#### **REVIEW PAPER**



# Error Metrics and Performance Fitness Indicators for Artificial Intelligence and Machine Learning in Engineering and Sciences

M. Z. Naser<sup>1,2</sup> · Amir H. Alavi<sup>3</sup>

Received: 29 July 2021 / Accepted: 15 November 2021 © The Author(s), under exclusive licence to Springer Nature Switzerland AG 2021

#### **Abstract**

Artificial intelligence (AI) and Machine learning (ML) train machines to achieve a high level of cognition and perform human-like analysis. Both AI and ML seemingly fit into our daily lives as well as complex and interdisciplinary fields. With the rise of commercial, open-source, and user-catered AI/ML tools, a key question often arises whenever AI/ML is applied to explore a phenomenon or a scenario: what constitutes a good AI/ML model? Keeping in mind that a proper answer to this question depends on various factors, this work presumes that a goodmodel optimally performs and best describes the phenomenon on hand. From this perspective, identifying proper assessment metrics to evaluate the performance of AI/ML models is not only necessary but is also warranted. As such, this paper examines 78 of the most commonly-used performance fitness and error metrics for regression and classification algorithms, with emphasis on engineering and sciences applications.

**Keywords** Error metrics · Machine learning · Regression · Classification

#### Introduction

Learning is the process of seeking knowledge [1]. We, as humans, can learn from our daily interactions and experiences because we have the ability to communicate, reason, and understand. With the rapid technological advancement in computer sciences, computational intelligence has led to the development of modern cognitive and evaluation tools [2, 3]. One such tool is machine learning (ML) which is often described as a set of methods that, when applied, can allow machines to learn/understand meaningful patterns from data repositories; while maintaining minimal human interaction [4]. More specifically, a "computer program is said to learn from experience E with respect to some class"

of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E" [5]. In other words, ML trains machines to understand real-world applications, use this knowledge to carry out pre-identified tasks with the goal of optimizing and improving the machines' performance with time and new knowledge. A closer look at the definition of ML infers that computers do not learn by reasoning but rather by algorithms.

From the perspective of this work, traditional statistical regression techniques are often used to carry out behavioral modeling wherein such techniques may suffer from large uncertainties, the need for the idealization of complex processes, approximation, and averaging widely varying prototype conditions. Furthermore, statistical analysis often assumes linear, or in some cases nonlinear, relationships between the output and the predictor variables, and these assumptions do not always hold true - especially in the context of engineering/real data. On the other hand, ML methods adaptively learn from experiences and extract various discriminators. One of the major advantages of ML approaches over the traditional statistical techniques is their ability to derive a relationship(s) between inputs and outputs without assuming prior forms or existing relationships. In other words, ML approaches are not confined to one particular space that requires the availability of physical

M. Z. Naser
mznaser@clemson.edu; m@mznaser.com

Amir H. Alavi alavi@pitt.edu

- School of Civil and Environmental Engineering and Earth Sciences, Clemson University, Clemson, SC 29634, USA
- Artificial Intelligence Research Institute for Science and Engineering (AIRISE) at Clemson University, Clemson, SC 29634, USA
- Department of Civil and Environmental Engineering, University of Pittsburgh, Pittsburgh, PA 15261, USA

Published online: 24 November 2021



representation but rather goes beyond that to explore hidden relations in data patterns [6-11].

While ML was initially developed for computer sciences, it is now an integral part of various fields including, energy/mechanical engineering [6–9], social sciences [10, 11], space applications [12, 13], among others [14–19]. Due to the availability of high-computationally powered machines and ease-of-access to data (thanks in part to the rise of Internet-of-Things and data-driven-applications), the utilization of ML into civil engineering, in general, and materials science, engineering in particular, has been duly noted in recent years [20–25].

An integral part of the wide spread of integrating ML into new research areas is due to the availability of userfriendly and easy-to-use software packages that simplifies the process of ML by utilizing pre-defined algorithms and training/validation procedure [26–30]. The availability of such tools, while facilitating ML analysis and providing new opportunities for researchers often unfamiliar with the ML fundamentals with means to easily carry out such analysis, could still be misused by providing a false sense of analysis interpretation [31]. Another concern of utilizing user-ready approaches to carry out ML analysis lies in the need for compiling proper observations (i.e. datapoints). In some classical fields (say material sciences, earthquake or fire engineering) where there is a limited number of observations due to expensive tests, or need for specialized instrumentation/facilities [32], then the use of ML may lead to a biased outcome - especially when combined with lack of expertise on ML [33, 34].

An examination of open literature raises a few questions: 1) are we developing accurate ML models? 2) are such models useful to our fields? 3) are we properly validating ML models? And 4) how to confidently answer "yes" to the aforementioned questions?

A distinction should be drawn in which we need to acknowledge that, we often apply existing ML algorithms to our problems rather than developing new algorithms. This acknowledgment goes hand in hand with that similar to applying other numerical tools such as the finite element method, to investigate the response of materials and structures (say concrete beams) under harsh environments (i.e. fire conditions) [35, 36]. From this perspective, we use an existing tool, say a finite element (FE) software (ANSYS [37], ABAQUS [38] etc.), to investigate how failure mechanism occurs in a concrete beam under fire. The accuracy of this FE model is often established through a validation procedure in which a comparison of predictions from the FE model (say temperature rise in steel rebars or mid-span deflection during a fire, or in some cases, point in time when the beam fails) is plotted against that measured in an actual fire test. If the comparison is deemed well, then the FE model is said to be valid and hence can be used to explore the effect of key response parameters (i.e. magnitude of loading, strength of concrete, intensity of fire etc.). From this perspective, the validity of an FE model is established if the variation between predicted results and measured observations is between 5–15%<sup>1</sup> [39].

Unlike the use of FE simulation, ML is often used in two domains: 1) to show the applicability of ML to understand a phenomenon [40, 41], and 2) to identify hidden patterns governing a phenomenon [33, 42]. In the first domain, ML is primarily used to show that an ML algorithm can replicate a phenomenon – or in other words, to validate the applicability of that particular ML algorithm to a material science problem (i.e. can deep learning be applied to predict the compressive strength of concrete given that information regarding the components in a concrete mix is available?). While works in this domain showcase the diversity of ML, these also provide an additional validation platform/case studies to already well-established algorithms. The contribution of such works to our knowledge base is to be thanked and acknowledged.

The second domain is where ML shines and can be proven as a powerful ally to researchers. This is because ML strives on data and is designed to explore hidden features and patterns. The integration of these two items has not been thoroughly applied into our fields and, if applied properly, cannot only open new opportunities but also revolutionize our perspective into our fields. Unfortunately, the open literature continues to lack works in this domain, and hence such works are to be encouraged.

Whether ML is used in the first or second domain, ML models need to be rigorously assessed [43, 44]. This is a critical key to ensure: 1) the validity of the developed ML model in understanding a complex phenomenon given a limited set of data points, and 2) proper extension of the same models towards new/future datasets. Traditionally, the adequacy of ML models is often established through performance fitness and error metrics (PFEMs). Performance and error measures are vital elements in the process of evaluating ML models/frameworks. These are defined as logical and/ or mathematical constructs intended to measure the closeness of actual observations to that expected (or predicted). In other words, PFEMs are used to establish an understanding of how predictions from a model compare to real (or measured) observations. Such metrics often relate to the variation between predicted and measured observations in terms of errors [45–47].

Diverse sets of performance metrics have been noted in the open literature i.e. correlation coefficient (R), root mean



<sup>&</sup>lt;sup>1</sup> One should note that the validation of an FE model is also governed by satisfying convergence criteria input in the FE software. More on this can be found elsewhere [37, 38].

squared error (RMSE), etc. In practice, one, a multiple, or a combination of metrics are used to examine the adequacy of a particular ML model. However, there does not seem to be a systematic view into which scenarios specific metrics are preferable to use. In order to bridge this knowledge gap, this work compiles the commonly-used PFEMs and highlights their use in evaluating the performance of regression and classification ML models.

## **Performance Fitness and Error Metrics**

This section presents the most widely-used PFEMS and highlights fundamentals, recommendations, and limitations associated with their use in assessing ML models.<sup>2</sup> In this work, PFEMs are grouped under two categories; traditional and modern. In this section, these reoccurring terms are used; *A*: actual measurements, *P*: predictions, *n*: number of data points.

# Regression

Regression ML methods deal with predicting a target value using independent variables. Some of these methods include artificial neural networks, genetic programing, etc. PFEMs grouped herein belong to a group of metrics that are based on methods to calculate point distance primarily using subtraction or division operations. These metrics contain fundamental operations, either A-P or P/A, and can be supplemented with absoluteness or squareness. These are the most widely-used metrics in literature. The simplest form of common PFEMs results from subtracting a predicted value from its corresponding actual/observed value. This is often straightforward, easy to interpret, and most of all yields the magnitude of error (or difference) in the same units as those measured and predicted and can indicate if the model overestimates or underestimates observations (by analyzing the sign of the reminder). One should remember that an issue could arise where due to the opposite between predictions and observations i.e. canceling positive and negative errors. In this scenario, a zero error could be calculated, indicating false accuracy.

This can be avoided by using an absolute error (i.e. |*A-P*|) which only yields non-negative values. Analogous to traditional error, the absolute error also maintains the same units of predictions (and observations), and hence is easily

relatable. However, due to its nature, the bias in absolute errors cannot be determined.

Similar to the same concept of absolute error, the squared error also mitigates mutual cancellation of errors. This metric can be continuously differentiable and thus facilitates optimization. However, this metric emphasizes relatively large errors (as opposed to small errors), unlike absolute error, and could be susceptible to outliners. The fact that the units of squared error is squared leads to unconventional units for error (i.e. squared days); which are not intuitive. Other metrics may also include logarithmic quotient error (i.e. ln(P/A)) as well as absolute logarithmic quotient error (i.e. ln(P/A)). Table 1 lists other commonly used metrics, together with some of their limitations and shortcomings as identified by surveyed studies.

