

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/341658398>

Locating Infinite Discontinuities in Computer Experiments

Article in *SIAM/ASA Journal on Uncertainty Quantification* · January 2020

DOI: 10.1137/18M1209076

CITATIONS

0

READS

26

3 authors, including:



Ying-Chao Hung

National Chengchi University

36 PUBLICATIONS 182 CITATIONS

[SEE PROFILE](#)



George Michailidis

University of Michigan

366 PUBLICATIONS 6,417 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Control of Network Systems [View project](#)



Continuous Space Reinforcement Learning [View project](#)

Locating Infinite Discontinuities in Computer Experiments*

Ying-Chao Hung[†], George Michailidis[‡], and Horace PakHai Lok[§]

Abstract. Identification of input configurations so as to meet a pre-specified output target under a limited experimental budget has been an important task for computer experiments. Such a task often involves the development of response models and design of experimental trials that rely on the models exhibiting continuity and differentiability properties. Motivated by two canonical examples in systems and manufacturing engineering, we propose a strategy for locating the boundary of the response surface in computer experiments, wherein on one side the response is finite, whereas on the other side is infinite, leveraging ideas from active learning and quasi-Monte Carlo methods. The strategy is illustrated on an example from computer networks engineering and one from precision manufacturing and shown to allocate experimental trials in a fairly effective manner. We conclude by discussing extensions of the proposed strategy to characterize other types of output discontinuity or non-differentiability in high-cost experiments, including jump discontinuities in the target output response or pathological structures such as kinks and cusps.

Key words. Computer experiments, Infinite discontinuity, Active learning, Support vector machines, Quasi-Monte Carlo methods

AMS subject classifications. 60G15, 62M20, 62K99, 65Y20, 91B74

1. Introduction. Computer experiments involve complex computer codes underlying some physical phenomenon or experiment. They act as surrogates to explore the relationship between the corresponding input factors and an output measure of interest. However, running such computer codes is usually computationally expensive in terms of CPU time [47]. A natural question arising is how to obtain a comprehensive understanding of output performance over the input domain by utilizing a limited number of experimental trials, which leads to how to best model the input-output relationship and select the corresponding computer code trials.

The Gaussian Stochastic Process (GASP) model has been the most popular technique for modeling the input-output relationship of computer experiments [40, 45, 46, 47], briefly summarized next. Denote the code output corresponding to K -dimensional inputs (locations) x_1, \dots, x_n by $y(x_1), \dots, y(x_n)$. The stationary GASP model models such pairs of inputs-outputs as $y(x_i) = \mu + z(x_i)$, where μ is the output mean and $z(x_i)$ is a spatial process with mean zero, constant variance σ^2 , and covariance $\text{Cov}(z(x_i), z(x_j)) = \sigma^2 R_{ij}$, where R_{ij} is a correlation function defined by a measure of spatial distance between inputs x_i and x_j . Further, it is assumed that $y(X)$ follows a multivariate normal distribution with mean vector $\mathbf{1}_n \mu$ and covariance matrix $\sigma^2 [R_{ij}]$. The benefit of adopting the GASP model is that

*Submitted to the editors DATE.

Funding: This work was supported in part by research grant MOST 106-2118-M-004-002-MY2 (YCH) and NSF IIS 1632730 and NSF DMS 1830175 (GM).

[†]Department of Statistics, National Chengchi University, Taipei 11605, Taiwan (hungy@nccu.edu.tw).

[‡]Department of Statistics and the Informatics Institute, University of Florida, Gainesville, FL 32611, USA (gmichail@ufl.edu).

[§]Department of Statistics, National Chengchi University, Taipei 11605, Taiwan (105354008@nccu.edu.tw).

36 the explicit form of the best linear unbiased predictor (BLUP) at a new input location can
37 be obtained in closed form [46, 47], thus allowing researchers to explore the experiment’s
38 output performance over the entire input space. There is a considerable body of work on
39 related models, including limit kriging [27], co-kriging [16], blind kriging [28], Bayesian treed
40 Gaussian process models [19] and scaled Gaussian stochastic process [20], just to name a few.
41 Based on the above-mentioned characteristics, there exists a number of design methods that
42 aim to select the experimental trials so as to minimize uncertainty in estimating the GASP
43 model estimation, as well as output predictions. Examples include a sequential design based
44 on a designated improvement function for contour estimation [44], a design over non-convex
45 regions based on multidimensional scaling to the geodesic distance [43], and a sequential
46 design strategy based on maximum mutual information [2], etc. A large portion of such
47 design strategies belongs to the framework of space filling designs [26, 47].

48 Note that GASP based modeling and the associated design techniques developed for com-
49 puter experiments automatically assume that the output response surface is *smooth* (*contin-*
50 *uous* and *differentiable*) [2]. This assumption provides a convenient mathematical formulation
51 and corresponding solution, but may not be suitable for certain settings as outlined next.
52 For example, a fundamental issue in many service engineering systems is to identify possible
53 input configurations so that some performance measure of interest does not increase in an un-
54 bounded manner to *infinity* [4, 8, 15, 34, 50]. This is related to the problem of system *stability*
55 and often needs to be validated through a long-run simulation. Another example is the laser
56 cutting epoxy film in manufacturing engineering, for which the goal is to detect the physical
57 cutting limits by characterizing a number control factors. Since the cutting limits refer to
58 the boundary of input configurations that produce “nonzero” outputs (or the boundary of
59 the no-cut region), by taking the reciprocal of the output value, the goal here is exactly the
60 same as identifying input configurations that produce finite outputs. Due to the complexity of
61 the laser cutting process, a computer surrogate model is often utilized to understand how the
62 control factors affect the outputs. Note that experimental trials needed to identify the input
63 configurations for such problems are often computationally expensive. A detailed description
64 of these two examples is provided in Section 2.

65 The broad objective of this paper is to develop methodology to identify input configura-
66 tions in computer experiments that trace the boundary that separates infinite valued outputs
67 from finite valued ones. Technically, it aims to locate the boundary in the input domain
68 (henceforth referred to as the set of *infinite discontinuities*) in an efficient and accurate man-
69 ner. As pointed out in the literature [19], the popular GASP model cannot predict well
70 non-continuous changes in the output (e.g. jumps). To overcome this limitation, we propose
71 to formulate the problem of locating infinite outputs as a classification one. In this case, input
72 configurations that lead to finite output values are labeled as belonging to class “+1”, while
73 those that lead to infinite output values are labeled as belonging to another class “−1”. Hence,
74 the problem becomes to identify the decision boundary separating the two classes. Based on
75 this, the proposed strategy leverages ideas from *active learning*, a semi-supervised machine
76 learning technique (see [33] and references therein), whose primary goal is to train a good
77 classifier for predicting the output labels by using a small amount of sampled data. Note that
78 rather than randomly obtaining a full training data set, active learning iteratively selects the
79 most “informative” instances and query their labels to train the classifier. Further, the use

80 of active learning for the problem at hand belongs to the *pool-based sampling* [36, 53] and is
 81 related to sequential experimental designs. Further, active learning in our strategy is based
 82 on the powerful formalism of support vector machines (SVM) [3, 10, 42, 48, 54] for label pre-
 83 dictions. However, the sampling technique is different from those employed by conventional
 84 methods in active learning. In order to obtain a good initial training set (i.e., the passive
 85 learning process) and accommodate well the training samples selected at later learning stages
 86 (i.e., the active learning process), our sampling strategy utilizes a robust quasi-Monte Carlo
 87 (QMC) method - called uniform design (UD) in the statistical literature [12, 13].

88 Note that conceptually the problem addressed in this paper shares similarities to the so-
 89 called regression discontinuity design (RDD) that has many applications in economics and
 90 engineering [25, 32, 35, 51]. The remainder of the paper is organized as follows. In [Section 2](#),
 91 a motivating example that addresses the stability issue of queueing service engineering sys-
 92 tems is described in detail. In [Section 3](#), the strategy based on the concept of active learning
 93 is introduced. The focus is on (i) how to choose an adequate initial training set and unlabeled
 94 instances at later learning stages from a pool of candidate instances; and (ii) how to develop
 95 sophisticated SVM based active learners for label prediction. Note that there are two itera-
 96 tively updated weight functions associated with data involved in this active learning process
 97 - one is designed to improve the accuracy of the SVM learner, while the other to represent
 98 the sample informativeness/importance based on which the training data can be best selected
 99 by UD. In [Section 4](#), the proposed strategy is illustrated on a queueing service engineering
 100 system and numerical evaluations are provided. A discussion and some concluding remarks
 101 are drawn in [Section 5](#). Finally, an accelerated algorithm for implementing UD in higher di-
 102 mensional spaces and the Hierarchical Mixing Linear SVMs (HMLSVM, see [55]) for locating
 103 infinite discontinuities with piecewise linear shapes are given in [Appendix A](#) and [Appendix C](#),
 104 respectively.

105 **2. Motivating Examples.** Next, we describe in considerable detail two motivating exam-
 106 ples from systems engineering and precision manufacturing.

2.1. Example 1: Stability of Queueing Systems. Consider a multi-class queueing system
 comprised of Q parallel, infinite capacity, and first-in-first-out (FIFO) queues, with each queue
 containing traffic from a different job class. Each job of class q arrives according to a random
 process, carries a random amount of workload (or service requirement), and waits in line for
 service. For generosity here we assume each arrival process has a *time-heterogeneous rate*
 (i.e., the instantaneous input rate can vary over time). Thus, it is reasonable to define a
 time-averaged arrival rate for each input class q , i.e.,

$$x_q = \lim_{T \rightarrow \infty} \frac{E[A_q(T)]}{T},$$

107 where $A_q(T)$ is the amount of workload arriving at queue q before time T , $q = 1, \dots, Q$.
 108 These time-averaged input rates can be collected in a vector $x = (x_1, \dots, x_Q)$, where each
 109 element x_q is assumed to be non-degenerate. At any point in time, the system can be in
 110 one of M ($M \geq 2$) service modes indexed by $m \in \{1, \dots, M\}$. When the system switches
 111 into the m -th service mode, the jobs in queue q receive service at a constant rate μ_{mq} (i.e.,
 112 the amount of work that can be processed in one time unit). Therefore, mode m is asso-

113 ciated with the service rate vector $U_m = (\mu_{m1}, \dots, \mu_{mQ})$. Note that the switching scheme
 114 between service modes automatically introduces intricate dependencies amongst the queues.
 115 For a general real-life system, we assume the switching time between any two service modes
 116 i and j (denoted by Δ_{ij}) is non-negligible, while no service is provided during each switching
 117 epoch. This particular queueing system is known as the Switched Processing System (SPS)
 118 and captures the essence of the fundamental resource allocation problem in many modern
 119 service engineering applications involving heterogeneous processors and multiple classes of
 120 job traffic flows; examples include network switches for routing traffic through the Internet,
 121 modern flexible manufacturing processes, and automatic call distributor (ACD) (see [1, 23]
 122 and references therein).

123 A fundamental issue for SPS studies is to characterize the allowable range of input rates
 124 (i.e. x) so that system can achieve a certain level of “stability” under a given service-
 125 mode switching policy. For example, denote by $W_q(t)$ the workload (or remaining service
 126 requirement) of queue q at time t , which is obviously a function of the input rate vector
 127 x . Since each arrival process is time-heterogeneous, system stability can be defined based
 128 on the time-averaged expected amount of workload [57]. That is, given an input-rate vector
 129 $x = (x_1, \dots, x_Q)$, the system is *stable* if

$$130 \quad (2.1) \quad \theta_q(x) = \lim_{T \rightarrow \infty} \int_0^T \frac{E[W_q(t)]}{T} dt < \infty \quad \text{for all } q = 1, \dots, Q.$$

131 On the other hand, the system is characterized as “unstable” if $\theta_q(x)$ is infinite for at least
 132 one queue q . Given a service-mode switching policy π , the maximum set of input-rate vectors
 133 x satisfying (2.1) is called the *stability region*.

134 Note that based on the definition in (2.1), the outer *boundary* of the stability region
 135 clearly represents the set of infinite discontinuities. However, identifying this boundary is
 136 mathematically intractable for even very small systems and hence in practice is approximated
 137 through computer simulations of the underlying system. To illustrate, let us consider a simple
 138 2-queue system with three service modes $U_1 = (0, 5)$, $U_2 = (6, 4)$, $U_3 = (7, 0)$ and assume
 139 that the switching time between any two service modes is constant, say, $\Delta_{ij} \equiv \Delta$ for all
 140 i, j . Jobs of each input arrive according to a compound Poisson process with exponentially
 141 distributed service requirements, having a rate of one. Suppose a well-known service-mode
 142 switching policy called the MaxProduct is employed [1, 23, 24], the estimated stability regions
 143 are shown in Figure 1 for two selected switching times $\Delta = 0.5$ and 1.0, where stability is
 144 validated based on (2.1) by simulating the system at a set of superimposed input-rate grids
 145 over an \mathbb{R}^2 region $[7, 0] \times [0, 5]$. Note that for comparison purpose, the maximum stability
 146 region (i.e. when $\Delta = 0$) is also included in Figure 1.

147 Note that for a system with the stability region shown in Figure 1, we can simply define
 148 $\theta(x) = \theta_1(x)$, $\theta_2(x)$, or $\theta_1(x) + \theta_2(x)$. Thus, by (2.1) $\theta(x)$ is characterized as “finite” when
 149 the input-rate vector $x = (x_1, x_2)$ is inside the stability region. Specifically, $\theta(x)$ becomes
 150 substantially larger when x gets fairly close to the boundary of the stability region and jumps
 151 to “infinity” when x reaches or crosses the boundary. The later case naturally leads to an
 152 output model characterized by jump discontinuities. The main challenge of identifying the
 153 stability region based on simulation is the associated computational cost, since for each input
 154 configuration the system needs to be simulated for a large number of events (job arrivals,

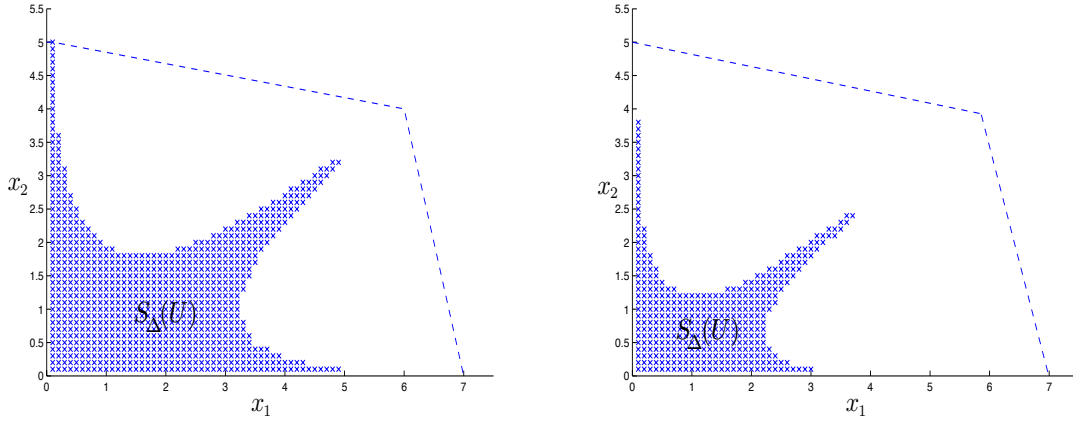


Figure 1. The estimated stability regions $S_{\Delta}(U)$ with service-mode switching times $\Delta = 0.5$ (the dotted area in the left panel) and $\Delta = 1.0$ (the dotted area in the right panel) under the MaxProduct policy. The polygon with the boundaries described by the dashed lines refers to the maximum stability region with $\Delta = 0$.

