

LOCAL VORONOI TESSELLATIONS FOR ROBUST MULTI-WAVE CALIBRATION OF COMPUTER MODELS

Wenzhe Xu,^{1,*} Daniel B. Williamson,² & Peter Challenor³

¹Department of Mathematical Sciences University of Exeter, UK, EX4 4PY

²Department of Mathematical Sciences University of Exeter, UK, EX4 4PY and The Alan Turing Institute The British Library, London, UK, NW1 2DB

³Department of Mathematical Sciences University of Exeter, UK, EX4 4PY and The Alan Turing Institute The British Library, London, UK, NW1 2DB

*Address all correspondence to: Wenzhe Xu, Department of Mathematical Sciences, University of Exeter, UK, EX4 4PY, E-mail: wx229@exeter.ac.uk

Original Manuscript Submitted: 17/05/2020; Final Draft Received: mm/dd/yyyy

History matching using Gaussian process emulators is a well-known methodology for the calibration of computer models. It attempts to identify the parts of input parameter space that are likely to result in mismatches between simulator outputs and physical observations by using emulators. These parts are then ruled out. The remaining “Not Ruled Out Yet (NROY)” input space is then searched for good matches by repeating the history matching process. An easily neglected limitation of this method is that the emulator must simulate the target NROY space well, else good parameter choices can be ruled out. We show that even when an emulator passes standard diagnostic checks on the whole parameter space, good parameter choices can easily be ruled out. We present novel methods for detecting these cases and a Local Voronoi Tessellation method for a robust approach to calibration that ensures that the true NROY space is retained and parameter inference is not biased.

KEY WORDS: History matching, Uncertainty quantification, Calibration, Gaussian Processes, Emulator Diagnostics

1 1. INTRODUCTION

2 Computer models typically solve physical equations, such as coupled systems of PDEs in order to learn
3 about features of the real-world. Calibration of computer models broadly involves using partial and imper-
4 fect observations of the real world in order to learn which settings of the model’s input parameters lead to
5 outputs that are consistent with real-world observations given relevant uncertainties such as measurement
6 error and model discrepancy (we define these terms in Section 2.2).

7 When a computer model is inexpensive, it can be embedded in an MCMC or optimization algorithm
8 for calibration directly (see e.g. [1]). However, many computer models are expensive and/or take a long
9 time to run. For example, climate models may take days or weeks of run time on supercomputers [2].
10 When it is not possible to run the model often enough to calibrate directly, a small, carefully chosen, set
11 of model runs, often termed a ‘design’ or ‘ensemble’, can be run and used to construct an ‘emulator’ or
12 ‘surrogate’; an inexpensive statistical model used to approximate the computer model [3,4].

13 Bayesian calibration [5,6] and history matching [7–10] are both extensively used methods for calibrat-
14 ing with an emulator. Bayesian calibration places a probability distribution over a ‘best input’ and updates
15 this distribution using model runs and observations of the real process. Rather than relying on making
16 distributional assumptions, history matching attempts to identify the parts of the input parameter space
17 that are likely to result in mismatches between computer outputs and observations by iteratively remov-
18 ing those regions of parameter space in which we are virtually certain that there are no good matches.
19 Previous research has applied history matching to many fields, including oil reservoir modelling [11,12],
20 epidemiology [13,14] galaxy formation [15–17] climate systems, [1,18–23]. In this paper we focus on history
21 matching. For discussions comparing the two approaches, see [24,25] and the discussion in [15].

22 Beginning with an emulator, history matching uses a distance function to rule out input choices that
23 lead to outputs that are ‘too far’ from observations. The distance, called *implausibility*, is computed ac-
24 cording to a norm that standardises according to all relevant uncertainties, including the uncertainty con-
25 tributed by the emulator. We define this formally in Section 2. Large implausibility regions are ruled out,
26 whilst small implausibility regions may either be good parameter choices or “too uncertain to tell”, hence
27 they are termed “Not Ruled Out Yet” (NROY). Following an initial history match (wave 1), a new design
28 is then constructed within NROY space and more accurate emulators are constructed within this region
29 for the purposes of cutting out further space.

1 A currently unexplored limitation of history matching can occur when the emulator is unable to simu-
 2 late the unknown target NROY space effectively, even if it seems to pass all standard emulator diagnostic
 3 checks [26]. Poor simulation may result in true NROY space being ruled out without any indication for
 4 the analyst that this has occurred. For simulators that are constantly under development, such as climate
 5 models, this could be a costly mistake that causes parameterizations or even computational methods and
 6 hardware to be needlessly revisited, even though the model was already fit for purpose.

7 This paper will discuss factors that contribute towards ruling out good parameter choices and then
 8 will present a novel Local Voronoi Tessellation design that can be used for robust multi-wave calibration
 9 of computer models that ensures the true NROY space is retained without biasing the parameter inference.
 10 Section 2 will present a review of emulation and history matching methodologies. It will further illustrate
 11 these with a numerical example to demonstrate that, even when an emulator validates well on the whole
 12 parameter space, good parameter choices can still be ruled out. Section 3 presents the novel detecting
 13 method we have developed, as well as a local Voronoi Tessellation method for robust history matching.
 14 Section 4 will outline the application of the study methods to two illustrative examples, as well as to the
 15 output of the French climate model, IPSL-CM [27,28].

16 2. EMULATION AND HISTORY MATCHING

17 2.1 Emulation

18 Gaussian process emulators are used to approximate expensive computer simulators whilst quantifying
 19 uncertainty in the approximation [29,30]. A Gaussian Process (GP) is a stochastic process. Any finite num-
 20 ber of random variables from the Gaussian process has a joint Gaussian distribution [31]. Assuming f
 21 represents the complex computer model with input parameters x , an emulator for $f(x)$ can be constructed
 22 by fitting a mean function $m(x)$ and a correlation function $c(x, x')$ so that

$$f(x)|\beta, \sigma^2, \nu, \delta \sim \mathcal{GP}(m(x), \sigma^2 c(x, x'; \delta, \nu)), \quad (1)$$

23 where $m(x) = h^T(x)\beta$, $\beta = (\beta_1, \beta_2, \dots, \beta_k)$ is a vector of unknown regression coefficients, $h(x)$ is a q -
 24 length vector of regression functions, σ^2 is a hyperparameter that controls the scaling of the process, δ is
 25 a vector of correlation length parameters that used to define the correlation function. The nugget term,
 26 ν , is a small number added on the principal diagonal of the design correlation matrix and is often used

1 to represent observation noise, account for uncertainty in inactive inputs or to avoid numerical instability
 2 during computation. The nugget can be specified or it may need to be trained along with the other hyper-
 3 parameters (for example, when inputs are varied in the design, yet not included in the correlation function)
 4 [32]. With a zero nugget, the emulator interpolates the model runs at the training locations.

5 The choice of correlation function, $c(x, x')$, is one of the key elements of the Gaussian Process. Standard
 6 choices include the power exponential correlation function [21] and the Matérn correlation function [31],
 7 both depending on correlation length parameters (which we will refer to as δ throughout), and only on
 8 the distance $|x - x'|$, rather than the individual locations x and x' . Such correlation functions, also called
 9 kernels, are known as stationary and are used in the majority of cases. An overview of methods for fitting
 10 non-stationary Gaussian processes is given in [33].

11 Let the computer model run at n points, $X = (x_1, \dots, x_n)^T \in \mathcal{X}$ at which we obtain training runs
 12 $\mathbf{F} = (f(x_1), \dots, f(x_n))$, where \mathcal{X} is a p -dimensional input space. Given the emulator hyperparameters and
 13 \mathbf{F} , the posterior for $f(x)$ is

$$f(x)|\mathbf{F}, X, \beta, \delta, \nu, \sigma^2 \sim \mathcal{GP}(m^*(x), c^*(x, x')), \quad (2)$$

with the posterior expectation $m^*(x)$

$$m^*(x) = h(x)^T \beta + c(x, X) \mathbf{A}^{-1} (\mathbf{F} - h(X)^T \beta),$$

and posterior variance $c^*(x, x')$

$$c^*(x, x') = c(x, x') - c(x, X) \mathbf{A}^{-1} c(X, x'),$$

14 and $\mathbf{A} = k(X, X)$. There are often different approaches to handling hyperparameters β, δ and σ^2 . [29]
 15 adopted a maximum likelihood method to fit the hyperparameters. A drawback of this is that the hyper-
 16 parameters are usually highly confounded leading to a ridge on the likelihood surface for large δ and σ^2 .
 17 One way out of this is to specify δ . [30] do this and then propose a 'non-informative' prior $P(\beta, \sigma^2) \propto \sigma^{-2}$.
 18 β and σ^2 can then be integrated out leading to a posterior predictive student-t process emulator. [34] use
 19 an informative Normal-inverse gamma prior for β and σ^2 . Fitting GPs via Full Bayes Markov chain Monte
 20 Carlo (MCMC) methods with a benefit that prior distributions can be used to penalise the ridge on the
 21 likelihood surface [6,33,35].

