

Evaluating the Suitability of Regression-Based Emulators of Building Performance in Practice: A Test Suite

Parag Rastogi ^{*,1, 2, 3}, Mohammad Emtiyaz Khan², and Marilynne Andersen³

¹arbnco Ltd., Glasgow, UK

²RIKEN Center for Advanced Intelligence Project, Tokyo, Japan

³LIPID laboratory, École Polytechnique Fédérale de Lausanne, Switzerland

DD MMM 2021

1 Abstract

Building Performance Simulation (BPS), a useful tool to assess the operational performance of buildings and systems, can often be computationally expensive. The use of BPS is cumbersome for problems where the speed of response is important, e.g., real-time control, uncertainty quantification, parametric exploration, or stock modelling. Emulators, such as those based on regression, offer a faster substitute, but their reliability can be questionable. This paper proposes seven tests to check if an emulator *is a suitable replacement for simulation* in practice. The tests are categorized using four criteria: accuracy, speed, generalisability, and ease of use. The tests can be included in the process of setting up an emulator-based workflow. A use case is provided for emulators based on linear and non-linear regression (Gaussian Process models). This work aims to enable a practitioner to reliably conduct performance assessment for buildings using emulators.

Keywords— non-linear regression, building simulation, test suite, regression model

1 Introduction

Building Performance Simulation (BPS) is a tool to quantify the impact of decisions about a building's design, specification, or operation on its performance, i.e., changes in thermal or visual conditions, energy use, or other physical quantities. This is useful when measured data cannot be obtained from a real building and its systems, or when the systems being modelled are too complex for manual calculations. BPS provides estimates of building performance under hypothetical conditions such as a future climate, changes in building operations, or the impact of retrofits. The ensembles of physics-based equations that make up BPS are usually deterministic: given a set of inputs, simulators will always give the same outputs. These outputs are also precise, i.e., simulators do not typically estimate the uncertainty in outputs. A BPS program can often be computationally-expensive and many design exercises require hundreds of runs for each decision, both of which are difficult to implement in practice and slow down decision-making. Important use cases where the simulator's speed of response is consequential include:

- Monte Carlo (MC) sampling for uncertainty or sensitivity quantification – simulating a building or system with several plausible values of an unknown or poorly characterised input like weather,
- parametric design exploration – testing the impact of several different variations of the

*Corresponding author: contact@paragrastogi.com

57 design or specification of a component like
 58 window sizes,

- 59 • stock modelling – estimating the performance
 60 of large groups of buildings, like modelling
 61 housing archetypes at an urban or national
 62 scale,
- 63 • early design phase exploration – when several
 64 consequential decisions must be made while it
 65 is impractical to set up a full simulation for
 66 each decision since too many aspects of the
 67 final building design are unknown.

68 To be practical in time-bound decision-making
 69 workflows, either individual simulations must be
 70 sped up or computational power increased suffi-
 71 ciently. While big datasets have been analysed
 72 with supercomputers for some applications [e.g., 1],
 73 most users only have access to single computers or
 74 small-scale cloud services like JEPlus¹ or NREL’s
 75 PAT². This means that individual simulations
 76 must be sped up. Options exist in most simula-
 77 tion programs to use simpler variants of underlying
 78 algorithms, these achieve only limited efficiencies.
 79 One option for speeding up individual estimates is
 80 the use of so-called *emulators*, alternative models
 81 that estimate the same quantities as the original
 82 BPS. How an emulator may be judged to be a suit-
 83 able replacement for a simulator is the problem ad-
 84 dressed in this paper.

85 Emulators replace, and usually simplify, the
 86 complex ensemble of physics-based equations that
 87 make up a simulator. The inputs for emulators
 88 need not be exactly the same as those used by
 89 physics-based simulators. For regression models,
 90 these inputs are called features. That is, those
 91 quantifiable features or characteristics of a dataset
 92 that can be used to characterise the variety or vari-
 93 ance seen in the dataset. If the dataset is then rep-
 94 resentative of the physical system being examined,
 95 the features can be said to describe the physical
 96 system as well.

97 While the use of emulators is motivated by the
 98 need to reduce computational burden, an emulator
 99 must accurately represent the physical behaviour

¹www.jeplus.org

²http://nrel.github.io/OpenStudio-user-documentation/reference/parametric_analysis_tool_2/

of the building systems being modelled to be use-
 ful. That is, the difference between predictions of
 performance by an emulator and simulator should
 be acceptably low. The objective of this work is to
 provide a set of tests that can be used to deter-
 mine whether a simpler, faster mathematical model
 (an emulator) is a suitable and viable substitute for
 BPS.

This paper lays the groundwork for a test suite
 to evaluate emulators for a given problem, along
 with its application to a simple problem of predict-
 ing energy use. This could be extended with the
 development of a catalogue of results using com-
 mon design problems and emulators showing, for
 example, how emulators will improve solutions for
 some problems and not for others, where they are
 more applicable, and whether the selection of spe-
 cific emulators/algorithms can be generalised to a
 class of problems. The original code used for this
 paper is available online³, and is open for imple-
 mentation in tools.

To reduce computational time without com-
 promising the usefulness of performance estimates,
 any replacement for BPS must be (i) **accurate**,
 (ii) **fast**, (iii) **generalisable**, and (iv) **easy** to use
 and setup. These criteria can be used to judge
 whether an emulator is a sufficiently useful replace-
 ment for a simulator. In this paper, we provide
 a suite of seven tests that can be used to evalu-
 ate emulators against these criteria. Namely, the
 model should fulfil the following conditions, further
 developed in Section 3.1:

- (1) error compared to simulator outputs is accept-
ably low,
- (2) performance improves with more training data
(lower error),
- (3) performance improves with more complex or
varied training data,
- (4) performance does not degrade (error does not
increase) too much for a specific test case,
- (5) performance is consistent across different test
data sets,
- (6) emulator is computationally cheaper than a
simulator,

³www.github.com/author/repository

144 (7) performs appreciably better (lower error) than
145 a simpler model.

146 These tests can be used to both identify whether
147 a specific emulator is usable or not (pass/fail) and
148 to compare different emulators, e.g., ranking candi-
149 dates by error on the test data. The tests are
150 not a substitute for expert judgement and may give
151 conflicting answers, e.g., a model may generalise
152 poorly (test 5) but perform very well for a specific
153 test case (test 4), a model may show very low error
154 (tests 2-5) but have an unacceptably high compu-
155 tational burden (test 6), etc. It is also possible for a
156 model to fail on a particular test, e.g., failing to per-
157 form consistently well across different test data sets
158 (test 5). It is difficult to offer a generalisable rule in
159 these cases, and users must consider the nature of
160 the problem. What is important to a given problem
161 determines the importance of a given test: gener-
162 alisability, accuracy, speed, or ease. An emulator
163 should do well on tests 6 and 7, since there is little
164 point to replacing a simulator with a more expens-
165 ive emulator or using a complex emulator when a
166 simple one will do.

167 In the next section, we describe possible candid-
168 ates for emulators, existing work, and introduce
169 regression models. After that, Section 3 lays out
170 the context for how and when these emulators are
171 useful, and the example dataset used in this pa-
172 per. This example dataset consists of four subsets
173 named Breadth (B), Depth (D), Home (H), and
174 Urban (U). Each subset will be used to demonstrate
175 a different test. In Section 4, we show a use case:
176 applying a class of models known as Gaussian Pro-
177 cess (GP) Regression to the example dataset. Note
178 that the use of GP regression or the specific data-
179 set have no bearing on the test suite itself; they
180 are only convenient examples chosen in this paper
181 to illustrate the applicability of the tests. Finally,
182 we conclude with a discussion of possible use cases
183 and limitations of this approach. We discuss how
184 the tests are generally applicable, and the use case
185 is meant to serve as an example for how the pro-
186 posed test suite may work in practice.

187 2 Background

188 In this section we outline the mathematical back-
189 ground for the models and tests, as well as existing

work and how it relates to this paper.

BPS is best characterised as a non-linear, stochastic, causal system [2]. Linear regression models are popular emulators [e.g. 3] since they are easy to fit, use and interpret, but they could be inaccurate since the simulator they are trying to model is a non-linear system itself. For such non-linear systems, emulators based on non-linear regression are more appropriate.

2.1 Simulators

A typical simulation uses, as inputs, parameters such as building geometry; building envelope characteristics; Heating, Ventilation, and Air Conditioning (HVAC) system specifications; operation schedules and control strategies. We denote a set of inputs with a vector θ containing all these parameters. To assess the performance of a building design θ , a BPS simulation also requires plausible operating conditions that the building might experience, e.g., local weather and internal heat gains from lighting, occupants and equipment loads. We denote these operating conditions by a vector \mathbf{z} . The set of inputs, therefore, is equal to θ and \mathbf{z} which we denote by $\mathbf{x} := \{\theta, \mathbf{z}\}$.

Given an input vector \mathbf{x} , the goal of BPS is to estimate performance indicators, such as temperature trends, comfort indicators, energy demand, etc. In this paper, we use a common BPS problem: predicting the energy performance of a building as the sum of the hour-by-hour energy demand (power draw) over the year. This gives us a scalar energy output which we denote by y . Denoting the simulator by a function f that takes \mathbf{x} as input and outputs y , we can express the Simulator as:

$$y = f_s(\mathbf{x}), \quad \text{where } \mathbf{x} = \{\theta, \mathbf{z}\}. \quad (1)$$

Since the output depends on \mathbf{z} , the choices of operating conditions is extremely important. The future operating conditions are unknown, but can still be obtained using other sources of information. For example, future weather conditions can assumed to be random draws from a probability distributions $p(\mathbf{z})$ which can be estimated using past weather data [4, ch. 2]. A reliable prediction of performance can then be obtained using the Monte

Carlo (MC) method:

$$\hat{y}_{MC} := \frac{1}{N} \sum_{n=1}^N f_s(\boldsymbol{\theta}, \mathbf{z}^{(n)}), \quad (2)$$

where $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)}$ are N operating conditions drawn from a statistical model $\hat{p}(\mathbf{z})$ obtained using past weather data. The above quantity gives a reliable estimate of the energy performance because it accounts for a wide variety of operating conditions. A design process built upon such estimates results in a robust building design that is less likely to fail under extreme operating conditions.

This approach, the MC method, which might be time-consuming to complete with a simulator since each simulation may itself take hours. Even though multiple simulations may be run in parallel, the design process itself is sequential and iterative, necessitating the repetition of the N simulations for each design parameter vector $\boldsymbol{\theta}$ investigated. A more common practice is to use ‘typical’ (average, median, representative) values for \mathbf{z} and only perform one or two simulations. This would result in performance estimates that are low variance but are heavily biased towards typical operating conditions. A design process relying on such estimate will be less robust, since it might miss important operating conditions under which a building may perform poorly or even break down. Emulators are models of the simulator which can predict the outputs of simulations quickly and, therefore, enable the use of computationally-expensive techniques in several situations.

2.2 Emulators

Fitting or training emulators requires overcoming three principal challenges: the requirement of a large, varied, and representative database for training; the time and effort to specify the form and compute the parameters of the models; and inflexibility in real-world application, usually indicated by an inability to predict well on test sets [5, 6].

The emulators we discuss here are based on regression models and take the following form:

$$\hat{y} = f_e(\mathbf{x}), \quad (3)$$

where the emulator is represented by the function $f_e(\cdot)$, which must be estimated, and \hat{y} is an

estimate of the simulator output y . The best estimate of f_e is that function in the set of functions \mathcal{F} that minimizes a cost function, e.g., Mean Squared Error (MSE),

$$f_e^* = \arg \min_{f_e \in \mathcal{F}} \mathbb{E}_{p(y, \mathbf{x})} \left[(y - f_e(\mathbf{x}))^2 \right], \quad (4)$$

where $p(y, \mathbf{x})$ is the joint distribution of y and \mathbf{x} . Since this distribution is unknown, we approximate the expectation using the sample mean over input-output pairs for y and \mathbf{x} observed in practice. The set \mathcal{F} is usually the set of all continuous and differentiable functions.

The dataset used to build a regression model should be independent and identically distributed (*i.i.d.*), though this is difficult to achieve in practice. One naive method is to use past measurements of \mathbf{z} (if available, e.g., past weather data) as samples from $p(\mathbf{z})$ and run the simulator with these to obtain samples from the distribution of y (e.g., Figure 9). Specifically, given N measurements $\mathbf{z}^{(n)}$ for $n = 1, 2, \dots, N$, we can calculate the corresponding energy outputs $y^{(n)}$ by running the simulator.

