$$

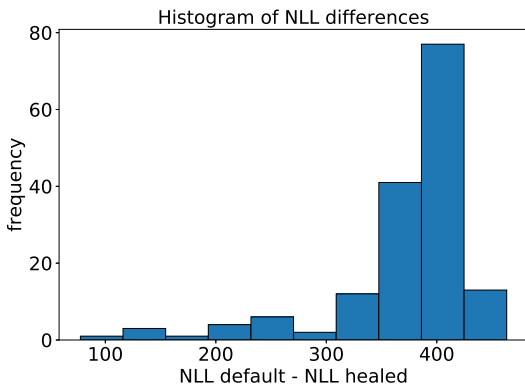


Figure 6: Distribution of the differences of NLL values.

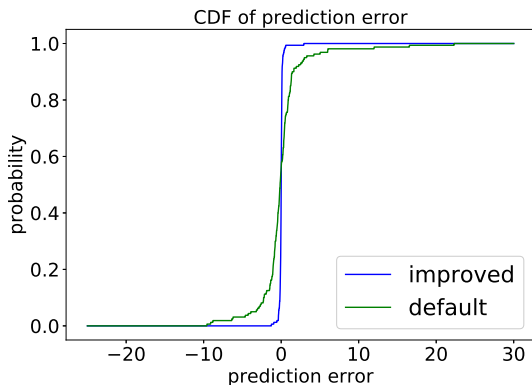


Figure 7: ECDFs of the prediction errors at points removed.

## 5 Concluding remarks

- Off-the-shelf GP implementations should be implemented carefully.
- The numerical noise on the likelihood should not be overlooked.
- Adaptive jitter cannot be considered as a do-it-all solution.
- The ML estimation can be significantly improved using some simple and effective strategies.
- This study intends to encourage practitioners to develop more robust GP implementations.

DCE Reading Group @ ATI

Numerical issues in maximum likelihood parameter estimation for  
**Gaussian process interpolation**

Subhasish Basak<sup>1</sup>, Sébastien Petit<sup>1,2</sup>, Julien Bect<sup>1</sup> & Emmanuel Vazquez<sup>1</sup>

March 10, 2021

1. Laboratoire des Signaux et Systèmes, CNRS, CentraleSupélec, Univ. Paris-Saclay
2. Safran Aircraft Engines, France