1 2.2 History matching

2 Like Bayesian calibration, History matching requires a ‘best input’ assumption linking the model and
3 reality via

$$y = f(x^*) + \eta, \quad (3)$$

4 where y represents reality, x^* is the ‘best input’ of the computer model and η is the model discrepancy,
5 which is independent of $f(x)$ and x^* [7]. To learn about x^* , we have observations z with unknown mea-
6 surement error e and

$$z = y + e, \quad (4)$$

7 where e has mean zero and is independent of η .

8 The implausibility measure is used to calculate the distance between the output of the model and the
9 observations so that we can rule out parameter settings that are too far from the observations. For a single
10 output at a given value x , implausibility is defined as

$$\mathcal{I}(x) = \frac{|z - \mathbf{E}[f(x)]|}{\sqrt{\text{Var}[z - \mathbf{E}[f(x)]]}}, \quad (5)$$

11 where $\mathbf{E}[f(x)]$ is the emulator prediction. Under (3) and (4)

$$\text{Var}[z - \mathbf{E}[f(x)]] = \text{Var}[f(x)] + \text{Var}[e] + \text{Var}[\eta]. \quad (6)$$

12 For r outputs, the implausibility can be calculated as the maximum across all outputs [11], the second or
13 third largest [25], or via a multivariate version of (5) [e.g.24].

14 Large values of $\mathcal{I}(x)$ indicate that we are confident that $f(x)$ is too far from the observations and so
15 can be ruled out. The space that has not yet been ruled out, “Not Ruled Out Yet” (NROY) space, \mathcal{X}^1 , is
16 defined as

$$\mathcal{X}^1 = \{x \in \mathcal{X} | \mathcal{I}(x) \leq T\}, \quad (7)$$

17 where T is a selected threshold. A common choice is $T = 3$ based on the three sigma rule [36]. We define
18 the NROY space found by the computer model directly (without an emulator) as “true” NROY space or
19 target NROY in order to compare with the NROY space found using an emulator. Target NROY space \mathcal{X}^*

1 is defined as

$$\mathcal{X}^* = \left\{ x \in \mathcal{X} \mid \frac{|z - f(x)|}{\sqrt{\text{Var}[e] + \text{Var}[\eta]}} \leq T \right\}. \quad (8)$$

2 History matching does not seek to identify \mathcal{X}^* using a single set of computer model evaluations, but
 3 through iteratively designed experiments known as ‘waves’ [15]. In the first wave, the emulator is con-
 4 structed based on an ensemble at a set of points which cover the whole input space \mathcal{X} , then history match-
 5 ing attempts to rule out space from the initial space through (7). Using an ensemble at a new set of points
 6 $X^1 \in \mathcal{X}^1$ (defined in equation(7)), a new emulator can be constructed, and a second wave of history match-
 7 ing will further reduce the input space. In general, at wave k , a new set of points x^k is drawn from the wave
 8 $k-1$ NROY space, $x^k \in \mathcal{X}^{k-1}$ and are used to construct a new emulator. The wave k NROY space is defined
 9 as

$$\mathcal{X}^k = \{x \in \mathcal{X}^{k-1} \mid \mathcal{I}^k(x) \leq T_k\}, \quad (9)$$

10 where T_k is a selected threshold for wave k and $\mathcal{I}^k(x)$ is the implausibility function (e.g. (5)) defined on
 11 \mathcal{X}^{k-1} and with the wave k emulator. Moreover, for different approaches developed for multi-wave designs
 12 \mathbf{x}_k , please see [25] and [37].

13 As NROY space is iteratively reduced, the majority of the runs of the simulator are made closer to
 14 the target space, \mathcal{X}^* . This ensures the density of points in important regions of the model input space
 15 are greater than using Bayesian calibration (assuming an equivalent budget of model runs), ensuring our
 16 emulators are more accurate where it counts. In the later waves, we are more likely to believe Normality
 17 assumptions implicit in using Gaussian processes, a smaller value of threshold T_k can be adopted for
 18 history matching.

19 **2.3 The importance of diagnostics**

20 Before cutting areas of parameter space, diagnostics must be used to assess the adequacy of an emulator.
 21 [26] present a variety of diagnostics that compare Gaussian process emulator predictions and simulation
 22 outputs at validation points. One example is to look at standardised prediction errors, $D_i(f(x_i))$, calculated
 23 via

$$D_i(f(x_i)) = \frac{f(x_i) - \text{E}[f(x_i)]}{\sqrt{\text{Var}[f(x_i)]}}, \quad (10)$$

1 for a set of runs $f(x_1), \dots, f(x_m)$ left out of the training set. Note often these left out runs are actually those
 2 used to fit the emulator, with one at a time left out and the emulator refit (these are Leave One Out (LOO)
 3 errors). [26] state that standardised large errors (larger than 2) suggest there could be a conflict between
 4 the emulator and the simulator, though we should expect 5% of points to fail this test if we have not
 5 been underconfident, and we may have extrapolation issues on the input space boundaries. In practice,
 6 if less than 5% of the errors are large and there is no systematic problem (e.g. all large errors are in the
 7 same region of parameter space) an emulator is considered to have “validated”. It would then be used in
 8 history matching, calibration or for any other purpose by practitioners. Many applications using GPs now
 9 fit hundreds or thousands of emulators simultaneously [38,39], making detailed examination, beyond a
 10 quick check to see if the number of large errors is “about 5%”, impractical.

11 Whilst the test described above may be adequate to assess the global performance of an emulator,
 12 when history matching the primary concern should be the local performance of the emulator within target
 13 NROY, \mathcal{X}^* . When the emulator at a given wave is incorrect outside target NROY, the worst thing that could
 14 happen is that a poor parameter choice is retained. However, this could still be removed by a future wave
 15 with a more accurate emulator. However, an emulator that is inaccurate within true NROY could lead to
 16 good choices being irrecoverably ruled out. As far as we are aware, this concern has not been addressed
 17 within the literature.

To fix ideas, we consider history matching on the 1D function considered by [40]. The function has the equation

$$y(x) = \sin(30(x - 0.9)^4) \cos(2(x - 0.9)) + \frac{(x - 0.9)}{2}.$$

18 We use a 10-run maximin Latin Hypercube (LHC) [41] to design the runs to train the emulators used in
 19 this example.

20 Figure 1 shows the emulator performance and the results of the first wave of history matching, com-
 21 pared to the true NROY space found by the 1D model directly. The true function and the corresponding
 22 emulator posterior mean with uncertainty is shown in the top right panel. We can see that the region
 23 $[0, 0.4]$ is hard to predict by comparing the true function and emulator prediction. From the leave one out
 24 diagnostic plot (top left) we can see that the emulator has failed at one point, but that this single failure
 25 wouldn't be deemed serious enough to invalidate the emulator. The results of history matching with this
 26 emulator are shown in the bottom left and right panels (we set the threshold as 3). Comparing the true
 27 model calibration results with the emulator calibration results, we find nearly one third of the true NROY

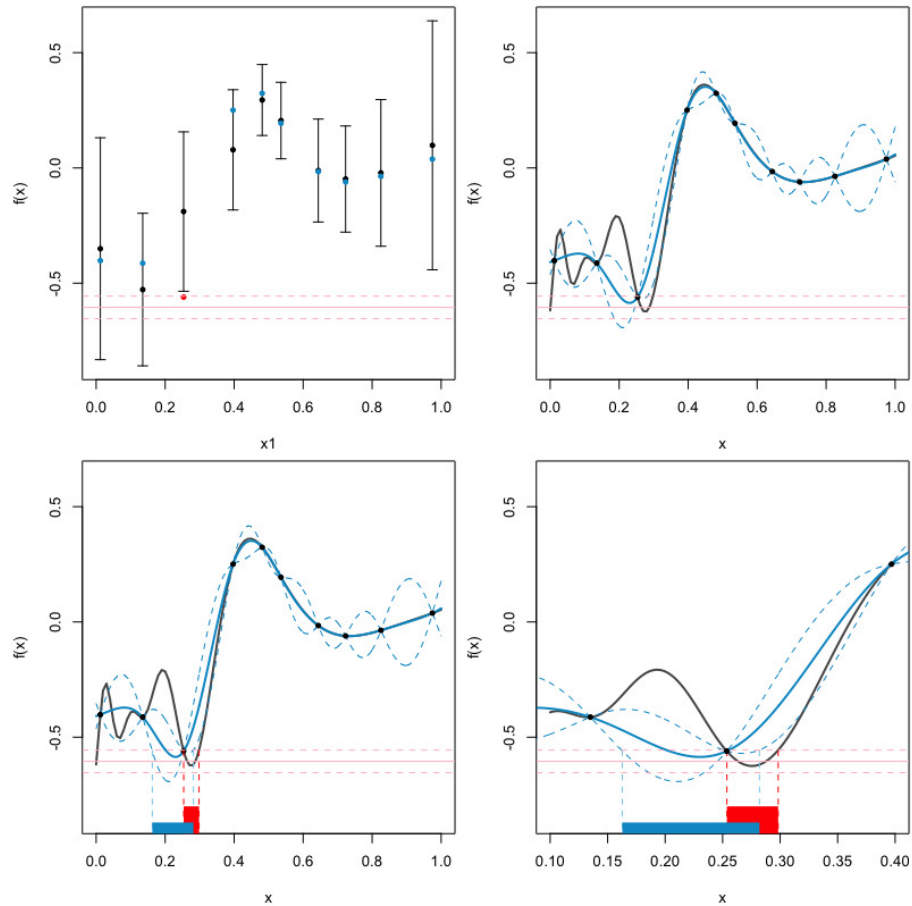


FIG. 1: *Top left:* Leave One Out diagnostic plot against x . The emulator prediction and two standard deviation intervals are given in black. The true function values are in blue if they lie within two standard deviation prediction intervals, or red otherwise. The pink line and the pair of red dotted lines present the observation with observation error and discrepancy in all 4 panels. *Top right:* Emulator performance for the 1D model. The true function is represented by the black curve and ten black points are inputs used to train the emulator. The blue line represents the emulator posterior mean, and the blue dotted lines give the two standard deviation prediction intervals. *Bottom left:* History matching results and the true NROY region. The blue interval defining the NROY space after first wave, the red Interval defining true NROY \mathcal{X}^* . *Bottom right:* As with bottom left but enlarged over the NROY regions.

