# A Comparison of Machine Learning Methods
## for Predicting the Compressive Strength of Field-Placed Concrete

*M.A. DeRousseau*[1], *E. Laftchiev*[2], *J.R. Kasprzyk*[1], *B. Rajagopalan*[1], *W.V. Srubar III*[1,†]

[1] Department of Civil, Environmental, and Architectural Engineering, University of Colorado Boulder,

ECOT 441 UCB 428, Boulder, Colorado 80309-0428 USA,

[2] Mitsubishi Electric Research Labs, 201 Broadway FL8, Cambridge, MA 02139

[†] Corresponding Author, T +1 303 492 2621, F +1 303 492 7317, E wsrubar@colorado.edu

## Abstract

This study evaluates the efficacy of machine learning (ML) methods to predict the compressive strength of field-placed concrete. We employ both field- and laboratory-obtained data to train and test ML models of increasing complexity to determine the best-performing model specific to field-placed concrete. The ability of ML models trained on laboratory data to predict the compressive strength of field-placed concrete is evaluated and compared to those models trained exclusively on field-acquired data. Results substantiate that the random forest ML model trained on field-acquired data exhibits the best performance for predicting the compressive strength of field-placed concrete; the RMSE, MAE, and $R^2$ values were 730 psi, 530 psi, and .51, respectively. We also show that hybridization of field- and laboratory-acquired data for training ML models is a promising method for reducing common over-prediction issues encountered by laboratory-trained models that are used in isolation to predict the compressive strength of field-placed concrete.

**Keywords:** concrete; compressive strength; machine learning; prediction; statistical modeling

## 1. Introduction

The 28-day compressive strength of concrete is a critical design parameter for reinforced concrete structures [1]. Empirical prescriptive- and performance-based mixture design methodologies remain the conventional means to obtain concrete mixture design proportions that meet minimum 28-day compressive strength requirements. However, numerical approaches for predicting the 28-day compressive strength of concrete are emerging in the literature. Accurate numerical estimation of the 28-day compressive strength of concrete is desirable because more precise prediction (1) provides assurance of concrete quality, (2) reduces the number of concrete batches that are needed to be tested to meet strength targets, and (3) enables a reduction in factors of safety. Recent computational studies have demonstrated the ability of advanced statistical modeling techniques to numerically predict concrete compressive strength for laboratory-mixed concrete, termed *laboratory concrete* herein [2]–[13].

35    However, prediction of the 28-day compressive strength of concrete placed in the field on an actual

36    construction site, termed *field concrete* herein, remains a challenge for the concrete industry due to

37    variable environmental conditions and other uncertainties encountered during mixture proportioning,

38    transport, placing, curing, and finishing.

39    *1.1 Prediction challenges for field concrete mixtures*

40    Estimating the 28-day compressive strength of concrete is a multifaceted problem. Complex physical and

41    chemical interactions occur between concrete constituents, which, in turn, affect compressive strength.

42    Therefore, nonlinear mathematical models are advantageous for accurately capturing all phenomena. As

43    an example, consider the following physically intuitive correlations: compressive strength decreases

44    (nonlinearly) as the water-to-cement ratio (w/c) increases [14], [15]; increasing air content for improved

45    workability and freeze-thaw resistance also reduces compressive strength [16]. Other correlations have

46    not been as intuitively deduced to-date. For example, it is well known that the proportion of coarse-to-fine

47    aggregate affects compressive strength, but the relationship has not been precisely determined due to

48    confounding factors, such as particle size distribution, aggregate angularity, and water demand. Coarse

49    aggregate, for example, may vary in nominal size, grading, chemical composition, shape, surface texture,

50    and absorptivity [17]; these properties can impact the strength of the interfacial bonds between the

51    aggregate and mortar, which, in turn, affect the compressive strength of concrete. Furthermore, the

52    addition of supplementary cementitious materials (SCMs), like fly ash, slag, and silica fume, also

53    introduce new, complex, and nonlinear relationships to compressive strength because of complex factors,

54    such as fineness, chemical variability, and pozzolanic reactivity [18], [19]. Additionally, the fineness and

55    mineral composition of fly ash and slag can be highly variable, depending on the original industrial

56    source and additional processing steps [20].

57        The conditions of the job site at which field concrete is mixed and placed are also highly variable and

58    lead to high variability in field compressive strength compared to laboratory concrete. For instance, it is

59    commonplace for the environmental conditions at construction sites to be loosely controlled. Here,

60    temperature, humidity, and inclement weather can all affect concrete curing and the final compressive

61    strength [21], [22]. Such variabilities do not exist in laboratory concrete mixing, which suggests that

62    accurate prediction of the compressive strength of field concrete is a more challenging problem compared

63    to compressive strength prediction of laboratory concrete.

64    *1.2 Machine learning methods for compressive strength prediction*

65    Because of the physical limitations described above, there is growing interest in predicting concrete

66    compressive strength using machine learning (ML) models for both field and laboratory concrete

67    mixtures [23], [24]. ML models predict compressive strength (*i.e.*, the target variable) from the types and

68    quantities of the mixture ingredients (*i.e.*, the input variables). Using pairs of data of the form [input

69 variables, target variable], a model is trained from a collected dataset and learns the relationship between

70 the target and input variables without constraint on prior intuitive understanding. The vast majority of this

71 type of research has been performed on laboratory concrete, which, as discussed, suggests limitations on

72 the actual usefulness of these models for predicting the compressive strength of field concrete, given the

73 myriad of convoluting factors.

74     Prior research in ML methods for compressive strength prediction has been limited to testing ML

75 methods using laboratory data to determine best-possible prediction models for concrete compressive

76 strength. A particularly popular ML algorithm is artificial neural networks (ANNs). The first study of

77 ANNs by Yeh *et al*. [25] employed ANNs on a dataset of over 1000 laboratory concrete mixture designs.

78 Since then, other researchers have reported ANN studies with coefficients of determination ($R^2$) of up to

79 0.999 [2]–[8], [10], [26]–[28]. However, a significant number of ANN studies employ less than 100

80 experimental data points, which may not sufficiently sample the predictor variable space. While ANNs

81 are a flexible and powerful ML method, it suffers from the need to train a large number of parameters.

82 For small datasets (as is common for field concrete), ANNs can quickly overfit the data, which leads to

83 strong training set performance but poor generalization performance on new datasets. Other ML methods

84 that appear in the literature include support vector machines (SVM) [25], [26] and decision tree-based

85 models [13], [29]. Studies that employ these methods are less common than ANN studies, due to the

86 historical alignment of compressive strength prediction and ML methods.

87     Some narrower-scope prediction studies that used ML have focused on modeling concrete mixtures

88 that contain particular mixture ingredients, such as fly ash [28], blast-furnace slag [30], recycled

89 aggregate [31], silica fume [32], and metakaolin [33]. This body of research generates models that are

90 useful for predicting compressive strength when specific constituents are included. However, this

91 approach narrowly tailors the model to the particular dataset and, thus, is less useful when either mixture

92 ingredients or external conditions (possibly unmeasured) may change.

93     A recent study by Young *et al.* considered field concrete data and compared the predictive

94 performance of four ML models for predicting both *field* and *laboratory concrete* [23]. This study found

95 that variance can be significantly better explained in the laboratory concrete dataset, which is compatible

96 with the idea that *laboratory concrete* has fewer uncontrolled variables. The study determined that the

97 four ML methods investigated exhibited equivalent predictive performance for *field concrete* – a

98 somewhat unintuitive result, given that the four methods employed do not share common assumptions

99 about the underlying data. In addition, it is also of note that the laboratory and field datasets contained

100 different mixture ingredients (*i.e.*, input variables). For example, the laboratory concrete dataset included

101 blast-furnace slag, while the field concrete dataset did not, making an apples-to-apples comparison

102 difficult between models for both laboratory and field concrete.

103 *1.3. Innovative contribution/knowledge gaps*

104 Despite a large body of research in this area of study, the challenge of training a ML model for accurate

105 prediction of concrete compressive strength remains relevant. More specifically, two significant gaps

106 exist in the literature. First, prior studies are not well-grounded in best-practice methods of the ML

107 community. The standard procedure in ML is to generate a pipeline of methods that increase in

108 complexity [34]. The reason for this is two-fold: (1) while powerful, ML methods often search a large

109 model space and may miss simple solutions recognized by the researcher and (2) the failure of simpler

110 models is typically caused by a failure in model assumptions that reveals previously hidden details about

111 the data interactions and non-linear behavior observed in the system. These failures can thus be used to

112 inform the appropriate choice of ML tools for further development. Second, consensus on the best model

113 architectures for predicting the compressive strength of field concrete has not yet been reached.

114 To this end, this study aims to address the aforementioned knowledge gaps and is particularly focused

115 on approaches for accurate prediction of the compressive strength of *field concrete*. First, we employ the

116 standard ML procedure of testing models of increasing complexity in order to determine the best-

117 performing model for field concrete. This procedure enables us to build on past research by discussing

118 *why* certain ML methods are particularly well-suited for the concrete compressive strength prediction

119 problem. The field concrete dataset in this study contains 1681 concrete mixtures and was collected by

120 the Colorado Department of Transportation (CDOT). The laboratory concrete dataset in this study was

121 obtained from the University of California, Irvine Machine Learning Repository, which contains data for

122 more than 1000 mixtures [35].

123 Following the analysis of the field concrete models trained on the field concrete dataset, we evaluate

124 the ability of ML models learned on laboratory concrete data to predict the compressive strength of field

125 concrete mixtures. For this analysis, we perform the same ML procedures for the laboratory data and

126 select the best-performing model. This model is then used to predict the compressive strength of field

127 concrete mixtures, and the relative model performance is analyzed. It was hypothesized that the

128 laboratory ML model performance would be unsatisfactory for predicting field compressive strength

129 compared to that of models trained exclusively on field concrete data. Finally, this work includes an

130 analysis of laboratory data-trained models that are supplemented with varying percentages of field data in

131 order to determine if such hybridized datasets can improve performance the predictive capabilities of

132 laboratory concrete models.

