

1 Article

2 Comparative Evaluation of Machine Learning Models for Pre- 3 dicting Compressive Strength of Concrete Made with Soap Fac- 4 tory Wastewater

5 Asad Wahab^{1*}, Tausif Junaid Khan¹, Touqeer Ali Rind¹, Maaz Khan¹, Muhammad Faarid Shah¹, Muhammad Faisal
6 Javed¹

7 1 Department of Civil Engineering (DCvE), Ghulam Ishaq Khan (GIK) Institute of Engineering Sciences and
8 Technology, Topi, Khyber Pakhtunkhwa; gcv2520@giki.edu.pk, tausif.junaid@giki.edu.pk,
9 touqeer.ali@giki.edu.pk, gcv2463@giki.edu.pk, gcv2560@giki.edu.pk, arbabfaisal@giki.edu.pk

10 * Correspondence: gcv2520@giki.edu.pk

11 Abstract

12 The reuse of industrial wastewater in concrete production offers a sustainable alter-
13 native to reducing freshwater consumption in construction. This study investigates the
14 prediction of compressive strength of concrete prepared with soap factory wastewater us-
15 ing machine learning techniques. A dataset containing key mix design parameters was
16 analyzed using five machine learning models: Random Forest, Gradient Boosting, Deci-
17 sion Tree, k-Nearest Neighbors (kNN), and Linear Regression. Model performance was
18 evaluated using MSE, RMSE, MAE, MAPE, and R^2 with 10-fold cross-validation. The re-
19 sults show that Random Forest achieved the best predictive performance ($R^2 = 0.883$), fol-
20 lowed by Gradient Boosting ($R^2 = 0.872$), while kNN and Linear Regression showed lower
21 accuracy. Feature importance analysis indicates that cement content, water–cement ratio,
22 and curing age are the most influential parameters affecting compressive strength. These
23 findings demonstrate the potential of machine learning for predicting the strength of
24 wastewater-based concrete and supporting sustainable construction practices.

25 **Keywords:** Machine learning, Compressive strength, Gradient Boosting, Sustainable con-
26 struction, Mix design optimization

27

28 1. Introduction

29 One of the most widely used construction materials is concrete due to its dura-
30 bility, versatility and affordability. One of the most significant mechanical properties that dictate
31 the performance of structure and service life is the compressive strength of concrete. A
32 number of mix design parameters play a critical role in this strength such as cement con-
33 tent, water/ cement ratio, proportions of aggregates, dosage of superplasticizer, and cur-
34 ing age [1], [2]. Traditionally, determination of compressive strength is done in the form
35 of experimental testing in the laboratory, which is time consuming, resource consuming
36 and not always feasible in the context of quick mix optimization [3]. As the scale and com-
37 plexity of the current construction projects continue to increase, there is also an increased

38 need of efficient and reliable predictive techniques that are capable of estimating the
39 strength of concrete without engaging in cumbersome laboratory experimentation [4].

40 Over the last several years, machine learning (ML) techniques have received a lot of
41 attention because of their capabilities to represent rather complicated and highly nonlin-
42 ear correlations between concrete mix components and compressive strength [5]. It has
43 been established in previous studies that existing statistical and linear regression models
44 could not fully model the nonlinear behavior of cement-based materials [6]. To alleviate
45 these limitations, artificial intelligence-related methods like artificial neural network
46 (ANN), support vector machine (SVM), and decision-tree algorithms have been highly
47 studied. Yeh [7] indicated that ANN models are much better predictors of the compressive
48 strength of concrete compared to linear regression models, whereas Chou et al. [8] re-
49 vealed that SVM performs better in the prediction of strengths.

