

Assessing Convergence in Gaussian Process Surrogate Model Optimization

Nicholas R. Grunloh and Herbert K. H. Lee

Abstract

Identifying convergence in numerical optimization is an ever-present, difficult, and often subjective task. The statistical framework provided by Gaussian Process surrogate model optimization provides useful secondary measures for tracking optimization progress; however the identification of convergence via these criteria has often provided only limited success and often requires a more subjective analysis. Here we use ideas originally introduced in the field of Statistical Process Control to define convergence in the context of an robust and objective convergence heuristic. The Exponentially Weighted Moving Average (EWMA) chart provides an ideal starting point for adaptation to track convergence via the EWMA convergence chart introduced here.

Keywords: Computer Model, Derivative-free Optimization, Emulator, Expected Improvement.

1 Introduction

Black-box derivative-free optimization has a wide variety of applications, especially in the realm of computer simulations [?, ?]. When dealing with computationally expensive computer models, a key question is that of convergence of the optimization. Because each function evaluation is expensive, one wants to terminate the optimization as early as possible. However for complex simulators, the response surface may be ill-behaved and optimization routines can easily become trapped in a local mode, so one needs to run the optimization sufficiently long to achieve a robust solution. In this paper, we provide an automated method for assessing convergence of Gaussian Process surrogate model optimization by bringing in elements of Statistical Process Control.

Our motivating example is a hydrology application, the Lockwood pump-and-treat problem [10], discussed in more detail in Section 3.1, wherein contamination in the ground-water near the Yellowstone River is remediated via a set of treatment wells. The goal is to minimize the cost of running the wells while ensuring that no contamination enters the river. The contamination constraint results in a complicated boundary that is unknown in advance and requires evaluation of the simulator, and

thus finding the global constrained minimum is a difficult problem where it is easy for optimization routines to, at least temporarily, get stuck in a local minimum. Without knowing the answer in advance, how do we know when to terminate the optimization routine?

The context of this paper is Gaussian Process surrogate model optimization, a statistical modeling approach to derivative-free numerical optimization that constructs a fast approximation to the expensive computer simulation using a statistical surrogate model [7]. Analysis of the surrogate model allows for efficient exploration the objective solution space. Typically a Gaussian Process (GP) surrogate model is chosen for its robustness, relative ease of computation, and its predictive framework [12]. Arising naturally from the GP predictive distribution [13, 4], the maximum Expected Improvement (EI) criterion has shown to be a valuable criterion for guiding the exploration of the objective function [17, 8]; furthermore the EI shows promise for use as a convergence criterion [cite](#).

[Literature cite](#) recommends considering the EI as a convergence criterion for surrogate model optimization; as of yet, little work has been done to describe what convergence of these algorithms actually looks like in the context of the EI criterion. However, the basic idea behind the use of EI as a convergence criterion is that convergence should occur when the surrogate model produces low expectations for discovering new optimum; that is to say, globally small EI values should be associated with convergence of the algorithm. Thus an obvious stopping rule first defines some lower EI threshold, then claims convergence upon the first instance of the maximum EI value falling below this threshold [example \[3\]](#), [cite](#). This use of EI as a convergence criterion is analogous to the standard approach in the numerical optimization literature of setting a threshold for some value that is monitored and declaring convergence when the threshold is obtained (e.g., the vanishing step sizes of a Newton-Raphson algorithm). However, applying this same threshold strategy to the convergence of surrogate model optimization has not yet been adequately justified. In fact, this use of EI ignores the nature of the EI criterion as a random variable, and oversimplifies the stochastic nature of convergence in this setting. Thus it is no surprise that this treatment of the EI criterion can result in an inconsistent stopping rule as demonstrated in Figure (1).



Figure 1: Three ELAI series plotted alongside an example convergence threshold value shown as a dashed line at -10.

