

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

## **Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs**

M.Z. Naser, PhD, PE

<sup>1</sup>School of Civil & Environmental Engineering and Earth Sciences (SCEEES), Clemson University, USA

<sup>1</sup>Artificial Intelligence Research Institute for Science and Engineering (AIRISE), Clemson University, USA

E-mail: [mznaser@clemson.edu](mailto:mznaser@clemson.edu), Website: [www.mznaser.com](http://www.mznaser.com)

### **Abstract**

Machine learning (ML) has been shown to bypass key limitations of traditional methods (i.e., physical testing and numerical simulations) and hence presents itself as an attractive and effective technology in engineering. While the integration of ML brings exciting opportunities, it also brings unique challenges. One such challenge is related to the heavy reliance of ML on large datasets and computing facilities – for algorithm development, training, and storage. To realize Green ML (GML), this paper argues that ML users are to be cognizant of the hidden costs of energy consumption and subsequent carbon emissions arising from ML modeling, and hence they are ethically bound to apply ML responsibly. In this pursuit, a series of simple and exotic ML algorithms are examined, and their performance on a large dataset (~8,000 observations) is documented on five fronts; predictive performance, model size, training time, energy consumption, and equivalent carbon emissions. In addition, this work also examines the influence of algorithm architecture, processing language, number of features, as well as dataset size on model predictivity and energy consumption. Findings from this investigation infer that a 23-99% reduction in energy consumption and carbon emissions can be attained (while maintaining a comparable level of accuracy) by adopting simple as opposed to exotic ML models. The same findings have also led to the development of two new metrics that can tie the predictivity (i.e., level of accuracy) to the amount of energy consumed per algorithm. These metrics can be used to compare model performance in a similar manner to that traditionally used to assess the accuracy of ML predictions, thereby integrating energy-based awareness as a dimension of model comparison.

**Keywords:** Machine learning (ML); Structural engineering; Energy; Carbon emissions.

### **1. Introduction**

Ongoing advancements in machine learning (ML) have made it possible to extend this technology to various industries. One such industry is structural engineering which also happens to rely on numerical simulations for multi-scale problems [1]. While traditional simulation techniques (such as the finite element (FE) method) continue to be favorable given the familiarity of structural engineers with the fundamentals of such techniques, recent works have noted a rising interest in ML [2,3]. The outcome of such recent works, together with others [4–7], notes the potential of ML as a complementary technology that can offer novel solutions to structural engineers.

More specifically, a collection of structural engineering-based ML approaches have been reviewed in [8–10] with a common convergence of how ML can bypass the limitations of traditional methods and hence advocate for its adoption as we near the era of *Construction 4.0*. In addition, the application of ML is seen to be fruitful across a number of sub-disciplines belonging to

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

structural engineering. For example, ML models are reported to be able to accurately predict the properties of construction materials [11,12], detecting damage in structures [13,14], evaluating structural performance of elements [15,16], components [17,18], and structures [19]. In addition, ML has also been used as an aid tool for laying out blueprints and structural plans [20–22]. Recent trends in existing literature show that ML is being integrated into complex and unique problems as well [23,24].

Given the current interest in adopting ML, a look toward current practices [4–6,9,25–33] reveals a few interesting observations; 1) we continue to lack a unified procedure to apply ML to our problems<sup>1</sup>, 2) the application of ML is a user-derived operation wherein individual expertise influence and drive model development and deployment, and 3) ML is primarily applied to well-defined problems of a smaller search space as that commonly faced in other domains (e.g., medicine, or finance, etc.). These observations imply that we are yet to witness an upcoming boom in ML within structural engineering – which is likely to take place in the coming few years.

Therefore, it is of merit to proactively set the stage to mold the integration of ML in this domain. While it is virtually impractical to deliver one solution that fits all early into adopting ML, a few essential items can be emphasized and discussed. Such items may include the need for transparency, interpretability, fairness, and energy efficiency. The latter is the focus and motivation of this work.

With the growing need for larger yet accurate models, structural engineers will be expected to develop complex models. Such models are often tied to intricate architectures/topologies and are likely to require a tremendous amount of energy for training, development, storage, and deployment. To put this into perspective, two recent works have estimated the cost of training various deep learning models in terms of carbon emissions and monetary costs and reported estimates varying between 6.0-150,000 Kg and between \$41.0 to \$3.2 million per model [34,35]. The reader is to note that a human being, on average, contributes to about 5,000 Kg of carbon emissions on an annual basis, and 150,000 Kg of carbon emissions is equivalent to that emitted through the lifetime of five fuel-based vehicles [34].

On a larger scale, 2.3% of global carbon emissions are attributed to the information and communication technology sector, as estimated in a recent report by the Global e-Sustainability Initiative [36]. For comparison, the cement and aviation industry generates carbon emissions of about 5% and 1.7%, respectively [37,38]. Given the rise in recent calls to cut carbon emissions by 50% within the next ten years to negate the escalating rates of climate change, one means to align with such calls is to ensure considerate ML energy expenditure [35].

On a different front, exotic or highly advanced versions of existing ML models provide improvements in accuracy, among others; however, such improvements may or may not positively reflect upon the consumed resources to attain such improvements. One particular example that fits the former is the well-known computer vision model released in 2015 known as *ResNet*. The improved version of this model, *ResNeXt*, not only required 35% more computational resources to

---

<sup>1</sup> It is worth noting that the same is also true in other domains as well.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

train than *ResNet* but only achieved a 0.5% improvement in accuracy [39]. A second example of the latter is that related to the deep learning network, *SqueezeNet*, which attained the same accuracy as *AlexNet* (the winner of The ImageNet Large Scale Visual Recognition Challenge (ILSVRC) in 2012, with fifty times fewer parameters and a squeezed size of 0.5 MB as compared to 240 MB for *AlexNet* [40].

It is true that improved accuracy comes in handy; however, for *some applications*, minimal improvement in accuracy may not warrant the additional energy-based resources needed to realize such accuracy. For instance, predicting if a concrete mixture can develop strength of 38 or 40 MPa may not *significantly* alter the practical application of such a mixture since both strengths remain within the range of normal strength concrete. In this example, pushing a typical ML model to chase accuracy metrics of unity may not be warranted. Similarly, training an exotic model (e.g., an ensemble or a deep learning model) to predict the above when a simple model (i.e., a decision tree) can be trained to yield a comparable prediction accuracy may also be unnecessary. On a parallel front, developing a highly complex neural network to classify if structural members will undergo damage or not may lead to high energy consumption as compared to developing a leaner classifier – one that is based on logistic regression (assuming that both models can deliver a comparable performance).

Given the rise in ML adoption, which is also expected to initiate a ML-based race within the construction industry, one can appreciate the environmental and societal benefits of arriving at energy efficient ML models. Realizing energy efficient ML models, or strategies that enable cognizant and reduced-order modeling, can be of merit since these are the main tasks that ML users can control in our domain – as opposed to tackling other dimensions to this problem such as developing new algorithms or hardware, or storage, or tuning of cloud computing capabilities etc. While advancements in other dimensions are indeed necessary to realize a more energy efficient ML; practically speaking, investments in new hardware-like facilities by structural engineering firms are often expected to last for a few years, and hence embracing a user-driven approach is seen of equal importance, if not of more importance, to accelerate adopting Green ML (GML). This work argues that adopting cognizant and reduced-order modeling strategies can be seen as a sustainable approach to ML adoption. The same exercise will educate structural engineers to be inherently efficient and is expected to have long-lasting benefits.

This paper examines cognizant and reduced-order strategies that can be practiced by structural engineers to minimize energy consumption arising as a byproduct of developing ML models. In addition, this paper carries out a series of comparisons between simple and exotic ML algorithms to report on their predictive performance, model size, energy consumption, and storage need. More specifically, this work explores the influence of algorithm architecture, processing language, number of features, as well as dataset size<sup>2</sup>. The selected algorithms comprise; Gradient Boosted

---

<sup>2</sup> Given the wide possibilities and approaches of tuning algorithms (especially those of different topologies), it is deemed necessary to establish the following rationale early into this work. Thus, the reader is to realize that in order to enable reproduction of this analysis and to allow reporting true benchmarks on models' performance, all of the examined algorithms were applied in their default settings. As such, the goal of this work is *not* to arrive at a "solution" to the examined phenomenon, *nor* to identify the most "suited" algorithm to solve such a phenomenon, but rather to

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.jclepro.2022.135334>.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

Trees (GBT, including those of Extreme and Light Gradient nature), Keras Deep Residual Neural Network (KDNN), TensorFlow Deep Learning (TDL), Vowpal Wabbit (VW), and Random Forest (RF). The performance of these algorithms was compared in classifying if a given concrete mixture can attain specified design strength at in-situ conditions by assessing a large tabular dataset.

## 2. Materials and methods

### 2.1 Selected Machine Learning Algorithms

The selected algorithms in this work are briefly described herein, and their full description can be found in their respective references, as well as in [41–49]. These algorithms are among the most widely used algorithms in this research area, as noted in the following recent review papers [10,50].

#### 2.1.1 Variants of Gradient Boosted Trees (GBT)

A GBT trains a simple tree-like model and then uses the first model's error as a feature to build successive models. Focusing on errors obtained from a feature when building successive models reduces the overall error in the model. Two variants of GBT are the XGBoost and Light Gradient Boosting Machine (LGBM). The XGboost improves the performance of a GBT via a weighted quantile sketch (an approximation algorithm that determines how to split candidates in a tree) and the sparsity-aware split finding (which works on sparse data, as well as data with missing values). The XGboost uses a pre-sorted algorithm and a histogram-based algorithm for computing the best split. This algorithm was first published by Chen and Guestrin [43], and more details on this algorithm can be found in such work.

On the other hand, the LGBM algorithm by Microsoft [51] introduces two techniques to improve the performance of a GBT. These techniques are gradient-based one-side sampling (which identifies the most informative observations and skips those less informative) and exclusive feature bundling (which groups features in a near-lossless way). The LGBM is an improvement over both the XGBoost and traditional GBTs.

The code of the used XGBoost can be found online at [44,45]. This algorithm incorporates the following pre-tuned settings; learning rate = 0.05, maximum tree depth = 3.0, subsample feature = 0.5, number of boosting stages = 6,250, and minimum interval for early stopping = 200. The LGBM algorithm can also be found at [52] with the following default settings: learning rate = 0.05, maximum depth = "none", number of boosting stages = 6,250, and minimum interval for early stopping = 200. Finally, two GBTs were used (one of Python origin and another of R origin), which can also be found herein in their default settings [46,47]<sup>3</sup>.

---

report on findings from each of the examined algorithms (as obtained from their default settings) from an energy point of view. A full diagnostic test to explore the effect of hardware, explicit tuning of algorithms, as well as other algorithms that were not showcased herein, is deemed cumbersome to fit into one work and hence interested readers are invited to extend this work beyond its original message. The author hopes that the collective effort within this community and in the coming years will generate in depth insights to expedite the adoption of GML.

<sup>3</sup> This analysis acknowledges that each algorithm was developed with certain default settings as per the developers of each algorithm. Hence, algorithm settings were not set to be similar in GBT variants since: 1) other algorithms such as KDNN, TFDL, etc. also have unique settings, and 2) it is not possible to ensure that all algorithms will have similar

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

### 2.1.2 Keras Deep Residual Neural Network (KDNN)

Keras is a high-level library for developing neural networks [53]. In a residual network, a direct connection exists, linking data points to the outputs. Such a connection smoothens the loss function and enables better network optimization. In the used KDNN, default settings of a learning rate of 0.03 was used, along with a *Sigmoid* activation function, *Adam* optimizer, and different layer architectures (1 layer of 64 units, 1 layer of 1536 units, 2 layers of 64, and 64, and 3 layers of 256, 128, and 64 units). KDNN can be readily found at [48].

### 2.1.3 TensorFlow Deep Learning (TFDL)

A TFDL is an open-source and free neural network-based model that uses Deep Learning and is hosted by Google [54]. The used algorithm in its default settings (neurons in each layer = 100, number of training examples = 128, optimizer = *Adam*, adaptive learning rate, early stopping window = 5, and activation function of *ReLU*) can be found at [49].

