

# Machine learning can predict setting behavior and strength evolution of hydrating cement systems

Tandr  Oey (<sup>a\*</sup>), Scott Jones (<sup>b</sup>), Jeffrey W. Bullard (<sup>b</sup>), and Gaurav Sant (<sup>a,c,d,e</sup>)

## Abstract

Setting and strength development of ordinary portland cement (OPC) binders is a complex process that involves multiple interacting chemical reactions, which result in the formation of a solid microstructure. A long-standing yet elusive goal of the cementitious materials community has been to establish a basis for prediction of the properties and performance of concrete using knowledge of the chemical and physical attributes of its components – OPC, sand, stone, water, and chemical admixtures – together with the environmental conditions under which they react. Machine learning provides a *data-driven* basis for the estimation of properties, and has recently been applied to estimate the 28 d (compressive) strength of concrete simply from knowledge of its mixture proportions [1]. Building on this success, the current work uses a diverse dataset of different ASTM C150 cements, the chemical composition and other attributes of which have been measured. Machine learning (ML) estimators were trained with this dataset to estimate both paste setting time and mortar strength development as a function of the OPC composition and fineness. The ML estimation errors are typically similar to or lower than the measurement repeatability of the relevant ASTM test methods. ML therefore can be used to estimate the influence of binder composition and fineness on the engineering properties of cementitious systems. This creates new opportunities to apply data intensive methods to optimize concrete formulations under multiple constraints of cost, CO<sub>2</sub> impact, and performance attributes.

**Keywords:** cement composition, fineness, strength, setting, machine learning

---

<sup>a</sup> Laboratory for the Chemistry of Construction Materials (LC<sup>2</sup>), Department of Civil and Environmental Engineering, University of California, Los Angeles, CA, USA

<sup>b</sup> Engineering Laboratory, National Institute of Standards and Technology, Gaithersburg, MD, USA

<sup>c</sup> Institute for Carbon Management, University of California, Los Angeles, CA, USA

<sup>d</sup> Department of Materials Science and Engineering, University of California, Los Angeles, CA, USA

<sup>e</sup> California Nanosystems Institute, University of California, Los Angeles, CA, USA

29 **1.0 Introduction**

30 The hydration of ordinary portland cement (PC) entails multiple concurrent chemical reactions  
31 [2]. These reactions cause extensive changes in phase assemblage and microstructure, which in  
32 turn determine the time-dependent evolution of concrete properties and performance such as  
33 setting time and compressive strength. Mature (28 d) compressive strength is the metric most  
34 commonly used to specify and qualify a concrete for structural design [3]. Multi-scale  
35 simulations suggest the need to couple microstructural and mechanical models as a means to  
36 predict time-dependent mechanical properties [4]. However, these approaches are still severely  
37 limited by gaps in knowledge of OPC's hydration process and its constituent mechanisms, and  
38 are generally unable to forecast the evolution of properties and performance unless they are  
39 experimentally, and narrowly, calibrated to the specific system of interest [5].

40  
41 In the absence of knowledge needed to predict cement hydration rates and associated changes in  
42 properties, data-driven machine learning (ML) methods offer an attractive, mechanism-agnostic  
43 approach for estimating engineering properties such as the 28 d compressive strength of concrete  
44 [6–15]. Young *et al.* [1] have recently demonstrated that ML can, when trained on enough data,  
45 make reasonable estimations of the 28 d compressive strength of field-produced concretes as a  
46 function of its attributes such as water-to-cement mass ratio (w/c), aggregate content, and the  
47 content and type of mineral and chemical admixtures. Such results demonstrate the potential of  
48 ML approaches for predicting concrete performance because the data that were used therein were  
49 obtained for concrete produced under the *relatively uncontrolled* conditions of diverse  
50 construction sites. The predictions could likely be made even more accurate by including site-  
51 specific variables such as temperature and humidity changes with time. However, the study was  
52 limited to concretes produced with Type I/II PC.

53  
54 To supplement and extend existing models, the current study takes another step toward truly  
55 predictive models of concrete properties by applying ML methods to estimate the effects of OPC  
56 characteristics, such as chemical composition and fineness, on target performance characteristics  
57 such as paste setting time and mortar compressive strength. In addition, a tentative lower bound  
58 on the number of data records that are required for future estimation of other concrete properties  
59 is established. Special focus is paid to identify potential technical barriers faced by ML methods  
60 to identifying *general trends* among thousands of data points and, more importantly, to  
61 accurately predict the properties of any *one* material of interest.

62  
63 **2.0 Background and Methods**

64  
65 **2.1 Machine learning algorithms**

66 Young *et al.* showed that bootstrap-aggregated (or bagged) decision tree ensembles can  
67 accurately estimate the 28 d compressive strength of concrete when trained on large datasets with  
68 potentially high inherent variability [1]. These rule-based estimators identify logical splits in  
69 data, partitioning the input space into a tree of decision nodes that are traversed until arriving at a  
70 final prediction of the target, called a leaf node. A simple operation, such as the multiplication of  
71 the input by a constant, produces the output estimation from each leaf node. A collection, or  
72 ensemble of trees are constructed, each tree being trained on different data sets and attributes,  
73 and their results are then averaged to produce the final prediction of the target [16]. This study  
74 focuses on three different decision tree ensembles because of their ability to estimate field

75 concrete compressive strength [1]. The first method is a bagged\* tree ensemble, which bootstrap  
76 samples  $n$  different subsets of the training data with replacement to train  $n$  trees. Other than the  
77 random sampling from the training data, the method is deterministic in the sense that the decision  
78 nodes are chosen from among all attributes using a deterministic function such as information  
79 gain or Gini index [17]. In addition, the threshold value for splitting at a decision node is chosen  
80 to be that which optimizes that deterministic function. The second method, a random forest  
81 ensemble, differs from the bagged tree ensemble in that it selects the attribute chosen for each  
82 decision node from among a randomly chosen small subset of the attributes. The third method,  
83 called extra† trees, is the same as a random forest except that the threshold value for splitting a  
84 decision node is also chosen at random instead of being prescribed by optimization of a  
85 thresholding function [18]. Other ML estimators besides these three tree ensembles were also  
86 examined, including basic linear regression and K-nearest neighbor (K-NN) regression [19–21].  
87 The tree ensembles provided the highest prediction accuracy for every attribute, although the  
88 results of the other regression methods are also shown for comparison. All the algorithms used  
89 for estimator construction are regressors from the scikit-learn library, and can be accessed and  
90 downloaded, along with their documentation, at <http://scikit-learn.org/stable/> [19].