The fitting and testing of regression models to training datasets consists of four steps, as outlined in Figure 1: (A) picking a subset of data of size $N_{\text{train}} < N$; (B) estimating the hyper-parameters ($\boldsymbol{\psi}$) with this training dataset; (C) picking an additional set of data of size N_{pred} for prediction; (D) using the model to predict on the test dataset and calculate error of prediction.

The structure of linear and non-linear functions, especially the GP regression models used in this paper, is discussed in Appendix A.

2.3 Existing Work

The existing BPS literature, including our previous work on GP regression [7], focussed on proving the utility of a specific regression method or emulator type to tackle a specific problem in BPS. A large variety of approaches have been proposed for a variety of outputs, and these can be broadly divided into two classes: ‘grey box’ models and regression-based models. The so-called grey box or ‘reduced order’ models, which use simplified physics to approximate performance [8]. These are simple to use but inflexible since a single grey-box model is

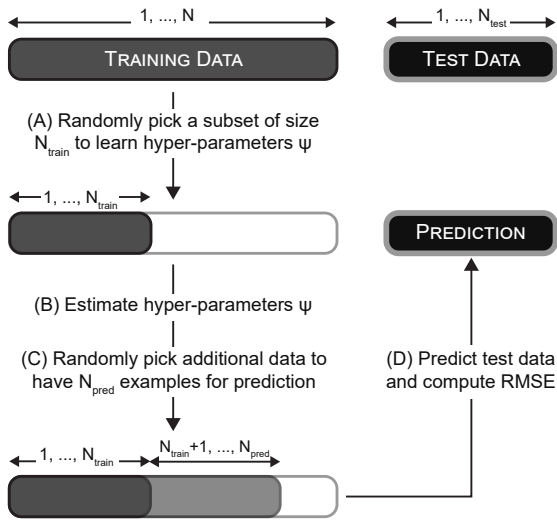


Figure 1: A schematic of the method used in this paper to train and test regression models. (A)-(D) represent steps in the procedure, while the boxes represent data sets. When testing on real-world data, we use the term validation dataset to denote the test dataset.

299 appropriate only for the specific system it approx- 300 imates. The second class of emulators is based on 301 regression models, constructed by fitting the model 302 to a database of inputs and outputs obtained from 303 a simulator or measured data [e.g., 3, 9]. Confus- 304 ingly, regression models are also sometimes referred 305 to as ‘reduced order’ models.

306 Most existing work on emulators to supplement 307 BPS uses regression-based models. The published 308 work largely focusses on characterising the rela- 309 tionship of some performance metric like energy 310 consumption for heating or cooling to paramet- 311 ers/properties that may be controlled by designers, 312 like insulation levels and window area [4, 8, 10, 11]. 313 Examples of these kinds of studies include those 314 that address:

- 315 (i) general uncertainty and sensitivity analysis for 316 performance analysis and what-if analyses [4, 317 12–14],
- 318 (ii) computational cost, such as energy optimisa- 319 tion for large-scale retrofits [9], grid-scale de- 320 mand prediction from buildings [15], bench- 321 marking [16],

- (iii) lack of sufficient information to run simulat- 322 ors reliably, e.g., prediction of the potential 323 to harvest solar energy for neighbourhoods 324 [17, 18], statistical evaluation of the energy 325 performance of different office designs [19–24], 326
- (iv) lack of certainty about future weather, e.g., 327 prediction of indoor conditions using a small 328 number of measured parameters [25], correlat- 329 ing probabilistic climate projections with of- 330 fice cooling demand and overheating analyses 331 in the UK [26–28], development of “climate 332 change amplification coefficients” to estimate 333 resilience [29], 334
- (v) calibration of building energy models, and 335 fault detection and control [30–32]. 336

3 Method 337

In this section we described the proposed tests in 338 detail and the dataset used to show how they might 339 be used in practice. These tests are based on gen- 340 eral principles of statistical learning outlined in 341 texts such as Hastie et al. [5], Rasmussen and Wil- 342 liams [33] 343

3.1 Proposed Test Suite 344

As described in Section 1, we use the following four 345 criteria to judge if, for a given problem, an emulator 346 is a good-enough replacement for a simulator: ac- 347 curacy, speed, generalisability, and ease of use. The 348 logic for these criteria relates to the use cases and 349 issues discussed in the preceding sections. There 350 are usually trade-offs between accuracy and gener- 351 alisability, accuracy and speed, and between ease 352 of use and the others. In addition, speed and ease 353 of use are somewhat subjective and based on the 354 problem at hand. 355

The tests that can be used to evaluate these cri- 356 teria are all described here in terms of *error*, i.e., the 357 difference between values (simulation outputs of inter- 358 est) predicted by an emulator and those output 359 by a simulator for the same input. We will also use 360 the concept of test and train datasets, i.e., data- 361 sets used to fit (train) a regression model and test 362 its performance. These tests are described for use 363 in the process of training a regression model for 364 a given problem. When an emulator is ‘deployed’ 365

366 for use in a given problem, it cannot be tested ex-
 367 cept by running a small number validation simula-
 368 tions. The use of error to train a regression model
 369 is outlined later in Figure 1, section 2.2, and ap-
 370 pendix A.2. We begin by summarising the tests
 371 here and describing their relevance to predicting
 372 building performance.

373 3.1.1 Accuracy

374 Accuracy implies that for a given set of inputs, the
 375 emulator predicts the simulator output well. The
 376 need to predict simulator outputs accurately is ob-
 377 vious: if an emulator is not sufficiently accurate,
 378 it cannot replace the simulator. Given that the
 379 simulator itself is an estimate of actual perform-
 380 ance during the lifetime of a building, introducing
 381 an unacceptably large additional error will degrade
 382 the utility of simulation-aided decision-making. We
 383 propose that the accuracy of a model during the
 384 training process can be assessed using the follow-
 385 ing tests.

Test 1: Error on a validation dataset is less than
 some acceptable tolerance of error, i.e.,

$$f_e^*(y - \hat{y}) \leq \varepsilon, \quad (5)$$

386 where y is the output of a simulator and \hat{y} is the
 387 output of an emulator for the same input, $f(\cdot)$ is
 388 some function to aggregate the differences between
 389 emulator and simulator outputs like Root Mean
 390 Square Error (RMSE) or Mean Absolute Error
 391 (MAE), and ε is some acceptable tolerance for the
 392 error.

393 The error tolerance is a decision for the practi-
 394 tioner, and it may vary based on the context. For
 395 example, small random errors in the performance of
 396 individual buildings in a large stock energy model,
 397 like the Breadth dataset, will make little difference
 398 to evaluating the effectiveness of large-scale applica-
 399 tion of retrofit measures. On the other hand, error
 400 tolerance for the Home dataset, consisting of
 401 a single-family home, would be considerably smal-
 402 ler. The tighter the design requirements and lower
 403 the average consumption, in general, the lower this
 404 tolerance would be.

Test 2: Error improves with increased training
 data, i.e., as the number of observations available

for training increase, the error on the validation
 set reduces. This can be expressed as an inverse
 correlation:

$$f_e^*(y - \hat{y}) \propto \frac{1}{N_{train}} \quad (6)$$

where N_{train} is the number of observations in the
 training dataset.

Since the training dataset must be obtained from
 a simulator, with its attendant computational cost
 and effort, the increased investment must be justi-
 fied by improvement in prediction. When applying
 these tests to our sample dataset in Section 4, we
 will show how the return on investment can dimin-
 ish as the size of the training dataset increases.

3.1.2 Generalisability

Test 3: Error improves with more complex or
 varied training data:

$$f_e^*(y - \hat{y}) \propto \frac{1}{\sigma_{\mathbf{x}}^2} \quad (7)$$

where $\sigma_{\mathbf{x}}^2$ is the variance of the features or inputs
 corresponding to the observations in the training
 dataset.

This test is a check against over-fitting to data-
 set representing a narrow set of inputs, potentially
 unrepresentative of the problem. For example, a
 model trained entirely on variations in window-to-
 wall ratio while everything else is held constant,
 like part of the Home dataset, is unlikely to estim-
 ate changes in insulation levels accurately. If the
 prediction on the test set improves as the training
 dataset includes more building or system options,
 designs, or scenarios, then the user has more confid-
 ence that the emulator will more accurately predict
 over the variety of designs and scenarios.

The training set must represent the problem to
 be explored. When a problem is narrowly-defined,
 i.e., only a limited aspect of design or uncertainty
 is to be explored, this test may not matter.

Test 4: Error does not degrade too much when
 moving from predicting on a general dataset during
 training to a validation set that is more specific to
 the problem at hand. In the context of this paper,
 this would mean training on the Breadth dataset,
 and predicting on the Depth dataset. When the

440 change/increase in error is too much depends on
 441 the magnitude of the initial training error relative
 442 to that of the average prediction and the problem
 443 being studied. If an emulator fails Test 1 on a spe-
 444 cific validation set for example, i.e., the error sur-
 445 passes the tolerance set by the user, the emulator
 446 is likely unsuitable.

A dataset is more specific if it deals with only one aspect of a given design problem or exercise. For example, after training a model on variations in layouts and several building systems (Breadth dataset), we task the emulator to predict only on variations of one system or component in a specific building (Depth). If the initial error,

$$f_e^*(y_1 - \hat{y}_1) \leq \varepsilon,$$

447 then,

$$f_e^*(y_2 - \hat{y}_2) \leq \varepsilon, \quad (8)$$

448 where y_1 , y_2 are outputs from two different valid-
 449 ation datasets and y_2 is the result of testing on a
 450 validation dataset \mathbf{x}_2 more specific to a problem
 451 than \mathbf{x}_1 , the corresponding validation dataset for
 452 y_1 .

453 **Test 5:** Error does not degrade too much when
 454 testing on a more complex validation dataset, i.e.,
 455 one with more variety of inputs. If the initial error,

$$f_e^*(y_1 - \hat{y}_1) \leq \varepsilon,$$

456 then,

$$f_e^*(y_2 - \hat{y}_2) \leq \varepsilon, \quad (9)$$

457 such that,

$$\sigma_{\mathbf{x}_2}^2 > \sigma_{\mathbf{x}_1}^2,$$

458 where y_1 , y_2 are outputs from two different valid-
 459 ation sets, \mathbf{x}_1 and \mathbf{x}_2 respectively, and $\sigma_{\mathbf{x}}^2$ is the
 460 variance of an input validation set. Here, \mathbf{x}_1 may
 461 represent a single archetype building used for train-
 462 ing a model while \mathbf{x}_2 would represent a portfolio of
 463 buildings that should conform to the same arche-
 464 type but with a variety of designs and systems. For
 465 example, creating variations on one archetype office
 466 building by varying the properties of different sys-
 467 tems as in Depth dataset.

3.1.3 Speed 468

Test 6: Emulator is computationally cheaper 469
 than a simulator. Regardless of the computing 470
 infrastructure being used, an emulator should be 471
 cheaper, and therefore faster, to run in order to 472
 justify accepting the increased error in estimation 473
 of output. If an emulator performs particularly well 474
 on this test, it may also be suitable for implement- 475
 ation in Building Management (Automation) Sys- 476
 tems (BMS) controllers for Model Predictive Con- 477
 trol (MPC) and similar low-resource applications 478
 that require rapid response. The computational 479
 complexity of the emulator, 480

$$\mathcal{O}_e(aN_{\text{val}}^b) \leq \mathcal{O}_s(cN_{\text{val}}^d), \quad (10)$$

where $\mathcal{O}_e(aN_{\text{val}}^b)$ is the complexity of the emu- 481
 lator that scales with the number of observations 482
 in the validation dataset with some exponent, and 483
 $\mathcal{O}_s(cN_{\text{val}}^d)$ is the complexity of the simulator pre- 484
 dicting over the same observations. 485

3.1.4 Ease of use 486

Test 7: Emulator is appreciably better than sim- 487
 pler methods. This final test ensures that the 488
 simplest method that delivers adequate perform- 489
 ance is used. As discussed in Figures 10 and 11 490
 and appendix A.1, more complex models will gen- 491
 erally tend to overfit to a given dataset. The error 492
 of an emulator 493

$$f_{e1}^*(y - \hat{y}) \leq f_{e2}^*(y - \hat{y}), \quad (11)$$

where $f_{e1}^*(\cdot)$ is a higher-order model (more complex, 494
 more parameters) than $f_{e2}^*(\cdot)$. Models with fewer 495
 parameters are both cheaper to train, so would also 496
 be quicker to train and deploy (Equation (10)). 497

3.2 Data 498

To demonstrate the use of the test suite proposed 499
 in this paper, we generated a large dataset by sim- 500
 ulating numerous combinations of buildings and 501
 weather conditions from four different simulation 502
 exercises, labelled ‘Breadth’, ‘Depth’, ‘Home’, and 503
 ‘Urban’. This dataset does not comprehensively 504
 represent all the possible use cases for emulators 505
 discussed in Section 1, and it does not need to. In- 506
 stead, the complete dataset gives enough variety of 507

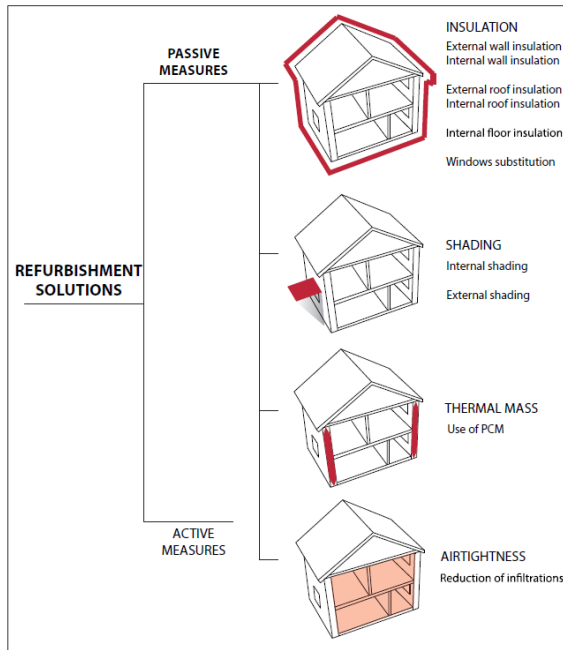


Figure 2: *The envelope variations simulated for the single-family home (H). Original figure in Rastogi [4].*

508 use cases to demonstrate the test suite proposed in
 509 this paper. Specific subsets are also used for in-
 510 dividual tests, e.g., Breadth and Depth for Test 4
 511 (Section 4.4).

