
Determining Convergence for Bayesian Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

Bayesian optimization routines may have theoretical convergence results, but determining whether a run has converged in practice can be a subjective task. This paper provides a framework inspired by statistical process control for monitoring an optimization run for convergence. An Exponentially Weighted Moving Average chart is adapted for automated convergence analysis.

Keywords: Derivative-free Optimization, Computer Simulation, Emulator, Expected Improvement, EWMA

1 Introduction

Bayesian optimization aims to find a global optimum of a complex function that may not be analytically tractable, and where derivative information may not be readily available (Mockus, 1989; Brochu et al., 2010). A common application is for computer simulation experiments Gramacy (2020). Because each function evaluation may be expensive, one wants to terminate the optimization algorithm 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. So far there has been little work on assessing convergence for Bayesian optimization. In this paper, we provide an automated method for determining convergence of surrogate model-based optimization by bringing in elements of statistical process control.

Among the wide variety of Bayesian optimization approaches, we focus on those that are based on a statistical surrogate model, such as a Gaussian process (Santner et al., 2003). We further focus on approaches based on Expected Improvement (EI) (Schonlau et al., 1998), although our methods are generalizable for other acquisition functions.

There have been a few hints in the literature that monitoring EI directly could be used to assess convergence (Jones et al., 1998). Taddy et al. (2009) considers the use of the improvement distribution for identifying global convergence. The basic idea is that convergence should occur when the surrogate model produces low expectations for discovering a new optimum; that is to say, globally small EI values should be associated with convergence of the algorithm. Thus a simplistic stopping rule might first define some lower EI threshold, then claim convergence upon the first instance of an EI value falling below this threshold, as seen in Diwale et al. (2015). This use of EI as a convergence criterion is analogous to other standard convergence identification methods in numerical optimization (e.g., the vanishing step sizes of a Newton-Raphson algorithm). However, applying this same threshold strategy to the convergence of Bayesian 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).

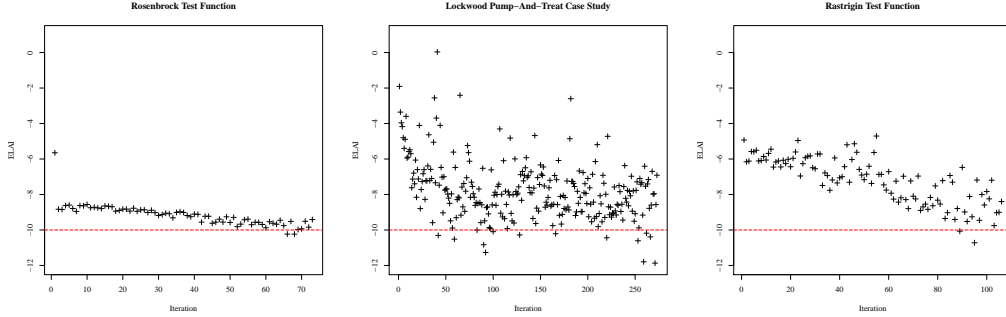


Figure 1: Three Expected Log-normal Approximation to the Improvement series (more details in Section 4) plotted alongside an example convergence threshold value shown as a dashed line at -10.

Because EI is strictly positive but decreasingly small, we find it more productive to work on the log scale, using a log-normal approximation to the improvement distribution to generate a more appropriate convergence criterion, as described in Section 3.2. Figure (1) represents three series of the Expected Log-normal Approximation to the Improvement (ELAI) values from three different optimization problems. We will demonstrate later in this paper that convergence is established near the end of each of these series. These three series demonstrate the kind of diversity observed among 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 described in Section 4.3) 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 values 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.

Since the Improvement function is itself random, attempting to set a lower threshold bound on the EI, without consideration of the underlying EI distribution through 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 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 4 provides synthetic and real examples, and then we provide some conclusions in the final section.

2 Bayesian Optimization via Expected Improvement

Bayesian optimization attempts to solve problems of the form

$$x^* = \operatorname{argmin}_{x \in \mathcal{X}} f(x),$$

where f is an objective function (often not available in analytical form) and $x \in \mathcal{X} \subset \mathbb{R}^d$. \mathcal{X} may be defined via constraints. Without loss of generality, we frame all optimizations as minimizations in this paper, as maximization can be recovered by minimizing the negative of the function. Bayesian optimization proceeds by iteratively developing a statistical surrogate model of the objective function f , and using predictions from the statistical surrogate to choose the next point to evaluate based

on some criterion. A common choice of surrogate model is the Gaussian process (GP) [Gramacy \(2020\)](#); [Pourmohamad and Lee \(2021\)](#), as it combines flexibility with smoothness.