Most of the works conducted so far in the areas of engineering applications only utilized a few of the above PFEMs [20, 33, 61, 62, 72–92]. The bulk of the reviewed works continue to incorporate traditional metrics such as R,  $R^2$ , MAE, MAPE, and RMSE as primary indicators of adequacy of the regression-based ML models. This seems to stem from our familiarity with these indicators, as opposed to others; such as Golbraikh and Tropsha's [58] criterion, QSAR model by Roy and Roy [59], Frank and Todeschini [60], and specifically designed objective functions, often used in the realms of other fields and data sciences. It should be noted that out of the reviewed studies, the works of Gandomi et al. [90], Golafshani and Behnood [40] as well as Cheng et al. [62] applied a multi-criteria verification process that incorporated the use of traditional as well as modern PFEMs. Utilizing multi-criteria is not only beneficial to ensure the validity of a particular ML model but is also recommended to overcome some of the identified limitations of traditional metrics in Table 1 and hence should be encouraged.

#### Classification

In ML, classification refers to categorizing data into distinct classes. This is a supervised learning approach where machines learn to classify observations into binary or multiclasses. Binary classes are those with two labels (i.e. positive vs. negative etc.), and multi-classes are those having more than two labels (i.e. types of concrete e.g., normal strength, high strength, high performance etc.). Classification algorithms may include logistic regression, k-nearest neighbors, support vector machines, etc. [93, 94].

The performance of classifiers is often listed in a confusion matrix. This matrix contains statistics about actual and predicted classifications and lays the fundamental foundations necessary to understand accuracy measurements for a specific classifier. Each column in this matrix signifies predicted instances, while each row represents actual instances.



<sup>&</sup>lt;sup>2</sup> It should be noted that other works have used a different classification for PFEMs [2]. Botchkarev [2] went even further to survey the most preferred metrics reported by researchers during the 1980–2007 era and also explored multiplication and addition point distance methods.

literature
oben
from
cted
s colle
els as
pou 1
regression
Ē
Σ
for
<b>I</b> s
PE
nseq
duommc
t of c
List
Table 1

No	) Metric	Definition	Formula	Remarks
1	Error (E)	The amount by which an observation differs from its actual value	E = A - P	<ul><li>Intuitive</li><li>Easy to apply</li><li>Works with numeric data</li></ul>
2	Mean error (ME)	The average of all errors in a set	$ME = \frac{\sum_{i=1}^{n} E_i}{n}$	<ul> <li>May not be helpful in cases where positive and negative predictions cancel each other out</li> <li>Works with numeric data</li> </ul>
$\mathcal{S}$	Mean Normalized Bias (MNB)	Associated with observation-based minimum MNB = threshold	$MNB = \frac{\sum_{i=1}^{n} E_i / A_i}{n}$	<ul><li>Biased towards overestimations</li><li>Works with numeric data</li></ul>
4	Mean Percentage Error (MPE)	Computed average of percentage errors	$MPE = \frac{\sum_{i=1}^{n} E_i/A_i}{n/100}$	<ul> <li>Undefined whenever a single actual value is zero</li> <li>Works with numeric data</li> </ul>
ĸ	Mean Absolute Error (MAE)*	Measures the difference between two continuous variables	$MAE = \frac{\sum_{i=1}^{n}  E_i }{n}$	<ul> <li>Uses a similar scale to input data [48]</li> <li>Can be used to compare series of different scales</li> <li>Works with numeric data</li> </ul>
9	Mean Absolute Percentage Error (MAPE)*	Measures the extent of error in percentage terms	$MAPE = \frac{100}{n} \sum_{i=1}^{n}  E_i  /  A_i $	<ul> <li>Commonly-used as a loss function [49]</li> <li>Cannot be used if there are actual zero values</li> <li>Percentage error cannot exceed 1.0 for small predictions</li> <li>There is no upper limit to percentage error in predictions that are too high</li> <li>Non-symmetrical (adversely affected if a predicted value is larger or smaller than the corresponding actual value) [49]</li> <li>Works with numeric data</li> </ul>
7	Relative Absolute Error (RAE)	Expressed as a ratio comparing the mean error to errors produced by a trivial model	$RAE = \sum_{i=1}^{n}  E_i / A_i - A_{mean} $	<ul> <li>E<sub>i</sub> ranges from zero (being ideal) to infinity</li> <li>Works with numeric data</li> </ul>
∞	Mean Absolute Relative Error (MARE)	Measures the average ratio of absolute error to random error	$MARE = \frac{1}{n} \sum_{i=1}^{n}  E_i  /  A_i $	<ul> <li>Sensitive to outliers (especially of low values)</li> <li>Division by zero may occur (if actuals contain zeros)</li> <li>Works with numeric data</li> </ul>
6	Mean Relative Absolute Error (MRAE)	Ratio of accumulation of errors to cumulative error of random error	$MRAE = \sum_{i=1}^{n}  E_i / A_i - A_{mean} $	• For a perfect fit, the numerator equals to zero [50] • Works with numeric data



Tab	Table 1 (continued)			
N <sub>o</sub>	o Metric	Definition	Formula	Remarks
10	Geometric Mean Absolute Error (GMAE)*	Defined as the n-th root of the product of error values	$GMAE = \sqrt[q]{\prod_{i=1}^n  E_i }$	<ul> <li>GMAE is more appropriate for averaging relative quantities as opposed to arithmetic mean [51]</li> <li>This metric can be dominated by large outlies and minor errors (i.e. close to zero)</li> <li>Works with numeric data</li> </ul>
11	Fractional Absolute Error (FAE)	Evaluates the absolute fractional error	$FAE = \frac{1}{n} \sum_{i=1}^{n} \frac{2 \times  E_i }{ A_i  +  P_i }$	• Works with numeric data
12	Mean Squared Error (MSE)	Measures the average of the squares of the errors	$MSE = \frac{\sum_{i=1}^{n} E_i^2}{n}$	<ul> <li>Scale dependent [52]</li> <li>Values closer to zero present adequate state</li> <li>Heavily weights outliers</li> <li>Highly dependent on fraction of data used (low reliability) [53]</li> <li>Works with numeric data</li> </ul>
13	Root Mean Squared Error (RMSE)	Root square of average squared error	$RMSE = \sqrt{\frac{\sum_{i=1}^{n} E_i^2}{n}}$	<ul> <li>Scale dependent</li> <li>A lower value for RMSE is favorable</li> <li>Sensitive to outliers</li> <li>Highly dependent on fraction of data used (low reliability) [53]</li> <li>Works with numeric data</li> </ul>
14	Sum of Squared Error (SSE)	Sums the squared differences between each observation and its mean	$SSE = \sum_{i=1}^{n} E_i^2$	<ul> <li>A small SSE indicates a tight fit [54]</li> <li>Works with numeric data</li> </ul>
15	Relative Squared Error (RSE)	Normalizes total squared error by dividing by the total squared error	$RSE = \sum_{i=1}^{n} E_i^2 / \left( A_i - A_{mean} \right)^2$	<ul> <li>A perfect fit is achieved when the numerator equals to zero [50]</li> <li>Works with numeric data</li> </ul>
16	Root Relative Squared Error (RRSE)	Evaluates the root relative squared error between two vectors	$RRSE = \sqrt{\sum_{i=1}^{n} E_i^2 / \left(A_i - A_{mean}\right)^2}$	<ul> <li>Ranges between zero and 1, with zero being ideal [50]</li> <li>Works with numeric data</li> </ul>
17	Geometric Root Mean Squared Error (GRMSE)	Evaluates the geometric root squared errors	$GRMSE = \sqrt[2n]{\prod_{i=1}^{n} E_i^2}$	<ul> <li>Scale dependent</li> <li>Less sensitive to outliners than RMSE [52]</li> <li>Works with numeric data</li> </ul>
18	: Mean Square Percentage Error (MSPE)*	Evaluates the mean of square percentage errors	$MSPE = \frac{\sum_{i=1}^{n} ( E_i / A_i )^2}{n/100}$	<ul><li>Non-symmetrical [49]</li><li>Works with numeric data</li></ul>
19	Root Mean Square Percentage Error (RMSPE)*	Evaluates the mean of squared errors in percentages	$RMSPE = \sqrt{\frac{\sum_{i=1}^{n} ( E_i / A_i )^2}{n/100}}$	<ul> <li>Scale independent</li> <li>Can be used to compare predictions from different datasets</li> <li>Non-symmetrical [49]</li> <li>Works with numeric data</li> <li>An extension of RMSE</li> </ul>