155 switching between service modes epochs, job service and departures), so that the long term
 156 average defined in equation (2.1) can be estimated accurately. This computational time grows
 157 fast for larger systems compromising of more queues and more switching modes, or for input
 158 configurations close to the boundary of the stability region.

159 **2.2. Example 2: Limits of Laser Cutting Process.** Laser cutting epoxy film has replaced
 160 the conventional die cutting film and is being widely used for circuit board attachment to
 161 carriers and housings in defense electronics, commercial radio frequency (RF) microwave and
 162 microelectronics industries. An important issue in such manufacturing process is to detect
 163 the physical cutting limits by characterizing the control factors such as laser power, beam
 164 radius, pulse duration, focal position, and Rayleigh length. Denote the cutting width at the
 165 bottom of the material by $y(x_1, x_2)$, where x_1 represents the pulse duration (e.g. in μs), x_2
 166 represents the laser power (e.g. in watts), while other factors are kept constant. The goal is to
 167 identify the control factor settings (x_1, x_2) so that cut through does not occur (i.e. the cutting
 168 width $y(x_1, x_2) = 0$). However, identifying the limits of the cutting process is a challenging
 169 task since (i) the relationship between the cutting width and other control factors may not be
 170 trivial; and (ii) each experimental trial is physically expensive [32]. Thus, simulation of such
 171 complex processes along with an adequate surrogate model is often used to understand how
 172 the control factors affect the outputs.

173 Suppose that the experiment is emulated, say, with a proper transformation of scales, by
 174 the following functional relationship (see [32] for a similar function):

$$175 \quad (2.2) \quad y(x_1, x_2) = \begin{cases} \frac{1}{100} \{ [(x_1-7)^2 + (x_2-7)^2 - 62]^2 + [(x_1-7)^2 - 0.5(x_1-7) - 0.5(x_2-7) - 1.5]^2 \} + 10 & \text{if } x_1^2 + x_2^2 > 80, \\ 0 & \text{if } x_1^2 + x_2^2 \leq 80. \end{cases}$$

176 The interpolated surface and the contour plot of function (2.2) over the 2-dimensional square
 177 domain $[0, 14] \times [0, 14]$ are depicted in Figure 2. As can be seen, the boundary of the blue-color

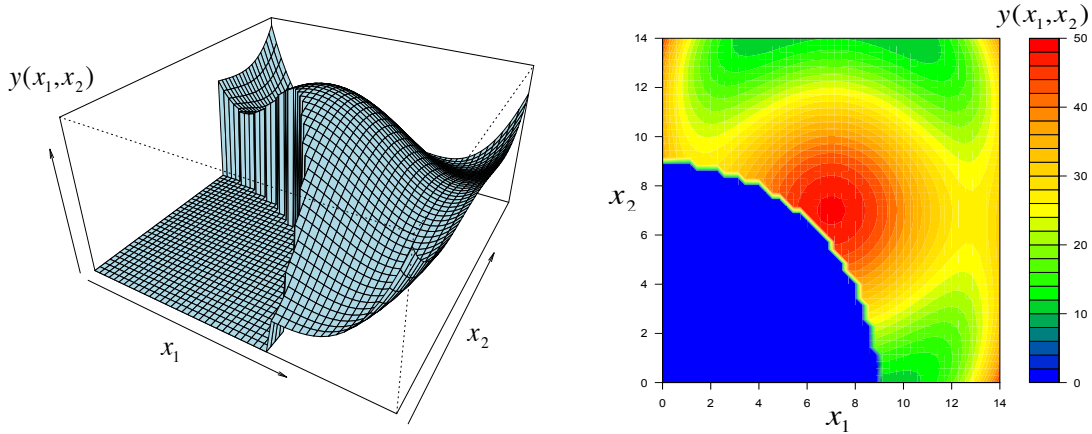


Figure 2. The interpolated surface (left panel) and the contour plot (right panel) of function (2.2) based on $50 \times 50 = 2,500$ grids superimposed over the input domain $[0, 14] \times [0, 14]$. The blue-color area in the contour plot represents the settings of “no cut” (i.e. $y(x_1, x_2) = 0$) in the input domain.

178 area in the right panel of Figure 2 represents the set of discontinuities in the input domain,
 179 which is the focus of this work. Also note that identifying the limits (i.e. discontinuities) of
 180 the cutting process can be thought as a “dual” problem to that given in Example 1, since one
 181 can consider the reciprocal of $y(x_1, x_2)$ as the response in Example 1. In words, by letting
 182 $\theta(x) = 1/y(x_1, x_2)$, the input points that produce the “zero” output are equivalent to the
 183 infinite discontinuities described in Example 1.

184 **3. Locating Infinite Valued Discontinuities via Active Learning.** Since the GASP model
 185 is not suitable for finding the infinite valued discontinuities, we propose next an alternative
 186 strategy based on the concept of active learning. The strategy is also sequential and formulates
 187 the primary problem as a classification one. Further, it includes two main components: (i)
 188 a quasi-Monte Carlo method called uniform design (UD) for sampling the training data (or
 189 allocating experimental trials); and (ii) an SVM classifier for estimating the discontinuities.
 190 Note that there are two well-designed weight functions associated with these two components.
 191 The first one represents the importance of data to be sampled, while the second one represents
 192 the importance of sampled data for training the SVM classifier. We start with introducing
 193 the sampling technique based on UD, while active learning with SVM is discussed afterwards.

194 **3.1. Sampling Based on Quasi-Monte Carlo Methods.** The quasi-Monte Carlo (QMC)
 195 method is a deterministic version of the Monte Carlo method. It yields a faster convergence
 196 rate for numerical integration and has proved to be a robust technique for exploring relation-
 197 ships between input factors and experiments outputs. The key idea of QMC is to choose a set
 198 of n points $\mathcal{P} = \{p_1, \dots, p_n\}$ that are uniformly distributed on a bounded input domain D ,
 199 $D \subset \mathbb{R}^K$. Let $Z(n)$ be the collection of all possible sets $\{p_1, \dots, p_n\}$ on D . A uniform design
 200 (UD) is an efficient QMC that seeks the best set $\mathcal{P} \in Z(n)$ such that a measure of uniformity

201 $M(\mathcal{P})$ is minimized. That is,

$$202 \quad (3.1) \quad \mathcal{P}^* = \arg \min_{\mathcal{P} \in Z(n)} M(\mathcal{P}).$$

203 The most popular choice for the measure of uniformity is the so-called “discrepancy” [13, 14,
 204 21]. In this study, we employ a measure named the *central composite discrepancy* (CCD),
 205 which is developed to overcome the limitations of conventional space filling designs [9, 38].
 206 The CCD is attractive from the following perspectives: (i) the solution is robust to the output
 207 model and has a faster convergence rate than traditional random sampling procedures (i.e. the
 208 property of quasi-Monte Carlo methods); (ii) it can be applied to any shape of design domains;
 209 (iii) the solution is invariant to rotations; (iv) the optimal design \mathcal{P}^* can be obtained to explore
 210 a designated output response by placing different weights at the input configurations. Note
 211 that (iv) is the key element of our proposed framework, since all allocated experimental points
 212 may not be equally important in locating the desired discontinuities. The formulation of the
 213 CCD with weights is given next.

214 At a particular learning stage s , let $f_s(x)$ be the weight function such that $f_s(x) > 0$ for
 215 all $x \in D$ and assume $\int_D f_s(x) dx = 1$ (this can be done by a simple normalization). With this
 216 formulation, the weight function $f_s(x)$ can also be viewed as the probability density function
 217 of all data points defined on D . To find a set of n points $\mathcal{P} = \{p_1, \dots, p_n\}$ in D so that they
 218 have a “good representation” for $f_s(x)$, we define

$$219 \quad (3.2) \quad WCCD_{f_s, p}(n, \mathcal{P}) = \left\{ \frac{1}{v(D)} \int_D \frac{1}{2^K} \sum_{k=1}^{2^K} \left| \frac{N(D_k(x), \mathcal{P})}{n} - F(D_k(x)) \right|^p dx \right\}^{1/p},$$

220 where $v(D)$ is the volume of D , $N(D_k(x), \mathcal{P})$ is the number of points allocated in the subregion
 221 $D_k(x)$ (note that D is partitioned into 2^K subregions at each x), $F(D_k(x)) = \int_{D_k(x)} f_s(x) dx$ is
 222 the proportion of points expected to be allocated in $D_k(x)$, $k = 1, \dots, 2^K$. Therefore, a good
 223 representation for $f_s(x)$ will be the set of points \mathcal{P}^* that minimizes (3.2). Note that solving \mathcal{P}^*
 224 is an NP-hard problem as the number of allocated design points goes to infinity. In practice,
 225 one can superimpose a set of grid points Z over the input domain D and approximate the
 226 solution by using sophisticated search algorithms [6, 9, 38]. The solution of such a discretized
 227 optimization problem is known as the nearly uniform design (NUD), which minimizes the
 228 approximation of $WCCD_{f_s, p}(n, \mathcal{P})$,

$$229 \quad (3.3) \quad \widehat{WCCD}_{f_s, p}(n, \mathcal{P}) = \left\{ \frac{1}{|Z|} \sum_{z \in Z} \frac{1}{2^K} \sum_{k=1}^{2^K} \left| \frac{N(D_k(z), \mathcal{P})}{n} - F(D_k(z)) \right|^p \right\}^{1/p},$$

230 where $\mathcal{P} \in Z(n) \subset Z$, and $|Z|$ represents the cardinality of Z . It should be noted that the
 231 solution \mathcal{P}^* here is deterministic and is different from the solution obtained by stochastic
 232 sampling methods (e.g. random sampling).

233 **3.2. Acceleration of the Sampling Procedure.** Next, we discuss how to accelerate the
 234 associated sampling procedure. Suppose a UD of size n needs to be found based on N

235 grid points (i.e. let $|Z| = N$) superimposed over the input domain D . Under the natural
 236 assumption that $K < n \ll N$, an exhaustive search requires computational time of order
 237 $O(N^n)$ for finding \mathcal{P}^* . To reduce the computation load, Chuang and Hung (2010) [9] proposed
 238 a Switching Algorithm and showed that an NUD can be constructed within the computational
 239 time of order $O(N^{(2+p)})$. Although the algorithm provides a good approximation to \mathcal{P}^* in
 240 lower dimensions, it may still be computationally intensive for many problems with multi-
 241 dimensional inputs. To see this, let us look at a rectangular domain with $K = 3$, $N = 10^3$, $n =$
 242 6 and $p = 2$ (wherein an ℓ_2 -distance is considered). Note that even for such a simple example,
 243 exhaustive search requires computation of order 10^{18} , while the Switching Algorithm requires
 244 computation up to order of 10^{12} , which represents a big improvement in absolute terms, but
 245 not still sufficient for most applications.

246 In order to implement UD with a more manageable time line, we introduce next an
 247 accelerated algorithm for finding an NUD. Let ψ_i be a non-degenerate subset of $\{1, \dots, K\}$
 248 and denote by Z^{ψ_i} the projection of Z onto the corresponding axes in a lower $|\psi_i|$ -dimensional
 249 space, $|Z^{\psi_i}| \geq n$. The idea of the proposed accelerated algorithm is to first locate the design
 250 points based on the lower $|\psi_i|$ -dimensional uniformity. Then, by augmenting the selected
 251 locations in the lower $|\psi_i|$ -dimensional space, we sequentially increase the dimensionality of
 252 Z^{ψ_i} so that the design points are relocated step by step to achieve better uniformity in the
 253 primary input space. It should be highlighted that the acceleration algorithm is developed
 254 to search a flexible and bigger area than that of the well-known Latin Hypercube Designs
 255 (LHD) [11], while retaining the property of uniformity and a fairly fast computation speed. A
 256 detailed description of the acceleration algorithm and discussion of its performance are given
 257 in [Appendix A](#).

258 We next introduce how to develop an adequate response model $g(x)$ at each design stage,
 259 so as to improve the estimation of the discontinuities based on the concept of active learning.

260 **3.3. Active Learning with SVM.** Let us recall the example of the queueing system intro-
 261 duced in [Section 2](#) and unify some notations. Denote the input rate vector by x , the labeled
 262 response, say, $\theta(x) = \theta_1(x)$ by $y(x)$, and the stability region $S_\Delta(U)$ by $T(x)$, respectively. To
 263 overcome the difficulty of modeling infinite outcomes, we formulate the problem as a binary
 264 classification one. That is, label $y(x) = +1$ represents the case that the input point x is inside
 265 the target region (or equivalently $\theta(x) < \infty$), while label $y(x) = -1$ represents the case that x
 266 is outside the target region (or $\theta(x) = \infty$). Therefore, the objective of locating/estimating the
 267 boundary of $T(x)$ becomes that of identifying the decision boundary between design points
 268 with labels $y(x) = +1$ and -1 in the primary data space, wherein the allocated experimental
 269 trials are treated as training samples. Since labeling data is time consuming or costly, instead
 270 of obtaining the training samples at random, active learning selects “informative” samples so
 271 as to maximize the classification accuracy with less training data. We focus on active learning
 272 based on an SVM classifier since: (i) SVM is a proven powerful classification tool capable of
 273 dealing with different shapes of decision boundaries; and (ii) the distance of a sample x to
 274 the decision boundary of $T(x)$ can be easily computed (see below) and subsequently used as
 275 a measure of information about input x .

