
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 (??). A common application is for computer simulation experiments ?. 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.

Our motivating example is a hydrology application, the Lockwood pump-and-treat problem (Matott et al., 2011), discussed in more detail in Section 4.3, wherein contamination in the groundwater 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. Finding the global constrained minimum is a difficult problem where it is easy for optimization routines to temporarily get stuck in a local minimum. Without knowing the answer in advance, how does one know when to terminate the optimization algorithm?

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

37 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
38 threshold strategy to the convergence of Bayesian optimization has not yet been adequately justified.
39 In fact, this use of EI ignores the nature of the EI criterion as a random variable, and oversimplifies
40 the stochastic nature of convergence in this setting. Thus it is no surprise that this treatment of the
41 EI criterion can result in an inconsistent stopping rule as demonstrated in Figure (1).
42

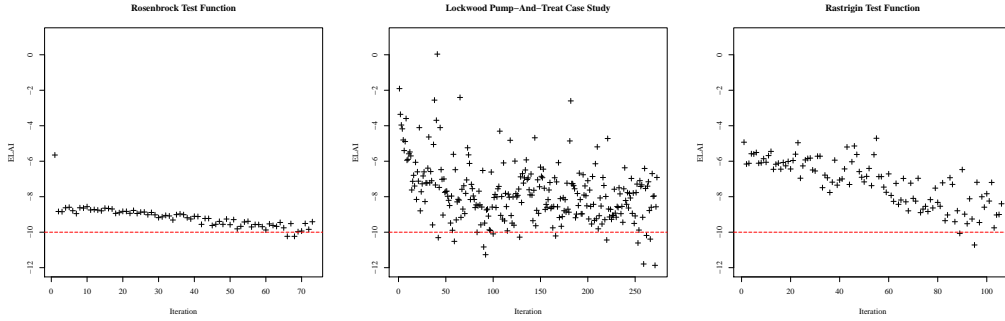


Figure 1: Three Expected Log-normal Approximation to the Improvement series (more details in Section 4) plotted alongside an example convergence threshold value shown as a dashed line at -10.

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

64 Since the Improvement function is itself random, attempting to set a lower threshold bound on the EI,
65 without consideration of the underlying EI distribution through time, over-simplifies the dynamics
66 of convergence in this setting. Instead, we propose taking the perspective of Statistical Process
67 Control (SPC), where a stochastic series is monitored for consistency of the distribution of the most
68 recently observed values. In the next section, we review the surrogate model approach and the use
69 of EI for optimization. In Section 3, we discuss our inspiration from SPC and how we construct
70 our convergence chart. Section 4 provides synthetic and real examples, and then we provide some
71 conclusions in the final section.

72 2 Bayesian Optimization via Expected Improvement

73 Bayesian optimization attempts to solve problems of the form

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

74 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
 75 defined via constraints. Without loss of generality, we frame all optimizations as minimizations in
 76 this paper, as maximization can be recovered by minimizing the negative of the function. Bayesian
 77 optimization proceeds by iteratively developing a statistical surrogate model of the objective func-
 78 tion f , and using predictions from the statistical surrogate to choose the next point to evaluate based
 79 on some criterion. A common choice of surrogate model is the Gaussian process (GP) ??, as it
 80 combines flexibility with smoothness.

81 In many cases the assumption of a globally smooth f with a homogeneous uncertainty structure
 82 can provide an effective and parsimonious model. However, in other problems, f may have sharp
 83 boundaries, f may show different levels of smoothness across its domain, or numerical simulators
 84 may have variable stability in portions of the domain. In this paper, we use treed Gaussian pro-
 85 cesses Gramacy and Lee (2008), a generalization of a standard GP that uses treed partitioning of
 86 the domain, fitting separate hierarchically-linked stationary GP surfaces to separate portions of f .
 87 The partitioning scheme is fit via a reversible jump MCMC algorithm and averaging over the full
 88 parameter space to provide smooth predictions except where the data call for a discontinuous pre-
 89 diction. We use the R package `tgp` (Gramacy, 2007; Gramacy and Taddy, 2010). While the treed
 90 GPs provide additional modeling flexibility, we emphasize that the approach of this paper can be
 91 applied to standard GPs as well as any surrogate model that provides both predictions and predictive
 92 uncertainty.