## 91 92 **2.2 Data collection and preprocessing**

93 Two datasets were utilized. The first dataset was provided by the Cement and Concrete  
94 Reference Laboratory (CCRL) Proficiency Sample Program, which issues four OPCs each year  
95 for comprehensive physical and chemical testing by nearly 200 different laboratories. This  
96 dataset consists of measurements of 48 attributes of a given OPC sample (see Table 1), as  
97 established by ASTM test methods [22]. The second dataset is a compilation of different industry  
98 survey data supplied by the Portland Cement Association (PCA) and the National Institute of  
99 Standards and Technology (NIST), formerly the National Bureau of Standards (NBS). This  
100 dataset comprises 2211 different PCs characterized by an unknown number of testing institutions  
101 using standard test methods. It also includes the averages‡ of 19 of the 48 attributes for each of  
102 the CCRL cements (marked in bold in Table 1). Two other attributes, normal consistency and  
103 final setting time, were also reported in the majority of records available, and so were also  
104 considered in this study (italicized in Table 1). The bolded entries in the “Chemical Tests”  
105 column of Table 1 were used as inputs to the final ML estimators, along with Blaine fineness,  
106 while the bolded and italicized entries in the “Physical Tests” column of Table 1, with the  
107 exception of Blaine fineness, were used as targets for ML prediction using these estimators.

108  
109 Prior to use as inputs and targets in the machine learning estimators, the data were preprocessed  
110 to remove obvious errors and to ensure they would be compatible with all the ML algorithms  
111 used. First, on an attribute-by-attribute basis, unphysical or meaningless values were deleted.  
112 Among these were percentages outside the range of 0 % to 100 %§ and unphysical values such as  
113 negative setting time or compressive strength. Second, a filter was applied to each attribute to  
114 delete any outliers, which we defined according Chauvenet’s criterion [23] as more than four

---

\* The term “bagged” is a portmanteau of the terms “bootstrap” and “aggregated”.

† The term “extra” is a portmanteau of the terms “extremely” and “randomized”.

‡ Use of averages was necessary to ensure that no cement was over-represented in the input data to ML models, as this is known to negatively impact ML estimator performance.

§ Any percentage values in excess of 100 % or below 0 % were retained only if they were physically meaningful. For example, negative percentages in autoclave expansion measurements correspond to shrinkage.

115 standard deviations from the mean of that attribute across all cements. The mean(s) were  
 116 recalculated after those outliers were removed and the filter was reapplied, the process being  
 117 repeated until no more outliers were identified. Less than 0.05 % of the data were discarded by  
 118 this filtering for any given attribute, and the process of omitting outliers required only three  
 119 iterations. Afterward, duplicate records (that is, identical cements) were deleted and any missing  
 120 attributes were replaced by mean imputation, setting each missing value to the mean for the  
 121 appropriate attribute as determined using data from the other cements. This is the simplest of all  
 122 methods of data imputation, used in situations when data are missing completely at random, i.e.,  
 123 when the absence of a value is unrelated to the state of the system or values of other variables.  
 124

125 Of the two datasets, that from the CCRL contains the greater number of records, nearly 31 000,  
 126 and has a more comprehensive list of potential attributes to be used as inputs or targets for ML  
 127 estimators. However, that dataset is also missing more data, contains many more duplicates  
 128 (consisting of only about 200 unique cements), and consequently was unable to train any ML  
 129 estimators as accurately as the composite survey dataset. The CCRL data were incorporated in  
 130 the composite survey dataset, however, by using the mean value of each attribute for each  
 131 cement instead of the individual records. Randomly shuffling the order of data records proved  
 132 essential for effectively training the ML estimators regardless of the algorithm used. This  
 133 indicates that the ranges of attribute values are not homogeneously distributed across the  
 134 different surveys in the compilation, and that leaving the data grouped by survey alone  
 135 introduces an inadvertent bias in the sampling of input attributes toward one particular study.  
 136 Therefore, random shuffling as implemented herein is an effective and necessary way to  
 137 ameliorate that artifact.  
 138

139 **Table 1:** The cement attributes provided in the datasets, and the ASTM standards [22] (in square brackets) used to  
 140 measure them. The boldfaced entries are reported consistently for nearly all cements in the full dataset, and italicized  
 141 entries are reported in at least 50 % of the records in the dataset. Other entries were not consistently reported and  
 142 were excluded from inputs to ML estimators. All boldfaced and italicized entries listed under “Physical Tests,” with  
 143 the exception of Blaine Fineness, were utilized as target attributes in this study, and as such were also excluded from  
 144 inputs to ML estimators. All other entries that were excluded from inputs to ML estimators were verified to be of  
 145 minimal importance to estimator performance, as outlined in Section 3.2.

<b>Chemical Tests</b>	<b>Physical Tests</b>
<b>SiO<sub>2</sub> (mass %) [C114]</b>	<i>Paste Normal Consistency (%) [C187]</i>
<b>Al<sub>2</sub>O<sub>3</sub> (mass %) [C114]</b>	<b>Vicat Paste Initial Set (minutes) [C191]</b>
<b>Fe<sub>2</sub>O<sub>3</sub> (mass %) [C114]</b>	<i>Vicat Paste Final Set (minutes) [C191]</i>
<b>CaO (mass %) [C114]</b>	Gillmore Initial Set (minutes) [C266]
<b>C<sub>3</sub>S (mass %) [C150]</b>	Gillmore Final Set (minutes) [C266]
<b>C<sub>2</sub>S (mass %) [C150]</b>	False Set (%) [C451]
<b>C<sub>3</sub>A (mass %) [C150]</b>	Autoclave Expansion (%) [C151]
<b>C<sub>4</sub>AF (mass %) [C150]</b>	Air Content (%) [C185]
<b>Free CaO (mass %) [C114]</b>	Air Content Mixing Water (%) [C185]
<b>MgO (mass %) [C114]</b>	Air Content Mixture Flow (%) [C185]
<b>SO<sub>3</sub> (mass %) [C114]</b>	<b>3 Day Mortar Compressive Strength (MPa) [C109]</b>
<b>Na<sub>2</sub>O (mass %) [C114]</b>	<b>7 Day Mortar Compressive Strength (MPa) [C109]</b>
<b>K<sub>2</sub>O (mass %) [C114]</b>	<b>28 Day Mortar Compressive Strength (MPa) [C109]</b>
<b>Loss on Ignition (mass %) [C114]</b>	Mortar Compressive Strength Mixture Flow (%) [C109]
<b>Insoluble Residue (mass %) [C114]</b>	<b>Blaine Fineness (m<sup>2</sup>/kg) [C204]</b>