1				
No	o Metric	Definition	Formula	Remarks
20	Normalized Root Mean Squared Error (NRMSE)**	Normalizes the root mean squared error	$NRMSE = \frac{\sqrt{\sum_{i=1}^{n} E_i^2}}{A_{mean}}$	<ul> <li>Can be used to compare predictions from different datasets [55]</li> <li>Works with numeric data</li> <li>An extension of RMSE</li> </ul>
21	Normalized Mean Squared Error (NMSE)	Estimates the overall deviations between measured values and predictions	$NMSE = \frac{\sum_{i=1}^{n} E_i^2}{n}$ $variance = \frac{\sum_{i=1}^{n} E_i^2}{\sum_{i=1}^{n} mean}^2$	<ul> <li>Biased towards over-predictions [56]</li> <li>Works with numeric data</li> <li>An extension of MSE</li> </ul>
22	Coefficient of Determination (R <sup>2</sup> )	The square of correlation	$\left(\frac{1}{n}\right)^{2}/\sum_{i=1}^{n}\left(A_{i}-A_{mean}\right)^{2}$	<ul> <li>R<sup>2</sup> values close to 1.0 indicate strong correlation</li> <li>Can be used in predicting material properties</li> <li>Works with numeric data</li> <li>Related to R</li> </ul>
23	Correlation coefficient (R)	Measures the strength of association between variables	$R = \frac{\sum_{i=1}^{n} (A_i - \overline{A}_i) (P_i - \overline{P}_i)}{\sqrt{\sum_{i=1}^{n} (A_i - \overline{A}_i)^2 \sum_{i=1}^{n} (P_i - \overline{P}_i)^2}}$	• R > 0.8 implies strong correlation [57] • Does not change by equal scaling • Can be used in predicting material properties • Works with numeric data
24	Mean Absolute Scaled Error (MASE)	Mean absolute errors divided by the mean absolute error	$\frac{\sum_{i=1}^{n} \frac{E_{i}}{N_{i} 100}}{n/100} / (\frac{1}{n} - 1) \sum_{i=1}^{n} \left  A_{i} - A_{i-1} \right $	<ul> <li>Scale independent</li> <li>Stable near zero [52]</li> <li>Works with numeric data</li> </ul>
25	Golbraikh and Tropsha's [58] criterion		At least one slope of regression lines $(k \text{ or } k')$ between the regressions of actual $(A_i)$ against predicted output $(P_i)$ or $P_i$ against $A_i$ through the origin, i.e. $A_i = k \times P_i \underbrace{q_i d_i P_i = k' A_i}_{P_i = M_i A_i X_i P_i}_{P_i = M_i A_i X_i P_i}$ ss $k' = \frac{\sum_{i=1}^n (A_i X_i P_i)}{\sum_{i=1}^n A_i A_i X_i P_i}$ ss $m = \frac{P_i - M_i P_i}{R^2}$ $n = \frac{P_i - M_i P_i}{R^2}$	<ul> <li>• k and k' need to be close to 1 or at least within the range of 0.85 and 1.15</li> <li>• m and n are performance indexes and their absolute value should be lower than 0.1</li> <li>• Works with numeric data</li> </ul>
26	OSAR model by Roy and Roy [59]		$\times (1 - \sqrt{ R^2 - R_o^2 })$ $= \sum_{n=1}^{n-A_o} \frac{1}{\sqrt{2}}, A_i^o - k \times P_i R_i^o = 1 - \sum_{n=1}^{n} \frac{(A_i - A_i^o)^2}{\sqrt{2}}$	• $R_m$ is an external predictability indicator. $R_m > 0.5$ implies a good fit $P_{n+1} = P_n $
27	Frank and Todeschini [60]		Recommend maintaining a ratio of $3-5^{L_{i=1}(A_i-A_i)}$ between the number of observations and input parameters	
58	Objective function by Gandomi et al. [61]	A multi-criteria metric	Function = (No. Training —No. Validation) RMSE_Training +MAE_Learner_Fluid Fluid Flu	்திழ் furktion and edge with the further further further high specification and service of the further by the further furthe



a	lable 1 (continued)			
No	o Metric	Definition	Formula	Remarks
53	Reference index (RI) by Cheng et al. [62]	A multi-criteria metric that uniformly accounts for RMSE, MAE and MAPE	$RI = \frac{RMSE + MAE + MAPE}{3}$	<ul> <li>Each fitness metric is normalized to achieve the best performance</li> <li>Works with numeric data</li> <li>An extension of RMSE, MAE and MAPE</li> </ul>
30	Scatter index (SI) [63]	Applied to examine whether RMSE is good or not	$SI = \frac{\sqrt{\sum_{i=1}^{n} \binom{p_{max(i)} - P_{max(i)}}{n}}}{P_{max(i)}}$ where, $n = number$ of data sets used during the training phase $.P_{max(p)} = mean$ actual observations data	• SI is RMSE normalised to the measured data mean • If SI is less than one, then estimations are acceptable • Works with numeric data • "excellent performance" when SI<0.1, a "good performance" when 0.1 < SI<0.2, a "fair performance" when 0.2 < SI<0.3, and a "poor performance" when SI>0.3
31	Synthesis index (SyI) [64]	Comprehensive performance measure a based on MAE, RMSE, and MAPE a	$SyI = \frac{1}{n} \sum_{i=1}^{n} \left( \frac{P_i - P_{mini}}{P_{max,i} - P_{mini}} \right)$ where, $n = number$ of performance measures; and $P_i = ith$ performance measure	• The SI ranged from 0 to 1; an SI value close to 0 indicated a highly accurate predictive model • Works with numeric data
32	Relative root mean squared error (RRMSE) [65]	Present percentage variation in accuracy	$RRMSE = \sqrt{\frac{1}{n}\sum (A - P)^2}$	<ul> <li>Lower RRMSE values result in more accurate model predictions</li> <li>Works with numeric data</li> </ul>
33	Performance index (PI) [65]	Performance index to evaluate predictivity of a model	$PI = \frac{RRMSE}{1+R}$	<ul> <li>Lower PI values result in more accurate model predictions</li> <li>Works with numeric data</li> </ul>
34	420-index [66]	Performance index to evaluate predictivity of a model within 20% variation	$a_{20-index} = \frac{m_{20}}{M}$ where, $m_{20}$ is the number of samples with the ratio of experimental value over predicted value falling from 0.8 to 1.2 and M is the number of samples in the dataset	<ul> <li>Presents the number of samples with the difference between the predicted value and experimental value within ± 20%</li> <li>Works with numeric data</li> </ul>
35	Fractional bias (FB) [67]	Measure of the shift between the observed and predicted values	$FB = rac{2\sum_{i=1}^{n}(A-P)}{\sum_{i=1}^{n}(A+P)}$	<ul> <li>Dimensionless metric, which is convenient for comparing the results from studies involving different scales</li> <li>Symmetrical and bounded; values for the fractional bias range between -2.0 (extreme underprediction) to +2.0 (extreme overprediction)</li> <li>Perfect model has FB of zero</li> <li>Works with numeric data</li> </ul>
36	Relative index of agreement (RD) [68]	A standardized measure of the degree of model prediction error	$RD = 1 - \frac{\sum_{i=1}^{N} (\frac{A_i - P}{A})}{\sum_{i=1}^{N} (\frac{(P - A_i + A_i - A_i)}{A})}^{2}$	A value of 1.0 indicates a perfect match, and zero indicates no agreement at all     Overly sensitive to extreme values     Works with numeric data



No Metric 37 Nash–Sutcliffe coefficient (NSE) [69]			
37 Nash–Sutcliffe coefficient (NSE) [69]	Definition	Formula	Remarks
	A metric often used in flow predictions	NSE = 1 - $\left[\sum_{k=1}^{N} \frac{(A-P)^2}{(A-\overline{A})^2}\right]$	<ul> <li>NSE = 1 indicates perfect correspondence</li> <li>NSE = 0 indicates that the model simulations have the same explanatory power as the mean of the observations</li> <li>NSE &lt; 0 indicates that the model is a worse predictor than the mean of the observations</li> <li>Works with numeric data</li> </ul>
38 Kling-Gupta efficiency (KGE) [70]	A metric often used in flow predictions	KGE = $1 - \sqrt{(r - 1)^2 + (\alpha - 1)^2 + (\beta - 1)^2}$ , where, $r$ is the linear correlation between the predicted and actuals. $\alpha$ is the magnitude of the variability calculated as the standard deviation in predictions divided by the standard deviation in actuals. $\beta$ is the bias term calculated as the predictions means divided by the actual mean. $N$ is the number of dataset over the training and testing phases	KGE=1 indicates perfect agreement between actuals and predictions     KGE<0 indicates that the mean of actuals provides better estimate than predictions     For other values of KGE, please refer to [71]     Works with numeric data



<sup>\*\*</sup>can be normalized by standard deviation of actual observations

<sup>\*\*\*</sup>The reader is encouraged to review the cited references for full details on specific metrics.

This matrix was identified to be the "go-to" metric used in studies examining materials science and engineering problems [22, 95–98]. However, there are other PFEMs that can be used to evaluate classification models, and these, along with others, are listed in Table 2. Similar to Table 1, Table 2 also lists some of the remarks and limitations pointed out by surveyed works. In this table, *P* (denotes number of real positives), *N* (denotes number of real negatives), *TP* (denotes true positives), *TN* (denotes true negatives), *FP* (denotes false positives), and *FN* (denotes false negatives).

# **Closing Remarks**

Our confidence in the accuracy of predictions obtained from ML algorithms heavily relies on the availability of actual observations and proper PFEMs. From this point of view, it is unfortunate that observations relating to the engineering discipline continue to be 1) limited in size, and 2) lack completeness. The lack of such observations is often related to limitations in conducting full-scale tests, the need for specialized equipment, and a wide variety of tested samples. For instance, one can think of how normal strength concrete mixes can significantly vary from one study to another simply due to variation in raw materials, mix proportions, and casting/curing procedures, etc.

Combining the above two points with the notion of simply "applying ML" to understand a given phenomenon (say flexural strength of beams) without a thorough validation is deemed to fail. In fact, in many instances, researchers noted the validity of a specific ML model by reporting its performance against traditional PFEMs, only to be later identified that such a model does not properly represent actual observations - despite having good fitness. This can be avoided by adopting a rigorous validation procedure [121, 122]. Unfortunately, many of the published studies in the area of ML application in engineering do not include multi-criteria/ additional validation phases and simply rely on conventional performance metrics such as R or  $R^2$  of the derived models. Furthermore, adopting a set of PFEMs does not negate the occurrence of some common issues, most notably, overfitting, biasedness etc. As such, an analysis that utilizes ML should also consider some of the following techniques e.g. use of independent test datasets, varying degrees of crossvalidation etc.