50 Random Forest and Gradient Boosting algorithms are ensemble learning algorithms
51 that have displayed better predictive behavior because of their capabilities to decrease
52 over-fitting and enhance generalization [9], [10]. The method of the Random Forest algo-
53 rithm proposed by Breiman [9] is a randomized and trained set of decision trees that are
54 integrated to increase robustness. Gradient Boosting as a model suggested by Friedman
55 [10], creates sequence of models to correct the past errors in prediction hence it is espe-
56 cially powerful in regression models. Recent studies have used these ensemble methods
57 to model concrete compressive strength and have reported a greater coefficient of deter-
58 minism (R^2) and a smaller error measure than the non-ensembled models like Decision
59 Trees and Linear Regression [11], [12]. The above conclusions indicate that ensemble-
60 based ML models can be successfully used when forecasting the properties of concrete
61 that depend on several interacting variables.

62 In addition to the development of predictive models, the issue of sustainability has
63 triggered the desire to explore the nature of alternative resources and greener concrete
64 production. The construction sector is a major consumer of freshwater which puts a lot of
65 pressure on the environment, especially in areas that experience water shortage [13]. As
66 such, reuse of factual wastewater in the mixing of concrete has been suggested as a sus-
67 tainable approach. Experimental experiments have already been conducted to determine
68 the impact of wastewater in textile, sugar, and chemical industries on concrete perfor-
69 mance and have shown that with the proper usage of wastewater, compressive strength
70 can be acceptable or even better [14], [15]. Nevertheless, the nature of wastewater can be
71 affecting concrete properties in a complex way, emphasizing the importance of prediction
72 modeling methods.

73 The wastewater of soap factories is full of surfactants, organic compounds, and dis-
74 solved salts that may affect the cement hydration and the microstructural formation. Ex-
75 perimentally, Zoyem et al. [16] showed that compressive strength can be similar using
76 partial replacement of mixing water with soap factory wastewater when controlled con-
77 ditions are considered. Although these findings are encouraging, most of the current re-
78 search is based on ex-experimental trial-and-error techniques limiting the scalability and

79 implementation in practice. Furthermore, the use of machine learning methods in using
80 wastewater-based concrete data is sparse in literature.

81 Hence, there is an obvious research gap in nominating machine learning models with
82 experimental data of industrial wastewater-modified concrete. This research fills this gap
83 by using and comparing five machine learning models, i.e., Gradient Boosting, Random
84 Forest, Decision Tree, k-Nearest Neighbors (kNN), and Linear Regression to predict the
85 compressive strength of concrete containing soap factory wastewater. These models are
86 measured in terms of several statistical metrics that focus on the most efficient predictive
87 method and contribute to the sustainable and data-driven concrete mix design.

88 **2. Methodology**

89 *2.1 Data Collection*

90 The dataset used in this study was obtained from a previously published experi-
91 mental investigation that examined the influence of soap factory wastewater on the com-
92 pressive strength of concrete (Zoyem et al.). In that study, several concrete mixtures were
93 prepared by partially replacing conventional mixing water with industrial wastewater in
94 different proportions. The experimental program documented key mix design parame-
95 ters, including cement content, aggregate proportions, water–cement ratio, and curing
96 age. Compressive strength tests were performed at different curing periods to evaluate
97 the mechanical performance of the mixtures. The resulting dataset provides detailed in-
98 formation on mix composition and corresponding compressive strength values, which
99 were used in this study to develop machine learning models for predicting concrete
100 strength and assessing the potential use of industrial wastewater as a sustainable alterna-
101 tive in concrete production.

102 *2.2 Data Preprocessing*

103 Prior to model development, exploratory statistical analysis was carried out to under-
104 stand the distribution and variability of the dataset variables. Descriptive statistics of the
105 input parameters and target variable are presented in Table 1, including the mean, mode,
106 median, dispersion, minimum, and maximum values.

107 The cement content shows a mean value of 276.50 kg/m³ and a median of 266.00 kg/m³,
108 indicating moderate variability across the concrete mixtures. The water content has a
109 mean value of 182.98 kg/m³ and relatively low dispersion, suggesting that water content
110 remained comparatively consistent among most samples.

111 The water–cement ratio shows noticeable variation across the dataset, reflecting dif-
112 ferences in mix design proportions and its important role in strength development. Simi-
113 larly, the superplasticizer content exhibits greater variability than most other input varia-
114 bles, indicating that admixture dosage was not uniform across all mixtures.