Carbon Dioxide (mass %) [C114]	Wagner Fineness (m <sup>2</sup> /kg) [C115]
Limestone (mass %) [C114]	Sieve Fineness (% passing) [C430]
ZnO (mass %) [C114]	0 Day Heat of Solution (cal/g) [C186]
Mn <sub>2</sub> O (mass %) [C114]	7 Day Heat of Solution (cal/g) [C186]
P <sub>2</sub> O <sub>5</sub> (mass %) [C114]	28 Day Heat of Solution (cal/g) [C186]
TiO <sub>2</sub> (mass %) [C114]	7 Day Heat of Hydration (cal/g) [C186]
Cl (mass %) [C114]	28 Day Heat of Hydration (cal/g) [C186]
	Mortar Bar Expansion (%) [C1038]
	Mortar Bar Mixing Water (%) [C1038]
	Mortar Bar Flow [C1038]

146

### 147 **2.3 Estimator optimization**

148 Numerous machine learning estimators were constructed and applied to predict initial set  
 149 (minutes), 3 d compressive strength (MPa), 7 d compressive strength (MPa), and 28 d  
 150 compressive strength (MPa). These targets were chosen because the first three affect the  
 151 scheduling of construction operations, and the 28 d strength is both an input for structural design  
 152 and a specification criterion. Each estimator was trained and tested on the combined datasets  
 153 with the performance of each estimator being evaluated using several error metrics. Both training  
 154 and testing were conducted on different portions of data using a standard low-bias resampling  
 155 procedure called k-Fold Cross-Validation\* [24,25]. The data records were randomly split into k =  
 156 10 “folds,” nine of which were used to train the estimator, and one of which was used to evaluate  
 157 the estimator after training. The process was then repeated nine additional times, each time using  
 158 a different fold as the test set, and the remaining nine folds as the training set.

159

160 The estimators used in this study are sensitive to the magnitude of the attributes in the sense that  
 161 they will be biased to assign more importance to attributes with inherently greater values. For  
 162 example, merely changing the units of Blaine fineness of the powder from m<sup>2</sup>/kg to cm<sup>2</sup>/g  
 163 increases the numerical value by a factor of ten and can influence the accuracy of the estimators  
 164 even though the physical data are the same. To address this kind of artifact, after the training and  
 165 testing sets were identified and separated, the data for each attribute were rescaled to a standard  
 166 normal distribution (mean = 0, variance = 1). This step was taken after the separation of the  
 167 training and testing sets to avoid data leakage (i.e., the unintentional passing of information  
 168 about the test set to the training set) which could potentially happen if the combined testing and  
 169 training data were rescaled together.

170

171 Estimator optimization was performed by determining extremal values of one of three objective  
 172 functions that characterize the overall fidelity of the predictions to the actual values in the testing  
 173 set. The objective functions are the root mean square error (RMSE), the mean absolute  
 174 percentage error (MAPE), and the coefficient of determination (R<sup>2</sup>):

175

176

177

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (P_i - A_i)^2}{n}} \quad (1)$$

---

\* Cross-validation is necessary to evaluate how machine-learning estimators are likely to perform when making predictions on previously unseen data: a portion of the data are taken as a training set and used to train and optimize the model, and the remainder of the data are withheld as a testing set to evaluate the model’s performance.

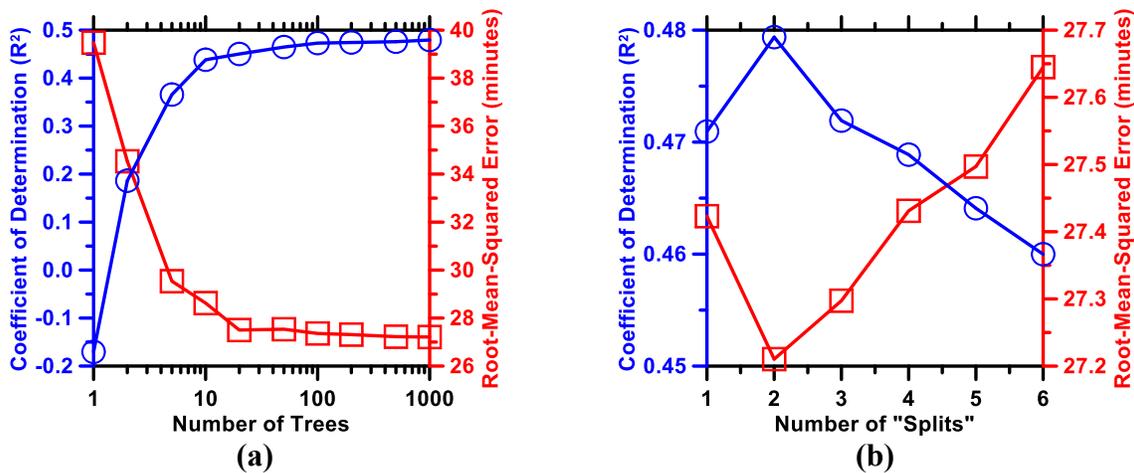
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203

$$\text{MAPE} = \frac{100}{n} \sum_{i=1}^n \left| \frac{P_i - A_i}{A_i} \right| \quad (2)$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (P_i - \bar{A})^2}{\sum_{i=1}^n (A_i - \bar{A})^2} \quad (3)$$

where  $n$  is the number of records in the testing set,  $P_i$  and  $A_i$  are the predicted and actual target value of the  $i^{\text{th}}$  record in the testing set, respectively, and  $\bar{A}$  is the arithmetic average of the actual target values. RMSE and MAPE indicate the average departure of estimated values from actual values, whereas  $R^2$  is the fraction of the variance of the target values that is predictable from the attributes using the model. As described in Section 3.1, low RMSE and MAPE values may still be achieved even when the data are relatively scattered and the  $R^2$  value is low. This has also been observed previously [1] and suggests that an error-based metric such as MAPE is a better test of estimator performance than  $R^2$  because it can be compared more directly with the acceptable range of physical test values for attributes such as setting time or strength [22].

Each of the machine learning estimators were finalized by optimizing their estimation performance via hyperparameter tuning. This procedure varied both the number of trees used in random forest estimators and the number of attributes considered per tree split when partitioning the input space. The results of this hyperparameter tuning, shown in Figure 1, indicate that estimator performance improves only marginally beyond a certain number of trees. Consequently, the final estimators reported here employ only 100 trees to avoid over-fitting of the training data, and the extremely-random forest estimators employ only two attributes per “split,” for similar reasons. These fully-optimized ML estimators are a substantial improvement (roughly a two-fold reduction in MAPE) over prior work [1], by merit of their consideration of cement composition.



**Figure 1:** The results of a representative parameter tuning exercise for the extremely random forest estimators constructed to estimate initial setting time, showing: (a) A plateau in estimator performance with increasing number of trees (i.e., in each case using two attributes to determine each partitioning of the input space), and (b) A modest optimum of two splits is observed when using 1000 trees.

