

# Friends don't let friends over-fit: Benchmarking open-source prediction methods on building electrical meter data

Clayton Miller

*Building and Urban Data Science (BUDS) Lab, Dept. of Building, School of Design and Environment (SDE), National University of Singapore (NUS)*

---

## Abstract

Prediction is a common machine learning (ML) models used on sub-hourly building energy consumption data. This process is valuable for anomaly detection, load profile-based control, energy plant systems control, and measurement and verification procedures. Literally hundreds of building energy prediction techniques have been developed over the last three decades, yet there is still no consensus on which techniques are the most effective for various building types. In addition, many of the techniques developed are proprietary and unavailable to the general research community. This paper outlines a library of open source regression techniques from the *Scikit-Learn* Python library and describes the process of applying them to open hourly electrical meter data from 482 non-residential buildings from data from the *Building Data Genome Project*. The results illustrate that there is no *one size-fits-all* modeling solution and that various types of temporal behavior are difficult to capture using machine learning. This framework and methodology is designed to be a *baseline* implementation for other building energy data prediction methods developed by commercial providers or the wider research community. The benchmark data set can also be expanded with numerous other building performance data from a wider representation of buildings from around the world. The use of a baseline data set in future prediction research results in comparability and reproducibility of techniques in the built environment domain.

---

\*Corresponding author email: clayton@nus.edu.sg, Phone: +65 81602452

*Keywords:* Building energy prediction, Building performance prediction, Performance prediction, Machine learning, Smart meters, Artificial neural networks, Support vector machines, Deep Learning

---

## 1. Introduction

Machine learning prediction models are highly impacting all facets of industry and science. They are being developed to diagnose illnesses, drive cars, suggest purchases to potential customers and mine the human genome. The built environment has the opportunity to leverage the same algorithms and techniques to improve efficiency and to create new business models [1]. Building performance analysis has dozens of uses for temporal prediction of electricity, heating and cooling energy. Prediction is often made both on the short-term (hours or days ahead) or long-term (weeks, months or years ahead). Short-term prediction are generally used for real-time HVAC control and efficiency of upcoming hours [2], scheduling and management of power stations and demand response schemes [3], and the analysis of residential metering and sub-metering [4], in addition to many other applications. Long-term prediction is used for the evaluation of energy conservation measures through a baseline model generation [5] and capacity expansion and planning [cite here]. Figure 1 illustrates the measurement and verification procedure using long-term energy prediction models. A period of baseline energy consumption is used to create a machine learning prediction model to evaluate how much energy a building would use in a *status quo* baseline mode. A energy conservation measure (ECM) is installed and the difference between the baseline is the avoided energy consumption or demand. This process is crucial for the implementation of energy savings implementations as it gives building owners, designers, and contractors a means of evaluating success of such measures.

### 1.1. Contemporary building energy prediction

An updated comprehensive review of building performance prediction studies is available that describes the models, techniques, input features, and uses for energy prediction in buildings [7]. This study reviewed research that implemented the most common prediction modelling techniques as applied to building energy prediction. These models include Support Vector Machines (SVM), Artificial Neural Networks (ANN), Ensemble Methods, and various

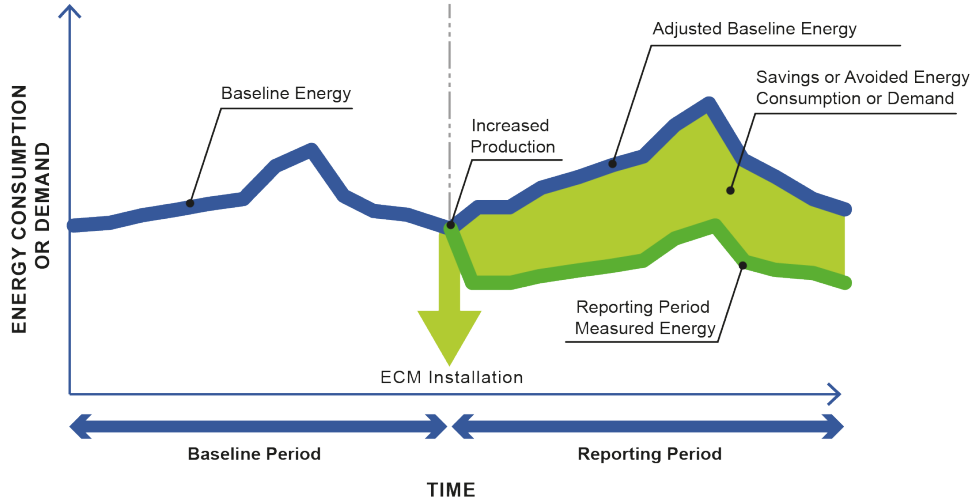


Figure 1: Prediction models as a comparison to energy savings interventions (Used with permission from the EVO - IPMVP [6])

other methods. A majority of the literature focuses on single building or a small set of buildings case studies.

The results of the studies with a low number of samples are problematic from a machine learning standpoint as the techniques are mostly applied to a single building or a small set of buildings. Training a machine learning model on such a small set of data results in a solution that is grossly over-fitted for a specific scenario. Such over-fitting results in models that are inappropriate to apply to data from the wider building stock. Over-fitting is a well-understood challenge in the machine learning community. This concept is related the bias-variance trade off that is illustrated in Figure 3. A model with high bias suffers from a lack of complexity to capture the behavior that is occurring in reality in a meaningful way. This model would be considered *underfitting* the data. On the other hand, a model that is overly complex suffers from high variance; this model is able to capture all the detailed behavior in the data that is used to train the model, but the model does not perform as well with new data sets or future data.