The coarse aggregate and fine aggregate contents remain comparatively stable, with mean values of 964.83 kg/m³ and 770.49 kg/m³, respectively. Their relatively low dispersion suggests that aggregate proportions were maintained within a narrower range than cementitious and admixture-related parameters.

The age of testing has a mean value of 44.06 days, while both the mode and median are 28 days, which corresponds to the standard curing duration commonly used for compressive strength evaluation. This indicates that most specimens were tested at 28 days, although additional curing ages were also included to capture longer-term strength development.

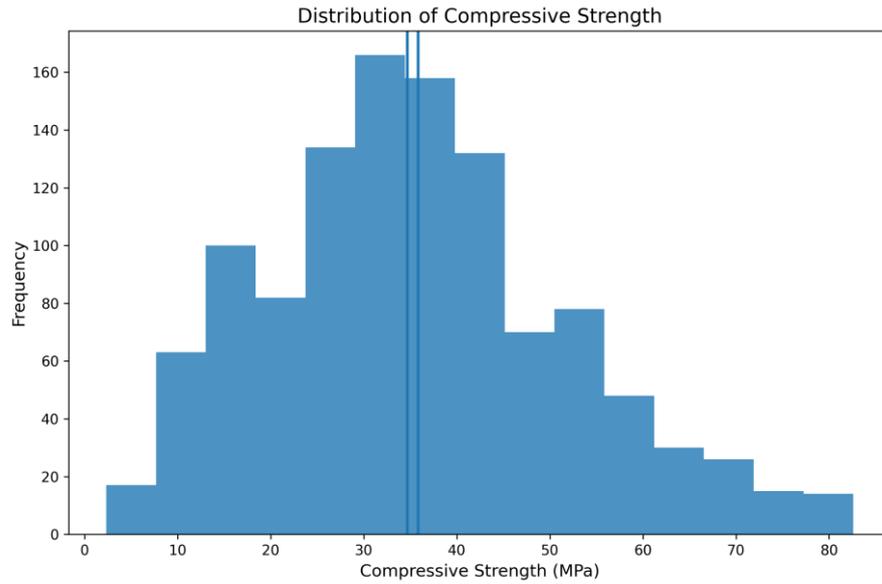
The compressive strength values show substantial variation, with a mean of 35.83 MPa, demonstrating that the dataset covers a broad range of concrete performance levels. As shown in Figure 1, the strength distribution is concentrated mainly in the medium-strength range, with fewer observations at very low and very high strength values. This spread provides sufficient variability for training and evaluating machine learning models.

Overall, the statistical analysis indicates that the dataset contains meaningful variation in the principal mix design variables, particularly cement content, water–cement ratio, superplasticizer dosage, and curing age, all of which are known to influence concrete compressive strength. Before model training, the input variables were standardized to ensure fair learning across features with different numerical scales.

A correlation heatmap was generated to visualize relationships among the dataset variables. The analysis in figure 2 indicates that cement content and curing age exhibit positive relationships with compressive strength, while the water–cement ratio tends to show a negative relationship. These trends align with established principles of concrete mix design and confirm that the dataset contains meaningful interactions suitable for machine learning modeling.