204  
205  
206

### 3.0 Results and discussion

### 207 3.1 Estimation accuracy for a given target is comparable to ASTM repeatability limits

208 Among the algorithms examined, ensembles of decision trees consistently produced the lowest  
 209 errors, as shown in Table 2. Of the tree ensembles, the extra trees estimator most accurately  
 210 estimated every primary attribute as measured by MAPE. The error metrics are not much greater  
 211 in magnitude than the reported repeatability of the corresponding ASTM test methods, reported  
 212 as a coefficient of variation, though there is no standard ML error parameter that would enable  
 213 more direct comparisons [22]. For example, the MAPE for 7 d compressive strength predictions  
 214 by extra trees estimators is 6.58 %, less than twice the single-operator coefficient of variation of  
 215 the measurement using ASTM C109 (3.8%). Similarly, the variability in initial set time for that  
 216 estimator, 25.5 minutes, is considerably less than the acceptable range of two successive  
 217 measurements using ASTM C191 (34 minutes). This suggests that, for cement compositions  
 218 covered by ASTM C150, ensemble machine learning approaches may reliably estimate the  
 219 *average properties* and performance of paste / mortar formulations nearly as well or better than  
 220 they can be repeatably measured in the lab.

221  
 222 **Table 2:** The results of 10-fold cross-validation using the following error metrics: root mean square error (RMSE),  
 223 coefficient of determination ( $R^2$ ), and mean absolute percentage error (MAPE). The input attributes were SiO<sub>2</sub> (mass  
 224 %), Al<sub>2</sub>O<sub>3</sub> (mass %), Fe<sub>2</sub>O<sub>3</sub> (mass %), CaO (mass %), SO<sub>3</sub> (mass %), and Blaine fineness (m<sup>2</sup>/kg), as determined by  
 225 attribute importance in Section 3.2.

Target Attributes for Ten-Fold Cross-Validation												
Estimator	Initial Set Time <sup>b</sup>			3 Day Strength <sup>c</sup>			7 Day Strength <sup>c</sup>			28 Day Strength <sup>c</sup>		
	RMSE (min)	R <sup>2</sup>	MAPE (%)	RMSE (MPa)	R <sup>2</sup>	MAPE (%)	RMSE (MPa)	R <sup>2</sup>	MAPE (%)	RMSE (MPa)	R <sup>2</sup>	MAPE (%)
Linear	29.6	0.392	17.7	3.26	0.676	9.01	3.59	0.573	7.90	4.01	0.305	7.06
K-NN <sup>a</sup>	27.8	0.437	15.9	3.25	0.691	8.32	3.48	0.614	7.27	3.77	0.394	6.35
Decision Tree Ensemble Estimators:												
Bagged	26.2	0.524	15.0	2.79	0.766	7.29	3.18	0.668	6.67	3.50	0.489	5.91
Random	25.6	0.541	14.9	2.82	0.763	7.35	3.15	0.674	6.68	3.48	0.497	5.87
<b>Extra</b>	<b>25.5</b>	<b>0.547</b>	<b>14.7</b>	<b>2.82</b>	<b>0.762</b>	<b>7.29</b>	<b>3.14</b>	<b>0.675</b>	<b>6.58</b>	<b>3.44</b>	<b>0.506</b>	<b>5.79</b>
Boosted <sup>a</sup>	29.0	0.417	17.7	3.44	0.646	10.0	3.63	0.567	8.19	3.89	0.368	6.81
Gradient <sup>a</sup>	26.9	0.495	15.7	2.89	0.749	7.74	3.30	0.642	6.93	3.59	0.460	6.13

226 <sup>a</sup> K-nearest neighbors, boosted decision trees, and gradient boosted decision trees were also used, among other  
 227 estimators (not shown), as they are likely to perform similarly to bagged decision trees. None performed better  
 228 for these target attributes. For details regarding the implementation, see <http://scikit-learn.org/stable/>.

229 <sup>b</sup> ASTM C191.

230 <sup>c</sup> ASTM C109.

231

### 232 3.2 Higher errors for late-age strength suggest missing data attributes

233 Table 2 shows that estimator performance for predicting compressive strength is progressively  
 234 poorer at later ages, regardless of the estimator used. For example, the RMSE of the extremely  
 235 randomized forest estimator increases from 2.82 MPa at 3 d to 3.14 MPa and 3.44 MPa at 7 d  
 236 and 28 d, respectively. Despite the somewhat poorer estimator performance for 28 d strength  
 237 compared to earlier times, both the MAPE and RMSE for 28 d strength estimates are modestly  
 238 better than those determined by Young *et al.* [1] for industrially produced concretes using similar  
 239 estimators, likely due to more detailed knowledge of mixture and material characteristics in the  
 240 current study (cement composition, fineness). In any case, the greater errors at later ages may  
 241 indicate that the available datasets are missing some important attributes that influence  
 242 compressive strength at later ages.

243

244 One possible reason for this decrease in accuracy at later ages may be inconsistent or poorly  
 245 controlled curing conditions in practice, the effects of which would become progressively more  
 246 important with time. It is impossible to assess the likelihood of that possibility based on the data  
 247 alone, however, because there are no requirements in ASTM C109 to report the imposed degree  
 248 of control over curing temperature or moisture conditions. A second possible reason for  
 249 increased error is air entrainment in some subset of the measurements, given that ASTM C109  
 250 allows the user to decide whether or not the sample will contain entrained air – macroscopic air  
 251 voids stabilized by chemical admixtures to improve freeze-thaw resistance – requiring a lower  
 252 water-cement mass ratio (w/c) of 0.460 than the value of 0.485 required for samples without air  
 253 entrainment. Finally, differences in water content may play a significant role in poor estimator  
 254 performance for initial setting time measured using ASTM C191, wherein the mixture must be  
 255 prepared with “normal consistency” as measured by ASTM C187, which is the empirically  
 256 determined water content required to achieve a prescribed paste stiffness after 30 s of mixing  
 257 with 0.65 kg of cement powder (varying from about 22 % to 30 % of the powder mass among  
 258 different PCs). Therefore, ML estimation of normal consistency has also been investigated, as it  
 259 may serve as a proxy for w/c and is available in some of the compiled survey data.

261 **Table 3:** Results of 10-fold cross-validation for the final machine learning estimators of secondary targets with  
 262 partial data records, evaluated using the same error metrics given in Table 1. The best-performing estimator (lowest  
 263 MAPE) is marked in bold. The number of available data points used in each estimator is also reported.

Target Attribute for Ten-Fold Cross-Validation								
Estimator	Final Set				Normal Consistency			
	RMSE (min)	R <sup>2</sup>	MAPE (%)	Data Points	RMSE (%)	R <sup>2</sup>	MAPE (%)	Data Points
Linear	60.1	0.422	18.4	1144	1.04	0.292	2.96	1447
K-NN	59.6	0.432	17.6	1144	0.935	0.427	2.29	1447
<b>Trees:</b>								
Bagged	55.5	0.505	16.6	1144	0.920	0.446	2.28	1447
Random	55.5	0.505	16.6	1144	0.935	0.427	2.29	1447
<b>Extra</b>	<b>54.7</b>	<b>0.513</b>	<b>16.4</b>	1144	<b>0.894</b>	<b>0.471</b>	<b>2.23</b>	1447
Boosted	57.6	0.461	17.8	1144	1.11	0.193	3.31	1447
Gradient	57.9	0.461	17.5	1144	0.999	0.358	2.49	1447

264 <sup>a</sup>ASTM C191.