133 **2. Machine Learning (ML) Methods**

134 As discussed in the introduction, this paper builds a pipeline of ML methods with increasing complexity,

135 such that the underlying structure in the training data can be stepwise analyzed. First, in Section 2.1, we

136 describe the ML methods used in the pipeline. We introduce *linear methods* (*i.e.*, linear regression,

137     polynomial regression), *transformed linear methods* (*i.e.*, kernelized support vector regression, kernelized

138     Gaussian process regression), and *non-linear methods* (*i.e.*, regression trees, boosted trees, random

139     forest). In general, simple models are introduced first, and subsequent models increase in complexity. The

140     simplest methods (*e.g.*, linear regression) tend to require the most assumptions about the underlying data

141     structure, and the most complex methods (*e.g.*, boosted trees) require few assumptions about the

142     underlying structure of the data. Second in Section 2.2, we analyze the utility of predictive models trained

143     on laboratory concrete data for predicting field concrete strength. Third, in Section 2.3, we introduce the

144     performance measures used to evaluate the effectiveness of each model: the coefficient of determination

145     ($R^2$), root mean squared error (RMSE), and mean absolute error (MAE). Last, in Section 2.4, overfitting is

146     discussed, which occurs when a model not only captures the desired qualities in the data, but also begins

147     to exactly model the training data itself. An overfitted model is undesirable because it lowers the

148     predictive performance on "unseen" testing data. In other words, overfitted models do not generalize well

149     to real-world cases. In this analysis, we describe and utilize nested cross-validation as a means reduce

150     overfitting. Reserved testing data is used for final determination of the best-performing model.

151     *2.1 ML Methods*

152     All models were created in the R Project for Statistical Computing [36]; in addition, Table 1 lists the ML

153     methods employed in this study, as well as the specific package and function used for model training. For

154     each ML method, we discuss parameter tuning and the intuitive meaning of the parameters.

155     **Table 1.** ML models and corresponding R packages used in this study.

| Model Type | R Package | R Function |
|---|---|---|
| *Linear Methods* | | |
| Linear regression | stats | lm |
| Polynomial regression | stats | lm |
| *Transformed Linear Methods* | | |
| Kernelized support vector regression | kernlab | ksvm |
| Kernelized Gaussian process | kernlab | gausspr |
| *Non-Linear Methods* | | |
| Regression trees | rpart | rpart |
| Random forest | randomForest | randomForest |
| Boosted trees | gbm | gbm |

156

157     **2.1.1 Linear Regression**

158     The simplest model to apply and analyze is linear regression. In addition to providing useful

159     understanding of the data, linear regression also serves as a good baseline from which other techniques

160     can be evaluated. Linear regression is a model that describes the output (target) variable as a linear

161     combination of the predictor variables [37]. This linear combination is a hyperplane in N-dimensional

162     space, where N is the number of coefficients in the model. The model solution is the hyperplane that

163     minimizes the squared error between the observed output and the predicted output. Mathematically, the

164     solution is described as:

165 $$\hat{y} = \boldsymbol{x}^T\boldsymbol{\beta},$$ Eq. 1

166     where $\boldsymbol{x}$ is the input vector, $\boldsymbol{\beta}$ is the N-dimensional vector of coefficients (parameters) for the linear

167     model, and $\hat{y}$ is the predicted output variable from the model. The underlying assumption in linear

168     regression is that the relationship between the predictor variables and the output variable is linear.

169     Moreover, the model assumes that predictor variables are independent from one another, and the resulting

170     residuals, the difference between the predicted and observed output variables, are both homoscedastic

171     (*i.e.*, have constant variance) and normally distributed. When these assumptions are violated, it indicates

172     that a linear model is not appropriate. When such violations occur, it is reasonable to use transformations

173     on the input data to try to reduce or eliminate the violation in assumptions. Failure of such methods to

174     improve the resulting model error and reduce violation of the assumptions means that the dataset requires

175     more complex non-linear models.

176     **2.1.2 Multivariate Polynomial Regression**

177     Multivariate polynomial regression (called *polynomial regression* in this study) uses nth degree

178     polynomials of the input variables to predict the output variable. Polynomial regression is a generalization

179     of Eq. 1; however, each *x* term may be: (1) an original predictor variable (*e.g.*, $x_1$), (2) a pure higher-order

180     term of one predictor variable (*e.g.*, $x_1^4$), or (3) an interaction term between two or more predictor

181     variables (*e.g.*, $x_1 x_2^2$)

182       A generalized example of an expanded second-order polynomial solution with two predictor variables

183     (for simplicity) is described by:

184 $$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{11} x_1^2 + \beta_{22} x_2^2 + \beta_{12} x_1 x_2$$ Eq. 2

185     The transformation of the predictor variables allows for modeling of higher-order relationships and

186     modeling interactions between the input variables; Eq. 2, for example, shows a parabolic relationship.

187     When the original predictor variables are transformed, they are called "features." This term, also

188     commonly applied to all input variables of the models, denotes the fact that the inputs have been

189     transformed from their original space. In this analysis, polynomials up to third-order are employed, where

190     third order is chosen due to limits in computational power. Since polynomial regression is a form of linear

191     regression, the same assumptions are required—more specifically, independence of the input features,

192     homoscedasticity of the residuals, and normality of the residuals. Note, however, these assumptions apply

193     to the transformed features and not the original data space.

194     **2.1.3 Kernalized Regression Methods**

195     Kernalized regression methods utilize two mathematical concepts applied in tandem – a transformation of

196     the predictor variables and the pairing of the new predictors with a regression method. These pairings can

197     then be analyzed in order to determine which (if any) kernel and regression assumptions fit the data well.

198         Kernels are a set of transformations that can be used to map the original predictor variable space to a

199     high-dimensional feature space [34]. Here, this mapping is more complex than the polynomial mappings

200     in the previous section, and all mappings are the result of extensive previous research effort [38]. Each

201     kernel has its own set of *tuning parameters* that must be optimized. This paper compares four kernel

202     transformations, including: linear kernel, radial basis function (RBF) kernel, sigmoid kernel, and

203     polynomial kernels (up to order 4). The model order of the polynomial kernels is only limited by the

204     available computational power.

205         Kernel transformations have the form:

206     
$$k(\boldsymbol{x}, \boldsymbol{x}') = \langle \phi(\boldsymbol{x}), \phi(\boldsymbol{x}') \rangle \qquad \text{Eq. 3}$$

207     where $k$ is the kernel function, $\boldsymbol{x}$ and $\boldsymbol{x}'$ are N-dimensional input vectors (N is the number of predictor

208     variables), and $\phi$ is a mapping from m dimensions to an m-dimensional space. Note that $\langle \phi(\boldsymbol{x}), \phi(\boldsymbol{x}') \rangle$

209     denotes the inner product between the two mappings and can be thought of as a measure of similarity

210     between the two transformed vectors. The kernel tuning parameters are optimized in tandem with the

211     optimization of a regression model. This optimization is discussed below. Table 2 provides the kernel

212     transformation equations and kernel tuning parameters used in this study.

213             **Table 2.** Kernel Transformation equations and tuning parameters

| Name | Kernel Transformation | Tuning Parameters |
|---|---|---|
| Linear | $k(\boldsymbol{x}, \boldsymbol{x}') = \langle \boldsymbol{x}, \boldsymbol{x}' \rangle$ | n/a |
| Radial Basis Function | $k(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\gamma \lVert \mathbf{x} - \boldsymbol{x}' \rVert^2\right)$ | $\gamma$ |
| Polynomial | $k(\boldsymbol{x}, \boldsymbol{x}') = (\gamma \langle \boldsymbol{x}, \boldsymbol{x}' \rangle + r)^d$ | $\gamma$, r, and d |

214

215     The transformed variables (features) can be utilized with any regression method. The concept here is that

216     parameters for both the kernel transformations and the regression methods are tuned simultaneously such

217     that the cross-validated model error is minimized. When there are multiple tuning parameters, a grid

218     search technique is employed in order to find near-optimal parameter values. In this paper, two kernalized

219     regression methods are tested:  support vector regression and Gaussian process regression.

220         Support vector regression (SVR) is a version of support vector machines (SVM) used for regression

221     purposes (rather than classification) [39]. The regression model generated by SVR depends on only a

222     subset of the dataset, and these data points are deemed *support vectors*. When an SVR model is trained,

223     support vectors are the points from the dataset that produce error values ($\epsilon$) larger than a prescribed

224     threshold value. SVR model training generates values for $\beta_m$ (the coefficients for the transformed support

225     vectors) and $\beta_0$ (the intercept). This occurs via minimization of Eq. 4 using gradient descent:

$$\text{min:} H(\beta_m, \beta_0) = \sum_{i=1}^{N} V(y_i - \hat{y}) + \frac{\lambda}{2} \sum \beta_m^2 \qquad \text{Eq. 4}$$

227     Here, $V(r)$ is the prescribed error measure, $y$ is the observed target variable, and $\lambda$ is a regularization

228     parameter that serves as a degree of importance given to large error values. When $\lambda$ increases, large errors

229     are more greatly penalized in the model; this parameter can be tuned using cross-validation. In this study

230     the SVR is paired with the aforementioned kernel transformations in order to examine the utility of

231     transformations of the predictor variables.