512 These combinations results in approximately
 513 600,000 data points, details of which are in Table 1.
 514 Each simulation calculated the space heating and
 515 cooling requirement for the given combination of
 516 building design and weather. The total energy for
 517 heating and cooling was divided in each case by
 518 the area of the building (normalised by area) to
 519 make the outputs of all the simulations compar-
 520 able. These simulations are also described in detail
 521 in separate publications and summarised in Rastogi
 522 [4, sec. 4.2]⁴.

523 The size of the datasets used in this paper is
 524 not an indication of the minimum sizes required
 525 for each test. As the figures show in Section 4, the
 526 result for a test is usually obvious using a fraction
 527 of the data shown in this paper for demonstration.

528 The Breadth (B) and Depth (D) case studies
 529 are from the United States Department of Energy

(USDOE) commercial buildings reference database 530
 [34]. The Breadth dataset consists of 16 different 531
 building types. This dataset is a representation of 532
 a national or regional building ‘stock’, i.e., a repres- 533
 entative sample of buildings. The Depth set consists 534
 of simulations on one of the building types in 535
 the USDOE database: the ‘medium office’. This 536
 dataset is a simulation of a design exercise vary- 537
 ing envelope properties. These properties were: U- 538
 value (practically, changing the thickness of insu- 539
 lation material), thermal mass (nominal quantity 540
 of internal mass), shading (obtained by varying the 541
 depths of overhangs and fins), permeability (chang- 542
 ing infiltration levels), and transparency (chang- 543
 ing the Window-to-Wall ratio). 544

545 The Urban case (U) is composed of a set of build- 546
 ings constructed over a century (1900-2010) in the 547
 centre of Geneva, Switzerland [35]. This case was 548
 chosen because it includes the urban context sur- 549
 rounding the buildings. We expect the urban con- 550
 text to add noise to the data, since we have not 551
 included any features that explicitly describe the 552
 influence of the surroundings. The buildings were 553
 all modelled in at least two variants: with the ori- 554
 ginal envelope and with an envelope upgraded to 555
 the latest Swiss standards for infiltration, insula- 556
 tion, etc.

557 Finally, the single-family home case (H) is an ex- 558
 ample of a very simple simulation study, one where 559
 we expect the response to be characterized well- 560
 enough by a linear regressor. The changes to the 561
 house are described in Figure 2. The simulation 562
 model of the house is based on an actual home in 563
 north-central Germany [36].

564 The weather data used is of three types: recor- 565
 ded data from the Integrated Surface Database⁵, 566
 typical years [38, 39], and synthetic weather time 567
 series [4, 40, 41].

⁴The data may be downloaded from <https://doi.org/10.5281/zenodo.291858>

⁵<https://www.ncdc.noaa.gov/isd>

Table 1: List of datasets. The size indicated here may be slightly different from the amount of data used in the scripts due to the presence of invalid data entries. Less than 5% of the entries were invalid in any dataset. Mean values for annual sum of Heating and Cooling loads were calculated over all valid values in kWh/m².

Name	Size	Description	Ref.	Mean Heat	Mean Cool
Breadth (B)	88,242	USDOE commercial reference buildings (all building types)	[4, 34]	131.63	48.97
Home (H)	77,934	Single-family home, Central European construction	[4, 36, 37]	23.41	74.84
Depth (D)	445,334	Variations on the medium office building from the USDOE database		102.90	30.98
Urban (U)	6,003	Mixed-use buildings in Geneva, Switzerland (with surrounding buildings)	[35]	134.27	0.00

4 Results

In this section we show the application of the proposed test suite to judge the suitability of GP regression models with linear and non-linear kernels, using the dataset described in Section 3.2. The kernels are listed and described in Table 2.

The fitting and testing of GP models to large training datasets consists of four steps, as outlined in Figure 1: (A) picking a subset of data of size $N_{\text{train}} < N$; (B) estimating the hyper-parameters (ψ) with this training dataset; (C) picking an additional set of data of size N_{pred} for prediction; (D) using the model to predict on the test dataset and calculate error of prediction. Before beginning the fitting procedure, we set aside a portion of the total data available to us as ‘test’ data of size N_{test} (approx. 60% of the dataset). The amount of training data N_{train} used to estimate the hyper-parameters (ψ) was varied from 50 to 4000 for each dataset and model (except for models fit to the Urban dataset individually). For a given training size N_{train} , we draw that many observations from the large overall training set and use it to train a model. We repeat this process 100 times to obtain an empirical distribution of the RMSE of each model’s predictions on the test set. For Test 3 (Section 4.3), results from fitting to datasets much larger than 4000 are presented (e.g., Figure 5). In this case, for runs with $N_{\text{train}} > 2000$, only $N_{\text{train}} = 2000$ was used to calculate hyper-

parameters, while a separate subset of the training dataset, N_{pred} , was used to fit the model (see Appendix A.4.2 for details). In this case, both N_{pred} and N_{train} are training datasets. This allowed for models fit to much larger datasets than our computer could handle if the entire process were carried out with these very large datasets.

We use Root Mean Square Error (RMSE), the square root of the MSE term from Equation (4), to quantify errors. An advantage of RMSE is that it has the same units as the original outputs, in this example kWh/m², which makes it easier to understand and judge magnitudes of error. Readers are invited to compare plotted errors against the means of the datasets given in Table 1, and against the overall mean of 51.49 kWh/m² and 64.31 kWh/m² for annual heating and cooling loads respectively. In the figures, we present the distribution of the RMSE calculated over a hundred subsets of the test set as an additional check on the reliability of a model: if the distribution of RMSE over subsets of a test set is too wide, the model does not reliably represent the range of possible values of inputs to be tested, i.e., the test dataset. The test dataset is entirely separate from the training dataset (Figure 1). The RMSE plotted in each graph below was calculated solely on the test dataset, i.e., test error.

Each test is presented as it would be applied to the problem of quantifying the uncertainty in predicted energy use of a building design for space conditioning due to lack of knowledge about fu-

Table 2: List of models compared in this study.

Model	Description
Mean	Mean of the outputs y_n
Lin-ISO	Linear model with isometric kernel
Lin-ARD	Linear model with automatic relevance determination
NonLin-ISO	Squared-exponential isometric kernel
NonLin-ARD	Squared-exponential kernel with automatic relevance determination

629 ture weather conditions. Each design is intended
630 to represent a choice available to a designer at construction or renovation, and the output of interest
631 is the impact of the design choice on annual whole-
632 building energy performance over the lifetime of the
633 building. The uncertainty would be quantified using
634 the MC method, sampling plausible weather scenarios
635 and simulating them for a given building design to
636 obtain a reliable estimate of energy performance
637 (Equation (2)). While there is no definitive number
638 of operating conditions (weather files) that must
639 be simulated for an estimate to be reliable, we have
640 found that stable results could be obtained with
641 about a hundred simulations with random weather
642 conditions per building design. Depending on the
643 complexity of the building designs and systems used
644 in this paper, each simulation took 15 minutes to
645 almost 2 hours.

647 4.1 Test 1: Error on validation set

648 The acceptable level of error depends on the problem
649 being explored. Fitting regression models always
650 involves a trade-off between minimising the error
651 on the available training dataset and ensuring that
652 this does not overfit the model to the specific
653 dataset. This is a form of the bias-variance trade-
654 off common to all statistical learning approaches
655 [5], as discussed in Section 2.1.

656 4.2 Test 2: Error with larger training dataset

658 This test is unambiguous and straightforward for
659 this example: the validation error improves with
660 size of training set for all cases described in Table 1.
661 For example, see the changes in prediction error
662 when predicting heating loads on the Breadth and

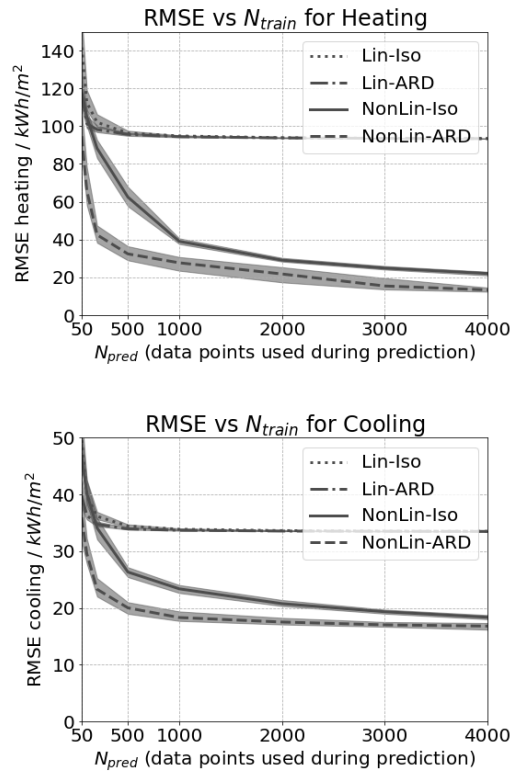


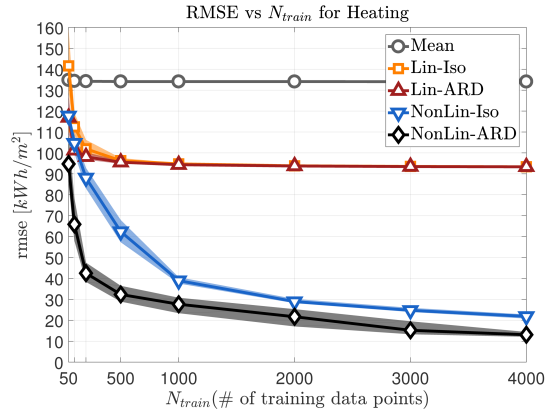
Figure 3: RMSE for heating [top] and cooling [bottom] models fit to the Breadth database. The (test) errors are calculated on approx. 52,940 points, and plotted against the size of the data set used at prediction (same as the dataset used to learn hyper-parameters in this case, i.e., $N_{pred} == N_{train}$). The lines indicate median errors, and the filled areas are bounded by the 25th and 75th percentiles.

663 Home datasets, plotted in Figure 4 for each model
 664 type. Compare the errors with those from the
 665 simplest possible predictive model: the mean of the
 666 training outputs, \bar{y}_{train} .