In many cases the assumption of a globally smooth f with a homogeneous uncertainty structure can provide an effective and parsimonious model. However, in other problems, f may have sharp boundaries, f may show different levels of smoothness across its domain, or numerical simulators may have variable stability in portions of the domain. In this paper, we use treed Gaussian processes [Gramacy and Lee \(2008\)](#), a generalization of a standard GP that uses treed partitioning of the domain, fitting separate hierarchically-linked stationary GP surfaces to separate portions of f via a reversible jump MCMC algorithm and averaging over the full parameter space to provide smooth predictions except where the data call for a discontinuous prediction. We use the R package `tgpp` ([Gramacy, 2007](#); [Gramacy and Taddy, 2010](#)). While the treed GPs provide additional modeling flexibility, we emphasize that the approach of this paper can be applied to standard GPs as well as any surrogate model that provides both predictions and predictive uncertainty.

2.1 Expected Improvement

Bayesian optimization requires an acquisition function that guides the choice of a new function evaluation at each iteration. There are a wide variety of suggestions for acquisition functions. A large family of options is based on Expected Improvement. The EI criterion predicts how likely a new minimum is to be observed, at new locations of the domain, based upon the predictive distribution of the surrogate model. EI is built upon the improvement function ([Jones et al., 1998](#)):

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

where f_{\min} is the smallest function value observed so far. EI is the expectation of the improvement function with respect to the posterior predictive distribution of the surrogate model, $\mathbb{E}[I(\mathbf{x})]$. EI rewards candidates both for having a low predictive mean, as well as high uncertainty (where the function has not been sufficiently explored), thus balancing global exploration and local exploitation. By definition the improvement function is always non-negative and the posterior predictive $\mathbb{E}[I(\mathbf{x})]$ is strictly positive. The EI criterion is available in closed form for a stationary GP. For other models the EI criterion can be quickly estimated using Monte Carlo posterior predictive samples at given candidate locations.

2.2 Optimization Procedure

Optimization can be viewed as a sequential design process, where locations are selected for evaluation on the basis of how likely they are to decrease the objective function, i.e., based on the EI. Optimization begins by initially collecting a set, \mathbf{X} , of locations to evaluate the true function, f , to get an initial fit of the statistical surrogate model, using $f(\mathbf{X})$ as observations of the true function. Based on the surrogate model, a set of candidate points, $\tilde{\mathbf{X}}$, are selected from the domain and the EI criterion 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 using the updated $f(\mathbf{X})$. The optimization procedure carries on in this way until convergence. The key contribution of this paper is an automated method for checking convergence, which we develop in the next section.

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 $\arg\max_{\tilde{\mathbf{x}}_i} \mathbb{E}[I(\tilde{\mathbf{x}}_i)]$ 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 book (Shewhart, 1931) 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. Thus Shewhart can express his idea of control as the expected behavior of random observations from this sampling distribution. For example, an \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 known or estimated sampling distribution for the mean, $\bar{x}_j \sim N\left(\mu, \frac{\sigma^2}{n}\right)$. By considering confidence intervals on this sampling distribution we can draw explicit boundaries (i.e., control limits) to identify when the process is in control and when it is not. Observations violating our expectations (falling outside of the control limits) indicate an out-of-control state. Since neither μ nor σ^2 are typically known, it is common to collect an initial set of data from which point estimates of μ and σ^2 may establish an initial standard for control that is further refined as the process proceeds. This logic relies upon the typical asymptotic results of the central limit theorem (CLT), and care should be taken to verify the relevant assumptions required.

It is important to note that we are not performing traditional SPC in this context, as the EI criterion will be stochastically decreasing as an optimization routine proceeds. Only when convergence is reached will the EI series look approximately like an in-control process. Thus our perspective is completely reversed from the traditional SPC approach—we start with a process that is out of control, and we determine convergence when the process stabilizes and becomes locally in control. An alternative way to think about our approach is to consider performing SPC backwards in time on our EI series. Starting from the most recent EI observations and looking back, we declare convergence if the process starts in control and then becomes out of control. This pattern generally appears only when the optimization has progressed and reached a local mode without other prospects for a global mode. If the optimization were still proceeding, then the EI would still be decreasing and the final section would not appear in control.