### 2.1.4 Vowpal Wabbit (VW)

Vowpal Wabbit is a fast and out-of-core algorithm learner capable of streaming data (which comes in handy in large databases that cannot be supported by existing hardware). VW initially started at Yahoo and then moved to Microsoft. VW has a learning rate of 0.1, *logistic* loss function, power on the learning rate decay = 0.5, and can be found herein [55].

### 2.1.5 Support Vector Machines (SVM)

The Support Vector Machine (SVM) algorithm aims to identify a line or a boundary (i.e., hyperplane) in an  $n$ -dimensional space to classify data into separate classes [56]. This plane is found by maximizing the distance between data points of each class to allow for confident classification [27]. SVM uses a special form of mathematical function defined as kernels ( $k$ ). A kernel function transforms inputs into the required form. The used SVM utilizes penalty parameter = 12915, Gamma parameter = 0.208, and can be found at [57].

### 2.1.6 Logistic Regression (LR) and Elastic Net (EN)

The logistic regression model is a generalized linear model that applies a binomial distribution to fit regression models to binary response variables. The applied LR model herein had the following settings; Sigma =  $10^{-6}$ , fit intercept = *True*, and tolerance = 0.0001. The Elastic Net (EN) classifier is an extension of LR that applies L1 (estimate the median of the data) and L2 (estimate the mean of the data) prior as regularizers. Both models were taken from [58].

### 2.1.7 Random Forest (RF)

The Random Forest (RF) algorithm is an ensemble learner that forms a series of decision trees [59]. In a classification problem, the majority of predictions, as compared against all trees, are

---

settings (given that some settings which may exist for a particular algorithm, yet may not exist for another/all algorithms). As mentioned earlier, this analysis does *not* seek to identify the most “suited” algorithm, but rather to report on findings from each of the examined algorithms from an energy perspective.

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.jclepro.2022.135334>.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

used to consolidate the final outcome. Two RFs were used (one of Python origin and another of R origin), which can also be found herein in their default settings herein [60] and [61], respectively.

## 2.2 Description of Dataset

This section describes the examined dataset to be used in this work. This dataset was selected given its tabular nature, which is also likely to be present in most structural engineering problems [10,50]. As such, all observations were numeric, and no missing data points were present. This allows for somewhat of a fair field for comparison for the selected ML algorithms, given each's tendency to handle missing and/or categorical data differently (which may skew the primary focus of this investigation).

This dataset comprises compressive strength and mixture proportion of about 8,000 concrete mixtures as measured from in-situ conditions and reported by Young et al. [62,63]. The compressive strength of each concrete mixture at 28 days was measured following the ASTM C39 standard. In addition, mixture proportions including water, cement, and fly ash contents (in  $\text{kg/m}^3$  of concrete), water-reducing admixture (WRA), and air-entraining admixture contents (AEA in  $0.01\text{kg/kg}$  of cementitious material), coarse and fine aggregate contents (in  $\text{kg/m}^3$  of concrete), and fresh air content (in volume %) were reported. The focus of this dataset was to predict if a given concrete mixture will be able to develop its design compressive strength of concrete at in-situ conditions. As such, each measured mixture was compared against that of the design strength. If the attained strength is equal to or exceeds that of the design strength, then the mixture is deemed satisfactory. If not, then the mixture is labeled “*under-designed*”.

After cleansing the dataset of outliers and measurements with features of missing values, 6,647 measurements were labeled as *satisfactory*, and 1,099 measurements were labeled as *under-designed*. This cleansing process was carried out using a new approach that has been recently published by the author, wherein unsupervised and supervised learning methods are used to identify and remove outliers (see [64] for full details on this approach)<sup>4</sup>. Figure 1 shows additional insights into the statistical distributions of each feature used in this dataset.

---

<sup>4</sup> From a practical view, and for simplicity, the reader can assume that this dataset has two classes (Class 1 contains 6,647 satisfactory measurements and Class 2 with 1,099 under designed measurements).

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.jclepro.2022.135334>.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*. <https://doi.org/10.1016/j.jclepro.2022.135334>.

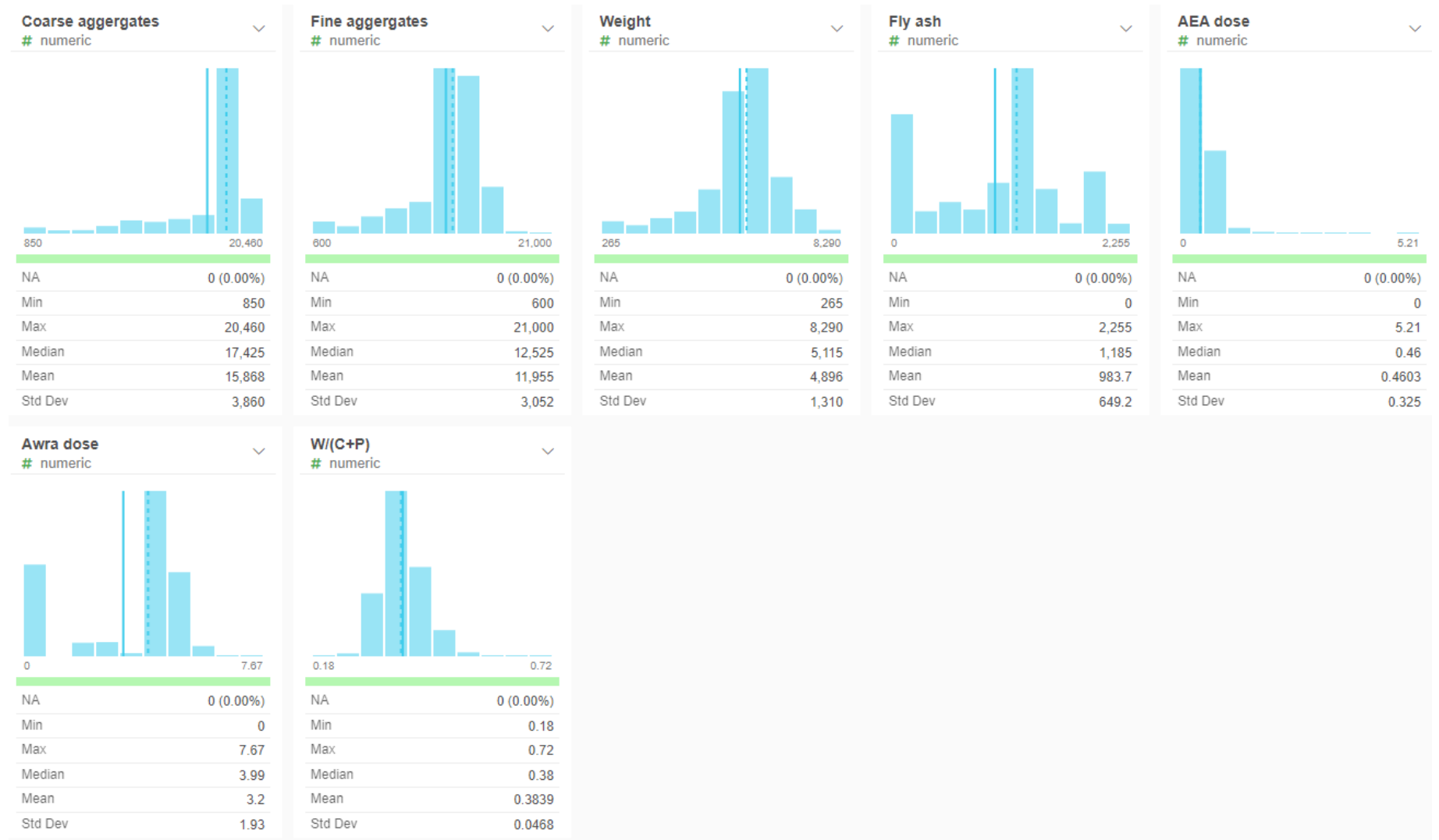


Fig. 1 Details on the dataset

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*. <https://doi.org/10.1016/j.jclepro.2022.135334>.

### 2.3 Model Development

Given the large range of features noted in the previous section, it is thought to apply data processing methods such as feature normalization [65]. In this process, the collected data is transformed such that the standard deviation is set equal to unity and the mean equal to zero [66]. Then, additional steps were taken to ensure the proper performance of the model. For a start, the dataset was randomly shuffled to eliminate the influence of neighboring measurements and data points (especially those taken within one in-situ job). Then, the shuffled dataset is split into two sets – a training and a testing set. The split of choice was 70:30 as per recommendations of recent works [67,68].

The training set was further processed via 10-fold cross-validation technique which became handy in ensuring the proper distribution of samples and hence allowed the algorithms to be examined at different subsets. In this technique, the training set is split into 10 equal-sized subsets such that 9 subsets are used for training, and the remainder of these is kept for validation [69]. All subsets were the same across all algorithms.

In addition to adopting a 10-fold cross-validation technique, model predictions were also checked against a series of performance metrics that measure the closeness of a predicted outcome to that predicted by a given ML model [70,71]. In this work, metrics that are commonly used by the structural engineering domain are selected [29,72–74]. Three classification metrics are used in this classification-based problem. These metrics are accuracy (ACC), Area under the ROC curve (AUC), Log Loss Error (LLE). Additional details on each metric are shown herein.

$$ACC = \frac{TP+TN}{P+N} \quad (1)$$

where,  $P$ : predictions,  $N$ : number of real negatives,  $TP$ : number of true positives,  $TN$ : number of true negatives,

- Evaluates the ratio of the number of correct predictions to the total number of samples.
- Presents performance at a single class threshold only.
- Assumes equal cost for errors [75].

$$AUC = \sum_{i=1}^{N-1} \frac{1}{2} (FP_{i+1} - FP_i) (TP_{i+1} - TP_i) \quad (2)$$

where,  $FP$  number of false positives,  $FN$  number of false negatives.

- A value of unity indicates an accurate prediction.

$$LLE = - \sum_{c=1}^M A_i \log P \quad (3)$$

where,  $M$ : number of classes,  $c$ : class label,  $y$ : binary indicator (0 or 1) if  $c$  is the correct classification for a given observation.

- Penalizes for being confident in the wrong prediction.
- Has a probability between zero and unity.
- A lower value for log loss is favorable.



Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

### 3. Results and Discussion

This section evaluates the performance of all of the selected ML models upon the presented dataset from two main perspectives: predictivity and estimated energy consumption. In addition, this section highlights the influence of the dataset size, the number of features used in model development, ML model architecture, and programming language in more detail and from each of the aforementioned two perspectives. Finally, a discussion on the estimated carbon emissions emitted per algorithm is also presented. It is worth noting that a private cloud computing service was obtained for this analysis, and hence all algorithms were granted access to the same computation and hardware resources. As such, the influence of such resources on model performance is not discussed as it is normalized.

#### 3.1 Insights from a Predictivity Perspective

The selected ML models were first evaluated from a predictivity point of view (as in how well the predictions of each model match with the measured observations). Thus, predictions from all models were compared utilizing the ACC, AUC, and LLE performance metrics. Figure 2 and Table 1 articulate the outcome of this comparison. Table 1 shows that most, if not all, models seem to capture the phenomenon on hand adequately. The same table also shows that all models achieved comparable performance across all metrics and against training, validation, and testing splits. While some variations exist between all models in the reported metrics, a look into Fig. 2 infers that these variations can be considered minor (as per the scale of the horizontal axis).

For example, the lowest ACC value obtained against the full dataset was 0.946 by the KDNN with one layer of 64 units during its training procedure (which also happens to be the only value outside of 0.983 (or 98.3%) in ACC or AUC metrics. Beyond that, all models performed comfortably within the 98.6-100.0% range. In the case of the AUC metric, the lowest value was noted by the EN model with 99.1%, with all other models reaching the range of 99.2-100.0%. All models also performed well when their performance was compared using the LLE metric. Except for the KDNN with one layer of 64 units, all models scored within 0.05, with the majority scoring in the range of 0.01-0.004 (note that a score of 0.0 implies a perfect model).