Because EI is strictly positive but decreasingly small, we find it more productive to work on the log scale, using a lognormal approximation to the improvement distribution to generate a more appropriate convergence criterion, as described in more detail in Section 3.2. Figure (1) represents three series of the Expected Lognormal Approximation to the Improvement (ELAI) convergence criterion values from three different optimization problems that will be demonstrated later in this paper, where it will be shown that convergence is established near the end of each of these series. These three series demonstrate various ELAI convergence behaviors, and illustrate the difficulty in assessing convergence. In the left-most panel, optimization of the Rosenbrock test function results in a well-behaved series of ELAI values, demonstrating a case in which the simple threshold stopping rule can accurately identify convergence. However the center panel (the Lockwood problem) demonstrates a failure of the threshold stopping rule, as this ELAI series contains much more variance, and thus small ELAI values are observed quite regularly. In the Lockwood example a simple threshold stopping rule could falsely claim convergence within the first 50 iterations of the algorithm. The large variability in ELAI with occasional large values indicates that the optimization routine sometimes briefly settles into a local minimum but is still exploring and is not yet convinced that it has found a global minimum. This optimization run appears to have converged only after the larger ELAI values stop appearing and the variability has decreased. Thus one might ask if a decrease in variability, or small variability, is a necessary condition for convergence. The right-most panel (the Rastrigin test function) shows a case where convergence occurs by meeting the threshold level, but where variability has increased, demonstrating that a decrease in variability is not a necessary condition.

As the Improvement function is itself a random variable, attempting to set a lower threshold bound on the EI, without consideration of the underlying EI distribution over time, over-simplifies the dynamics of convergence in this setting. Instead, we propose taking the perspective of Statistical Process Control (SPC), where a stochastic series is monitored for consistency of the distribution of the most recently observed values. In the next section, we review the statistical surrogate model approach and the use of EI for optimization. In Section 3, we discuss our inspiration from SPC and how we construct our convergence chart. Section ?? provides synthetic and real examples, and then we provide some conclusions in the final section.

2 Gaussian Process Surrogate Model Optimization

The primary motivation for the use of surrogate modeling in optimization is to manage a computationally challenging objective function with the use of a fast and relatively simple functional working model (i.e. the surrogate model) of the problem function. The surrogate model serves as a tool for using function evaluations, to infer the expected behavior of the objective function and thus determine where further optima may exist with minimal evaluation the objective function itself. Surrogate modeling is therefore useful for optimizing large computer simulations experiments, where each function evaluation may consume considerable computational resources. GP surrogate modeling considers the objective function $f \sim \text{GP}(m(\mathbf{x}), C(\mathbf{x}, \mathbf{x}'))$, and thus aims to minimize the number of actual function evaluations by using the GP predictive surface to interpolate between function evaluations. Inference on such GP models seems to strike an equitable balance between relative ease of computation and efficient learning of the true objective function behavior. As a result GP surrogate models have become a standard tool for analysis of computer emulation experiments [11, 12].

As with most popular optimization strategies for optimizing functions over a real domain, GP surrogate optimization can make use of the assumption of some cohesive smoothness of the objective function. This smoothness relates points close in space, with similar expectations for the objective value, providing the primary mechanism by which optimization may proceed in most numerical algorithms. GPs can directly express this assumption of smoothness with the a’priori choice of an appropriate covariance function which smoothly relates points through their relative positions in

the domain. The choice of a covariance function to represent a smooth objective function is the typical modeling choice for optimization in computer emulation [12], although theoretically GPs may accommodate many common covariance structures [1, 15].

In many cases the assumption of a globally smooth f with a homogeneous uncertainty structure can provide an effective and parsimonious model. However for the sake of providing a flexible surrogate model, it is desirable to have the ability to loosen these restrictions in cases when f may have inherently sharp boundaries, or numerical simulators have variable stability in portions of the domain. Gramacy and Lee [6] use the idea of allowing subpopulations of flexibility via a treed partitioning of the domain, fitting stationary GP surfaces to separately stationary portions of f . The domain is recursively sub-partitioned and separate hierarchically linked GP models are fit within each sub-partition. The partitioning scheme is fit via a reversible jump MCMC algorithm, jumping between models with differing partitioning schemes. By partitioning the domain in this way it allows parsimonious surrogate models in simple objective function cases and quite flexible surrogate models when the the objective function display's complex behavior. For further explanation of partitioned Gaussian process models as well as notes on implementing such models in R, see the R package `tgpp` [5, 7].