In order to ensure fruitful use of ML, it is our duty to seek proper application of ML. Besides, one of the major concerns about the ML-based models is their robustness under a wide range of conditions [123]. A robust ML model should not only provide reasonable PFEMs but should also be capable of capturing the underlying physical mechanisms that govern the investigated system [124]. An essential approach to verify the robustness of the ML models is

to perform parametric and sensitivity analyses [123, 125]. These types of analyses ensure that the ML predictions are in sound agreement with the system's real behavior and physical processes rather than being merely a combination of the variables with the best fit on the data. Another item to consider is to develop a user-friendly phenomenon-specific recommendation system wherein novice users who apply preidentified PFEMs are selected to evaluate the performance of a given problem (say using  $R^2$  in a regression problem etc.).

The reader is to remember that the addition of one example to showcase recommended or important PFEMs negates the purpose of this paper (which is to compile commonly used performance metrics and list their key characteristics into one document to provide interested researchers in carrying out a ML analysis with a starting point to select proper performance metrics). Providing a comparison for all of the reviewed metrics will significantly extend this work beyond its scope and may not be feasible at the moment. We feel that this is best suited for a series of more in-depth reviews wherein metrics for classification and regression problems can be separately evaluated and reviewed under well-designed problems and a variety of conditions to ensure fairness and unbiasedness to come in the near future.

It is our intention to not specifically identify a measure (or a set of measures) due to the wide range of problems (as well as the quality of data) that a scientist could face. Please note that other researchers (which are quoted herein) also followed a similar approach.

- o "Although some methods clearly perform better or worse than other methods on average, there is significant variability across the problems and metrics. Even the best models sometimes perform poorly, and models with poor average performance occasionally perform exceptionally well." [126].
- "It is clearly difficult to convincingly differentiate ML algorithms (and feature reduction techniques) on the basis of their achievable accuracy, recall and precision."[127].
- q "Different performance metrics yield different tradeoffs that are appropriate in different settings. No one metric does it all, and the metric optimized to or used for model selection does matter."[102].

### **Conclusions**

Based on the information presented in this note, the following conclusions can be drawn.

ML is expected to rise into a key analysis tool in the coming few years; especially within material scientists and structural engineers. As such, the integration of ML is



 Table 2
 List of the commonly-used PFEMs for ML classification models as collected from open literature

2				
No	Metric	Definition	Formula	Remarks
_	True Positive Rate (TPR) or Sensitivity or Recall	Measures the proportion of actual positives that are correctly identified as positives	$TPR = \frac{TP}{P} = \frac{TP}{TP + FN} = 1 - FNR$	<ul> <li>Describes the proportion of actual positives that are correctly identified</li> <li>Does not account for indeterminate results</li> <li>Works with categorial data</li> </ul>
2	True Negative Rate (TNR) or Specificity or selectivity	Measures the proportion of actual negatives that are correctly identified negatives	$TNR = \frac{TN}{N} = \frac{TN}{TN + FP} = 1 - FPR$	<ul> <li>Describes the proportion of actual negatives that are correctly identified</li> <li>Works with categorial data</li> </ul>
8	Positive Predictive Value (PPV) or Precision	The proportions of positive observations that are true positives	$PPV = \frac{TP}{TP+FP} = 1 - FDR$	<ul> <li>Has an ideal value of 1 and the worst value of zero</li> <li>Works with categorial data</li> </ul>
4	Negative Predictive Value (NPV)	The proportions of negative observations that are true positives	$NPV = \frac{TN}{TN + FN} = 1 - FOR$	<ul> <li>Has an ideal value of 1 and the worst value of zero</li> <li>Works with categorial data</li> </ul>
S	False Positive Rate (FPR)	Measures the proportion of positive cases in that are correctly identified as positives	$FPR = \frac{FP}{N} = \frac{FP}{FP + TN} = 1 - TNR$	<ul> <li>Describes proportion of negative cases incorrectly identified as positive cases</li> <li>Works with categorial data</li> </ul>
9	False Discovery Rate (FDR)	Expected proportion of false observations	$FDR = \frac{FP}{FP+TP} = 1 - PPV$	<ul> <li>Describes proportion of the individuals with a positive test result for which the true condition is negative</li> <li>Works with categorial data</li> </ul>
7	False Omission Rate (FOR)	Measures the proportion of false negatives that are incorrectly rejected	$FDR = \frac{FN}{FN + TPN} = 1 - NPV$	<ul> <li>Describes proportion of the individuals with a negative test result for which the true condition is positive</li> <li>Works with categorial data</li> </ul>
∞	Positive likelihood ratio (LR+)	Evaluates the change in the odds of having a diagnosis with a positive test	$LR + = \frac{TPR}{FPR}$	<ul> <li>Measures the ratio of TPR (sensitivity) to the FPR (1 – specificity)</li> <li>Presents the likelihood ratio for increasing certainty about a positive diagnosis</li> <li>Works with categorial data</li> </ul>
6	Negative likelihood ratio (LR-)	Evaluates the change in the odds of having a diagnosis with a negative test	$LR - = \frac{FNR}{TNR}$	<ul> <li>Describes the ratio of FNR to TNR (specificity)</li> <li>Works with categorial data</li> </ul>
10	Diagnostic odds ratio (DOR)	Measures the effectiveness of a (diagnostic) test	$DOR = \frac{LR+}{LR-} = \frac{TP/FP}{FN/TN}$	<ul> <li>Often used in binary classification</li> <li>Works with categorial data</li> </ul>
11	Accuracy (ACC)	Evaluates the ratio of number of correct predictions to the total number of samples	$ACC = \frac{TP + TN}{P + N} = \frac{TP + TN}{TP + TN + FP + FN}$	<ul> <li>Presents performance at a single class threshold only</li> <li>Assumes equal cost for errors [96]</li> <li>Works with categorial data</li> </ul>



Processing   Pro					
F <sub>1</sub> score Harmonic mean of the precision and recall $r_1 = \frac{2PP_1 + 2PP_2}{PP_1 + PP_2} = \frac{2TP}{2PP_2 + PP_2}$ Matthews Correlation Coefficient (MCC) Measures the quality of binary classifications $MCC = \frac{2PP_2 + PP_2}{(D^2 + P^2 N^2 + P^2 N^2 + PP_2 N^2 N^2 + PP_2 N^2 N^2 + PP_2 N^2 + PP$	No	Metric	Definition	Formula	Remarks
Matthews Correlation Coefficient (MCC) Measures the quality of binary classifications $MCC = \frac{TPATN - FPAZN}{(TP + FPZN + FP$		F <sub>1</sub> score	Harmonic mean of the precision and recall	$F_1 = \frac{2PPV \times TPR}{PFV + TPR} = \frac{2TP}{2TP + FP + FN}$	<ul> <li>Describes the harmonic mean of precision and sensitivity</li> <li>Focuses on one class only</li> <li>Biased to the majority class [99]</li> <li>Works with categorial data</li> </ul>
Bookmaker Informedness (BM) or Youden's Evaluates the discriminative power of the $BM = TPR + TNR - 1$ test [101] test [101]  Markedness (MK) Measures trustworthiness of positive and $MK = PPV + NPV - 1$ negative predictions	13	Matthews Correlation Coefficient (MCC)	Measures the quality of binary classifications analysis	$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TN + FP)(TN + FP)(TN + FN)}}$	<ul> <li>Measures the quality of binary and multiclass classifications</li> <li>Can be used in classes with different sizes</li> <li>When MCC equals +1 → perfect prediction, →0 equivalent to a random prediction and → -1 false prediction</li> <li>Considered as a balanced measures as it involves values of all the four quardants of a confusion matrix [100]</li> <li>Works with categorial data</li> </ul>
Markedness (MK) Measures trustworthiness of positive and $MK = PPV + NPV - 1$ negative predictions	41	Bookmaker Informedness (BM) or Youden's J statistic		BM = TPR + TNR - 1	<ul> <li>Describes the probability of an informed decision (vs. a random guess)</li> <li>Has a range between zero and 1 (being ideal)</li> <li>Considers both real positives and real negatives</li> <li>Takes into account all predictions [102]</li> <li>Works with categorial data</li> <li>Counterpart of recall</li> <li>It is also suitable with imbalanced data</li> <li>It does not change concerning the differences between the sensitivity and specificity [101]</li> </ul>
• Works with categorial data	15	Markedness (MK)		MK = PPV + NPV - 1	Measures trustworthiness of positive and negative predictions by a model [103]     Considers both predicted positives and predicted negatives     Counterpart of precision     Specifies the probability that a condition is marked by the predictor (as opposed to luck/chance) [104]     Sensitive to data changes (not suitable for imbalanced data) [101]     Works with categorial data



