

MACHINE LEARNING APPROACHES FOR HIGH-PRECISION CONCRETE STRENGTH PREDICTION IN INDUSTRIAL CEMENT PRODUCTION

G. Vidyulatha¹, S. Ajay², V. Snehith rao², Sk. Arif², K. Sai Krishna²

¹Assistant Professor, ²UG Student, ^{1,2}Department of Computer Science and Engineering (CSBS),

^{1,2}Sree Dattha Institute of Engineering and Science, Sheriguda, Ibrahimpatnam, 501510, Telangana.

To Cite this Article

G. Vidyulatha, S. Ajay, V. Snehith Rao, Sk. Arif, K. Sai Krishna, "Machine Learning Approaches For High-Precision Concrete Strength Prediction In Industrial Cement Production", Journal of Science Engineering Technology and Management Science, Vol. 02, Issue 08, August 2025, pp: 556-566, DOI: <http://doi.org/10.64771/jsetms.2025.v02.i08.pp556-566>

Submitted: 15-07-2025

Accepted: 21-08-2025

Published: 28-08-2025

ABSTRACT

In India, infrastructure plays a pivotal role in economic development, yet building collapses have increased due to structural weaknesses. According to the National Crime Records Bureau, over 1,200 building collapses were reported annually between 2015 and 2020. Traditionally, civil engineers rely on manual formulas, which lack precision. The objective is to predict the compressive strength of buildings. Before machine learning, building strength was predicted using empirical formulas, standard material codes (e.g., IS 456:2000), and engineering judgment based on concrete mix ratios, curing time, and load tests. These methods involved trial-and-error concrete mix design, compressive strength tests (e.g., cube testing), and regression through manually plotted graphs. Traditional systems are time-consuming, cost-intensive, and less adaptable to varying construction scenarios. They lack predictive capabilities, struggle with large datasets, and ignore hidden patterns between influencing variables. These drawbacks lead to inaccurate assessments, delayed results, and increased risks in structural safety decisions. This research aims to overcome traditional limitations by introducing machine learning for predictive accuracy. ML models like XGBoost can handle non-linear data, detect complex correlations, reduce human error, and generate quick, real-time predictions. These improvements significantly enhance safety, cost-efficiency, and performance monitoring in modern building projects. The proposed system leverages machine learning models to predict the compressive strength of building materials based on historical data inputs such as cement, slag, ash, water, superplasticizer, coarse/fine aggregates, and curing age. By applying XGBoost Regressor and Random Forest Regressor, the system learns from existing datasets to predict outcomes with high accuracy. Compared to traditional methods, this automated approach provides faster results, handles large datasets effectively, and reduces manual intervention. XGBoost, in particular, yields superior performance making it ideal for real-world deployment in construction sites, improving quality assurance, and aiding in smart civil engineering practices.

Keywords: XGBoost, Random Forest, civil engineering, construction safety, predictive modeling, data-driven approach, building materials, infrastructure, India, quality assurance, smart construction, structural integrity.

This is an open access article under the creative commons license
<https://creativecommons.org/licenses/by-nc-nd/4.0/>



1. INTRODUCTION

Concrete is one of the most widely used construction materials globally, with its compressive strength being a critical parameter for ensuring its structural integrity. In the cement manufacturing industry, accurately predicting the compressive strength of concrete is essential to maintain quality, optimize material usage, and reduce production costs. Traditional methods for predicting concrete compressive strength—relying on empirical formulas and time-consuming physical testing—are increasingly inadequate in the face of modern demands for speed, accuracy, and efficiency in cement manufacturing. These methods fail to account for complex, real-world variables such as ingredient variability, curing conditions, and environmental influences. This project addresses the critical need for a faster, data-driven solution by integrating machine learning (ML) to predict concrete strength in real-time, reducing production delays, material waste, and testing costs. Motivated by India's rapid urbanization and its \$1.4 trillion infrastructure vision by 2030, the research aims to modernize traditional practices, cut 28-day curing delays, and improve sustainability by optimizing material use—particularly important in an industry responsible for 8% of global CO₂ emissions.