232      The second regression method that is employed with the kernel-transformed data is Gaussian process

233     regression (GP). GP can be thought of as the Bayesian interpretation of linear regression. Rather than

234     assuming that the relationship between the predictor variables and the target variable has the prescribed

235     linear functional form (e.g. $\hat{y} = \boldsymbol{x}^T \boldsymbol{\beta}$), GP simply assumes that the data can be represented as a sample

236     from a multivariate Gaussian distribution and that the mean of this distribution is zero. This approach is

237     "less parametric" in the sense that the model is more loosely defined. Using GP, the predictions of the

238     target variable are made using the conditional probability, $p(y_*|\boldsymbol{y})$. In short: given the data, how likely is

239     a certain prediction for $y_*$? Here, note the subtle difference between $\hat{y}$ and $y_*$. $\hat{y}$ represents a predicted

240     target variable from a model, and $y_*$ represents a distribution of possible outputs from the model. In the

241     case of GP, $\hat{y}$ is the expected value of $y_*$, $E[y_*] = E[y_*|\boldsymbol{y}]$.

242      Given the assumed Gaussian distribution, the matrix of all predictor variables in the dataset ($\boldsymbol{X}$), the

243     output vector ($\boldsymbol{y}$) and the new matrix of data inputs, the goal is to make a prediction on the new set of data

244     points ($\boldsymbol{x}_*$). The derived conditional distribution has the form,

$$y^*|\boldsymbol{y} \sim N(\boldsymbol{K}_* \boldsymbol{K}^{-1}\boldsymbol{y}, \ \boldsymbol{K}_{**} - \boldsymbol{K}_* \boldsymbol{K}^{-1}\boldsymbol{K}_*^T), \qquad \text{Eq. 5}$$

246     where $\boldsymbol{K}$, $\boldsymbol{K}_*$, and $\boldsymbol{K}_{**}$ are the covariance matrices resulting from $k(\boldsymbol{x}, \boldsymbol{x}')$, $k(\boldsymbol{x}_*, \boldsymbol{x}')$, and $k(\boldsymbol{x}_*, \boldsymbol{x}_*')$,

247     respectively. The prediction, $\hat{y}$, is the expected value of this distribution, which can be reduced to the
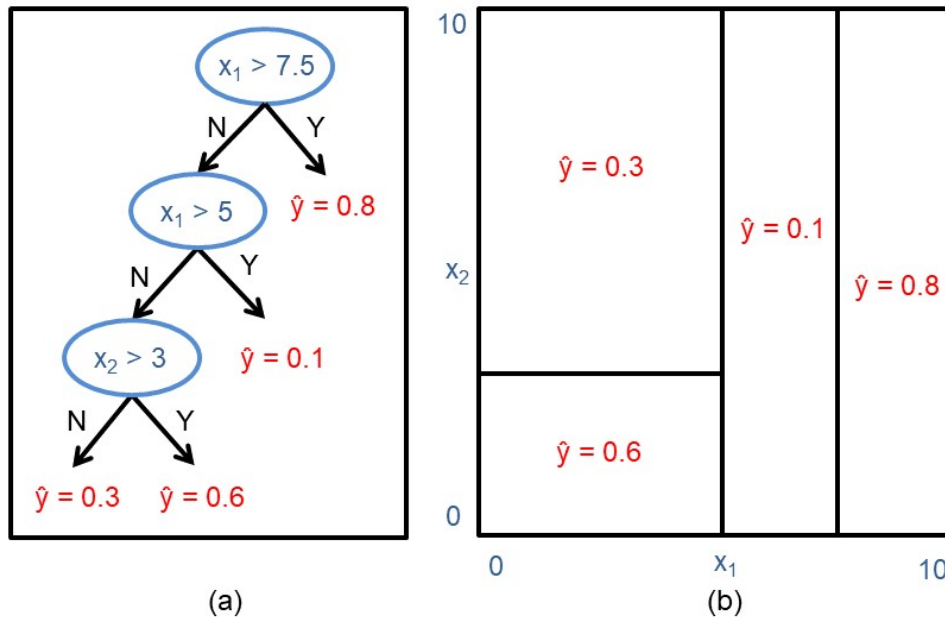
248     equation below.

$$\hat{y} = \boldsymbol{K}_* \boldsymbol{K}^{-1}\boldsymbol{y} \qquad \text{Eq. 6}$$

250     Since GP employs only the assumption of a Gaussian distribution and the covariance matrices for model

251     formulation, no tuning parameters are necessary for this regression method beyond those required for the

252     choice of kernel. The performance of GP allows us to assess the veracity of the Gaussian distribution

253     assumption for the data under multiple different transformations of the predictor variables. If none of the

254     above regression methods can adequately model the output variable, models with no linearity assumptions

255     (e.g. regression trees, artificial neural networks) are reasonable model options consider.

## 2.1.4 Regression Trees

The goal of a regression tree is to generate partitions in the predictor variables such that the target variable can be predicted based on the partitions among the input variables. Figure 1a provides a simple illustration of regression tree "nodes" (*i.e.,* partition rules) and "leaves" (*i.e.,* terminal nodes that lead to one output value). For instance, in the example provided in Figure 1a, there are two predictor variables ($x_1$ and $x_2$). The "root node" (the uppermost blue ellipse) is a rule that partitions the data along $x_1$. For this node, if $x_1$ is greater than 7.5, then the predicted output (in red) is 0.8. However, if the value of $x_1$ is less than 7.5, then one must proceed to the next node in the tree. This process continues until a predicted output variable is reached. For the same example regression tree, Figure 1b demonstrates that a regression tree partitions the predictor variables into rectangular spaces; the and the predicted output is the same value throughout each of these rectangular cells.



(a)　(b)

**Figure 1: (a)** Diagram of an example regression tree model with two predictor variables, $x_1$ and $x_2$. **(b)** This diagram shows the same decision tree using the two predictor variables as axes. It helps visualized the rectangularity of the target variable predictions when simple regression trees are employed. Within each rectangle the predicted target variable would be the same.

Training a regression tree is performed by selecting partitions in succession using a criterion of variance reduction in the target variable [40]. Since each successive partition is always chosen such that the variance of the target variable is reduced, regression trees are prone to overfitting the data. To prevent overfitting, a variety of regularization techniques can be employed. This study minimizes the cost

276　　complexity function, which places a penalty for each additional node that is selected for the model. As

277　　shown below, the cost complexity function $R_\alpha(T)$ has two terms that influence its value:

278　　$$R_\alpha(T) = R(T) + \alpha * f(T),$$　　　　　　　　　　　　　　　Eq. 7

279　　where $R(T)$ is the training error, $f(T)$ is the number of leaves in the regression tree, and $\alpha$ is the

280　　regularization parameter that is determined via cross-validation [41]. In Section 2.3 the cross-validation

281　　procedure used in this study is thoroughly discussed.

282　　　　Regression trees have the advantage that they do not assume linearity in the data, and, therefore, no

283　　complex data transformations are needed. Overall, this approach is simpler than linear methods, but it

284　　requires careful consideration so as not to overfit the data. Regression trees also implicitly select

285　　variables, which means that a trained regression tree will show variables that have more importance for

286　　predicting the target variable in earlier nodes in the tree. Lastly, regression trees are interpretable and can

287　　provide some insight on the dataset being analyzed.

288　　　　A disadvantage of simple regression trees is that they suffer from model instability; in other words,

289　　small changes to the dataset might create a completely different set of partitions, and, consequently does

290　　not lead to the best-performing model. For this reason, more complicated tree-based methods are often

291　　considered that are more stable. Random forest and boosted trees are examples of more complex tree-

292　　based methods that aim to reduce this instability and are discussed in the subsequent sections.

293　　**2.1.5 Random Forest**

294　　Random forest is a method that builds an ensemble of regression trees in order to reduce the instability of

295　　individual trees. Random forest utilizes two strategies for improving the instability issue. First, it employs

296　　the concept of "bootstrap aggregation" (sampling with replacement) in order to generate many similar

297　　datasets that were sampled from the same original dataset. These datasets each lead to an individual tree

298　　within the ensemble. Second, it incorporates randomness during tree-learning in order to reduce the

299　　correlation between each tree within the ensemble. For instance, when generating new nodes (for

300　　individual trees within the random forest), only a subset of the original predictor variables is selected as

301　　the set of candidate variables on which to partition the data. The variable value that minimizes variance in

302　　the output from these randomly selected predictors is the variable selected for that node. This process is

303　　repeated for all nodes in a regression tree and then for all regression trees in the random forest. For a

304　　random forest model, the tuning parameters are: the number of randomly selected predictors ($k$), the

305　　number of individual trees that are trained ($n$), and the tree depth ($d$) [42].

306　　　　The advantage of the random forest method is that it significantly reduces the instability of simple

307　　regression trees. Furthermore, this method has been shown to minimize correlation between trees

308　　compared to other tree-ensemble methods (e.g. "bagging trees" that use only bootstrap aggregation and

309  not random variable selection) [40]. One disadvantage of random forests is their reduction in

310  interpretability compared to simple regression trees; random forests cannot be easily visualized and

311  individual trees are often not good predictive models on their own. However, variable importance plots

312  can reveal the relative importance of predictor variables.

313  **2.1.6 Boosted Trees**

314  Like random forest, boosted trees are an ensemble method for dealing with the instability and poor

315  predictive performance of simple regression trees. Generally, the concept of "boosting" is an ensemble

316  strategy that can be used to improve weak learning algorithms (e.g. regression trees) [43], [44]. Boosting

317  can be applied to any weak learning algorithm but is commonly utilized for regression trees. The main

318  concept of boosting is to build a model using the weak learning algorithm. Then another model is learned

319  on the *residuals* from the first model. This step of model-building on the previous model's residuals is

320  repeated for a set number of iterations. Therefore, a boosted tree is simply a model where the weak

321  learning algorithm used in each iteration is a regression tree.