3.2 Expected Log-normal 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, in which case, the EI is far from normal. Additionally, this right skew becomes exaggerated as convergence approaches, due to the decreasing trend in the EI criterion. These issues naturally suggest 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 would consider the log of the improvement distribution. However due to the Monte Carlo sample-based implementation of the Gaussian process, it is not uncommon to obtain at least one sample that is computationally indistinguishable from zero in double precision. Thus simply taking the log of the improvement samples can result in numerical failure, particularly as convergence approaches, even though the quantities are theoretically strictly positive. Despite this numerical inconvenience, the distribution of the improvement samples is often very well approximated by the log-normal distribution.

We avoid the numerical issues by using a model-based approximation. With the desire to model $\mathbb{E}[\log I] \sim N\left(\mu, \frac{\sigma^2}{n}\right)$, we switch to a log-normal perspective. Recall that if a random variable $X \sim \text{Log-}N(\theta, \phi)$, then another random variable $Y = \log(X)$ is distributed $Y \sim N(\theta, \phi)$. Furthermore, if ω and ψ are, respectively, the mean and variance of a log-normal sample, then the mean, θ , and variance, ϕ , of the associated normal distribution are given by the following relation.

$$\theta = \log\left(\frac{\omega^2}{\sqrt{\psi + \omega^2}}\right) \quad \phi = \log\left(1 + \frac{\psi}{\omega^2}\right). \quad (2)$$

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 log-normal, improvement samples, then use relation (2) to directly compute ω as the Expectation under the Log-normal Approximation to the Improvement (ELAI). The ELAI value is useful for assessing convergence because of the reduced right skew of the log of the posterior predictive improvement distribution. Additionally, the ELAI 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 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 depends.

3.3 Exponentially Weighted Moving Average

The Exponentially Weighted Moving Average (EWMA) control chart (Lucas and Saccucci, 1990; Scrucca, 2004) elaborates on Shewhart's original notion of control by viewing the repeated sampling process in the context of a moving average smoothing of series data. Pre-convergence ELAI evaluations tend to be variable and overall decreasing, and so do not necessarily share distributional consistency among all observed values. Thus 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 can emphasize recent observations and shift 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 $i \in \{1, 2, 3, \dots\}$ the EWMA statistic is expressed as $Z_i = \lambda Y_i + (1 - \lambda)Z_{i-1}$. Here $\lambda \in (0, 1]$ is a smoothing parameter that defines the weight assigned to the most recent observation. The recursive expression of the statistic ensures that all subsequent weights geometrically decrease.

Box et al. (1997) describes a method for computing optimal choices of λ by minimizing the sum of squared forecasting deviations (S_λ). Through this analysis of S_λ , as seen in Figure (3), it is 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.

Identifying convergence in this setting now requires the computation of control limits on the EWMA statistic. As in the simplified \bar{x} -chart, defining the control limits for 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 (1990) show that we can write $\sigma_{Z_i}^2$ in terms of σ_Y^2 .

$$\sigma_{Z_i}^2 = \sigma_Y^2 \left(\frac{\lambda}{2 - \lambda} \right) [1 - (1 - \lambda)^{2i}] \quad (3)$$

Thus if $Y_i \stackrel{i.i.d.}{\sim} N\left(\mu, \frac{\sigma^2}{n}\right)$, then the sampling distribution for Z_i is $Z_i \sim N\left(\mu, \sigma_{Z_i}^2\right)$. Furthermore by choosing a confidence level through choice of a constant c , the control limits based on this sampling distribution are seen in Eq. (4).

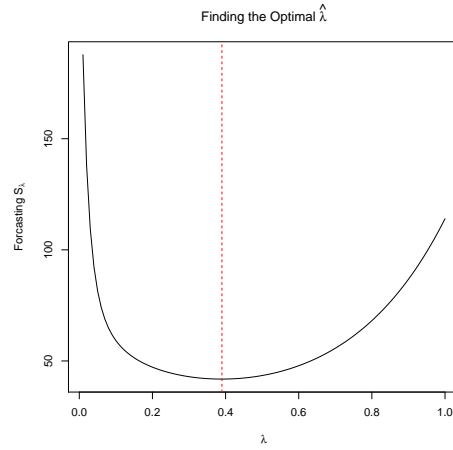


Figure 3: S_λ as calculated for ELAI values derived under the Rastrigin test function. $\hat{\lambda}$ is shown by the vertical dashed line.