276 Denote the training data set (i.e. the existing experimental trials) at a particular learning
 277 stage by $\{(x_1, y_1), \dots, (x_n, y_n)\}$, where $x_i = (x_{i1}, \dots, x_{iK})$ is a K -dimensional input vector

278 and $y_i = y(x_i) \in \{-1, +1\}$ is the corresponding output label. In order to obtain flexible
 279 decision boundaries (as shown in Figure 1), the SVM algorithm projects the training data
 280 into a higher dimensional feature space \mathcal{H} using a mapping function $\Phi(\cdot)$. With this mapping,
 281 it is shown that the optimal solution depends on data merely through the inner products in \mathcal{H}
 282 [54], that is, on functions of the form $\langle \Phi(x_i), \Phi(x_j) \rangle$. Hence, a computationally less expensive
 283 approach is to use a kernel function $K(x_i, x_j) = \langle \Phi(x_i), \Phi(x_j) \rangle$ instead of explicit mappings
 284 $\Phi(x_i)$ and $\Phi(x_j)$. This is known as the *kernel trick* - a commonly used technique in machine
 285 learning [49].

286 Since different data points may have different contributions to the training of a classifier,
 287 we use the weighted version of SVM (termed as WSVM) for active learning. Let W_i represent
 288 the weight assigned to each training sample x_i , $i = 1, \dots, n$. To obtain the WSVM based on
 289 all weighted training data (x_i, y_i, W_i) with a mapping function $\Phi(\cdot)$, one solves the constrained
 290 convex optimization problem [59]:

$$\begin{aligned}
 291 \quad & \underset{w, b}{\text{minimize}} \quad \frac{1}{2} w' w + C \sum_{i=1}^n W_i \xi_i \\
 292 \quad (3.4) \quad & \text{subject to} \quad \begin{cases} y_i (\langle w, \Phi(x_i) \rangle + b) \geq 1 - \xi_i, & i = 1, \dots, n, \\ \xi_i \geq 0, & i = 1, \dots, n. \end{cases}
 \end{aligned}$$

293 where $\langle w, \Phi(x) \rangle + b$ represents the hyperplane separating the y_i 's in the feature space \mathcal{H} , w
 294 is the coefficient vector, b controls the offset of the decision boundary from the origin, ξ_i are
 295 known as *slack variables* that penalize wrong classifications of y_i , and $C > 0$ controls the
 296 tradeoff between classification accuracy and the margin between two bounding planes. The
 297 standard approach for solving (3.4) is to formulate its dual given by [59]:

$$\begin{aligned}
 298 \quad & \underset{\alpha_i}{\text{maximize}} \quad \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i, x_j) \\
 299 \quad (3.5) \quad & \text{subject to} \quad \begin{cases} 0 \leq \alpha_i \leq C W_i, & i = 1, \dots, n, \\ \sum_{i=1}^n \alpha_i y_i = 0, \end{cases}
 \end{aligned}$$

300 where α_i are the associated Lagrange multipliers. Solving (3.5) yields the decision function

$$301 \quad (3.6) \quad g(x) = \text{sign}(\langle w, \Phi(x) \rangle + b) = \text{sign} \left(\sum_{i=1}^n \alpha_i y_i K(x_i, x) + b \right),$$

302 where $w = \sum_{i=1}^n \alpha_i y_i \Phi(x_i)$ (readers can refer to LIBSVM [5] for how to obtain the WSVM
 303 learner in practice). Note that a popular choice of the kernel function is the radial basis
 304 function (RBF), which has the form $K(x_i, x_j) = \exp\{-\gamma \|x_i - x_j\|^2\}$, $\gamma > 0$. The RBF
 305 kernel offers the following two advantages in practice: (i) it can produce nonlinear decision
 306 boundaries by mapping data into high (or infinite) dimensional space; and (ii) it has rela-
 307 tively low complexity for model selection. Therefore, it is particularly suitable for exploring
 308 high-dimensional data structures. In practice, the two parameters (C, γ) can be chosen by
 309 performing an appropriate grid search in R_+^2 and utilizing the idea of cross-validation (CV) so

310 that the “prediction error rate” is minimized [22, 29, 31]. However, there are some alternative
 311 methods that have been shown to achieve better efficiency and/or accuracy than traditional
 312 grid search based on CV. For example, the optimal value of γ can be found by minimizing the
 313 Fisher discriminant function [56], by considering parameter selection as a recognition problem
 314 [58], and by optimizing simultaneously the overall accuracies (for within-class samples) and
 315 kappa accuracies (for between-class samples) [37], just to name a few.

316 Let us denote the decision function obtained at learning stage s by $g_s(x)$, $x \in D$. Analogous
 317 to the idea of boosting [17, 18], the weight for each sampled point x_i is chosen as

$$318 \quad (3.7) \quad W_i = 1 + M_i^s = 1 + \sum_{t=1}^{s-1} I\{y_i \cdot g_t(x_i) < 0\},$$

319 where M_i^s represents the total number of times that x_i was wrongfully labeled (i.e. the
 320 decision function yields a different label from y_i , such as $y_i = +1$ and $g_t(x_i) < 0$, or $y_i = -1$
 321 and $g_t(x_i) > 0$) in the past $s - 1$ stages, $M_i^1 = 0$ for all i . Note that if M_i^s is large, the
 322 weight W_i at stage s is set to be large, while W_i is set to be small if M_i^s is small. Intuitively,
 323 this will force the decision function $g_s(x)$ to accommodate all the training samples based on
 324 their classification history, thus improving the overall classification accuracy. Once $g_s(x)$ is
 325 obtained at a particular learning stage s , the boundary of $T(x)$ can be simply estimated by

$$326 \quad (3.8) \quad \partial\tilde{T}_s(x) = \{x \in D : g_s(x) = 0\}.$$

327 Note that the goal of obtaining $\partial\tilde{T}_s(x)$ is to quantify the “informativeness” (or importance)
 328 of all unlabeled samples, based on which a suitable weight function can be incorporated into
 329 UD for choosing the next-stage instance from the pool.

330 Next, we proceed to introduce how the weight function $f_s(x)$ in (3.2) is defined in accor-
 331 dance with the quantified measure of information. For each instance x , compute its minimum
 332 Euclidean distance to the estimated boundary of the target region, that is,

$$333 \quad (3.9) \quad d(x, \partial\tilde{T}_s(x)) = \min_{\delta \in \partial\tilde{T}_s(x)} \|x - \delta\|.$$

334 If x is unlabeled and has a smaller $d(x, \partial\tilde{T}_s(x))$, we say that it is more informative for training
 335 the SVM classifier and thus more likely to be selected at the next stage. This idea is quite
 336 similar to that of the Simple Margin querying method [53], which chooses the unlabeled sample
 337 closest to the decision boundary in \mathcal{H} . Based on this argument, the weight function $f_s(x)$ for
 338 implementing UD is given by

$$339 \quad (3.10) \quad f_s(x) = \exp \left\{ -\frac{d(x, \partial\tilde{T}_s(x))}{\beta(s)} \right\}.$$

340 As can be seen in (3.10), the weight of an instance x becomes exponentially larger when it
 341 gets closer to the estimated decision boundary and becomes exponentially smaller when it
 342 gets farther away from the boundary. This will push the next-stage learning to query an
 343 instance near the estimated decision boundary $\partial\tilde{T}_s(x)$. Further, $\beta(s)$ is chosen as a *decreasing*

344 function of learning stage s , which is similar to the concept of *learning rate* [41]. The goal of
 345 this setup is to accelerate the process of querying unlabeled instances near the “true” decision
 346 boundary, thus improving the estimation efficiency at subsequent stages. To normalize the
 347 weight function, we simply set

$$348 \quad (3.11) \quad f_s^*(x) = \frac{f_s(x)}{\int_{x \in D} f_s(x) dx}.$$

349 *Remark:*

350 (a) For nonlinear decision boundaries the computation of $d(x, \partial\tilde{T}_s(x))$ may not be trivial,
 351 especially when the data dimension is high. A computationally cheaper way is to represent
 352 the input domain D as a pool of superimposed grid points Z (as done in [Subsection 3.1](#)) so
 353 that $\partial\tilde{T}_s(x)$ can be approximated by a subset of Z , say,

$$354 \quad (3.12) \quad \partial\tilde{T}_s(z) = \{z \in Z : -\epsilon < g_s(z) < \epsilon\},$$

355 where $\epsilon > 0$ is known as the *tolerance level* and chosen by the designer. Therefore, $d(x, \partial\tilde{T}_s(x))$
 356 can be approximated by

$$357 \quad (3.13) \quad d(x, \partial\tilde{T}_s(z)) = \min_{z \in \partial\tilde{T}_s(z)} \|x - z\|.$$

358 To evaluate the quality of the approximation based on (3.12), let us denote $|Z| = N$ by the
 359 total number of equal-size grids z superimposed over a typical K -dimensional input domain
 360 $D = [0, d]^K$. The value of N reflects the total computational cost one can afford so as to
 361 implement the quasi-Monte Carlo method. By definition, we then have

$$362 \quad (3.14) \quad N = (d/l)^K,$$

363 and thus $l = dN^{-(1/K)}$, where l is the side-length of each K -dimensional grid $z \in Z$. Note
 364 that if we use (3.12) to approximate the decision boundary obtained in (3.8), the maximum
 365 error corresponds to the “diameter” of each superimposed grid z , which is given by

$$366 \quad (3.15) \quad \sqrt{K}l = \sqrt{K}dN^{-(1/K)}.$$

367 It is easy to see that the approximation is influenced by the number of input dimensions (K),
 368 domain size (d) and maximum affordable computational cost (N).

369 (b) Instead of using $d(x, \partial\tilde{T}_s(x))$ to formulate the weight function $f_s(x)$ in (3.10), for compu-
 370 tational simplicity, one may consider the following two alternatives: (i) consider the distance
 371 between x and the obtained decision boundary $\langle w, \Phi(x) \rangle + b = 0$ in \mathcal{H} , which has a simple
 372 mathematical expression $|\langle w, \Phi(x) \rangle + b| / \|w\|$. Then, assign the weight for each x as

$$373 \quad (3.16) \quad f_s(x) = \exp \left\{ -\frac{|\langle w, \Phi(x) \rangle + b| / \|w\|}{\beta(s)} \right\};$$

374 (ii) consider the predictive value $g_s(x)$ and assign the weight for each x as

$$375 \quad (3.17) \quad f_s(x) = \exp \left\{ -\frac{\|g_s(x)\|}{\beta(s)} \right\}.$$

376 Note that assigning the weight based on the decision value is straightforward since $g_s(x) = 0$ on
 377 the decision boundary. Further, the decision values are readily available from software pack-
 378 ages (e.g. the R package e1071). Thus, for computational simplicity purposes, we strongly
 379 recommend assigning the weight for each x by using (3.17).

380

381 We summarize the detailed steps of our strategy for estimating the target region with
 382 nonlinear boundaries in Algorithm 3.1. Note that we denote by \mathcal{P}_s^* the training set obtained
 383 at learning stage s , where $|\mathcal{P}_s^*| = n_s$, $s = 1, 2, \dots$

384

Algorithm 3.1 Active Learning with SVM

- 1: Determine an appropriate training domain D associated with a set of superimposed grid points Z . Set the learning stage to $s = 1$ and determine the size of the initial training set n_1 . Find the initial training set \mathcal{P}_1^* of size n_1 that minimizes (3.3) based on Z , where $p = 2$ and the weight function is initially placed as $f_0(x) \equiv 1$. Train a kernel WSVM classifier based on \mathcal{P}_1^* with the weights $W_i \equiv 1$ and obtain the decision boundary $\partial\tilde{T}_1(z)$ by (3.12).
 - 2: **while** the stopping criterion is not met **do**
 - 3: Update the weight functions $f_s(x)$ and $f_s^*(x)$ by (3.17) and (3.11), respectively. Set $s = s + 1$.
 - 4: Select an unlabeled instance z^* from the pool $Z \setminus \mathcal{P}_{s-1}^*$ so that $\mathcal{P}_s^* = \mathcal{P}_{s-1}^* \cup \{z^*\}$ minimizes (3.3). Query the label of z^* .
 - 5: Update the weights W_i for all instances in \mathcal{P}_s^* by (3.7) and train a kernel WSVM classifier based on \mathcal{P}_s^* . Obtain a new decision boundary $\partial\tilde{T}_s(z)$.
 - 6: **end while**
 - 7: **return** $\partial\tilde{T}_s(z)$
-

385

We highlight a few important issues regarding the implementation of Algorithm 3.1.

386

- A suitable input domain D should be considered, so as to accurately and efficiently estimate the desired decision boundary. This often relies on the practitioner's knowledge about the response model. If possible, one can choose D in the form of hyper-rectangles, so as to accelerate the implementation of Algorithm 3.1.
- The number of grid points superimposed on D depends on the computational complexity of the proposed framework and the executing processor's performance. For example, if the input domain is high dimensional and/or the executing processor has moderate efficiency, a smaller number of grid points in Z should be considered.
- In cases where there is lack of knowledge regarding the output label, it is recommended to allocate a larger portion of the experimental budget to the initial training set (i.e., for the passive learning process) so as to better accommodate the model/sample uncertainty and improve the estimation efficiency. A rule of thumb for this proportion is about 20 ~ 50%, which was suggested by a large number of experimental scenarios.
- The choice of $\beta(s)$ in (3.10) controls the speed of convergence for the estimated decision boundary - it is basically a positive decreasing function of the learning stage s . However, if $\beta(s)$ decays too fast, the algorithm will be forced to stop at an earlier

387

388

389

390

391

392

393

394

395

396

397

398

stage (since almost all the experimental trials are allocated near the estimated decision boundaries). This will apparently accelerate the convergence of the algorithm, but lose the opportunity of exploring the output model at some unlabeled input locations. Due to the sampling feature of the proposed framework (i.e. only one query at each learning stage), it is recommended that $\beta(s)$ be chosen to have a moderately slow decay. For example, one may consider

$$\beta(s) = a_0 \times (c)^{-s},$$

399 where a_0 is a positive constant and $c \gtrsim 1$ (i.e. c is a positive constant greater than
400 (but close to) one). Note that such $\beta(s)$ apparently has a slower decay than the
401 exponential function. Numerical results also show that it performs well for various
402 types of computer experiments (see Section 4).

- 403 • The algorithm stops usually when the experimental budget is expended or the esti-
404 mated decision boundary has achieved a stable status.
- 405 • Note that one can speed up the learning of Algorithm 3.1 by selecting a batch of
406 samples at each learning stage. Such sampling process refers to the so-called batch-
407 mode active learning [52].