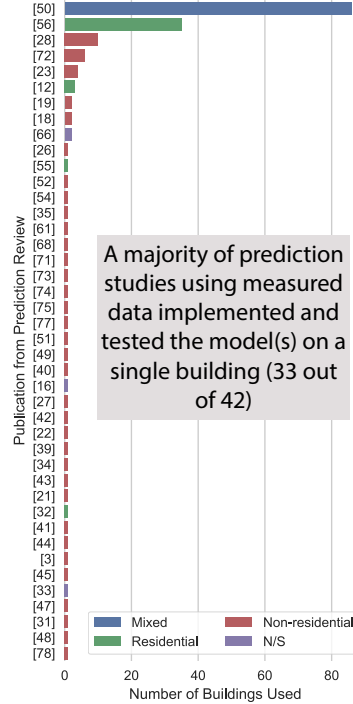


Figure 2: In the most recent prediction review study [7], a majority of publications developed and tested a machine learning modeling framework on data from a single building. These studies have the tendency to create very complex models for those particular buildings

Three key studies were completed in the last four years which finally took the process of machine learning for prediction to a diverse set of buildings to attempt to tackle the issue of overfitting. The first used 400 randomly selected buildings and applied six common, open prediction models for the purpose of creating a measurement and verification baseline [5]. The next study used 537 buildings and applied ten prediction models to further understand which model performed best among a large set of building [9]. The final study focused on the evaluation of what percentage of the building stock are appropriate to be used with *automated methods* as developed in the previous two studies [10].

Overall, these studies tested a large set energy prediction methods on a large set of buildings and this method was a step in the direction of general-

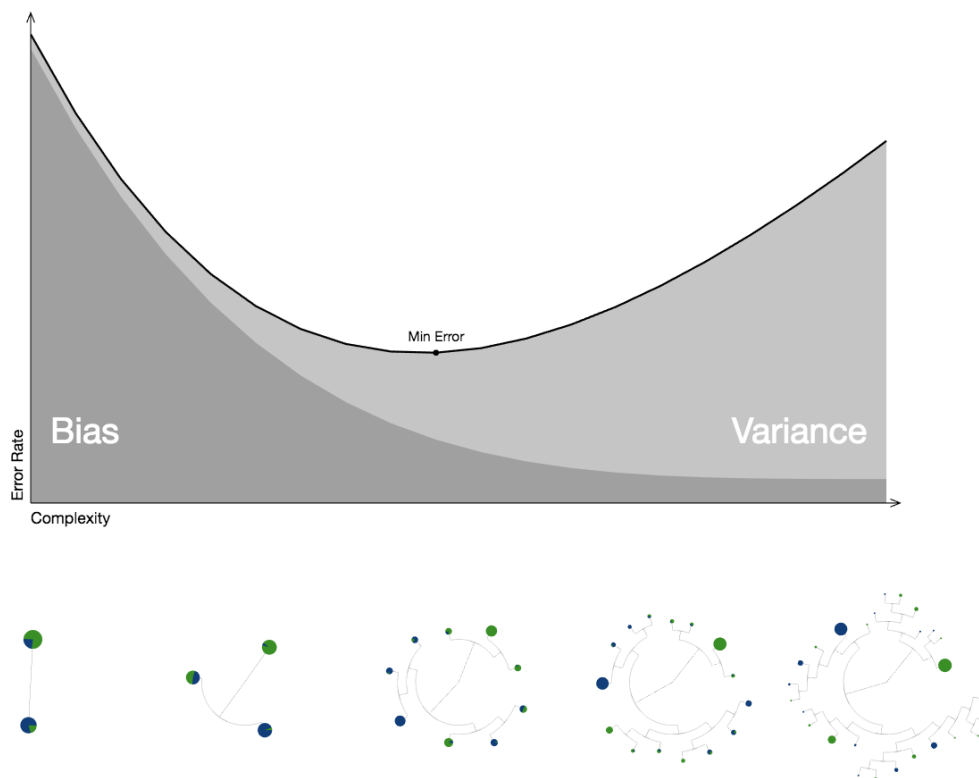


Figure 3: Bias-variance trade-off - The graph shows how as complexity increases the *bias* of the model drops as does the error rate. At a certain point the model becomes so complex that variance, or over-fitting becomes the key issue and the error rates increase again. The tree-based models below the chart show the increasing complexity of models of that type. Graphic adapted with permission from [8]

60 izability of energy prediction models. However, these studies are problematic  
 61 due to the lack of access to the exact models or data the researchers used.  
 62 Access to these aspects of their studies would allow future techniques to be  
 63 compared directly through application of the old machine learning methods  
 64 on new data or application of new machine learning methods on their old  
 65 data. These studies give the community an understanding of what models  
 66 work better than others in the context of the bias-variance trade-off, but they  
 67 do not provide the ability to test new models and techniques.

68 1.2. *The importance of benchmarking - A remedy for a crisis of over-fitting*

69 Despite the advancement of machine learning and prediction for perfor-  
70 mance data for buildings, a major barrier to wide-spread dissemination is  
71 that the techniques are not easily reproducible. Engineers, data scientists  
72 and researchers should be really asking themselves *does my machine learn-*  
73 *ing technique actually scale across hundreds of buildings? And is it actually*  
74 *faster or more accurate? How do we actually compare, each individual tech-*  
75 *nique against previously created methods?*

76 The time-series data mining community identified this problem as early  
77 as 2003: “Much of this work has very little utility because the contribution  
78 made”...“offer an amount of improvement that would have been completely  
79 dwarfed by the variance that would have been observed by testing on many  
80 real world data sets, or the variance that would have been observed by chang-  
81 ing minor (unstated) implementation details.” [cite keogh]

82 Figure 4 illustrates the conventional way that machine learning is gen-  
83 erally implemented in the built environment domain. Most of the existing  
84 building performance data science studies rely on each individual researcher  
85 creating their own methods, finding a case study data set and determining  
86 efficacy on their own. Not surprisingly, most of those researcher find positive  
87 results. However, the ability to compare those results to other publications  
88 and techniques is limited.