2.1 Expected Improvement

The EI criterion is fundamentally based on the improvement criterion [cite](#) which evaluates how possible it may be to encounter new minima at a given location based on the predictive surrogate model. The improvement function takes the following form,

$$I(\mathbf{x}) = \max \left\{ (f_{min} - f(\mathbf{x})), 0 \right\}. \quad (1)$$

By considering the expectation of $I(\mathbf{x})$, candidate locations are not only rewarded for having a low predictive mean, but the $\mathbb{E}[I(\mathbf{x})]$ also rewards poorly explored locations due to the high uncertainty of $I(\mathbf{x})$ in these places. Notice that by definition the $I(\mathbf{x})$ function is always non-negative, however the GP posterior predictive $\mathbb{E}[I(\mathbf{x})]$ is a strictly positive criterion. Considering the MCMC inferential setting of our GP surrogate model, the EI criterion can be quickly computed by using pos-

terior predictive $I(\mathbf{x})$ samples at given candidate locations to empirically approximate the $\mathbb{E}[I(\mathbf{x})]$ calculation.

2.2 Optimization Procedure

The idea for optimization, in this context, is to only evaluate the objective function at locations that have a good chance of providing a new minimum. Optimization begins by initially collecting a set, \mathbf{X} , of locations to evaluate the true function, f , to gather a basic impression of f . A statistical surrogate model is then fitted with $f(\mathbf{X})$ as observations of the true function. Using the surrogate model, a set of candidate points, $\tilde{\mathbf{X}}$, are selected from the domain and the EI cri-

terion is calculated among these points. The candidate point that has the highest EI is then chosen as the best candidate for a new minimum and thus, it is added to \mathbf{X} . The objective function is evaluated at this new location and the surrogate model is refit based on the updated $f(\mathbf{X})$. The optimization procedure carries on in this way until convergence.

Figure 2: Optimization Procedure

- 1) Collect an initial set, \mathbf{X} .
- 2) Compute $f(\mathbf{X})$.
- 3) Fit surrogate based on evaluations of f .
- 4) Collect a candidate set, $\tilde{\mathbf{X}}$.
- 5) Compute EI among $\tilde{\mathbf{X}}$
- 6) Add $\operatorname{argmax}_{\tilde{\mathbf{x}}_i} \mathbb{E}[I(\tilde{\mathbf{x}}_i)]$ to \mathbf{X} .
- 7) Check convergence.
- 8) If converged exit. Otherwise go to 2).

3 EWMA Convergence Chart

3.1 Statistical Process Control

In Shewhart’s seminal 1931 book [14] on the topic of control in manufacturing, Shewhart explains that a phenomenon is said to be in control when, “through the use of past experience, we can predict, at least within limits, how the phenomenon may be expected to vary in the future.” This notion provides an instructive framework for thinking about convergence because it offers a natural way to consider the distributional characteristics of the EI as a proper random variable. In its most simplified form, SPC considers an approximation of a statistic’s sampling distribution as repeated sampling occurs

in time. For example, the \bar{x} -chart tracks the mean of repeated samples (all of size n) through time so as to expect the arrival of each subsequent mean in accordance with the typical sampling distribution for the mean, $\bar{x}_j \sim N\left(\mu, \frac{\sigma^2}{n}\right)$. Shewhart expresses his idea of control, in this case, as the expected behavior of random observations from this sampling distribution. By considering confidence intervals on this sampling distribution we can easily draw explicit boundaries (i.e. control limits) to identify which samples are in control, and which are not. Observations violating our expectations (i.e. observations that fall outside of the confidence interval/beyond the control limits) indicate an out-of-control state. Since neither μ nor σ^2 are typically known, it is of primary importance to use the data carefully to form accurate approximations of these values, thus establishing a standard for control. Furthermore, this logic relies upon the typical asymptotic results of the central limit theorem (CLT), and special care should always be taken to satisfy its requirements.

3.2 Expected Lognormal Approximation to the Improvement (ELAI)

For the sake of obtaining a robust convergence criterion to track via SPC, it is important to carefully consider properties of the improvement distributions which generate the EI values. The improvement criterion is strictly positive but decreasingly small, thus the improvement distribution is often strongly right skewed, and the EI is often far from normal. Additionally, this right skew becomes exaggerated as convergence approaches, due to the decreasing trend in the EI criterion. Together these characteristics of the improvement distribution give the EI criterion inconsistent behavior for tracking convergence via a typical \bar{x} -chart.