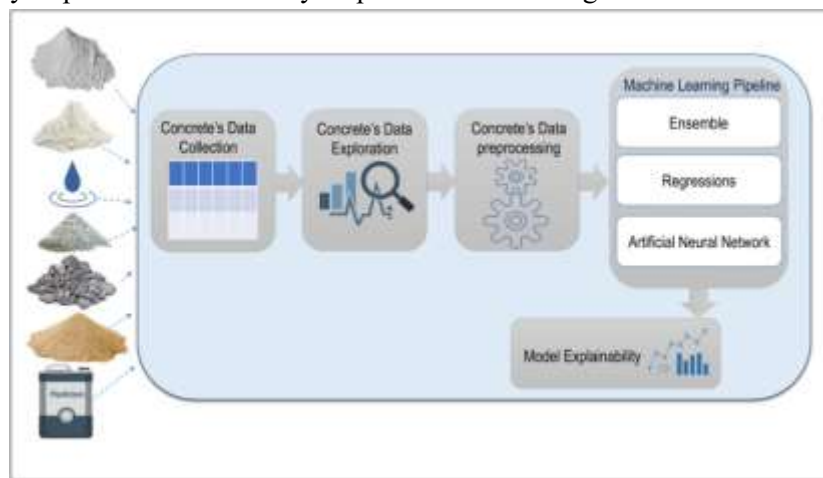


Figure 1: Generalised block diagram.

With applications spanning cement manufacturing, construction quality control, material optimization, and R&D, this ML-driven approach enhances structural reliability, reduces environmental impact, and supports smarter, greener infrastructure development nationwide.

2. LITERATURE SURVEY

Concrete is the most frequently used construction material globally due to its versatility, durability, and cost-effectiveness [1]. Its mechanical properties, particularly compressive strength, are critical for ensuring the safety and longevity of structures. Accurate prediction of concrete's compressive strength is essential for mix design optimization, quality control, and compliance with engineering standards [2]. Traditional empirical methods for estimating compressive strength often involve extensive laboratory testing and simplistic models that may not capture the complex interactions among the multitude of variables in concrete mixtures. This complexity has led researchers to explore advanced computational techniques, particularly machine learning (ML), to model and predict concrete behavior more accurately [3]. In recent years, ML algorithms have gained prominence in civil engineering applications due to their ability to model nonlinear relationships and handle large datasets. These algorithms learn patterns from historical data and can make accurate predictions based on input features, which makes them suitable for predicting the properties of various concrete types, including those modified with Supplementary Materials such as fly ash, nano-silica, recycled aggregates, and other industrial by-products. Several studies have applied ML models to predict concrete compressive strength with notable success. Alghrairi et al. [4] developed nine ML models to estimate the compressive strength of lightweight concrete modified with nanomaterials. Among these,

the gradient-boosted trees (GBT) model outperformed others by achieving a coefficient of determination (R^2) of 0.90 and a root mean square error (RMSE) of 5.286 MPa. The study highlighted that water content was the most influential factor affecting compressive strength predictions and emphasized the critical role of the water-to-cement ratio in concrete mix design. Similarly, Ding et al. [5] investigated ML models to predict the compressive strength of alkali-activated cementitious materials using solid waste components. They employed six ML algorithms, including support vector machine (SVM), random forest (RF), radial basis function neural network (RBF), and long short-term memory network (LSTM). The SVM model achieved the highest performance with an R^2 of 0.9054 and a normalized root mean square error of 0.0997. In addition to the evaluation of prediction accuracy, feature importance analysis using SHapley Additive exPlanations (SHAP) revealed key influencing factors such as calcium oxide content, water-to-binder ratio, silicon dioxide content, modulus of water glass, and aluminum trioxide content. Ekanayake et al. [6] addressed the “black-box” nature of ML models by employing SHAP to interpret predictions of concrete compressive strength. Utilizing tree-based algorithms including XGBoost and light gradient boosting machine (LGBM), they achieved high accuracy with an R-value of 0.98. The SHAP analysis provided insights into feature importance and confirmed that age and cement content were the most influential features. This approach demonstrated that ML models could capture complex relationships among variables and lead to enhanced trust among domain experts.