$$CL_i = \mu \pm c\sigma_{Z_i} = \mu \pm c \frac{\sigma}{\sqrt{n}} \sqrt{\left(\frac{\lambda}{2-\lambda}\right) [1 - (1-\lambda)^{2i}]} \quad (4)$$

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, and thus the control limits widen until the limiting case, $i \rightarrow \infty$, where the control limits are defined by $\mu \pm c \sqrt{\frac{\lambda\sigma^2}{(2-\lambda)n}}$.

Our aim in applying the EWMA framework in this context is to recognize the fundamental notion of control that EWMA enforces in the newly arriving EI values, as optimization proceeds. Convergence often arises as a subtle shift of the EI distribution into place. In this context a more traditional \bar{x} chart will often overlook convergence as a subtle random fluctuation, when in fact it is often this subtle signal that we aim to pick-up. EWMA is among the better techniques for recognizing such subtly shifting means [Aerne et al. \(1991\)](#); [Zou et al. \(2009\)](#), while maintaining the capability to detect abrupt shifts in mean. As convergence approaches the newly arriving Y_i begin to fit into the *i.i.d.* EWMA framework and the Z_i increasingly begin to fall within the EWMA control limits. EWMA's recognition of such a controlled region in the newly arriving ELAI values, indicates the notion of distributional consistency that is necessary for defining convergence for stochastic measures of convergence, such as EI.

3.4 The Control Window

The final structural feature of the EWMA convergence chart for identifying convergence is the *control window*, which 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. To assess convergence, the EWMA statistic is 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. First, it serves to dichotomize the series for evaluating subsets of the Y_i for distributional consistency. Second, 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 of the control window, w , may vary from problem to problem based on the difficulty of optimization in each case. A reasonable way of choosing w is to consider the number of observations necessary to establish a standard of control. In this setting w is a kind of sample size, and as such the choice of w will naturally increase as the variability in the ELAI series increases. Just as in other sample size calculations, the choice of an optimal w must consider the cost of poor inference (premature identification of convergence) associated with underestimating w , against the cost of over sampling (continuing to sample after convergence has occurred) associated with overestimating w . Providing a default choice of w is somewhat arbitrary without careful analysis of the particulars of the objective function behavior and the costs of each successive objective function evaluation.

For the purpose of exploring the behavior of w in examples presented here, we use the following procedure for educating the choice of w . We hand tune w for two informative known example functions (i.e., Rosenbrock and Rastrigin). From exploration of w in known examples, it is clear that w needs to increase directly with ELAI variance. Furthermore, if one considers the form of sample size calculations based on classical power analysis, sample size increases directly proportional with the sample variance. Thus we linearly extrapolate the choice of w for the Lockwood case study based on a default starting value of 30 (based on sampling conventions) with a slope term structured to make use of the proportionality of w with the observed ELAI variance (\hat{v}) so that $\hat{w} = \frac{\Delta w}{\Delta V(\text{ELAI})} \hat{v} + 30$.

3.5 Identifying Convergence

In identifying convergence, we not only desire that the ELAI series reaches a state of control, but we desire that the ELAI series demonstrates a move from a state of pre-convergence to a consistent state of convergence. To recognize the move into convergence we combine the notion of the control window with the EWMA framework to construct the so called, *EWMA Convergence Chart*. Since

we expect EI values to decrease upon convergence, the primary recognition of convergence is that new ELAI values demonstrate values that are consistently lower than initial pre-converged values.

First, we require that all exponentially weighted Z_i values inside the control window fall within the control limits. This ensures that the most recent ELAI values demonstrate distributional consistency within the bounds of the control window. Second, since we wish to indicate a move from the initial pre-converged state of the system, we require at least one point beyond the initial control window to fall outside the defined EWMA control limits. This second rule suggests that the new ELAI observations have established a state of control which is significantly different from the previous pre-converged ELAI observations. Jointly enforcing these two rules implies convergence based on the notion that convergence enjoys a state of consistently decreased expectation of finding new minima in future function evaluations.