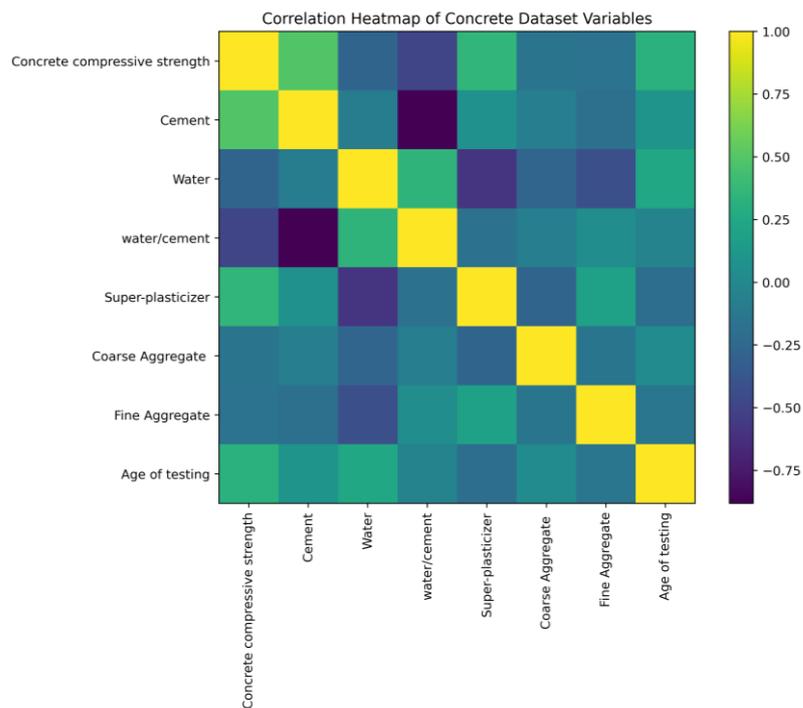
Table 1: Descriptive Statistics

Feature Statistics Name	Mean	Mode	Median	Dispersion	Minimum	Maximum
Cement	276.50	362.60	266.00	0.37	102.00	540.00
Water	182.98	192.00	185.70	0.12	121.75	247.00
water/cement	0.76	0.45	0.70	0.42	0.27	0.88
Super-plasticizer	6.41	0.00	6.70	0.90	0.00	32.20
Coarse Aggregate	964.83	932.00	966.80	0.09	708.00	1145.00
Fine Aggregate	770.49	594.00	777.50	0.10	594.00	992.60
Age of testing	44.06	28.00	28.00	1.37	1	365.00
Compressive strength	35.83	33.39	34.67	0.45	2.33	82.60



142
143

Figure 1: Distribution of compressive strength



144

145

146

Figure 2: Correlation heatmap showing relationships between input variables and compressive strength.

147

2.3 Machine Learning Models

148

2.3.1 Gradient Boosting

149

Gradient Boosting is an ensemble learning technique that builds multiple weak models (typically decision trees) sequentially, where each new model corrects the errors of the previous ones. It minimizes a loss function using gradient descent [10].

150

151

Equation:

152

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x)$$

153

2.3.2 Random Forest

Random Forest is an ensemble of multiple decision trees where each tree is trained on a random subset of the data. The final prediction is obtained by averaging (for regression) or majority voting (for classification) [9].

Equation:

$$\hat{y} = \frac{1}{N} \sum_{i=1}^N T_i(x)$$

2.3.3 Decision Tree

A Decision Tree splits the data into branches based on feature values, forming a tree-like structure. The split is determined by using criteria like Gini impurity or mean squared error (MSE) for regression [21].

Equation (for regression using MSE):

$$\text{Split} = \arg \min_s \sum_{i=1}^n (y_i - \bar{y}_s)^2$$

2.3.4 k-Nearest Neighbors (kNN)

kNN is a non-parametric model that predicts a value based on the average (for regression) or majority class (for classification) of its nearest neighbors [22].

Equation:

$$\hat{y} = \frac{1}{k} \sum_{i=1}^k y_i$$

2.3.5 Linear Regression

Linear Regression models the relationship between input features and output by fitting a linear equation to the data [6] [23].

Equation:

$$y = \beta_0 + \sum_{i=1}^n \beta_i x_i + \epsilon$$

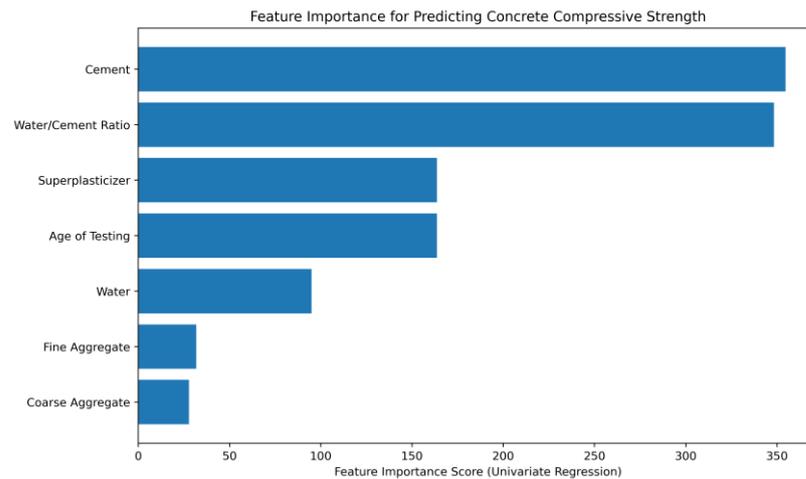
2.4 Training and Testing

To evaluate the predictive performance of the machine learning models, a 10-fold cross-validation approach was employed. In this method, the dataset is divided into ten equal subsets. During each iteration, nine subsets are used for model training while the remaining subset is used for validation. This process is repeated ten times so that each subset is used once as a validation set. The performance metrics are then averaged across all folds to obtain a reliable estimate of model accuracy. This cross-validation strategy reduces the risk of overfitting and ensures that the models are evaluated on multiple data partitions, providing a more robust assessment of predictive capability. The procedure was applied consistently to all five models: Gradient Boosting, Random Forest, Decision Tree, k-Nearest Neighbors (kNN), and Linear Regression.

2. Results and Discussion

3.1 Feature Importance Analysis

191 Feature importance analysis was conducted to determine the relative influence of in-
192 put variables on compressive strength prediction. The results in figure 3 indicate that ce-
193 ment content and water–cement ratio are the most influential parameters, followed by
194 superplasticizer content and curing age. In contrast, fine and coarse aggregate contents
195 show comparatively lower importance, indicating that variations in these parameters con-
196 tribute less to the predictive performance of the machine learning models. These findings
197 align with established concrete mix design principles where cement content and water–
198 cement ratio are primary factors governing strength development.



199
200 Figure 3: Feature importance ranking of input variables influencing compressive strength prediction
201 based on Univariate Regression

202 3.2 Model Evaluation

203 The performance comparison of the five machine learning models used for predicting
204 the compressive strength of concrete is presented in Table 2. Among the evaluated models,
205 Random Forest achieved the best predictive performance, exhibiting the lowest error val-
206 ues with MSE = 30.417, RMSE = 5.515, MAE = 3.954, and MAPE = 0.140, along with the
207 highest coefficient of determination ($R^2 = 0.883$). These results indicate that Random Forest
208 is highly effective in capturing the complex nonlinear relationships between the concrete
209 mix design parameters and compressive strength.

210 The Gradient Boosting model also demonstrated strong predictive capability, with
211 slightly higher error values (MSE = 33.250, RMSE = 5.766, MAE = 4.240, MAPE = 0.147)
212 and an R^2 value of 0.872, indicating that ensemble learning methods perform well for this
213 prediction task. The Decision Tree model showed moderate performance with an R^2 value
214 of 0.824, although its prediction errors were higher compared to the ensemble-based mod-
215 els.

216 In contrast, the k-Nearest Neighbors (kNN) and Linear Regression models showed
217 comparatively lower prediction accuracy. The kNN model produced higher error values
218 and an R^2 of 0.718, while Linear Regression performed the worst among all models, with
219 the largest prediction errors and the lowest R^2 value of 0.553. This suggests that simple

linear models are not sufficient to capture the complex nonlinear interactions among the mix design variables affecting compressive strength.

Overall, the results demonstrate that ensemble learning models, particularly Random Forest and Gradient Boosting, outperform simpler machine learning approaches in predicting the compressive strength of wastewater-based concrete mixtures.