No Metric 16 Averag				
	Metric	Definition	Formula	Remarks
	Average Class Accuracy (ACA)	Measures the average accuracy of predictions in a class	$ACA = W\left(\frac{TP}{TP+FP}\right) + (1-W)\left(\frac{TN}{TN+FP}\right)$ where $0 < W < 1$	<ul> <li>Used with unbalanced data</li> <li>Choosing a good weighting factor a priori [99]</li> <li>When W&gt; 0.5, minority class accuracy contributes more than majority class</li> <li>Presents performance at a single class threshold</li> <li>Works with categorial data</li> </ul>
17	Receiver Operating Characteristic (ROC)	Plots the diagnostic ability of a binary classifier system as its discrimination threshold is varied	The ROC curve is plotted such that TPR is on vertical axis and FPR is on the horizontal axis (the line TPR=FPR represents a random guess of a specific class) [105]	<ul> <li>Characterizes tradeoff between hit rate and false alarm rate</li> <li>Designates the relationship between sensitivity and specificity [106]</li> <li>Takes a value between zero and 1 to relate the probability distribution to a single state [107]</li> <li>A threshold of zero ensures highest sensitivity and 1 ensures best specificity</li> <li>Can be used to estimate cost ratio (slope of line tangent to ROC curve)</li> <li>Should be used in datasets with roughly equal numbers of observations for each class [108, 109]</li> <li>Works with categorial data</li> </ul>
81	Area under the ROC curve (AUC)	Measures the two-dimensional area underneath the entire ROC curve	$AUC = \sum_{i=1}^{N-1} \frac{1}{2} (FP_{i+1} - FP_i) (TP_{i+1} - TP_i)$ or $AUC = \frac{1}{2} w(h+h'),$ where, $w = width$ , and $h$ and $h' = heights$ of the sides of a trapezoid histogram	<ul> <li>Not dependent on a single class threshold</li> <li>Associated with increased training times</li> <li>Works with categorial data</li> </ul>
19	Precision-Recall curve	Plots the tradeoff between precision and recall for different thresholds	Plots precision (in the vertical axis) and the recall (in the horizontal axis) for different thresholds	<ul> <li>Applicable in cases of moderate to large class imbalance [108]</li> <li>Used in binary classification</li> </ul>
20	Log Loss Error (LLE)	Measures the where the prediction input is a probability value	Where, $M$ : $a_{i}$ $a_{i}$ $a_{j}$ $a_{i}$ $a_{j}$ $b_{i}$ $a_{i}$ $b_{i}$	<ul> <li>Measures the uncertainty of the probabilities by comparing predictions to the true labels</li> <li>Penalizes for being too confident in wrong prediction</li> <li>Has probability between zero and 1</li> <li>A log loss of zero indicates a perfect model</li> <li>Works with categorial data</li> </ul>
21	Hinge Loss Error (HLE)		$HLE = max(0, 1 - q \cdot y)$ where, $q = \pm 1$ and y: classifier score	Linearly penalize incorrect predictions     Primarily used in support vector machine



Tab	Table 2 (continued)			
<sup>N</sup> o	Metric	Definition	Formula	Remarks
22	Wilcoxon-Mann-Whitney (WMW) test [99]		$WMW = \frac{\sum_{i \in Minorclass} \sum_{i \in Migorclass} I_{nom.}(P_i.P_i)}{ Minorclass  \times  Minorclass  \times$	<ul> <li>Used in scenarios with unbalanced data</li> <li>The indicator function I<sub>wmw</sub> returns 1 if</li> <li>P<sub>i</sub> &gt; P<sub>j</sub> and P<sub>i</sub> ≥ 0 or 0 if otherwise</li> </ul>
23	Fitness Function Amse (FFA) [99]	Measures pattern difference between input and output	$FFA = \frac{1}{K} \sum_{c=1}^{K} \left( 1 - \frac{\sum_{i=1}^{W_c} (1 - sig(P_c)) - T_c)}{N_c \times 2} \right)^2,$ $sig(x) = \frac{2}{1 + e^{-x}} + 1$ where, $P_{ci}$ output of a classifier evaluated on the ith example, $N_c$ : number of examples, $K$ : number of classes, $T_c$ : target values (equals to -0.5 and 0.5 for majority and minority classes, respectively)	<ul> <li>Used in scenarios with unbalanced data</li> <li>Appropriate for genetic programing</li> <li>Needs to be scaled to a range of [-1, 1] and hence the need for sigmoid function</li> <li>FFA = 1 presents an ideal scenario</li> </ul>
24	Fitness Function Incr (FFI) [99]		$\begin{cases} Eq(D_{c_i}, P_{c_i}) \\ c \in Minority cl \\ c \in Majority \end{cases}$ herwise	<ul> <li>Used in scenarios with unbalanced data</li> <li>Sasigns incremental rewards to predictions</li> <li>clateat fall further away from the class boundary</li> </ul>
25	Fitness Function Correlation (FFC)		$Eq(p,q) = \begin{cases} 1, & \text{if } p = q \\ \frac{1}{\sqrt{\sum_{k=1}^{K} \frac{1}{N_k} (q_1 - p_1)^2}}, & \text{otherwise} \end{cases}$ $Ff_k C = \int_{V_k} \frac{\left(\sum_{k=1}^{K} \frac{1}{N_k} (q_1 - p_1)^2\right)}{\sum_{k=1}^{K} \frac{1}{N_k} (q_1 - p_1)^2},$ $W_t^{there} \sum_{k=1}^{K_k} \frac{1}{N_k} \sum_{i=1}^{K_k} \frac{1}{N_k} \frac{1}{N_$	<ul> <li>Appropriate for genetic programming</li> <li>Ranges [0, 1] (zero being worst fitness)</li> <li>Used in scenarios with unbalanced data</li> </ul>
26	Fitness Function Distribution (FFD)	Measures the distance between class distributions as a function of class separability	$FFD = \frac{ \mu_{\min} - \mu_{\min} }{\sum_{i=1}^{N_c} \sigma_{\min}^{n} + \sigma_{\min}} \times I_{2r}(2, \mu_{\min}, \mu_{\min})}$ $\mu_c = \frac{\sum_{i=1}^{N_c} \sigma_{\min}^{n} + \sigma_{\min}}{\sum_{i=1}^{N_c} (P_{ci} - \mu_c)^2}.$ where, $\mu_c^{c}$ and $\sigma_c$ : $\mu_c^{cm}$ and standard deviation of the class distribution, respectively,	<ul> <li>Used in scenarios with unbalanced data</li> <li>Treats predictions as independent distributions</li> <li>Measures separability (i.e. distance between class distributions) [110] – high separability (no overlap) and this distance turns large (go to +∞)</li> <li>Uses I<sub>n</sub> to enforce zero class threshold</li> </ul>
27	Canberra Metric (CM)	Measures the distance between pairs of points in a vector space	$CM = \sum_{i=1}^{n} \frac{ E_i }{A_i + P_i}$	
78	Wave Hedges Distance (WHD)		$WHD = \sum_{i=1}^{n} \frac{ E_i }{\max(A_i, P_i)}$	Normalizes the difference of each pair of coefficients with its maximum [111-113]



No Metric	Definition	Formula	Remarks
29 Lift [114]	Measures the performance of a model at predicting or classifying cases	$LIFT=rac{\%oftwapositivesxaboverheelnreshold}{\%oftkuasetaboverheelnreshold}$	Measures betterness of a classifier than a baseline classifier that randomly predicts positives     Threshold is set as a static fraction of the positive dataset     Lift and Accuracy do not always correlate well
30 Mean Cross Entropy (MXE)	Measures the performance of a model where the output is a probability between zero and one	$MXE = -\frac{1}{N} \sum True \times ln(Predicted) + (1 - Tru)$ (The assumptions are that $Predicted \in [0, 1]$ and $True \in \{0, 1\}$ )	$MXE = -\frac{1}{N}\sum True \times ln(Predicted) + (1 - True)$ WithitpizipseMASEagives the maximum likeli- (The assumptions are that $Predicted \in [0, 1]$ hood [102] and $True \in [0, 1]$ )
		1. Order cases 1–100 by their predicted in the same bin 2. Evaluate the percentage of true positives 3. Calculate the mean prediction for true positives 4. Calculate the mean prediction calibration error for this bin (using the absolute value of the difference between the observed frequency and the mean) 5. Repeat steps 1–4 for cases 2–101, 3–102, etc 6. CAL is calculated as the mean of these binned calibration errors [102]	• Lengthy procedure
32 Precision-recall break-even point	Point at which the precision-recall-curve intersects the bisecting line	Precision = Recall	Defines the point when precision and recall are equal
33 Average precision (AP)	Combines recall and precision for ranking	$AP = \sum_{n} (Recall_{n} - Recall_{n-1}) Percision_{n}$	• Describes the weighted mean of precision in each threshold with the increase in recall from the previous threshold used
34 Balanced accuracy [115]	Calculates the average of the correctly identified proportion of individual classes	Defined as the average of recall obtained on each class	Used in binary and multiclass classification problems     Accommodates imbalanced datasets
35 Brier score (BS)	Measures the accuracy of probabilistic-based predictions	$BS = \frac{1}{N} \sum_{i=1}^{N} (f_i - A_i)^2$ in which $f_j$ is the probability that was forecast, $A_i$ the actual outcome of the event at instance i	Measures the mean squared difference between the predicted probability and the actual outcome     Takes on a value between zero and 1 (the lower the score is, the better the predictions)     Composed of refinement loss and calibration loss     Appropriate for binary and categorical outcomes     Inappropriate for ordinal variables