1 space is ruled out.

2 In an application with many emulators being used to cut a high dimensional parameter space using
 3 many metrics, such critical cases may often occur and be difficult to catch. A trainable nugget would enable
 4 simpler functions to fit the data [35], and this might alleviate the problem in some cases, particularly if we
 5 have achieved what looks like an acceptable fit by over-fitting. We would generally use a trainable nugget
 6 when building emulators, for these reasons. However, in most cases where we see this pathology, the
 7 emulator fits well across the parameter space, but near the true NROY there is an issue which would not
 8 normally raise a diagnostic flag. In these cases, it is likely that the overall fit is good, but that there is some

1 local non-stationarity near where the function behaves like the data. In such cases, if a trainable nugget had
 2 not already been used, one would be unlikely to solve the problem and may lead to reduction in global
 3 performance (i.e. we may still have a good emulator from a validation perspective, but might rule out less
 4 space given that there will be a larger posterior variance). With or without a trainable nugget, we should
 5 still expect 5% of predicted points to lie outside our prediction intervals. If most or all of these occur near
 6 true NROY, we may still rule out good parts of that space by mistake.

7 The history matching literature usually recommends only ruling space out if 3 or more outputs have
 8 large implausibility, and this might insure against ruling out true NROY in some cases (if we have more
 9 than 1 or 2 metrics). However, a poor emulator near the true NROY region may often be a feature of the
 10 design that can appear for emulators of all metrics, and flagging this issue in a key region might make us
 11 wary of trusting emulators for other metrics in that region. Larger cutoff thresholds are sometimes used in
 12 earlier waves to retain more space until we are more confident about ruling out. If this is done routinely, it
 13 may still be that true NROY is ruled out using the larger threshold. If this is done to ensure that all points
 14 where there may be an issue are not ruled out, the cutoff may have to be set so high as to ensure that no
 15 space would be ruled out at all. In the following section we offer a method for detecting the type of history
 16 matching failure we have highlighted and then use it to present a method of robust history matching that
 17 makes use of the emulator we have in the regions where it performs well.

18 3. METHODOLOGY

19 3.1 Detection

20 Suppose we have already fitted a GP emulator. We may use our GP to compute the standardised errors,
 21 $D_i(f(x_i))$, given in equation(10), for design \mathbf{X} . Those errors that would normally be flagged as too large
 22 are grouped into what we term the ‘failed’ set,

$$\mathbf{X}_F = \{x_i \in \mathbf{X} | D_i(f(x_i)) > T_F\} \subseteq \mathbf{X}. \quad (11)$$

23 Here T_F is a threshold which is usually set as 2 (or even 1.96 with the argument that if the emulator were
 24 a good fit, 5% of these points should raise a flag). We treat \mathbf{X}_F as a set of candidate points near which the
 25 emulator could be failing in such a way as to cut out regions of target NROY, \mathcal{X}^* . This could only happen
 26 if a point were inside or close to \mathcal{X}^* and if, at the same time, the flag indicated a failure that would mean it

1 (and neighbouring space) was ruled out. For each point $x_m \in \mathbf{X}_F$, $f(x_m)$ is compared with the observation
 2 z to discover whether the emulator failure points are near \mathcal{X}^* and we form a set of ‘doubt points’ that
 3 could be close enough to target NROY to cause an issue. The doubt points set, \mathbf{X}_D , is defined using (8) via

$$\mathbf{X}_D = \{x_m \in \mathbf{X}_F \mid |z - f(x_m)| \leq T \sqrt{\text{Var}[e] + \text{Var}[\eta]}\}, \quad (12)$$

4 where T is the selected threshold, e is the observation error and η is the discrepancy. We define the set of
 5 remaining points \mathbf{X}_N , $\mathbf{X}_N = \mathbf{X} \setminus \mathbf{X}_D$.

6 Standard history matching can be applied directly if \mathbf{X}_D is empty. Otherwise, in principle, with ex-
 7 isting methods we might have to seek to add further runs from the computer model and/or find a more
 8 complex or bespoke emulation. The latter may not always be easy or even possible. Emulation and history
 9 matching are increasingly methods being adopted by modellers in order to calibrate their own models.
 10 Developing a bespoke emulator using a tailored kernel or mean function may be possible for UQ experts
 11 in any given problem, but it raises barriers to wide application in general that may not be necessary. Fur-
 12 ther model runs near \mathbf{X}_D will likely enable standard methods to work well and fix the issue in many
 13 cases. However, in applications like climate modelling where running the model requires specialist equip-
 14 ment (e.g. supercomputers) and scientist time, it often the case that runs need to be done in batches and
 15 time/budget constraints mean that only a small number of batches will be available. Wasting one of these
 16 just to improve part of an emulator may sacrifice a whole potential wave of history matching.

17 Our method is based on the notion that the emulator works well enough in most of the parameter
 18 space so that it can be used for history matching anyway. However, in regions of space near \mathbf{X}_D , it would
 19 be safer not to remove space at all, and to resample that space in the next wave. Essentially, we will add
 20 further runs of our simulator to correct the errors in this region, but we will first cut out all of the space
 21 that can safely be cut out with the existing emulator. The goal then, when \mathbf{X}_D is not empty, is to separate
 22 the whole input space into two regions, one containing \mathbf{X}_D and the other containing \mathbf{X}_N in such a way as
 23 to ensure that history matching in the latter region only will not discard parts of \mathcal{X}^* .

24 **3.2 Local Voronoi Tessellation**

25 There are several different approaches that can be employed to partition the input space into two distinct
 26 regions. One conventional approach is to use a classification method [42]. Logistic regression may be seen

1 as an obvious choice for classification [43]. However, in these problems \mathbf{X}_D may contain only 1 or 2 points,
 2 meaning that our training data is a highly imbalanced dataset which is hard to use to train a logistic
 3 regression, or any similar model-based classifier. The typical results of these attempts tend to put most or
 4 all points in the ‘normal’ region and fail to adequately capture the doubt region.

5 A machine learning method Synthetic Minority Over-sampling Technique (SMOTE) can be used for
 6 classification on imbalanced datasets [44]. SMOTE uses synthetic data generation to increase the number
 7 of samples in the minority class so that the data set becomes balanced. SMOTE first finds the n -nearest
 8 neighbours in the minority class for each of the samples in the class, then random samples are generated
 9 on the lines between the neighbours. Though promising, SMOTE requires at least two points in \mathbf{X}_D which
 10 in many instances will not apply.

11 A Voronoi tessellation is a partitioning of a space into convex cells called Voronoi regions [45]. Sup-
 12 posing that $\mathbf{X} = (x_1, \dots, x_n)^T \in \mathcal{X}$ is a set of centres of an n -cell Voronoi tessellation, a Voronoi region, \mathcal{V}_i ,
 13 is defined as the set of points in \mathcal{X} , whose ‘nearest’ point is x_i , so that

$$\mathcal{V}_i = \left\{ x \in \mathcal{X} \mid d(x, x_i) \leq d(x, x_j) \right\}, \forall j \in \{1, \dots, n\} \setminus i. \quad (13)$$

14 where $d(x, x_i)$ is commonly defined as the Euclidean distance. When history matching, our correlation
 15 function, $c(x, x')$, provides an appropriate notion of distance between inputs. Our n inputs $x \in \mathcal{X}$ can be
 16 used as the centres of a Voronoi tessellation. We cannot use the correlation directly (as the distance between
 17 the two points increases, the correlation decreases), therefore we define a Voronoi Tessellation \mathcal{V}_i with the
 18 emulator posterior correlation function as

$$\mathcal{V}_i = \left\{ x \in \mathcal{X} \mid |c^*(x, x_i)| \geq |c^*(x, x_j)| \right\}, \forall j \in \{1, \dots, n\} \setminus i. \quad (14)$$

19 Finding a Voronoi tessellation can be computationally challenging when the design is large or when
 20 the input dimensions become even moderately sized (e.g. > 4). However, we do not need to map the
 21 whole parameter space. Our goal is to find the local Voronoi tiles that cover \mathbf{X}_D , ensuring that all possible
 22 values we have not run but might doubt our emulator for near true NROY are included. Specifically, a
 23 local Voronoi tessellation will partition the input space into a doubt region, $\mathcal{X}_D \supseteq \mathbf{X}_D$, and normal region,
 24 $\mathcal{X}_N \supseteq \mathbf{X}_N$, by finding \mathcal{X}_D .