Despite these advancements, a persistent limitation in the existing literature is the inadequate exploration of feature interactions and their cumulative impact on model predictions. Most studies emphasize achieving high predictive accuracy without thoroughly investigating how input variables interact within the models. For instance, Paudel et al. [7] compared the performance of non-ensemble and ensemble ML models in predicting the compressive strength of concrete containing fly ash. The study identified age, cement content, and water content as the most influential features but lacked a comprehensive analysis of feature interactions. Similarly, Song et al. [8] employed ML algorithms, including gene expression programming (GEP), artificial neural network (ANN), decision tree (DT), and bagging regressor, to predict the compressive strength of concrete with fly ash admixture. While the study confirmed that the selection of input parameters and regressors significantly affects the accuracy of predicted outcomes, it did not extensively explore feature interactions. Tran et al. [9] evaluated the compressive strength of concrete made with recycled concrete aggregates using six ML models. The GB_PSO model achieved the highest prediction accuracy with an R^2 of 0.9356. Feature importance analysis revealed that cement content and water content were the most important factors affecting compressive strength. However, the study primarily focused on individual feature importance rather than the interactions between variables. Ahmad et al. [10] compared supervised ML algorithms, including ANN, AdaBoost, and boosting, to predict the compressive strength of geopolymer concrete containing high-calcium fly ash. This study demonstrated the potential of ensemble methods in capturing complex patterns in data, which can lead to more accurate predictions. Nevertheless, it did not explore the interactions among input features. Anjum et al. [11] applied ensemble ML methods, including gradient boosting, RF, bagging regressor, and AdaBoost regressor, to estimate the compressive strength of fiber-reinforced nano-silica modified concrete. SHAP analysis revealed that the coarse aggregate to fine aggregate ratio had a stronger negative correlation with compressive strength, while specimen age positively affected it. The study highlighted the importance of considering the interaction and effects of input parameters but did not provide a detailed feature interaction analysis. Ullah et al. [12] predicted the compressive strength of sustainable foam concrete using individual and ensemble ML approaches, including SVM, RF, bagging, boosting, and a modified ensemble learner. The study suggested that ensemble learners significantly enhance the performance and robustness of ML models but did not explore feature interactions in depth. Moreover, Kumar and Pratap [13] investigated the use of ML models to predict the compressive

strength of high-strength concrete and focused on the influence of superplasticizer, sand, and water content. The study acknowledged the significant influence of superplasticizer on compressive strength but lacked a comprehensive analysis of feature interactions. Nguyen et al. [14] proposed a machine learning approach using multivariate polynomial regression and automated feature engineering to predict the compressive strength of ultra-high-performance concrete (UHPC). While this study provided insights into feature interactions, it was specific to UHPC and did not address broader concrete types.

3. PROPOSED SYSTEM

The proposed project for predicting concrete compressive strength using machine learning begins with uploading a structured dataset—typically from sources like the Yeh concrete dataset—containing inputs such as cement, water, aggregates, and curing age, along with compressive strength values. The data undergoes preprocessing to handle missing values, remove duplicates, and apply feature scaling using standardization to ensure consistent input ranges. Exploratory Data Analysis (EDA) follows, employing statistical summaries, visualizations, and correlation heatmaps to uncover patterns and key predictors influencing strength. A baseline prediction is established using the Random Forest Regressor (RFR), which captures nonlinear relationships and is evaluated using metrics like MAE, MSE, and R^2 . Building on this, the XGBoost Regressor (XGBR) is proposed as a more advanced model, leveraging gradient boosting to refine accuracy and capture complex variable interactions more effectively. Both models are compared using standardized metrics and visualizations, often showing XGBR's superior performance in reducing prediction error. Finally, the trained XGBR model is applied to new test data, offering real-time predictions of compressive strength, demonstrating practical utility for cement manufacturers by enabling faster, cost-effective, and more accurate quality control in concrete production.