89 [here go into the efforts in thi benchmarking in computer science and  
90 machine learning]

91 Using a large, consistent benchmark data set from hundreds (or thou-  
92 sands) of buildings, a researcher can determine how well their methods ac-  
93 tually perform across a heterogeneous data set. If multiple researcher use  
94 the same data set, then there can be meaningful comparisons of accuracy,  
95 speed and ease-of-use. The purpose of this paper is to establish an example  
96 of machine learning benchmarking for the purposes of energy forecasting and  
97 prediction.

98 The creation of benchmarking data sets and methods sets the research  
99 community towards an environment in which there is a measure of account-  
100 ability in the claims made by new algorithms. These new techniques will  
101 likely be implemented on the test case(s) developed in the course of a re-  
102 search study, but will need to also be applied as much as possible to bench-  
103 marking methods in order to show which type of improvement is occurring  
104 in the literature and the quantification of improvement.

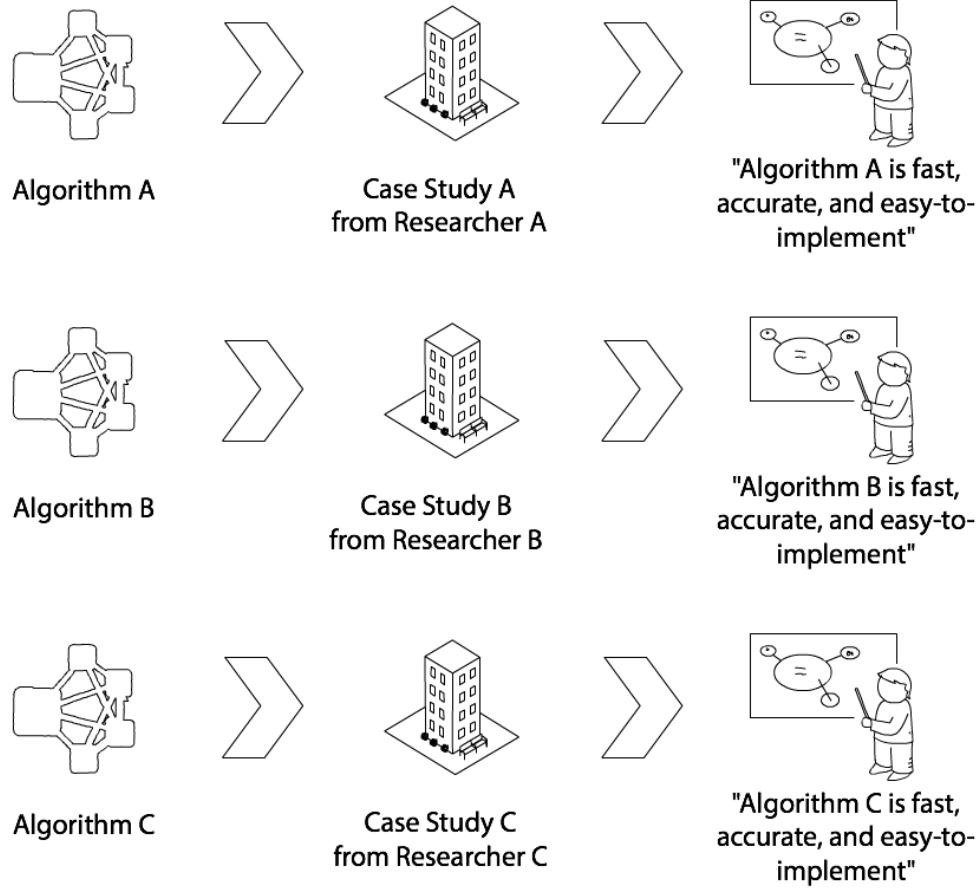


Figure 4: Old anecdotal machine learning testing and implementation methodology that results in a lack of comparative opportunities for new techniques

105 1.2.1. *The original benchmark - Great Energy Predictor Shootout I and II*

106 In the non-residential building research domain, there is a single set of

107 examples of a benchmarking data set that was utilized for several machine

108 learning studies. This set is the *Great Energy Predictor Shootout* compe-

109 titions held in the early 1990's. The first competition included the use of

110 a single building's electrical, cooling, and heating meters in a competition

111 where the participants were asked predict a single month of data using three

112 months of training data [11]. This competition resulted in several publication