408 It should be mentioned that if we know a priori that the target region has a piecewise
409 linear boundary (e.g. the stability region when $\Delta = 0$ in Figure 1 or target regions comprising
410 of a number of unidentifiable linear constraints), the strategy of using the RBF kernel SVM
411 may not be plausible when the size of the training data is small. An alternative strategy is to
412 employ a technique analogous to the so-called Hierarchical Mixing Linear SVMs (HMLSVMs)
413 introduced in [55]. The idea of HMLSVMs is quite similar to that of the hierarchical classifi-
414 cation tree, but it assigns different weights to the training data and uses the linear SVM (thus
415 known as the WLSVM) as the decision function to split the nodes. The detailed description of
416 how this technique is implemented at each learning stage of our proposed framework is given
417 in Appendix C.

418 **4. Performance Assessment.** In this section, we illustrate the proposed strategy on the
419 two motivating examples and evaluate its performance. All numerical results presented here
420 were performed by using the statistical software package R (version 3.2.5) and executed on a
421 3.2 GHz Intel® Core™ i5-6500 processors with 16 GB of memory under the operating system
422 of Microsoft Windows 7 64-bit Service Pack 1 (SP1). The SVM techniques were implemented
423 by utilizing the R package “e1071” and “wSVM”, while all the tuning parameters used in the
424 proposed framework were determined based on a number of sensitivity tests.

425 **4.1. Queueing Example.** Let us consider the queueing example with a service-mode
426 switching time $\Delta = 0.5$ introduced in Subsection 2.1. The “true” stability region was ob-
427 tained by simulating the system at $50 \times 70 = 3,500$ superimposed grid points, as shown in
428 the left panel of Figure 1. We now demonstrate how to estimate the stability region based
429 on active learning with SVM, as given in Algorithm 3.1. Since no evidence indicates that
430 the target region has a linear (or piecewise linear) boundary, the RBF kernel is a reasonable
431 choice for fitting the WSVM model. Note that the initial training set \mathcal{P}_1^* is obtained by UD,
432 whose size is chosen to be $n_1 = 10$ (see Figure 3). For fitting the RBF kernel WSVM model,

433 the parameters (C, γ) are chosen to minimize the prediction error (based on the 10-fold cross-
 434 validation) based on the set of candidate grid points (e^i, e^j) , where $i, j = -10, -9, \dots, 9, 10$.
 435 As a result, the grid search yields the best parameters $(C, \gamma) = (1096.633, 0.0183)$. Since
 436 all computations now are based on the superimposed grid points, at each learning stage the
 437 estimated boundary of the target region has the form of (3.12), wherein ϵ is chosen to be 0.01.
 438 We then update the weight function by using (3.17), wherein we choose $\beta(s) = 100 \times (1.1)^{-s}$
 439 so that it has a moderately slow decay as the design stage moves. The estimated boundaries
 440 $\partial\tilde{T}_s(z)$ and the contour lines of the associated weight functions $f_s(x)$ at learning stages $s = 8$
 441 ($n_s = 17$), $s = 19$ ($n_s = 28$), $s = 46$ ($n_s = 55$), $s = 66$ ($n_s = 75$) are depicted in Figure 3 and
 Figure 4, respectively.

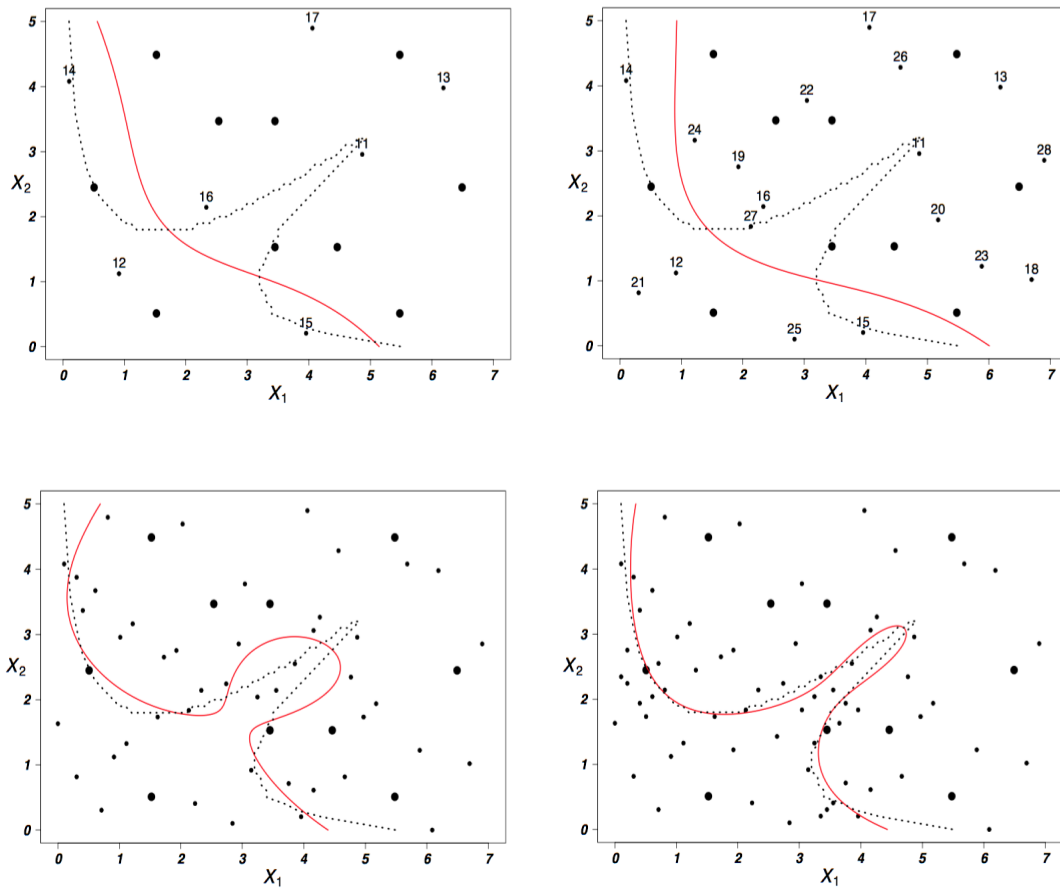


Figure 3. The estimated boundary $\partial\tilde{T}_s(z)$ (the red curve) based on WSVM at the learning stage $s = 8$ (upper left), $s = 19$ (upper right), $s = 46$ (lower left) and $s = 66$ (lower right), respectively. The dotted lines represent the true stability region, the big dots represent the selected initial training data, while the small numbered dots represent the selected samples at later learning stages.

442

443

As can be seen from Figure 3, the estimated target region gets closer to the true one as

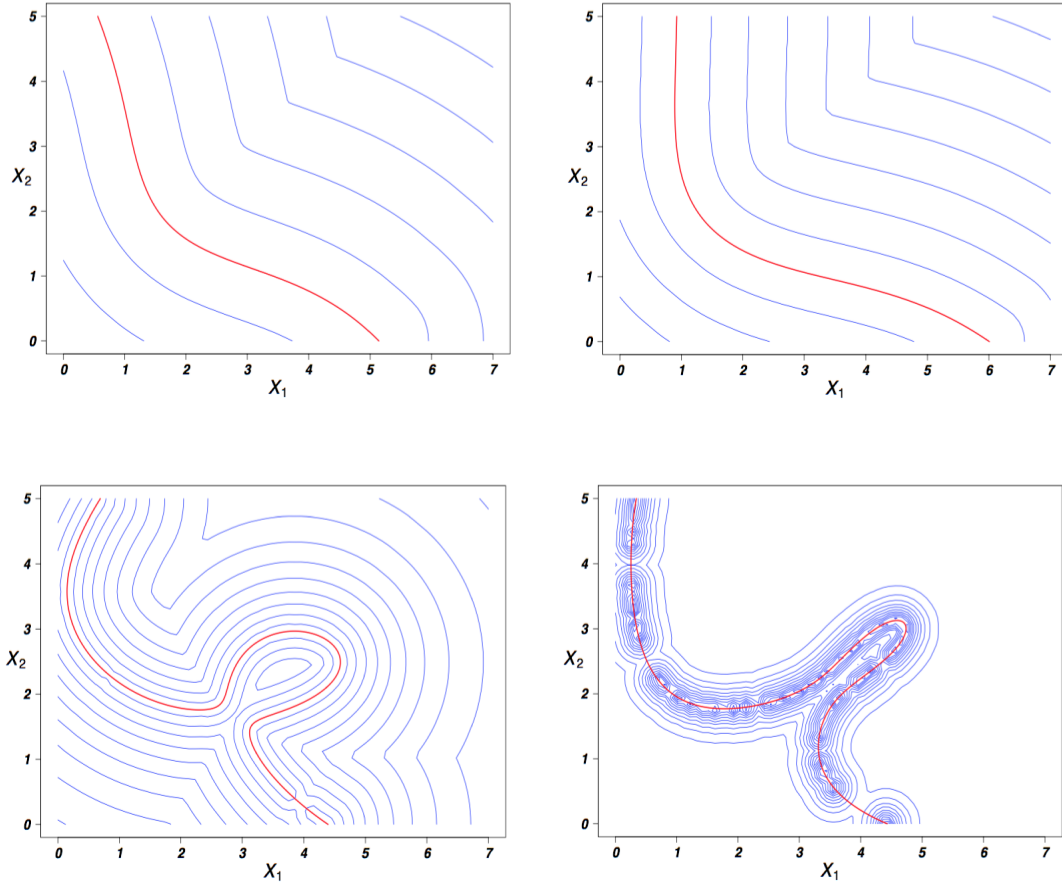


Figure 4. The contour lines of the weight functions $f_s(x)$ at the learning stages $s = 8$ (upper left), $s = 19$ (upper right), $s = 46$ (lower left) and $s = 66$ (lower right). The red line represents the weight function on the estimated boundary.

444 the learning stage increases. In addition, the initially selected training set uniformly spreads
 445 out the experimental trials over the input domain, while thereafter it tends to allocate the
 446 experimental trials closer and closer to the boundary of the estimated target region. Note
 447 that an adequate choice of the control process $\beta(s)$ in the weight function $f_s(x)$ allows us
 448 to allocate experimental trials in the areas not been explored at early design stages. These
 449 phenomena are also supported by the contour plots of the weight function shown in [Figure 4](#).

450 To evaluate the efficiency of [Algorithm 3.1](#), we compute the label error rate of the estimated
 451 target region based on the superimposed 3,500 grid points and compare with that obtained by
 452 (i) passive learning with SVM (i.e. allocate all experimental trials at once based on the UD);
 453 (ii) GASP modeling technique along with a sequential sampling strategy for contour estimation
 454 (a modification of the method introduced in [44], see [Appendix B](#) for detailed steps); and (iii)
 455 the GASP modeling technique based on pure UD (i.e. allocate all experimental trials at once).

456 Note that for (ii), we take the reciprocal of the primary response (i.e. $y(x)=1/\theta(x)=1/\theta_1(x)$)
 457 for building the GASP model. Thus, the contour of interest has height “close to zero” (but
 458 not equal to zero) and can be estimated in a stepwise manner by using a stage-dependent
 459 improvement function. The numerical results with respect to different sizes of training samples
 460 are depicted in Figure 5, where each GASP model is estimated by using the R package “GPfit”
 [39].

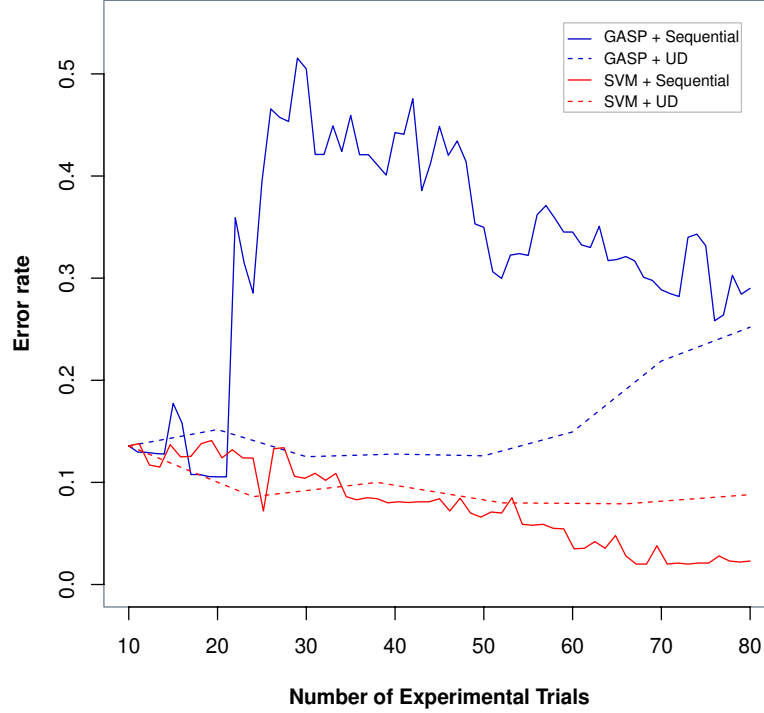


Figure 5. The error rates for estimating the stability region based on passive and active learning with SVM (red lines) and the GASP model for contour estimation (blue lines). The solid lines refer to the results by using the sequential sampling strategy; the dashed lines refer to the results by sampling at once based on UD.

461

462 As can be seen from Figure 5, the error rate of the proposed strategy based on active
 463 learning reduces in a steady manner to zero as the experimental budget increases. Note that
 464 the convergence may not be particularly fast, which is due to the spiky shape of the target
 465 region in the middle part of the experimental domain. However, the estimation error rate
 466 is no more than 6% when the total number of experimental trials exceeds 52. On the other
 467 hand, the error rate based on the GASP model is significantly larger. The estimation error
 468 rate remains high (above 26%), even when the number of experimental trials is increased
 469 to 80. Such a slow convergence rate may be induced by over-allocating experimental trials
 470 outside the target region, which is a consequence of utilizing the sampling strategy based
 471 on improvement function. Not surprisingly, the strategy of allocating all experimental trials
 472 at once based on UD does not work well for estimating the desired boundary. Even if we
 473 increase the size of training data to 80, the estimation error rate is always above 9% for both
 474 the GASP and SVM modeling methodologies. In conclusion, the strategy of employing active

475 learning with SVM outperforms that of employing passive learning and the GASP model for
 476 estimating the desired input boundary of infinite valued outcomes.