322     Unlike random forest in which all trees are of the same importance, boosted trees are hierarchical,

323  meaning that each tree layer is constructed recursively. The tuning parameters for boosted trees are: the

324  number of trees, the interaction depth (maximum number of nodes per tree), the minimum number of

325  observations per node (a stopping criteria used to prevent trees that have only one observation at each

326  leaf), and the shrinkage rate (the rate at which the impact of each additional tree is reduced).

327     Boosted trees are similar to random forest in their advantages and disadvantages. Boosted trees tend to

328  have high predictive performance on highly nonlinear datasets and can be successful on problems where

329  there is unequal importance of predictor variables [34]. One disadvantage of boosted trees is that this

330  method has low interpretability; it is difficult to gain much intuition of the patterns that the model has

331  learned or to determine why a boosted tree model is successful (or not) at predicting the target variable.

332  This means that a strong ML pipeline must be used to train boosted trees to ensure that the approach has

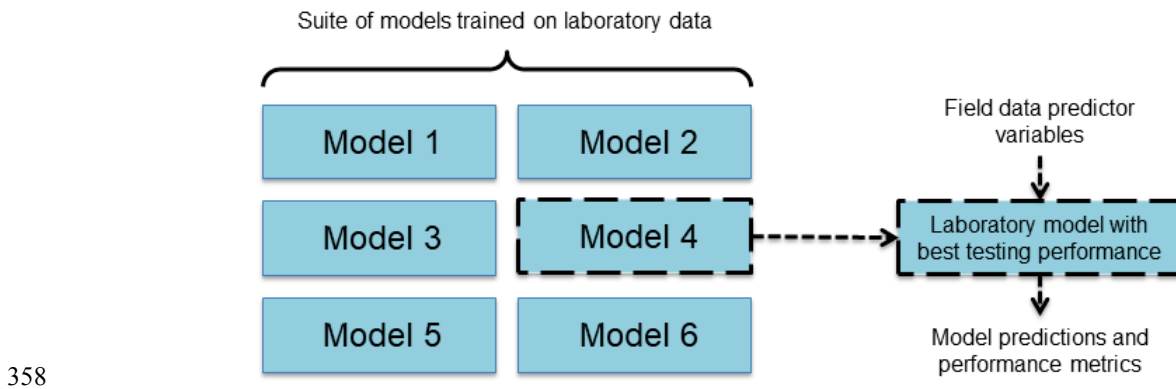333  not overfit the data.

334

335  *2.2 Testing of Laboratory and Hybrid Models for Field Concrete Strength Prediction*

336  **2.2.1 Laboratory Models**

337  As was discussed in the introduction, many studies in the literature have developed ML models for

338  predicting concrete compressive strength using laboratory concrete datasets. While these laboratory

339  models report high predictive performance [2], [4]–[8], [10], [26], [27], it has not yet been tested whether

340  they are useful for predicting the compressive strength of field concrete. A significant novelty of this

341  study is that laboratory models are tested to determine if they are, in fact, useful for predicting

342  compressive strength when presented with other datasets – namely, field concrete data.

343  One issue preventing the direct testing of laboratory ML models from the literature is the use of
344 concrete age as a predictor variable. In other studies, age is a convenient predictor variable because it can
345 explain a high percentage of variance in compressive strength data. In other words, removing age as a
346 predictor and using only the final compressive strength as the output causes the compressive strength
347 problem to be significantly more difficult (*i.e.*, model performance measures tend to be poorer). In this
348 analysis, the desired model output is the final compressive strength (approximated by the 28-day strength)
349 of a concrete mixture as a function of only the quantities of the mixture ingredients. Due to this difference
350 between the prediction problem described herein and that of the literature, laboratory models for
351 predicting the 28-day compressive strength using laboratory concrete data have been trained specifically
352 for this study. Model utility is examined via the process described below and illustrated in Figure 2.
353  First, the aforementioned suite of ML models (*i.e.,* linear regression, polynomial regression, kernel
354 regression, tree-based models) is trained and tested using the laboratory data described in Section 3. The
355 model with the best testing performance is selected. Then, the predictor variables from the field data are
356 used as inputs and the performance measures and diagnostic plots for this new data shall be reported and
357 analyzed.

358



359 **Figure 2.** Process for testing the predictive capability of laboratory models using field concrete data. The
360 dotted outline indicates the laboratory model that is selected based on its performance measures.

361 **2.2.2 Models Trained on Hybrid Data**
362 It was hypothesized that the previously-described laboratory model will not satisfactorily predict the
363 compressive strength of real concrete mixtures. Thus, an analysis of models trained on hybrid data (*i.e.,* a
364 dataset that is composed of both field laboratory data) is conducted to determine whether they can
365 improve predictive performance compared to "pure" laboratory models.
366  Models trained on hybrid data are potentially valuable because there is an inherent tradeoff between
367 the use of laboratory and field models for predicting real concrete compressive strength. On one hand,
368 laboratory data is the cheapest and most accessible data to acquire. It is also the best method for exploring

369    new and exotic concrete mixtures that are uncommon in industry. However, laboratory compressive

370    strength data has the disadvantage that it does not reflect the full set environmental variables experienced

371    by field concrete. Accordingly, it is expected that ML models trained on hybrid data may have the

372    potential to improve the predictive performance of laboratory models.

373        In this novel hybrid approach, a percentage, $\alpha$, of the hybrid dataset is composed of the field data, and

374    the rest is composed of the laboratory dataset. This procedure is used to determine if small amounts of

375    field data can improve model performance. In order to determine the effect of variable amounts of field

376    data, different $\alpha$ values are utilized (10%, 20%, 30%, 40%, and 50%). The model building process occurs,

377    as follows, for each value of $\alpha$:

378    For each $\alpha$:

379      1.  Sort the field dataset in the order of lowest compressive strength to highest compressive strength

380          and partition this sorted dataset into quintiles.

381      2.  In order to ensure the field data portion of the hybrid data is well-sampled, randomly sample (in

382          equal number) the appropriate number for points from the quintiles of the sorted field dataset.

383          Randomly sample from the field dataset the appropriate number of points.

384      3.  Use this hybrid data to train a cross-validated ML model. (The selection of ML model is

385          determined by the best performing laboratory model.)

386      4.  Use the remaining, unsampled field data to determine the average testing performance of the

387          hybrid model. The performance measures described in the following section are reported.

388      5.  Repeat steps 1-4 five times to find average performance measures for each $\alpha$.

389    *2.3 Performance Measures*

390    When training statistical, data-driven models, it is necessary to have a method to quantify the model

391    performance so that hyperparameter tuning can be iterated to select the best possible model. There are

392    several established metrics for determining predictive performance, each with advantages and

393    disadvantages, which will be discussed below. Common quantitative performance measures common to

394    regression modeling (rather than classification modeling) include the coefficient of determination ($R^2$),

395    root mean square error (RMSE), and mean absolute error (MAE) [45], [46]. These metrics, coupled with

396    model diagnostic plots and visualization of predicted versus observed output values, provide a

397    comprehensive picture of a model's performance.

398        $R^2$ is a measure of the proportion of the variance in the data that is explained by the model.

399    Accordingly, $R^2$ is the ratio shown in Eq. 8, where $y_i$ is the observed value from the data, $\hat{y}_i$ is the

400    predicted value from the model, and $\bar{y}$ is the average output from the data.

401
$$R^2 = \frac{\sum_i (\hat{y}_i - \bar{y})^2}{\sum_i (y_i - \bar{y})^2} \qquad\qquad \text{Eq. 8}$$

402     The value of $R^2$ ranges from zero to one, with higher values indicating a better ability to explain the

403     variance in the data with the model. However, $R^2$ is a measure of correlation, not accuracy, and should be

404     used with other performance measures because it is dependent on the variance of the output variable.

405         The root mean square error (RMSE) indicates how concentrated the data is around the model fit. The

406     RMSE is measured on the same scale as the output variable, and is always positive due to the squared

407     residuals in its calculation. Using the RMSE accentuates the effect of outliers in the error metric. This

408     means that if median error of the model (usually captured by the mean absolute error) is low, the RMSE

409     of the model can still be large due to the inability to model some outliers in the data. Given observed

410     values, $y_i$, predicted values, $\hat{y}_i$, and $n$ observed values RMSE is calculated as:

411 $$RMSE = \sqrt{\frac{\sum_{i=1}^{n}(\hat{y}_i - y_i)^2}{n}} \qquad \text{Eq. 9}$$

412     The mean absolute error (MAE) is a measure of prediction accuracy of a model that uses the absolute

413     value of the errors rather than a squared value. The use of the absolute value reduces the influence of very

414     large errors on the measure of performance. Thus, MAE is a measure of the median error of the model

415     and is complimentary to the use of $R^2$ and RMSE.

416 $$MAE = \frac{\sum_{i=1}^{n}|\hat{y}_i - y_i|}{n} \qquad \text{Eq. 10}$$

417     Like RMSE, MAE is measured on the same scale as the output variable, and a lower value indicates a

418     better model fit. In addition, MAE values for a model are typically smaller than the RMSE value for the

419     same model. In this paper we chose to use both RMSE and MAE in order to report both the median and

420     mean error of each model.

421     *2.4 Cross-validation*

422     A critical issue to consider when training and comparing statistical and machine learning models is the

423     prevention of *overfitting*. Overfitting is a problem for ML models that have a high capacity to learn non-

424     linear relationships and are trained on datasets that do not contain a sufficiently large variance of the data

425     (*i.e.,* on datasets that are not rigorously sampled). When using iterative training methods such as grid

426     search, a model is particularly prone to overfitting if the same data is used for the training and validation

427     datasets. In this case, the resulting performance measures would indicate that the model has good

428     predictive performance, but when these models are tested on new data, poor performance is observed.