Considering the optimization procedure outlined in Figure (2), the check for convergence indicated in step 7) amounts to computing new EWMA Z_i values, and control limits, from the inclusion of the most recent observation of the improvement distribution, and checking if the subsequent set of Z_i satisfy both of the above rules of the EWMA convergence chart. Satisfying one, or none, of the convergence rules indicates insufficient exploration and further iterations of optimization are required to gather more information about the objective function.

4 Examples

We first look at two synthetic examples from the optimization literature, where the true optimum is known, so we can be sure we have converged to the true global minimum. We tune the EWMA Convergence Charts for each of these synthetic examples, then extrapolate the choice of w to provide a real world example from hydrology.

4.1 Rosenbrock

The Rosenbrock function (Rosenbrock, 1960), $f(x_1, x_2) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$, was an early test problem in the optimization literature. It combines a narrow, flat parabolic valley with steep walls, and thus it can be difficult for gradient-based methods. Convergence is non-trivial to assess, because optimization routines can take some time to explore the relatively flat, but non-convex, valley floor for the global minimum. Here we focus on the region $-2 \leq x_1 \leq 2$, $-3 \leq x_2 \leq 5$. While the region around the mode presents some minor challenges, this problem is unimodal, and thus represents a relatively easier optimization problem in the context of Bayesian optimization, with a well-behaved convergence process.

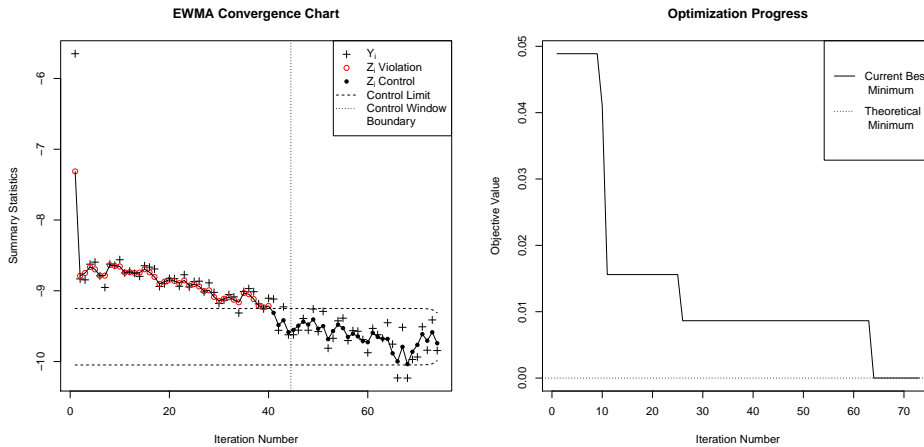


Figure 4: Rosenbrock function: Convergence chart on the left, optimization progress on the right.

We estimate λ via the minimum S_λ estimator, $\hat{\lambda} \approx 0.5$. Due to the relative simplicity of this problem we find that $w = 30$ results in a well behaved convergence pattern with a final ELAI variance of

0.35. Figure 4 shows the result of surrogate model optimization at convergence, as assessed by our method. The right panel shows the best function value (y -axis) found so far at each iteration (x -axis), and verifies that we have found the global minimum. The left panel shows the convergence chart, with the control window to the right of the vertical line, and the control limits indicated by the dashed lines. Iteration 74 is the first time that all EWMA points, in the control window, are observed within the control limits, and thus we declare convergence. This declaration of convergence comes after the global minimum has been found, but not too many iterations later, just enough to establish convergence. Note that the EWMA points generally trend downward until the global minimum is found at iteration 63.

4.2 Rastrigin

The 2- d Rastrigin function is a commonly used test function for evaluating the performance of global optimization schemes such as genetic algorithms (Whitley et al., 1996), $f(x_1, x_2) = \sum_{i=1}^2 [x_i^2 - 10 \cos(2\pi x_i)] + 2(10)$. The global behavior of Rastrigin is dominated by the spherical function, $\sum_i x_i^2$, however Rastrigin has been oscillated by the cosine function and vertically shifted so that it achieves a global minimum value of 0 at the lowest point of its lowest trough at (0, 0). We focus on the domain $-2.5 \leq x_i \leq 2.5$. This function is highly multimodal, and the many similar modes present a challenge for identifying convergence. The multimodality of this problem increases the variability of the EI criterion, and thus represents a moderately difficult optimization problem.