477 **4.2. The Laser Cutting Process.** Let us consider the laser cutting problem introduced
 478 in [Subsection 2.1](#), where the process has four control factors and the experiment is emulated
 479 by the following functional relationship

480 (4.1)

$$481 \quad y(x_1, x_2, x_3, x_4)$$

$$482 = \begin{cases} \frac{[(x_1+x_2-7)^2+(x_3+x_4-7)^2-62]^2 + [(x_1+x_2-7)^2-0.5(x_1+x_2-7)-0.5(x_3+x_4-7)-1.5]^2}{100} + 10 & \text{if } (x_1+x_2)^2+(x_3+x_4)^2 > 80, \\ 0 & \text{if } (x_1+x_2)^2+(x_3+x_4)^2 \leq 80. \end{cases}$$

483 Suppose now we are interested in identifying the limits of the cutting process (i.e. the input
 484 boundary that results in $y(x) = 0$) over the 4D hypercube domain $D = [0, 7]^4$ and 10^4 grid
 485 points are superimposed over D . The prediction error rates for the methods based on GASP
 486 and SVM with two sampling strategies are given in [Figure 6](#), where the remaining simulation
 487 setups are the same as those given in [Subsection 4.1](#). As can be seen from [Figure 6](#), the
 488 error rate of the proposed strategy based on active learning drops down stably to zero as the
 489 experimental budget increases. It is no more than 5% when the total number of experimental
 490 trials exceeds 50. On the other hand, the error rate based on GASP with the sequential
 491 sampling strategy appears to be significantly larger and the decay is relatively slow. The
 492 estimation error rate remains high (above 17%), even when the number of experimental trials
 493 is increased to 80. Interestingly, estimation does not benefit much by using the sequential
 494 sampling strategy based on contour extraction. In summary, the strategy of utilizing active
 495 learning with SVM (i) performs slightly better than that of utilizing passive learning (i.e.
 496 sampling at once by utilizing UD); and (ii) significantly outperforms both sampling strategies
 497 based on GASP.

498 We next examine the performance of the proposed method for larger-scale systems. In
 499 order to reduce the computational burden increased by high dimensionality, here we provide
 500 some useful implementation guidelines.

501

502 *Guidelines for Conducting Experiments for Large-scale Systems*

503

504 • Consider a set of rather “coarse grid” Z points over the experimental domain D . For ex-
 505 ample, here we consider 10 grid points (0.7, 1.4, 2.1, 2.8, 3.5, 4.2, 4.9, 5.6, 6.3, 7.0) for each
 506 dimension x_i . Thus, there are 10^4 grid points in D . A quick numerical examination shows
 507 that there are 2,082 grid points with the output $y(x) = 0$.

508

509 • Select a simpler version of [\(3.3\)](#) to implement the quasi-Monte Carlo method. For example,
 510 by choosing $K = 0$ in [\(3.3\)](#), we consider “uniformity” merely through one subregion $D_1(z)$ (i.e.
 511 the hyper-rectangular subregion from 0 to a grid point z) instead of through 2^K subregions.

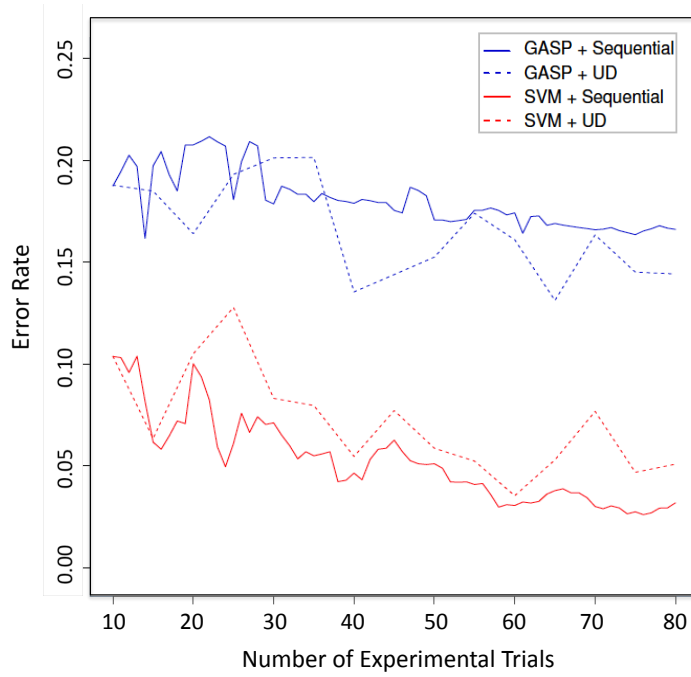


Figure 6. The error rates for estimating the 4D laser cutting boundaries based on passive and active learning with SVM (red lines) and the GASP model for contour estimation (blue lines). The solid lines refer to the results by using the sequential sampling strategy; the dashed lines refer to the results by sampling at once based on UD.

512 That is, the NUD minimizes

$$513 \quad (4.2) \quad \widehat{WCCD}_{f_s, p}(n, \mathcal{P}) = \left\{ \frac{1}{|Z|} \sum_{z \in Z} \left| \frac{N(D_1(z), \mathcal{P})}{n} - F(D_1(z)) \right|^p \right\}^{1/p},$$

514 where $\mathcal{P} \in Z(n) \subset Z$, $|Z|$ represents the cardinality of Z , and $F(D_1(z)) = \sum_{g \in D_1(z)} f_s^*(g)$ is
 515 the proportion of points expected to be allocated in $D_1(z)$, where $\sum_{z \in Z} f_s^*(z) = 1$ (summa-
 516 tion of all normalized weights is equal to one). This simplified measure is the hybrid of the
 517 so-called L_p -discrepancy and F -discrepancy [9], based on which the optimal design requires
 518 significantly less computation load. For simplicity, we can choose $p = 1$ in (4.2).

519

520 • To further reduce the computation load of finding the best new design point, we can im-
 521 plement the accelerated algorithm shown in Appendix A. The idea is to augment on the
 522 lower-dimensional design and increase one dimension each step to find the location of a new
 523 design point.

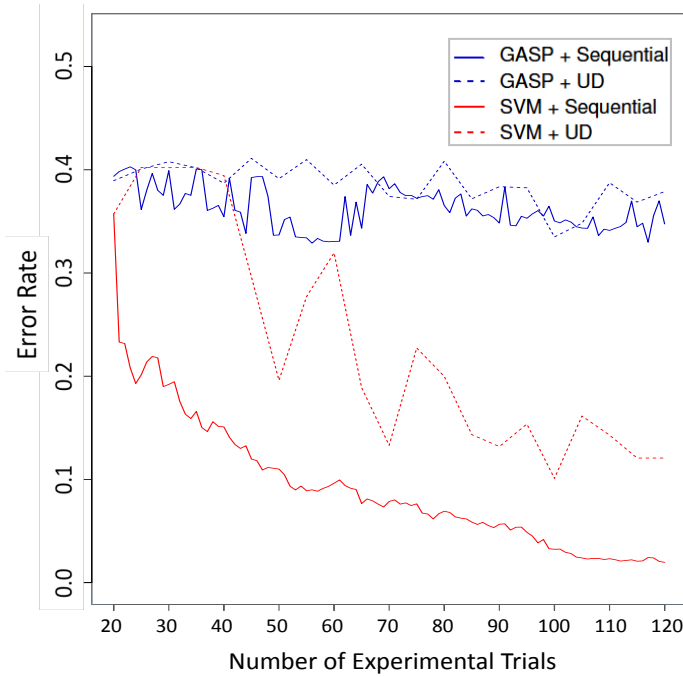
524

525 Let us now consider the laser cutting problem with eight control factors and the experiment

526 is emulated by the following functional relationship

$$\begin{aligned}
 & y(x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8) \\
 & = I\{(\sum_{i=1}^4 x_i)^2 + (\sum_{i=5}^8 x_i)^2 > 120\} \times \\
 (4.3) \quad & \frac{1}{100} \left\{ [(\sum_{i=1}^4 x_i - 7)^4 + (\sum_{i=5}^8 x_i - 7)^2 - 62]^2 + [(\sum_{i=1}^4 x_i - 7)^2 - 0.5(\sum_{i=1}^4 x_i - 7) - 0.5(\sum_{i=5}^8 x_i - 7) - 1.5]^2 \right\} + 10.
 \end{aligned}$$

530 Analogously, we are interested in identifying the limits of the cutting process (i.e. the input
 531 boundary that results in $y(x) = 0$) over the 8D hypercube domain $D = [0, 3]^8$. To accelerate
 532 the implementation of the proposed procedure by using the above guidelines, we consider 3
 533 coarse grid points (1.0, 2.0, 3.0) for each input dimension x_i . Thus, there are 6, 561 grid points
 534 superimposed over D , among which 2, 299 grid points have the output value $y(x) = 0$. Further,
 535 the quasi-Monte Carlo method is implemented by using a simpler version of the discrepancy
 536 measure (say, (4.2) with $p = 1$) and the accelerated search algorithm shown in Appendix A,
 537 with an initial design of size $n_0 = 20$. The prediction error rates for the methods based on
 GASP and SVM with two different sampling strategies are given in Figure 7.



538 **Figure 7.** The error rates for estimating the 8D laser cutting boundaries based on passive and active
 539 learning with SVM (red lines) and the GASP model for contour estimation (blue lines). The solid lines refer
 540 to the results by using the sequential sampling strategy; the dashed lines refer to the results by sampling once
 541 based on UD.
 542
 543

538 It can be seen that Figure 7 reveals similar patterns to those shown in Figure 6. The
 539 error rate of the proposed method based on active learning and SVM drops down quickly to
 540 zero, as the experimental budget increases. It is no more than 5% when the total number of
 541 experimental trials exceeds 100. On the other hand, the error rates based on GASP with both
 542 sampling strategies appear to be significantly larger and their decay is rather slow. It is also
 543

544 worth noting that for the method based on SVM, active learning significantly outperforms
545 passive learning (i.e. sampling at once based on UD) in terms of both estimation efficiency
546 and stability.

547 **5. Conclusion and Discussion.** Although the GASP model is particularly popular for
548 output analysis of computer experiments, it nevertheless exhibits certain drawbacks in ad-
549 dressing the problem under study: (i) it cannot predict well outputs exhibiting sharp changes
550 in nearby locations due to its underlying assumption of continuity and differentiability [19];
551 (ii) it may encounter serious computational issues (i.e. parameter estimation) for large designs
552 [30]; (iii) the associated design methods may not provide efficient allocations for estimating
553 the desired outputs. In this article, we propose a strategy for locating infinite discontinuities
554 in computer experiments leveraging ideas from active learning. The strategy formulates the
555 problem as a binary classification one and employs active learning with an SVM classifier and
556 UD to predict the output discontinuities and allocate the experimental trials in a more effec-
557 tive manner. Numerical results show that, given a limited experimental budget, the proposed
558 strategy is superior to the modeling strategy based on GASP, as well as strategies based on
559 a pure uniform design (i.e. with all experimental trials allocated at once by UD) in terms of
560 both efficiency and prediction accuracy.

561 It should be mentioned that the sampling technique here is different from the conventional
562 methods in active learning, for which the initial training set (for passive learning) is usually
563 obtained by random sampling or LHD. In general, the quality and cost of the initially selected
564 training set in active learning are beyond consideration. However, we believe that for computer
565 experiments significant attention needs to be paid to the selection of the initial design, since it
566 can take up a large portion of the total experimental budget and thus significantly influences
567 the output analysis. The sampling method employed in our strategy belongs to the class of
568 quasi-Monte Carlo methods and has proven more robust than random sampling and other
569 space filling designs. Further, with the placement of a well-designed weight function, the UD
570 can be automatically incorporated into the sampling process of active learning. For practical
571 implementation, we have also provided an accelerated algorithm so that the solution of UD can
572 be fairly well approximated within a more reasonable computational time. Such acceleration
573 is particularly useful and necessary for designs with a large number of input factors. Finally,
574 another learning technique called HMLSVMs is introduced so as to better deal with infinite
575 valued discontinuities that lead to piecewise linear shapes.

576 We believe that the proposed strategy can be extensively used to identify other types of
577 output discontinuity or non-differentiability in computer or other physical experiments. For
578 example, to identify (finite) *jump discontinuities* one can examine the magnitude of response
579 $y(x)$ in the neighborhood of each input point x ; to identify pathological structures like *kinks*
580 and *cusps* (which are essential in manufacturing processes and computer graphics) one may
581 examine the “gradients” of response in the neighborhood of each input point x and determine
582 suitable class labels (note that the gradients for kinks and cusps are not continuous). However,
583 these problems will require a more delicate/complicated formulation.

584 **Appendix A. Accelerated Algorithm for Finding an NUD.**

585 The accelerated algorithm for constructing an NUD is shown as follows.

586 The following theorem describes how much computation can be saved by [Algorithm A.1](#)

Algorithm A.1 Acceleration for Finding an NUD

- 1: Consider a set of superimposed grid points Z over D . Select $\psi_0 \subset \{1, \dots, K\}$, find the optimal design \mathcal{P}_0^* in the selected $|\psi_0|$ -dimensional space (e.g. by using the exhaustive search or Switching Algorithm). Set $i = 0$.
 - 2: **while** $|\psi_{i+1}| < K$ **do**
 - 3: Select $\psi_{i+1} \subset \{1, \dots, K\} \setminus \psi_i$, set $\psi_{i+1} = \psi_{i+1} \cup \psi_i$.
 - 4: Find the optimal design \mathcal{P}_{i+1}^* in $Z^{\psi_{i+1}}$ by augmenting the design \mathcal{P}_i^* in the lower dimensional input domain Z^{ψ_i} .
 - 5: Set $i = i + 1$.
 - 6: **end while**
 - 7: **return** \mathcal{P}_{i+1}^*
-

587 for finding an NUD.

588 **Theorem A.1 (Computational time of Algorithm A.1).** *If the number of superimposed grid*
 589 *points in each dimension of D is taken to be $N^{\frac{1}{K}}$ and $|\psi_i| = d(i+1)$ for each iteration i (where*
 590 *d is divisible by K), the accelerated exhaustive search has the computational time $O(N^{\frac{nd}{K}})$,*
 591 *while the accelerated Switching Algorithm has the computational time up to $O(N^{\frac{d}{K}(2+p)})$.*

Proof. Note that $|\psi_i| = d(i+1)$ indicates that at each iteration i , an additional d -dimensional subspace of D is aggregated with the current design space so as to find the optimal design \mathcal{P}_i^* . Therefore, for fixed values of K , d and n , the accelerated exhaustive search requires the computational time

$$\binom{K}{d} \binom{N^{\frac{d}{K}}}{n} + \binom{K-d}{d} \binom{N^{\frac{d}{K}}}{n} + \dots + \binom{d}{d} \binom{N^{\frac{d}{K}}}{n} = O(N^{\frac{nd}{K}}).$$

On the other hand, the accelerated Switching Algorithm requires the computational time up to

$$(K/d) \cdot O(N^{\frac{d}{K}(2+p)}) = O(N^{\frac{d}{K}(2+p)}).$$

592 Note that in practice the accelerated algorithm would have a dramatic reduction on the
 593 computational time of finding an NUD (though there may be a minor loss of uniformity in the
 594 primary space). To see this, let us recall the early example with $K = 3$, $d = 1$, $n = 6$, $p = 2$
 595 and $N = 10^3$. The accelerated exhaustive search requires merely 10^6 computations for finding
 596 an NUD, while the search requires 10^{18} computations without acceleration. The accelerated
 597 Switching Algorithm requires up to 10^4 computations for finding an NUD, while the search
 598 requires up to 10^{12} computations without acceleration.