429         To prevent overfitting, ML learning methods and pipelines can employ several strategies. The strategy

430     employed herein, is called nested cross-validation (nested CV), which splits the data into "training",

431     "validation", and "testing" datasets. In the "inner CV loop", the performance measures are approximately

432     optimized by fitting a model to each of several training datasets. Subsequently, the performance measures

433     are directly optimized by selecting hyperparameters with each validation dataset. In the "outer CV loop",

434 the testing error is estimated by averaging test set scores for several dataset splits. In order to prevent data

435 leakage, it is critical that the trained models have never been exposed to the testing data.

436 When performing CV, the selection of the sizes of the training, validation, and testing sets is critical

437 because this choice affects the bias/variance tradeoff for a given statistical model. To strike a balance

438 between bias and variance error, this paper uses five folds (*i.e.*, partitions) for both the inner and outer CV

439 loops which can generate a favorable bias/variance tradeoff according to the literature [34]. This choice

440 results in 25 validation scores and 5 testing performance measure scores for each model.

441 **3. Datasets**

442 In this study, two datasets – field and laboratory concrete compressive strength data - are used. The field

443 dataset is from the Colorado Department of Transportation (CDOT); it has 1681 mixture designs and

444 corresponding compressive strength values. The mixture constituent variables in this dataset include

445 masses of cement, fly ash, water, water-reducing admixtures (WRA), coarse aggregate, fine aggregate,

446 and percent air entrainment. The laboratory dataset was obtained from the Machine Learning Repository

447 at the University of California, Irvine [35]. This dataset contains over 1000 mixture designs and

448 corresponding compressive strength values. However, it originally contained some mixtures that included

449 blast-furnace slag (a mixture ingredient not included in the field dataset) as well as some mixtures in

450 which the compressive strength was measured earlier than 28 days of curing. In order to reconcile these

451 differences, only mixtures that do not include blast furnace slag and that measure compressive strength

452 after 28 days are included in this analysis. This decision reduced the number of usable mixtures to 311.

453 One last discrepancy is that the laboratory dataset does not report air entrainment values. It is not clear

454 which of the following is true: a constant amount of air was entrained, no air was entrained, or variable

455 amounts of air were entrained but not reported. Notably, this discrepancy does not prevent model training

456 for either dataset. However, when the best laboratory predictive model is used to predict field

457 compressive strength, the air entrainment predictor cannot be utilized.

458

459 Table 1 provides a statistical summary of the two datasets. The laboratory dataset has been converted

460 to US customary units for ease of comparison. Note also that both datasets have used the Absolute

461 Volume Method for proportioning concrete mixtures, which generates weights of ingredients on a cubic

462 yard basis; this means that ingredient quantities are comparable between datasets.

463

464 **Table 1.** Statistical summary of laboratory- and field-acquired datasets.

| Dataset | Statistic | Cement | Fly Ash | Coarse Aggregate | Sand | Water | Air | WRA | Strength |
|---------|-----------|--------|---------|------------------|------|-------|-----|-----|----------|
| | Units | lbs/yd$^3$ | lbs/yd$^3$ | lbs/yd$^3$ | lbs/yd$^3$ | lbs/yd$^3$ | Vol. % | oz/yd$^3$ | psi |
| **Lab** | Mean | 501 | 113 | 1678 | 1332 | 307 | - | 149 | 5357 |

15

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | Median | 487 | 161 | 1689 | 1330 | 314 | - | 154 | 5362 |
| | Min | 227 | 0 | 1350 | 1001 | 236 | - | 0 | 1239 |
| | Max | 910.2 | 337 | 1896 | 1593 | 384 | - | 761 | 11602 |
| **Field** | Mean | 540 | 106 | 1697 | 1256 | 265 | 6.6 | 28 | 5938 |
| | Median | 528 | 120 | 1725 | 1250 | 265 | 5.8 | 24 | 5820 |
| | Min | 395 | 0 | 430 | 445 | 142 | 0 | 0 | 3400 |
| | Max | 900 | 250 | 2240 | 2250 | 392 | 9.6 | 305 | 13040 |

## 4. Results and Discussion

In this analysis, we evaluate the predictive performance of the aforementioned ML models. The values for RMSE, MAE, and $R^2$ for all models are reported in Figure 3. Low values for RMSE and MAE, and high values for $R^2$ indicate better model performance, respectively. For simplicity of discussion, RMSE is used as the primary metric of performance. In addition, both the testing and validation performance is reported, which facilitates the discussion on overfitting in the models. These performance measures are plotted as boxplots to illustrate the range and variance of the error. The set of errors for each model is determined using a nested five-fold cross-validation, with five testing values and twenty-five training values for each model. Each model's performance from a methodological standpoint is discussed in the sections to follow. The methodological and architectural reasons for each model's performance are also examined.

**Figure 3.** Boxplots of the three cross-validated performance measures – **(a)** RMSE, **(b)** $R^2$, and **(c)** MAE for all ML models. Both the training and testing performance measures are reported. The abbreviations are as follows: linear regression (Linear), polynomial regression (Poly), support vector regression (SVR), Gaussian process regression (GP), regression tree (RT), boosted tree (BT), random forest (RF). Kernels are referred to as follows: second-order polynomial (Poly 2), third-order polynomial (Poly 3), fourth-order polynomial (Poly 4), radial basis function (RBF).

17

483    *4.1 Linear Regression*

484    Linear regression is the first model tested in this analysis. This model assumes that the predictors are

485    independent and the residuals are homoscedastic and normally distributed. The performance of linear

486    regression is used as a baseline for comparing model performance and for determining what other models

487    may be more appropriate for the data. For the linear regression model, Bayesian information criterion

488    (BIC) – a parsimonious model selection criterion – is employed to select important predictor variables. Of

489    the seven mixture ingredients, BIC selects five of these as predictor variables (cement, fly ash, water, air,

490    and WRA); this model has a mean testing RMSE, MAE, and $R^2$ of 803 psi, 582 psi, and 0.40,

491    respectively. Of note is the relatively low value of $R^2$, which indicates that a linear model is only able to

492    capture 40% of the variance in the data.

493    There are two possible reasons for the poor performance of this model. One reason is that there are

494    strong predictor variables that were not measured in the dataset. Consequently, the model does not have

495    all necessary information and is unable to perform well. A second possible reason is that the data does not

496    fit the linear assumption of the model, that is, the assumption that the predictors are linearly to produce an

497    output. These possibilities are further evaluated below in diagnostic plots.

498    Four diagnostic plots are shown in Figure 4. Figure 4a shows a plot of the residuals versus the

499    predicted outputs; significant deviation of the smoothed red line indicates non-constant error variances

500    and outliers. For this model, the smoothed average of the error variances indicates nearly constant error

501    variance. The quantile-quantile (Q-Q) plot (Figure 4b) diagnoses the normality of the residuals. Normal

502    residuals (in the statistical sense) lie along the dotted line; however, this figure indicates that there is some

503    deviation from normality of the residuals among higher residual values. Figure 4c is a scale-location plot,

504    which illustrates whether the homoscedasticity assumption is violated. For this plot, the residuals are

505    standardized (to have a mean of zero and a variance of one) and the absolute value is taken. This plot

506    shows that there is a slight increase in error variance with increasing compressive strength, which is

507    indicative of minor heteroscedasticity. Lastly, Figure 4d shows the standardized residuals against their

508    leverage, which is helpful for indicating if particular points more strongly influence the regression. In this

509    case, a few outlier points more highly influence the regression. However, the figure also plots contours of

510    the Cook's distance measure, which measures the effect of deleting a given observation. Cook's distance

511    is increased by both leverage and large residuals. Since no points have a Cook's distance greater than 0.5,

512    there is no great concern about large residuals also having too great of leverage over the fit.

513    One conclusion from the model diagnostics is that there are only minor assumption violations (non-

514    normality of residuals and heteroscedasticity). Despite this result, the linear model retains poor predictive

515    performance, which indicates that there are unmeasured variables needed for predicting compressive

516  strength. Nevertheless, it is reasonable to investigate the use of other types of models to determine if

517  improved performance can be achieved.



519  **Figure 4.** Model diagnostic plots: **(a)** Residuals versus predicted plot to check for non-constant error

520  variance for both positive and negative residuals, **(b)** Quantile-quantile plot to check normality of

521  residuals, **(c)** Scale-location plot to inspect homoscedasticity, and **(d)** Residuals versus leverage plot to

522  determine if any outliers severely impact the regression equation. The blue lines represent the smoothed

523  average for each model dignostic.

524  **4.1.1 Polynomial Regression**

525  Polynomial regression introduces higher order terms and interaction terms between variables, which can

526  sometimes improve model performance because they approximate unobserved phenomena. Here, the

527  polynomial regression has potential because the linear regression analysis indicates a lack of the

528  necessary predictors for improving model performance. In this analysis, polynomial regression is

529  employed for second order and third order terms to determine if there is a physical basis for higher order

530  variables or interaction terms.

531      One key aspect of polynomial regression is that the method acts like a feature selection method. In

532  other words, a set of polynomial features is created, and then the features with the largest reduction in

533  RMSE are kept for the final model. This is the method by which interaction terms are discovered. During

534  the experiments in this paper, the following terms were discovered and included in the model:

535  $(Water) \times (WRA) \times (Air)$ and $(Cement)^2 \times (Fly\ ash)$. The first feature is somewhat intuitive; it is

536  expected that some interaction between water and WRA would be relevant. However, it is somewhat less

537  intuitive that air content is also a part of this feature. The second feature is intuitive because it is expected

19

538     that fly ash and cement would interactively have an impact on concrete compressive strength.