1 We define a local Voronoi tessellation $\mathcal{X}_D = \bigcup_{\{i:x_i \in \mathbf{X}_D\}} \mathcal{V}_i$, with

$$\mathcal{V}_i = \left\{ x \in \mathcal{X} \mid |c^*(x, x_i)| \geq |c^*(x, x_j)| \right\}, \forall j \text{ s.t. } x_j \in \mathbf{X}_N.$$

2 3.3 Local augmentation

3 The local Voronoi tessellation, \mathcal{X}_D , represents the union of convex sets around the doubt points. Given that
 4 the emulator failed to predict the doubt points, but was able to predict the surrounding normal points, we
 5 can deduce that there is a region between each normal point and each doubt point where the emulator is
 6 reliable (it predicts the normal points well) and a region near the doubt points where it is not. Though \mathcal{X}_D
 7 will contain much of this region, if not all, there is no guarantee that it should contain the whole badly
 8 performing region. We therefore include an augmentation step to ensure that as much of the region where
 9 the emulator cannot be trusted (near target NROY) is included in \mathcal{X}_D .

10 For any design point x_i , the design point x_j with the largest value of $c^*(x_i, x_j)$ is the point with the
 11 most influence on x_i . For $x_i \in \mathbf{X}_N$, we want to ensure that their most influential points are not doubt points
 12 where we do not trust our emulator as this would indicate a possibility that some part of the region bor-
 13 dering \mathcal{X}_D and near to x_j is unreliable. Our augmentation step adds all points from \mathbf{X}_N with this property
 14 to \mathbf{X}_D before arriving at a final \mathcal{X}_D .

Specifically, for each $x_i \in \mathbf{X}_N$, let

$$x_{k(i)} = \arg \max_{k:x_k \in \mathbf{X}} c^*(x_i, x_k).$$

15 Let $\mathbf{X}_{D'} = \{x_i : x_{k(i)} \in \mathbf{X}_D\}$, the set of points in \mathbf{X}_N whose most influential point is a doubt point. We then
 16 augment the doubt set by $\mathbf{X}_{D'}$ so that $\mathbf{X}_D = \mathbf{X}_D \cup \mathbf{X}_{D'}$, and compute the local Voronoi tessellation on the
 17 augmented set as before.

18 3.4 Robust history matching

19 Having isolated a region of parameter space, $\mathcal{X}_N = \mathcal{X} \setminus \mathcal{X}_D$ in which we feel confident enough in our
 20 emulator to rule out parameter space, we can history match in just that region. Specifically, we define

$$\mathcal{X}'_N = \{x \in \mathcal{X}_N : \mathcal{I}(x) \leq T\} \tag{15}$$

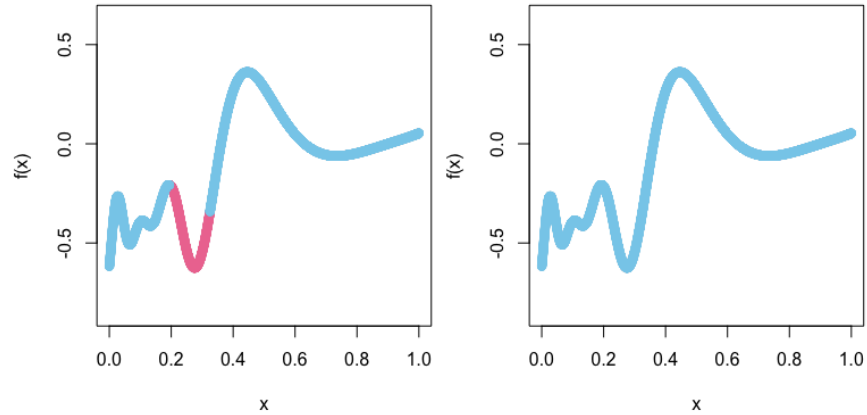


FIG. 2: Comparison between logistic regression classification (right), and Voronoi Tessellation with the GP emulator correlation prior function (left). The blue part is the normal region which can be employed in history matching. The red part is the retained doubt region.

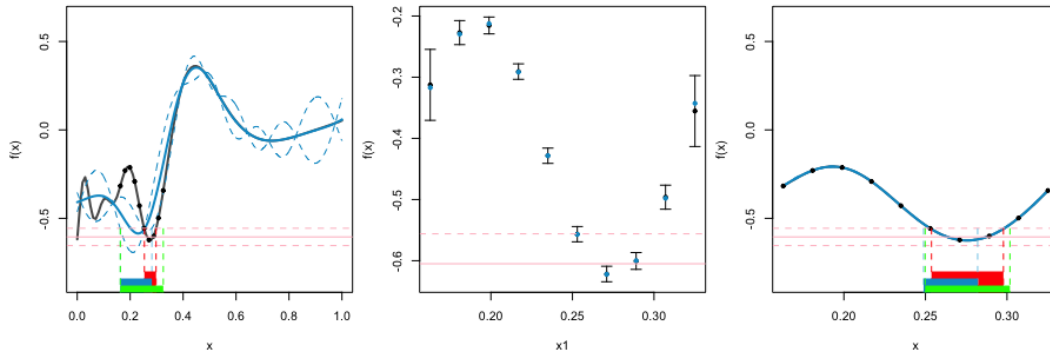


FIG. 3: The 1-dimensional model multi-wave calibration result. Left: history matching with our method first wave result, the red interval defining the true NROY space, the blue interval defining the NROY space by standrad history matching and the green interval defining the the NROY space by our method. Centre: leave One Out diagnostic plot against x for second wave emulator. Right: history matching second wave result.

1 with $\mathcal{I}(x)$ as in equation(5). The NROY space \mathcal{X}^1 after wave 1 is defined as

$$\mathcal{X}^1 = \mathcal{X}_D \cup \mathcal{X}'_N. \quad (16)$$

2 4. NUMERICAL EXAMPLES

3 4.1 The 1-dimensional function

4 We apply the methodology of the last section to the 1-dimensional function from Section 2.3 and a 5-
5 dimensional function described below. We use the R package `DiceKriging` [46] to construct the emula-
6 tors throughout.

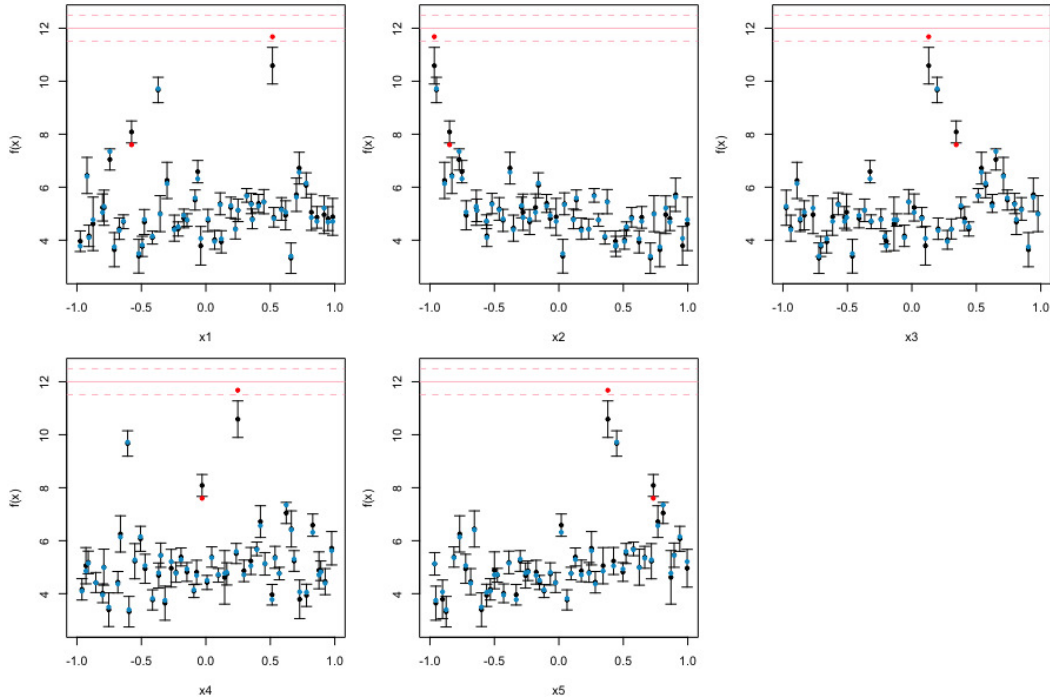


FIG. 4: Leave one out diagnostic plots. Each panel represents leave one out predictions from an emulator against one of the 5 inputs. Black points and error bars are from the emulator posterior mean and two standard deviation prediction intervals. The true function values are in blue if they lie within two standard deviation prediction intervals, or red otherwise.