These issues naturally suggest releasing the bound at 0 by modeling transformations of the improvement, rather than directly considering the improvement distribution on its own. One of the simplest of the many possible helpful transformations in this case could simply consider the log of the improvement distribution. However due to the MCMC sample-based implementation of the Gaussian Process, and the desire for a large number of samples from the improvement distribution, it is not uncommon to obtain at least one sample, that in double precision, is computationally indistinguishable from 0. Thus simply taking the log of the improvement samples can be computationally undefined, particularly as convergence approaches. Despite this numerical inconvenience, the distribution of the improvement samples is often very well approximated by the Lognormal distribution.

Rather than simply taking the log of the improvement samples, to determine the more robust statement that $\mathbb{E}[\log I] \approx N\left(\mu, \frac{\sigma^2}{n}\right)$, it is computationally useful to consider the following approximate model-based perspective. Recall that if a random variable $X \sim \text{Log-N}(\omega, \nu)$, then another random variable $Y = \log(X)$ is distributed $Y \sim N(\omega, \nu)$. Furthermore, if m and v are, respectively, the mean and variance of a lognormal sample, then the mean, ω , and variance, ν , of the associated normal distribution are given by the following relation.

$$\omega = \log\left(\frac{m^2}{\sqrt{v + m^2}}\right) \quad \nu = \log\left(1 + \frac{v}{m^2}\right). \quad (2)$$

Using this relation we do not need to transform any of the improvement samples. We compute the empirical mean and variance of the unaltered, approximately lognormal, improvement samples, then use relation (2) to directly compute ω as the Expectation under the Lognormal Approximation to the Improvement (ELAI). The ELAI convergence criterion is a useful convergence criterion in this case because of the reduced right skew of the log of the improvement distribution, and the ELAI convergence criterion serves as a computationally robust approximation of the $\mathbb{E}[\log I]$ under reasonable log-normality of the improvements. Furthermore, both the $\mathbb{E}[\log I]$ and ELAI convergence criterion are distributed approximately normally in repeated sampling. This construction allows for more consistent and accurate use of the fundamental theory on which our SPC perspective requires.

3.3 Exponentially Weighted Moving Average

The EWMA control chart is based upon Shewhart's original notion of control, however the EWMA control chart views the typical repeated sampling process in the context of moving average smoothing of series data. Since preconvergence ELAI convergence criterion evaluations do not necessarily share consistency among all observed values, a weighted series perspective was chosen to follow the moving average of the most recent ELAI observations while still smoothing with some memory of older evaluations. EWMA achieves this robust smoothing behavior, relative to shifting means, by assigning exponentially decreasing weights to successive points in a rolling average among all of the points of the series, thus the EWMA emphasizes recent observations and shifts the focus of the moving average to the most recent information while still providing shrinkage towards the global mean of the series.

If Y_i is the current ELAI value, and Z_i is the EWMA statistic associated with this current value, then the initial value Z_0 is set to Y_0 and for $i > 0$ the EWMA statistic is expressed as,

$$Z_i = \lambda Y_i + (1 - \lambda)Z_{i-1}. \quad (3)$$

Above, λ is a smoothing parameter that defines the weight (i.e. $0 < \lambda \leq 1$) assigned to the most recent observation, Y_i . The recursive expression of the statistic ensures that all subsequent weights geometrically decrease as they move back through the series.

Typical values of λ can range from $0.1 \leq \lambda \leq 0.3$, with a default value of λ around 0.2, as described by Box et al. [2]. Large values of λ assign more weight to recent observations in the series, allowing for a more flexible fit for unstable series. However, the choice of a large λ may over-fit the Z_i to noise in the Y_i . It is thus desirable to choose the smallest λ which still provides good forecasts of future observations in the series. Box et al. [2, p. 87] explains

how to choose an optimal value for λ by choosing the $\hat{\lambda}$ which minimizes the sum of squared forecasting deviations (S_λ) for each new observation. Through this analysis of S_λ , as seen in Figure (3), is it evident that EWMA charts can be very robust to reasonable choices of λ , due to the small first and second derivatives of S_λ for a large range of sub-optimal choices of λ around $\hat{\lambda}$. In fact, Figure (3) shows that for $\lambda \in [0.2, 0.6]$, S_λ stays within 10% of its the minimum possible value.