539     Promisingly, polynomials of order two and three decrease the *training* RMSE compared to the linear

540     model by 2.0% and 2.8%, respectively. Given this trend, it's likely that the RMSE of this model will

541     decrease given unlimited computational power.

542      However, it is critical to also analyze the testing error. The testing error values for polynomial orders

543     two and three are higher than the training error by 40.6% and 123.8%. This result suggests that the

544     polynomial regression models are too flexible and overfit the data as the polynomial order grows. Thus

545     this model type is not suitable for compressive strength prediction in concrete.

546     *4.2 Kernel Transformations and Regression*

547     A different approach to discovering interactions and modeling unobserved phenomena is to use non-linear

548     transformations of the data. Some of these are commonly known as kernel transformations. This section

549     will survey techniques in using kernel transformations.

550     **4.2.1 Support Vector Regression**

551     Solving the regression problem using kernel transformations, support vector regression is a popular

552     technique that has shown good results in the literature. In this paper, an array of kernels was tested in

553     cross-validation. These kernels include the RBF kernel, and polynomial kernels (2, 3, and 4).

554     One of the major goals of adaptive regression techniques like SVR is to discover any underlying structure

555     in the data. Of the tested kernels, the RBF kernel has the greatest reduction in RMSE compared to linear

556     regression. Here, RBF SVR reduces the average RMSE by 2.9%. In contrast, the linear and polynomial

557     kernels (orders 2, 3, and 4) reduce this error by -0.6%, -0.1%, 1.1%, and 0.2%, respectively.

558     From this result, it is inferred that the RBF kernel generates the optimal hyperplane for linearly separable

559     patterns among the tested kernels. The minimal improvement from polynomial kernels implies that the

560     regression curve is not well-modeled by a polynomial.

561      The performance of SVR with RBF demonstrates that transformation of the predictor variables

562     improves upon the linear regression baseline model. However, as will be demonstrated in section 3.3,

563     further improvements in performance can be made with other models. One possible explanation for this

564     behavior is that SVR can suffer from the curse of dimensionality in the sense that all terms in the

565     transformed space are given equal weight, so the kernel cannot adapt itself to focus on the critical

566     "subspaces" of the data [34]. Hastie et al. illustrates this concept via a prediction problem with four

567     standard normal features (*i.e*., "real" features) with a polynomial decision boundary and six Gaussian

568     random features (*i.e.,* "noise" features) [34]. Although applying a polynomial kernel with SVR reduces

569     the test error, the real features are drowned out by the noise features. In the example, kernelized SVR is

570     unable to perform as well compared to when the real features are the only modeled features. We

571     hypothesize that this behavior is also true in this case; the noise of irrelevant variables essentially

572 overpowers the predictive capability of SVR to capture the true underlying behavior of field compressive

573 strength.

**4.2.2 Gaussian Process Regression**

575 As is displayed in Figure 3 the GP training and testing performance show that the RBF kernel also

576 generates the highest performance for GP for the kernels utilized in this study. Compared to the linear

577 regression baseline, the GP with RBF-transformed data decreases the average testing RMSE by 3.6%.

578 Utilizing the linear and polynomial (orders 2, 3, and 4) transformations, the reduction in RMSE is -0.1%,

579 1.0%, 2.5%, and 1.6% respectively. With these results, we can conclude that the same transformation

580 (RBF) generates the hyperplane most suitable for use in both SVR and GP.

581 Moreover, this analysis shows that GP is preferred over SVR for this type of data due to its improved

582 performance measures. We hypothesize that GP is a better-performing method (compared to SVR) due its

583 further relaxation of the linearity assumption. Unlike GP, SVR retains the assumption that a

584 transformation of the predictor space causes the data to be linearly separable. GP, on the other hand,

585 makes predictions based on the maximum likelihood of an output given the data, normal parameter

586 distributions, and penalty term that minimize the prediction error. The improved performance of GP over

587 SVR indicates that model performance improves when no linearity assumption exists.

*4.3 Tree-based Models*

**4.3.1 Simple Regression Trees**

590 Unlike the aforementioned techniques, tree-based methods assume that the predictor variables may be

591 partitioned repeatedly and that each final partition generates a different output value. For the simplest

592 tree-based method (regression trees), the average testing RMSE indicates an increase of 6.9% compared

593 to linear regression. We hypothesize that this result is due to the instability of regression trees. In other

594 words, the constructed nodes for a tree may change significantly if the input training sample is slightly

595 changed. Figure 3 illustrates the decreased performance of this model for all three metrics: RMSE,

596 MAPE, and $R^2$.

597 Although the testing performance of the simple regression tree indicates it should not be used for

598 prediction, the results of the model can be used to better understand the relative importance of certain

599 variables for determining concrete compressive strength. In Figure 7, the nodes (e.g. Cement < 569 lbs.)

600 and terminal node predictions (e.g. 4868 psi) are illustrated in the regression tree graph. Values of cement

601 are the first and second nodes, as well as multiple nodes lower in the tree, which indicate the importance

602 of cement quantity as a discriminating predictor variable for this tree. The next most important variable is

603 the quantity of fly ash, which, like cement, has positive correlation with strength. All of the mixture

604 ingredients appear in nodes in the tree, indicating that all are valuable for prediction.

605

**Figure 5.** This figure represents the best-performing simple regression tree graph for compressive
strength prediction. Final predictions from each terminal node are in shown in ellipses.

**4.3.2 Boosted Trees**

Boosted methods are used to reduce the instability of single trees. In this paper, the ensemble tree model
reduced the average testing RMSE by 13.2% compared to the simple regression tree and by 6.9%
compared to linear regression. For this dataset, boosted trees are the second best method for prediction
based on the three performance measures. Notably, the average training RMSE for boosted trees (749 psi)
is slightly lower than that of the random forest model (751 psi). However, the random forest model has
the lower testing RMSE by 5.4%. Despite the nested cross-validation routine, it appears that the boosted
tree model is slightly overfitted due to the higher value of testing RMSE compared to the training RMSE.
Recall from section 2.1.6, that this method iteratively builds regression trees on the residuals from each
consecutive tree. We hypothesize that the model has learned noise in the residuals rather than signal in the
data, which has lead to lower testing performance.

**4.3.3 Random Forest**

Like boosted trees, the random forest model reduces the instability of simple regression trees by utilizing
an ensemble of trees that utilize bootstrap aggregation and random variable selection. Consequently, the
model decreases the average testing RMSE by 9.4% compared to the linear model. It also improves upon
the testing RMSE of the simple regression tree by 20.0%. Furthermore, the average testing error is
slightly lower than the average validation error (730 psi versus 739 psi) indicating that it is unlikely that

625 the random forest model is overfitted. These testing and validation performance measures indicate that

626 random forest is the best method for predicting compressive strength with this dataset. It has the lowest

627 RMSE and MAE as well as the highest R^2 value (730 psi, 530 psi, and .51, respectively).

628     This result may be due to the ability of tree-based methods to learn inconsistent variable importance in

629 the data. In other words, each tree, trained on a subset of the data might learn a slightly different set of

630 variable importance weights. In aggregate, the random forest can then better predict the target variable.

631 An example of inconsistent variable importance can be seen in Figure 5; for mixtures with cement

632 quantities of less than 569 pounds, the next most important variable for determining strength is fly ash. In

633 contrast, above 569 pounds, the next splitting criterion is an even higher quantity of cement. Not only do

634 random forest models have the ability to learn inconsistent variable importance, they also reduce the

635 instability of individual trees and reduce the potential for overfitting [47].
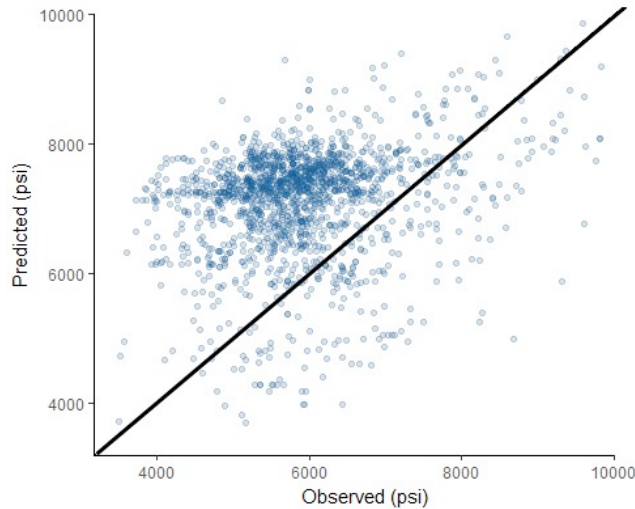
636 *4.4 Prediction of Field Compressive Strength with Laboratory and Hybrid Models*

637 **4.4.1 Models Trained on Laboratory Data**

638 As was discussed in Section 2, many studies in the literature have developed ML models for predicting

639 concrete compressive strength using laboratory datasets. While these laboratory models report high

640 predictive performance, it is relevant to consider whether they are useful for predicting field concrete

641 strength.

642     Consequently, in this study, a suite of ML models (*i.e.,* linear regression, polynomial regression,

643 kernel regression, tree-based models) is trained and tested using the laboratory data described in Section

644 3. Among those tested, the highest-performing model for the laboratory dataset is the random forest

645 model, in which the number of random variables selected at each node was 3, and the number of trees was

646 550 trees; this model achieves a testing $R^2$ value of 0.80.