265 <sup>b</sup>ASTM C187.

267 **3.3 Estimation of secondary targets suggests a limited ability to account for missing attributes**

268 Among the other attributes in the dataset besides initial set and compressive strength, both  
 269 normal consistency and final setting time were reported frequently enough to construct viable  
 270 estimators. Estimators for these secondary targets, results of which are given in Table 3, were  
 271 indeed about as accurate as those for primary targets in Table 2. However, in contrast to the  
 272 primary targets, the errors in estimating normal consistency are significantly higher than the  
 273 tolerances listed in its associated ASTM C187 test method. Nevertheless, the normal consistency  
 274 estimators have the lowest MAPE of any estimator used in this study. ASTM C187 uses OPC  
 275 pastes prepared with normal consistency, so the estimator’s ability to capture the dependence of  
 276 normal consistency on composition and fineness may explain why ML estimators are able to  
 277 predict initial and final setting times from those same attributes despite the fact that the w/c used  
 278 can be different for each cement. In other words, cement details such as fineness are able to at  
 279 least *somewhat* capture this indicator of “water demand” of a cement, but there likely are other  
 280 powder characteristics – perhaps microscale texture or grinding aid type or dose – that affect

281 normal consistency but are currently not being measured by standard test methods. This example  
282 highlights both a limitation of, and an opportunity for, ML methods: they can estimate certain  
283 aspects of concrete performance from routinely collected data, but they can also identify other  
284 performance attributes, the systematic estimation of which requires additional or perhaps  
285 qualitatively different material characterization. Similarly, as taken up in the next section, it is  
286 helpful for understanding to identify which currently-measured attributes contribute most  
287 strongly to the quality of ML estimations of different targets.

288

### 289 ***3.4 Selective omission identifies six attributes necessary for estimation of set and strength***

290 One can evaluate the relative importance of the different attributes in determining estimator  
291 performance in predicting the primary targets (initial set, compressive strength) by eliminating  
292 them one at a time from the training set. The corresponding increase in MAPE was used as a  
293 quantitative measure of attribute importance, as shown in Figure 2(a). Unsurprisingly, cement  
294 fineness is by far the most influential input attribute, followed by the oxides of sulfur, calcium,  
295 aluminum, silicon, and iron. Similar attribute rankings were obtained for all targets estimated.  
296 This is reassuring because (i) available surface area is well known to be a key factor that affects  
297 cement reaction rates and water demand, (ii) calcium and aluminum bearing cement phases such  
298 as tricalcium silicate ( $C_3S^*$ ) and tricalcium aluminate ( $C_3A$ ), are known to be the most reactive  
299 cement phases, and (iii) proper sulfation of a cement is empirically known to influence setting  
300 and early-age strength gain. Predictions showed only marginal improvement upon inclusion of  
301 any other other attributes from Table 1 besides these six, such as minor oxides (Mg, Na, K), loss  
302 on ignition, or free lime content. Whether added alone or in combination with other such  
303 attributes, none affected the MAPE by more than 0.1 %. Replacing the four major oxides with  
304 the Bogue estimates of the four major clinker phases also did not improve estimator  
305 performance, which is understandable because the Bogue estimates are merely linear functions  
306 of the oxide proportions.

307

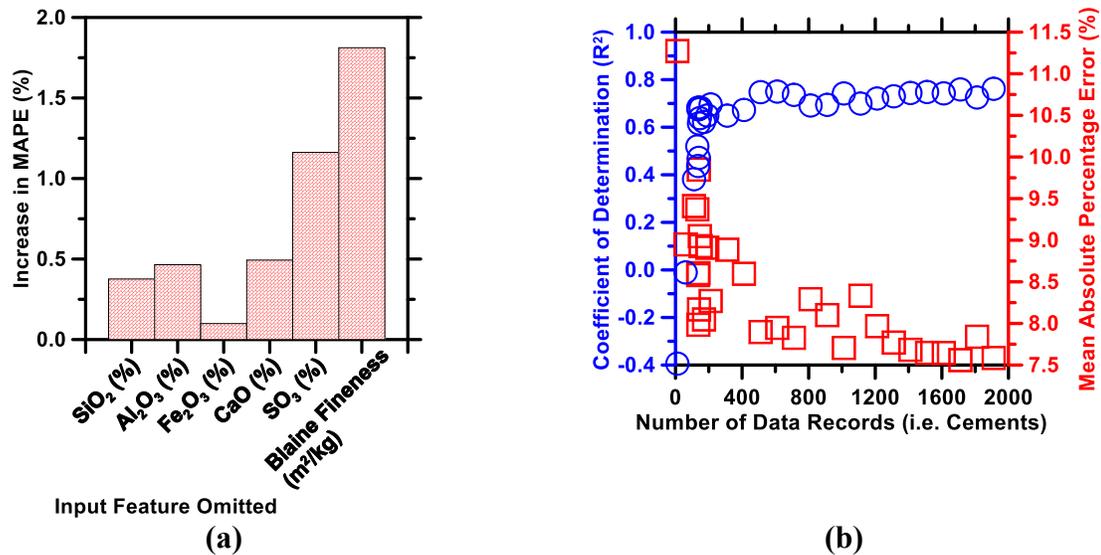
### 308 ***3.5 Random omission identifies a tentative lower bound on data needed to train estimators***

309 Having now established the minimum attributes necessary for predicting the primary targets, we  
310 now turn attention to determining the minimum number of data records needed to make accurate  
311 target estimates. This measure of robustness of the different ML algorithms, when applied to  
312 these datasets, can be evaluated by retraining them with a sparse subset of the data. Specifically,  
313 learning curves were constructed by randomly omitting data records from the input, as illustrated  
314 in Figure 2(b). For convenience in terminology, we define “data-sufficiency” as the minimum  
315 number of data records at which the learning curves plateau. Figure 2(b) shows that the  
316 estimators approach peak performance, at least with respect to  $R^2$ , with less than 10 % of the  
317 available dataset; those trained with a random selection of at least 200 of the 2211 total available  
318 data records performed within about 1 % of the MAPE of the same estimators that were given  
319 access to the full training set. This suggests the viability of applying such estimators even for  
320 relatively smaller datasets and is an encouraging sign that these methods can also be used  
321 reliably even with limited field data. However, the error metrics frequently used to evaluate the  
322 quality of ML estimators, such as MAPE, are not necessarily suitable for the direct comparison  
323 between estimator accuracy *on average* and the ability of the estimator to make consistently  
324 accurate predictions of engineering properties of particular cement systems.