It is interesting to note that for the example series used in Figure (3), the optimal $\hat{\lambda}$ is found to be about 0.4, which obviously falls outside of the standard values for λ mentioned above. The typical use of EWMA in

SPC begins with the premise of a relatively stable series and attempts to identify new out-of-control observations which would indicate some change in the data generating process. The use of EWMA to identify convergence is markedly different from the typical SPC usage. For identifying convergence

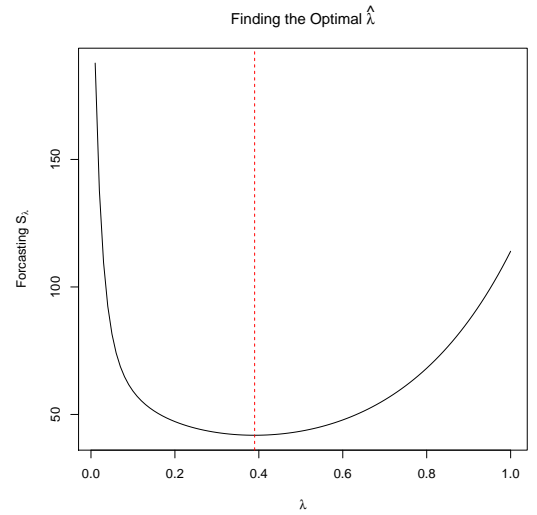


Figure 3: S_λ as calculated over a fine grid of possible λ values for ELAI values derived from optimization of the Rastrigin test function. The optimal forecasting $\hat{\lambda}$ is shown by the vertical the dashed line.

the series begins in an out-of-control state and we wish to identify when the series has fallen into control (i.e. convergence). As a result, ELAI values for tracking convergence are inherently less stable than typical SPC applications of EWMA. Often the optimal $\hat{\lambda}$, for identifying convergence, may fall above the recommended upper limit for λ . Thus, in this context it is useful for to borrow some of the machinery used in SPC, but the ultimate analysis of the data should not explicitly be considered SPC, as the perspective of the data does not fully align with that of SPC.

For identifying convergence it is of primary importance to define the control limits on the EWMA statistic in this setting. As in the simplified \bar{x} -chart, defining the control limits in the EWMA setting amounts to considering an interval on the sampling distribution of interest. In the EWMA case we are interested in the sampling distribution of the Z_i . Assuming that the Y_i are *i.i.d.* then Lucas and Saccucci [9] show that we can write $\sigma_{Z_i}^2$ in terms of σ_Y^2 .

$$\sigma_{Z_i}^2 = \sigma_Y^2 \left(\frac{\lambda}{2 - \lambda} \right) [1 - (1 - \lambda)^{2i}] \quad (4)$$

Thus if the $Y_i \stackrel{i.i.d.}{\sim} N\left(\mu, \frac{\sigma^2}{n}\right)$ the sampling distribution for Z_i is $Z_i \sim N(\mu, \sigma_{Z_i}^2)$. Furthermore by choose a confidence level through choice of a constant c , the control limits based on this sampling distribution are seen in Eq. (5).

$$\begin{aligned} \text{CL}_i &= \mu \pm c \sigma_{Z_i} \\ &= \mu \pm c \frac{\sigma}{\sqrt{n}} \sqrt{\left(\frac{\lambda}{2 - \lambda} \right) [1 - (1 - \lambda)^{2i}]} \end{aligned} \quad (5)$$

Notice that since $\sigma_{Z_i}^2$ has a dependence on i , the control limits do as well. Looking back through the series brings us away from the focus of the moving average, at i , and thus the control limits widen until the limiting case as, $i \rightarrow \infty$, the control limits are defined by $\mu \pm c \sqrt{\frac{\lambda \sigma^2}{(2 - \lambda)n}}$.

At first glance it is not clear that the Y_i are in fact *i.i.d.* Indeed the early iterations of the convergence processes seen in Figure (1) certainly do not display *i.i.d.* Y_i . However as the series approaches convergence, the Y_i eventually do enter a state of control see Figure (4). For these controlled Y_i an *i.i.d.* assumption is very reasonable. The realization of such a controlled region

of the series defines the notion of consistency which is part of what allows for the identification of convergence here.