Table 2: Models' Performance

Model	MSE	RMSE	MAE	MAPE	R ²
Random Forest	30.417	5.515	3.954	0.140	0.883
Gradient Boosting	33.250	5.766	4.240	0.147	0.872
Decision Tree	45.550	6.749	4.720	0.166	0.824
k-Nearest Neighbors (kNN)	73.166	8.554	6.660	0.255	0.718
Linear Regression	115.831	10.762	8.400	0.318	0.553

3.2 Prediction Performance Analysis

Figure 4 presents the relationship between the actual compressive strength values and the values predicted by the five machine learning models. In each plot, the dashed diagonal line represents the ideal prediction line, where the predicted values perfectly match the actual compressive strength. The closer the data points lie to this line, the more accurate the model predictions.

Among the evaluated models, Random Forest shows the best prediction performance, with an R² value of 0.883. The predicted values closely follow the ideal line, indicating a strong agreement between actual and predicted compressive strength values. The Gradient Boosting model also demonstrates high prediction accuracy (R² = 0.872), with most data points concentrated near the reference line.

The Decision Tree model shows moderate prediction capability (R² = 0.824), although a slightly larger spread of points around the reference line indicates higher prediction errors compared to the ensemble models. In contrast, the k-Nearest Neighbors (kNN) model exhibits greater dispersion (R² = 0.718), suggesting weaker predictive performance. The Linear Regression model performs the worst (R² = 0.553), as many predicted values deviate from the ideal line, particularly at higher compressive strength levels.

Overall, the graphical analysis confirms that ensemble learning models, particularly Random Forest and Gradient Boosting, provide more accurate predictions of concrete compressive strength compared with simpler machine learning approaches.