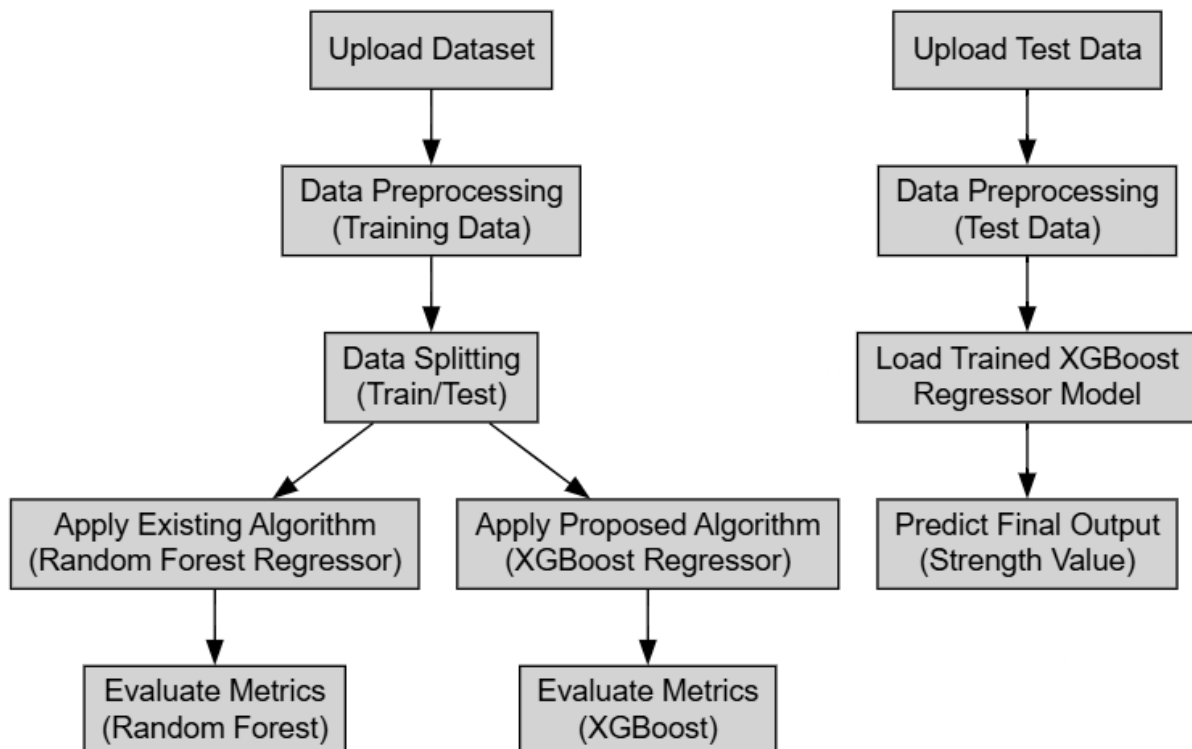


Figure 2: Architectural block diagram of the proposed system.

The data processing phase begins by loading the raw concrete compressive strength dataset into a Pandas DataFrame, followed by checks for missing values and duplicate rows. While this dataset typically contains no null entries, the process ensures data integrity by removing any duplicates—identical rows of all input features and target values—that could bias the model during training. The

dataset is then split into predictors (cement, slag, fly ash, water, superplasticizer, aggregates, and age) stored in X, and the target variable (compressive strength in MPa) stored in y. To ensure compatibility with regression models like Linear, Lasso, Ridge, and SVR, a StandardScaler is applied, standardizing the features to have zero mean and unit variance, thus preventing any feature with a larger numerical scale from disproportionately influencing model learning. Finally, the dataset is split into training and testing subsets using an 80-20 ratio via `train_test_split`. This approach ensures that 80% of the data is used for model training (X_train, y_train) and 20% is reserved for performance evaluation (X_test, y_test), with automatic shuffling applied to maintain randomness and prevent any ordering bias—ensuring a fair assessment of the model's generalization capabilities.

Proposed Algorithm – XGBoost Regressor

The **XGBoost Regressor** (Extreme Gradient Boosting Regressor) is an advanced implementation of gradient boosting designed for high performance, efficiency, and scalability. It is particularly well-suited for structured/tabular data and is widely used in regression, classification, and ranking problems. XGB Regressor builds an ensemble of decision trees sequentially, where each new tree corrects the errors made by the previous trees by focusing more on the residuals. It utilizes both first and second-order derivatives (gradient and hessian) of the loss function to make optimized decisions at each split, leading to faster convergence and improved accuracy. XGBoost includes regularization terms in its objective function, which helps prevent overfitting and enhances generalization compared to traditional boosting algorithms. In this project, XGB Regressor is employed as the proposed algorithm for predicting concrete compressive strength more accurately.

How XGB Regressor Works

XGB Regressor works by adding trees one at a time, where each tree is trained on the residuals (errors) of the predictions made so far. Initially, the model starts with a base prediction (usually the mean of the target values). Then, at each iteration, it fits a new tree to the negative gradient of the loss function (i.e., the residuals). This process continues for a specified number of iterations or until convergence. The final prediction is the sum of the outputs of all individual trees. Unlike simple gradient boosting, XGBoost applies advanced regularization (L1 and L2), parallel computation for efficiency, and sophisticated handling of missing values, making it robust and faster.

Architecture of XGB Regressor

1. **Input Layer** – Receives multiple numerical features from the dataset (e.g., cement, water, etc.).
2. **Initial Prediction** – Starts with a base prediction, often the mean of target values.
3. **Gradient Calculation** – Computes the gradient (error/residuals) of the loss function.
4. **Tree 1 Construction** – Builds the first decision tree to fit the residuals from initial prediction.
5. **Prediction Update** – Updates prediction by adding the output of Tree 1 (weighted by learning rate).
6. **New Residual Calculation** – Calculates new residuals using updated predictions.
7. **Next Tree Construction** – Builds Tree 2 to fit new residuals and adds it to the model.
8. **Iterative Learning** – Repeats steps 5 to 7 for a set number of boosting rounds.
9. **Regularization** – Applies L1/L2 regularization to avoid overfitting.
10. **Final Output** – Final prediction is the sum of all trees' contributions.