93 2.1 Expected Improvement

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

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

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

107 2.2 Optimization Procedure

108 Optimization can be viewed as a sequential de-
 109 sign process, where locations are selected for
 110 evaluation on the basis of how likely they are
 111 to decrease the objective function, i.e., based
 112 on the EI. Optimization begins by initially col-
 113 lecting a set, \mathbf{X} , of locations to evaluate the
 114 true function, f , to gather a basic impression
 115 of f . A statistical surrogate model is then fit-
 116 ted with $f(\mathbf{X})$ as observations of the true func-
 117 tion. Using the surrogate model, a set of can-
 118 didate points, $\tilde{\mathbf{X}}$, are selected from the domain
 119 and the EI criterion is calculated among these
 120 points. The candidate point that has the high-
 121 est EI is then chosen as the best candidate for a
 122 new minimum and thus, it is added to \mathbf{X} . The
 123 objective function is evaluated at this new location and the surrogate model is refit based on the
 124 updated $f(\mathbf{X})$. The optimization procedure carries on in this way until convergence. The key con-

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 $\arg\max_{\tilde{x}_i} \mathbb{E} [I(\tilde{x}_i)]$ to \mathbf{X} .
- 7) Check convergence.
- 8) If converged exit. Otherwise go to 2).

tribution of this paper is an automated method for checking convergence, which we develop in the next section.

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 $\mathbb{E}[\log I] \approx N\left(\mu, \frac{\sigma^2}{n}\right)$, 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 ,

176 and variance, ϕ , of the associated normal distribution are given by the following relation.

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

177 Using this relation we do not need to transform any of the improvement samples. We compute the
 178 empirical mean and variance of the unaltered, approximately log-normal, improvement samples,
 179 then use relation (2) to directly compute ψ as the Expectation under the Log-normal Approximation
 180 to the Improvement (ELAI). The ELAI value is useful for assessing convergence because of the
 181 reduced right skew of the log of the posterior predictive improvement distribution. Additionally,
 182 the ELAI serves as a computationally robust approximation of the $\mathbb{E}[\log I]$ under reasonable log-
 183 normality of the improvements. Furthermore, both the $\mathbb{E}[\log I]$ and ELAI are distributed approx-
 184 imately normally in repeated sampling. This construction allows for more consistent and accurate
 185 use of the fundamental theory on which our SPC perspective depends.

186 3.3 Exponentially Weighted Moving Average

187 The Exponentially Weighted Moving Average (EWMA) control chart (Lucas and Saccucci, 1990;
 188 Scrucca, 2004) elaborates on Shewhart's original notion of control by viewing the repeated sam-
 189 pling process in the context of a moving average smoothing of series data. Pre-convergence ELAI
 190 evaluations tend to be variable and overall decreasing, and so do not necessarily share distributional
 191 consistency among all observed values. Thus a weighted series perspective was chosen to follow
 192 the moving average of the most recent ELAI observations while still smoothing with some memory
 193 of older evaluations. EWMA achieves this robust smoothing behavior, relative to shifting means,
 194 by assigning exponentially decreasing weights to successive points in a rolling average among all
 195 of the points of the series. Thus the EWMA can emphasize recent observations and shift the focus
 196 of the moving average to the most recent information while still providing shrinkage towards the
 197 global mean of the series.

198 If Y_i is the current ELAI value, and Z_i is the EWMA statistic associated with this current value,
 199 then the initial value Z_0 is set to Y_0 and for $i \in \{1, 2, 3, \dots\}$ the EWMA statistic is expressed
 200 as, observation, Y_i . The recursive expression of the statistic ensures that all subsequent weights
 201 geometrically

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

202 Above, λ is a smoothing parameter that defines the weight (i.e. $0 < \lambda \leq 1$) assigned to the most
 203 recent decrease as they move back through the series.