647     Subsequent to the random forest model selection, the predictor variables from the field data have been

648 used as inputs in the laboratory random forest model to determine how well the model can predict

649 compressive strength of real concrete. The predicted output is plotted versus the observed field strength

650 value in Figure 6. Points near the 1:1 line would indicate a high-performing model. This plot shows that

651 despite its high performance using laboratory data, the laboratory model is not able to predict field

652 strength to a high degree of accuracy; the RMSE for the field data is 1655 psi. Furthermore, this plot

653 illustrates that, overall, the laboratory model tends to over-predict compressive strength. It is likely that

654 this effect is due to the ideal curing conditions in the laboratory setting, which would tend to generate

655 higher compressive strength values than if the same mixture was cured under highly variable

656 environmental conditions.

657

**Figure 6.** Predicted versus observed plot for field compressive strength predictions using the random
forest laboratory model, which illustrates the models tendency to overpredict strength.
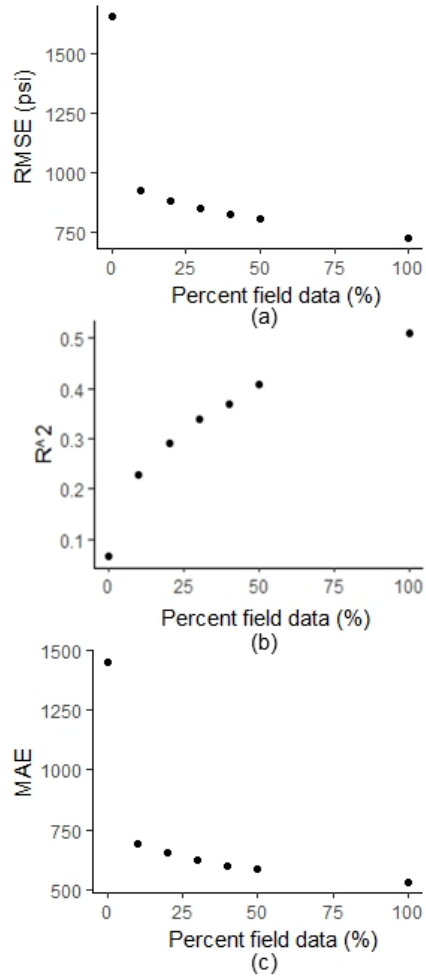
### 4.4.2 Models Trained on Hybrid Data

As was described in Section 2.3.2, models employing hybrid training data are explored in order to
determine if small amounts of field data can improve the performance of laboratory ML models for
predicting compressive strength of field concrete. In this analysis, α values of 10%, 20%, 30%, 40%, and
50% replacement percentages are selected via the quintile sampling method discussed in section 2.3.2.
The remaining, unused field data is to determine the average testing performance of each hybrid model.

As was hypothesized, the inclusion of small percentages of field data significantly reduces the RMSE
MAE and increases the $R^2$ (compared to a pure laboratory model). As is shown in Figure 7, the most
significant model improvements occur with the addition of the initial 10% of field data, which reduces the
RMSE by 43.0%. However, continued performance improvements occur with the additional
supplementation of field data driving the models. Furthermore, Figure 8 illustrates via predicted vs.
observed scatter plots how the addition of field data improves predictive performance. A model
comprised of 100% field data, which was analyzed in Section 4.3, is the standard with which the hybrid
models are compared in terms of the extent to which predictive performance could be improved. This
analysis illustrates that ML modeling of hybrid training data is a promising area of research that improves
upon the downsides of field models and laboratory models being used in isolation.

Future research in this area may explore different ML methods (*i.e.*, models other than random forest)
or other hybridization strategies for utilizing hybrid training data. In addition, it may be of interest to
focus this modeling procedure on concretes with exotic mixture ingredients, which inherently have been
rarely employed in industry, and thus, have few data points with which to model compressive strength.

680
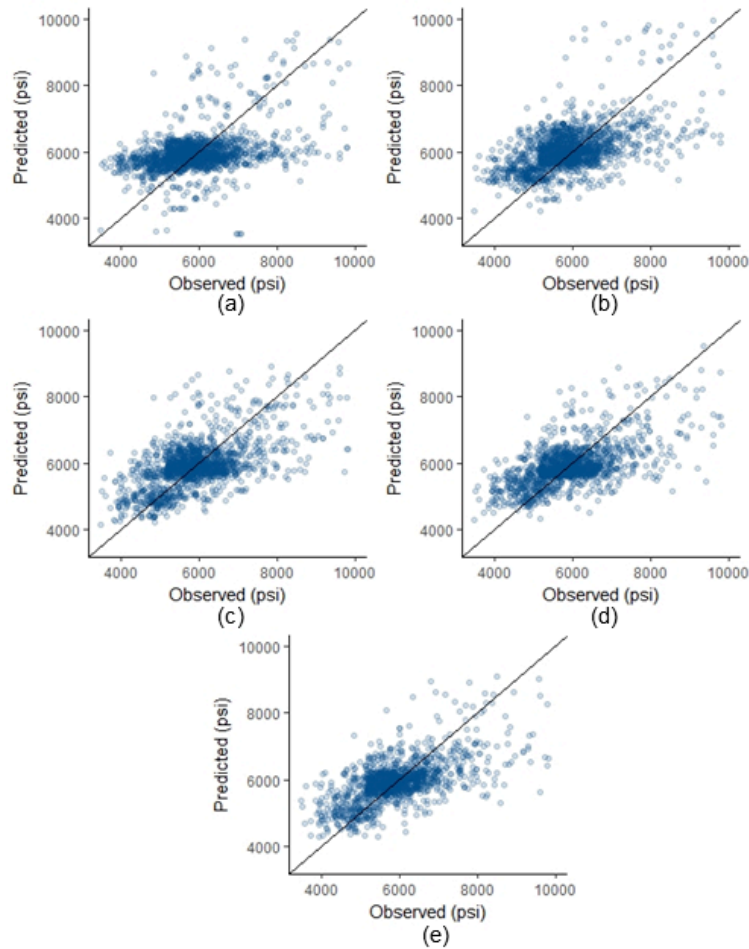
24

681

**Figure 7.** Graphs illustrating the continued improvement in **(a)** RMSE, **(b)** $R^2$, and **(c)** MAE as additional
field data is supplied to the model.

**Figure 8.** Scatter plots of predictive versus observed for ML models trained on hybrid data with the following percentages of field data: **(a)** 10%, **(b)** 20%, **(c)** 30%, **(d)** 40%, **(e)** 50%. Points lying near the one-to-one line indicate better model performance.

## 5. Conclusions

The goal of this work was to specifically analyze the compressive strength behavior of *field concrete* as a function of mixture ingredient quantities. Furthermore, this work trained and tested a variety of ML models for predicting compressive strength of field concrete mixtures and determined which ML models are best suited for the data. By analyzing the performance measures and a variety of diagnostic plots, the reasons for differing performance for field concrete ML models have been elucidated. For instance, from the linear regression model diagnostics, it was found that there are only very minor violations of linearity assumptions; this result indicated it is likely that important predictor variables are missing from the data. Further manipulation of the predictor space via polynomial regression and kernel transformation indicated that a transformed predictor space can improve predictive capability (via a 4% reduction in testing RMSE). Moreover, it was found that nonlinear models, specifically random forest, generated the best

699 performance measures, which is attributed to its full rejection of linear assumptions and ability to learn
700 inconsistent variable importance in the data.

701     It was also confirmed that, at the current time, the most accurate prediction of compressive strength of
702 field concrete is achieved with models trained on field concrete data; however, ML models that employ
703 hybrid training data show promise for significantly improving predictive performance of laboratory
704 concrete models even when only small amounts of field concrete data are available. For instance, it was
705 found that, when only 10% of the training data were from field concrete, the RMSE was reduced by 43%.
706 Moving forward, this research could be extended to explore other ML models with the hybridized
707 approach or applications when it is desirable to explore modeling of exotic concrete mixtures and
708 ingredients.

709     Broadly, the results of this research support two main conclusions: (1) Prediction of field concrete
710 strength requires the application of nonlinear ML models using field-specific data. In particular, advanced
711 tree-based models, such as random forest, are the highest-performing, even when field data is relatively
712 less abundant than laboratory data. (2) Although there is value in testing and statistical-model training for
713 the strength prediction of laboratory concrete, these models should not be used for stand-alone prediction
714 of field concrete strength, because they do not capture the many convoluting factors of field concrete
715 placement and curing. However, ML models that employ hybrid training data can significantly improve
716 the predictive performance compared to laboratory concrete ML models that are used in isolation.

717

## 6. Acknowledgments

723

## References

725 [1]    ACI Committee 318, "Building code requirements for reinforced concrete," American Concrete
726       Institute, ACI 318, 2014.
727 [2]    M. Alshihri, A. Azmy, and M. El-Bisy, "Neural networks for predicting compressive strength of
728       structural light weight concrete," *Constr. Build. Mater.*, vol. 23, no. 6, pp. 2214–2219, 2009.
729 [3]    A. Oztas, M. Pala, E. Ozbay, E. Kanca, N. Caglar, and M. A. Bhatti, "Predicting the compressive
730       strength and slump of high strength concrete using neural network," *Constr. Build. Mater.*, vol. 20,
731       no. 9, pp. 769–775, 2006.

732 [4] C. Bilim, C. D. Atiş, H. Tanyildizi, and O. Karahan, "Predicting the compressive strength of ground
733 granulated blast furnace slag concrete using artificial neural network," *Adv. Eng. Softw.*, vol. 40, no.
734 5, pp. 334–340, May 2009.

735 [5] H.-G. Ni and J.-Z. Wang, "Prediction of compressive strength of concrete by neural networks,"
736 *Cem. Concr. Res.*, vol. 30, no. 8, pp. 1245–1250, Aug. 2000.