No Metric	Definition	Formula	Remarks
36 Cohen's kappa (CK) [116]	Measures interrater (agreement) reliability	$\kappa = (p_o - p_e)/(1 - p_e)$ where, $p_o$ ; empirical probability of agreement on the label assigned to any sample, $p_e$ ; expected agreement when both annotators assign labels randomly and this is estimated using a per-annotator empirical prior over the class labels	Measures inter-annotator agreement     Expresses the level of agreement between two annotators [117]     Ranges between -1 and 1. The maximum value means complete agreement
37 Hamming loss (HL)	Fraction of the wrongly identified labels	$HL = \frac{1}{m} \sum_{i=1}^{m} 1_{p_{\frac{i+\lambda_i}{m}}}$	<ul> <li>Describes fraction of labels that are incorrectly predicted</li> <li>Optimal value is zero [118]</li> </ul>
		Finness(T) = $Q(T) + \alpha * R(T) + \beta * Cost(T)$ where, $Q(T)$ : accuracy, $R(T)$ : sum of $R(T_i)$ in all multi-tests of the $T$ tree, $Cost(T)$ : sum of the costs of attributes constituting multi-tests. The default parameters values are: $1 - 0$ and $\beta = -0.5$ , $R(T_i) = \frac{ X_i }{ X_i } * \sum_{j=1}^{ m_i -1} r_{ij}$ where, $X$ : learning set, $X_i$ : instances in i-th node, and $ m_i $ : size of a multi-test $Cost(T_i) = \frac{ X_i }{ X_i } * C(\alpha_{ij})$ where: $\alpha_{ij}$ : $\frac{ X_i }{ X_i } * C(\alpha_{ij})$ where: $\alpha_{ij}$ : $\frac{ X_i }{ X_i } * C(\alpha_{ij})$	Used for fitting decision trees     This function needs to be maximized to achieve high performance
39 F2 score [120]	Measured as the weighted average of precision and recall	$F_{\beta} = 1 + \beta 2 \times \frac{precision \times recall}{(\beta 2 \times precision) + recall}$ where: $\beta = 2$	Used in genetic programming and medical fields Computes a weighted harmonic mean of Precision and Recall  Learning about the minority class
40 Distance score (D score) [120]		$\begin{array}{ll} D_{xc} &= \frac{2 \times C1 \times C2}{C1 + C2} \\ where \sum_{i=0}^{N_{eng}} sig(P_{M_{eng}}) \times  T-sig(P_{M_{eng}}) }{C1 = \sum_{i=0}^{N_{eng}} sig(P_{M_{eng}}) \times  T-sig(P_{M_{eng}}) } \times func(1,P_{M_{eng}}) & \text{Distance score (D score) which learns ab} \\ sig(x) &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \frac{N_{eng}}{M_{eng}} \times func(1,P_{M_{eng}}) & \text{Distance score (D score) which learns ab} \\ C2 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \frac{N_{eng}}{M_{eng}} \times func(1,P_{eng}) & \text{Distance score (D score) which learns ab} \\ C2 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \frac{N_{eng}}{M_{eng}} \times func(1,P_{eng}) & \text{Distance score (D score) which learns ab} \\ C2 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \frac{N_{eng}}{M_{eng}} \times func(1,P_{eng}) & \text{Distance score} \\ C3 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \frac{N_{eng}}{M_{eng}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{eng})^{N_{eng}}} \times func(1,P_{eng}) & \text{Distance score} \\ C4 &= \sum_{i=0}^{N_{eng}} \frac{2}{sig(N_{$	$\begin{aligned} & D_{xc} = \frac{2 \times CI \times C2}{CI + C2} \\ & where: \frac{C_{1} + C2}{CI + C2} \\ & = \frac{C_{1} + C2}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times func(1, P_{Maji}) \\ & = \frac{1}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times func(1, P_{Maji}) \\ & sig(x) = \frac{1}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times func(1, P_{Maji}) \\ & = \frac{1}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times f(P_{Maji}) \\ & = \frac{1}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times f(P_{Maji})} \times f(P_{Maji}) \\ & = \frac{1}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times f(P_{Maji})} \times f(P_{Maji}) \\ & = \frac{1}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times f(P_{Maji})} \times f(P_{Maji}) \\ & = \frac{1}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times f(P_{Maji})} \times f(P_{Maji}) \\ & = \frac{1}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times f(P_{Maji})} \times f(P_{Maji})} \times f(P_{Maji}) \\ & = \frac{1}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times f(P_{Maji})} \times f(P_{Maji})} \times f(P_{Maji}) \\ & = \frac{1}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times f(P_{Maji})} \times f(P_{Maji})} \times f(P_{Maji}) \\ & = \frac{1}{\sum_{i=0}^{N_{opt}} sig(P_{Maji})} \times f(P_{Maji})} \times f($
*The reader is encouraged to review the cited references for full details on specific metrics	ences for full details on specific metrics		

The reader is encouraged to review the cited references for full details on specific metrics.



to be thorough and proper. Hence, the need for proper validation procedure.

A variety of performance metrics and error metrics exists for regression and classification problems. This work recommends the utilization of multi-fitness criteria (where a series of metrics are checked on one problem) to ensure the validity of ML models as these metrics may overcome some of the limitations of induvial metrics. Such metrics can be of independent nature to each other such as,  $R^2$ , RSME, and  $a_{20-index}$ .

The performance of the existing metrics and future fitness functions can be further improved through systematic collaboration between researchers of interdisciplinary backgrounds. For example, efforts are invited to identify and recommend metrics suitable for specific problems and datasets.

Future works should be directed towards documenting and exploring performance metrics for other types of learnings such as unsupervised learning and reinforcement learning. This is ongoing research need that is to be addressed in the coming years.

**Data Availability** No data, models, or code were generated or used during the study.

#### **Declarations**

Conflict of interest none.

## References

- Mahdavi S, Rahnamayan S, Deb K (2018) Opposition based learning: A literature review. Swarm Evol Comput. https://doi. org/10.1016/j.swevo.2017.09.010
- Botchkarev A (2019) A new typology design of performance metrics to measure errors in machine learning regression algorithms. Interdiscip J Information Knowledge Manag 14:045–076. https://doi.org/10.28945/4184
- Bishop C (2007) Pattern Recognition and Machine Learning. Technometrics. https://doi.org/10.1198/tech.2007.s518
- Fu G-S, Levin-Schwartz Y, Lin Q-H, Zhang D (2019) Machine Learning for Medical Imaging. J Healthc Eng. https://doi.org/10. 1155/2019/9874591
- Michalski, R. S., Carbonell, J. G., & Mitchell TM (1983)
   Machine learning: An artificial intelligence approach.
- Majidifard H, Jahangiri B, Buttlar WG, Alavi AH (2019) New machine learning-based prediction models for fracture energy of asphalt mixtures. Meas J Int Meas Confed. https://doi.org/10. 1016/j.measurement.2018.11.081
- Hu X, Li SE, Yang Y (2016) Advanced Machine Learning Approach for Lithium-Ion Battery State Estimation in Electric Vehicles. IEEE Trans Transp Electrif. https://doi.org/10.1109/ TTE.2015.2512237

- 8. Voyant C, Notton G, Kalogirou S, et al (2017) Machine learning methods for solar radiation forecasting: A review. Renew. Energy
- Shukla R, Singh D (2017) Experimentation investigation of abrasive water jet machining parameters using Taguchi and Evolutionary optimization techniques. Swarm Evol Comput. https://doi.org/10.1016/j.swevo.2016.07.002
- Hindman M (2015) Building Better Models: Prediction, Replication, and Machine Learning in the Social Sciences. Ann Am Acad Pol Soc Sci. https://doi.org/10.1177/0002716215570279
- Grimmer J (2014) We are all social scientists now: How big data, machine learning, and causal inference work together. In: PS - Political Science and Politics
- Naser M, Chehab A (2018) Materials and design concepts for space-resilient structures. Prog Aerosp Sci 98:74–90. https:// doi.org/10.1016/j.paerosci.2018.03.004
- Rashno A, Nazari B, Sadri S, Saraee M (2017) Effective pixel classification of Mars images based on ant colony optimization feature selection and extreme learning machine. Neurocomputing. https://doi.org/10.1016/j.neucom.2016.11.030
- Jordan MI, Mitchell TM (2015) Machine learning: Trends, perspectives, and prospects. Science 349:255–260. https://doi.org/10.1126/science.aaa8415
- Seitllari A (2014) Traffic Flow Simulation by Neuro-Fuzzy Approach. In: Second International Conference on Traffic. Belgrade, pp 97–102
- Naser MZ (2019) AI-based cognitive framework for evaluating response of concrete structures in extreme conditions.
   Eng Appl Artif Intell 81:437–449. https://doi.org/10.1016/J.
   ENGAPPAI.2019.03.004
- Li X, Qiao T, Pang Y et al (2018) A new machine vision realtime detection system for liquid impurities based on dynamic morphological characteristic analysis and machine learning. Meas J Int Meas Confed. https://doi.org/10.1016/j.measu rement.2018.04.015
- Oleaga I, Pardo C, Zulaika JJ, Bustillo A (2018) A machinelearning based solution for chatter prediction in heavy-duty milling machines. Meas J Int Meas Confed. https://doi.org/10. 1016/j.measurement.2018.06.028
- Shanmugamani R, Sadique M, Ramamoorthy B (2015) Detection and classification of surface defects of gun barrels using computer vision and machine learning. Meas J Int Meas Confed. https://doi.org/10.1016/j.measurement.2014.10.009
- Naser MZ (2019) Properties and material models for common construction materials at elevated temperatures. Constr Build Mater 10:192–206. https://doi.org/10.1016/j.conbuildmat. 2019.04.182
- Raccuglia P, Elbert KC, Adler PDF et al (2016) Machinelearning-assisted materials discovery using failed experiments. Nature. https://doi.org/10.1038/nature17439
- 22. Alavi AH, Hasni H, Lajnef N et al (2016) Damage detection using self-powered wireless sensor data: An evolutionary approach. Meas J Int Meas Confed. https://doi.org/10.1016/j.measurement.2015.12.020
- Farrar CR, Worden K (2012) Structural Health Monitoring: A Machine Learning Perspective
- Mcfarlane C (2011) The city as a machine for learning. Trans Inst Br Geogr. https://doi.org/10.1111/j.1475-5661.2011. 00430.x
- Chan J, Chan K, Yeh A (2001) Detecting the nature of change in an urban environment: A comparison of machine learning algorithms. Photogramm. Eng. Remote Sensing
- 26. King DE (2009) Dlibml: A Machine Learning Toolkit. J Mach
- 27. Collobert R, Kavukcuoglu K, Farabet C (2011) Torch7: A Matlab-like Environment for Machine Learning