3.4 The Control Window

The final structural feature of the EWMA convergence chart for identifying convergence is the so called *control window*. The control window contains a fixed number, w , of the most recently observed Y_i . Only information from the w points currently residing inside the control window is used to calculate the control limits, but the EWMA statistic is still computed for all Y_i values. Initially, the convergence algorithm is allowed to fill the control window, by collecting an initial set of w observations of the Y_i . As new observations arrive, the oldest Y_i value is removed from the control window, thus allowing for the inclusion of a new Y_i .

The purpose of the control window is two fold. Firstly it serves to dichotomizes the series for evaluating subsets of the Y_i for distributional consistency. Secondly it offers a structural way for basing the standard for consistency (i.e. the control limits) only on the most recent and relevant information in the series.

The size, w , of the control window is an important parameter for correctly identifying convergence. The size of the control window, w , Because w may vary from problem to problem it is ultimately left as a tuning parameter of the system. Choosing the correct value of w presents an interesting decision problem since underestimating the size of the control window may lead to premature identification of convergence, however if w is too large, we compute unnecessary objective function evaluations. As a general trend, harder optimization problems require larger values of w since the EI criterion follows a less structured decreasing pattern as new modes are discovered at irregular patterns.

3.5 Identifying Convergence

4 Examples

5 Conclusion

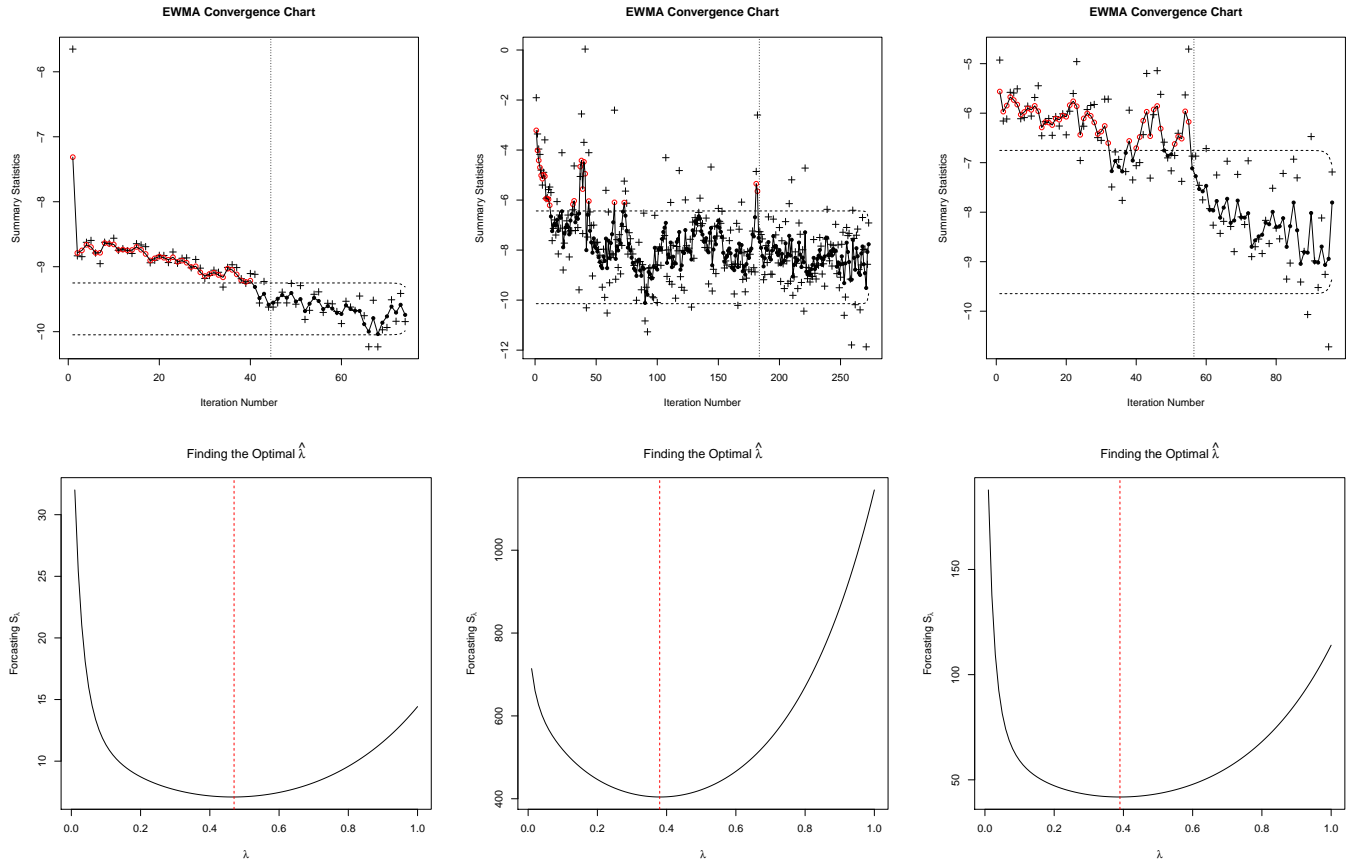


Figure 4: *left* : Rosenbrock | *center* : Lockwood | *right* : Rastrigin

References

- [1] Petter Abrahamsen. *A review of Gaussian random fields and correlation functions*. Norsk Regnesentral/Norwegian Computing Center, 1997.
- [2] George E. P. Box, Alberto Luceño, and María Del Carmen Paniagua-Quiñones. *Statistical Control by Monitoring and Adjustment*. Wiley, New York, NY, 1997.
- [3] Sanket Sanjay Diwale, Ioannis Lymperopoulos, and Colin Jones. Optimization of an airborne wind energy system using constrained gaussian processes. In *IEEE Multi-Conference on Systems and Control*, 2014.