1 The 1D model of Section 2.3 only has 1 doubt point. The doubt region highlighted by our local Voronoi
 2 tessellation is highlighted in red on the left hand panel of Figure 2. The right hand panel shows that the
 3 logistic regression classifier fails to identify any doubt region due to the unbalanced design.

4 We present our robust history match of this function in Figure 3. The wave 1 result is shown in the
 5 left panel with the green interval defining wave 1 NROY space. A second wave is performed with 10
 6 randomly selected runs within NROY space and the leave one outs are shown in the middle panel of
 7 Figure 3, highlighting that there are no doubt points. The right panel shows the second wave results. We
 8 see that all of the target NROY space is retained.

9 4.2 A 5-dimensional function

In order to examine the performance of our method in higher dimensions we look at the 5D function

$$f(\mathbf{x}) = \sqrt{x_1} + \frac{1}{\sqrt{x_2}} + x_3 + \sin(x_4) + \exp(x_5).$$

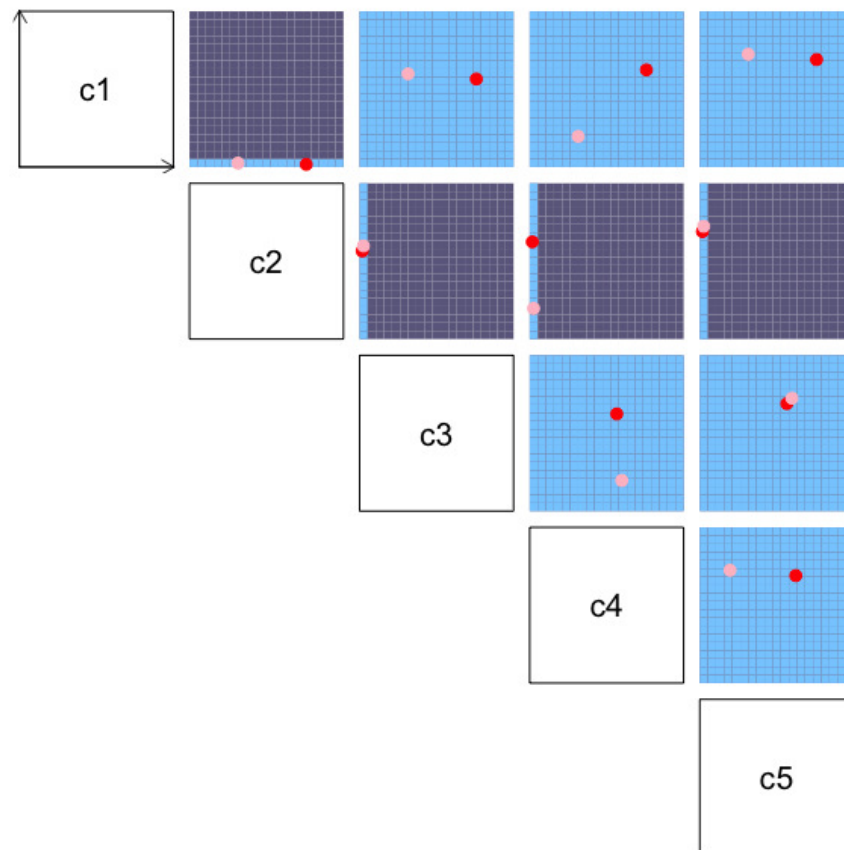


FIG. 5: Local Voronoi cell plots over each two parameters. The red point is the doubt point and the pink points are selected by augmentation step. The light blue region is the Local Voronoi cell of the doubt points which is the doubt region.

- 1 Note that this function tends to infinity as x_2 tends to zero which may happen in many physical models.
- 2 We use a 50 member maximin Latin Hypercube (LHC) to select points for wave 1 and use the function
- 3 evaluations to construct an emulator.
- 4 Leave one out diagnostics against each input are presented in Figure 4. By eye, we see that the emulator
- 5 has individual large errors near the observations, which might indicate that the emulator does not simulate
- 6 the target 'NROY' space effectively. Using equation(12) we identified 1 doubt point. The local Voronoi
- 7 tessellation plot for all inputs is presented in Figure 5. The red point is the doubt point and the pink point
- 8 is selected by the augmentation step. The light blue range is the Local Voronoi tessellation. We apply the
- 9 robust history matching algorithm described in Section 3.4, retaining the local Voronoi tiles as part of
- 10 NROY space and applying the usual constraint to the rest of the space.

TABLE 1: Standard vs robust history matching with top row as the percentage of the original space as NROY and the bottom the percentage of target NROY retained.

	Standard HM wave 1	Robust HM wave 1	Robust HM wave 2	Robust HM wave 3
Retained NROY volume	0.6373%	4.6660%	3.9132%	0.2251 %
Retained target NROY %	24.755%	99.643%	99.643%	99.449%

1 We compare our robust method with standard history matching in Figure 6. In these density plots,
 2 each pixel on any panel represents the proportion of points behind that pixel in the other 3 dimensions of
 3 the parameter space that is NROY. The scale corresponds to the colours in the upper triangles, whilst plots
 4 on the lower triangle mirror the upper triangle but with independent scales so as to reveal any structure
 5 hidden by the comparative colour scheme.

6 The top left panel in Figure 6 shows the target NROY space and the top right panel in Figure 6 shows
 7 the wave 1 NROY space following standard history matching. The first wave has incorrectly cut out a large
 8 corner of the target region (low x_1 and low x_2, x_3, x_4 and the lower half of x_5). Wave 1 of robust history
 9 matching, shown in Figure 6 (bottom left panel), does not have this issue and cuts out less parameter
 10 space overall (as expected). We continue to perform robust history matching for 2 further waves, though
 11 in waves 2 and 3, there were no doubt points, so these waves are the same as standard history matching
 12 (but from a different wave 1). The wave 3 NROY space is shown in Figure 6 (bottom right panel).

13 Table 1 shows the volume of NROY space as a percentage of the original space (top row) and the
 14 percentage of target NROY retained following each wave of history matching (bottom row). Target NROY
 15 is 0.09% of the original space. Though standard history matching cuts more space than our robust method
 16 in wave 1, it cuts out nearly 75% of target NROY, whilst we only cut 0.2%. After 2 further waves of history
 17 matching, we have still retained the target NROY, but have reduced our NROY to 0.17% of the original
 18 space.

19 This example shows a case where history matching can be non-robust in 5 dimensions and that our
 20 robust history matching effectively enables us to continue the analysis, without having to run a new wave
 21 1. We now show a case from our work tuning climate models where this issue has presented itself and how
 22 our method performs.