As expected, and as Fig. 2 and Table 1 show, not a particular model scored the highest in all metrics. As such, a new composite metric was developed to enable a unified and more accessible comparison between all models developed in this study (from a predictivity perspective). This composite metric was created by combining ACC, AUC, and LLE. The synthetic metric acknowledges that higher values of ACC and AUC and lower values of LLE infer good predictivity. This metric is listed below, and Fig. 2d shows that the best performing model based on this metric is the XGboost model, followed by RF and LGBM<sup>5</sup>.

$$\text{Predictivity } (P) = \sum_i^n \frac{ACC \times AUC}{LLE} \quad (4)$$

---

<sup>5</sup> Despite this comparison, the reader is reminded that the goal of this work is not to identify the best performing ML model but rather to showcase that in many instances, a variety of ML models can inherently be good candidates for solving a particular problem.

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.jclepro.2022.135334>.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

where,  $i$  = training, validation, and testing. Higher values of this metric are favorable; for example, for a hypothetical case of ACC and AUC of unity, and LLE = 0.001, this metric yields 1000.

### 3.1.1 Influence of dataset size

Given the large size of our dataset (~8,000 observations), it is thought to explore the influence of reducing the size of the dataset upon the predictivity of all selected models (on training, validation, and testing split by examining the ACC, AUC, LLE and P metrics). Thus, only the training split is changed from 70% to 25% while controlling all other aspects of model development (i.e., the analysis is re-run by reducing the number of used observations from 70% of the whole dataset to 25% with the remainder (45%) of data used in the initial training omitted from the analysis). The outcome of this analysis is shown in Fig. 2 and Table 1.

When Figs. 2a, 2b, 2c, and 2d are compared against Figs. 2e, 2f, 2g, and 2h, one can see the following with regard to the performance of models in the full and reduced dataset; 1) the variation between models performance in terms of ACC and AUC is very comparable; with ML models performing slightly better in the reduced dataset, 2) LLE metric performance in the whole dataset is better than the reduced dataset (especially in the case of KDNN variants), 3) the performance of models in terms of the composite metric is higher in the case of the full dataset which reflects the larger degree of changes in LLE as to those observed in ACC and AUC, 4) despite XGBoost remaining the highest-ranked model, minor changes to overall models ranking occurred, and 5) from a practical engineering point of view, the performance of all models is adequate which may infer the suitability of all models from a predictivity perspective.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*. <https://doi.org/10.1016/j.jclepro.2022.135334>.

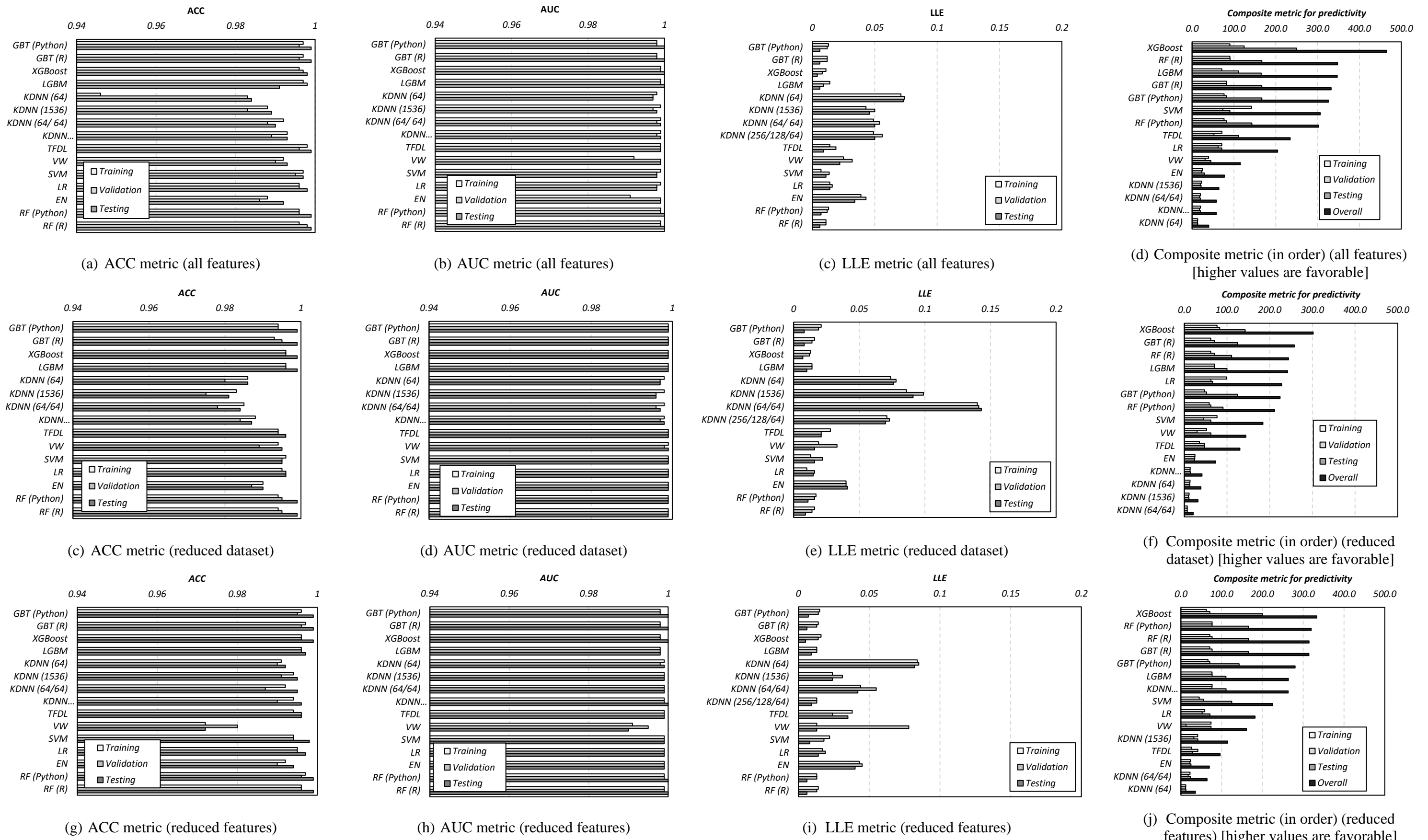


Fig. 2 Comparison of metrics

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*. <https://doi.org/10.1016/j.jclepro.2022.135334>.

Table 1 Evaluation of ML models.

| Cases            | Model   | ACC**             | AUC               | LLE               | Model size (MB)   | Training time (Sec) | Prediction time (sec)*** |       |
|------------------|---------|-------------------|-------------------|-------------------|-------------------|---------------------|--------------------------|-------|
| Full features    | GBT     | Python            | 0.997/0.996/0.999 | 0.998/0.998/1.000 | 0.013/0.012/0.006 | 0.337               | 389                      | 0.523 |
|                  |         | R                 | 0.997/0.996/0.999 | 0.998/0.998/1.000 | 0.012/0.012/0.006 | 0.259               | 391                      | 2.567 |
|                  | XGBoost |                   | 0.996/0.997/0.998 | 0.999/0.999/1.000 | 0.011/0.008/0.004 | 1.239               | 437                      | 0.478 |
|                  | LGBM    |                   | 0.997/0.998/0.991 | 0.999/0.999/1.000 | 0.014/0.009/0.006 | 1.130               | 610                      | 0.428 |
|                  | KDNN    | 64*               | 0.946/0.983/0.984 | 0.998/0.997/0.997 | 0.071/0.074/0.073 | 0.271               | 677                      | 0.851 |
|                  |         | 1536              | 0.988/0.983/0.989 | 0.999/0.997/0.998 | 0.043/0.051/0.046 | 4.757               | 197                      | 0.895 |
|                  |         | 64/64             | 0.992/0.988/0.990 | 0.999/0.998/0.999 | 0.049/0.054/0.050 | 0.290               | 677                      | 0.835 |
|                  |         | 256/128/ 64       | 0.993/0.989/0.993 | 0.999/0.998/0.999 | 0.049/0.056/0.050 | 0.446               | 677                      | 1.300 |
|                  | TFDL    |                   | 0.998/0.996/0.999 | 0.999/0.999/0.999 | 0.014/0.019/0.009 | 1.818               | 318                      | 0.778 |
|                  | VW      |                   | 0.992/0.990/0.993 | 0.992/0.999/0.999 | 0.025/0.032/0.022 | 0.223               | 696                      | 0.404 |
|                  | SVM     |                   | 0.997/0.995/0.997 | 0.999/0.998/0.998 | 0.007/0.014/0.011 | 2.162               | 233                      | 0.490 |
|                  | LR      |                   | 0.996/0.996/0.998 | 0.999/0.998/0.998 | 0.014/0.016/0.014 | 0.152               | 259                      | 0.430 |
|                  | EN      |                   | 0.988/0.986/0.992 | 0.991/0.999/0.999 | 0.039/0.043/0.034 | 0.223               | 322                      | 0.489 |
|                  | RF      | Python            | 0.996/0.996/0.999 | 0.999/0.999/1.000 | 0.013/0.012/0.007 | 0.626               | 366                      | 0.427 |
| R                |         | 0.996/0.998/0.999 | 0.999/0.999/1.000 | 0.011/0.011/0.006 | 1.460             | 386                 | 2.657                    |       |
| Reduced dataset  | GBT     | Python            | 0.994/0.994/0.999 | 0.999/0.999/0.999 | 0.021/0.019/0.008 | 0.201               | 402                      | 0.212 |
|                  |         | R                 | 0.993/0.995/0.999 | 0.999/0.999/0.999 | 0.016/0.014/0.008 | 0.144               | 366                      | 2.581 |
|                  | XGBoost |                   | 0.996/0.996/0.999 | 0.999/0.999/0.999 | 0.013/0.012/0.007 | 0.573               | 446                      | 0.259 |
|                  | LGBM    |                   | 0.996/0.996/0.999 | 0.999/0.999/0.999 | 0.014/0.014/0.010 | 1.452               | 246                      | 0.255 |
|                  | KDNN    | 64                | 0.986/0.980/0.986 | 0.998/0.997/0.997 | 0.074/0.078/0.076 | 0.184               | 785                      | 0.745 |
|                  |         | 1536              | 0.983/0.975/0.981 | 0.998/0.996/0.996 | 0.086/0.099/0.091 | 1.816               | 857                      | 0.580 |
|                  |         | 64/64             | 0.985/0.978/0.984 | 0.998/0.996/0.997 | 0.140/0.141/0.143 | 0.204               | 842                      | 0.625 |
|                  |         | 256/128/ 64       | 0.988/0.984/0.987 | 0.998/0.997/0.998 | 0.071/0.073/0.070 | 0.359               | 750                      | 0.556 |
|                  | TFDL    |                   | 0.994/0.994/0.996 | 0.999/0.999/0.999 | 0.028/0.021/0.021 | 3.881               | 310                      | 0.586 |
|                  | VW      |                   | 0.994/0.989/0.995 | 0.999/0.998/0.999 | 0.019/0.033/0.016 | 0.136               | 410                      | 0.224 |
|                  | SVM     |                   | 0.996/0.995/0.995 | 0.999/0.999/0.999 | 0.013/0.022/0.016 | 2.067               | 198                      | 0.271 |
|                  | LR      |                   | 0.995/0.996/0.996 | 0.999/0.999/0.999 | 0.010/0.016/0.015 | 0.066               | 226                      | 0.249 |
|                  | EN      |                   | 0.990/0.987/0.990 | 0.999/0.999/0.999 | 0.040/0.040/0.041 | 0.137               | 440                      | 0.249 |
|                  | RF      | Python            | 0.994/0.995/0.999 | 0.999/0.999/0.999 | 0.017/0.016/0.011 | 0.336               | 333                      | 0.334 |
| R                |         | 0.994/0.995/0.999 | 0.999/0.999/0.999 | 0.016/0.014/0.009 | 0.929             | 492                 | 2.670                    |       |
| Reduced features | GBT     | Python            | 0.996/0.995/0.999 | 0.998/0.998/1.000 | 0.015/0.014/0.007 | 0.333               | 269                      | 0.367 |
|                  |         | R                 | 0.997/0.996/0.999 | 0.998/0.998/1.000 | 0.014/0.013/0.006 | 0.258               | 179                      | 2.595 |
|                  | XGBoost |                   | 0.996/0.996/0.999 | 0.998/0.998/1.000 | 0.016/0.014/0.005 | 0.933               | 458                      | 0.595 |
|                  | LGBM    |                   | 0.996/0.996/0.997 | 0.998/0.998/0.998 | 0.013/0.013/0.009 | 5.886               | 337                      | 0.482 |
|                  | KDNN    | 64                | 0.991/0.990/0.992 | 0.999/0.998/0.999 | 0.084/0.085/0.082 | 0.221               | 417                      | 0.835 |
|                  |         | 1536              | 0.994/0.991/0.995 | 0.999/0.999/0.999 | 0.024/0.031/0.024 | 2.880               | 217                      | 0.773 |
|                  |         | 64/64             | 0.992/0.987/0.995 | 0.999/0.999/0.999 | 0.044/0.055/0.042 | 0.239               | 364                      | 0.756 |
|                  |         | 256/128/ 64       | 0.994/0.990/0.996 | 0.999/0.999/1.000 | 0.013/0.013/0.009 | 0.390               | 406                      | 0.815 |
|                  | TFDL    |                   | 0.994/0.996/0.996 | 0.999/0.999/0.999 | 0.038/0.024/0.035 | 4.061               | 362                      | 0.675 |
|                  | VW      |                   | 0.972/0.980/0.972 | 0.991/0.995/0.990 | 0.013/0.078/0.013 | 0.176               | 527                      | 0.392 |
|                  | SVM     |                   | 0.994/0.994/0.998 | 0.999/0.999/0.999 | 0.022/0.018/0.008 | 2.115               | 184                      | 0.450 |
|                  | LR      |                   | 0.995/0.995/0.997 | 0.999/0.999/0.999 | 0.017/0.019/0.014 | 0.151               | 134                      | 0.480 |
|                  | EN      |                   | 0.992/0.990/0.994 | 0.999/0.999/0.999 | 0.043/0.045/0.040 | 0.176               | 290                      | 0.400 |
|                  | RF      | Python            | 0.997/0.996/0.999 | 0.999/0.999/1.000 | 0.013/0.013/0.006 | 0.579               | 103                      | 0.459 |
| R                |         | 0.996/0.996/0.999 | 0.999/0.999/1.000 | 0.014/0.013/0.006 | 1.403             | 171                 | 2.661                    |       |