204 Box et al. (1997) describes a method for comput-
 205 ing optimal choices of λ by minimizing the sum of
 206 squared forecasting deviations (S_λ). Through this
 207 analysis of S_λ , as seen in Figure (3), it is evident
 208 that EWMA charts can be very robust to reasonable
 209 choices of λ , due to the small first and second deriva-
 210 tives of S_λ for a large range of sub-optimal choices
 211 of λ around $\hat{\lambda}$. In fact, Figure (3) shows that for
 212 $\lambda \in [0.2, 0.6]$, S_λ stays within 10% of its the mini-
 213 mum possible value.

214 Identifying convergence in this setting now requires
 215 the computation of control limits on the EWMA
 216 statistic. As in the simplified \bar{x} -chart, defining the
 217 control limits for the EWMA setting amounts to con-
 218 sidering an interval on the sampling distribution of
 219 interest. In the EWMA case we are interested in the
 220 sampling distribution of the Z_i . Assuming that the
 221 Y_i are *i.i.d.* then Lucas and Saccucci (1990) show
 222 that we can write $\sigma_{Z_i}^2$ in terms of σ_Y^2 .

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

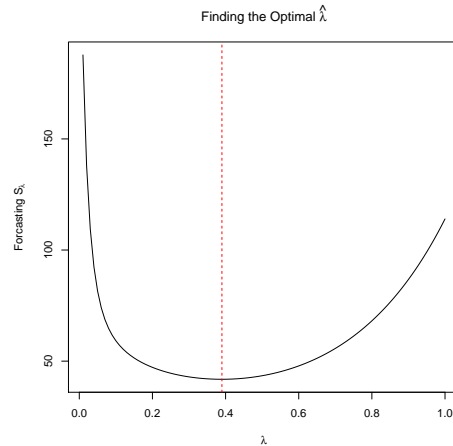


Figure 3: S_λ as calculated for ELAI values derived under the Rastrigin test function. $\hat{\lambda}$ is shown by the vertical dashed line.

Thus if $Y_i \overset{i.i.d.}{\sim} N\left(\mu, \frac{\sigma^2}{n}\right)$, 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. (5).

$$CL_i = \mu \pm c\sigma_{Z_i} = \mu \pm c \frac{\sigma}{\sqrt{n}} \sqrt{\left(\frac{\lambda}{2-\lambda}\right) [1 - (1-\lambda)^{2i}]} \quad (5)$$

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, $i \rightarrow \infty$, 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 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, however 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. Rosebrock and Rastrigin). From exploration of w in known examples, it is clear that w increases 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, $\hat{w} = \frac{\Delta w}{\Delta V(\text{ELAI})} \hat{v} + 30$. Further exploration of the exact form of an estimator of w is left to be discovered, although the connection of w with sample size calculations is a promising line of research in itself.

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 to 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](#)) 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. It is generalizable to higher dimensions, but we use the two-dimensional version here. Convergence is non-trivial to assess, because optimization routines can take a while 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 GP surrogate model optimization. This example illustrates a well-behaved convergence process.

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.