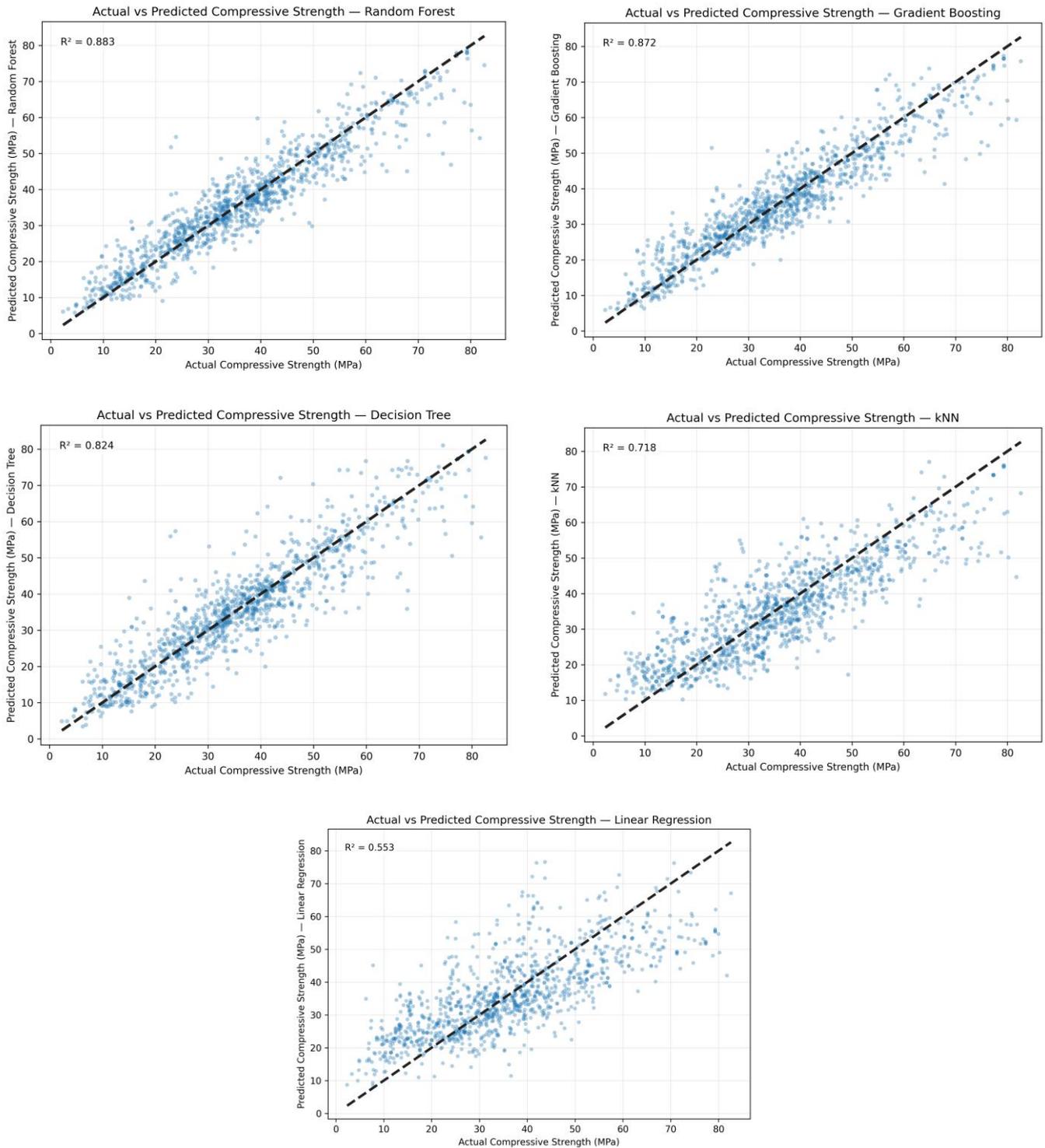


Figure 4: Actual vs predicted compressive strength for the evaluated machine learning models.

4. Conclusions

This study evaluated the ability of five machine learning models—Random Forest, Gradient Boosting, Decision Tree, k-Nearest Neighbors (kNN), and Linear Regression—to predict the compressive strength of concrete produced using soap factory wastewater. The models were assessed using performance metrics including MSE, RMSE, MAE, MAPE, and the coefficient of determination (R^2).

251

252

253

254

255

256

257

258

259 Among the evaluated models, Random Forest achieved the best predictive
260 performance with the highest R^2 value (0.883) and the lowest prediction errors. Gradient
261 Boosting also demonstrated strong performance, while Decision Tree showed moderate
262 accuracy. In contrast, kNN and Linear Regression exhibited lower prediction accuracy,
263 indicating limited ability to capture the nonlinear relationships between mix design
264 parameters and compressive strength.

265 Overall, the results suggest that ensemble learning models, particularly Random
266 Forest and Gradient Boosting, are effective tools for predicting compressive strength of
267 wastewater-based concrete mixtures, supporting more efficient mix design and
268 sustainable construction practices.

269 5. Patents

270 No Patents have resulted from the work reported in this manuscript.

271 Author Contributions:

272 The conceptualization, methodology development, data collection, formal analysis, model
273 testing, and data validation, initial drafting and final paper preparation were carried out by Asad
274 Wahab, Tauseef Junaid Khan, Touqeer Ali Rind. Maaz Khan and Muhammad Faarid Shah played
275 a key role in the initial drafting, methodology development, and critical analysis of the results. Mu-
276 hammad Faisal Javed supported the supervision, manuscript preparation, including reviewing and
277 editing.

278 All authors have read and agreed to the published version of the manuscript.

279 **Institutional Review Board Statement:** Not applicable. This study did not involve humans or ani-
280 mals and therefore did not require ethical approval.

281 **Informed Consent Statement:** Not applicable. This study did not involve human participants.

282 **Data Availability Statement:** Data will be made available upon request.

283 **Funding:** This research received no external funding.

284 **Acknowledgments:** The authors gratefully acknowledge Ghulam Ishaq Khan Institute of Engineer-
285 ing Science and Technology (GIKI EST), for providing support and facilities in conducting this re-
286 search.

287 **Conflicts of Interest:** The authors declare that there is no conflict of interest regarding this study
288 and affirm that the work is original, without any form of plagiarism. All sources of information have
289 been properly cited and acknowledged.