113 based on the top set of performers in the competition. A second competition

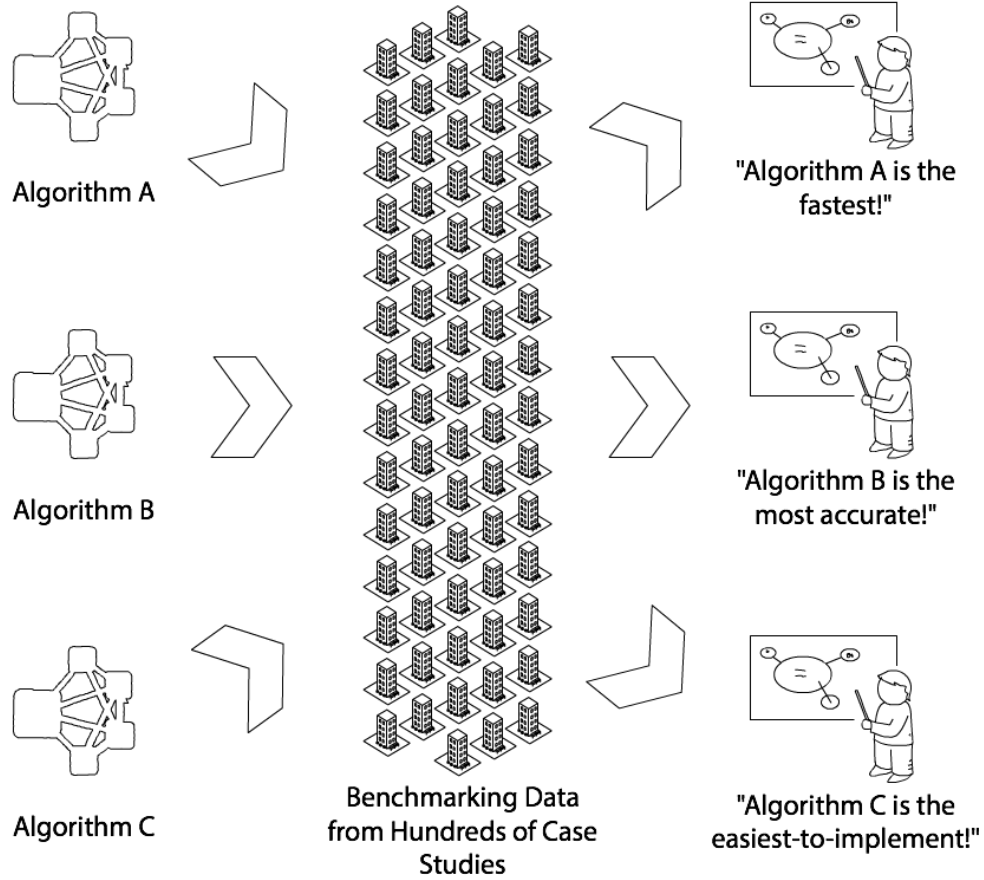


Figure 5: Using a common bench-marking data set enables the comparison of algorithms developed in the research community

114 was held, the aptly-named *Great Energy Predictor Shootout II* that asked  
 115 participants to predict the energy savings for a energy savings project [12].  
 116 These data sets have been since been used as a benchmarking comparison  
 117 data set on two studies; one focused on residential modelling using SVM [ref]  
 118 and the other to predict consumption of office buildings in Greece [ref].



### 119 1.3. Tackling the over-fitting problem - open benchmarking data sets and open 120 models

121 This paper outlines a demonstration of how benchmarking of performance  
122 prediction models can be accomplished on an open data set using open-  
123 source prediction techniques. Initially, a review of prediction considerations  
124 for building energy performance is reviewed. Next, a framework for imple-  
125 mentation is outlined using the Building Data Genome Project data set. The  
126 results are showcased and interpreted in the context of the different primary  
127 use types and operations considerations. Finally, a discussion of the limita-  
128 tions and guidelines for implementation of future techniques is presented.

## 129 2. Review of building energy prediction input considerations

130 To initiate the prediction model discussion, an overview of the conven-  
131 tional energy performance prediction considerations is covered in this section.  
132 These aspects of machine learning-based modelling for buildings are most  
133 prominent in non-residential buildings such as offices, educational facilities,  
134 laboratories, and health-care. These categories are the key considerations  
135 when developing prediction models for buildings.

### 136 2.1. Daily, weekly and seasonal schedules

137 Buildings operate on several types of schedules. A majority of non-  
138 residential building have *occupied* and *unoccupied* periods that usually co-  
139 incide with daytime and nighttime. These diurnal cycles are generally one  
140 of the best indicators of the building use type - i.e.: offices are open from  
141 9am to 6pm and hotels are most active from 6pm to 10am. Most buildings  
142 also have a weekly schedule; the default being certain behavior on weekdays  
143 and a different behavior on weekends. Finally, many buildings have seasonal  
144 changes such as when educational buildings have certain behavior during  
145 a regular session versus during breaks or holiday seasons. The concept of  
146 scheduling in non-residential building is most often related to the predefined  
147 schedules programmed into the automation systems in the building, unlike  
148 the human behavior that is discussed later. These schedules are determined  
149 usually by the operations and maintenance policy or by the energy manage-  
150 ment group within an organization. Many buildings have very predictable  
151 schedules, while others are much more volatile. Modelling such behavior can  
152 often be done using time-series methods that find auto-correlation behavior  
153 or by using date/time features as inputs to prediction models.

154 [insert the references from the temporal mining - uses dayfilter and other  
155 prominent daily and weekly clustering]

## 156 2.2. *Human behavior*

157 The concept of *human behavior* as an influence is similar to the previously-  
158 discussed schedules, however there is a more stochastic element to these be-  
159 havior. Buildings that are more influenced by occupant behavior generally  
160 have demand response-based control systems that use sensors, cameras or  
161 other detection methods to modulate systems only when humans are present  
162 or using the space for a specific purpose. Sometimes humans even have the  
163 ability to control spaces using various types of interfaces with the building,  
164 although this is less common in non-residential buildings. Modelling occu-  
165 pant behavior is considered more complex than schedule as human behavior  
166 can be less systematic, thus auto correlation-based time-series methods are  
167 less effective.

## 168 2.3. *Weather*

169 A major energy consuming component for most buildings are heating,  
170 ventilation, and air-conditioning (HVAC) systems. Intuitively, these systems  
171 tend to use more energy as the outdoor conditions get hotter (in the case of  
172 cooling) or colder (in the case of heating). Often this relationship is linear,  
173 but it can also be non-linear based on the HVAC system type and operation  
174 policy. The degree of influence of weather varies greatly among the building  
175 stock. It is influenced by the percentage of internal load vs. envelope-based  
176 load, HVAC system type, climate, and other factors.