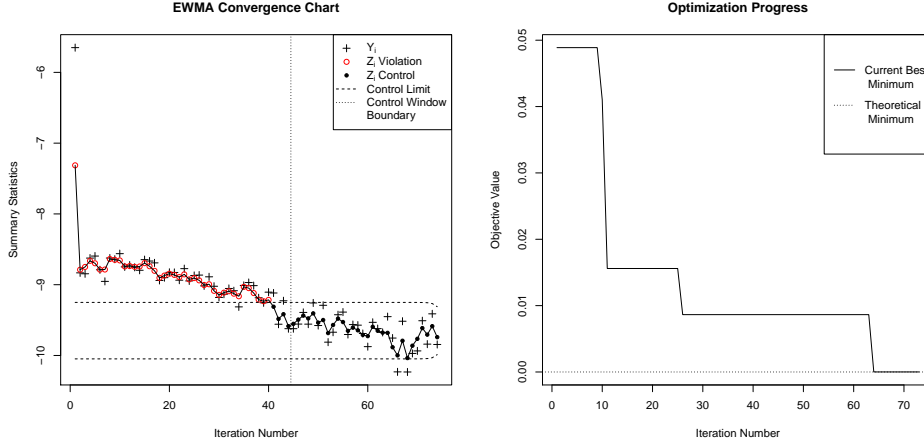


Figure 4: Rosenbrock function: Convergence chart on the left, optimization progress on the right.

4.2 Rastrigin

The Rastrigin function is a commonly used test function for evaluating the performance of global optimization schemes such as genetic algorithms (Whitley et al., 1996). 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).

Rastrigin is generalizable to arbitrarily many dimensions, but to develop intuition, this example considers Rastrigin over the 2 dimensional square $-2.5 \leq x_i \leq 2.5$. Rastrigin is a highly multimodal function, and as such, 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 in this context. It should be noted that by increasing the size of the search domain, either by increasing the bounds of the search square and/or increasing the dimension of the domain would make this example considerably more difficult.

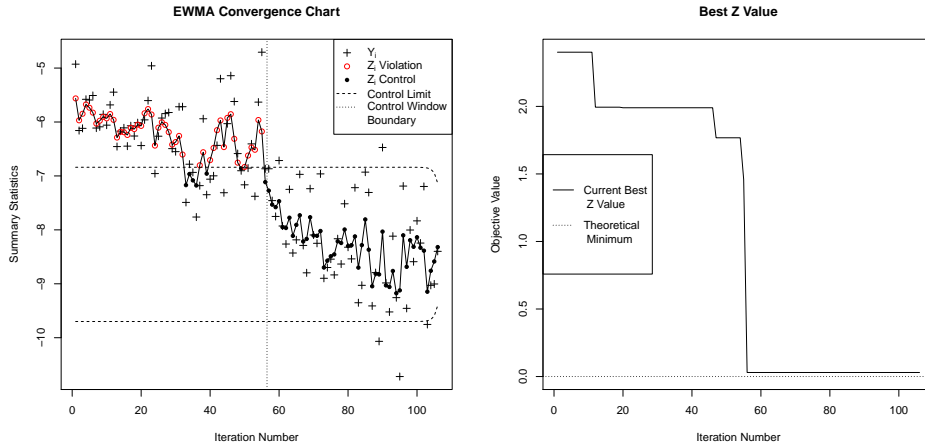


Figure 5: Rastrigin function: Convergence chart on the left, optimization progress on the right.

$\hat{\lambda}$ in this example is calculated to be about 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. The added noise of the ELAI series here comes from the regular discovery of dramatic new modes as optimization proceeds. Due to the increased complexity of Rastrigin relative to Rosenbrock, a larger w is needed to recognize convergence in the presence of increased noise in

the ELAI criterion. In this application $w = 50$ 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 105 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. This is done for the sake of developing 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 in the practical optimization setting of pump and treat optimization problems as formulated by Mayer et al. (2002). Specifically we consider the Lockwood pump and treat problem, originally presented by Matott et al. (2011).

The Lockwood pump and treat case study considers an industrial site, along the Yellowstone River in Montana, with groundwater contaminated by chlorinated solvents. If left untreated, this contaminated groundwater may contaminate the Yellowstone river, as dictated by the hydrology of the system. In order to control this contaminated groundwater, a total of six pumps, situated over two plumes of the contaminated groundwater are used to redirect groundwater away from the river to a treatment facility. Due to the cost of running these pumps, it is desirable to determine how to best allocate the pumping effort among these pumps so as to determine the lowest cost pumping strategy to protect the river. Pumping each of these six wells at different rates can drastically change the groundwater behavior, and thus a numerical simulation of the system is required to predict the behavior of the system at a given set of pumping rates.