Advantages of XGB Regressor

The XGB Regressor has several key advantages that make it superior to many traditional machine learning models. Firstly, it delivers high predictive accuracy due to its robust gradient boosting framework that continuously corrects errors. XGBoost includes regularization techniques (both L1 and L2), which prevent overfitting and make the model generalize better on unseen data. It is highly

efficient and scalable, with support for parallelization and distributed computing, allowing it to train quickly even on large datasets. Moreover, it has built-in handling of missing values, which is particularly useful in real-world data scenarios. The algorithm also includes cross-validation during training, which enhances model reliability. XGBoost supports various custom objective functions and evaluation metrics, making it highly flexible for different kinds of regression problems. Its tree pruning, column subsampling, and early stopping features help in optimizing performance and computational cost. These advantages make it particularly suitable for the proposed solution to predict concrete strength more accurately and efficiently compared to existing models like Random Forest.

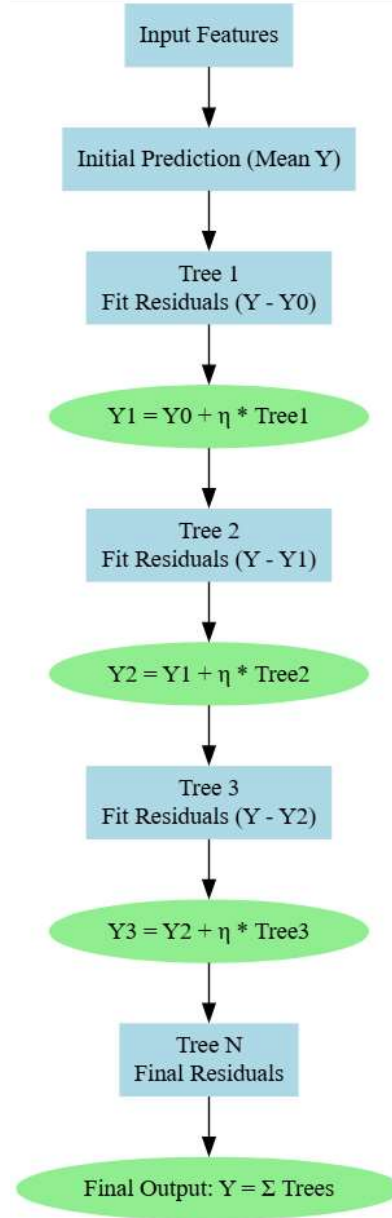


Figure 3: Working of XG boost regressor.

4. RESULTS

Figure 4 presents a tabular view of the dataset used for predicting concrete compressive strength. It includes columns such as cement, slag, flyash, water, superplasticizer, coarseaggregate, fineaggregate, age, and csMPa, with each row representing a unique concrete mix. The values, measured in units like kg/m³ for materials and days for age, provide the raw data for analysis. This visual establishes the dataset's structure, serving as the foundation for all subsequent steps in the project.

	cement	slag	flyash	water	superplasticizer	coarseaggregate	fineaggregate	age	csMPa
0	540.0	0.0	0.0	162.0	2.5	1040.0	676.0	28	79.99
1	540.0	0.0	0.0	162.0	2.5	1055.0	676.0	28	61.89
2	332.5	142.5	0.0	228.0	0.0	932.0	594.0	270	40.27
3	332.5	142.5	0.0	228.0	0.0	932.0	594.0	365	41.05
4	198.6	132.4	0.0	192.0	0.0	978.4	825.5	360	44.30

Fig. 4: Sample Dataset.

	count	mean	std	min	25%	50%	75%	max
cement	1030.0	281.167864	104.506364	102.00	192.375	272.900	350.000	540.0
slag	1030.0	73.895825	86.279342	0.00	0.000	22.000	142.950	359.4
flyash	1030.0	54.188350	63.997004	0.00	0.000	0.000	118.300	200.1
water	1030.0	181.567282	21.354219	121.80	164.900	185.000	192.000	247.0
superplasticizer	1030.0	6.204660	5.973841	0.00	0.000	6.400	10.200	32.2
coarseaggregate	1030.0	972.918932	77.753954	801.00	932.000	968.000	1029.400	1145.0
fineaggregate	1030.0	773.580485	80.175980	594.00	730.950	779.500	824.000	992.6
age	1030.0	45.662136	63.169912	1.00	7.000	28.000	56.000	365.0
csMPa	1030.0	35.817961	16.705742	2.33	23.710	34.445	46.135	82.6

Fig. 5: Dataset Description.