599 To evaluate the numerical performance of **Algorithm A.1**, we consider a simple design area
 600 $D = [0, 1] \times [0, 1]$. Suppose there are 100 candidate grids superimposed over D (so $100^{1/2} = 10$
 601 grids for each dimension), and we choose $d = 1$ and $p = 2$. The CCD are then computed by
 602 utilizing (3.3) (with F chosen as an identity function) based on the four algorithms: Exhaustive
 603 Search, Switching, Accelerated Exhaustive Search, and Accelerated Switching. The results
 604 for $n = 1, \dots, 10$ are shown in **Figure 8**. Note that the CCD for the Switching algorithm
 605 is computed by taking the average of 1,000 initial designs based on random sampling, while

606 the CCD for the Accelerated Switching algorithm is computed by taking the average of all
 607 solutions based on a set of NUDs found in the first dimension (i.e. a set of \mathcal{P}_1^*). As can
 608 be seen, for this small design area both the Switching and Accelerated Switching algorithms
 609 approximate very well the optimal design based on Exhaustive Search. This provides a strong
 numerical evidence that the proposed accelerated algorithm is adequate for finding an NUD.

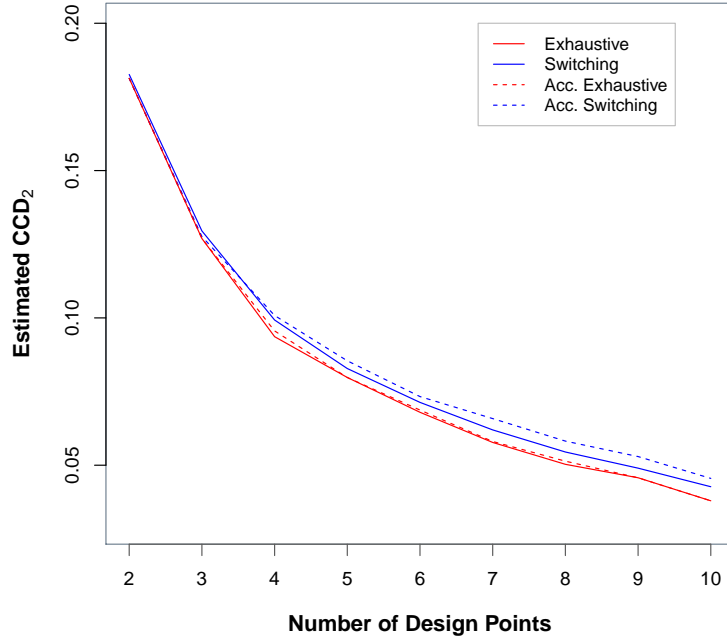


Figure 8. The estimated CCD (with $p = 2$) for the Exhaustive, Switching, Accelerated Exhaustive and Accelerated Switching algorithms with respect to the number of allocated experimental trials. Here the input domain is $D = [0, 1]^2$.

610

611 We next consider the laser cutting example in [Subsection 4.2](#) with an 8D hypercube domain
 612 $D = [0, 3]^8$, while 3 grid points (0.5, 1.5, 2.5) are considered for each input dimension so as
 613 to evaluate the performance of the above algorithms. Thus, there are $N = 3^8 = 6,561$ grid
 614 points superimposed over D . To accelerate the computation of [Algorithm A.1](#), we choose
 615 $d = 2$ (augmenting two dimensions one at a time to find the optimal design), $K = 0$ and
 616 $p = 1$ in [\(3.3\)](#), i.e., a simplified version of WCCD with merely one subregion $D_1(z)$ and
 617 the L_1 -norm distance is considered. Since now the design is unweighted (i.e. F is chosen
 618 as an identity function), it is clear that $F(D_1(z)) = |D_1(z)|/|D|$, which corresponds to the
 619 proportion of the volume taken up by the subregion $D_1(z)$. Under such setup, the WCCD
 620 reduces to the form of so-called L_1 -discrepancy. [Figure 9](#) yields the values of L_1 -discrepancy
 621 associated with the number of design points allocated in D based on the above algorithms. As
 622 can be seen, for this 8D domain with coarse grid points, the trend of discrepancy measure is
 623 somewhat different from that shown in [Figure 8](#) - it does not drop down monotonically as we

624 increase the number of design points. However, numerical values show that both the Switching
 625 and the Accelerated Switching algorithms approximate fairly well the optimal design based
 626 on Exhaustive Search. The empirical evidence renders strong support that our proposed
 accelerated algorithm performs well in finding an NUD for larger scale systems.

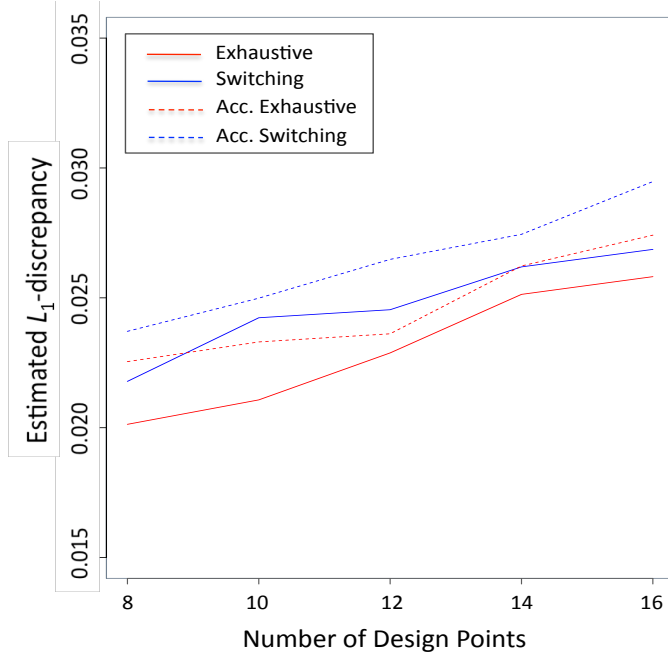


Figure 9. The estimated L_1 -discrepancy for the Exhaustive, Switching, Accelerated Exhaustive and Accelerated Switching algorithms with respect to the number of allocated design points. Here the input domain is $D = [0, 3]^8$.

627

628 **Appendix B. A sequential sampling strategy based on GASP and contour estimation.**

629

630 Suppose one would like to estimate a contour $S(a) = \{x : y(x) = a\}$ (the set of input
 631 points with the response value a) based on the GASP model, a useful strategy is to start with
 632 a relatively small experimental design (called initial design) and then sequentially allocate
 633 experimental trials near the estimated contour and/or input locations with large prediction
 634 uncertainty [44]. Specifically, consider the following *improvement function*

635 (B.1)
$$I(x) = \epsilon^2(x) - \min\{(y(x) - a)^2, \epsilon^2(x)\},$$

636 where $\epsilon(x) = \alpha s(x)$ for some $\alpha > 0$ and $y(x) \sim N(\hat{y}(x), s^2(x))$. Here $\hat{y}(x)$ and $s^2(x)$ represent
 637 the best linear predictor and the associated mean squared error at any non-sampled point
 638 x , respectively. Therefore, at each design stage one tends to select an untried point that
 639 maximizes $I(x)$ over the input space. However, due to the uncertainty of un-sampled design
 640 points, the expected value of the improvement function is considered instead. The follow-
 641 ing mathematical expression, which was given in [7], provides a simpler way for calculating

642 $E[I(x)]$:

$$643 \quad E[I(x)] = [\alpha^2 s^2(x) - (\hat{y}(x) - a)^2 - s^2(x)] \cdot [\Phi(u_2) - \Phi(u_1)] \\ 644 \quad (B.2) \quad + s^2(x) [u_2 \phi(u_2) - u_1 \phi(u_1)] + 2s(x)(\hat{y}(x) - a) [\phi(u_2) - \phi(u_1)],$$

645 where $u_1 = (a - \hat{y}(x) - \alpha s(x))/s(x)$, $u_2 = (a - \hat{y}(x) + \alpha s(x))/s(x)$, $\Phi(\cdot)$ and $\phi(\cdot)$ represent the
646 cdf and pdf of a standard normal random variable, respectively. Thus, in the standard contour
647 estimation problem a not yet tested point x is chosen to maximize (B.2) at each design stage.

648 Note that the goal is to locate the ‘‘boundary’’ for the set of input points x with $y(x) = 0$,
649 which is different from identifying the contour line of height $a = 0$ associated with a contin-
650 uous response surface. Thus, we need to modify accordingly the above sequential sampling
651 procedure. Suppose there are N_j available experimental trials at a particular design stage j
652 and we denote by x_1, \dots, x_{N_j} and $y(x_1), \dots, y(x_{N_j})$ the input locations and their associated
653 outputs, respectively. After fitting the GASP model based on the N_j experimental trials, the
654 estimated contour (i.e. target boundary) is given by

$$655 \quad (B.3) \quad \partial \tilde{T}_j(x) = \{x \in D : 0 < \hat{y}(x) \leq a_j\},$$

656 where

$$657 \quad (B.4) \quad a_j = \frac{1}{2} \cdot \min\{y(x_i) : y(x_i) > 0, i = 1, \dots, N_j\}.$$

658 We next introduce a *stage-dependent improvement function*

$$659 \quad (B.5) \quad I_j(x) = \epsilon^2(x) - \min\{(y(x) - a_j)^2, \epsilon^2(x)\},$$

660 and thus the next-stage design point is chosen as

$$661 \quad (B.6) \quad x^* = \arg \max_{x \notin \{x_1, \dots, x_{N_j}\}} E[I_j(x)],$$

662 where the computation of $E[I_j(x)]$ can be done by using (B.2).

663 Since $y(x) = 0$ for all $x \notin T(x)$, the sampling strategy based on the concept of improve-
664 ment function tends to over-allocate experimental trials outside the target region (this is also
665 validated by numerical investigations). Unfortunately, certain portions of such trials are not
666 particularly informative for estimating the boundary $\partial T(x)$ (e.g. data far away from the
667 boundary). The goal of utilizing a_j instead of 0 in (B.5) is to reduce the amount of over-
668 sampled data outside the target region $T(x)$. Intuitively, the approximation of $\partial T(x)$ should
669 then be done in a more efficient manner. Now we proceed to prove that the estimated $\partial \tilde{T}_j(x)$
670 converges to the true boundary $\partial T(x)$ as $N_j \rightarrow \infty$. It is clear that the GASP model is not
671 able to fit well the primary output surface $y(x)$, which is neither continuous nor differentiable
672 on $\partial T(x)$. However, as we learn that $y(x)$ becomes continuous on $\partial T(x)$ after a reciprocal
673 transformation (though it is still not differentiable), it suffices to show [44]:

$$674 \quad (B.7) \quad \lim_{j \rightarrow \infty} \sup_{x \in D \setminus \partial T(x)} E[I_j(x)] = 0.$$

Let us now break the input domain D into two parts: one for $x \in T(x)$ and the other for $x \notin T(x)$ (both not including the boundary $\partial T(x)$). Since the output surfaces $y(x)$ for both parts are continuous and differentiable, the proof of Theorem 1 in [44] yields

$$\sup_{x \in D \setminus \partial T(x)} s^2(x) \rightarrow 0 \text{ as } N_j \rightarrow \infty,$$

and thus

$$\sup_{x \in D \setminus \partial T(x)} I_j(x) \rightarrow 0 \text{ as } N_j \rightarrow \infty.$$

675 This then implies (B.7) and the proof of convergence is complete. The detailed steps of the
676 above procedure based on a set of superimposed grid points Z are shown in Algorithm B.1.

Algorithm B.1 Sampling based on contour estimation and GASP

- 1: Determine an appropriate experimental domain D associated with a set of superimposed grid points Z . Set the design stage $j = 1$ and determine the size of the initial design n_1 . Allocate an adequate initial design \mathcal{P}_1^* (e.g. by UD) based on Z and obtain the experimental outcomes $y(z)$ for all $z \in \mathcal{P}_1^*$. Compute a_1 by (B.4) and fit the GASP model based on all $z \in \mathcal{P}_1^*$. Compute $E[I_1(z)]$ by (B.2) for all $z \in Z \setminus \mathcal{P}_1^*$.
 - 2: **while** the stopping criterion is not met **do**
 - 3: Select the next-stage design point z^* which maximizes $E[I_j(z)]$ over the set $Z \setminus \mathcal{P}_j^*$.
 - 4: Obtain the experimental outcome for z^* and set $\mathcal{P}_{j+1}^* = \mathcal{P}_j^* \cup \{z^*\}$.
 - 5: Compute a_{j+1} by (B.4) and fit the GASP model based on all $z \in \mathcal{P}_{j+1}^*$.
 - 6: Set $j = j + 1$.
 - 7: **end while**
 - 8: **return** $\partial \bar{T}_{j+1}(z)$
-

677 **Appendix C. Active Learning with HMLSVMs.**

678 We briefly introduce how to develop a learner based on HMLSVMs by utilizing the idea of
679 classification tree. Let $[I]$ be a node of the hierarchical tree and denote the root node by $[0]$,
680 the left child node of $[I]$ by $[I]^+$, the right child node of $[I]$ by $[I]^-$, respectively. Let $g_{[I]}(x)$
681 be the decision function of fitting an LSVM at node $[I]$ and denote the associated “order of
682 split” in the hierarchical tree by i . We simply write $g_{[I]}(x) = g^{(i)}(x)$ and, it is clear that
683 $g_{[0]}(x) = g^{(1)}(x)$. For the classification tree having m splits, we denote the set of all decision
684 functions by $F^m = \{g^{(1)}(x), \dots, g^{(m)}(x)\}$. We next introduce how to split the nodes so as to
685 grow the decision tree.