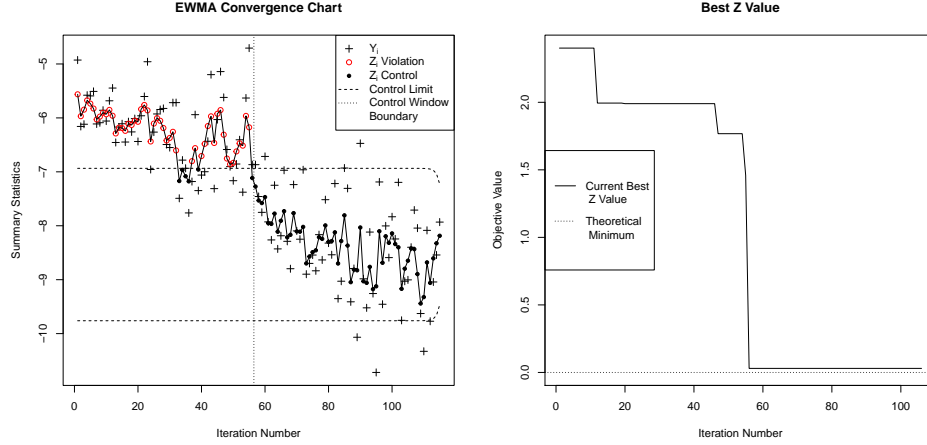


Figure 5: Rastrigin function: Convergence chart on the left, optimization progress on the right.

We estimate $\hat{\lambda} \approx 0.4$. The decreased value of $\hat{\lambda}$, relative to Rosenbrock, increases the smoothing capabilities of the EWMA procedure, as a response to the increased noise in the ELAI series. A larger w is needed to recognize convergence in the presence of increased noise in the ELAI criterion; $w = 60$ was found to work well, with a final ELAI variance of 1.71.

Figure (5) shows the convergence chart (left) and the optimization progress of the algorithm (right) after 115 iterations of optimization. Although the variability of the ELAI criterion increases as optimization proceeds, large ELAI values stop arriving after iteration 55, coincidentally with the surrogate model's discovery of the Rastrigin's main mode, as seen in the right panel of Figure (5). Furthermore notice that optimization progress in Figure (5, right) demonstrates that convergence in this case does indeed represent approximate identification of the theoretical minimum of the function, as indicated by the dashed horizontal line at the theoretical minimum.

4.3 Lockwood Case Study

The previous examples have focused on analytical functions with known minima, helping develop an intuition for tuning the EWMA convergence chart parameters and to ensure that our methods correspond to the identification of real optima. Here we apply the EWMA convergence chart on the Lockwood pump and treat problem, originally presented by Matott et al. (2011). This case study

331 considers an industrial site along the Yellowstone River in Montana, with groundwater contaminated
 332 by chlorinated solvents. Six pumps extract contaminated groundwater to attempt to prevent contam-
 333 ination of the river. The objective is to minimize the cost of running the pumps while preventing
 334 contamination of the river, and a computer simulator is used to compute the objective function.

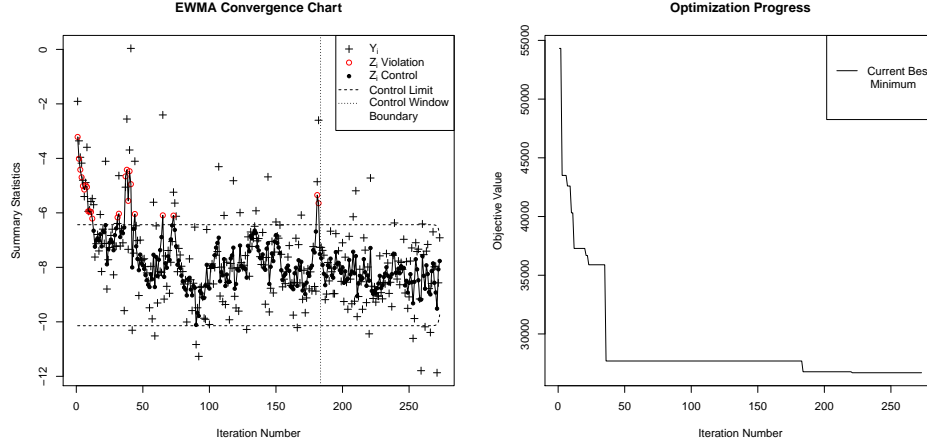


Figure 6: Lockwood Case-study: Convergence chart on the left, optimization progress on the right.