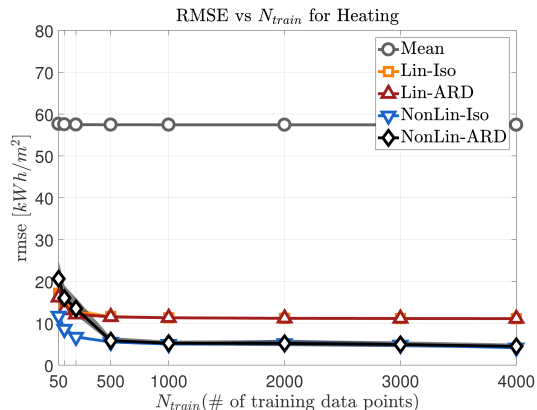
667 4.3 Test 3: Error improves with 668 training data variance

669 For this test, we combined all datasets to represent
 670 the case where a complex, varied problem is to be
 671 modelled and large number of simulation runs are
 672 available. Ideally, we would like to use all the data
 673 available for training. However, learning hyper-
 674 parameters is infeasible for large numbers of data
 675 points because each step in the learning procedure
 676 requires the inversion of an $N_{\text{train}} \times N_{\text{train}}$ matrix.
 677 We use a simple trick to make use of the additional
 678 data, outlined in Figure 1 and based on the pro-
 679 posal in Chalupka et al. [42]. We learn the hyper-
 680 parameters on a dataset of size $N_{\text{train}} = 2,000$,
 681 but during prediction we use a much larger data-
 682 set $N_{\text{pred}} = 2000, \dots, 12000$. Since prediction re-
 683 quires only one matrix inversion, the latter step is
 684 still feasible for dataset sizes of about 10,000 on the
 685 hardware we used for our study.

686 In Figure 5, we present the results of fitting and
 687 testing a non-linear GP model to the whole dataset
 688 (i.e., all subsets described in Table 1). Thus the dif-
 689 ferences between the results presented in Figure 3
 690 and those presented in Figure 5 are that, firstly,
 691 in the latter figure we use a larger dataset that
 692 is more representative of real-world outcomes, and
 693 secondly, we show the additional advantage of us-
 694 ing a larger dataset for prediction during training
 695 ($N_{\text{pred}} > 0$ in Step C of Figure 1). In Figure 3,
 696 the RMSE values were plotted against the num-
 697 ber of data points used to learn hyper-parameters
 698 and subsequently fit the model to the same data-
 699 set. In Figure 5, the RMSE values are plotted
 700 against the size of the prediction dataset (differ-
 701 ent from the fixed number used in learning hyper-
 702 parameters). The solid lines (with shaded curves)
 703 show the RMSE obtained when 2,000 data points
 704 are used for learning but a larger set is used dur-
 705 ing prediction. We see that RMSE decreases as
 706 the number of data points used for prediction is
 707 increased. Linear models are not presented here
 708 since this procedure makes no difference to their
 709 performance.



(a) Breadth



(b) Home

Figure 4: Evolution of RMSE for heating loads with increasing N_{train} . Non-linear models perform better than linear models for all datasets, and also show more improvement with increased training data. A separate, larger dataset was not used for prediction in this case since the sizes of N_{train} were still tractable.

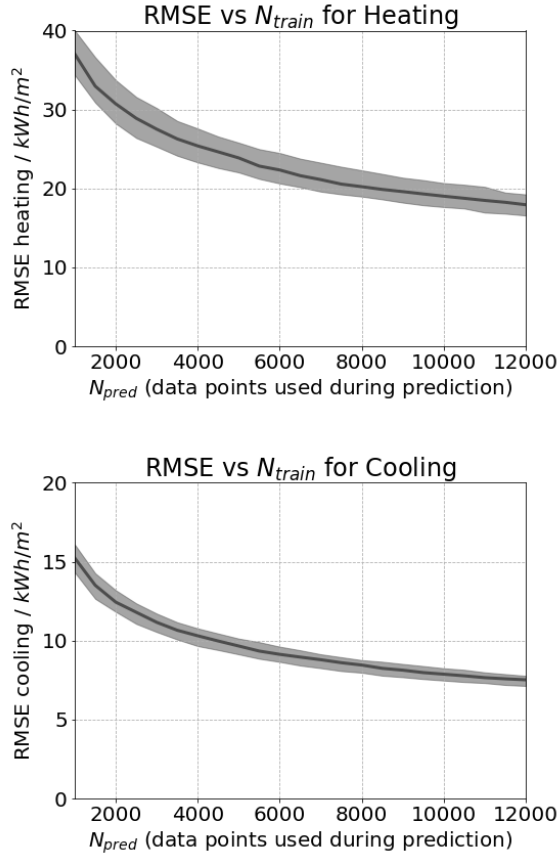
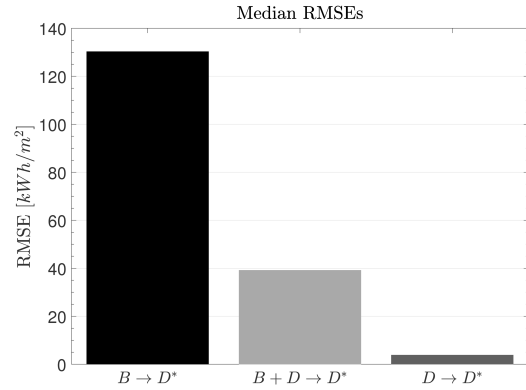


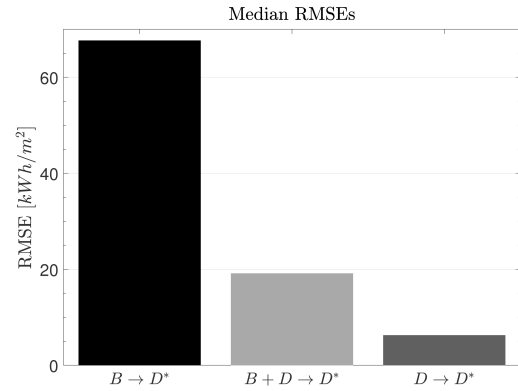
Figure 5: These figures show the RMSE for heating and cooling loads for the combined dataset, plotted against the size of the training data set used at prediction (Appendix A.4.2). The lines indicate median errors, and the filled areas are bounded by the 25th and 75th percentiles. The non-linear models outperform linear models, and the predictions improve with size of training data set.

We see that the use of non-linear GP models on a large dataset is both feasible and accurate. We will now check the performance of these models on two sets of problems: predicting for a specific building (individual) or on a variety of buildings together.

4.4 Test 4: Error on a specific validation set



(a)



(b)

Figure 6: Heating [top] and cooling [bottom] predictions when predicting on a specific building. $B \rightarrow D$ means training hyper-parameters on the Breadth dataset and testing on Depth, $B + D \rightarrow D$ implies training on a combination of the two and predicting on Depth, while $D \rightarrow D$ implies both training and testing on exclusively the Depth dataset.

In many applications, the designer might be interested in predicting the performance of only a

719 specific type of building. In that case, it is possible
720 that an emulator trained on a variety of building
721 types may not perform well. We show that a non-
722 linear model performs well if it has ‘seen’ enough
723 buildings that are similar to the one we want to pre-
724 dict. The results presented in this section establish
725 the satisfactory performance of the GP regression
726 models used here for this test.

727 We consider the task of predicting the perform-
728 ance of a building in the Depth dataset, the results
729 of which are shown in Figure 6. This dataset con-
730 tains only a specific type of building: a medium-
731 sized office. We used 60% of the Depth dataset as
732 the test set ($N_{\text{test}, D^*} = 247, 304$). We trained three
733 models: the first using only the Breadth dataset
734 ($N_{\text{train}} = N_{\text{pred}} = 1000$), referred to as ‘ $B \rightarrow D^*$ ’ in
735 Figure 6; the second using only the Depth dataset
736 ($N_{\text{train}} = N_{\text{pred}} = 200$), ‘ $D \rightarrow D^*$ ’; and the third
737 using a combination of the Breadth and Depth
738 datasets (1000 data points from Breadth and 200
739 data points from Depth), ‘ $B + D \rightarrow D^*$ ’.

740 We see that when using only the Breadth data-
741 set for training and prediction, the model performs
742 badly. Adding some points from the Depth data-
743 set significantly improves the performance of the
744 model, since this addition reduces the influence of
745 those points in the Breadth set that do not come
746 from the office building. Using only the Depth
747 dataset gives the best performance. This shows
748 that training an emulator with a small number of
749 examples for a specific building is enough to pre-
750 dict well for that building ($N_{\text{train}} = 200$ in this
751 case). When learning on a variety of buildings
752 and predicting on a specific one, the performance
753 is poor with naive selection of points during pre-
754 diction [left-most bar in each graph]. Performance
755 improves by adding simulations from the specific
756 building [middle bar]. Performance is best when
757 using training data only from that building [right-
758 most bar].

759 A GP model predicts on a new test point by
760 correlating the test inputs to training inputs. The
761 kernel is supposed to encode the influence of differ-
762 ent training points in predicting the test point: the
763 more closely related a group of training points is to
764 a test point, the more influence they ought to have
765 on the prediction. We find that this is not the case
766 for our study when we try to use training points
767 exclusively from the Breadth dataset to predict on

a test set from the Depth dataset. However, the
768 results improve dramatically if new training points
769 are added from the Depth case. This suggests that
770 explicitly encoding, perhaps with a categorical vari-
771 able representing building type/usage, the ‘close-
772 ness’ of a new test point to a subset of training
773 points, should improve prediction. 774

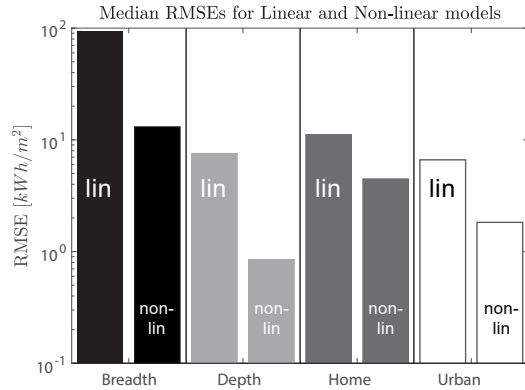
4.5 Test 5: Error on a varied valida- 775 tion set 776

777 We now present results to show that the non-linear
778 emulators presented here can be trained to obtain
779 accurate predictions for a variety of datasets, e.g.,
780 different buildings in different climates.

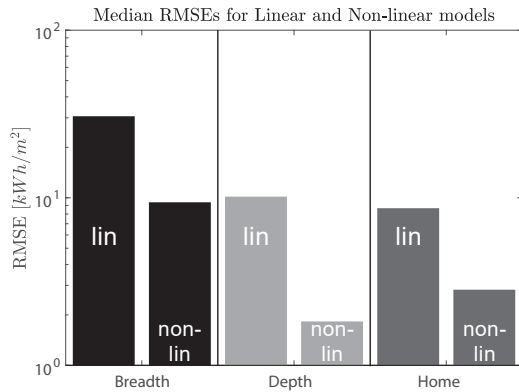
781 Figure 7 shows the results for all datasets separ-
782 ately using $N_{\text{train}} = N_{\text{pred}} = 4000$ data points, of
783 which 60% are used as the test set for each case. We
784 present results for two models: Linear automatic
785 relevance determination (ARD) (*lin*) and Squared
786 Exponential ARD (*non-lin*). The RMSE obtained
787 from using the non-linear model is of the same or-
788 der of magnitude for all datasets (1-10 kWh/m²)
789 and uniformly better than RMSE from linear mod-
790 els. These results demonstrate the flexibility of the
791 non-linear model compared to an equivalent linear
792 model.

4.6 Test 6: Computational Expense 793 and Time 794

795 Figure 8 shows the distribution of computer time
796 taken by each simulation in the dataset used for
797 this paper. A single iteration of the Monte Carlo
798 method for a single design, i.e., about 100 simula-
799 tions, would take between 25 and 200 hours (0.9e5
800 to 7.2e5 s). In addition, the time taken for each iter-
801 ation would not change, since the simulations can-
802 not be reused. For comparison, GP regression takes
803 1 microsecond (1e-6 s) to provide one estimate, so
804 each MC iteration would take about 10 milliseconds
805 (1e-3 s). This comparison is for run times, i.e., as-
806 suming that the regression model has already been
807 fitted. As we saw in Sections 4.1 to 4.4, reaching
808 a satisfactory error rate in this example requires a
809 training dataset of at least 1000-1500 simulations.
810 Assuming a fresh start for each problem, a worst
811 case scenario where no knowledge from previous
812 simulations is transferable to a new problem, would



(a) Heating



(b) Cooling

Figure 7: Median RMSE at $N_{train} = 4000$. The y-axis uses a log scale. For each dataset, the bar on the left is for the linear model and on the right for the non-linear model. Non-linear models perform than linear models for all datasets. The Urban case was modelled for Geneva, where buildings do not typically include cooling systems (air-conditioning).

require the user to use the simulator for some part 813
of the experiment. 814

When accounting for the cost of using an emulator, 815
both the cost of obtaining training data (usually 816
from the simulator) and of fitting the model 817
should be included. This means that emulators are 818
not suitable for short and quick design exercises un- 819
less that exercise is part of or similar to a problem 820
for which a model has been trained. This creates a 821
strong incentive for the use of pre-trained libraries 822
of models suitable for specific problems, especially 823
in situations where there is insufficient information 824
to create simulation models (Section 2.3). 825

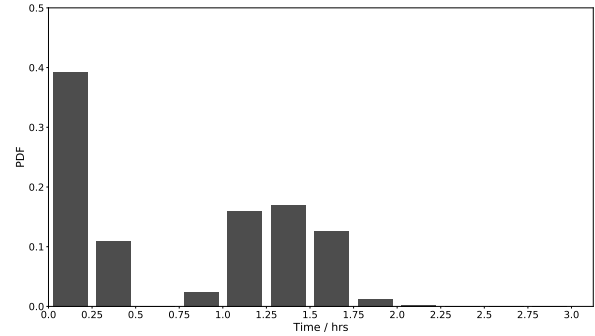


Figure 8: Probability Distribution Function (PDF) of the approximate time taken to run one simulation in a simulation exercise carried out by the authors using EnergyPlus v8.8 (<https://energyplus.net/>).