\*Training/validation/testing. \*\*Number of units in each layer. \*\*\*To score 1,000 observations.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*. <https://doi.org/10.1016/j.jclepro.2022.135334>.

### 3.1.2 Influence of number of features

The influence of the number of features used in model development was also investigated. As such, the importance of features in all models (based on the full dataset analysis) was measured via the SHAP method [76] and presented in Fig. 3. This figure clearly shows that the most reoccurring features among all models are coarse aggregates, fine aggregates, and weight. Thus, the features in the full dataset were reduced to only these three features (while keeping the number of observations (~8,000) and analysis procedure the same).

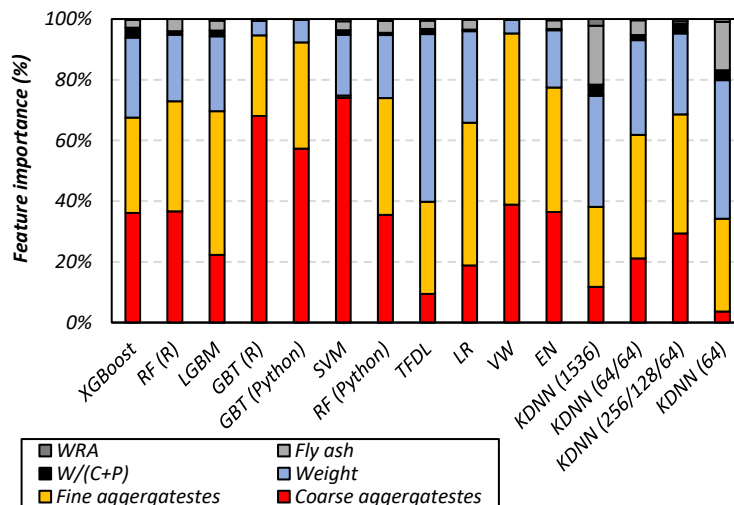


Fig. 3 Feature importance as measured in all models

A look into the ACC, AUC and LLE metrics reveals that the performance of all models is comparable, if not improved, to that of the case of the analysis that incorporates all features. The VW model shows odd performance with a slight reduction in predictivity, which could be related to this model being heavily reliant upon the two aggregate features with little to nothing allocated for other features. Figure 2j shows the ranking of all models as per the reduced number of features. As one can see, the XGBoost remains the highest ranking model, followed by the RF and GBT variants.

### 3.1.3 Influence of programming language

Results from the conducted analysis can also be used to investigate the influence of adopting similar models but from different programming languages (see Fig. 3 and Fig. 4). In this analysis, ML models were used of different languages RF (Python and R), as well as GBT (Python and R). For a start, the examined variants realized similar feature importance scores for most involved features, which implies consistency. Looking at the final outcome of the predictivity-based examination shows that the R version of both models performed well ahead of the Python version.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*. <https://doi.org/10.1016/j.jclepro.2022.135334>.

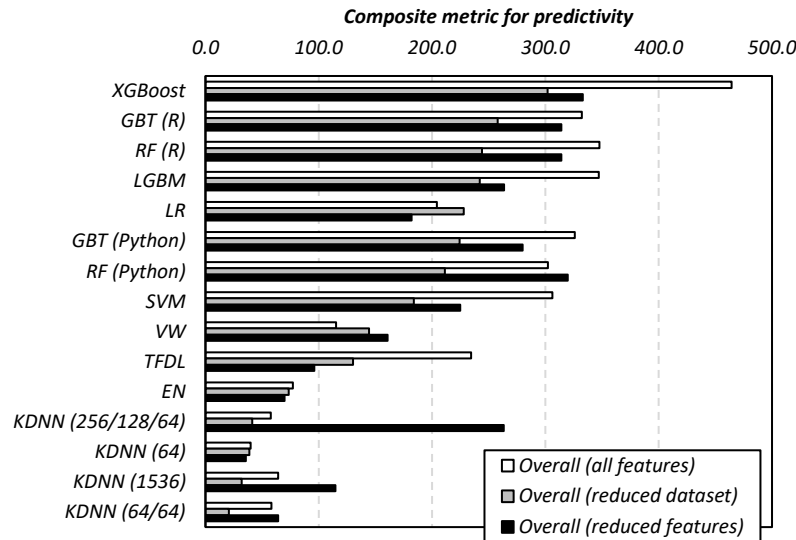


Fig. 4 Comparison of predictivity metric

### 3.1.4 Influence of model architecture

Figures 3 and 4 can also be used to explore the influence of model architecture. In this exercise, four architectures of the KDNN model were investigated by changing the number of units and layers in each model. The outcome of this investigation shows that these four models achieved consistent feature importance measurements (especially in the case of fine aggregates and weight). In addition, these models performed at the bottom of all other models in terms of predictivity (however, the reader must be reminded that all models performed adequately from a practical point of view).

The highest performing model of the KDNN variant in the case of the full dataset was that of one layer and 1536 units, followed by two layers (64/64), three layers (256/128/64), and one layer of 64 units. Figure 4 shows that the predictivity of the KDNN (256/128/64) variant significantly improved, bypassing all observed improvements of other models when applied to the case of reduced features. Furthermore, KDNN (64) seems to be the most stable variant in terms of predictivity. Acknowledging the limited number of programming languages and model architectures examined herein, it would be interesting to extend this investigation to other cases.

### 3.2 Insights from an Energy Perspective

The performance of the selected ML models was also examined from an energy perspective wherein three metrics were explored, namely, model size (in MB), training time (in seconds), and prediction time to score 1,000 rows (in seconds). The outcome of this examination is shown in Fig. 5 and Table 1. As one can see, the largest and smallest model size for the case of the full dataset was obtained from a KDNN (with 1536 units<sup>6</sup>) and LR, respectively. Overall, all ML models were

<sup>6</sup> This model also happens to be the fastest training model.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

of size within 0.152-4.757 MB and completed training within 197-696 seconds. On the other hand, prediction times varied between 0.404-2.657 seconds (with VW and RF being at these extremes, respectively).

Noting how faster training and prediction times and smaller model size (for the same workstation or cloud service) can be tied to lower energy expenditure, then all ML models were to be compared on this basis [77–79]. A new metric is then developed to account for the above rationale. This metric combines model size with training and prediction time such that:

$$\text{Estimated Energy } (E) = MS \times (TT + PT) \quad (5)$$

where, MS: model size, TT: training time, and PT: prediction time. Smaller values are favorable with a hypothetical minimum value = 1.0 MB × 10 sec = 10 MB.sec.

Figure 5 shows the results of applying this new metric. It is clear that the most energy efficient model is that of LR, followed by EN and the R version of GBT. In addition, once all models were normalized by the KDNN (with 1536 units), which happens to score the highest in the case of the full dataset analysis, one can appreciate the reduction of energy consumed by LR, which is estimated at 4.0% of that of the KDNN (with 1536 units).

### 3.2.1 Influence of dataset size

Surprisingly, reducing the dataset size was not always reflected by reducing model size, training time, or prediction time (see Fig. 5). While in large models such as KDNN (with 1536 units), this reduction results in a comparable reduction in model size, the same resulted in a significant increase in training time. In some instances, reducing the dataset size did not seem to affect model size much (i.e., SVM, KDNN (64, and 64/64)), nor prediction time (GBT (R), RF (R)). In all cases, the best performing models in terms of energy expenditure yielded the lowest composite metric, e.g., LR, GBT (R), and EN.

### 3.2.2 Influence of number of features

Figure 5 shows that the selected ML models behaved differently in response to feature reduction. For example, the model size of LGBM significantly increased despite a 40-50% reduction in training time. TFDL also underwent a similar performance to that seen in LGBM from a model size point of view. Finally, some models seem to have higher stability across all scenarios (i.e., XGBoost, LR, EN, among others, as seen in Fig. 5d).

### 3.2.3 Influence of programming language

Unlike results from the conducted analysis in the case of predictivity, the influence of different programming languages was not as apparent when compared from an energy consumption perspective. In this event, GBT (of both Python and R versions) seems to achieve a lower energy consumption than that of the RF variants. In fact, the RF (R) variant is noted to consume about 64% times more energy than the Python variant (as opposed to the GBT (Python) variant needing 48% more energy than the R variant) when compared based on the full dataset. Collectively, the Python variants achieved a smaller model size and shorter prediction times than the R variants. When compared to the energy consumption of LR, both R and Python variants consumed a more

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.jclepro.2022.135334>.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

considerable amount of energy in the range of 2.31-30.78 more than that of LR (with the Python variants scoring in the range of 2.95-7.50 and the R variants being on the higher side of the aforementioned range).

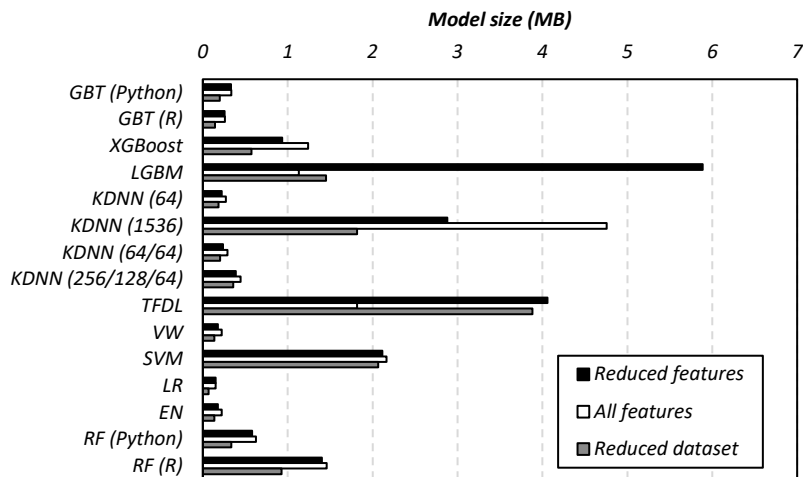
### 3.2.4 Influence of model architecture

The four previously discussed architectures of the KDNN model were also investigated to evaluate the influence of model architecture from an energy consumption perspective. As expected, our analysis indicates that simple architectures seem to consume less energy than those of deeper and heavier nature. More specifically, KDNN (64) and KDNN (1536) scored an E metric of 196.57, and 941.39 (for full dataset), 144.58, and 1557.37 (for reduced dataset), and 92.34, and 627.19 (for reduced features). Similarly, KDNN (64) and KDNN (64/64) scored an E metric of 196.57 and 196.57 (for full dataset), 144.58, and 171.90 (for reduced dataset), and 92.34, and 87.18 (for reduced features). A close examination of the trends shown in Fig. 5 reveals that energy consumption is positively correlated to the total number of units in KDNN variants. Noting the LR scored the lowest score as per the E metric, then a comparison between KDNN variants and LR shows that these variants require 4.66-23.87%, 4.29-30.88%, and 9.68-104.31% more energy for the whole dataset, reduced dataset, and reduced features, respectively. To echo the previous section, future works are invited to further examine the performance of different architectures and programming languages.