290 Abbreviations

291 The following abbreviations are used in this manuscript:

ANN	Artificial Neural Network
LR	Linear Regression
MSE	Mean Squared Error
RMSE	Root Mean Squared Error
MAE	Mean Absolute Error

292 References

- 293 [1] A. M. Neville, *Properties of Concrete*, 5th ed. London, UK: Pearson, 2011.
- 294 [2] S. B. Singh, P. Munjal, and N. Thammishetti, "Role of water-cement ratio on strength development of cement mortar," *J. Build.*
295 *Eng.*, vol. 4, pp. 94-100, 2015.
- 296 [3] S. Popovics, *Strength and Related Properties of Concrete*, New York, NY, USA: Wiley, 1998.
- 297 [4] S. C. Chapra and R. P. Canale, *Numerical Methods for Engineers*, 7th ed. New York, NY, USA: McGraw-Hill, 2015.

- 298 [5] M. Behnood and E. M. Golafshani, "Machine learning approaches for concrete strength prediction," *Constr. Build. Mater.*, vol.
299 270, 2021.
- 300 [6] M. Zain and S. M. Abd, "Multiple regression model for compressive strength prediction of concrete," *J. Appl. Sci.*, vol. 9, pp. 155–
301 160, 2009.
- 302 [7] I.-C. Yeh, "Modeling of strength of high-performance concrete using artificial neural networks," *Cem. Concr. Res.*, vol. 28, no. 12,
303 pp. 1797–1808, 1998.
- 304 [8] J.-S. Chou, C.-F. Tsai, and Y.-H. Pham, "Predicting concrete strength using machine learning techniques," *Autom. Constr.*, vol.
305 20, pp. 471–479, 2011.
- 306 [9] L. Breiman, "Random forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.
- 307 [10] J. H. Friedman, "Greedy function approximation: A gradient boosting machine," *Ann. Stat.*, vol. 29, no. 5, pp. 1189–1232, 2001.
- 308 [11] H. Nguyen et al., "Concrete compressive strength prediction using ensemble learning," *J. Build. Eng.*, vol. 32, 2020.
- 309 [12] M. Kamath et al., "Machine learning algorithms for high-performance concrete," *J. Eng. Des. Technol.*, 2022.
- 310 [13] V. Harshit, A. K. Rizwan, and K. K. Iqbal, "Sustainable use of wastewater in concrete construction," *J. Build. Eng.*, 2021.
- 311 [14] A. K. Padmini, K. Ramamurthy, and M. S. Mathews, "Influence of wastewater on properties of concrete," *Build. Environ.*, vol.
312 44, pp. 1761–1767, 2009.
- 313 [15] K. S. Al-Jabri et al., "Reuse of industrial wastewater in concrete," *Cem. Concr. Compos.*, vol. 33, pp. 603–608, 2011.
- 314 [16] Z. G. Mathurin et al., "Prediction of compressive strength of concrete made with soap factory wastewater using machine learn-
315 ing," *Model. Earth Syst. Environ.*, vol. 8, pp. 5625–5638, 2022.
- 316 [17] ASTM C39/C39M-21, Standard Test Method for Compressive Strength of Cylindrical Concrete Specimens, ASTM International,
317 2021.
- 318 [18] H. Chen et al., "Predicting compressive strength of cement-based materials," *PLoS One*, 2018.
- 319 [19] L. K. A. Sear et al., "Abrams' law and water–cement ratio," *Constr. Build. Mater.*, vol. 10, pp. 221–226, 1996.
- 320 [20] F. Khademi and K. Behfarnia, "ANN and regression models for concrete strength prediction," *Constr. Build. Mater.*, 2016.
- 321 [21] L. Rokach and O. Maimon, *Data Mining with Decision Trees*, Singapore: World Scientific, 2014.
- 322 [22] T. Cover and P. Hart, "Nearest neighbor pattern classification," *IEEE Trans. Inf. Theory*, vol. 13, no. 1, pp. 21–27, 1967.
- 323 [23] A. Al-Shamiri et al., "Modeling compressive strength of high-strength concrete," *Constr. Build. Mater.*, vol. 208, pp.
324 204–219, 2019.