4.7 Test 7: Simplest emulator 826

Non-linear models outperform linear models in all 827
of the results presented here (Figures 3 to 5), and 828
both model types outperform the mean of the training 829
set outputs. The predictions also improve with 830
size of training data set. Non-linear GP-based emu- 831
lators perform equally well when predicting on a 832
diverse set of buildings (Breadth) as a dataset con- 833
sisting of a single building (Depth, Home) or small 834
set of very similar buildings (Urban). 835

5 Discussion 836

In this section we discuss the applicability of this 837
approach, limitations of data-based approaches for 838
emulators, and practical issues around selecting 839
data for training. 840

841 5.1 Generalisability of Approach

842 Designers and analysts of buildings and their systems develop professional judgement and intuition
843 about the physics of the systems they study. This allows experts to identify common errors in numerical
844 simulations and improves the reliability of results. The approach proposed in this paper pre-
845 supposes some knowledge of how appropriate datasets may be acquired. Data-based methods are not
846 substitutes for knowledge about the physics of the problem. A tool like regression, which does not
847 use the physics of BPS, will require users to learn new skills in handling and interpreting statistical
848 learning models. An example of this is the difficulty of physically interpreting the dimensions of
849 the (mathematical) space of features, i.e., the numerical representations of the values of different
850 building properties or characteristics along numerical axes. When collecting data on complex problems,
851 the characteristics or features of interest may be many, which makes it difficult to maintain intuition
852 about which designs (combinations of features) are similar to others in the feature space.
853 This cannot easily be overcome with human judgement, a problem we encountered in the creation of
854 the dataset used here as well. Thus, the use of these techniques will benefit from the development of services
855 and tools which suggest methods for efficient data-gathering and training. While some progress
856 has been made recently in moving more simulation programs to cloud-based services that remove
857 much of the complexity of setting up models for the user, the possibility of augmenting these services
858 with regression models to improve their utility for computationally-intensive problems remains to be
859 explored.

860 The tests were demonstrated here using a single output: energy use for space conditioning. The use
861 of these tests could be more complicated for multiple outputs of interest, e.g., comfort and energy
862 use. Using multiple outputs could be handled with additional ranks or weights for different priorities,
863 combining the results of the tests for different outputs for a single decision. However, given that the
864 tests are comparing outputs from the *same physical system*, it may not always be the case that
865 the results of applying these tests to different outputs would be different. Additionally, since the user
866 would already have invested the effort to fit mod-

els at that point, they can also use different models for each output. There is no reason to suppose that
different emulators that work for different outputs would deliver inconsistent decisions.

847 5.2 Limitations of Data-driven Methods

848 Data-driven methods are not fail-proof; the model learnt on one dataset does not necessarily translate
849 perfectly to another (Sections 3.1.2, 4.4 and 4.5). An emulator does not incorporate any knowledge
850 about the physics of the problem being simulated, which means that emulators are, by construction,
851 usually less flexible than the simulator. Regression inputs may not be representative or may not explain
852 the variation in the data properly, which leads to inaccurate predictions and the inability to generalise.
853 Finally, the dataset used for training might itself have a systematic bias. That is, a dataset that does
854 not represent the problem properly, or is not a good proxy for real-world problems.

855 The best regression model contains just the right predictive inputs, for the selection of which there are
856 no fixed rules that will apply to every problem. The automatic relevance determination (ARD) procedure
857 [33, sec. 5.1, and references therein] used in this paper allows the user to begin with a large set of
858 input variables that might be important to a problem, letting the GP estimation procedure trim that number.
859 However, ARD does not necessarily follow the physics of the problem either and may not, therefore, generalise
860 to other problems. It is important to include all of the design parameters (input variables) which are
861 expected to be relevant in the modelling of the energy performance (output). This is both to ensure
862 a good fit to data and relevance to the design problem. Including too many inputs, however, makes it
863 difficult to obtain a good fit with a manageable size of training dataset (the so-called *curse of dimensionality*).

864 5.3 Collecting Data for Training

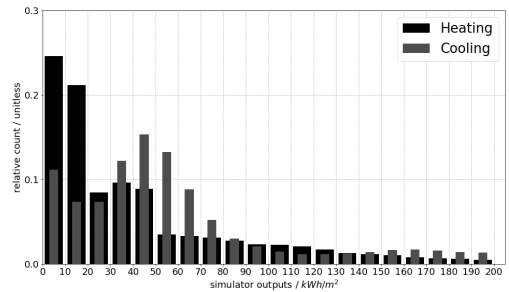
865 The quality of a dataset is determined by how faithfully it represents the true distribution of data
866 $p(\mathbf{x}, y)$ (assuming this distribution exists). However, since the true distribution is usually unknown,

935 it is impossible to measure the quality of a data-
 936 set objectively. To obtain a good-quality dataset
 937 of simulations of building designs, we can rely on
 938 designers (domain experts) who may, for example,
 939 choose realistic building designs and weather con-
 940 ditions from their portfolio and experience. The
 941 quality of dataset is, therefore, defined by the goal
 942 of the designer. For example, a designer interested
 943 in predicting performance over a variety of design
 944 parameters and weather conditions might need to
 945 acquire a dataset that contains similar examples in
 946 the training set. On the other hand, a designer who
 947 is interested in studying only a particular type of
 948 building might want to limit the dataset to that
 949 particular building. The results in this paper show
 950 that if the training dataset is too general, i.e., it
 951 contains too many examples dissimilar to the one
 952 considered during prediction, then the emulator
 953 does not perform well. This underlines the need
 954 to train separate models on different design prob-
 955 lems. Thus, the objectives of Test 3 are often in
 956 conflict with the objectives of Tests 4 and 5.

957 In principle, there should be no correlation with,
 958 or effect of, building type, usage, or location on the
 959 input variables θ . However, if the training sample is
 960 taken from realistic buildings in a given climate, the
 961 distribution of the input variables will be influenced
 962 by prevailing architectural idioms and other cul-
 963 tural and practical factors. For example, sampling
 964 many houses from a single region will not necessar-
 965 ily cover all possible values of wall conductance (U-
 966 value), since houses from a region typically follow
 967 local trends and laws concerning insulation level
 968 [e.g., 43]. Similarly, while modern office buildings
 969 may have up to 85% window-to-wall ratio (WWR),
 970 homes with the same proportion of window area are
 971 rare. At the other extreme, a value of less than 10%
 972 WWR is theoretically possible for any kind of build-
 973 ing, but windowless buildings are so rare as to be
 974 statistically insignificant. Sampling from a particu-
 975 lar type of building does not necessarily mean that
 976 the distribution of input variables will be identical
 977 to sampling from a different type/usage. In addi-
 978 tion, each type/usage has cultural or regional limits
 979 on the values of input variable seen in practice, e.g.,
 980 buildings will probably not include a layer of insu-
 981 lation in the walls when it is not appropriate for a
 982 given climate.

983 An example of selection bias is in the Breadth

984 dataset used in this paper (Figure 9). There is
 985 a preponderance of lower values of annual heating
 986 and cooling energy usage because the set has more
 987 moderate climates than extreme ones. The climates
 988 were selected based on data availability [4], prior-
 989 itising cities with several years of recorded data.
 990 These tended to be urban areas with major airports
 991 in continuous operation for decades either due to
 992 large, established populations or strategic reasons.
 993 The cities in the database have a combined popu-
 994 lation in excess of 200 million, though future work
 995 should incorporate weather from a wider selection
 996 of world climates.



997 **Figure 9:** *Distribution of simulator outputs (heating and cooling loads) from the overall dataset. There are more moderate loads than extreme ones in the dataset.*

997 These limitations in data quality and represent-
 998 ativeness can be overcome by including a large
 999 amount of simulated data and/or updating the
 1000 models using measured data. The flexibility of
 1001 regression-based emulators, as demonstrated in this
 1002 paper, means that different datasets can be easily
 1003 integrated into the model to improve its results.
 1004 This is in contrast to the original simulator, where
 1005 the results of one simulation do not have any im-
 1006 pact on the results of another.

1007 6 Conclusion

1008 This paper has proposed a new test suite for stand-
 1009 ardising the evaluation of emulators as suitable re-
 1010 placements for building performance simulators in
 1011 a variety of use cases, especially uncertainty quan-
 1012 tification. The use of emulators is promising for
 1013 applications where the speed of response from an
 1014 evaluation is important, provided the emulators are

1015 sufficiently accurate. Thus, the test suite is presented in the context of evaluating an emulator using
1016 four criteria: accuracy, generalisability, speed, and
1017 ease of use. We do not propose a specific model
1018 or class of models for the dataset used here, or
1019 any real-world problem exemplified by this dataset.
1020 Rather, we show how emulators may be evaluated
1021 in a given context, regardless of the structure of the
1022 problem or the dataset used to characterise it.
1023

1024 As an example of how the test suite may be used,
1025 we showed that non-linear models noticeably and
1026 consistently outperform linear models in emulating
1027 example specific and broad datasets. In all cases,
1028 the non-linear regression models show a Root Mean
1029 Square Error (RMSE) between 10-15% of the mean
1030 model (output from a model which consists of only
1031 one term: the mean of the training data). The
1032 GP regression models are able to predict well on a
1033 dataset consisting of a variety of buildings as well
1034 as a dataset consisting of a specific building. We
1035 find that the predictive performance of non-linear
1036 GP regression models is stable and repeatable. We
1037 showed procedures to use large datasets for learning
1038 and predicting with the same models on unseen
1039 data.

1040 Not all steps of a typical process require simulation,
1041 as designers make several decisions based
1042 on meeting existing laws, user needs, and functional
1043 requirements. The use of numerical simulation has
1044 expanded considerably with the advent
1045 of simulation tools or workflows offered as a service
1046 to non-specialists looking to carry out specific
1047 analyses [e.g., 44, 45], better diffusion of numerical
1048 and computational skills, and better interoperability
1049 between the models created by different professions.
1050 However, over-reliance on simulation tools for
1051 prediction rather than comparative what-if analyses,
1052 and excessive trust in results based on testing
1053 under limited operational conditions, e.g., typical
1054 weather files, can lead to a severe gap between
1055 expected and actual performance. The quantification
1056 of a possible cause of this gap can be partially
1057 addressed through the use of regression-based
1058 emulators. The use of these emulators can, in turn,
1059 become more systematic and widespread with the
1060 adoption of standard operating procedures such as
1061 the test suite proposed in this paper.

A Regression Models 1062

This appendix discusses the mathematical background of linear and non-linear regression models. 1063
Details are also included on the structure of GP regression models, how they may be fit to data and 1064
used, and a practical workaround for big datasets. 1065
1066
1067

A.1 Linear and Non-linear Regression Models 1068 1069

Non-linear regression models are more flexible and 1070
have the potential to model the output of a non-linear system more accurately. They can also account for the complex interactions of the large 1071
number of inputs that determine the outputs of a simulation. However, fitting non-linear models is computationally challenging, especially when a 1072
large amount of data is available. Ironically, a large amount of data is almost essential to obtain a good performance of the non-linear model, otherwise they might over-fit the data in hand [5, 6]. Another issue is that specification of non-linear models is difficult and requires a lot of effort and domain expertise. That is, for some problems where adequate data cannot be obtained within the budgetary or time allocation of an exercise, or the number of properties or factors for each test subject is limited because of the quality of data, non-linear models may not work [e.g., 3]. Some of these challenges are described in Section 3.2, exemplified by the data collected for this paper. 1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090

Non-linear regression models use more parameters than linear models, which increases both the time to build the model and the size of the training data set required to calculate the parameters of a model. For example, a non-linear model such as an Artificial Neural Network (ANN) could have millions of parameters, which means that the dataset required to estimate all parameters must be of the same order of magnitude. In addition to the larger number of parameters, the space of possible non-linear functions that can fit a given dataset is also larger. Thus, linear models, with fewer possible functions and fewer parameters to specify those functions, are simpler and easier to fit. 1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104

Linear models are not flexible enough to estimate non-linear systems such as a building performance simulator. This means that the estimates of lin- 1105
1106
1107

ear models are precise but may be inaccurate. The flexibility of non-linear models, on the other hand, means that the data from a non-linear system can be estimated more accurately. However, this flexibility could lead to over-fitting the model to the dataset at hand. The practical consequence of this would be a failure to predict well on real-world data different from that included in the training dataset.