Figure 5 summarizes the statistical properties of the dataset's columns, including cement, water, and csMPa. It displays metrics such as mean, standard deviation, minimum, and maximum values for each feature. For example, it shows average cement content in kg/m³ and typical strength in MPa across samples. This description quantifies the dataset's range and distribution, offering a clear overview of its numerical characteristics.

This below figure 6 illustrates the exploratory data analysis through visualizations like pair plots, distribution graphs, and a correlation heatmap. It reveals relationships between variables—such as cement and csMPa—and highlights distributions, like the skewness of age. The heatmap quantifies correlations, emphasizing strong predictors of strength. This comprehensive analysis guides feature selection and model development by exposing data patterns.

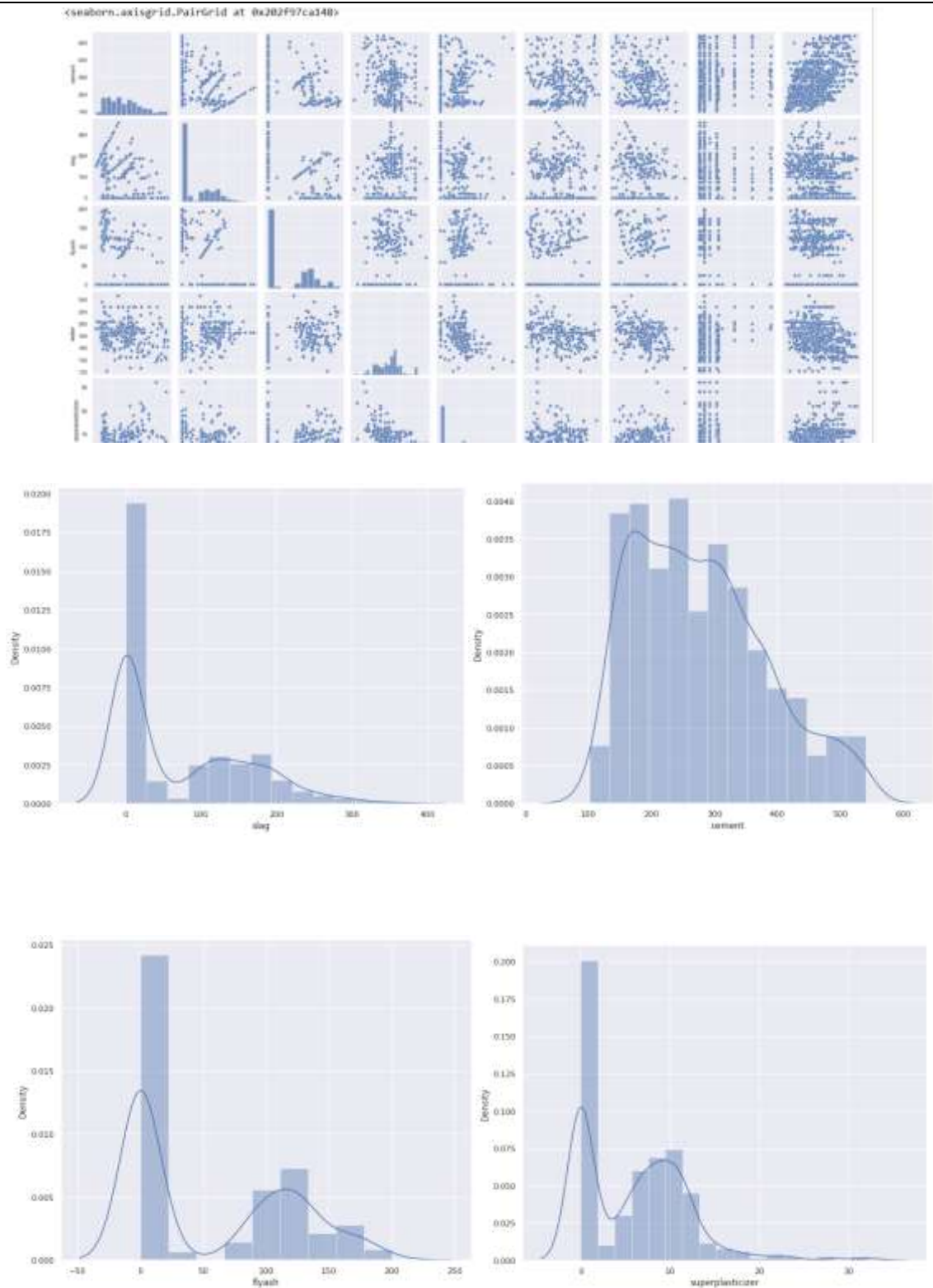


Fig. 6: EDA plots.