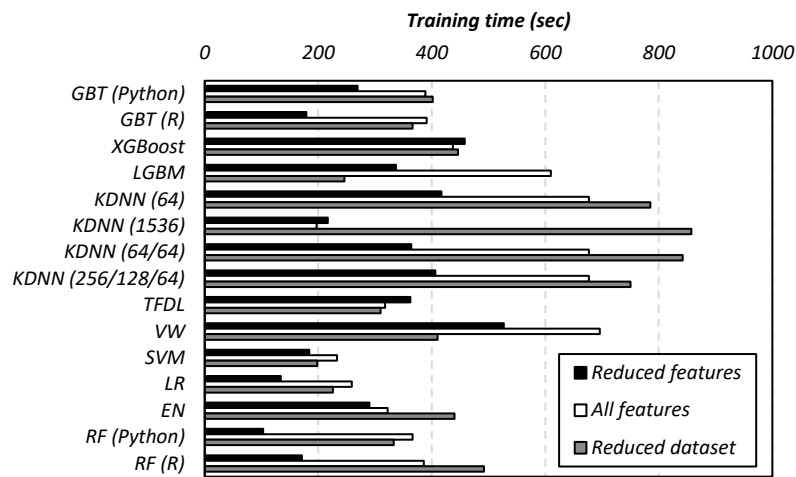


Please cite this paper as:

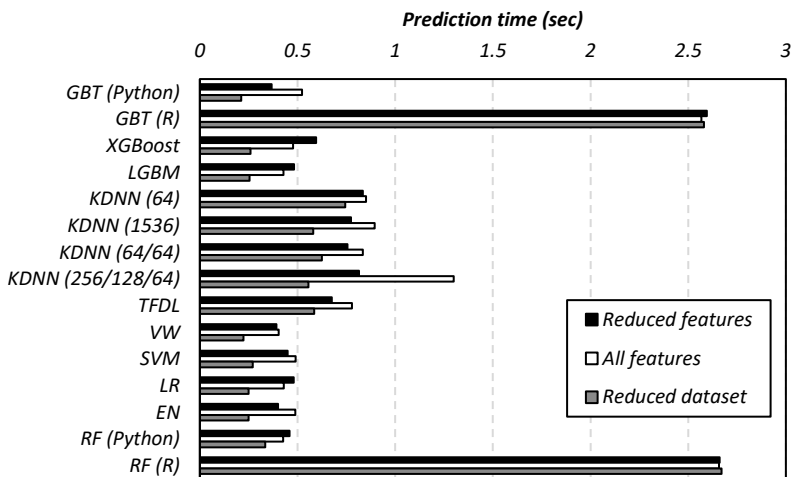
Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*. <https://doi.org/10.1016/j.jclepro.2022.135334>.



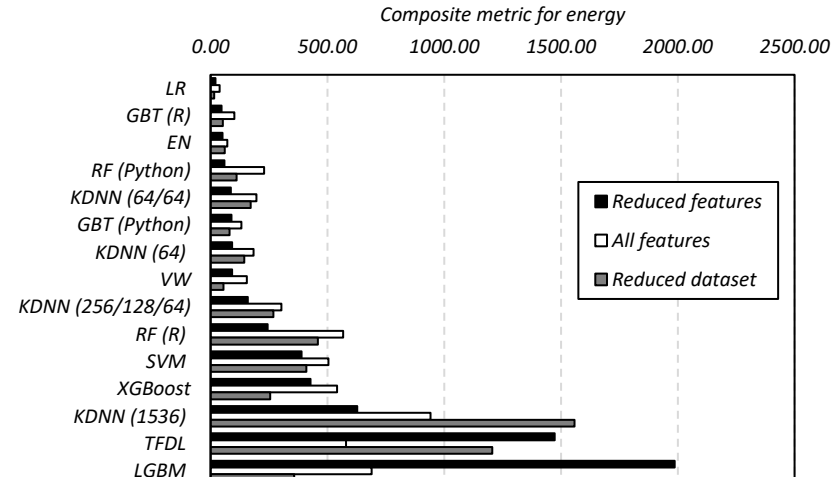
(a) Model size (MB)



(b) Training time (sec)



(c) Prediction time (sec)



(d) Composite metric [lower values are favorable]

Fig. 5 Comparison of metrics

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

### *3.3 Insights from a Holistic Perspective*

In order to tie the predictivity of ML models to their energy consumption, this section presents a visual aid employing a combined chart (see Fig. 6). This chart divides the models into four clusters and identifies models according to their combination of predictivity and energy consumption.

It is clear from Fig. 6 that RF (Python), GBT (Python), and GBT (R) consistently lay in the quadrant of high predictivity/low energy consumption. In addition, KDNN (1536) was found to reside in the low predictivity and high energy consumption, and the rest of the models were split into low predictivity/low energy consumption and high predictivity/high energy consumption. More exotic models, including XGBoost, LGBM, TFDL, and SVM, managed to cluster into one group (high predictivity/high energy consumption), and TFDL fell into the low predictivity/high energy consumption on two occasions.

It is worth noting that the XGBoost model achieved the highest predictivity in all cases and scored 0.58 on energy consumption (a slight 8% increase over the cut-off of 50%) in the case of the full dataset (while being on the lower energy side in the other two cases). As per the previous section, the LR model scored the lowest energy consumption in all cases.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*. <https://doi.org/10.1016/j.jclepro.2022.135334>.

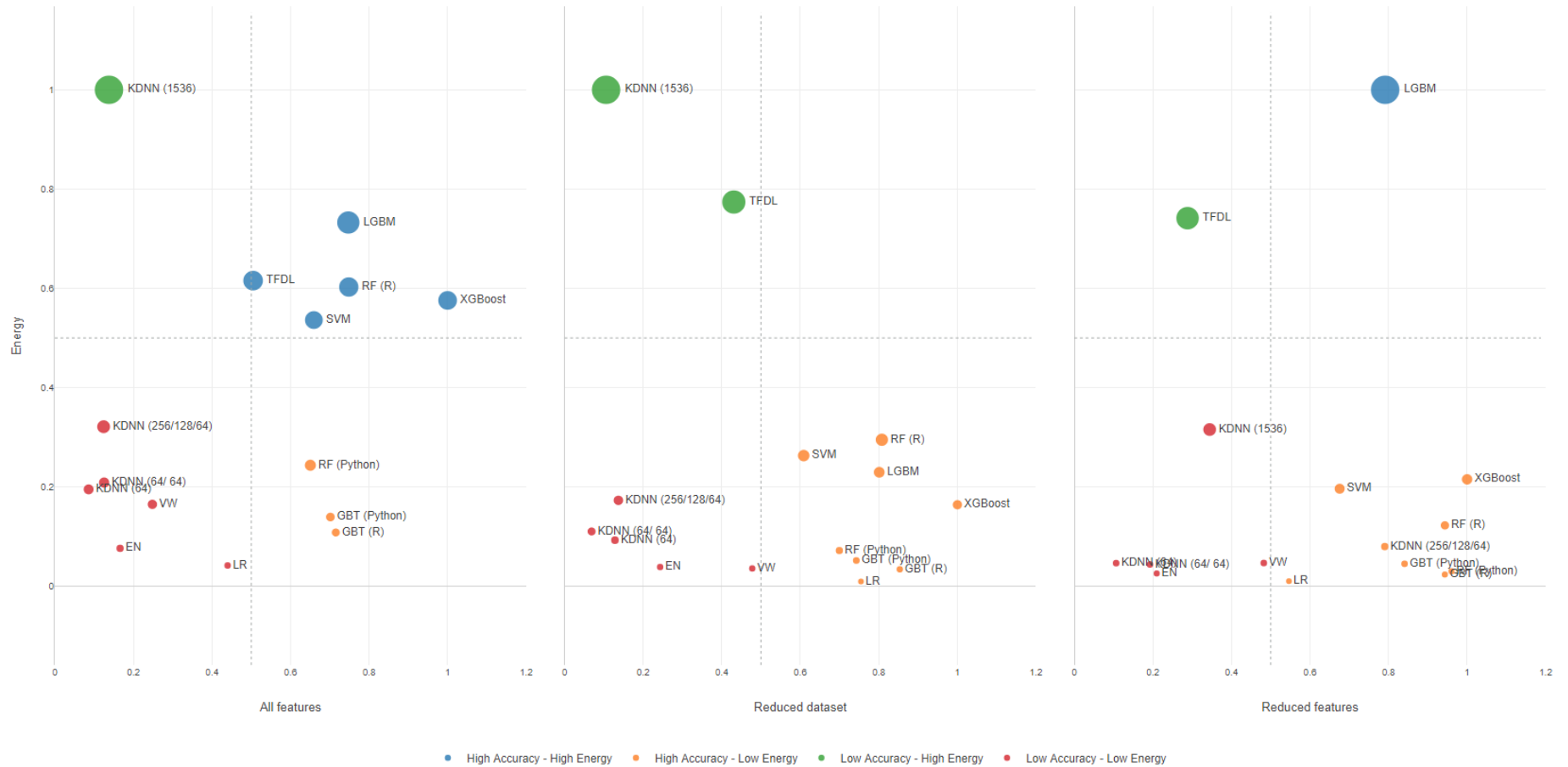


Fig. 6 Comparison of ML models when normalized by the highest attainable predictivity and energy consumption [size of the bubble represents the magnitude of estimated carbon emissions]

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

A look into Fig. 6 further shows interesting observations with regard to the programming language used and model architecture. For example, both GBT variants landed on the orange cluster (high predictivity/low energy consumption). The R version of RF deviated from this trend, given its high predictivity score. In addition, most KDNN variants (with the exception of KDNN (1536)) were clustered in the red cluster (low predictivity and low energy consumption), with the common notion of all variants being on the low predictivity side. The trend of varying programming languages seems to be negligible in the case of the reduced dataset and reduced features. The variation of model architecture remains consistent in the cases of the reduced dataset and reduced features.

### 3.3.1 A rough estimation of carbon emission

The above investigation can be further used to roughly estimate the amount of carbon emissions generated to train the selected models, and then this estimation can be tied to the predictivity and energy consumption of each model. One must realize that the open literature does not seem to contain a straightforward method to relate energy consumption of model development to carbon emissions as model development is governed by complex factors and is highly dependent upon the combination of software/hardware (e.g., type and model of CPU, GPU, processing units, etc.), geographical location of workstations (wherein some regions advocate for carbon neutrality by imposing a “carbon tax” for carbon generation and capture while others do not), and resources used in generating electricity (renewable resources vs. nonrenewable resources), etc. [34,78,80]. As such, our analysis will be based on a hypothetical scenario that applies available values as collected from their original resources.

According to the U.S. Energy Information Administration, 0.42 Kg of CO<sub>2</sub> is emitted per one kWh [81]. A study by the Lawrence Berkeley National Laboratory, USA, [82] estimated the unit annual energy consumption (AEC) for desktops to be 194 kWh/yr (with a median = 125 kWh/yr), with 20% of desktops consuming more than 300 kWh/yr and high-end computing desktops reaching 600 kWh/yr. While most structural engineers will be using an above average desktop in their modeling or a cloud computing service, this analysis opts to adopt the above estimate of 300 kWh/yr. As such, a typical computer is expected to generate  $0.42 \times 300 = 126$  Kg of CO<sub>2</sub> annually.

While the number of structural engineers in the world is not well established, Barter [83] notes that there are close to 50,000 structural engineers in the USA. Assuming that 2.5% of these engineers are proficient and dedicated ML users who use the above desktops, then applying the above calculation yields  $0.025 \times 50,000 \times 126 = 126,000$  Kg of CO<sub>2</sub> annually. This is equivalent to emissions generated from 27.4 passenger vehicles driven for one year (about 509,619.7 km) as per the published estimations made by the United States Environmental Protection Agency (EPA) [84].