---

\* Conventional cement chemistry notation is used: C = CaO, S = SiO<sub>2</sub>, A = Al<sub>2</sub>O<sub>3</sub>.

325



**Figure 2:** Representative evaluations of estimator performance shown for the extremely random forest estimators constructed to estimate 3 d compressive strength which highlight (a) Attribute importance as determined by an increase in MAPE upon omission of a given input attribute, and (b) so-called “learning curves” for the estimator showing the minimum number of input records required to construct an adequate estimator.

326

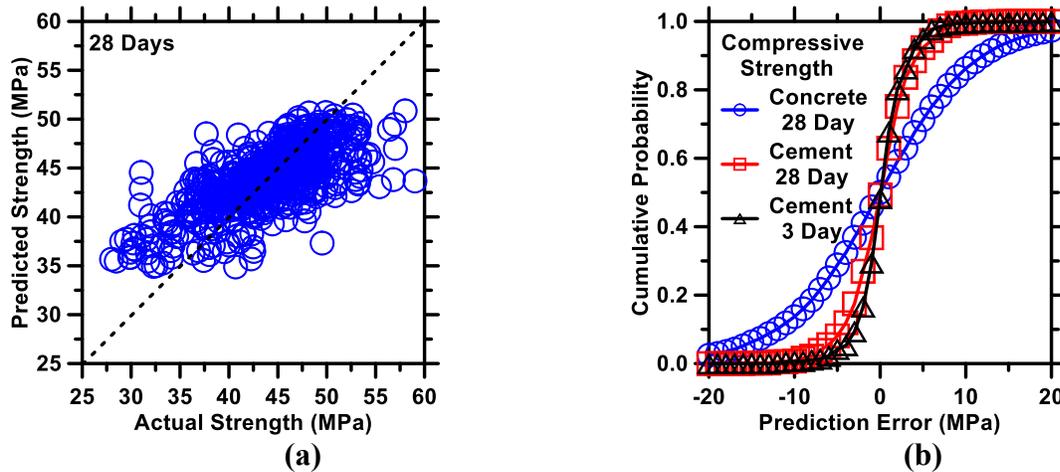
### 3.6 New evaluative metrics are needed to properly reflect estimator prediction accuracy

The three objective functions used to score the estimator performance in this study, which are among the most commonly used scoring metrics in other machine learning efforts, reflect the estimator’s performance on average for the entire dataset, which comprises many cements. However, indicators of average error such as RMSE and MAPE do *not* indicate the estimator’s accuracy in predicting the target value of any *particular* cement in the testing set. Just as a significant fraction of a normally distributed population lies outside one standard deviation of its mean, so does a given estimator produce individual errors much greater than the RMSE for a significant fraction of the cements. As an example, Figure 3(a) shows the individual predictions of 28 d compressive strength made by an extremely random forest regressor with 500 trees applied to a testing set after training. The predicted value for each data record is plotted against the actual target value for that record. The RMSE for this estimator is less than 5 MPa, but the maximum error for any particular cement could be as high as 20 MPa and corresponds to a relative error of about 50 %.

341

To view the situation in a different way, the absolute prediction errors for 28 d strength of individual cements were collected in a histogram with 1 MPa bin widths. The histogram was converted into a normalized probability density plot, the positive portion of which is shown in Figure 3(b). For comparison, the same figure shows the corresponding histograms for 3-d compressive strength obtained in this study and for 28-d concrete strength obtained by Young *et al.* [1]. The errors have an approximately normal distribution with a peak near 2 MPa and a standard deviation of approximately 3.6 MPa. A tolerance interval for an ML estimator can then be established in a similar manner to the ASTM standard test methods. For example, given that the 28 d strength errors in Figure 3(b) are approximately normally distributed with a mean of 2 MPa and a standard deviation of 3.6 MPa, there is a 95 % probability that 90 % of the predictions will be no less than 6.2 MPa below the actual value and no more than 10.2 MPa

353 above it. A tolerance interval this large is far from ideal. However, for comparison the interval  
 354 for similar estimations from concrete mixture proportions by Young *et al.* [1] comes in at about  
 355  $\pm 15.5$  MPa. As illustrated in Figure 3(b) by comparison to predictions on concrete, as well as 3 d  
 356 strength, the current results clarify both the substantial improvement achieved by inclusion of  
 357 attributes such as cement composition, as well as the potential future improvements that may  
 358 arise from inclusion of additional attributes such as curing conditions.  
 359

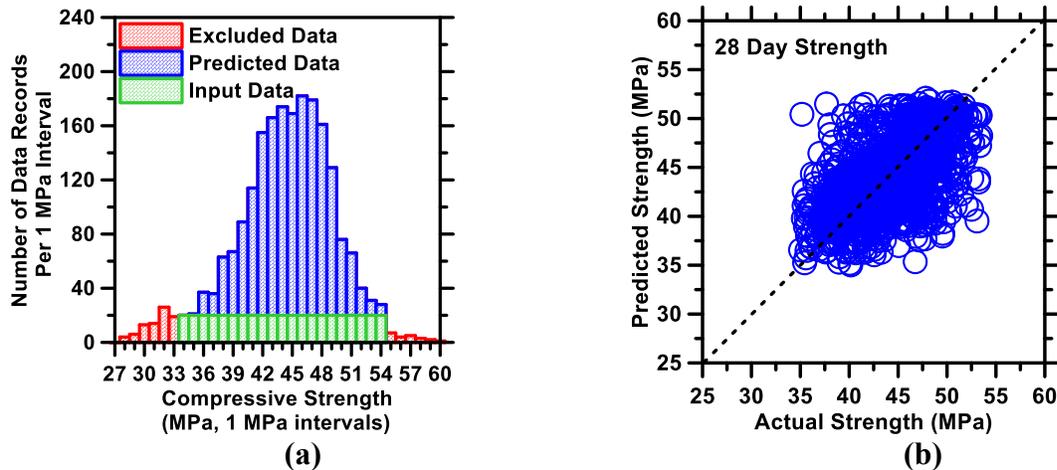


**Figure 3.** The prediction results of an optimized 500-tree extremely random forest regression estimator, shown as (a) predicted vs actual strength values with a dashed line of identity provided to guide the eye, and (b) the normalized cumulative probability distribution of a prediction by the estimator having a given error. Also shown for comparison are distributions for a similar estimator applied to prediction of 3 d compressive strength of mortars (this study) and 28 d compressive strength of concretes (Young *et al.* [1]).