This problem of generalisability of models could arise from over-fitting to a small amount of data, or an unrepresentative dataset. In the context of BPS, such an unrepresentative dataset could consist, for example, of only one building type or weather (context/location). This would make the model inaccurate for other building types or locations. The use of a larger, varied dataset can reduce this problem to some extent (see Figures 10 and 11 for an example), since it would make it more likely that the model would see examples of more real-world situations.

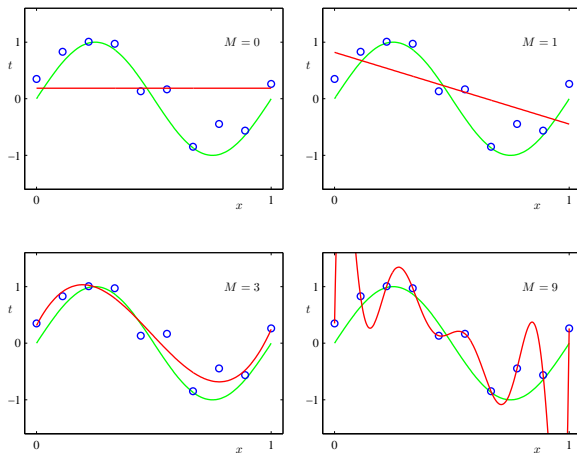


Figure 10: A generic representation of the tendency of non-linear models to over-fit data from Bishop [6]. In general, the more complex a model is, the more it will over-fit the data at hand. The green curve is used to generate the data and the red curves show a polynomial fit. Polynomials of progressively higher degrees ($M = 0, \dots, 9$) fit the data better but are probably over-fitting.

1128 A.2 Linear Models

The function $f_e(\cdot)$ from Equation (3) may take the form of a linear model:

$$\hat{y}_l = f_l(\mathbf{x}) := \boldsymbol{\beta}^T \mathbf{x}, \quad (12)$$

where \hat{y}_l is the prediction of the output at input \mathbf{x} obtained using a linear function $f_l(\mathbf{x})$ which is specified using $\boldsymbol{\beta}$, a real-valued parameter vector (of the same size as \mathbf{x}).

The parameter $\boldsymbol{\beta}$ is unknown but can be estimated using a *dataset* of the input-output pairs, e.g., obtained by running many BPS simulations on a plausible set of building designs $\boldsymbol{\theta}$ and its operating conditions \mathbf{x} . In a standard machine-learning framework, we first collect a large amount of such data: $\mathcal{D} := \{y_n, \mathbf{x}_n\}_{n=1}^N$ where n denotes the n 'th BPS simulation. Given a dataset \mathcal{D} , we may use the standard training-testing framework developed in statistics and machine learning [5] to estimate $\boldsymbol{\beta}$. In this framework, first, the N observations are split into two mutually-exclusive sets: training and testing. We denote the training set by \mathcal{D}_{train} which contains N_{train} number of observations. Similarly, we denote the test set by \mathcal{D}_{test} which may contain N_{test} . In this paper we use the term *validation* set to denote the dataset used for real-world testing of the model. By construction, $N = N_{train} + N_{test}$. The training set is used to *train* the linear model, i.e., to estimate $\boldsymbol{\beta}_*$ by minimizing a cost function, e.g., a mean-square error as shown below:

$$\boldsymbol{\beta}^* = \arg \min_{\boldsymbol{\beta}} \frac{1}{N_{train}} \sum_{n=1}^{N_{train}} (y_n - \boldsymbol{\beta}^T \mathbf{x}_n)^2, \quad (13)$$

This gives us a linear model $f_l^*(\mathbf{x}) := \boldsymbol{\beta}_*^T \mathbf{x}$ which can be used to predict the new inputs. The test set is then used to assess the *goodness-of-fit* of the estimator by computing the following cost,

$$\mathcal{L}(\hat{f}_l^*) = \frac{1}{N_{test}} \sum_{n=1}^{N_{test}} (y_n - \boldsymbol{\beta}_*^T \mathbf{x}_n)^2, \quad (14)$$

This is the *test error* which could be a faithful estimate of the real-world prediction error of the model when N_{test} is fairly large and representative of the real-world problem.

An advantage of a linear model is that training is easy. Equation (13) has a closed-form solution which can be obtained by using the ordinary least-squares method. This method scales well for medium-sized datasets and can also be extended to large datasets by using iterative methods such as stochastic gradient descent [46]. Another advantage of the linear model is that it is fairly straightforward to specify and interpret. An entry in the

1146 parameter β is a direct indicator of how important
 1147 the corresponding entry in input \mathbf{x} is for the linear
 1148 model to predict well. Unfortunately, linear models
 1149 are not good models of the simulator since the sim-
 1150 ulator is a non-linear model itself. As a result the
 1151 test error $\mathcal{L}(\hat{f}_i^*)$ is usually quite large, except for
 1152 the simplest problems as discussed in Section 2.3
 1153 above.

1154 A.3 Non-linear models

1155 Estimating the \hat{f}_e^* of Equation (4) is easier for linear
 1156 regression models since closed-form methods like
 1157 least-squares may be used [5, 6]. However, estimat-
 1158 ing the same quantity for non-linear models re-
 1159 quires the use of iterative methods, which is time-
 1160 consuming. The time to train parameters rises rap-
 1161 idly with the number of parameters to be estimat-
 1162 ed. As discussed in Section 2.3, several types
 1163 of non-linear model types have been proposed for
 1164 BPS. The recently-concluded ASHRAE Great En-
 1165 ergy Predictor III challenge [47] alone saw 415 solu-
 1166 tions submitted. Next, we discuss the structure of
 1167 a non-linear model and the process of fitting it to a
 1168 given dataset using a general-purpose model type:
 1169 Gaussian Process (GP) regression.

1170 A.4 Gaussian Process Regression

We use the framework of GP regression to estimate the non-linear function f_{nl} that minimizes Equation (14). GP regression uses Bayes’ rule to compute the posterior distribution over f_{nl} given sample outputs y_n [33, ch. 2]. This approach works directly in the space of f_{nl} and avoids both a direct estimation of β and also a direct specification of $\phi(\mathbf{x})$. Instead, we specify a ‘kernel’ function which defines the inner product of ϕ as

$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \Sigma \phi(\mathbf{x}_j), \quad (15)$$

where \mathbf{x}_i and \mathbf{x}_j are two inputs in our observation set. In practice, a kernel function is easier to specify than ϕ , even though it is sometimes unintuitive. For example, a linear model f_i can be specified by choosing the linear kernel

$$k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \Sigma \mathbf{x}_j. \quad (16)$$

The non-linear model used in this paper is a squared exponential function (SqE) kernel function

$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma_f^2 \exp \left[-\frac{1}{2}(\mathbf{x}_i - \mathbf{x}_j)^T \Sigma (\mathbf{x}_i - \mathbf{x}_j) \right], \quad (17)$$

where $\sigma_f^2 > 0$ is the signal variance. This kernel is also referred to as the radial basis function (RBF) kernel in the context of Artificial Neural Networks.

It is possible that emulator f_{nl} is not able to model the output y_n perfectly and, in that case, we can assume that there is noise in the estimation, i.e., $y_n = f_{nl}(\mathbf{x}_n) + \varepsilon_n$. Following the standard practice in GP regression, we assume that ε_n are independent Gaussian random variables with zero mean and noise variance σ_n^2 . Specifying this non-linear model requires estimation of the noise variance σ_n^2 , the signal variance σ_f^2 , and Σ . Collectively, these quantities are referred to as ‘hyper-parameters’ of the GP model, and we denote the set of hyper-parameters by ψ [33].

A.4.1 Fitting and Using a GP Model

Building a GP-based emulator requires two tasks. The first task is to estimate the hyper-parameters. This is called ‘learning’. The second task is to compute $\hat{f}_{nl}(\mathbf{x}_*)$ given a new input \mathbf{x}_* and the estimated hyper-parameters. This is called ‘prediction’.

We first give details of the prediction task. We wish to compute the *predictive* distribution of the output (here: annual energy use, denoted by y_*) at a new input (here: the set of features that define a building and weather conditions for that particular year, denoted by \mathbf{x}_*) present in the test data, i.e., the distribution $p(y_* | \mathbf{x}_*, \mathcal{D}, \psi)$ where $\mathcal{D} = \{y_1, \mathbf{x}_1, y_2, \mathbf{x}_2, \dots, y_{N_{\text{train}}}, \mathbf{x}_{N_{\text{train}}}\}$ and ψ is the set of hyper-parameters. For GP regression, this distribution is a Gaussian and has a closed form expression. This follows from the property that any finite number of samples drawn from a Gaussian Process are jointly Gaussian, giving the following expression for the distribution of $\mathbf{y} := [y_1, y_2, \dots, y_{N_{\text{train}}}]^T$ and the y_* corresponding to \mathbf{x}_* :

$$\begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma_n^2 \mathbf{I} & \mathbf{k}_* \\ \mathbf{k}_*^T & k_{**} + \sigma_n^2 \end{bmatrix} \right) \quad (18)$$

where \mathbf{K} is a matrix whose (i, j) ’th entry is $k(\mathbf{x}_i, \mathbf{x}_j)$, \mathbf{k}_* is a vector whose i ’th entry is

$k(\mathbf{x}_i, \mathbf{x}_*)$, $k_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$, and \mathbf{I} is an identity matrix of size $N_{\text{train}} \times N_{\text{train}}$. Using the above equation, we can write the expression for the distribution of y_* given \mathbf{y} by using the conditional distribution for a Gaussian distribution [33, pg. 16]:

$$p(y_* | \mathbf{x}_*, \mathcal{D}, \boldsymbol{\psi}) := \mathcal{N}(\mu_*, \sigma_*^2), \quad (19)$$

where $\mu_* := \mathbf{k}_*^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y}$
and $\sigma_*^2 := k_{**} - \mathbf{k}_*^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{k}_*$.

The computational complexity of these operations is $\mathcal{O}(N_{\text{train}}^3)$, due to the inversion of the matrix $\mathbf{K} + \sigma_n^2 \mathbf{I}$. That is, the number of operations required increases by the cube of the number of elements, so adding 2 data points, for example, would require 8 additional operations.

The distribution $p(y_* | \mathbf{x}_*, \mathcal{D}_t, \boldsymbol{\theta})$ depends on the specification of $\boldsymbol{\psi}$. We estimate $\boldsymbol{\psi}$ by maximizing the log-likelihood: $\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\psi})$. This is called the maximum likelihood estimation (MLE) method. The closed-form expression for the log-likelihood is

$$\log p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}) := -\frac{1}{2} \log |\mathbf{K} + \sigma_n^2 \mathbf{I}| - \frac{1}{2} \mathbf{y}^T (\mathbf{K} + \sigma_n^2 \mathbf{I})^{-1} \mathbf{y} - \frac{N_{\text{train}}}{2} \log(2\pi). \quad (20)$$

This can be optimised with a numerical optimization method [33]. However, every iteration requires a matrix inversion, which could be costly if the optimization takes too many iterations. We discuss a method to reduce the computation cost in Appendix A.4.2.

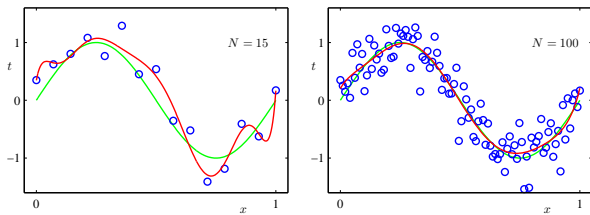
GP regression allows the specification of a ‘signal noise’, i.e., the variance of the uncertainty in the data itself (σ_n in Equation (18)). This noise variance may be fixed to some appropriate value or tuned along with the other hyper-parameters. Given that we are using simulated data, we expect the noise variance to be very low. However, the models are less stable when the signal noise is low because the covariance matrices are frequently non-invertible (Equation (20)). This is because the signal noise acts as a regulariser in this inversion of the covariance matrix \mathbf{K} when fitting a GP model. Therefore, a value of nearly zero for the signal noise foregoes the stability accorded by the regulariser. When we set a lower bound, $\sigma_n \geq 10^{-6}$, the ill-conditioning of the covariance matrices is reduced.