Now, revisiting results obtained from Table 1 and Fig. 2 clearly shows that all selected models can be used to accurately predict if a concrete mixture can develop its design strength at in-situ conditions. However, a look at Fig. 5 infers the amount of possible energy savings that can be attained by using models of low energy consumption.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*. <https://doi.org/10.1016/j.jclepro.2022.135334>.

For illustration purposes, we take the highest energy consumption model (KDNN (1536)) in the case of the full dataset as a base for normalization, then a comparison between all models used in this particular case study can be established. This comparison is shown in Table 2 and Fig. 7 and infers that adopting LR, GBT (R), XGBoost, or LGBM can lead to 96%, 89%, 42%, and 27% reduction in CO<sub>2</sub> as opposed to KDNN (1536). This clearly shows the significant reductions that can be attained by adopting simple models. Table 2 and Fig. 7 draw a similar comparison to the cases of the reduced dataset and the reduced features.

Table 2 Comparison between ML models\*

| Model                   | <i>P</i> metric | <i>E</i> metric | Normalized <i>P</i> | Normalized <i>E</i> | % Reduction of max. CO <sub>2</sub> emitted (per case) |
|-------------------------|-----------------|-----------------|---------------------|---------------------|--|
| <b>All features</b>     |                 |                 |                     |                     |  |
| LR                      | 204.3           | 39.43           | 0.44                | 0.04                | 0.96   |
| EN                      | 77.2            | 71.92           | 0.17                | 0.08                | 0.92   |
| GBT (R)                 | 332.3           | 101.93          | 0.72                | 0.11                | 0.89   |
| GBT (Python)            | 325.9           | 131.27          | 0.70                | 0.14                | 0.86   |
| VW                      | 115.4           | 155.30          | 0.25                | 0.16                | 0.84   |
| KDNN (64)               | 40.0            | 183.70          | 0.09                | 0.20                | 0.80   |
| KDNN (64/64)            | 58.3            | 196.57          | 0.13                | 0.21                | 0.79   |
| RF (Python)             | 302.2           | 229.38          | 0.65                | 0.24                | 0.76   |
| KDNN (256/128/64)       | 57.7            | 302.52          | 0.12                | 0.32                | 0.68   |
| SVM                     | 306.3           | 504.81          | 0.66                | 0.54                | 0.46   |
| XGBoost                 | <u>464.5</u>    | 542.04          | 1.00                | 0.58                | 0.42   |
| RF (R)                  | 347.6           | 567.44          | 0.75                | 0.60                | 0.40   |
| TFDL                    | 234.5           | 579.54          | 0.50                | 0.62                | 0.38   |
| LGBM                    | 347.1           | 689.78          | 0.75                | 0.73                | 0.27   |
| KDNN (1536)             | 64.0            | <u>941.39</u>   | 0.14                | 1.00                | 0.00   |
| <b>Reduced dataset</b>  |                 |                 |                     |                     |  |
| LR                      | 227.9           | 14.93           | 0.75                | 0.01                | 0.99   |
| EN                      | 257.8           | 53.08           | 0.85                | 0.03                | 0.97   |
| GBT (R)                 | 73.5            | 60.31           | 0.24                | 0.04                | 0.96   |
| GBT (Python)            | 211.3           | 112.00          | 0.70                | 0.07                | 0.93   |
| VW                      | 20.8            | 171.90          | 0.07                | 0.11                | 0.89   |
| KDNN (64)               | 224.3           | 80.84           | 0.74                | 0.05                | 0.95   |
| KDNN (64/64)            | 38.8            | 144.58          | 0.13                | 0.09                | 0.91   |
| RF (Python)             | 144.3           | 55.79           | 0.48                | 0.04                | 0.96   |
| KDNN (256/128/64)       | 41.4            | 269.45          | 0.14                | 0.17                | 0.83   |
| SVM                     | 244.0           | 459.55          | 0.81                | 0.30                | 0.70   |
| XGBoost                 | 183.8           | 409.83          | 0.61                | 0.26                | 0.74   |
| RF (R)                  | <u>302.0</u>    | 255.71          | 1.00                | 0.16                | 0.84   |
| TFDL                    | 32.0            | <u>1557.37</u>  | 0.11                | 1.00                | 0.00   |
| LGBM                    | 130.1           | 1205.38         | 0.43                | 0.77                | 0.23   |
| KDNN (1536)             | 241.9           | 357.56          | 0.80                | 0.23                | 0.77   |
| <b>Reduced features</b> |                 |                 |                     |                     |  |
| LR                      | 181.93          | 20.31           | 0.55                | 0.01                | 0.99   |

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*. <https://doi.org/10.1016/j.jclepro.2022.135334>.

|                   |               |                |      |      |      |
|-------------------|---------------|----------------|------|------|------|
| GBT (R)           | 314.03        | 46.85          | 0.94 | 0.02 | 0.98 |
| EN                | 69.85         | 51.11          | 0.21 | 0.03 | 0.97 |
| RF (Python)       | 319.65        | 59.90          | 0.96 | 0.03 | 0.97 |
| KDNN (64/64)      | 64.12         | 87.18          | 0.19 | 0.04 | 0.96 |
| GBT (Python)      | 279.91        | 89.70          | 0.84 | 0.05 | 0.95 |
| KDNN (64)         | 35.50         | 92.34          | 0.11 | 0.05 | 0.95 |
| VW                | 160.62        | 92.82          | 0.48 | 0.05 | 0.95 |
| KDNN (256/128/64) | 263.13        | 158.66         | 0.79 | 0.08 | 0.92 |
| RF (R)            | 314.11        | 243.65         | 0.94 | 0.12 | 0.88 |
| SVM               | 224.93        | 390.11         | 0.68 | 0.20 | 0.80 |
| XGBoost           | <u>332.93</u> | 427.87         | 1.00 | 0.22 | 0.78 |
| KDNN (1536)       | 114.73        | 627.19         | 0.34 | 0.32 | 0.68 |
| TFDL              | 96.02         | 1472.82        | 0.29 | 0.74 | 0.26 |
| LGBM              | 263.48        | <u>1986.42</u> | 0.79 | 1.00 | 0.00 |

\*Values used for normalization in each particular case are underlined and in *italic*.

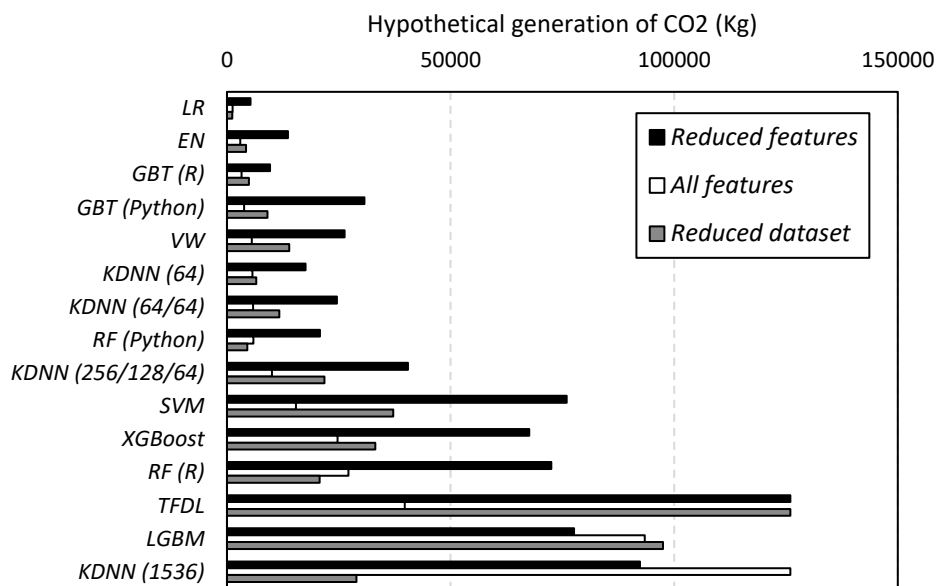


Fig. 7 Estimated carbon generated per ML model

### 3.4 Insights into Green ML

The presented analysis in this study highlights the need for advocating the adoption of responsible and Green ML early into fully integrating ML into the structural engineering domain. As such, one can be ready and prepared to maximize the benefits of ML while negating some of the complications that might arise (e.g., irresponsible energy consumption), which may delay harnessing the full potential of ML in the years to come.

The presented analysis can also be thought of as a first step towards a more articulated diagnostic investigation that dives in-depth to account for other dimensions that were not explored herein;

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

especially those related to software/hardware, storage, and upkeep of models, data transfer, data type and size, type of observations, geographical location of working stations, resources used in generating electricity, governance and policy aspects, etc. In a future analysis, further details on the number of structural engineers (and the ratio of those who are practicing ML in their works), preferred working stations/cloud services, and use of centralized vs. distributed computing, to name a few, are to be investigated. The open literature does contain a few works tackling the accuracy-energy trade-off front by exploring the influence of much more complex ML models (i.e., deep learning and natural language processing), a similar series of investigations on such models derivatives that can be deployed in structural engineering problems are of value to explore [85–87]. One particular work of interest is that carried out by García-Martín et al. [78]<sup>7</sup>, wherein this group presents options for hardware and software energy prediction methods – some of which can be used by structural engineers. These researchers also emphasized the need to investigate ML model energy consumption during the training phase and the inference phase.

In addition, a comparison between the energy consumption of traditional methods (such as FE modeling) to that of ML, when both are applied to the same problem, can indeed be interesting to investigate. While both methods can potentially use similar computing stations, it would be interesting to examine differences in setting up such stations to optimize their performance. Traditionally, FE models not only require large computational capacities but also necessitate the presence of large storage units, which can be a fraction of that required by a simple ML model for the same problem – especially if the FE model utilized finer mesh and nonlinear effects.

For example, a three-dimensional nonlinear FE model to predict the deformation history of reinforced concrete beams under fire conditions was developed by the author in earlier work. This model was composed of 50,125 elements and of a size of 64.58 MB [90]. The results from running this model for a typical beam yielded a file of a size of 700 MB. In a separate work [91], a ML model was also developed to predict the same phenomenon. This ML model was 305 KB. While both models were developed using different working stations approximately ten years apart, the truth is that both models still require storage. The American Council for an Energy-Efficient Economy estimates that it takes 5.12 kWh of electricity per gigabyte of transferred data. This implies that one GB (i.e., 1000 MB) emits  $0.42 \text{ Kg} \times 5.12 \text{ kWh} = 2.15 \text{ Kg}$  of CO<sub>2</sub>. Thus, carbon emissions arising solely from *storing* the FE model (without the result file), excluding those arising from model development and simulation time, are 215 times higher than that of the ML model.

Based on the presented analysis, the following are some of the big ideas that can be followed to create GML models:

- Care should be taken when developing new ML models. In all cases, engineers should consider the adequacy of simple models to tackle a problem instead of blindly selecting complex models. The same consideration can be applied to exploring issues arising from the accuracy-energy trade-off.

---

<sup>7</sup> Other works [88,89], have made solid advancements on this front and are worthy of review.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

- Engineers are to consider how accurate and detailed their developed model should be. In a similar analogy to FE modeling, the adoption of a very detailed model involves a significant increase in the resources allocated to developing such a model, which yields higher monetary and environmental costs.
- Establish and aim for a “functional” level of model performance. Chasing accuracy metrics does not guarantee optimal performance, as such metrics reflect upon the model’s performance against the available data and may not negate falling into model behavioral issues (e.g., overfitting, biasness, etc.).
- Explore pre-trained models and transfer learning techniques as a means to avoid the unnecessary retraining of models.
- Whenever possible, ML users are expected to share details of their conducted analysis and dataset and disclose energy performance metrics.
- Consider participating in carbon emission reduction initiatives by opting for computing services that employ environmental-friendly facilities.