The objective function, $f(\mathbf{x})$, to be minimized in this case, can be expressed as the sum of the pumping rates for each pump (a quantity proportional to the expense of running the pumps in USD), with additional large penalties associated with any contamination of the river.

$$f(\mathbf{x}) = \sum_{i=1}^6 x_i + 2[c_a(\mathbf{x}) + c_b(\mathbf{x})] + 20000[\mathbb{1}_{c_a(\mathbf{x}) > 0} + \mathbb{1}_{c_b(\mathbf{x}) > 0}] \quad (6)$$

Here $c_a(\mathbf{x})$ and $c_b(\mathbf{x})$ are outputs of a simulation, indicating the amount of contamination, if any, of the river as a function of the pumping rates, \mathbf{x} , for each of the six wells. Any amount of contamination of the river results in a large stepwise penalty which introduces a discontinuity into the objective function, at the contamination boundary. Each x_i is bounded on the interval $0 \leq x_i \leq 20,000$, representing a large range of possible management schemes. The full problem defines a six-dimensional optimization problem to determine the optimal rate at which to pump each well, so as to minimize the loss function defined in Eq (6). Since the loss function is defined over a large and continuous domain, and running the numerical simulation of the system is computationally expensive, this example presents an ideal situation for use with surrogate model based optimization.

By using the fitted values of w and the observed ELAI variance in each of the two previous examples we extrapolate an appropriate value of w for this case study based on an observed ELAI variance of 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. Again

λ was chosen via the minimum S_λ estimator to be $\hat{\lambda} \approx 0.4$ in this case. This level of smoothing is required here to reduce the noise in the ELAI criterion due to the large search domain, as well as the complicated contamination boundary among the six wells. Furthermore these features of the objective function complicate fit of the surrogate model and thus more function evaluations are required to produce an accurate model of f .

The convergence chart for monitoring the optimization of the Lockwood case study is shown in the left panel of Figure (6), as computed with $\hat{\lambda} \approx 0.4$ and $w = 93$. Convergence in this case does not occur with a dramatic shift in the mean level of the ELAI criterion, but rather convergence

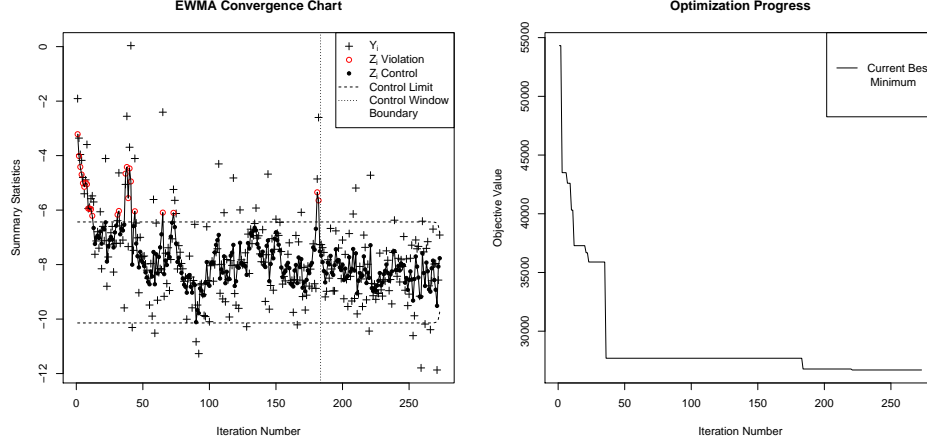


Figure 6: Lockwood Case-study: Convergence chart on the left, optimization progress on the right.

occurs as the series stabilizes after large ELAI values move beyond the control limit. Interestingly the last major spike in the ELAI series is observed alongside the discovery of the final major jump in the current best minimum value as seen at about iteration 180 in the right panel of Figure (6). The EWMA convergence chart identifies convergence as the EWMA statistic associated with this final ELAI spike eventually exits the control window at iteration 270. The solution shown here corresponds to $f(\mathbf{x}) \approx 26696$ at $\mathbf{x} \approx [0, 6195, 12988, 3160, 1190, 3163]$. This solution is well corroborated as a point of diminishing returns for this problem, by the analysis of Gramacy et al. (2015) on the same problem, as seen in their average EI surrogate modeling behavior.