## 177 2.4. *Non-routine events*

178 Non-routine events are disruptions in the systematic operation of a build-  
179 ing due to events such as a change in the operational schedule of the building,  
180 equipment breakdowns, and events or human behavior that is highly irregular  
181 as compared to past behavior. These events could be planned or unplanned  
182 by the operations and facilities management staff. Non-routine events are  
183 hard to predict with conventional input variables, but models can be tuned  
184 to quickly adapt to the change using approaches related to change-point  
185 prediction or other types of behavior change detection techniques.

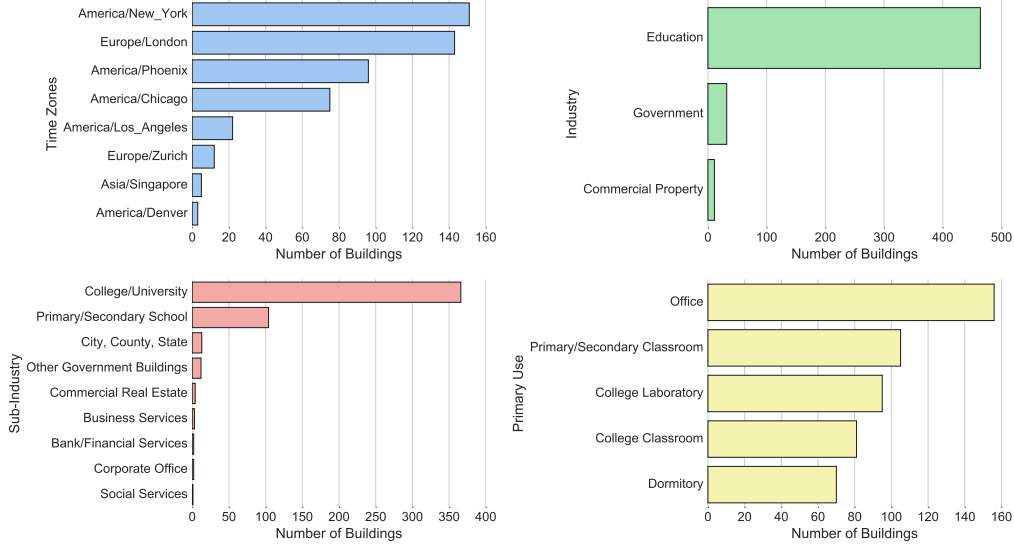


Figure 6: Meta-data breakdown of the *Building Data Genome Project* benchmarking data set used in this publication. Used with permission from (cite)

### 186 3. Benchmarking open source prediction methods on open building 187 energy data

188 The key to this benchmarking analysis is to use a diverse enough data  
189 set to illustrate the benefits of scalability and generalizability amongst the  
190 building stock. The regression testing framework for this paper is outlined in  
191 this section. The open hourly data from the Building Data Genome Project  
192 is used in this paper as a starting point. This data set includes data from  
193 over 500 buildings, mostly from educational institutions. Figure 6 illustrates  
194 the breakdown of four meta-data points from this data-set.

#### 195 3.1. Machine learning input variables

196 The input variables available as independent predictors of energy con-  
197 sumption are outlined in this section. Table 1 outlines the basic set of ma-  
198 chine learning input variables used in this benchmarking approach.

#### 199 3.2. Training and test data set scenario

200 For the purposes of the comparison of various open source techniques, a  
201 simplified training and testing scenario is utilized for the comparison. One

| Category | Variable                | Behavior Targeted                 |
|----------|-------------------------|-----------------------------------|
| Temporal | Time of Day             | Daily schedules                   |
|          | Day of Week             | Weekly schedules                  |
|          | Public Holiday Schedule | Holidays                          |
| Weather  | Schedule Type           | Seasonal schedules (summer, etc.) |
|          | Outdoor Air Temp.       | Sensible heating and cooling      |
|          | Outdoor Air Hum.        | Latent heating and cooling        |
| Meta     | Industry Sector         | General category of use           |
|          | Primary Use Type        | Specific category of use          |
|          | Floor area              | Size of building                  |
|          | Number of floors        | Height of building                |
|          | Climate zone            | Type of climate                   |
|          | Maximum occupancy       | Total number of people            |
|          | Cooling system type     | Typical cooling efficiency        |
|          | Heating system type     | Typical heating efficiency        |
|          | Performance rating      | Comparison to its peers           |

Table 1: Independent input features used in this benchmarking process

year of whole building hourly electrical meter data is available for all 482 buildings. Four different training and testing data scenarios are tested as seen in Figure 3.2. Scenarios 1-3 are made up of 3, 6, and 9 month training windows and a 3 month continuous testing window. Scenario 4 is made up of 3 month training windows followed by a 1 month test window that repeats itself three times. These scenarios provide a certain level of cross-validation that is realistic in the context of medium to long-term energy prediction in the built environment.

### 3.3. Open source regression models from Sci-kit Learn Python Library

The regression model catalogue from the Sci-kit Learn Library is used in this benchmarking process. We use these models as a starting point for the benchmarking process as many of these models have been developed for decades and are some of the most often used prediction models in the machine learning community. The larger array of more advanced models (e.g.: deep learning, specific energy-focused prediction models) is left outside the scope of this paper as these will be the *improved techniques* that the wider research community would likely be testing.