360  
 361 If ML estimators are to be used confidently for concrete mixture design and optimization, they  
 362 will need to achieve much lower tolerance intervals in their predictions than are indicated herein.  
 363 In statistical treatments such as those discussed above, the only way to reduce the probability that  
 364 a particular estimate is outside a tolerable limit is to significantly reduce the average error values  
 365 such as RMSE and MAPE, or to effectively tighten the distribution of errors about these average  
 366 values. The ways to decrease average error are to provide the estimator with data that more  
 367 uniformly span the range of possible values, to acquire better curated data, or to identify and  
 368 collect data on other attributes that may relate more meaningfully to the target being estimated.  
 369 Within the narrowly prescribed range of cement compositions and characteristics considered  
 370 herein, namely ASTM C150 PCs, the dataset would appear to be easily large enough to train the  
 371 estimators according to the plateau in learning curves demonstrated in Figure 2(b).  
 372 Consequently, the only feasible way to reduce the unexplained variance is to develop a means for  
 373 identifying relatively more inconsistent data within the currently applied dataset, or to  
 374 supplement the data with measurements of other material or processing characteristics that are  
 375 currently not being routinely captured including, but not restricted to, the types and dosages of  
 376 chemical admixtures, the particle size distribution of the OPC, clinker grinding parameters,  
 377 curing conditions, and data on the mineralogy, texture, and impurities in the individual cement  
 378 components.

379  
 380 **3.7 Under-sampling intermediate strength values reduces estimator bias**  
 381 The correlation between predicted and actual 28 d compressive strength values, as illustrated in  
 382 Figure 3(a), exhibits a distinct bias: low actual compressive strength values are consistently over-

383 predicted, while high values are consistently under-predicted. This suggests that such regression  
 384 estimators, including ensemble models such as extremely random forests, suffer from an  
 385 imbalance in the input data used to train them, specifically in that a scarcity of very low and very  
 386 high compressive strength values leads to less accurate predictions in these ranges. This issue has  
 387 been frequently addressed in the field of ML classification [26] by resampling, that is, omitting  
 388 or adding data records in the ML training set. Development of this practice for regression  
 389 estimators is only in its early stages [26], with primary interest so far in its ability to allow for  
 390 prediction of rare extremal values [27]. In the current case, where more accurate predictions  
 391 within a narrowly prescribed range are the goal, resampling methods provide a ready means to  
 392 reduce estimator bias by simply omitting a selection of the input data.  
 393



**Figure 4.** (a) The distribution of measured compressive strength values from the full dataset, with data that was used as input to train ML estimators, predictions to test ML estimators, and excluded data marked in green, blue, and red, respectively. (b) Prediction results of an optimized 500-tree extremely random forest regressor trained on an input set subject to under-sampling (as illustrated in part (a)), shown as predicted vs actual strength values with a dashed line of identity provided to guide the eye.

394  
 395 A tentative under-sampling procedure, developed specifically for the dataset under consideration,  
 396 demonstrates that the input of *fewer data* is preferable when predicting the compressive strength  
 397 of cement mortars (Figure 4b). The under-sampling in this case was conducted by analyzing the  
 398 distribution in actual compressive strength values (Figure 4(a)), divided arbitrarily into 1 MPa  
 399 intervals. About 90 % of the data records have compressive strengths between 34 MPa and 54  
 400 MPa. At least 20 data records were available within each 1 MPa interval in that range, but not  
 401 outside that range. Therefore, 20 data records were randomly selected from each 1 MPa interval  
 402 within the range of 34 MPa to 54 MPa, and the remainder of the records in that range were used  
 403 to test prediction accuracy. The input set constructed in this manner consisted of 420 data  
 404 records, more than enough to optimally train estimators according to Figure 2(b). Moreover, the  
 405 new restricted training set corrected the bias in 28 d strength predictions, as can be seen by  
 406 comparing Figure 4(b) with Figure 3(a).  
 407

408 The under-sampling procedure described above provides marginal improvements in previously  
 409 discussed average error metrics; for example,  $R^2$  for 28 d strength correlations increased from  
 410 0.506 to 0.582. However, the error in any specific prediction, as before, is still considerably  
 411 larger than that achieved by repeated experimental measurements. Nonetheless, this result  
 412 highlights an important guideline that should be taken into account, both when using existing

413 datasets and when acquiring new data with a broader array of attributes: prediction bias can be  
414 reduced when the training set contains data that are more evenly spread over the entire range of  
415 possible target values. The same principle might apply to imbalances also in specific attributes,  
416 which then reduce estimator performance but are not easily identifiable. This is likely most  
417 applicable to cases for which some of the input attributes are known to vary widely, like those of  
418 concrete mixture proportions, as opposed to the relatively well-bounded cement compositions  
419 considered currently. The potential applications of under-sampling and/or over-sampling across  
420 many attributes to improve the performance of ML regression estimators represents a significant  
421 area for future research, with particular relevance to cement and concrete-type materials.  
422

#### 423 **4. Summary and conclusions**

424 This study takes another important step toward predictive ML models of concrete properties by  
425 including the effects of OPC characteristics on the properties and performance of cement pastes  
426 and mortars. ML methods are applied to estimate 3 d, 7 d, and 28 d compressive strength and the  
427 time of initial set across numerous ASTM C150 compliant PCs – attributes that are typically  
428 measured in a laborious and time-intensive manner using standard test methods. At a minimum,  
429 accurate estimation of these properties by ML requires knowledge of the cement fineness and the  
430 mass fractions of the oxides of silicon, aluminum, iron, calcium, and sulfur. Additionally, a  
431 lower bound of approximately 200 data records for different cements is required to enable this  
432 nature of estimations, with estimator performance improving only marginally with provision of  
433 more data records, likely due to the relatively narrow range of cement compositions and  
434 finenesses that are included. This implies that suitably-trained ML approaches may be used even  
435 when limited data are available.  
436

437 A distinction of the dataset used in this study is that all the attributes and targets were measured  
438 following standard test methods that are intended to minimize the variability of measurement  
439 conditions. One advantage of this is that it enables the ML estimators to isolate and discover the  
440 influences of OPC powder characteristics on engineering performance without the complications  
441 of variability among other important parameters such as mixture proportions and curing  
442 temperature. In the field, these latter variables are not held constant and can have a decisive  
443 influence on concrete performance. However, prior applications of ensemble ML estimators to  
444 field concrete performance have demonstrated that realistic mixture proportioning, and  
445 production procedures and curing conditions can be accommodated and still yield reasonably  
446 accurate estimations of 28 d compressive strength [1]. Therefore, in a limited sense, this effort  
447 confirms the ability of ML methods to estimate how OPC powder characteristics affect binder  
448 properties, while outlining the limitations, such as the difference between an estimator's average  
449 accuracy and its accuracy in making single predictions. Tight tolerance intervals are a major goal  
450 in the ongoing effort to develop more comprehensive ML approaches to predicting the field  
451 performance of concrete with multicomponent binders. ML approaches are all the more desirable  
452 in this context, however, because they can, if provided with suitable and sufficient data, capture  
453 the effects of variable environmental conditions and curing practices on concrete properties and  
454 performance.  
455

#### 456 **Acknowledgements**

457 The authors acknowledge financial support for this research provided by the NIST Engineering  
458 Laboratory's Exploratory Research Program, the National Science Foundation (CAREER:

459 1253269, CMMI: 1562066), and the Henry Samueli Fellowship. The contents of this paper  
460 reflect the views and opinions of the authors, who are responsible for the accuracy of data  
461 presented herein. This research was conducted in the Laboratory for the Chemistry of  
462 Construction Materials (LC<sup>2</sup>) at the University of California, Los Angeles (UCLA) and the  
463 Inorganic Materials Group of the Materials and Structural Systems Division of the Engineering  
464 Laboratory at NIST. The authors gratefully acknowledge the support that has made these  
465 laboratories and their operations possible.