#### 4. Conclusions

This paper carries out a series of comparisons between simple and exotic ML algorithms to report on their predictive performance, model size, energy consumption, and storage need. More specifically, this work examines the influence of algorithm architecture, processing language, number of features, as well as dataset size on model predictivity, energy consumption, and expected carbon emissions. The following list of inferences can be drawn from the findings of this study:

- Structural engineers are expected to be cognizant of the accuracy-energy trade-offs when developing future ML models.
- Between 23-99% reduction in energy consumption and carbon emissions (while maintaining a similar level of accuracy) can be attained by adopting simple ML models and reduced-order modeling strategies.
- Of all examined ML models, the XGBoost ranked highest in terms of predictivity, and the LR ranked lowest in terms of energy consumption (and carbon emittance).
- There is an apparent influence on accuracy and energy consumption in similar ML models developed using different programming languages and ML models of the same origin but with different architectures. As such, care is needed when selecting such algorithms.
- Additional experiments specifically designed to examine the accuracy-energy trade-off in terms of system architecture, computing expenditure, etc., are warranted.
- Proactively setting the stage towards responsible and green ML is of merit to this community, and hence efforts in this direction are to be invited and appreciated.

#### Data Availability

Some or all data, models, or code that support the findings of this study are available from the corresponding author upon reasonable request.

#### Conflict of Interest

The author declares no conflict of interest.



Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

## References

- [1] M. Mahendran, Applications of Finite Element Analysis in Structural Engineering, Int. Conf. Comput. Aided Eng. (2015).
- [2] C.R. Farrar, K. Worden, Structural Health Monitoring: A Machine Learning Perspective, 2012. <https://doi.org/10.1002/9781118443118>.
- [3] B. D'Amico, R.J. Myers, J. Sykes, E. Voss, B. Cousins-Jenvey, W. Fawcett, S. Richardson, A. Kermani, F. Pomponi, Machine Learning for Sustainable Structures: A Call for Data, Structures. (2019). <https://doi.org/10.1016/j.istruc.2018.11.013>.
- [4] M.H. Rafiei, H. Adeli, A novel machine learning-based algorithm to detect damage in high-rise building structures, *Struct. Des. Tall Spec. Build.* (2017). <https://doi.org/10.1002/tal.1400>.
- [5] M.Z. Naser, Machine learning assessment of fiber-reinforced polymer-strengthened and reinforced concrete members, *ACI Struct. J.* (2020). <https://doi.org/10.14359/51728073>.
- [6] S.K. Babanajad, A.H. Gandomi, A.H. Alavi, New prediction models for concrete ultimate strength under true-triaxial stress states: An evolutionary approach, *Adv. Eng. Softw.* (2017). <https://doi.org/10.1016/j.advengsoft.2017.03.011>.
- [7] M. Naser, A Faculty's Perspective into Infusing Artificial Intelligence to Civil Engineering Education, *J. Civ. Eng. Educ.* (2022). [https://doi.org/10.1061/\(ASCE\)EI.2643-9115.0000065](https://doi.org/10.1061/(ASCE)EI.2643-9115.0000065).
- [8] Y. Wang, S. Sun, X. Chen, X. Zeng, Y. Kong, J. Chen, Y. Guo, T. Wang, Short-term load forecasting of industrial customers based on SVM and XGBoost, *Int. J. Electr. Power Energy Syst.* 129 (2021) 106830. <https://doi.org/10.1016/j.ijepes.2021.106830>.
- [9] Y. Xie, M. Ebad Sichani, J.E. Padgett, R. DesRoches, The promise of implementing machine learning in earthquake engineering: A state-of-the-art review, *Earthq. Spectra.* (2020). <https://doi.org/10.1177/8755293020919419>.
- [10] A. Tapeh, M.Z. Naser, Artificial Intelligence, Machine Learning, and Deep Learning in Structural Engineering: A Scientometrics Review of Trends and Best Practices, *Arch. Comput. Methods Eng.* (2022). <https://doi.org/10.1007/s11831-022-09793-w>.
- [11] A. Ashrafian, F. Shokri, M.J. Taheri Amiri, Z.M. Yaseen, M. Rezaie-Balf, Compressive strength of Foamed Cellular Lightweight Concrete simulation: New development of hybrid artificial intelligence model, *Constr. Build. Mater.* (2020). <https://doi.org/10.1016/j.conbuildmat.2019.117048>.
- [12] J. Gola, J. Webel, D. Britz, A. Guitar, T. Staudt, M. Winter, F. Mücklich, Objective microstructure classification by support vector machine (SVM) using a combination of morphological parameters and textural features for low carbon steels, *Comput. Mater. Sci.*

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.

<https://doi.org/10.1016/j.jclepro.2022.135334>.

- (2019). <https://doi.org/10.1016/j.commatsci.2019.01.006>.
- [13] A. Diez, N.L.D. Khoa, M. Makki Alamdari, Y. Wang, F. Chen, P. Runcie, A clustering approach for structural health monitoring on bridges, *J. Civ. Struct. Heal. Monit.* 6 (2016) 429–445. <https://doi.org/10.1007/s13349-016-0160-0>.
- [14] H. Hasni, A.H. Alavi, N. Lajnef, M. Abdelbarr, S.F. Masri, S. Chakrabartty, Self-powered piezo-floating-gate sensors for health monitoring of steel plates, *Eng. Struct.* (2017). <https://doi.org/10.1016/j.engstruct.2017.06.063>.
- [15] A. Ashteyat, Y.T. Obaidat, Y.Z. Murad, R. Haddad, Compressive strength prediction of lightweight short columns at elevated temperature using gene expression programming and artificial neural network, *J. Civ. Eng. Manag.* (2020). <https://doi.org/10.3846/jcem.2020.11931>.
- [16] Y. Panev, P. Kotsovinos, S. Deeny, G. Flint, The Use of Machine Learning for the Prediction of fire Resistance of Composite Shallow Floor Systems, *Fire Technol.* (2021). <https://doi.org/10.1007/s10694-021-01108-y>.
- [17] A.N. Tarawneh, A.N. Tarawneh, A.N. Tarawneh, B.E. Ross, B.E. Ross, B.E. Ross, T.E. Cousins, T.E. Cousins, T.E. Cousins, Shear Behavior and Design of Post-Installed Anchors in Thin Concrete Members, *ACI Struct. J.* (2020). <https://doi.org/10.14359/51723508>.
- [18] A.A.H. Alwanas, A.A. Al-Musawi, S.Q. Salih, H. Tao, M. Ali, Z.M. Yaseen, Load-carrying capacity and mode failure simulation of beam-column joint connection: Application of self-tuning machine learning model, *Eng. Struct.* (2019). <https://doi.org/10.1016/j.engstruct.2019.05.048>.
- [19] F. Fu, Fire induced progressive collapse potential assessment of steel framed buildings using machine learning, *J. Constr. Steel Res.* (2020). <https://doi.org/10.1016/j.jcsr.2019.105918>.
- [20] W. Liao, Y. Huang, Z. Zheng, X. Lu, Intelligent generative structural design method for shear wall building based on “fused-text-image-to-image” generative adversarial networks, *Expert Syst. Appl.* 210 (2022) 118530. <https://doi.org/10.1016/J.ESWA.2022.118530>.
- [21] M. Freischlad, M. Schnellenbach-Held, A machine learning approach for the support of preliminary structural design, *Adv. Eng. Informatics.* (2005). <https://doi.org/10.1016/j.aei.2005.07.001>.
- [22] J.C.P. Cheng, W. Chen, K. Chen, Q. Wang, Data-driven predictive maintenance planning framework for MEP components based on BIM and IoT using machine learning algorithms, *Autom. Constr.* (2020). <https://doi.org/10.1016/j.autcon.2020.103087>.
- [23] O. Avci, O. Abdeljaber, S. Kiranyaz, M. Hussein, M. Gabbouj, D.J. Inman, A review of

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

- vibration-based damage detection in civil structures: From traditional methods to Machine Learning and Deep Learning applications, *Mech. Syst. Signal Process.* (2021).  
<https://doi.org/10.1016/j.ymsp.2020.107077>.
- [24] M.Z. Naser, Mechanistically Informed Machine Learning and Artificial Intelligence in Fire Engineering and Sciences, *Fire Technol.* (2021) 1–44.  
<https://doi.org/10.1007/s10694-020-01069-8>.
- [25] H. Adeli, Neural networks in civil engineering: 1989-2000, *Comput. Civ. Infrastruct. Eng.* (2001). <https://doi.org/10.1111/0885-9507.00219>.
- [26] H. Sun, H. V. Burton, H. Huang, Machine learning applications for building structural design and performance assessment: State-of-the-art review, *J. Build. Eng.* 33 (2021) 101816. <https://doi.org/10.1016/j.jobe.2020.101816>.
- [27] A. Çevik, A.E. KURTOĞLU, M. Bilgehan, M.E. Gülşan, H.M. Albegmpri, Support vector machines in structural engineering: A review, 2015.  
<https://doi.org/10.3846/13923730.2015.1005021>.
- [28] E.M. Golafshani, A. Behnood, Estimating the optimal mix design of silica fume concrete using biogeography-based programming, *Cem. Concr. Compos.* 96 (2019) 95–105.  
<https://doi.org/10.1016/J.CEMCONCOMP.2018.11.005>.
- [29] V. V. Degtyarev, Neural networks for predicting shear strength of CFS channels with slotted webs, *J. Constr. Steel Res.* (2021). <https://doi.org/10.1016/j.jcsr.2020.106443>.
- [30] S. Mangalathu, H. V. Burton, Deep learning-based classification of earthquake-impacted buildings using textual damage descriptions, *Int. J. Disaster Risk Reduct.* 36 (2019) 101111. <https://doi.org/10.1016/j.ijdr.2019.101111>.
- [31] Y. Pan, L. Zhang, Roles of artificial intelligence in construction engineering and management: A critical review and future trends, *Autom. Constr.* (2021).  
<https://doi.org/10.1016/j.autcon.2020.103517>.
- [32] M.Z. Naser, The Role of Computational Intelligence in Realizing Modern and Autonomous Fire Evaluation Methods, in: *Struct. Congr. 2020 - Sel. Pap. from Struct. Congr. 2020*, 2020. <https://doi.org/10.1061/9780784482896.059>.
- [33] M. Mohammadi, A. Al-Fuqaha, Enabling Cognitive Smart Cities Using Big Data and Machine Learning: Approaches and Challenges, *IEEE Commun. Mag.* (2018).  
<https://doi.org/10.1109/MCOM.2018.1700298>.
- [34] E. Strubell, A. Ganesh, A. McCallum, Energy and Policy Considerations for Modern Deep Learning Research, *Proc. AAAI Conf. Artif. Intell.* (2020).  
<https://doi.org/10.1609/aaai.v34i09.7123>.
- [35] R. Jackson, A roadmap to reducing greenhouse gas emissions 50 percent by 2030, 2019.  
<https://earth.stanford.edu/news/roadmap-reducing-greenhouse-gas-emissions-50-percent->

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.jclepro.2022.135334>.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

2030#gs.zij3zf.

- [36] G. SMARTer, GeSI SMARTer 2020: The Role of ICT in Driving a Sustainable Future GeSI SMARTer 2020 2 The Role of ICT in Driving a Sustainable Future Independent review by, 2020.
- [37] J. Farfan, M. Fasihi, C. Breyer, Trends in the global cement industry and opportunities for long-term sustainable CCU potential for Power-to-X, *J. Clean. Prod.* (2019).  
<https://doi.org/10.1016/j.jclepro.2019.01.226>.
- [38] H. Ritchie, M. Roser, Emissions by sector - Our World in Data, 2021.  
<https://ourworldindata.org/emissions-by-sector> (accessed April 28, 2021).
- [39] Y. Lu, A. Zhang, Q. Li, B. Dong, Beyond Finite Layer Neural Networks: Bridging Deep Architectures and Numerical Differential Equations | OpenReview, ICLR 2018 Conf. (2018). <https://openreview.net/forum?id=ryZ283gAZ> (accessed April 23, 2021).
- [40] F.N. Iandola, S. Han, M.W. Moskewicz, K. Ashraf, W.J. Dally, K. Keutzer, SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5MB model size, (2016).  
<https://doi.org/10.48550/arxiv.1602.07360>.
- [41] E.R. Ziegel, The Elements of Statistical Learning, *Technometrics*. (2003).  
<https://doi.org/10.1198/tech.2003.s770>.
- [42] N. Ketkar, N. Ketkar, Introduction to Keras, in: *Deep Learn. with Python*, 2017.  
[https://doi.org/10.1007/978-1-4842-2766-4\\_7](https://doi.org/10.1007/978-1-4842-2766-4_7).
- [43] T. Chen, C. Guestrin, XGBoost: A scalable tree boosting system, in: *Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Min.*, 2016. <https://doi.org/10.1145/2939672.2939785>.
- [44] Scikit, *sklearn.ensemble.GradientBoostingRegressor* — *scikit-learn 0.24.1 documentation*, (2020). <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingRegressor.html> (accessed February 9, 2021).
- [45] XGBoost Python Package, *Python Package Introduction* — *xgboost 1.4.0-SNAPSHOT documentation*, (2020).  
[https://xgboost.readthedocs.io/en/latest/python/python\\_intro.html#early-stopping](https://xgboost.readthedocs.io/en/latest/python/python_intro.html#early-stopping) (accessed February 10, 2021).
- [46] Scikit, *sklearn.ensemble.GradientBoostingClassifier* — *scikit-learn 0.24.1 documentation*, (n.d.). <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingClassifier.html> (accessed April 26, 2021).
- [47] A. Natekin, A. Knoll, Gradient boosting machines, a tutorial, *Front. Neurobot.* 7 (2013).  
<https://doi.org/10.3389/fnbot.2013.00021>.