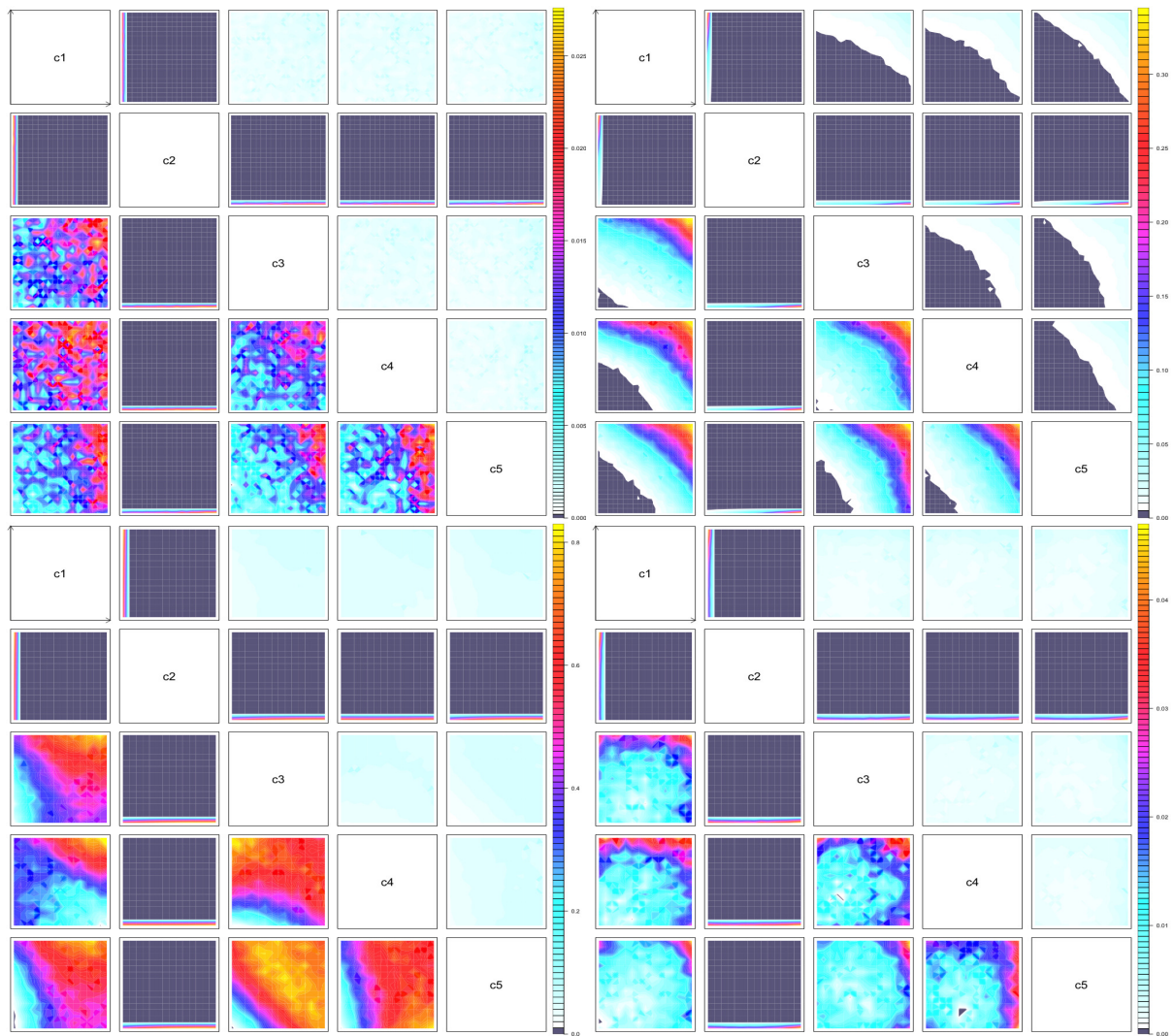


FIG. 6: NROY density plots for 2-D projections of NROY space. *Top left:* Target NROY space. *Top right:* Wave 1 NROY space following standard history matching. *Bottom left:* Wave 1 NROY space following robust history matching. *Bottom right:* Wave 3 NROY space after robust history matching. The scale corresponds to the colours in the upper triangles, whilst plots on the lower triangle mirror the upper triangle but with independent scales so as to reveal any structure hidden by the comparative colour scheme (the change from light blue, blue to red indicates that the density is rising).

1 5. APPLICATION: PROCESS-BASED TUNING OF CLIMATE MODELS

2 As part of the ANR (The French National Research Agency) funded project HIGH-TUNE, developers of
3 the French climate models CNRM-CM and IPSL-CM are developing tools to automatically tune bound-
4 ary layer cloud parameterisations within their models based on history matching to high resolution Large
5 Eddy Simulations. Our collaboration involves providing methods to emulate and history match to a large
6 number of process-based metrics quickly and automatically, so that the modellers can use the tools in-
7 dependently. With multiple unsupervised history matches, it is important that our methods are robust to
8 possible ensemble issues, and so the method we describe in this paper should be part of our set of tools.
9 We illustrate its importance through an example of a metric that fails our tests in IPSL-CM.

10 IPSL-CM is an atmosphere model that is used to predict planetary atmospheres, including the Earth
11 and other celestial bodies (Mars, Titan, Venus), as well as regional climate, process studies [2,27,28]. We run
12 a single column version of the model and perturb 5 cloud parameters chosen by the modellers. The model
13 is run for a particular boundary layer case (in this case SANDU, capturing transitions from cumulus to
14 stratocumulus clouds) with the idea of seeing which parameter choices lead to reasonable representations
15 of clouds in these region types (compared to high-resolution simulations). Parameter ranges were deter-
16 mined by the project, and in our analysis we have mapped the parameters onto $[-1, 1]$ for fitting emulators
17 and history matching.

18 We generate a 30-member design as the first 2 LHCs in a 150-member extended LHC composed of
19 10, 15-member LHCs following [22] (each additional LHC ensures that the composite design is orthogonal
20 and fills the space in each extension phase). Leave one out diagnostic plots for our fitted emulator are
21 presented in the top of Figure 7. To history match, we use an observation of 12.18, the observation error
22 variance and discrepancy variance are both set as 0.0006.

23 There are 2 failed points near the observation, which might indicate that the emulator does not sim-
24 ulate the target NROY space well. Using equation(12) we identify 1 doubt point and another doubt point
25 defined by our augmentation step. Since the target NROY is unknown in the climate model, in order to
26 fairly compare our method with standard history matching, we use the remaining 120 data points (from
27 our 150 member LHC) as validation data. The validation results are shown on Figure 7 in the middle and
28 bottom rows. In this small data set, we have 11 points in target NROY space, the standard history matching
29 misses one target point after wave 1, our method retains all the true NROY. In order to fairly compare, we
30 do three waves with each methodology.

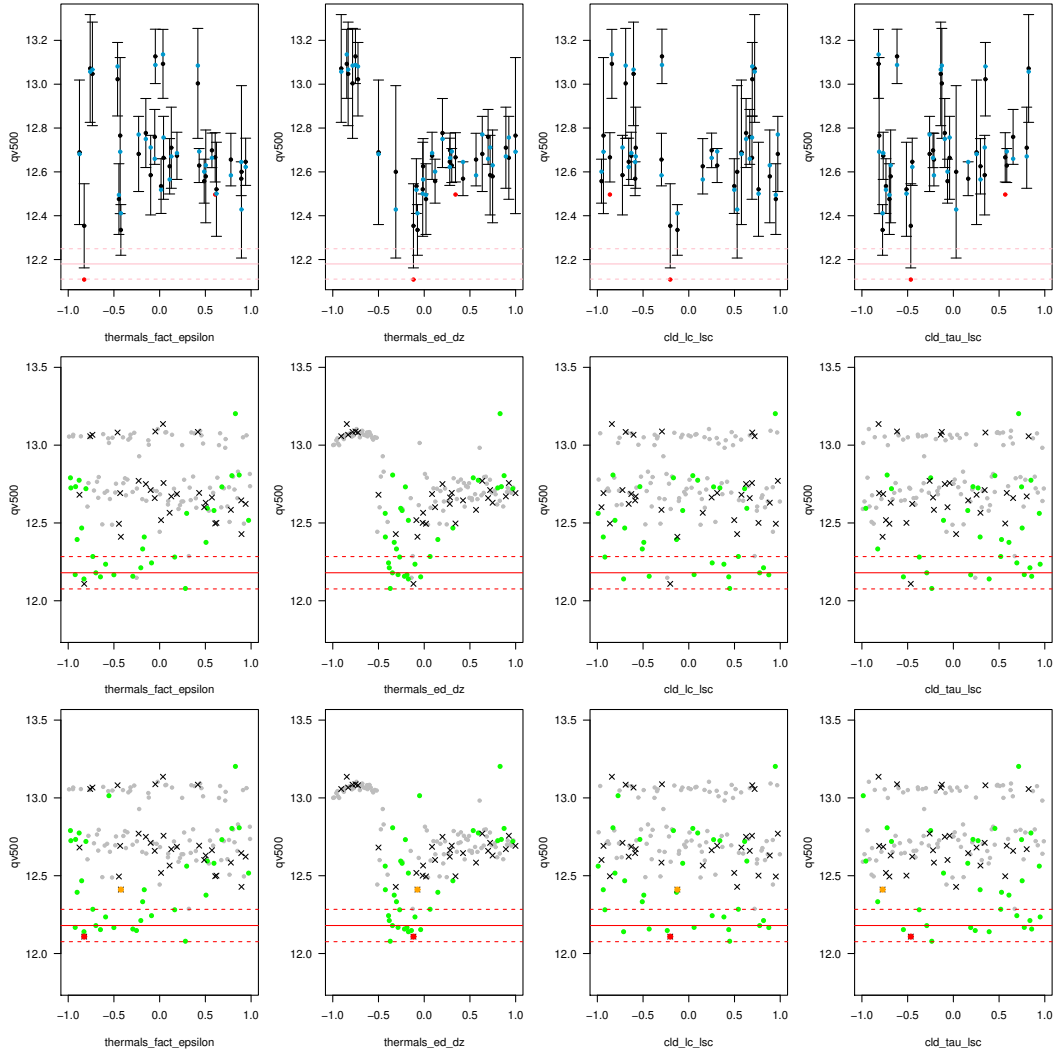


FIG. 7: *Top:* Leave one out diagnostic plots. Each panel represents one left-out latin hypercube predicted by the emulator, black points and error bar are from the emulator posterior mean and two standard deviation prediction intervals. The true function values are in blue if they lie within two standard deviation prediction intervals, or red otherwise. The observation with observation error are in solid and dotted red lines respectively. *Middle:* Validation results after wave 1 following standard history matching. All the points are model runs with the emulator training data presented in black. The validation data are green if they are retained in the NROY after wave 1 history matching, or grey otherwise. *Bottom:* Validation results after wave1 following robust history matching. The colours are as for the middle row with the red point being the original doubt point and the orange point, the doubt point selected by our augmentation step.

1 The NROY density plots (upper triangle) and minimum implausibility plots (lower triangle) are pre-
 2 sented in Figures 8, 9, 10 and 11 (upper right), for each pair of parameters. For the NROY density plots,
 3 each pixel on any panel represents the proportion of points behind that pixel in the other 2 dimensions of
 4 the parameter space that is NROY. The right scale corresponds to the colours in the NROY density plots.