466

## 467 **References**

468

- 469 [1] B.A. Young, A. Hall, L. Pilon, P. Gupta, G. Sant, Can the compressive strength of concrete  
470 be estimated from knowledge of the mixture proportions?: New insights from statistical  
471 analysis and machine learning methods, *Cement and Concrete Research*. 115 (2019) 379–  
472 388.
- 473 [2] J.W. Bullard, H.M. Jennings, R.A. Livingston, A. Nonat, G.W. Scherer, J.S. Schweitzer, K.L.  
474 Scrivener, J.J. Thomas, Mechanisms of cement hydration, *Cement and Concrete Research*.  
475 41 (2011) 1208–1223.
- 476 [3] S. Mindess, J.F. Young, D. Darwin, *Concrete*, Prentice Hall, 2003.  
477 <http://www.bcin.ca/Interface/openbcin.cgi?submit=submit&Chinkey=302193>.
- 478 [4] H.M. Jennings, J.W. Bullard, J.J. Thomas, J.E. Andrade, J.J. Chen, G.W. Scherer,  
479 Characterization and modeling of pores and surfaces in cement paste, *Journal of Advanced*  
480 *Concrete Technology*. 6 (2008) 5–29.
- 481 [5] J.J. Thomas, J.J. Biernacki, J.W. Bullard, S. Bishnoi, J.S. Dolado, G.W. Scherer, A. Luttge,  
482 Modeling and simulation of cement hydration kinetics and microstructure development,  
483 *Cement and Concrete Research*. 41 (2011) 1257–1278.
- 484 [6] I.-C. Yeh, Modeling of strength of high-performance concrete using artificial neural  
485 networks, *Cement and Concrete Research*. 28 (1998) 1797–1808.
- 486 [7] J.-S. Chou, C.-K. Chiu, M. Farfoura, I. Al-Taharwa, Optimizing the prediction accuracy of  
487 concrete compressive strength based on a comparison of data-mining techniques, *Journal of*  
488 *Computing in Civil Engineering*. 25 (2010) 242–253.
- 489 [8] K.O. Akande, T.O. Owolabi, S. Twaha, S.O. Olatunji, Performance comparison of SVM and  
490 ANN in predicting compressive strength of concrete, *IOSR Journal of Computer*  
491 *Engineering*. 16 (2014) 88–94.
- 492 [9] M.F. Zarandi, I.B. Türksen, J. Sobhani, A.A. Ramezani pour, Fuzzy polynomial neural  
493 networks for approximation of the compressive strength of concrete, *Applied Soft*  
494 *Computing*. 8 (2008) 488–498.
- 495 [10] U. Atici, Prediction of the strength of mineral admixture concrete using multivariable  
496 regression analysis and an artificial neural network, *Expert Systems with Applications*. 38  
497 (2011) 9609–9618.
- 498 [11] J. Kasperkiewicz, J. Racz, A. Dubrawski, HPC strength prediction using artificial neural  
499 network, *Journal of Computing in Civil Engineering*. 9 (1995) 279–284.
- 500 [12] H.-G. Ni, J.-Z. Wang, Prediction of compressive strength of concrete by neural networks,  
501 *Cement and Concrete Research*. 30 (2000) 1245–1250.
- 502 [13] A. Öztaş, M. Pala, E. Özbay, E. Kanca, N. Caglar, M.A. Bhatti, Predicting the compressive  
503 strength and slump of high strength concrete using neural network, *Construction and Building*  
504 *Materials*. 20 (2006) 769–775.

- 505 [14] M.H. Rafiei, W.H. Khushefati, R. Demirboga, H. Adeli, Supervised Deep Restricted  
506 Boltzmann Machine for Estimation of Concrete., *ACI Materials Journal*. 114 (2017).
- 507 [15] I.B. Topcu, M. Sarıdemir, Prediction of compressive strength of concrete containing fly ash  
508 using artificial neural networks and fuzzy logic, *Computational Materials Science*. 41 (2008)  
509 305–311.
- 510 [16] Z.Q. John Lu, The elements of statistical learning: data mining, inference, and prediction,  
511 *Journal of the Royal Statistical Society: Series A (Statistics in Society)*. 173 (2010) 693–694.
- 512 [17] L. Breiman, Bagging predictors, *Machine Learning*. 24 (1996) 123–140.
- 513 [18] P. Geurts, D. Ernst, L. Wehenkel, Extremely randomized trees, *Machine Learning*. 63 (2006)  
514 3–42.
- 515 [19] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P.  
516 Prettenhofer, R. Weiss, V. Dubourg, Scikit-learn: Machine learning in Python, *Journal of*  
517 *Machine Learning Research*. 12 (2011) 2825–2830.
- 518 [20] E. Fix, J.L. Hodges Jr, Discriminatory analysis-nonparametric discrimination: consistency  
519 properties, California Univ Berkeley, 1951.
- 520 [21] T. Cover, Estimation by the nearest neighbor rule, *IEEE Transactions on Information Theory*.  
521 14 (1968) 50–55.
- 522 [22] ASTM International, Annual Book of ASTM Standards, (2012).
- 523 [23] J.O. Irwin, On a criterion for the rejection of outlying observations, *Biometrika*. (1925) 238–  
524 250.
- 525 [24] G. McLachlan, K.-A. Do, C. Ambroise, Analyzing microarray gene expression data, John  
526 Wiley & Sons, 2005.
- 527 [25] J. Brownlee, Machine Learning Mastery with Python, Machine Learning Mastery Pty Ltd.  
528 (2016) 100–120.
- 529 [26] B. Krawczyk, Learning from imbalanced data: open challenges and future directions,  
530 *Progress in Artificial Intelligence*. 5 (2016) 221–232.
- 531 [27] L. Torgo, P. Branco, R.P. Ribeiro, B. Pfahringer, Resampling strategies for regression, *Expert*  
532 *Systems*. 32 (2015) 465–476.  
533