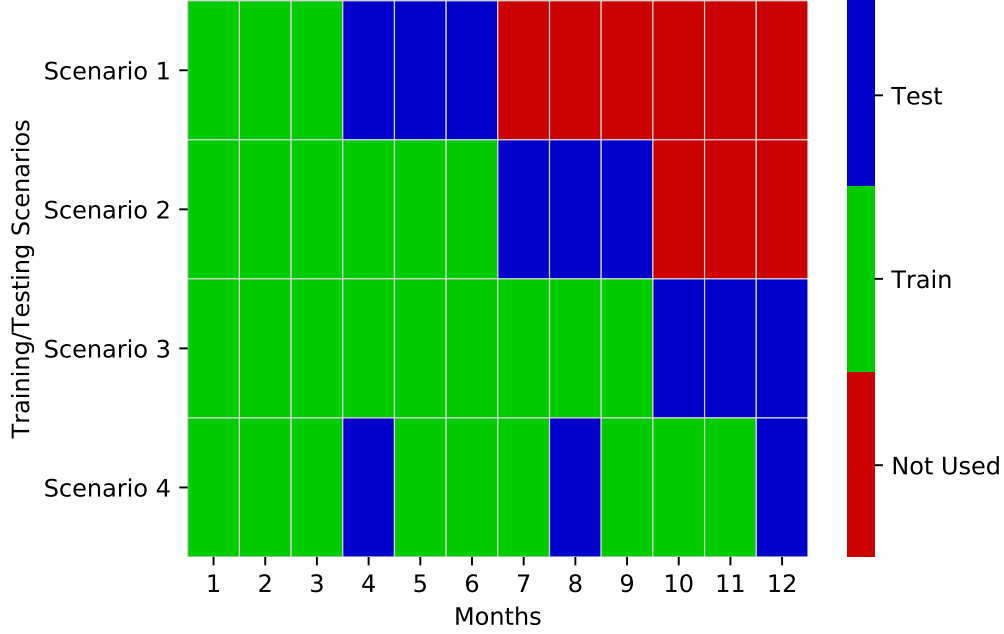


Figure 7: Overview of the four training/testing scenarios

219 This study focuses on general forecasting methods as a foundation for  
 220 comparison for more building domain specific regression or prediction models.  
 221 These models do not take into consideration the auto-correlation aspect of  
 222 prediction or the building context-specific nature of the built environment.

### 223 3.4. Accuracy metrics

224 The three metrics used in this analysis to evaluate model fit are the  
 225 MAPE, NMBE, and CVRSME metrics.

## 226 4. Results

### 227 4.1. Model fit overview

228 Figure 8 illustrates 12 Sci-Kit learn models applied to all the building use  
 229 types using the MAPE and CVRSME metrics. In general, laboratories have  
 230 the highest accuracy across all of the models and in general. This situation  
 231 is due to laboratories being the most *systematically schedule-driven* of all

the building types. Laboratories often have large equipment that is operated continuously or in set time schedules throughout the course of an entire year. University classrooms and offices behave in similar ways across the models as these two use types are often similar and many of these buildings are mixed use types. University dormitories have the fourth highest accuracy for most of the models, however they are better fits for models such as the Huber Regressor or the TheilSen Regressor. Finally, the Primary School buildings are the worst performers among the building use types. These buildings are dependent more on human behavior within their annual schedule phases, resulting in the lack of predictability using the methods and models outlined in this simple example.

Appendix A provides a much more detailed breakdown of each of the building use types according to all three error metrics as well as an analysis of each of the four training/testing scenarios.

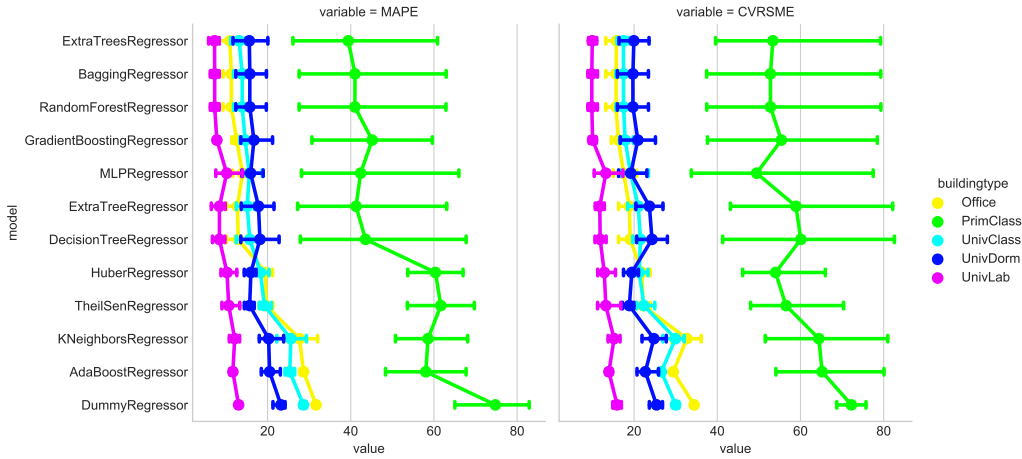


Figure 8: Overview of models fit metrics MAPE and CVRSME for each building use type across 12 Sci-kit Learn models

## 5. Conclusion

This publication gives readers a simple example of mainstream prediction techniques applied to a large, diverse data set. This example is designed to illustrate how future techniques can and should be applied to open data sets.

### 250 5.1. Limitations

251 The key limitation of this analysis is congruent with the limitation of  
252 machine learning in general: the models developed and the resulting metrics  
253 are limited by the range of training data utilized to build the models. The  
254 conclusions of the model comparison in this paper are only relevant to office,  
255 laboratories, classrooms, and dormitories from university campuses in the  
256 context of the geographical and operational environments of the Building  
257 Data Genome Project data set. Thus, the primary goal of this analysis is  
258 not to be comprehensive in capturing the behavior of the building stock, but  
259 to provide an example and data set that can be built upon with a more diverse  
260 set of data. *It is crucial that additional data from thousands (or millions) of*  
261 *other buildings are added over time for the methodology to achieve the main*  
262 *goal of a comprehensive benchmarking process.*

263 New algorithms coming into the public domain could be tested against  
264 this growing set of building data sets to quantify what improvements are  
265 being made in the domain.