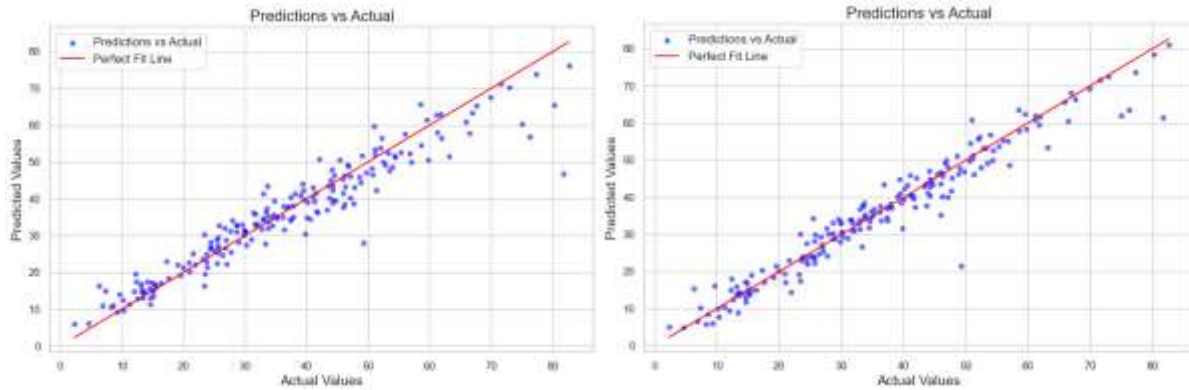


Fig. 7: Scatter plot for Actual vs Prediction on RFR and XGB Models.

Figure 7 displays two scatter plots comparing actual compressive strength (csMPa) against predictions from Random Forest Regressor (RFR) and XGBoost (XGB) models. Each plot shows data points along a diagonal line, with XGB clustering tighter than RFR, indicating higher accuracy. The X-axis represents actual values, and the Y-axis shows predicted values from test data. This visual comparison underscores XGB's superior performance over RFR in predicting concrete strength.

Model	MAE	MSE	r2 Score	RMSE (Cross Validated)
XGB Regressor	2.631153	17.495196	0.941355	9.324258
Random Forest Regressor	3.516771	27.372299	0.908247	10.080081

Table 1: Performance metrics of RFR and XGBoost Models.

This table 1 lists performance metrics for Random Forest Regressor (RFR) and XGBoost (XGB) models, including MAE, MSE, R^2 Score, and RMSE (Cross Validated). XGB records MAE of 2.631153, MSE of 17.495196, R^2 of 0.941355, and RMSE of 9.324258, while RFR shows MAE of 3.516771, MSE of 27.372299, R^2 of 0.908247, and RMSE of 10.080081. These values quantify XGB's lower error and higher fit compared to RFR. The table provides a concise benchmark for model evaluation.

	cement	slag	flyash	water	superplasticizer	coarseaggregate	fineaggregate	age	Predicated
0	540.0	0.0	0	162	2.5	1040.0	676.0	28	49.491047
1	540.0	0.0	0	162	2.5	1055.0	676.0	28	49.491047
2	332.5	142.5	0	228	0.0	932.0	594.0	270	46.377384
3	332.5	142.5	0	228	0.0	932.0	594.0	365	46.377384
4	198.6	132.4	0	192	0.0	978.4	825.5	360	46.377384
5	266.0	114.0	0	228	0.0	932.0	670.0	90	46.377384
6	380.0	95.0	0	228	0.0	932.0	594.0	365	46.377384
7	380.0	95.0	0	228	0.0	932.0	594.0	28	46.377384
8	266.0	114.0	0	228	0.0	932.0	670.0	28	46.377384

Fig. 8: Proposed XG Boost Model Prediction on test data.

Figure 8 showcases the XGBoost model's predictions on a separate test dataset, presented as a table or graph. It includes input features like cement and water alongside predicted csMPa values in megapascals. The output aligns with actual strength requirements, demonstrating practical applicability in cement manufacturing. This visual confirms the model's effectiveness in real-world scenarios.

5. CONCLUSION

The research successfully implements a machine learning-based system to predict concrete compressive strength, leveraging the XGBoost Regressor as the superior model over the Random

Forest Regressor. The dataset, comprising features like cement, water, and age, undergoes thorough preprocessing and exploratory analysis, revealing critical predictors of strength. XGBoost achieves a lower mean absolute error of 2.631153 MPa, mean squared error of 17.495196, and an R^2 score of 0.941355, outperforming Random Forest's metrics of 3.516771, 27.372299, and 0.908247, respectively. Cross-validated RMSE further confirms XGBoost's precision at 9.324258 compared to Random Forest's 10.080081. This system delivers accurate, efficient predictions, eliminating the delays of traditional 28-day testing and optimizing quality control in cement manufacturing. The implementation proves effective for real-time strength assessment, reducing production costs and enhancing material performance in the industry.