- Hall M, Frank E, Holmes G et al (2009) The WEKA data mining software. ACM SIGKDD Explor Newsl DOI 10(1145/1656274):1656278
- 29. Ramsundar B (2016) TensorFlow Tutorial. CS224d
- Zaharia M, Franklin MJ, Ghodsi A et al (2016) Apache Spark. Commun ACM. https://doi.org/10.1145/2934664
- 31. Korolov M (2018) AI's biggest risk factor: Data gone wrong | CIO. In: CIO. https://www.cio.com/article/3254693/ais-bigge st-risk-factor-data-gone-wrong.html. Accessed 5 Jul 2019
- Kodur VKR, Garlock M, Iwankiw N (2012) Structures in Fire: State-of-the-Art, Research and Training Needs. Fire Technol 48:825–839. https://doi.org/10.1007/s10694-011-0247-4
- Naser MZ (2019) Fire Resistance Evaluation through Artificial Intelligence - A Case for Timber Structures. Fire Saf J 105:1–18. https://doi.org/10.1016/j.firesaf.2019.02.002
- Domingos P (2012) A few useful things to know about machine learning. Commun ACM. https://doi.org/10.1145/2347736. 2347755
- Shakya AM, Kodur VKR (2015) Response of precast prestressed concrete hollowcore slabs under fire conditions. Eng Struct. https://doi.org/10.1016/j.engstruct.2015.01.018
- Kodur VKR, Bhatt PP (2018) A numerical approach for modeling response of fiber reinforced polymer strengthened concrete slabs exposed to fire. Compos Struct 187:226–240. https://doi.org/10.1016/J.COMPSTRUCT.2017.12.051
- 37. Kohnke PC (2013) ANSYS. In: @ ANSYS, Inc.
- 38. Abaqus 6.13 (2013) Abaqus 6.13. Anal User's Guid Dassault Syst
- Franssen JM, Gernay T (2017) Modeling structures in fire with SAFIR®: Theoretical background and capabilities. J Struct Fire Eng. https://doi.org/10.1108/JSFE-07-2016-0010
- Golafshani EM, Behnood A (2018) Automatic regression methods for formulation of elastic modulus of recycled aggregate concrete. Appl Soft Comput J. https://doi.org/10.1016/j.asoc.2017. 12.030
- Sadowski Ł, Nikoo M, Nikoo M (2018) Concrete compressive strength prediction using the imperialist competitive algorithm. Comput Concr. https://doi.org/10.12989/cac.2018.22.4.355
- Alavi AH, Gandomi AH, Sahab MG, Gandomi M (2010) Multi expression programming: A new approach to formulation of soil classification. Eng Comput 26:111–118. https://doi.org/10.1007/ s00366-009-0140-7
- Mirjalili S, Lewis A (2015) Novel performance metrics for robust multi-objective optimization algorithms. Swarm Evol Comput. https://doi.org/10.1016/j.swevo.2014.10.005
- Mishra SK, Panda G, Majhi R (2014) A comparative performance assessment of a set of multiobjective algorithms for constrained portfolio assets selection. Swarm Evol Comput. https://doi.org/10.1016/j.swevo.2014.01.001
- 45. Schmidt MD, Lipson H (2010) Age-fitness pareto optimization
- Cremonesi P, Koren Y, Turrin R (2010) Performance of Recommender Algorithms on Top-N Recommendation Tasks Categories and Subject Descriptors. RecSys
- Laszczyk M, Myszkowski PB (2019) Survey of quality measures for multi-objective optimization: Construction of complementary set of multi-objective quality measures. Swarm Evol Comput 48:109–133. https://doi.org/10.1016/J.SWEVO.2019.04.001
- Willmott CJ, Matsuura K (2005) Advantages of the mean absolute error (MAE) over the root mean square error (RMSE) in assessing average model performance. Clim Res. https://doi.org/10.3354/cr030079
- Makridakis S (1993) Accuracy measures: theoretical and practical concerns. Int J Forecast. https://doi.org/10.1016/0169-2070(93)90079-3
- Ferreira C (2001) Gene Expression Programming: a New Adaptive Algorithm for Solving Problems. Ferreira, C (2001) Gene

- Expr Program a New Adapt Algorithm Solving Probl Complex Syst 13
- (2016) Handbook of Time Series Analysis, Signal Processing, and Dynamics
- Hyndman RJ, Koehler AB (2006) Another look at measures of forecast accuracy. Int J Forecast. https://doi.org/10.1016/j.ijfor ecast.2006.03.001
- Shcherbakov MV, Brebels A, Shcherbakova NL et al (2013) A survey of forecast error measures. World Appl Sci J. https://doi. org/10.5829/idosi.wasj.2013.24.itmies.80032
- Bain LJ (1967) Applied Regression Analysis. Technometrics. https://doi.org/10.1080/00401706.1967.10490452
- Armstrong JS, Collopy F (1992) Error measures for generalizing about forecasting methods: Empirical comparisons. Int J Forecast. https://doi.org/10.1016/0169-2070(92)90008-W
- Poli AA, Cirillo MC (1993) On the use of the normalized mean square error in evaluating dispersion model performance. Atmos Environ Part A, Gen Top. https://doi.org/10.1016/0960-1686(93) 90410-Z
- Smith G (1986) Probability and statistics in civil engineering. Collins, London
- Golbraikh A, Shen M, Xiao Z et al (2003) Rational selection of training and test sets for the development of validated QSAR models. J Comput Aided Mol Des 17:241–253. https://doi.org/ 10.1023/A:1025386326946
- Roy PP, Roy K (2008) On some aspects of variable selection for partial least squares regression models. QSAR Comb Sci 27:302–313. https://doi.org/10.1002/qsar.200710043
- 60. Frank I, Todeschini R (1994) The data analysis handbook
- Gandomi AH, Yun GJ, Alavi AH (2013) An evolutionary approach for modeling of shear strength of RC deep beams. Mater Struct Constr. https://doi.org/10.1617/s11527-013-0039-z
- 62. Cheng MY, Firdausi PM, Prayogo D (2014) High-performance concrete compressive strength prediction using Genetic Weighted Pyramid Operation Tree (GWPOT). Eng Appl Artif Intell. https://doi.org/10.1016/j.engappai.2013.11.014
- Alwanas AAH, Al-Musawi AA, Salih SQ et al (2019) Loadcarrying capacity and mode failure simulation of beam-column joint connection: Application of self-tuning machine learning model. Eng Struct. https://doi.org/10.1016/j.engstruct.2019.05. 048
- Chou JS, Tsai CF, Pham AD, Lu YH (2014) Machine learning in concrete strength simulations: Multi-nation data analytics. Constr Build Mater. https://doi.org/10.1016/j.conbuildmat.2014.09.054
- Sadat Hosseini A, Hajikarimi P, Gandomi M et al (2021) Genetic programming to formulate viscoelastic behavior of modified asphalt binder. Constr Build Mater. https://doi.org/10.1016/j. conbuildmat.2021.122954
- Nguyen TT, Pham Duy H, Pham Thanh T, Vu HH (2020) Compressive Strength Evaluation of Fiber-Reinforced High-Strength Self-Compacting Concrete with Artificial Intelligence. Adv Civ Eng. https://doi.org/10.1155/2020/3012139
- Sultana N, Zakir Hossain SM, Alam MS, et al (2020) Soft computing approaches for comparative prediction of the mechanical properties of jute fiber reinforced concrete. Adv Eng Softw 149:. https://doi.org/10.1016/j.advengsoft.2020.102887
- Willmott CJ (1981) On the validation of models. Phys Geogr. https://doi.org/10.1080/02723646.1981.10642213
- Nash JE, Sutcliffe JV (1970) River flow forecasting through conceptual models part I - A discussion of principles. J Hydrol. https://doi.org/10.1016/0022-1694(70)90255-6
- Gupta HV, Kling H, Yilmaz KK, Martinez GF (2009) Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling. J Hydrol. https://doi.org/10.1016/j.jhydrol.2009.08.003



- Knoben WJM, Freer JE, Woods RA (2019) Technical note: Inherent benchmark or not? Comparing Nash-Sutcliffe and Kling-Gupta efficiency scores. Hydrol Earth Syst Sci. https://doi.org/10.5194/hess-23-4323-2019
- Cheng MY, Chou JS, Roy AFV, Wu YW (2012) High-performance Concrete Compressive Strength Prediction using Time-Weighted Evolutionary Fuzzy Support Vector Machines Inference Model. Autom Constr. https://doi.org/10.1016/j.autcon. 2012.07.004
- Yaseen ZM, Deo RC, Hilal A et al (2018) Predicting compressive strength of lightweight foamed concrete using extreme learning machine model. Adv Eng Softw. https://doi.org/10.1016/j.adven gsoft.2017.09.004
- Yang L, Qi C, Lin X et al (2019) Prediction of dynamic increase factor for steel fibre reinforced concrete using a hybrid artificial intelligence model. Eng Struct. https://doi.org/10.1016/j.engst ruct.2019.03.105
- Qi C, Fourie A, Chen Q (2018) Neural network and particle swarm optimization for predicting the unconfined compressive strength of cemented paste backfill. Constr Build Mater. https:// doi.org/10.1016/j.conbuildmat.2017.11.006
- Chou J-S, Chiu C-K, Farfoura M, Al-Taharwa I (2010) Optimizing the Prediction Accuracy of Concrete Compressive Strength Based on a Comparison of Data-Mining Techniques. J Comput Civ Eng. https://doi.org/10.1061/(asce)cp.1943-5487.0000088
- Deepa C, SathiyaKumari K, Sudha VP (2010) Prediction of the Compressive Strength of High Performance Concrete Mix using Tree Based Modeling. Int J Comput Appl. https://doi.org/ 10.5120/1076-1406
- Erdal HI (2013) Two-level and hybrid ensembles of decision trees for high performance concrete compressive strength prediction. Eng Appl Artif Intell. https://doi.org/10.1016/j.engappai. 2013.03.014
- Yan K, Shi C (2010) Prediction of elastic modulus of normal and high strength concrete by support vector machine. Constr Build Mater. https://doi.org/10.1016/j.conbuildmat.2010.01.006
- Rafiei MH, Khushefati WH, Demirboga R, Adeli H (2017) Supervised Deep Restricted Boltzmann Machine for Estimation of Concrete. ACI Mater J 114:. https://doi.org/10.14359/51689 560
- Yan K, Xu H, Shen G, Liu P (2013) Prediction of Splitting Tensile Strength from Cylinder Compressive Strength of Concrete by Support Vector Machine. Adv Mater Sci Eng. https://doi.org/10.1155/2013/597257
- Anoop Krishnan NM, Mangalathu S, Smedskjaer MM et al (2018) Predicting the dissolution kinetics of silicate glasses using machine learning. J Non Cryst Solids. https://doi.org/10.1016/j. jnoncrysol.2018.02.023
- Okuyucu H, Kurt A, Arcaklioglu E (2007) Artificial neural network application to the friction stir welding of aluminum plates. Mater Des. https://doi.org/10.1016/j.matdes.2005.06.003
- 84. Lim CH, Yoon YS, Kim JH (2004) Genetic algorithm in mix proportioning of high-performance concrete. Cem Concr Res. https://doi.org/10.1016/j.cemconres.2003.08.018
- Haghdadi N, Zarei-Hanzaki A, Khalesian AR, Abedi HR (2013)
   Artificial neural network modeling to predict the hot deformation behavior of an A356 aluminum alloy. Mater Des. https://doi.org/10.1016/j.matdes.2012.12.082
- Golafshani EM, Behnood A (2019) Estimating the optimal mix design of silica fume concrete using biogeography-based programming. Cem Concr Compos 96:95–105. https://doi.org/10. 1016/J.CEMCONCOMP.2018.11.005
- Naser MZ (2018) Deriving temperature-dependent material models for structural steel through artificial intelligence. Constr Build Mater 191:56–68. https://doi.org/10.1016/J.CONBUILDMAT. 2018.09.186