686 Denote the set of training data associated with a node $[I]$ by $\mathcal{T}_{[I]}$, where $\mathcal{T}_{[0]} = \{(y_1, x_1),$
687 $\dots, (y_n, x_n)\}$. The training data associated with the left child and right child node of $[I]$
688 are then denoted by $\mathcal{T}_{[I]^+} = \{x \in \mathcal{T}_{[I]} : g_{[I]}(x) > 0\}$ and $\mathcal{T}_{[I]^-} = \{x \in \mathcal{T}_{[I]} : g_{[I]}(x) < 0\}$,
689 respectively. In order to best split the training data in a node, we need to define the so-called
690 *host training center* and *guest training center*. Given the current tree with m splits, denote
691 the set of training samples having “wrong classifications” by \mathcal{W} . For any $x_i \in \mathcal{T}_{[0]}$, denote the
692 set of its k “nearest” training samples with the same labels as y_i by $N_k(x_i)$. The host training

693 center is then defined as

$$694 \quad (C.1) \quad x_h = \arg \max_{x_i \in \mathcal{W}} \sum_{x_j \in N_k(x_i)} \exp \left\{ -\frac{\|x_j - x_i\|^2}{\sigma^2} \right\} E_j,$$

695 where $E_j = \min \{1, \max\{0, 1 - y_j \cdot G^{(m)}(x_j)\}\}$ and σ^2 is a positive constant. By this defi-
696 nition, the selected host training center appears to have the largest density of misclassified
697 training samples over its k nearest neighbors. Suppose now $x_h \in \mathcal{T}_{[I]}$ for some node $[I]$, the
698 guest training center is then defined as

$$699 \quad (C.2) \quad x_g = \arg \max_{x_i \in \bar{N}_k(x_h)} \sum_{x_j \in N_k(x_i)} \exp \left\{ -\frac{\|x_i - x_j\|^2 + \|x_j - x_h\|^2}{\sigma^2} \right\},$$

700 where $\bar{N}_k(x_h)$ represents the set of x_h 's k nearest training samples, with different labels from
701 y_h . In short, the guest training center is a training sample that is close to the host training
702 center and at the same time has a large density of correctly classified training samples over its
703 k nearest neighbors. Note that such setup is quite similar to the concept of ‘‘local likelihood’’,
704 while the goal is to split node $[I]$ based on the two selected centers x_h and x_g .

705 To develop the decision function to split node $[I]$, we first define the weights of the training
706 samples in $N_k(x_h)$ and $N_k(x_g)$ as

$$707 \quad (C.3) \quad W_i = \begin{cases} \exp \left\{ -\frac{\|x_i - x_h\|^2}{\sigma^2} \right\}, & \text{if } x_i \in N_k(x_h); \\ \exp \left\{ -\frac{\|x_i - x_g\|^2}{\sigma^2} \right\}, & \text{if } x_i \in N_k(x_g). \end{cases}$$

708 As can be seen, the weight of each training sample in both sets becomes exponentially larger
709 as it gets closer to x_h or x_g and exponentially smaller as it gets farther away. The WLSVM
710 is then used to develop a decision function $g_{N_k}(x)$ based on the two sets $N_k(x_h)$ and $N_k(x_g)$.
711 Once we have $g_{N_k}(x)$, the decision function for all training samples in node $[I]$ is obtained by
712 mixing $g_{N_k}(x)$ with the decision function of its parent node $[I_p]$, that is,

$$713 \quad (C.4) \quad g_{[I]}(x) = \begin{cases} \max\{g_{[I_p]}(x), g_{N_k}(x)\}, & \text{if } [I] = [I_p]^-; \\ \min\{g_{[I_p]}(x), g_{N_k}(x)\}, & \text{if } [I] = [I_p]^+. \end{cases}$$

714 Denote by $g_{[I]}(x) = g^{(m+1)}(x)$, then after the $(m+1)$ -th split, the set of all decision functions
715 for the tree is given by

$$716 \quad (C.5) \quad F^{m+1} = F^m \cup \{g^{(m+1)}(x)\}.$$

Denote by $F_s^m = \{g_s^{(1)}(x), \dots, g_s^{(m)}(x)\}$ the set of all decision functions obtained by
HMLSVMs at a particular learning stage s , where

$$g_s^{(i)}(x) = \text{sign} \left(\langle w_s^{(i)}, x \rangle + b_s^{(i)} \right)$$

717 is the linear decision function for the i -th split at stage s , $i = 1, \dots, m$. Then, for each input
718 point x we can analogously compute the distance (3.9) and the weight (3.10) so that the

719 weighted UD can be utilized to sample an instance at the next stage. Note that due to the
 720 structure of the developed decision tree, the computation of (3.9) can be intricate as the size
 721 of the tree becomes large - even when the input dimension is low. However, for some cases
 722 where the estimated target regions have a “convex” shape, the computation can be done in a
 723 simpler way. To see this, for any input point x inside the estimated target region $\tilde{T}_s(x)$, its
 724 distance to the boundary can be directly computed by

$$725 \text{ (C.6)} \quad d(x, \partial\tilde{T}_s(x)) = \min_{i \in \{1, \dots, m\}} \frac{|\langle w_s^{(i)}, x \rangle + b_s^{(i)}|}{\|w_s^{(i)}\|}.$$

726 If an input point x is outside the estimated target region, one will have to identify which
 727 facet/vertex of $\partial\tilde{T}_s(x)$ is closest to x and then compute the distance accordingly. We summa-
 728 rize the detailed steps of our framework for estimating the target region with piecewise linear
 729 boundaries in Algorithm C.1.

730

Algorithm C.1 Active Learning with HMLSVMs

- 1: Determine an appropriate training domain D associated with a set of superimposed grid points Z . Set the learning stage $s = 1$ and determine the size of the initial training set n_1 . Find the initial training set \mathcal{P}_1^* of size n_1 that minimizes (3.3) based on Z , where $p = 2$ and the weight function is initially placed as $f_0(x) \equiv 1$. Train an HMLSVMs classifier based on \mathcal{P}_1^* with the weights $W_i \equiv 1$ in the root node [0] and (C.3) for later splits, obtain the decision boundary $\partial T_1(z)$ based on the set of resulting decision functions F_1^m .
 - 2: **while** the stopping criterion is not met **do**
 - 3: Update the weight functions $f_s(x)$ and $f_s^*(x)$ by (3.10) and (3.11), set $s = s + 1$.
 - 4: Select a sample $z^* \in Z \setminus \mathcal{P}_{s-1}^*$ so that $\mathcal{P}_s^* = \mathcal{P}_{s-1}^* \cup \{z^*\}$ minimizes (3.3). Query the label of z^* .
 - 5: Train an HMLSVMs classifier based on \mathcal{P}_s^* with the weights $W_i \equiv 1$ in the root node [0] and (C.3) for later splits, obtain the decision boundary $\partial\tilde{T}_s(z)$ based on the set of resulting decision functions F_s^m .
 - 6: **end while**
 - 7: **return** $\partial\tilde{T}_s(z)$
-

731 Figure 10 yields the decision tree and estimated stability region for the queueing system
 732 shown in Figure 1 (with $\Delta = 0$) based on HMLSVMs with two splits, where the host center
 733 x_h and guest center x_g for the second split were computed based on two nearest samples. To
 734 evaluate the performance of Algorithm C.1, the prediction error rate along with the design
 735 stage is shown in Figure 11. Note that the prediction error rate here is calculated based on
 736 $40 \times 40 = 1,600$ grid points and the control process $\beta(s) = 100 \times (1.3)^{-s}$, while for comparison
 737 purpose the result based on the GASP model and Algorithm B.1 is also attached. As can be
 738 seen, for such a small system the prediction error rate of our proposed method is below 2%
 739 as the number of experiments exceeds 27 and remains stable after that. On the other hand,
 740 the prediction error rate of the GASP model appears to have a relatively slow decay by using
 741 the same number of experimental trials.

742

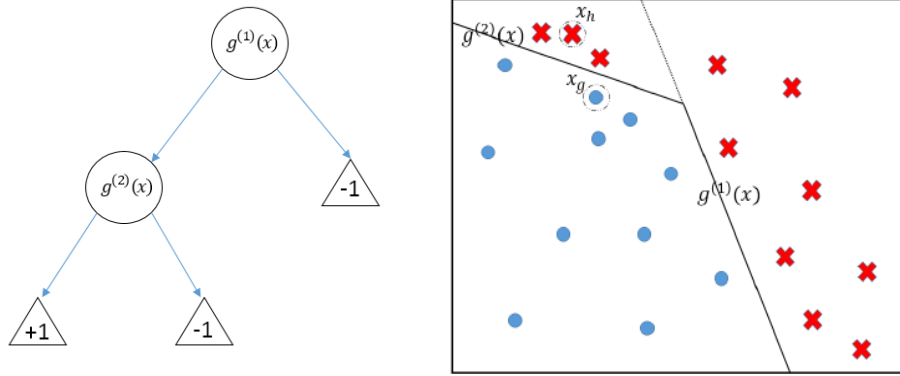


Figure 10. Illustration of the decision tree based on HMLSVMs with two splits (left panel) and the estimated stability region for the queueing system shown in Figure 1 (with $\Delta = 0$) along with the host center x_h and guest center x_g by choosing $k = 2$ (right panel). Here the red symbol “ \times ” represents the objects with a true class label “ -1 ”, while the blue symbol “ \circ ” represents objects with a true class label “ $+1$ ”.

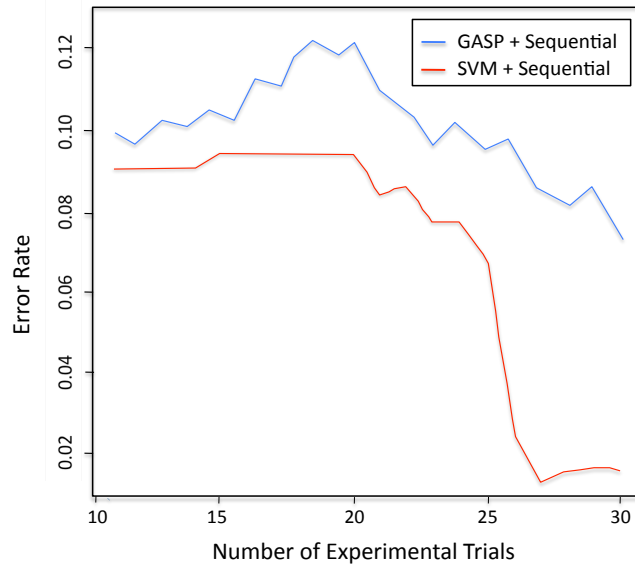


Figure 11. The error rates for estimating the piecewise linear boundaries based on passive and active learning with SVM (red line) and the GASP model for contour estimation (blue line). Note that both methods here utilize the associated sequential sampling strategy.

743 *Remark:* If we know a priori that the true target region is *convex*, one would not expect the
 744 decision tree to end up with a non-convex shape. To possibly avoid this issue, one can either
 745 (i) assign larger weights W_i to the training samples with label $y(x) = +1$ (i.e. points inside
 746 the true target region) when fitting the LSVM model or (ii) assign larger weights $f_s(x)$ to the
 747 input points inside the estimated target region so as to allocate the experimental trials based
 748 on WCCD. Intuitively, such strategy will push the decision boundary at each node of the tree
 749 to move “outward”, thus retaining as much as possible the convex shape of the estimated

750 target region.

751

752 **Acknowledgement.** The authors would like to thank Chia-Li Lin and Ruby Chiu-Hsing
753 Weng for their technical support and assistance with the work.