A.4.2 Using Big Datasets

Generally speaking, the more data a model sees to characterize a domain, the better it is able to predict on unseen data from that domain. Fitting a GP model involves the inversion of a matrix (the covariance matrix), whose size is $N \times N$, where N is the number of data points. This puts a limit on the size of dataset that can be considered for learning hyper-parameters or predicting. It is possible to work around this limitation by using the so-called ‘sparse’ methods, i.e., methods using sparse representations of the covariance matrix. However, these methods invariably reduce the predictive performance of the model. In this paper, we present a simple method to extend the amount of data considered, similar to a proposal in Chalupka et al. [42].

The training procedure consists of two steps: learning hyper-parameters through Maximum Likelihood Estimation (MLE) using some training data, and then predicting the output at test inputs using the same training data. These two steps correspond to Equations (19) and (20), respectively. The matrix \mathbf{K} from Equation (19) depends on training data, while the vector \mathbf{k}_* depends on the testing data. In our experiments, the estimates of hyper-parameters stabilize with about 1000-2000 training data points. If we continue with the policy of learning and predicting using the same dataset, we are restricted to models trained on about 5,000 points, because repeatedly inverting matrices (as part of the MLE step) of size 5,000 or more is impractical on the hardware available to us. A simple method to add more data was to increase the size of the dataset during prediction.

We modified the procedure to reduce run time by using sets of different sizes for learning and prediction (Figure 1). If learning is carried out on a smaller set of 2000 points ($N_{\text{train}} = 2000$), i.e., the matrix \mathbf{K} in Equation (20) is defined on a dataset of 2000 points, estimates of the hyper-parameters $\boldsymbol{\psi}$ are fixed relatively rapidly. Given this estimate of $\boldsymbol{\psi}$, we proceed to increase the data size to $N_{\text{pred}} > N_{\text{train}}$ for prediction using Equation (19). Since prediction involves only one matrix inversion, we were able to handle N_{pred} of size up to 12,000. This modification to the procedure gives a modest improvement in the validation error.



(a) Dataset of size $N = 15$ (b) Dataset of size $N = 100$

Figure 11: Using $N = 15$ [left] or $N = 100$ [right] data points to fit a polynomial of degree $M = 9$, shows that “...increasing the size of the data set reduces... overfitting” [6].

Acknowledgements

The authors would like to thank Dr Minu Agarwal, Dr Giorgia Chinazzo, Dr Georgios Mavromatidis, Margaux Peltier, Dr Masashi Sugiyama, Dr Norihiro Maeda, the RIKEN-AIP team, and the Energy Systems Research Unit (ESRU) team at University of Strathclyde. For part of this work, Parag was hosted at ESRU and RIKEN-AIP, and his work was funded by the Swiss National Science Foundation’s Postdoc.Mobility grant P2ELP2_168519.

References

[1] Joshua New, Jibonananda Sanyal, Bob Slatery, Anthony Gehl, William Miller, and Aaron Garrett. Big Data Mining for Assessing Calibration of Building Energy Models. *International Journal of Computer & Software Engineering*, 3(2), September 2018. ISSN 24564451. doi: 10.15344/2456-4451/2018/136. URL <https://www.graphyonline.com/archives/IJCSE/2018/IJCSE-136/>.

[2] J. A. Clarke. A vision for building performance simulation: a position paper prepared on behalf of the IBPSA Board. *Journal of Building Performance Simulation*, 8(2):39–43, March 2015. ISSN 1940-1493, 1940-1507. doi: 10.1080/19401493.2015.1007699. URL <http://www.tandfonline.com/doi/full/10.1080/19401493.2015.1007699>.

[3] Mahnameh Taheri, Parag Rastogi, Colin Parry, and Alan Wegienka. Benchmarking Building Energy Consumption Using Efficiency Factors. In *Proceedings of BS 2019*, pages 3863–3870, Rome, Italy, 2019. doi: 10.26868/25222708.

2019.210575. URL http://www.ibpsa.org/proceedings/BS2019/BS2019_210575.pdf.

[4] Parag Rastogi. *On the sensitivity of buildings to climate: the interaction of weather and building envelopes in determining future building energy consumption*. PhD, Ecole polytechnique fédérale de Lausanne, Lausanne, Switzerland, August 2016. URL <https://infoscience.epfl.ch/record/220971?ln=en>. doi:10.5075/epfl-thesis-6881.

[5] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition*. Springer, August 2009. ISBN 978-0-387-84858-7.

[6] Christopher M Bishop. *Pattern recognition and machine learning*. springer, 2006.

[7] Parag Rastogi, Mohammad Emtiyaz Khan, and Marilyne Andersen. Gaussian-Process-Based Emulators for Building Performance Simulation. In *Proceedings of BS 2017*, San Francisco, CA, USA, August 2017. IBPSA.

[8] Hai-xiang Zhao and Frédéric Magoulès. A review on the prediction of building energy consumption. *Renewable and Sustainable Energy Reviews*, 16(6):3586–3592, August 2012. ISSN 1364-0321. doi: 10.1016/j.rser.2012.02.049. URL <http://www.sciencedirect.com/science/article/pii/S1364032112001438>.

[9] Saleh Seyedzadeh, Farzad Pour Rahimian, Stephen Oliver, Sergio Rodriguez, and Ivan Glesk. Machine learning modelling for predicting non-domestic buildings energy performance: A model to support deep energy retrofit decision-making. *Applied Energy*, 279:115908, December 2020. ISSN 0306-2619. doi: 10.1016/j.apenergy.2020.115908. URL <http://www.sciencedirect.com/science/article/pii/S0306261920313702>.

[10] Aurélie Fouquier, Sylvain Robert, Frédéric Suard, Louis Stéphan, and Arnaud Jay. State of the art in building modelling and energy performances prediction: A review. *Renewable and Sustainable Energy Reviews*, 23:272–288, July 2013. doi: 10.1016/j.rser.2013.03.004. URL <http://linkinghub.elsevier.com/retrieve/pii/S1364032113001536>.

[11] Emilie Nault. *Solar Potential in Early Neighborhood Design. A Decision-Support Workflow Based*

- 1348 *on Predictive Models*. PhD Thesis, Ecole polytech- 1396
1349 nique fédérale de Lausanne, Lausanne, Switzer- 1397
1350 land, 2016.
- 1351 [12] Wei Tian. A review of sensitivity analysis 1398
1352 methods in building energy analysis. *Renew- 1399
1353 able and Sustainable Energy Reviews*, 20:411– 1400
1354 419, April 2013. doi: 10.1016/j.rser.2012.12. 1401
1355 014. URL [http://linkinghub.elsevier.com/](http://linkinghub.elsevier.com/retrieve/pii/S1364032112007101) 1402
1356 [retrieve/pii/S1364032112007101](http://linkinghub.elsevier.com/retrieve/pii/S1364032112007101). 1403
1404
- 1357 [13] Christina Johanna Hopfe. *Uncertainty and sensi- 1405
1358 tivity analysis in building performance simula- 1406
1359 tion for decision support and design optimization*. 1407
1360 PhD, Technische Universiteit Eindhoven, Eind- 1408
1361 hoven, The Netherlands, 2009. URL [http://www.](http://www.bwk.tue.nl/bps/hensen/team/past/Hopfe.pdf) 1409
1362 [bwk.tue.nl/bps/hensen/team/past/Hopfe.pdf](http://www.bwk.tue.nl/bps/hensen/team/past/Hopfe.pdf). 1410
- 1363 [14] Sten de Wit. *Uncertainty in predictions of 1411
1364 thermal comfort in buildings*. PhD, Delft Univer- 1412
1365 sity of Technology, Delft, The Netherlands, June 1413
1366 2001. URL [http://resolver.tudelft.nl/uuid:](http://resolver.tudelft.nl/uuid:a231bca8-ec81-4e22-8b34-4bafc062950e) 1414
1367 [a231bca8-ec81-4e22-8b34-4bafc062950e](http://resolver.tudelft.nl/uuid:a231bca8-ec81-4e22-8b34-4bafc062950e). 1415
- 1368 [15] Geoffrey K. F. Tso and Kelvin K. W. Yau. Pre- 1416
1369 dicting electricity energy consumption: A compar- 1417
1370 ison of regression analysis, decision tree and neural 1418
1371 networks. *Energy*, 32(9):1761–1768, September 1419
1372 2007. ISSN 0360-5442. doi: 10.1016/j.energy. 1420
1373 2006.11.010. URL [http://www.sciencedirect.](http://www.sciencedirect.com/science/article/pii/S0360544206003288) 1421
1374 [com/science/article/pii/S0360544206003288](http://www.sciencedirect.com/science/article/pii/S0360544206003288). 1422
- 1375 [16] Endong Wang, Zhigang Shen, and Kevin 1423
1376 Grosskopf. Benchmarking energy perform- 1424
1377 ance of building envelopes through a select- 1425
1378 ive residual-clustering approach using high di- 1426
1379 mensional dataset. *Energy and Buildings*, 1427
1380 75(Supplement C):10–22, June 2014. ISSN 1428
1381 0378-7788. doi: 10.1016/j.enbuild.2013.12.055. 1429
1382 URL [http://www.sciencedirect.com/science/](http://www.sciencedirect.com/science/article/pii/S0378778814000048) 1430
1383 [article/pii/S0378778814000048](http://www.sciencedirect.com/science/article/pii/S0378778814000048). 1431
- 1384 [17] Emilie Nault, Giuseppe Peronato, Emmanuel Rey, 1432
1385 and Marilyne Andersen. Review and critical ana- 1433
1386 lysis of early-design phase evaluation metrics for 1434
1387 the solar potential of neighborhood designs. *Build- 1435
1388 ing and Environment*, 92:679–691, October 2015. 1436
1389 ISSN 0360-1323. doi: 10.1016/j.buildenv.2015. 1437
1390 05.012. URL [http://www.sciencedirect.com/](http://www.sciencedirect.com/science/article/pii/S0360132315002243) 1438
1391 [science/article/pii/S0360132315002243](http://www.sciencedirect.com/science/article/pii/S0360132315002243). 1439
- 1392 [18] Emilie Nault, Parag Rastogi, Emmanuel Rey, and 1440
1393 Marilyne Andersen. The sensitivity of predicted 1441
1394 energy use to urban geometrical factors in vari- 1442
1395 ous climates. In *Proceedings of PLEA 2015*, 1443
1444 Bologna, Italy, September 2015. URL [http://](http://infoscience.epfl.ch/record/211101?ln=en) 1444
1445 infoscience.epfl.ch/record/211101?ln=en. 1446
- [19] Janelle S Hygh, Joseph F. DeCarolis, David B 1398
Hill, and S Ranji Ranjithan. Multivariate regres- 1399
sion as an energy assessment tool in early build- 1400
ing design. *Building and Environment*, 57:165– 1401
175, November 2012. doi: 10.1016/j.buildenv. 1402
2012.04.021. URL [http://dx.doi.org/10.1016/](http://dx.doi.org/10.1016/j.buildenv.2012.04.021) 1403
[j.buildenv.2012.04.021](http://dx.doi.org/10.1016/j.buildenv.2012.04.021). 1404
- [20] Somayeh Asadi, Shideh Shams Amiri, and Mo- 1405
hammad Mottahedi. On the development of multi- 1406
linear regression analysis to assess energy con- 1407
sumption in the early stages of building design. 1408
Energy and Buildings, 85:246–255, December 1409
2014. ISSN 0378-7788. doi: 10.1016/j.enbuild. 1410
2014.07.096. URL [http://www.sciencedirect.](http://www.sciencedirect.com/science/article/pii/S0378778814007154) 1411
[com/science/article/pii/S0378778814007154](http://www.sciencedirect.com/science/article/pii/S0378778814007154). 1412
- [21] Shideh Shams Amiri, Mohammad Mottahedi, and 1413
Somayeh Asadi. Using multiple regression ana- 1414
lysis to develop energy consumption indicators for 1415
commercial buildings in the U.S. *Energy and* 1416
Buildings, 109:209–216, December 2015. ISSN 1417
0378-7788. doi: 10.1016/j.enbuild.2015.09.073. 1418
URL [http://www.sciencedirect.com/science/](http://www.sciencedirect.com/science/article/pii/S0378778815303133) 1419
[article/pii/S0378778815303133](http://www.sciencedirect.com/science/article/pii/S0378778815303133). 1420
- [22] Torben Østergård, Rasmus L. Jensen, and Stef- 1421
fen E. Maagaard. Early Building Design: In- 1422
formed decision-making by exploring multidimen- 1423
sional design space using sensitivity analysis. *En- 1424
ergy and Buildings*, 142(Supplement C):8–22, May 1425
2017. ISSN 0378-7788. doi: 10.1016/j.enbuild. 1426
2017.02.059. URL [http://www.sciencedirect.](http://www.sciencedirect.com/science/article/pii/S0378778817306916) 1427
[com/science/article/pii/S0378778817306916](http://www.sciencedirect.com/science/article/pii/S0378778817306916). 1428
- [23] Xi Chen, Hongxing Yang, and Ke Sun. Devel- 1429
oping a meta-model for sensitivity analyses and 1430
prediction of building performance for passively 1431
designed high-rise residential buildings. *Applied* 1432
Energy, 194(Supplement C):422–439, May 2017. 1433
ISSN 0306-2619. doi: 10.1016/j.apenergy.2016. 1434
08.180. URL [http://www.sciencedirect.com/](http://www.sciencedirect.com/science/article/pii/S0306261916312892) 1435
[science/article/pii/S0306261916312892](http://www.sciencedirect.com/science/article/pii/S0306261916312892). 1436
- [24] Joshua Hester, Jeremy Gregory, and Randolph 1437
Kirchain. Sequential early-design guidance for res- 1438
idential single-family buildings using a probabil- 1439
istic metamodel of energy consumption. *Energy* 1440
and Buildings, 134(Supplement C):202–211, Janu- 1441
ary 2017. ISSN 0378-7788. doi: 10.1016/j.enbuild. 1442
2016.10.047. URL [http://www.sciencedirect.](http://www.sciencedirect.com/science/article/pii/S037877881631369X) 1443
[com/science/article/pii/S037877881631369X](http://www.sciencedirect.com/science/article/pii/S037877881631369X). 1444