335 By using the fitted values of w and the observed ELAI variance in each of the two previous examples
 336 we extrapolate an appropriate value of w for this case study based on an observed ELAI variance of
 337 2.86, resulting in an estimated w of $93 \approx \left(\frac{60-30}{1.71-0.35} \right) 2.86 + 30$, as discussed in Section 3.4. λ was
 338 chosen via the minimum S_λ estimator to be $\hat{\lambda} \approx 0.4$.

339 The convergence chart for monitoring the optimization of the Lockwood case study is shown in
 340 the left panel of Figure (6). Convergence in this case does not occur with a dramatic shift in the
 341 mean level of the ELAI criterion, but rather convergence occurs as the series stabilizes after large
 342 ELAI values move beyond the control limit. Interestingly the last major spike in the ELAI series
 343 is observed alongside the discovery of the final major jump in the current best minimum value as
 344 seen at about iteration 180 in the right panel of Figure (6). The EWMA convergence chart identifies
 345 convergence as the EWMA statistic associated with this final ELAI spike eventually exits the control
 346 window at iteration 270. The solution shown here corresponds to $f(x) \approx 26696$. This solution is
 347 corroborated as a point of diminishing returns by the analysis of [Gramacy et al. \(2015\)](#) on the same
 348 problem, as seen in their average EI surrogate modeling behavior.

349 5 Conclusion

350 Adapting the notion of control from the SPC literature, the EWMA convergence chart outlined here
 351 aims to provide an objective standard for identifying convergence in the presence of the inherent
 352 stochasticity of the improvement criterion in this setting. The examples provided here demonstrate
 353 how the EWMA convergence chart may accurately and efficiently identify convergence in the con-
 354 text of Bayesian optimization. We note that our approach could be applied with any optimization
 355 algorithm that allows computation of an expected improvement at each iteration.

356 As for any optimization algorithm, a converged solution may only be considered as good as the al-
 357 gorithm’s exploration of f . Thus poorly tuned strategies may never optimize f to their fullest extent,
 358 but the EWMA convergence chart presented here may still claim convergence in these cases. The
 359 EWMA convergence chart may only consider convergence in the context of the algorithm in which
 360 it is embedded, and thus should be interpreted as a means of identifying when a global algorithm
 361 has converged, and it is beneficial to stop iterating the routine and reflect upon the results.

References

- Aerne, L. A., Champ, C. W., and Rigdon, S. E. (1991). Evaluation of control charts under linear trend. *Communications in Statistics-Theory and Methods*, 20(10):3341–3349.
- Box, G. E. P., Luceño, A., and Paniagua-Quñones, M. D. C. (1997). *Statistical Control by Monitoring and Adjustment*. Wiley, New York.
- Brochu, E., Cora, V. M., and de Freitas, N. (2010). A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning. Technical Report 1012.2599, arXiv.
- Diwale, S. S., Lymperopoulos, I., and Jones, C. (2015). Optimization of an airborne wind energy system using constrained gaussian processes with transient measurements. In *First Indian Control Conference*, number EPFL-CONF-199719.
- Gramacy, R. B. (2007). tgp: an r package for bayesian nonstationary, semiparametric nonlinear regression and design by treed gaussian process models. *Journal of Statistical Software*, 19(9):1–46.
- Gramacy, R. B. (2020). *Surrogates: Gaussian Process Modeling, Design, and Optimization for the Applied Sciences*. Chapman & Hall/CRC.
- Gramacy, R. B., Gray, G. A., Le Digabel, S., Lee, H. K. H., Ranjan, P., Wells, G., and Wild, S. M. (2015). Modeling an augmented lagrangian for blackbox constrained optimization. *Technometrics*. to appear, preprint arXiv:1403.4890.
- Gramacy, R. B. and Lee, H. K. H. (2008). Bayesian treed gaussian process models with an application to computer modeling. *Journal of the American Statistical Association*, 103(483).
- Gramacy, R. B. and Taddy, M. (2010). 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(6):1–48.
- Jones, D. R., Schonlau, M., and Welch, W. J. (1998). Efficient global optimization of expensive black-box functions. *Journal of Global Optimization*, 13(4):455–492.
- Lucas, J. M. and Saccucci, M. S. (1990). Exponentially weighted moving average control schemes: properties and enhancements. *Technometrics*, 32(1):1–12.
- Matott, S. L., Leung, K., and Sim, J. (2011). Application of matlab and python optimizers to two case studies involving groundwater flow and contaminant transport modeling. *Computers & Geosciences*, 37(11):1894–1899.
- Mockus, J. (1989). *Bayesian Approach to Global Optimization*. Kluwer, Dordrecht.
- Pourmohamad, T. and Lee, H. K. H. (2021). *Bayesian Optimization with Application to Computer Experiments*. Springer, Cham, Switzerland.
- Rosenbrock, H. H. (1960). An automatic method for finding the greatest or least value of a function. *The Computer Journal*, 3(3):175–184.
- Santner, T. J., Williams, B. J., and Notz, W. (2003). *The design and analysis of computer experiments*. Springer, New York.
- Schonlau, M., Jones, D., and Welch, W. (1998). Global versus local search in constrained optimization of computer models. In *New Developments and applications in experimental design*, number 34 in IMS Lecture Notes - Monograph Series, pages 11–25. JSTOR.
- Scrucca, L. (2004). qcc: an r package for quality control charting and statistical process control. *R News*, 4/1:11–17.
- Shewhart, W. A. (1931). *Economic control of quality of manufactured product*. D. Van Nostrand Company, New York.