754

REFERENCES

- 755 [1] M. ARMONY AND N. BAMBOS, *Queueing dynamics and maximal throughput scheduling in switched*
756 *processing systems*, Queueing Systems, 44 (2003), pp. 209–252, [https://doi.org/10.1023/A:](https://doi.org/10.1023/A:1024714024248)
757 [1024714024248](https://doi.org/10.1023/A:1024714024248).
- 758 [2] J. BECK AND S. GUILLAS, *Sequential design with mutual information for computer experiments (mice):*
759 *Emulation of a tsunami model*, SIAM/ASA Journal on Uncertainty Quantification, 4 (2016), pp. 739–
760 766, <https://doi.org/10.1137/140989613>.
- 761 [3] B. E. BOSER, I. M. GUYON, AND V. N. VAPNIK, *A training algorithm for optimal margin classifiers*,
762 in Proceedings of the 5th Annual ACM Workshop on COLT, 1992, pp. 144–152, [https://doi.org/10.](https://doi.org/10.1145/130385.130401)
763 [1145/130385.130401](https://doi.org/10.1145/130385.130401).
- 764 [4] M. BRAMSON, *Stability of Queueing Networks*, vol. 5, Springer Verlag, Berlin, 2008, pp. 169–345, [https:](https://doi.org/10.1214/08-PS137)
765 [//doi.org/10.1214/08-PS137](https://doi.org/10.1214/08-PS137).
- 766 [5] C. C. CHANG AND C. J. LIN, *LIBSVM: A library for support vector machines*, ACM Transactions on
767 Intelligent Systems and Technology, 2 (2011), pp. 27:1–27:27. Software available at [http://www.csie.](http://www.csie.ntu.edu.tw/~cjlin/libsvm)
768 [ntu.edu.tw/~cjlin/libsvm](http://www.csie.ntu.edu.tw/~cjlin/libsvm).
- 769 [6] R. B. CHEN, Y. W. HSU, Y. HUNG, AND W. C. WANG, *Discrete particle swarm optimization for con-*
770 *structing uniform design on irregular regions*, Computational Statistics & Data Analysis, 72 (2014),
771 pp. 282–297, <https://doi.org/10.1016/j.csda.2013.10.015>.
- 772 [7] R. B. CHEN, Y. C. HUNG, W. C. WANG, AND S. W. YEN, *Contour estimation via two fidelity computer*
773 *simulators under limited resources*, Computational Statistics, 28 (2013), pp. 1813–1834, [https://doi.](https://doi.org/10.1007/s00180-012-0380-7)
774 [org/10.1007/s00180-012-0380-7](https://doi.org/10.1007/s00180-012-0380-7).
- 775 [8] W. J. CHEN, *Instability threshold and stability boundaries of rotor-bearing systems*, Journal of Engineering
776 for Gas Turbines and Power, 118 (1996), pp. 115–121, <https://doi.org/10.1115/1.2816526>.
- 777 [9] S. C. CHUANG AND Y. C. HUNG, *Uniform design over general input domains with applications to target*
778 *region estimation in computer experiments*, Computational Statistics & Data Analysis, 54 (2010),
779 pp. 219–232, <https://doi.org/10.1016/j.csda.2009.08.008>.
- 780 [10] C. CORTES AND V. N. VAPNIK, *Support vector networks*, Machine Learning, 20 (1995), pp. 273–297,
781 <https://doi.org/10.1023/A:1022627411411>.
- 782 [11] H. DETTE AND A. PEPELYSHEV, *Generalized latin hypercube design for computer experiments*, Techno-
783 metrics, 52 (2010), pp. 421–429, <https://doi.org/10.1198/TECH.2010.09157>.
- 784 [12] K. T. FANG AND D. K. J. LIN, *Uniform experimental designs and their applications in industry*, vol. 22,
785 Elsevier, Amsterdam, 2003, pp. 131–170, [https://doi.org/10.1016/S0169-7161\(03\)22006-X](https://doi.org/10.1016/S0169-7161(03)22006-X).
- 786 [13] K. T. FANG, D. K. J. LIN, P. WINKER, AND Y. ZHANG, *Uniform design: Theory and applications*,
787 Technometrics, 42 (2000), pp. 237–248, <https://doi.org/10.2307/1271079>.
- 788 [14] K. T. FANG AND Y. WANG, *Number-theoretic Methods in Statistics*, Chapman and Hall, London, 1994.
- 789 [15] Y. FANG AND K. A. LOPARO, *Stochastic stability of jump linear systems*, IEEE Transactions on Automatic
790 Control, 47 (2002), pp. 1204–1208, <https://doi.org/10.1109/TAC.2002.800674>.
- 791 [16] A. I. J. FORRESTER, A. SÓBESTER, AND A. J. KEANE, *Multi-fidelity optimization via surrogate modelling*,
792 in Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences, vol. 463,
793 Dec. 2007, pp. 3251–3269, <https://doi.org/10.1098/rspa.2007.1900>.
- 794 [17] Y. FREUND, *Boosting a weak learning algorithm by majority*, Information and Computation, 121 (1995),
795 pp. 256–285, <https://doi.org/10.1006/inco.1995.1136>.
- 796 [18] Y. FREUND AND R. E. SCHAPIRE, *A decision-theoretic generalization of on-line learning and an ap-*
797 *plication to boosting*, Journal of Computer and System Sciences, 55 (1997), pp. 119–139, [https:](https://doi.org/10.1006/jcss.1997.1504)
798 [//doi.org/10.1006/jcss.1997.1504](https://doi.org/10.1006/jcss.1997.1504).
- 799 [19] R. B. GRAMACY AND H. H. LEE, *Bayesian treed gaussian process models with an application to computer*
800 *modeling*, Journal of the American Statistical Association, 103 (2008), pp. 1119–1130, [https://doi.](https://doi.org/10.1198/016214508000000689)
801 [org/10.1198/016214508000000689](https://doi.org/10.1198/016214508000000689).
- 802 [20] M. GU AND L. WANG, *Scaled gaussian stochastic process for computer model calibration and prediction*,
803 May 2017, <https://arxiv.org/abs/1707.08215>.
- 804 [21] F. J. HICKERNELL, *A generalized discrepancy and quadrature error bound*, Mathematics of Computation,
805 67 (1998), pp. 299–322, <https://doi.org/10.1090/S0025-5718-98-00894-1>.
- 806 [22] C. M. HUANG, Y. J. LEE, D. K. J. LIN, AND S. Y. HUANG, *Model selection for support vector machines*

- 807 *via uniform design*, Computational Statistics & Data Analysis, 52 (2007), pp. 335–346, <https://doi.org/10.1016/j.csda.2007.02.013>.
- 808
- 809 [23] Y. C. HUNG AND C. C. CHANG, *Dynamic scheduling for switched processing systems with substantial*
810 *service-mode switching times*, Queueing Systems, 60 (2008), pp. 87–109, [https://doi.org/10.1007/](https://doi.org/10.1007/s11134-008-9088-3)
811 [s11134-008-9088-3](https://doi.org/10.1007/s11134-008-9088-3).
- 812 [24] Y. C. HUNG AND G. MICHAILIDIS, *Stability and control of acyclic stochastic processing networks with*
813 *shared resources*, IEEE Transactions on Automatic Control, 57 (2012), pp. 489–494, [https://doi.org/](https://doi.org/10.1109/TAC.2011.2164012)
814 [10.1109/TAC.2011.2164012](https://doi.org/10.1109/TAC.2011.2164012).
- 815 [25] G. IMBENS AND T. LEMIEUX, *Regression discontinuity designs: A guide to practice*, Journal of Econo-
816 *metrics*, 142 (2008), pp. 615–635, <https://doi.org/10.1016/j.jeconom.2007.05.001>.
- 817 [26] M. E. JOHNSON, L. M. MOORE, AND D. YLVIKAKER, *Minimax and maximin distance designs*, Journal
818 *of Statistical Planning and Inference*, 26 (1990), pp. 131–148, [https://doi.org/10.1016/0378-3758\(90\)](https://doi.org/10.1016/0378-3758(90)90122-B)
819 [90122-B](https://doi.org/10.1016/0378-3758(90)90122-B).
- 820 [27] V. R. JOSEPH, *Limit kriging*, Technometrics, 48 (2006), pp. 458–466, [https://doi.org/10.1198/](https://doi.org/10.1198/004017006000000011)
821 [004017006000000011](https://doi.org/10.1198/004017006000000011).
- 822 [28] V. R. JOSEPH, Y. HUNG, AND A. SUDJANTO, *Blind kriging: a new method for developing metamodels*,
823 *Journal of Mechanical Design*, 130 (2008), 0311021 (8 pages), <https://doi.org/10.1115/1.2829873>.
- 824 [29] A. KARATZOGLOU, D. MEYER, AND K. HORNIK, *Support vector machines in r*, Journal of Statistical
825 *Software*, 15 (2006), pp. 1–28, <https://doi.org/10.18637/jss.v015.i09>.
- 826 [30] C. G. KAUFMAN, D. BINGHAM, S. HABIB, K. HEITMANN, AND J. A. FRIEMAN, *Efficient emulators of*
827 *computer experiments using compactly supported correlation functions, with an application to cosmol-*
828 *ogy*, Annals of Applied Statistics, 5 (2011), pp. 2470–2492, <https://doi.org/10.1214/11-AOAS489>.
- 829 [31] S. S. KEERTHI AND C. J. LIN, *Asymptotic behaviors of support vector machines with gaussian kernel*,
830 *Neural Computation*, 15 (2003), pp. 1667–1689, <https://doi.org/10.1162/089976603321891855>.
- 831 [32] T. A. KHAWLI, U. EPELT, AND W. SCHULZ, *Advanced metamodeling techniques applied to multidimen-*
832 *sional applications with piecewise response*, vol. 9432, Springer International Publishing, Switzerland,
833 2015, pp. 93–104, https://doi.org/10.1007/978-3-319-27926-8_9.
- 834 [33] J. KREMER, K. S. PEDERSEN, AND C. IGEL, *Active learning with support vector machines*, Wiley
835 *Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 4 (2014), pp. 313–326,
836 <https://doi.org/10.1002/widm.1132>.
- 837 [34] P. R. KUMAR AND S. P. MEYN, *Stability of queueing networks and scheduling policies*, IEEE Transactions
838 *on Automatic Control*, 40 (1995), pp. 251–260, <https://doi.org/10.1109/9.341782>.
- 839 [35] D. LEE AND T. LEMIEUX, *Regression discontinuity designs in economics*, Journal of Economic Literature,
840 48 (2010), pp. 281–355, <https://doi.org/10.1257/jel.48.2.281>.
- 841 [36] D. LEWIS AND J. CATLETT, *Heterogeneous uncertainty sampling for supervised learning*, in Proceedings
842 *of the Eleventh International Conference on Machine Learning*, July 1994, pp. 148–156, <https://doi.org/10.1016/B978-1-55860-335-6.50026-X>.
- 843
- 844 [37] C. H. LI, C. T. LIN, AND B. C. KUO, *An automatic method for selecting the parameter of the RBF*
845 *kernel function to support vector machines*, in Proceedings of the IEEE International Conference
846 *on Geoscience and Remote Sensing Symposium*, July 2010, pp. 836–839, [https://doi.org/10.1109/](https://doi.org/10.1109/IGARSS.2010.5649251)
847 [IGARSS.2010.5649251](https://doi.org/10.1109/IGARSS.2010.5649251).
- 848 [38] D. K. J. LIN, C. SHARPE, AND P. WINKER, *Optimized u-type designs on flexible regions*, Computational
849 *Statistics & Data Analysis*, 54 (2010), pp. 1505–1515, <https://doi.org/10.1016/j.csda.2010.01.032>.
- 850 [39] B. McDONALD, P. RANJAN, AND H. CHIPMAN, *Gpfit: An r package for fitting a gaussian process model*
851 *to deterministic simulator outputs*, Journal of Statistical Software, 64 (2015), pp. 1–23, <https://doi.org/10.18637/jss.v064.i12>.
- 852
- 853 [40] M. D. MORRIS, T. J. MITCHELL, AND D. YLVIKAKER, *Bayesian design and analysis of computer ex-*
854 *periments: use of derivatives in surface prediction*, Technometrics, 35 (1993), pp. 243–255, <https://doi.org/10.2307/1269517>.
- 855
- 856 [41] K. P. MURPHY, *Machine Learning: A Probabilistic Perspective*, MIT Press, Cambridge, 2012.
- 857 [42] J. PLATT, *Probabilistic outputs for support vector machines and comparison to regularized likelihood*
858 *methods*, vol. 10, MIT Press, Cambridge, 1999, pp. 61–74, <https://doi.org/10.1.1.41.1639>.
- 859 [43] M. T. PRATOLA, O. HARARI, D. BINGHAM, AND G. E. FLOWERS, *Design and analysis of experiments*
860 *on nonconvex regions*, Technometrics, 59 (2017), pp. 36–47, <https://doi.org/10.1080/00401706.2015>.

- 861 1115674.
- 862 [44] P. RANJAN, D. BINGHAM, AND G. MICHAILIDIS, *Sequential experimental design for contour estima-*
863 *tion from complex computer codes*, *Technometrics*, 50 (2008), pp. 527–541, [https://doi.org/10.1198/](https://doi.org/10.1198/004017008000000541)
864 [004017008000000541](https://doi.org/10.1198/004017008000000541).
- 865 [45] J. SACKS, S. B. SCHILLER, AND W. J. WELCH, *Design for computer experiments*, *Technometrics*, 31
866 (1989), pp. 41–47, <https://doi.org/10.2307/1270363>.
- 867 [46] J. SACKS, W. J. WELCH, T. J. MITCHELL, AND H. P. WYNN, *Design and analysis of computer experi-*
868 *ments*, *Statistical Science*, 4 (1989), pp. 409–423, <https://doi.org/10.1214/ss/1177012420>.
- 869 [47] T. J. SANTNER, B. J. WILLIAMS, AND W. I. NOTZ, *The Design and Analysis of Computer Experiments*,
870 Springer Verlag, New York, 2003.
- 871 [48] B. SCHÖLKOPF, *Support Vector Learning*, Oldenbourg Verlag, Munich, 1997.
- 872 [49] J. SHAWE-TAYLOR AND N. CRISTIANINI, *Kernel Methods for Pattern Analysis*, Cambridge University
873 Press, New York, 2004, <https://doi.org/10.1017/CBO9780511809682>.
- 874 [50] M. A. L. THATHACHAR AND P. VISWANATH, *On the stability of fuzzy systems*, *IEEE Transactions on*
875 *Fuzzy Systems*, 5 (1997), pp. 145–151, <https://doi.org/10.1109/91.554461>.
- 876 [51] D. THISTLEWAITE AND D. CAMPBELL, *Regression-discontinuity analysis: An alternative to the ex post*
877 *facto experiment*, *Journal of Educational Psychology*, 51 (1960), pp. 309–317, [https://doi.org/10.](https://doi.org/10.1037/h0044319)
878 [1037/h0044319](https://doi.org/10.1037/h0044319).
- 879 [52] S. TONG AND E. CHANG, *Support vector machine active learning for image retrieval*, in *Proceedings*
880 *of the International Conference on Multimedia (MM)*, 2001, pp. 107–118, [https://doi.org/10.1145/](https://doi.org/10.1145/500141.500159)
881 [500141.500159](https://doi.org/10.1145/500141.500159).
- 882 [53] S. TONG AND D. KOLLER, *Support vector machine active learning with applications to text clas-*
883 *sification*, *Journal of Machine Learning Research*, 2 (2001), pp. 45–66, [https://doi.org/10.1162/](https://doi.org/10.1162/153244302760185243)
884 [153244302760185243](https://doi.org/10.1162/153244302760185243).
- 885 [54] V. N. VAPNIK, *Statistical Learning Theory*, Wiley, New York, 1998.
- 886 [55] D. WANG, X. ZHANG, M. FAN, AND X. YE, *Hierarchical mixing linear support vector machines for non-*
887 *linear classification*, *Pattern Recognition*, 59 (2016), pp. 255–267, [https://doi.org/10.1016/j.patcog.](https://doi.org/10.1016/j.patcog.2016.02.018)
888 [2016.02.018](https://doi.org/10.1016/j.patcog.2016.02.018).
- 889 [56] W. WANG, Z. XU, W. LU, AND X. ZHANG, *Determination of the spread parameter in the gaussian kernel*
890 *for classification and regression*, *Neurocomputing*, 55 (2003), pp. 643–663, [https://doi.org/10.1016/](https://doi.org/10.1016/S0925-2312(02)00632-X)
891 [S0925-2312\(02\)00632-X](https://doi.org/10.1016/S0925-2312(02)00632-X).
- 892 [57] J. WIELAND, R. PASUPATHY, AND B. W. SCHMEISER, *Queueing-network stability: simulation-based*
893 *checking*, in *Proceedings of the Winter Simulation Conference*, 2003, pp. 520–527, [https://doi.org/](https://doi.org/10.1109/WSC.2003.1261464)
894 [10.1109/WSC.2003.1261464](https://doi.org/10.1109/WSC.2003.1261464).
- 895 [58] Z. XU, M. DAI, AND D. MENG, *Fast and efficient strategies for model selection of gaussian support*
896 *vector machine*, *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 39
897 (2009), pp. 1292–1307, <https://doi.org/10.1109/TSMCB.2009.2015672>.
- 898 [59] X. YANG, Q. SONG, AND Y. WANG, *A weighted support vector machine for data classification*, *Journal*
899 *of Pattern Recognition and Artificial Intelligence*, 21 (2007), pp. 961–976, [https://doi.org/10.1142/](https://doi.org/10.1142/S0218001407005703)
900 [S0218001407005703](https://doi.org/10.1142/S0218001407005703).