5 Conclusion

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

As for any optimization algorithm, a converged solution may only be considered as good as the algorithm's exploration of f . Thus poorly tuned surrogate modeling strategies may never optimize f to their fullest extent, but the EWMA convergence chart presented here may still claim convergence in these cases. The EWMA convergence chart may only consider convergence in the context of the algorithm in which it is embedded, and thus should be interpreted as a means of identifying when an algorithm has converged rather than when the lowest minimum has been found. For poorly tuned surrogate modeling strategies the EWMA convergence chart may only identify that the algorithm has reached a point of diminishing returns, or that it has converged to a local mode; for well-tuned surrogate modeling strategies, this point should correspond with the realization of an optimal solution. In either case, the EWMA convergence chart identifies the moment at which it is beneficial to stop iterating the routine and reflect upon the results.

The EWMA convergence chart presented here is intended as a starting point for establishing an appropriate analysis of convergence for sequential surrogate modeling optimization algorithms. Details of the particular implementation of these methods may improve through further analysis of model usage and parameter estimation. The strategy presented here for transforming the improvement distribution via the Log-normal approximation to the improvement distribution has shown to be an empirically effective and computationally simple solution to better meet the assumptions of the EWMA control charting methodology. However, some applications may find it worthwhile to explore other transformations which could result in higher overall transformed signal to noise ratios, across a more broad set of improvement distributions. For example, rather than adopting the ELAI transformed estimate from the improvement distribution, it may be computationally feasible to apply

the two-parameter Box-Cox transformation (Box and Cox, 1964) to the improvement samples,

$$y_i^{(\lambda)} = \begin{cases} \frac{(y_i + \lambda_2)^{\lambda_1} - 1}{\lambda_1} & \lambda_1 \neq 0 \\ \log(y_i + \lambda_2) & \lambda_1 = 0 \end{cases} \quad (7)$$

thus alleviating any difficulties due to numerical truncation of the improvement samples at 0, while finding a flexible transformation to reduce skew. It should be noted that this approach adds additional computational expense, while our ELAI transformation requires minimal computation. Additionally the EWMA convergence chart could benefit from a more precise method for choosing the control window size parameter, w , although developing such a method would be a major research project itself; our empirical solution has worked well in practice. Although improvements to the details of these methods may exist, the fundamental consideration of the stochastic nature of convergence in this setting would remain, and SPC offers a nice framework for its inclusion.

References

- Abrahamsen, P. (1997). A review of gaussian random fields and correlation functions. Technical Report 917, Norsk Regnesentral/Norwegian Computing Center.
- 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. and Cox, D. R. (1964). An analysis of transformations. *Journal of the Royal Statistical Society. Series B (Methodological)*, 26:211–252.
- Box, G. E. P., Luceño, A., and Paniagua-Quinones, M. D. C. (1997). *Statistical Control by Monitoring and Adjustment*. Wiley, New York.
- 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., 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 Lee, H. K. H. (2012). Cases for the nugget in modeling computer experiments. *Statistics and Computing*, 22(3):713–722.
- 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.
- Kolda, T. G., Lewis, R. M., and Torczon, V. (2003). Optimization by direct search: New perspectives on some classical and modern methods. *SIAM Review*, 45:385–482.
- 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.