This is a preprint draft. The published article can be found at: <https://doi.org/10.1016/j.jclepro.2022.135334>.

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

- [48] Keras, GitHub - keras-team/keras: Deep Learning for humans, (2020).  
<https://github.com/keras-team/keras> (accessed February 9, 2021).
- [49] TensorFlow, GitHub - tensorflow/tensorflow: An Open Source Machine Learning Framework for Everyone, (2020). <https://github.com/tensorflow/tensorflow> (accessed February 9, 2021).
- [50] H.T. Thai, Machine learning for structural engineering: A state-of-the-art review, Structures. (2022). <https://doi.org/10.1016/j.istruc.2022.02.003>.
- [51] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, Q. Ye, T.Y. Liu, LightGBM: A highly efficient gradient boosting decision tree, in: Adv. Neural Inf. Process. Syst., 2017.
- [52] LightGBM, Welcome to LightGBM's documentation! — LightGBM 3.1.1.99 documentation, (2020). <https://lightgbm.readthedocs.io/en/latest/> (accessed February 9, 2021).
- [53] H. Li, Z. Xu, G. Taylor, C. Studer, T. Goldstein, Visualizing the loss landscape of neural nets, in: Adv. Neural Inf. Process. Syst., 2018.
- [54] M. Abadi, TensorFlow: learning functions at scale, ACM SIGPLAN Not. (2016).  
<https://doi.org/10.1145/3022670.2976746>.
- [55] VowpalWabbit, vowpalwabbit.sklearn — VowpalWabbit 8.11.0 documentation, (2021).  
[https://vowpalwabbit.org/docs/vowpal\\_wabbit/python/latest/vowpalwabbit.sklearn.html](https://vowpalwabbit.org/docs/vowpal_wabbit/python/latest/vowpalwabbit.sklearn.html) (accessed April 26, 2021).
- [56] B.E. Boser, I.M. Guyon, V.N. Vapnik, Training algorithm for optimal margin classifiers, in: Proc. Fifth Annu. ACM Work. Comput. Learn. Theory, 1992.  
<https://doi.org/10.1145/130385.130401>.
- [57] Scikit, sklearn.kernel\_approximation.Nystroem — scikit-learn 0.24.1 documentation, (2021). [https://scikit-learn.org/stable/modules/generated/sklearn.kernel\\_approximation.Nystroem.html](https://scikit-learn.org/stable/modules/generated/sklearn.kernel_approximation.Nystroem.html) (accessed April 27, 2021).
- [58] Scikit, lightning.classification.CDClassifier — lightning dev documentation, (2021).  
<http://contrib.scikit-learn.org/lightning/generated/lightning.classification.CDClassifier.html> (accessed April 27, 2021).
- [59] A. Liaw, M. Wiener, Classification and Regression by RandomForest, 2002.  
<https://www.researchgate.net/publication/228451484> (accessed April 8, 2019).
- [60] Scikit, sklearn.ensemble.RandomForestClassifier — scikit-learn 0.24.1 documentation, (2020). <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html> (accessed February 9, 2021).

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

- [61] L. Breiman, randomForest: Breiman and Cutler's Random Forests for Classification and Regression, *Mach. Learn.* 45 (2001) 5–32. <https://doi.org/10.1023/A:1010933404324>.
- [62] B.A. Young, A. Hall, L. Pilon, P. Gupta, G. Sant, Can the compressive strength of concrete be estimated from knowledge of the mixture proportions?: New insights from statistical analysis and machine learning methods, *Cem. Concr. Res.* 115 (2019) 379–388. <https://doi.org/10.1016/j.cemconres.2018.09.006>.
- [63] A. Yu, hyperparameter\_tune/Clean\_data.csv at master · huiziy/hyperparameter\_tune · GitHub, (2021).  
[https://github.com/huiziy/hyperparameter\\_tune/blob/master/Clean\\_data.csv](https://github.com/huiziy/hyperparameter_tune/blob/master/Clean_data.csv) (accessed March 18, 2021).
- [64] M.Z. Naser, Digital twin for next gen concretes: On-demand tuning of vulnerable mixtures through Explainable and Anomalous Machine Learning, *Cem. Concr. Compos.* 132 (2022) 104640. <https://doi.org/10.1016/J.CEMCONCOMP.2022.104640>.
- [65] S. Marsland, *Machine learning: An algorithmic perspective*, 2014. <https://doi.org/10.1201/b17476>.
- [66] D.C. Luor, A comparative assessment of data standardization on support vector machine for classification problems, *Intell. Data Anal.* (2015). <https://doi.org/10.3233/IDA-150730>.
- [67] A.U. Abubakar, M.S. Tabra, Prediction of Compressive Strength in High Performance Concrete with Hooked-End Steel Fiber using K-Nearest Neighbor Algorithm, *Int. J. Integr. Eng.* (2019). <https://doi.org/10.30880/ijie.2019.11.01.016>.
- [68] R. Biswas, B. Rai, P. Samui, S.S. Roy, Estimating concrete compressive strength using MARS, LSSVM and GP, *Eng. J.* (2020). <https://doi.org/10.4186/ej.2020.24.2.41>.
- [69] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, É. Duchesnay, E. Duchesnay, *Scikit-learn: Machine learning in Python*, *J. Mach. Learn. Res.* 12 (2011) 2825–2830.
- [70] M.D. Schmidt, H. Lipson, Age-fitness pareto optimization, in: 2010. <https://doi.org/10.1145/1830483.1830584>.
- [71] M. Laszczyk, P.B. Myszkowski, Survey of quality measures for multi-objective optimization: Construction of complementary set of multi-objective quality measures, *Swarm Evol. Comput.* 48 (2019) 109–133. <https://doi.org/10.1016/J.SWEVO.2019.04.001>.
- [72] A.H. Alavi, A.H. Gandomi, M.G. Sahab, M. Gandomi, Multi expression programming: A new approach to formulation of soil classification, *Eng. Comput.* 26 (2010) 111–118. <https://doi.org/10.1007/s00366-009-0140-7>.
- [73] M.Z.Z. Naser, A. Seitllari, Concrete under fire: an assessment through intelligent pattern

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

- recognition, *Eng. Comput.* 36 1–14. <https://doi.org/10.1007/s00366-019-00805-1>.
- [74] W.Z. Taffese, E. Sistonen, Machine learning for durability and service-life assessment of reinforced concrete structures: Recent advances and future directions, *Autom. Constr.* (2017). <https://doi.org/10.1016/j.autcon.2017.01.016>.
- [75] H. Huang, H. V. Burton, Classification of in-plane failure modes for reinforced concrete frames with infills using machine learning, *J. Build. Eng.* (2019).  
<https://doi.org/10.1016/j.jobe.2019.100767>.
- [76] S.M. Lundberg, S.I. Lee, A unified approach to interpreting model predictions, in: *Adv. Neural Inf. Process. Syst.*, 2017.
- [77] L.J. Ba, R. Caruana, Do deep nets really need to be deep?, in: *Adv. Neural Inf. Process. Syst.*, 2014.
- [78] E. García-Martín, C.F. Rodrigues, G. Riley, H. Grahn, Estimation of energy consumption in machine learning, *J. Parallel Distrib. Comput.* (2019).  
<https://doi.org/10.1016/j.jpdc.2019.07.007>.
- [79] J. Frankle, M. Carbin, The lottery ticket hypothesis: Finding sparse, trainable neural networks, *ArXiv*. (2018).
- [80] A.E.I. Brownlee, J. Adair, S.O. Haraldsson, J. Jabbo, Exploring the Accuracy-Energy Trade-off in Machine Learning, n.d. <https://www.kaggle.com/dinu1763/mortgage-loan-approval> (accessed April 28, 2021).
- [81] EIA, Frequently Asked Questions (FAQs) - U.S. Energy Information Administration (EIA), U.S. Energy Inf. Adm. (2021). <https://www.eia.gov/tools/faqs/faq.php?id=74&t=11> (accessed April 28, 2021).
- [82] L. Desroches, H. Fuch, J. Greenblatt, S. Pratt, Hcl. Willem, B. Beraki, M. Nagaraju, S. Price, S. Young, Computer usage and national energy consumption: Results from a field-metering study, 2014. <https://www.osti.gov/servlets/purl/1166988> (accessed April 28, 2021).
- [83] M. Barter, 10 Obstacles to Meaningful Licensing of Structural Engineers, *Struct. Mag.* (2014). <https://www.structuremag.org/?p=4039> (accessed April 28, 2021).
- [84] EPA, Greenhouse Gas Equivalencies Calculator, EPA. (2021).  
<https://www.epa.gov/energy/greenhouse-gas-equivalencies-calculator> (accessed April 28, 2021).
- [85] N.K. Jayakodi, A. Chatterjee, W. Choi, J.R. Doppa, P.P. Pande, Trading-off accuracy and energy of deep inference on embedded systems: A co-design approach, *IEEE Trans. Comput. Des. Integr. Circuits Syst.* (2018). <https://doi.org/10.1109/TCAD.2018.2857338>.
- [86] T.J. Yang, Y.H. Chen, J. Emer, V. Sze, A method to estimate the energy consumption of

Please cite this paper as:

Naser, M.Z. (2023). Do We Need Exotic Models? Engineering Metrics to Enable Green Machine Learning from Tackling Accuracy-Energy Trade-offs. *Journal of Cleaner Production*.  
<https://doi.org/10.1016/j.jclepro.2022.135334>.

- deep neural networks, in: Conf. Rec. 51st Asilomar Conf. Signals, Syst. Comput. ACSSC 2017, 2018. <https://doi.org/10.1109/ACSSC.2017.8335698>.
- [87] N.K. Jayakodi, S. Belakaria, A. Deshwal, J.R. Doppa, Design and optimization of energy-accuracy tradeoff networks for mobile platforms via pretrained deep models, *ACM Trans. Embed. Comput. Syst.* (2020). <https://doi.org/10.1145/3366636>.
- [88] A. Lacoste, A. Luccioni, V. Schmidt, T. Dandres, Quantifying the Carbon Emissions of Machine Learning, (2019). <https://doi.org/10.48550/arxiv.1910.09700>.
- [89] P. Henderson, J. Hu, J. Romoff, E. Brunskill, D. Jurafsky, J. Pineau, Towards the systematic reporting of the energy and carbon footprints of machine learning, *J. Mach. Learn. Res.* (2020).
- [90] R.A. Hawileh, M.Z. Naser, W. Zaidan, H.A. Rasheed, Modeling of insulated CFRP-strengthened reinforced concrete T-beam exposed to fire, *Eng. Struct.* 31 (2009) 3072–3079. <https://doi.org/10.1016/j.engstruct.2009.08.008>.
- [91] M.Z. Naser, AI-based cognitive framework for evaluating response of concrete structures in extreme conditions, *Eng. Appl. Artif. Intell.* 81 (2019) 437–449.  
<https://www.sciencedirect.com/science/article/pii/S0952197619300466> (accessed April 1, 2019).