- 407 Taddy, M. A., Lee, H. K. H., Gray, G. A., and Griffin, J. D. (2009). Bayesian guided pattern search
408 for robust local optimization. *Technometrics*, 51(4):389–401.
- 409 Whitley, D., Rana, S., Dzuber, J., and Mathias, K. E. (1996). Evaluating evolutionary algorithms.
410 *Artificial Intelligence*, 85(1-2):245–276.
- 411 Zou, C., Liu, Y., and Wang, Z. (2009). Comparisons of control schemes for monitoring the means
412 of processes subject to drifts. *Metrika*, 70(2):141–163.

413 Checklist

- 414 1. For all authors...
- 415 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
416 contributions and scope? [Yes]
- 417 (b) Did you describe the limitations of your work? [Yes]
- 418 (c) Did you discuss any potential negative societal impacts of your work? [N/A] Neg-
419 ative societal impacts in the field of optimization relate the particular objective func-
420 tions being optimized rather than the methods used to determine convergence of the
421 optimization algorithm itself.
- 422 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
423 them? [Yes]
- 424 2. If you are including theoretical results...
- 425 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 426 (b) Did you include complete proofs of all theoretical results? [N/A]
- 427 3. If you ran experiments...
- 428 (a) Did you include the code, data, and instructions needed to reproduce the main exper-
429 imental results (either in the supplemental material or as a URL)? [Yes] See supple-
430 mentary materials.
- 431 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
432 were chosen)? [Yes]
- 433 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
434 iments multiple times)? [No] While we have run replicates, it isn’t clear how to
435 measure uncertainty.
- 436 (d) Did you include the total amount of compute and the type of resources used (e.g.,
437 type of GPUs, internal cluster, or cloud provider)? [No] Details about computation
438 time will vary between particular implementations of different Bayesian optimization
439 algorithms. Compute used in the computation of the Bayesian optimization algorithm
440 itself is irrelevant to the scope of this paper, and identification of convergence requires
441 a negligible amount of compute.
- 442 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 443 (a) If your work uses existing assets, did you cite the creators? [Yes] The Lockwood
444 example uses an existing asset and the original author Matott et al. (2011) is cited.
- 445 (b) Did you mention the license of the assets? [No] The Lockwood code does not have
446 a public license. The Lockwood case study is merely used here as an example. The
447 details of the Lockwood code are not fundamental to our method of determining con-
448 vergence.
- 449 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
450 Example code is given in the supplemental material.
- 451 (d) Did you discuss whether and how consent was obtained from people whose data
452 you’re using/curating? [N/A] None.
- 453 (e) Did you discuss whether the data you are using/curating contains personally identifi-
454 able information or offensive content? [N/A] None.
- 455 5. If you used crowdsourcing or conducted research with human subjects...

- 456 (a) Did you include the full text of instructions given to participants and screenshots, if
457 applicable? [N/A]
- 458 (b) Did you describe any potential participant risks, with links to Institutional Review
459 Board (IRB) approvals, if applicable? [N/A]
- 460 (c) Did you include the estimated hourly wage paid to participants and the total amount
461 spent on participant compensation? [N/A]