- 1445 [25] Siyu Wu and Jian-Qiao Sun. Multi-stage regres- 1494
1446 sion linear parametric models of room temperat- 1495
1447 ure in office buildings. *Building and Environment*, 1496
1448 56:69–77, October 2012. doi: 10.1016/j.buildenv. 1497
1449 2012.02.026. URL [http://linkinghub.elsevier. 1498](http://linkinghub.elsevier.com/retrieve/pii/S0360132312000716)
1450 [com/retrieve/pii/S0360132312000716](http://linkinghub.elsevier.com/retrieve/pii/S0360132312000716). 1499
- 1451 [26] Sandhya Patidar, David P Jenkins, Gavin J Gib- 1500
1452 son, and P F G Banfill. Statistical techniques 1501
1453 to emulate dynamic building simulations for over- 1502
1454 heating analyses in future probabilistic climates. 1503
1455 *Journal of Building Performance Simulation*, 4 1504
1456 (3):271–284, 2011. doi: 10.1080/19401493.2010. 1505
1457 531144. 1506
- 1458 [27] D.P. Jenkins, M. Gul, and S. Patidar. Probabil- 1507
1459 istic future cooling loads for mechanically cooled 1508
1460 offices. *Energy and Buildings*, 66:57–65, Novem- 1509
1461 ber 2013. ISSN 0378-7788. doi: 10.1016/j.enbuild. 1510
1462 2013.07.040. URL [http://www.sciencedirect. 1511](http://www.sciencedirect.com/science/article/pii/S0378778813004313)
1463 [com/science/article/pii/S0378778813004313](http://www.sciencedirect.com/science/article/pii/S0378778813004313). 1512
- 1464 [28] David P Jenkins, Sandhya Patidar, P F G 1513
1465 Banfill, and Gavin J Gibson. Probabil- 1514
1466 istic climate projections with dynamic build- 1515
1467 ing simulation: Predicting overheating in dwell- 1516
1468 ings. *Energy and Buildings*, 43(7):1723–1731, 1517
1469 July 2011. doi: 10.1016/j.enbuild.2011.03. 1518
1470 016. URL [http://linkinghub.elsevier.com/ 1519](http://linkinghub.elsevier.com/retrieve/pii/S0378778811000946)
1471 [retrieve/pii/S0378778811000946](http://linkinghub.elsevier.com/retrieve/pii/S0378778811000946). 1520
- 1472 [29] David Coley and Tristan Kershaw. Changes in 1521
1473 internal temperatures within the built environ- 1522
1474 ment as a response to a changing climate. *Build- 1523*
1475 *ing and Environment*, 45(1):89–93, January 2010. 1524
1476 ISSN 0360-1323. doi: 10.1016/j.buildenv.2009. 1525
1477 05.009. URL [http://www.sciencedirect.com/ 1526](http://www.sciencedirect.com/science/article/pii/S0360132309001280)
1478 [science/article/pii/S0360132309001280](http://www.sciencedirect.com/science/article/pii/S0360132309001280). 1527
- 1479 [30] Yeonsook Heo and Victor M. Zavala. Gaus- 1528
1480 sian process modeling for measurement and ver- 1529
1481 ification of building energy savings. *Energy 1530*
1482 *and Buildings*, 53:7–18, October 2012. ISSN 1531
1483 0378-7788. doi: 10.1016/j.enbuild.2012.06.024. 1532
1484 URL [http://www.sciencedirect.com/science/ 1533](http://www.sciencedirect.com/science/article/pii/S037877881200312X)
1485 [article/pii/S037877881200312X](http://www.sciencedirect.com/science/article/pii/S037877881200312X). 1534
- 1486 [31] Michael C. Burkhart, Yeonsook Heo, and Vic- 1535
1487 tor M. Zavala. Measurement and verification 1536
1488 of building systems under uncertain data: A 1537
1489 Gaussian process modeling approach. *Energy 1538*
1490 *and Buildings*, 75:189–198, June 2014. ISSN 1539
1491 0378-7788. doi: 10.1016/j.enbuild.2014.01.048. 1540
1492 URL [http://www.sciencedirect.com/science/ 1541](http://www.sciencedirect.com/science/article/pii/S0378778814001091)
1493 [article/pii/S0378778814001091](http://www.sciencedirect.com/science/article/pii/S0378778814001091). 1542
- [32] Filippo Monari. *Sensitivity Analysis and Bayesian 1494*
Calibration of Building Energy Models. PhD, Uni- 1495
versity of Strathclyde, Glasgow, UK, 2016. 1496
- [33] Carl Edward Rasmussen and Christopher K. I. 1497
Williams. *Gaussian processes for machine learn- 1498*
ing. Adaptive computation and machine learning. 1499
MIT Press, Cambridge, Mass, 2006. ISBN 978-0- 1500
262-18253-9. 1501
- [34] Michael Deru, Kristin Field, Daniel Studer, Kyle 1502
Benne, Brent Griffith, Paul Torcellini, Bing Liu, 1503
Mark Halverson, Dave Winiarski, Michael Rosen- 1504
berg, Mehry Yazdani, Y. Joe Huang, and 1505
Drury B. Crawley. U.S. Department of Energy 1506
commercial reference building models of the na- 1507
tional building stock. Technical report, National 1508
Renewable Energy Laboratory (NREL), February 1509
2011. URL [http://digitalscholarship.unlv. 1510](http://digitalscholarship.unlv.edu/renew_pubs/44)
[edu/renew_pubs/44](http://digitalscholarship.unlv.edu/renew_pubs/44). 1511
- [35] Minu Agarwal, Parag Rastogi, Margaux Peltier, 1512
Luisa Pastore, and Marilyne Andersen. Exam- 1513
ining Building Design Decisions Under Long Term 1514
Weather Variability and Microclimate Effects: A 1515
Case Based Exploratory Study. In *Proceedings of 1516*
PLEA 2016, Los Angeles, USA, July 2016. 1517
- [36] Giorgia Chinazzo. Refurbishment of Existing En- 1518
velopes in Residential Buildings: assessing ro- 1519
bust solutions for future climate change. Master’s 1520
thesis, EPFL, Lausanne, Switzerland, 2014. URL 1521
<http://infoscience.epfl.ch/record/203438>. 1522
- [37] Parag Rastogi, Sönke Frederik Horn, and Mar- 1523
ilyne Andersen. Toward Assessing the Sensitivity 1524
of Buildings to Changes in Climate. In *Proceed- 1525*
ings of PLEA 2013, Munich, Germany, Septem- 1526
ber 2013. URL [http://infoscience.epfl.ch/ 1527](http://infoscience.epfl.ch/record/187507?ln=en)
[record/187507?ln=en](http://infoscience.epfl.ch/record/187507?ln=en). 1528
- [38] Jan Remund, Stefan Mueller, Stefan Kunz, and 1529
Christoph Schilter. METEONORM Handbook 1530
Part II : Theory. Technical report, Meteotest, May 1531
2012. URL [http://meteonorm.com/download/ 1532](http://meteonorm.com/download/software/mn70/)
[software/mn70/](http://meteonorm.com/download/software/mn70/). 1533
- [39] S Wilcox and W Marion. Users’ Manual for TMY3 1534
Data Sets. Technical report, National Renewable 1535
Energy Laboratory, May 2008. URL [http://www. 1536](http://www.nrel.gov/docs/fy08osti/43156.pdf)
[nrel.gov/docs/fy08osti/43156.pdf](http://www.nrel.gov/docs/fy08osti/43156.pdf). 1537
- [40] Parag Rastogi and Marilyne Andersen. Embed- 1538
ding Stochasticity in Building Simulation Through 1539
Synthetic Weather Files. In *Proceedings of BS 1540*

- 1541 2015, Hyderabad, India, December 2015. URL
1542 <http://infoscience.epfl.ch/record/208743>.
- 1543 [41] Parag Rastogi and Marilyne Andersen. Incorporating
1544 Climate Change Predictions in the Analysis
1545 of Weather-Based Uncertainty. In *Proceedings of*
1546 *SimBuild 2016*, Salt Lake City, UT, USA, Au-
1547 gust 2016. URL [http://infoscience.epfl.ch/](http://infoscience.epfl.ch/record/208743)
1548 [record/208743](http://infoscience.epfl.ch/record/208743).
- 1549 [42] Krzysztof Chalupka, Christopher KI Williams,
1550 and Iain Murray. A framework for evaluating ap-
1551 proximation methods for Gaussian process regres-
1552 sion. *Journal of Machine Learning Research*, 14
1553 (Feb):333–350, 2013. URL [http://www.jmlr.org/](http://www.jmlr.org/papers/v14/chalupka13a.html)
1554 [papers/v14/chalupka13a.html](http://www.jmlr.org/papers/v14/chalupka13a.html). 00071.
- 1555 [43] CIBSE. *Sustainability - CIBSE Guide L: 2020*.
1556 CIBSE guides. The Chartered Institution of Build-
1557 ing Service Engineers, London, UK, 2 edition, July
1558 2020. ISBN 978-1-912034-49-9.
- 1559 [44] J. A. Clarke. The role of building opera-
1560 tional emulation in realizing a resilient built en-
1561 vironment. *Architectural Science Review*, 61
1562 (5):358–361, September 2018. ISSN 0003-
1563 8628. doi: 10.1080/00038628.2018.1502157.
1564 URL [https://doi.org/10.1080/00038628.2018.](https://doi.org/10.1080/00038628.2018.1502157)
1565 1502157. Publisher: Taylor & Francis .eprint: ht-
1566 tps://doi.org/10.1080/00038628.2018.1502157.
- 1567 [45] J. A. Clarke. A simulation-based procedure
1568 for the holistic resilience testing of building
1569 performance. *IOP Conference Series: Earth*
1570 *and Environmental Science*, 329:012027, October
1571 2019. ISSN 1755-1315. doi: 10.1088/1755-1315/
1572 329/1/012027. URL [https://doi.org/10.1088/](https://doi.org/10.1088/1755-1315/329/1/012027)
1573 [1755-1315/329/1/012027](https://doi.org/10.1088/1755-1315/329/1/012027). Publisher:
1574 IOP Publishing.
- 1575 [46] Léon Bottou. Large-Scale Machine Learning with
1576 Stochastic Gradient Descent. In Yves Lecheval-
1577 lier and Gilbert Saporta, editors, *Proceedings of*
1578 *the 19th International Conference on Computa-*
1579 *tional Statistics (COMPSTAT'2010)*, pages 177–
1580 187, Paris, France, August 2010. Springer. URL
1581 <http://leon.bottou.org/papers/bottou-2010>.
- 1582 [47] Clayton Miller, Pandarasamy Arjunan, Anjukan
1583 Kathirgamanathan, Chun Fu, Jonathan Roth,
1584 June Young Park, Chris Balbach, Krishnan Gowri,
1585 Zoltan Nagy, Anthony Fontanini, and Jeff Haberl.
1586 The ASHRAE Great Energy Predictor III com-
1587 petition: Overview and results. *arXiv:2007.06933*
1588 *[cs]*, July 2020. URL [http://arxiv.org/abs/](http://arxiv.org/abs/2007.06933)
1589 [2007.06933](http://arxiv.org/abs/2007.06933). arXiv: 2007.06933.