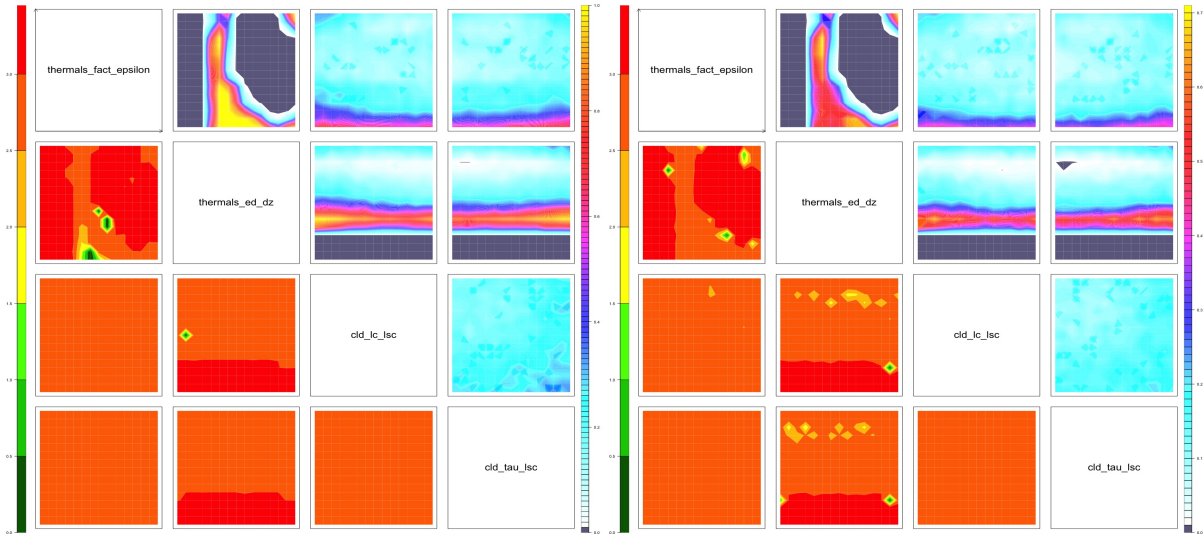


FIG. 8: Wave 1 NROY space for LMDZ-SANDU after robust history matching.

FIG. 9: Wave 3 NROY space for LMDZ-SANDU after robust history matching.

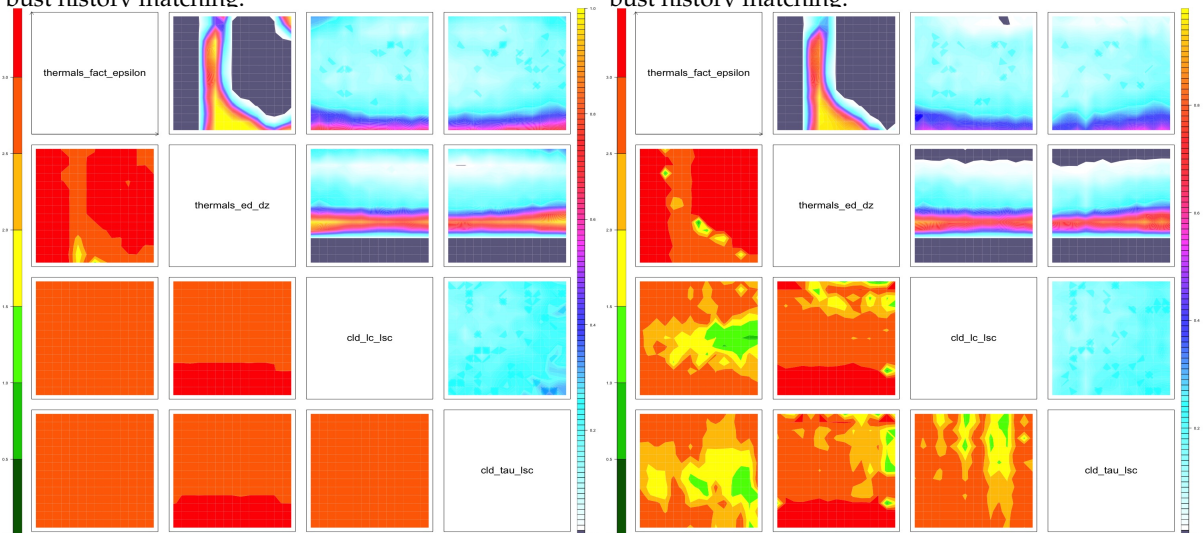


FIG. 10: Wave 1 NROY space for LMDZ-SANDU following standard history matching.

FIG. 11: Wave 3 NROY space for LMDZ-SANDU following standard history matching.

1 For the minimum implausibility plots, each pixel on any panel on the lower triangle represents the mini-
2 mum implausibility found in the remaining dimensions of parameters space behind that pixel. The colours
3 are given indicated by the scale on the left.

4 Only the first wave had doubt points, so wave 2 and wave 3 use standard history matching. Our
5 method retains more space in the first wave (around 1%). We can see from Figure 10 that this retained
6 space is in the centre of the space spanned by `thermals_fact_epsilon` and `thermals_ed_dz`. The wave
7 3 NROY density plot of robust history matching shows the doubt area is still in the NROY space, showing
8 that standard history matching incorrectly ruled out part of \mathcal{X}^* .

9 This application showed that incorrectly ruling out parameter space can occur in practice: in this case
10 when history matching climate model parameterisations. For climate models in particular, this mistake
11 could prove very costly as history matching is used to assess the quality of a given parameterisation or an
12 alternative. If good models are accidentally discarded, the parameterisation or even the resolution of the
13 model or its implementation might be needlessly changed, wasting time and resource for the modelling
14 centre.

15 6. DISCUSSION

16 In this paper we demonstrated a potential issue with history matching that occurs when emulators that
17 seem to validate well by most standard analyses, do not simulate the (unknown) target subspace well
18 enough. We showed that this can lead to good parts of parameter space being ruled out unintentionally,
19 and that existing methods, such as using nuggets, variable thresholds or only ruling out if multiple metrics
20 fail, were not designed specifically for such pathological cases and do not necessarily address the problem.
21 We developed a method for detecting these cases based on standard diagnostics. We then presented a
22 robust history matching method based on using a tailored local Voronoi tessellation designed to capture
23 the region where the emulator is not as good as it needs to be, and isolate it so that the rest of the input
24 space can be pruned as normal, without having to re-run the simulator.

25 We demonstrated the efficacy of our method in comparison to standard history matching for 2 numer-
26 ical examples designed to demonstrate the issue, and then applied the method to a process metric from
27 a single column version of the French climate model LMDZ. We showed that, unlike standard history
28 matching, our method manages to effectively cut parameter space whilst ensuring that the target space is
29 preserved.

1 Whilst it may be possible to observe the diagnostic issue we have highlighted and to offer a bespoke
2 history match for a particular quantity in any given application, this is not feasible in applications where
3 tens, hundreds or even thousands of emulators are built and are to be compared with observations [see,
4 e.g.15,38,39] as part of the calibration. We also want methods that do not require frequent intervention by
5 an experienced statistician. Hence our robust method provides a way to safely and automatically isolate
6 any regions of parameter space where it would be dangerous to history match with the current emulator,
7 but still allows the same emulator to be used appropriately without requiring a bespoke analysis.

8 **ACKNOWLEDGMENTS**

9 Daniel Williamson was funded by an Alan Turing Fellowship. The authors gratefully acknowledge the
10 support from Agence Nationale de la Recherche (ANR) (grant HIGH-TUNE ANR-16-CE01-0010), which
11 funded the runs of IPSL-CM. The authors thank the Isaac Newton Institute for Mathematical Sciences,
12 Cambridge, for support and hospitality during the Uncertainty Quantification programme where work on
13 this article was undertaken (EPSRC grant no EP/K032208/1).

14 **REFERENCES**

- 15 1. Gladstone, R.M., Lee, V., Rougier, J., Payne, A.J., Hellmer, H., Le Brocq, A., Shepherd, A., Edwards, T.L., Gregory,
16 J., and Cornford, S.L., Calibrated prediction of pine island glacier retreat during the 21st and 22nd centuries with
17 a coupled flowline model, *Earth and Planetary Science Letters*, 333:191–199, 2012.
- 18 2. Hourdin, F., Mauritsen, T., Gettelman, A., Golaz, J.C., Balaji, V., Duan, Q., Folini, D., Ji, D., Klocke, D., Qian, Y., ,
19 The art and science of climate model tuning, *Bulletin of the American Meteorological Society*, 98(3):589–602, 2017.
- 20 3. Sacks, J., Welch, W.J., Mitchell, T.J., and Wynn, H.P., Design and analysis of computer experiments, *Statistical sci-*
21 *ence*, pp. 409–423, 1989.
- 22 4. Santner, T.J., Williams, B.J., and Notz, W.I., *The design and analysis of computer experiments*, Springer Science & Busi-
23 *ness Media*, 2003.
- 24 5. Kennedy, M.C. and O’Hagan, A., Bayesian calibration of computer models, *Journal of the Royal Statistical Society:*
25 *Series B (Statistical Methodology)*, 63(3):425–464, 2001.
- 26 6. Higdon, D., Gattiker, J., Williams, B., and Rightley, M., Computer model calibration using high-dimensional out-
27 put, *Journal of the American Statistical Association*, 103(482):570–583, 2008.