- 465 Mayer, A. S., Kelley, C., and Miller, C. T. (2002). Optimal design for problems involving flow and
466 transport phenomena in saturated subsurface systems. *Advances in Water Resources*, 25(8):1233–
467 1256.
- 468 Rosenbrock, H. H. (1960). An automatic method for finding the greatest or least value of a function.
469 *The Computer Journal*, 3(3):175–184.
- 470 Sacks, J., Welch, W. J., Mitchell, T. J., and Wynn, H. P. (1989). Design and analysis of computer
471 experiments. *Statistical science*, 4:409–423.
- 472 Santner, T. J., Williams, B. J., and Notz, W. (2003). *The design and analysis of computer experi-*
473 *ments*. Springer, New York.
- 474 Schonlau, M., Jones, D., and Welch, W. (1998). Global versus local search in constrained opti-
475 mization of computer models. In *New Developments and applications in experimental design*,
476 number 34 in IMS Lecture Notes - Monograph Series, pages 11–25. JSTOR.
- 477 Scrucca, L. (2004). qcc: an r package for quality control charting and statistical process control. *R*
478 *News*, 4/1:11–17.
- 479 Shewhart, W. A. (1931). *Economic control of quality of manufactured product*. D. Van Nostrand
480 Company, New York.
- 481 Stein, M. L. (1999). *Interpolation of spatial data: some theory for kriging*. Springer, New York.
- 482 Taddy, M. A., Lee, H. K. H., Gray, G. A., and Griffin, J. D. (2009). Bayesian guided pattern search
483 for robust local optimization. *Technometrics*, 51(4):389–401.
- 484 Whitley, D., Rana, S., Dzuber, J., and Mathias, K. E. (1996). Evaluating evolutionary algorithms.
485 *Artificial Intelligence*, 85(1-2):245–276.
- 486 Zou, C., Liu, Y., and Wang, Z. (2009). Comparisons of control schemes for monitoring the means
487 of processes subject to drifts. *Metrika*, 70(2):141–163.

488 Checklist

489 The checklist follows the references. Please read the checklist guidelines carefully for information
490 on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
491 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
492 the appropriate section of your paper or providing a brief inline description. For example:

- 493 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 494 • Did you include the license to the code and datasets? **[No]** The code and the data are
495 proprietary.
- 496 • Did you include the license to the code and datasets? **[N/A]**

497 Please do not modify the questions and only use the provided macros for your answers. Note that the
498 Checklist section does not count towards the page limit. In your paper, please delete this instructions
499 block and only keep the Checklist section heading above along with the questions/answers below.

500 1. For all authors...

- 501 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
502 contributions and scope? **[TODO]**
- 503 (b) Did you describe the limitations of your work? **[TODO]**
- 504 (c) Did you discuss any potential negative societal impacts of your work? **[TODO]**
- 505 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
506 them? **[TODO]**

507 2. If you are including theoretical results...

- 508 (a) Did you state the full set of assumptions of all theoretical results? **[TODO]**

- 509 (b) Did you include complete proofs of all theoretical results? **[TODO]**
- 510 3. If you ran experiments...
- 511 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 512 mental results (either in the supplemental material or as a URL)? **[TODO]**
- 513 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 514 were chosen)? **[TODO]**
- 515 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 516 ments multiple times)? **[TODO]**
- 517 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 518 of GPUs, internal cluster, or cloud provider)? **[TODO]**
- 519 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 520 (a) If your work uses existing assets, did you cite the creators? **[TODO]**
- 521 (b) Did you mention the license of the assets? **[TODO]**
- 522 (c) Did you include any new assets either in the supplemental material or as a URL?
- 523 **[TODO]**
- 524 (d) Did you discuss whether and how consent was obtained from people whose data
- 525 you're using/curating? **[TODO]**
- 526 (e) Did you discuss whether the data you are using/curating contains personally identifi-
- 527 able information or offensive content? **[TODO]**
- 528 5. If you used crowdsourcing or conducted research with human subjects...
- 529 (a) Did you include the full text of instructions given to participants and screenshots, if
- 530 applicable? **[TODO]**
- 531 (b) Did you describe any potential participant risks, with links to Institutional Review
- 532 Board (IRB) approvals, if applicable? **[TODO]**
- 533 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 534 spent on participant compensation? **[TODO]**