### 266 5.2. Open Source Repository of Data and Implementation

267 To keep in the spirit of reproducibility, this publication is fully repro-  
268 ducible using the Building Data Genome Project.

### 269 5.3. Future work: the Great Building Energy Predictor Shootout 2019

270 A machine learning competition is in the planning phase to further extend  
271 the benchmark data set to over 4000 energy data streams from 1200+ build-  
272 ings from around the world. This competition will put potentially thousands  
273 of machine learning and building systems experts head-to-head to come up  
274 with the best set of methods to predict.

275 Key areas of potential improvement include:

- 276 1. Next generation time-series features - using the autocorrelated patterns  
277 in more advanced ways
- 278 2. Multi-variate learning - using a building's *peers* to predict future be-  
279 havior
- 280 3. Advanced models - using deep learning to leverage new and diverse  
281 types of behavior without feature engineering
- 282 4. Advanced ensembles - a voting ensemble that could significantly im-  
283 prove prediction for different building types
- 284 5. White or Gray Box model convergence - using physics-based models to  
285 inform the machine learning algorithm

## 6. References

- [1] A. Agrawal, J. Gans, A. Goldfarb, Prediction Machines: The simple economics of artificial intelligence, Harvard Business Press, 2018.
- [2] D. M. Solomon, R. L. Winter, A. G. Boulanger, R. N. Anderson, L. L. Wu, Forecasting energy demand in large commercial buildings using support vector machine regression, Department of Computer Science, Columbia University, Tech. Rep. CUCS-040-11 (2011).
- [3] C. E. Borges, Y. K. Penya, I. Fernández, J. Prieto, O. Bretos, Assessing tolerance-based robust short-term load forecasting in buildings, *Energies* 6 (2013) 2110–2129.
- [4] R. K. Jain, K. M. Smith, P. J. Culligan, J. E. Taylor, Forecasting energy consumption of multi-family residential buildings using support vector regression: Investigating the impact of temporal and spatial monitoring granularity on performance accuracy, *Applied Energy* 123 (2014) 168–178.
- [5] J. Granderson, P. N. Price, D. Jump, N. Addy, M. D. Sohn, Automated measurement and verification: Performance of public domain whole-building electric baseline models, *Applied Energy* 144 (2015) 106–113.
- [6] Efficiency Valuation Organisation, International performance measurement and verification protocol, Technical Report EVO 10000-1;, 2012.
- [7] K. Amasyali, N. M. El-Gohary, A review of data-driven building energy consumption prediction studies, 2018.
- [8] T. C. Yee, Stephanie, Model Tuning and the Bias-Variance Tradeoff, ????
- [9] J. Granderson, S. Touzani, C. Custodio, M. D. Sohn, D. Jump, S. Fernandes, Accuracy of automated measurement and verification (M&V) techniques for energy savings in commercial buildings, *Applied Energy* 173 (2016) 296–308.
- [10] J. Grandersona, S. Touzani, S. Fernandes, C. Taylor, Application of automated measurement and verification to utility energy efficiency program data, *Energy and Buildings* 142 (2017) 191–199.



317 [11] J. F. Kreider, J. S. Haberl, Predicting hourly building energy use: The  
 318 great energy predictor shootout – Overview and discussion of results  
 319 (????).

320 [12] J. S. Haberl, S. Thamilselan, Great energy predictor shootout II: Mea-  
 321 suring retrofit savings—overview and discussion of results, Technical Re-  
 322 port, American Society of Heating, Refrigerating and Air-Conditioning  
 323 Engineers~..., 1996.