- Naser MZ (2019) Properties and material models for modern construction materials at elevated temperatures. Comput Mater Sci 160:16–29. https://doi.org/10.1016/J.COMMATSCI.2018. 12.055
- Mousavi SM, Aminian P, Gandomi AH et al (2012) A new predictive model for compressive strength of HPC using gene expression programming. Adv Eng Softw. https://doi.org/10. 1016/j.advengsoft.2011.09.014
- Gandomi AH, Alavi AH, Sahab MG (2010) New formulation for compressive strength of CFRP confined concrete cylinders using linear genetic programming. Mater Struct Constr. https:// doi.org/10.1617/s11527-009-9559-y
- Mollahasani A, Alavi AH, Gandomi AH (2011) Empirical modeling of plate load test moduli of soil via gene expression programming. Comput Geotech. https://doi.org/10.1016/j. compgeo.2010.11.008
- 92. Erdal HI, Karakurt O, Namli E (2013) High performance concrete compressive strength forecasting using ensemble models based on discrete wavelet transform. Eng Appl Artif Intell. https://doi.org/10.1016/j.engappai.2012.10.014
- Pang B, Lee L, Vaithyanathan S (2002) Thumbs up? In: Proceedings of the ACL-02 conference on Empirical methods in natural language processing - EMNLP '02
- Galdi P, Tagliaferri R (2017) Data Mining: Accuracy and Error Measures for Classification and Prediction. In: Encyclopedia of Bioinformatics and Computational Biology
- Valença J, Gonçalves LMS, Júlio E (2013) Damage assessment on concrete surfaces using multi-spectral image analysis. Constr Build Mater. https://doi.org/10.1016/j.conbuildmat.2012.11. 061
- Huang H, Burton HV (2019) Classification of in-plane failure modes for reinforced concrete frames with infills using machine learning. J Build Eng. https://doi.org/10.1016/j.jobe.2019. 100767
- 97. Azimi SM, Britz D, Engstler M et al (2018) Advanced steel microstructural classification by deep learning methods. Sci Rep. https://doi.org/10.1038/s41598-018-20037-5
- Hore S, Chatterjee S, Sarkar S, et al (2016) Neural-based prediction of structural failure of multistoried RC buildings. Struct Eng Mech. https://doi.org/10.12989/sem.2016.58.3.459
- Bhowan U, Johnston M, Zhang M (2012) Developing new fitness functions in genetic programming for classification with unbalanced data. IEEE Trans Syst Man, Cybern Part B Cybern. https:// doi.org/10.1109/TSMCB.2011.2167144
- Boughorbel S, Jarray F, El-Anbari M (2017) Optimal classifier for imbalanced data using Matthews Correlation Coefficient metric. PLoS ONE. https://doi.org/10.1371/journal.pone.0177678
- Tharwat A (2018) Classification assessment methods. Appl. Comput. Informatics
- 102. Caruana R, Niculescu-Mizil A (2004) Data mining in metric space: an empirical analysis of supervised learning performance criteria. In: KDD-2004 - Proceedings of the Tenth ACM SIG-KDD International Conference on Knowledge Discovery and Data Mining
- Jurman G, Riccadonna S, Furlanello C (2012) A comparison of MCC and CEN error measures in multi-class prediction. PLoS ONE. https://doi.org/10.1371/journal.pone.0041882
- Powers DMW (2011) Evaluation: From Precision, Recall and F-Factor to ROC, Informedness, Markedness & Correlation. J Mach Learn Technol. 10.1.1.214.9232
- Bradley AP (1997) The use of the area under the ROC curve in the evaluation of machine learning algorithms. Pattern Recognit. https://doi.org/10.1016/S0031-3203(96)00142-2
- 106. Hanley JA, McNeil BJ (1982) The meaning and use of the area under a receiver operating characteristic (ROC) curve. Radiology. https://doi.org/10.1148/radiology.143.1.7063747



- Zhang Y, Burton HV, Sun H, Shokrabadi M (2018) A machine learning framework for assessing post-earthquake structural safety. Struct Saf. https://doi.org/10.1016/j.strusafe.2017.12.001
- 108. Davis J, Goadrich M (2006) The relationship between Precision-Recall and ROC curves. In: Proceedings of the 23rd international conference on Machine learning - ICML '06
- Bi, J.; Bennett KPP (2003) Regression Error Characteristic Curves. Proc Twent Int Conf Mach Learn
- Zhang M, Smart W (2006) Using Gaussian distribution to construct fitness functions in genetic programming for multiclass object classification. Pattern Recognit Lett. https://doi.org/10.1016/j.patrec.2005.07.024
- Kocher M, Savoy J (2017) Distance measures in author profiling. Inf Process Manag. https://doi.org/10.1016/j.ipm.2017.04.004
- Patel BV (2012) Content Based Video Retrieval Systems. Int J UbiComp. https://doi.org/10.5121/iju.2012.3202
- Giusti R, Batista GEAPA (2013) An empirical comparison of dissimilarity measures for time series classification. In: Proceedings 2013 Brazilian Conference on Intelligent Systems, BRACIS 2013
- 114. Vuk M, Curk T (2006) ROC Curve , Lift Chart and Calibration Plot. Metod Zv
- Brodersen KH, Ong CS, Stephan KE, Buhmann JM (2010) The balanced accuracy and its posterior distribution. In: Proceedings

   International Conference on Pattern Recognition
- Cohen J (1960) A Coefficient of Agreement for Nominal Scales.
   Educ Psychol Meas. https://doi.org/10.1177/001316446002000
   104
- Artstein R, Poesio M (2008) Inter-coder agreement for computational linguistics. Comput. Linguist.
- Destercke S (2014) Multilabel Prediction with Probability Sets:
   The Hamming Loss Case. In: Communications in Computer and Information Science
- Czajkowski M, Kretowski M (2019) Decision Tree Underfitting in Mining of Gene Expression Data. An Evolutionary Multi-Test Tree Approach. Expert Syst Appl. https://doi.org/10.1016/J. ESWA.2019.07.019

- Devarriya D, Gulati C, Mansharamani V, et al (2019) Unbalanced Breast Cancer Data Classification Using Novel Fitness Functions in Genetic Programming. Expert Syst Appl 112866. https://doi. org/10.1016/J.ESWA.2019.112866
- Bhaskar H, Hoyle DC, Singh S (2006) Machine learning in bioinformatics: A brief survey and recommendations for practitioners.
   Comput Biol Med. https://doi.org/10.1016/j.compbiomed.2005. 09.002
- Kohavi R (1995) A study of cross-validation and bootstrap for accuracy estimation and model selection. Proc 14th Int Jt Conf Artif Intell - Vol 2
- 123. Alavi AH, Gandomi AH (2011) Prediction of principal groundmotion parameters using a hybrid method coupling artificial neural networks and simulated annealing. Comput Struct. https://doi. org/10.1016/j.compstruc.2011.08.019
- Kingston GB, Maier HR, Lambert MF (2005) Calibration and validation of neural networks to ensure physically plausible hydrological modeling. J Hydrol. https://doi.org/10.1016/j.jhydr ol.2005.03.013
- Kuo YL, Jaksa MB, Lyamin AV, Kaggwa WS (2009) ANN-based model for predicting the bearing capacity of strip footing on multi-layered cohesive soil. Comput Geotech. https://doi.org/ 10.1016/j.compgeo.2008.07.002
- Caruana R, Niculescu-Mizil A (2006) An empirical comparison of supervised learning algorithms. In: ACM International Conference Proceeding Series. ACM Press, New York, USA, pp 161–168
- 127. Williams N, Zander S, Armitage G (2006) A preliminary performance comparison of five machine learning algorithms for practical IP traffic flow classification. Comput Commun Rev. https://doi.org/10.1145/1163593.s1163596

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