- [4] David Ginsbourger, Rodolphe Le Riche, Laurent Carraro, et al. A multi-points criterion for deterministic parallel global optimization based on gaussian processes. In *Journal of Global Optimization, in revision*. Citeseer, 2009.
- [5] Robert B. Gramacy. tgp: an r package for bayesian nonstationary, semiparametric nonlinear regression and design by treed gaussian process models. *Journal of Statistical Software*, 19(9):6, 2007.
- [6] Robert B. Gramacy and Herbert H. K. Lee. Bayesian treed gaussian process models with an application to computer modeling. *Journal of the American Statistical Association*, 103(483), 2008.
- [7] Robert B. Gramacy and Matthew Taddy. Categorical inputs, sensitivity analysis, optimization and importance tempering with tgp version 2, an r package for treed gaussian process models. *Journal of Statistical Software*, 33(i06), 2012.
- [8] Donald R. Jones, Matthias Schonlau, and William J. Welch. Efficient global optimization of expensive black-box functions. *Journal of Global optimization*, 13(4):455–492, 1998.
- [9] James M. Lucas and Michael S. Saccucci. Exponentially weighted moving average control schemes: properties and enhancements. *Technometrics*, 32(1):1–12, 1990.
- [10] Shawn L. Matott, Kenny Leung, and Junyoung Sim. Application of matlab and python optimizers to two case studies involving groundwater flow and contaminant transport modeling. *Computers & Geosciences*, 37(11):1894–1899, 2011.
- [11] Jerome Sacks, William J Welch, Toby J Mitchell, and Henry P Wynn. Design and analysis of computer experiments. *Statistical science*, pages 409–423, 1989.
- [12] Thomas J. Santner, Brian J. Williams, and William Notz. *The design and analysis of computer experiments*. Springer, 2003.
- [13] Matthias Schonlau, William J. Welch, and Donald R. Jones. Global versus local search in constrained optimization of computer models. *Lecture Notes-Monograph Series*, pages 11–25, 1998.
- [14] Walter A. Shewhart. *Economic control of quality of manufactured product*, volume 509. ASQ Quality Press, 1931.

- [15] Michael L Stein. *Interpolation of spatial data: some theory for kriging*. Springer, 1999.
- [16] Yan Su, Lianjie Shu, and Kwok-Leung Tsui. Adaptive ewma procedures for monitoring processes subject to linear drifts. *Computational statistics & data analysis*, 55(10):2819–2829, 2011.
- [17] Matthew A. Taddy, Herbert H. K. Lee, Genetha A. Gray, and Joshua D. Griffin. Bayesian guided pattern search for robust local optimization. *Technometrics*, 51(4):389–401, 2009.