737 [6] J. Zhang and Y. Zhao, "Prediction of Compressive Strength of Ultra-High Performance Concrete
738 (UHPC) Containing Supplementary Cementitious Materials," in *2017 International Conference on*
739 *Smart Grid and Electrical Automation (ICSGEA)*, 2017, pp. 522–525.

740 [7] S.-C. Lee, "Prediction of concrete strength using artificial neural networks," *Eng. Struct.*, vol. 25,
741 no. 7, pp. 849–857, Jun. 2003.

742 [8] S. Akkurt, S. Ozdemir, G. Tayfur, and B. Akyol, "The use of GA–ANNs in the modelling of
743 compressive strength of cement mortar," *Cem. Concr. Res.*, vol. 33, no. 7, pp. 973–979, Jul. 2003.

744 [9] I. B. Topcu, "Prediction of properties of waste AAC aggregate concrete using artificial neural
745 network," *Comput. Mater. Sci.*, vol. 41, no. 1, pp. 117–125, 2007.

746 [10] U. Atici, "Prediction of the strength of mineral admixture concrete using multivariable regression
747 analysis and an artificial neural network," *Expert Syst. Appl.*, vol. 38, no. 8, pp. 9609–9618, Aug.
748 2011.

749 [11] M. Rguig and M. El Aroussi, "High-Performance Concrete Compressive Strength Prediction Bsed
750 Weighted Support Vector Machines," *Int. J. Eng. Res. Appl.*, vol. 7, no. 1, pp. 68–75, Jan. 2017.

751 [12] B. G. Aiyer, D. Kim, N. Karingattikkal, P. Samui, and P. R. Rao, "Prediction of compressive
752 stregnth of self-compacting concrete using least square support vector machine and relevance vector
753 machine," *KSCE J. Civ. Eng.*, vol. 18, no. 6, pp. 1753–1758, 2014.

754 [13] C. Deepa, K. Sathiya Kumari, and V. Pream Sudha, "Prediction of the compressive strength of high
755 performance concrete mix using tree based modeling," *Int. J. Comput. Appl.*, vol. 6, no. 5, pp. 18–
756 24, 2010.

757 [14] D. A. Abrams, "Water-Cement Ratio as a Basis of Concrete Quality," *J. Proc.*, vol. 23, no. 2, pp.
758 452–457, Feb. 1927.

759 [15] S. Popovics, "Analysis of Concrete Strength Versus Water-Cement Ratio Relationship," *Mater. J.*,
760 vol. 87, no. 5, pp. 517–529, Sep. 1990.

761 [16] M. S. Mamlouk and J. P. Zaniewski, *Materials for Civil and Construction Engineers*, 2nd ed. Upper
762 Saddle River, NJ: Pearson Education, Inc., 2006.

763 [17] R. Kozul, "Effects of Aggregate Type, Size, and Content on Concrete Strength and Fracture
764 Energy," University of Kansas Center for Research, Inc., Lawrence, KS, SM Report No. 43, 1997.

765   [18] A. Fernández-Jiménez and A. Palomo, "Characterisation of fly ashes. Potential reactivity as alkaline

766        cements☆," *Fuel*, vol. 82, no. 18, pp. 2259–2265, Dec. 2003.

767   [19] A. A. Ramezanianpour and V. M. Malhotra, "Effect of curing on the compressive strength,

768        resistance to chloride-ion penetration and porosity of concretes incorporating slag, fly ash or silica

769        fume," *Cem. Concr. Compos.*, vol. 17, no. 2, pp. 125–133, Jan. 1995.

770   [20] J. Fox, "Fly Ash Classification - Old and New Ideas," presented at the 2017 World of Coal Ash

771        Conference, Lexington, KY, 2017.

772   [21] A. M. Zeyad, "Effect of curing methods in hot weather on the properties of high-strength

773        concretes," *J. King Saud Univ. - Eng. Sci.*, May 2017.

774   [22] O. Cebeci, "Strength of concrete in warm and dry environment," *Mater. Struct.*, pp. 270–272, 1987.

775   [23] B. A. Young, A. Hall, L. Pilon, P. Gupta, and G. Sant, "Can the compressive strength of concrete be

776        estimated from knowledge of the mixture proportions?: New insights from statistical analysis and

777        machine learning methods," *Cem. Concr. Res.*, Sep. 2018.

778   [24] M. A. DeRousseau, J. R. Kasprzyk, and W. V. Srubar, "Computational design optimization of

779        concrete mixtures: A review," *Cem. Concr. Res.*, vol. 109, pp. 42–53, Jul. 2018.

780   [25] I.-C. Yeh, "Optimization of Concrete Mix Proportioning Using Flattened Simplex-Centroid Mixture

781        Design and Neural Networks," *Eng. Comput.*, vol. 25, no. 179, pp. 179–190, 2009.

782   [26] M. Pala, E. Özbay, A. Öztaş, and M. I. Yuce, "Appraisal of long-term effects of fly ash and silica

783        fume on compressive strength of concrete by neural networks," *Constr. Build. Mater.*, vol. 21, no. 2,

784        pp. 384–394, Feb. 2007.

785   [27] G. Trtnik, F. Kavčič, and G. Turk, "Prediction of concrete strength using ultrasonic pulse velocity

786        and artificial neural networks," *Ultrasonics*, vol. 49, no. 1, pp. 53–60, Jan. 2009.

787   [28] İ. B. Topçu and M. Sarıdemir, "Prediction of compressive strength of concrete containing fly ash

788        using artificial neural networks and fuzzy logic," *Comput. Mater. Sci.*, vol. 41, no. 3, pp. 305–311,

789        Jan. 2008.

790   [29] Y. Ayaz, A. F. Kocamaz, and M. B. Karakoc, "Modeling of compressive strength and UPV of high-

791        volume mineral-admixtured concrete using rule-based M5 rule and treemodel M5P classifiers,"

792        *Constr. Build. Mater.*, vol. 94, pp. 235–240, 2015.

793   [30] C. Videla and C. Gaedicke, "Modeling Portland Blast-Furnace Slag Cement High-Performance

794        Concrete," *Mater. J.*, vol. 101, no. 5, pp. 365–375, Sep. 2004.

795   [31] F. Khademi, S. M. Jamal, N. Deshpande, and S. Londhe, "Predicting strength of recycled aggregate

796        concrete using Artificial Neural Network, Adaptive Neuro-Fuzzy Inference System and Multiple

797        Linear Regression," *Int. J. Sustain. Built Environ.*, vol. 5, no. 2, pp. 355–369, Dec. 2016.

798  [32]  M. Sarıdemir, "Prediction of compressive strength of concretes containing metakaolin and silica
799       fume by artificial neural networks," *Adv. Eng. Softw.*, vol. 40, no. 5, pp. 350–355, May 2009.

800  [33]  E. Güneyisi, M. Gesoğlu, Z. Algın, and K. Mermerdaş, "Optimization of concrete mixture with
801       hybrid blends of metakaolin and fly ash using response surface method," *Compos. Part B Eng.*, vol.
802       60, pp. 707–715, Apr. 2014.

803  [34]  T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining,*
804       *Inference, and Prediction, Second Edition*, 2nd ed. New York: Springer-Verlag, 2009.

805  [35]  "UCI Machine Learning Repository: Concrete Compressive Strength Data Set." [Online].
806       Available: https://archive.ics.uci.edu/ml/datasets/Concrete+Compressive+Strength. [Accessed: 19-
807       Dec-2017].

808  [36]  "R: The R Project for Statistical Computing." [Online]. Available: https://www.r-project.org/.
809       [Accessed: 06-Nov-2018].

810  [37]  G. A. F. Seber and A. J. Lee, *Linear Regression Analysis*. John Wiley & Sons, 2012.

811  [38]  T. Hofmann, B. Scholkopf, and A. Smola, "Kernel Methods in Machine Learning," *Ann. Stat.*

812  [39]  H. Drucker, C. J. Burges, L. Kaufman, A. J. Smola, and V. N. Vapnik, "Support Vector Regression
813       Machines," *Adv. Neural Inf. Process. Syst. 9*, pp. 155–161.

814  [40]  C. Strobl, J. Malley, and G. Tutz, "An introduction to recursive partitioning: rationale, application,
815       and characteristics of classification and regression trees, bagging, and random forests," *Psychol.*
816       *Methods*, vol. 14, no. 4, pp. 323–348, Dec. 2009.

817  [41]  L. Breiman, *Classification and Regression Trees*. Routledge, 2017.

818  [42]  P. Probst, M. Wright, and A.-L. Boulesteix, "Hyperparameters and Tuning Strategies for Random
819       Forest," *ArXiv180403515 Cs Stat*, Apr. 2018.

820  [43]  J. Friedman, T. Hastie, and R. Tibshirani, "Additive logistic regression: a statistical view of
821       boosting (With discussion and a rejoinder by the authors)," *Ann. Stat.*, vol. 28, no. 2, pp. 337–407,
822       Apr. 2000.

823  [44]  J. H. Friedman, "Greedy function approximation: A gradient boosting machine.," *Ann. Stat.*, vol. 29,
824       no. 5, pp. 1189–1232, 200110.

825  [45]  T. Chai and R. R. Draxler, "Root mean square error (RMSE) or mean absolute error (MAE)? –
826       Arguments against avoiding RMSE in the literature," *Geosci. Model Dev.*, vol. 7, no. 3, pp. 1247–
827       1250, Jun. 2014.

828  [46]  P. N. Chatur, A. R. Khobragade, and D. S. Asudani, "Effectiveness evaluation of regression models
829       for predictive data-mining," *Int. J. Manag. IT Eng.*, vol. 3, no. 3, pp. 465–483, Oct. 2013.

830  [47]  L. Breiman, "Random Forests - Random Features," University of California, Berkeley, CA,
831       Technical Report 567, Sep. 1999.