324 **Appendix A. Detailed Model Comparison Breakdowns**

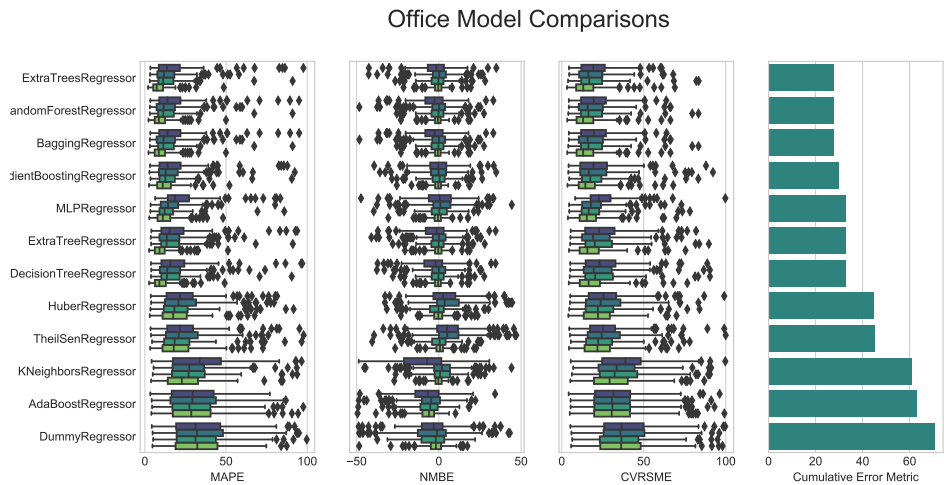


Figure A.9: Detailed Breakdown of Benchmarking Models on Office Buildings

Univ. Classroom Model Comparisons

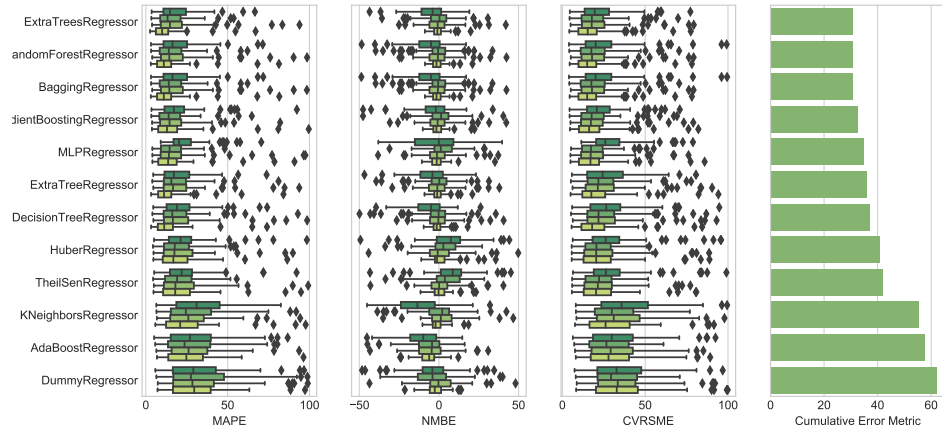


Figure A.10: Detailed Breakdown of Benchmarking Models on University Classroom Buildings

Univ. Lab Model Comparisons

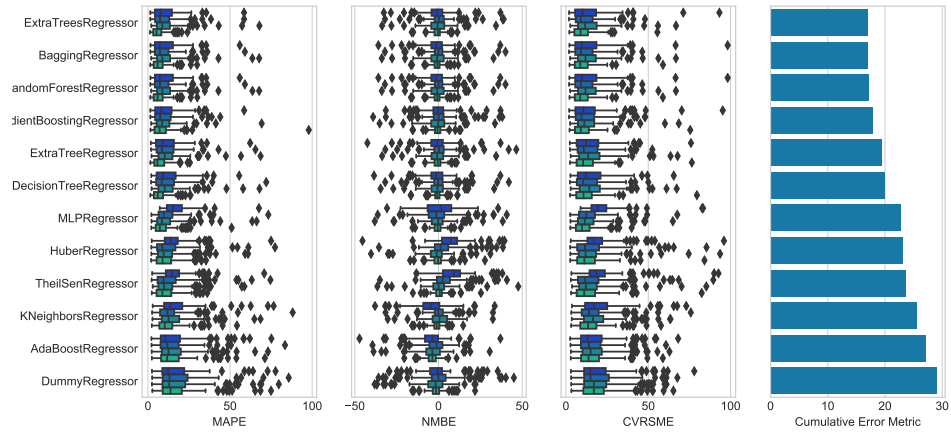


Figure A.11: Detailed Breakdown of Benchmarking Models on University Laboratory Buildings

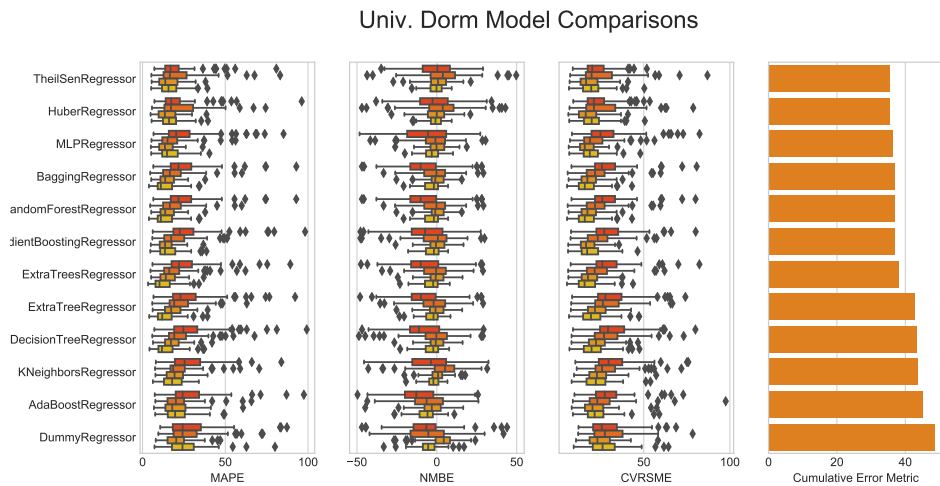


Figure A.12: Detailed Breakdown of Benchmarking Models on University Dormitory Buildings

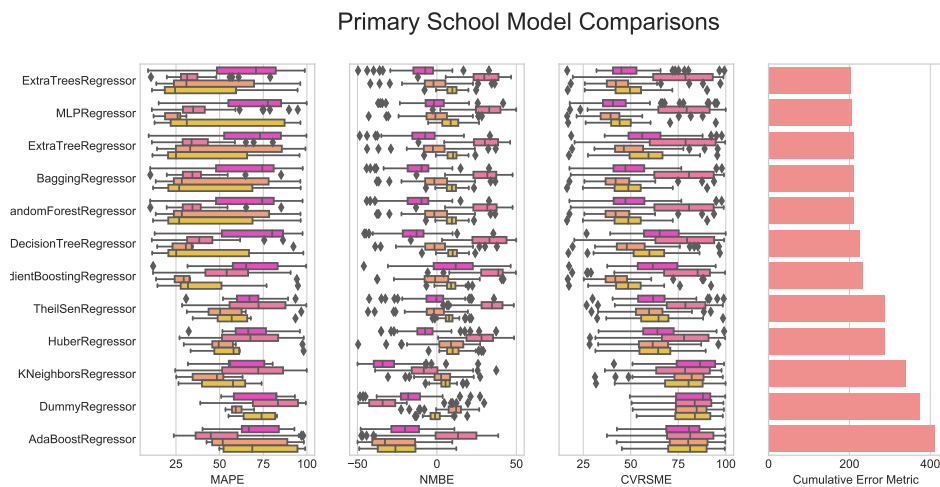


Figure A.13: Detailed Breakdown of Benchmarking Models on Primary School Buildings