REFERENCE

- [1] Griffiths, S.; Sovacool, B.K.; Furszyfer Del Rio, D.D.; Foley, A.M.; Bazilian, M.D.; Kim, J.; Uratani, J.M. Decarbonizing the Cement and Concrete Industry: A Systematic Review of Socio-Technical Systems, Technological Innovations, and Policy Options. *Renew. Sustain. Energy Rev.* 2023, *180*, 113291.
- [2] Young, B.A.; Hall, A.; Pilon, L.; Gupta, P.; Sant, G. Can the Compressive Strength of Concrete Be Estimated from Knowledge of the Mixture Proportions? New Insights from Statistical Analysis and Machine Learning Methods. *Cem. Concr. Res.* 2019, *115*, 379–388.
- [3] Li, Z.; Yoon, J.; Zhang, R.; Rajabipour, F.; Srubar, W.V., III; Dabo, I.; Radlińska, A. Machine Learning in Concrete Science: Applications, Challenges, and Best Practices. *npj Comput. Mater.* 2022, *8*, 127.
- [4] Alghairi, N.S.; Aziz, F.N.; Rashid, S.A.; Mohamed, M.Z.; Ibrahim, A.M. Machine Learning-Based Compressive Strength Estimation in Nanomaterial-Modified Lightweight Concrete. *Open Eng.* 2024, *14*, 20220604.
- [5] Ding, Y.; Wei, W.; Wang, J.; Wang, Y.; Shi, Y.; Mei, Z. Prediction of Compressive Strength and Feature Importance Analysis of Solid Waste Alkali-Activated Cementitious Materials Based on Machine Learning. *Constr. Build. Mater.* 2023, *407*, 133545.
- [6] Ekanayake, I.U.; Meddage, D.P.P.; Rathnayake, U. A Novel Approach to Explain the Black-Box Nature of Machine Learning in Compressive Strength Predictions of Concrete Using Shapley Additive Explanations (SHAP). *Case Stud. Constr. Mater.* 2022, *16*, e01059.
- [7] Paudel, S.; Pudasaini, A.; Shrestha, R.K.; Kharel, E. Compressive Strength of Concrete Material Using Machine Learning Techniques. *Clean. Eng. Technol.* 2023, *15*, 100661.
- [8] Song, H.; Ahmad, A.; Farooq, F.; Ostrowski, K.A.; Maślak, M.; Czarnecki, S.; Aslam, F. Predicting the Compressive Strength of Concrete with Fly Ash Admixture Using Machine Learning Algorithms. *Constr. Build. Mater.* 2021, *308*, 125021.
- [9] Quan Tran, V.; Quoc Dang, V.; Si Ho, L. Evaluating Compressive Strength of Concrete Made with Recycled Concrete Aggregates Using Machine Learning Approach. *Constr. Build. Mater.* 2022, *323*, 126578.
- [10] Ahmad, A.; Ahmad, W.; Chaiyasarn, K.; Ostrowski, K.A.; Aslam, F.; Zajdel, P.; Joyklad, P. Prediction of Geopolymer Concrete Compressive Strength Using Novel Machine Learning Algorithms. *Polymers* 2021, *13*, 3389.
- [11] Anjum, M.; Khan, K.; Ahmad, W.; Ahmad, A.; Amin, M.N.; Nafees, A. Application of Ensemble Machine Learning Methods to Estimate the Compressive Strength of Fiber-Reinforced Nano-Silica Modified Concrete. *Polymers* 2022, *14*, 3906.
- [12] Ullah, H.S.; Khushnood, R.A.; Farooq, F.; Ahmad, J.; Vatin, N.I.; Ewais, D.Y.Z. Prediction of Compressive Strength of Sustainable Foam Concrete Using Individual and Ensemble Machine Learning Approaches. *Materials* 2022, *15*, 3166.

- [13] Kumar, P.; Pratap, B. Feature Engineering for Predicting Compressive Strength of High-Strength Concrete with Machine Learning Models. *Asian J. Civ. Eng.* 2024, 25, 723–736.
- [14] Nguyen, N.-H.; Abellán-García, J.; Lee, S.; Vo, T.P. From Machine Learning to Semi-Empirical Formulas for Estimating Compressive Strength of Ultra-High-Performance Concrete. *Expert Syst. Appl.* 2024, 237, 121456.