- 1 7. Craig, P., Goldstein, M., Seheult, A., and Smith, J., Bayes linear strategies for matching hydrocarbon reservoir
2 history, *Bayesian statistics*, 5:69–95, 1996.
- 3 8. Vernon, I., Liu, J., Goldstein, M., Rowe, J., Topping, J., and Lindsey, K., Bayesian uncertainty analysis for complex
4 systems biology models: emulation, global parameter searches and evaluation of gene functions, *BMC systems*
5 *biology*, 12(1):1, 2018.
- 6 9. Gong, Z., DiazDelaO, F., and Beer, M., Sampling schemes for history matching using subset simulation, In *Pro-*
7 *ceedings for the 1st International Conference on Uncertainty Quantification in Computational Sciences and Engineering*,
8 2017.
- 9 10. Jackson, S.E., Vernon, I., Liu, J., and Lindsey, K., Bayesian uncertainty analysis establishes the link between the
10 parameter space of a complex model of hormonal crosstalk in arabidopsis root development and experimental
11 measurements, *Stat. Appl. Genet. Mol. Biol*, 2018.
- 12 11. Craig, P.S., Goldstein, M., Seheult, A.H., and Smith, J.A. Pressure matching for hydrocarbon reservoirs: a case
13 study in the use of bayes linear strategies for large computer experiments. In *Case studies in Bayesian statistics*, pp.
14 37–93. Springer, 1997.
- 15 12. Cumming, J.A. and Goldstein, M., Bayes linear uncertainty analysis for oil reservoirs based on multiscale computer
16 experiments, *O’Hagan, West, AM (eds.) The Oxford Handbook of Applied Bayesian Analysis*, pp. 241–270, 2010.
- 17 13. Andrianakis, I., Vernon, I.R., McCreesh, N., McKinley, T.J., Oakley, J.E., Nsubuga, R.N., Goldstein, M., and White,
18 R.G., Bayesian history matching of complex infectious disease models using emulation: a tutorial and a case study
19 on hiv in uganda, *PLoS computational biology*, 11(1):e1003968, 2015.
- 20 14. Andrianakis, I., McCreesh, N., Vernon, I., McKinley, T.J., Oakley, J.E., Nsubuga, R.N., Goldstein, M., and White,
21 R.G., Efficient history matching of a high dimensional individual-based hiv transmission model, *SIAM/ASA Jour-*
22 *nal on Uncertainty Quantification*, 5(1):694–719, 2017.
- 23 15. Vernon, I., Goldstein, M., Bower, R.G., , Galaxy formation: a bayesian uncertainty analysis, *Bayesian Analysis*,
24 5(4):619–669, 2010.
- 25 16. Bower, R.G., Vernon, I., Goldstein, M., Benson, A., Lacey, C.G., Baugh, C.M., Cole, S., and Frenk, C., The parameter
26 space of galaxy formation, *Monthly Notices of the Royal Astronomical Society*, 407(4):2017–2045, 2010.
- 27 17. Rodrigues, L.F.S., Vernon, I., and Bower, R.G., Constraints on galaxy formation models from the galaxy stellar
28 mass function and its evolution, *Monthly Notices of the Royal Astronomical Society*, 466(2):2418–2435, 2017.
- 29 18. Edwards, N.R., Cameron, D., and Rougier, J., Precalibrating an intermediate complexity climate model, *Climate*
30 *dynamics*, 37(7-8):1469–1482, 2011.
- 31 19. McNeall, D., Challenor, P.G., Gattiker, J., and Stone, E., The potential of an observational data set for calibration of

- 1 a computationally expensive computer model, 2013.
- 2 20. Williamson, D., Goldstein, M., Allison, L., Blaker, A., Challenor, P., Jackson, L., and Yamazaki, K., History matching
3 for exploring and reducing climate model parameter space using observations and a large perturbed physics
4 ensemble, *Climate dynamics*, 41(7-8):1703–1729, 2013.
- 5 21. Williamson, D. and Blaker, A.T., Evolving bayesian emulators for structured chaotic time series, with application
6 to large climate models, *SIAM/ASA Journal on Uncertainty Quantification*, 2(1):1–28, 2014.
- 7 22. Williamson, D., Blaker, A.T., Hampton, C., and Salter, J., Identifying and removing structural biases in climate
8 models with history matching, *Climate dynamics*, 45(5-6):1299–1324, 2015.
- 9 23. Salter, J.M. and Williamson, D., A comparison of statistical emulation methodologies for multi-wave calibration of
10 environmental models, *Environmetrics*, 27(8):507–523, 2016.
- 11 24. Salter, J.M., Williamson, D.B., Scinocca, J., and Kharin, V., Uncertainty quantification for spatio-temporal computer
12 models with calibration-optimal bases, *arXiv preprint arXiv:1801.08184*, 2018.
- 13 25. Williamson, D.B., Blaker, A.T., and Sinha, B., Tuning without over-tuning: parametric uncertainty quantification
14 for the nemo ocean model, *Geoscientific Model Development*, 10(4):1789, 2017.
- 15 26. Bastos, L.S. and O’Hagan, A., Diagnostics for gaussian process emulators, *Technometrics*, 51(4):425–438, 2009.
- 16 27. Voldoire, A., Sanchez-Gomez, E., y Méliá, D.S., Decharme, B., Cassou, C., Sénési, S., Valcke, S., Beau, I., Alias,
17 A., Chevallier, M., , The cnrm-cm5. 1 global climate model: description and basic evaluation, *Climate Dynamics*,
18 40(9-10):2091–2121, 2013.
- 19 28. Bony, S. and Dufresne, J.L., Marine boundary layer clouds at the heart of tropical cloud feedback uncertainties in
20 climate models, *Geophysical Research Letters*, 32(20), 2005.
- 21 29. Currin, C., Mitchell, T., Morris, M., and Ylvisaker, D., Bayesian prediction of deterministic functions, with applica-
22 tions to the design and analysis of computer experiments, *Journal of the American Statistical Association*, 86(416):953–
23 963, 1991.
- 24 30. Haylock, R. and O’Hagan, A., On inference for outputs of computationally expensive algorithms with uncertainty
25 on the inputs, *Bayesian statistics*, 5:629–637, 1996.
- 26 31. Rasmussen, C.E. and Williams, C.K., *Gaussian processes for machine learning*, Vol. 1, MIT press Cambridge, 2006.
- 27 32. Andrianakis, I. and Challenor, P.G., The effect of the nugget on gaussian process emulators of computer models,
28 *Computational Statistics & Data Analysis*, 56(12):4215–4228, 2012.
- 29 33. Volodina, V. and Williamson, D., Diagnostics-driven nonstationary emulators using kernel mixtures, *SIAM/ASA*
30 *Journal on Uncertainty Quantification*, 8(1):1–26, 2020.

- 1 34. Oakley, J.E. and O'Hagan, A., Probabilistic sensitivity analysis of complex models: a bayesian approach, *Journal of*
2 *the Royal Statistical Society: Series B (Statistical Methodology)*, 66(3):751–769, 2004.
- 3 35. Gramacy, R.B. and Lee, H.K., Cases for the nugget in modeling computer experiments, *Statistics and Computing*,
4 22(3):713–722, 2012.
- 5 36. Pukelsheim, F., The three sigma rule, *The American Statistician*, 48(2):88–91, 1994.
- 6 37. Coveney, S. and Clayton, R.H., Fitting two human atrial cell models to experimental data using bayesian history
7 matching, *Progress in biophysics and molecular biology*, 139:43–58, 2018.
- 8 38. Lee, L., Carslaw, K., Pringle, K., and Mann, G., Mapping the uncertainty in global ccn using emulation, *Atmospheric*
9 *Chemistry and Physics*, 12(20):9739–9751, 2012.
- 10 39. Gu, M., Berger, J.O., , Parallel partial gaussian process emulation for computer models with massive output, *The*
11 *Annals of Applied Statistics*, 10(3):1317–1347, 2016.
- 12 40. Xiong, Y., Chen, W., Apley, D., and Ding, X., A non-stationary covariance-based kriging method for metamodelling
13 in engineering design, *International Journal for Numerical Methods in Engineering*, 71(6):733–756, 2007.
- 14 41. Morris, M.D. and Mitchell, T.J., Exploratory designs for computational experiments, *Journal of statistical planning*
15 *and inference*, 43(3):381–402, 1995.
- 16 42. Bailey, K.D., *Typologies and taxonomies: an introduction to classification techniques*, Vol. 102, Sage, 1994.
- 17 43. Menard, S., *Applied logistic regression analysis*, Vol. 106, Sage, 2002.
- 18 44. Chawla, N.V., Bowyer, K.W., Hall, L.O., and Kegelmeyer, W.P., Smote: synthetic minority over-sampling technique,
19 *Journal of artificial intelligence research*, 16:321–357, 2002.
- 20 45. Gallier, J., Notes on convex sets, polytopes, polyhedra, combinatorial topology, voronoi diagrams and delaunay
21 triangulations, *arXiv preprint arXiv:0805.0292*, 2008.
- 22 46. Roustant, O., Ginsbourger, D., and Deville, Y., Dicekriging, diceoptim: Two r packages for the analysis of computer
23 experiments by kriging-based metamodelling and optimization, 2012.