

Statistical Analysis of Slump Flow Using Gene Expression Programming (GEP) for Self-Consolidated Concrete

Maaz Khan ^{1*}, Asad Wahab ¹, Muhammad Zikria Luqman¹, Muhammad Atta Ur Rehman¹, Waheed Ali Khoso¹, Muhammad Arsalan Khan¹, Touqeer Ali Rind¹

¹ Department of Civil Engineering, Ghulam Ishaq Khan Institute of Engineering Sciences and Technology, Topi, KPK, Pakistan. Email: gcv2463@giki.edu.pk; gcv2520@giki.edu.pk; gcv2466@giki.edu.pk; gcv2464@giki.edu.pk; waheed.ali@giki.edu.pk; arsalan.khan@giki.edu.pk; touqeer.ali@giki.edu.pk,

* Correspondence: gcv2463@giki.edu.pk

Abstract

Statistical analysis of the slump flow prediction by the application of Gene Expression Programming on the data points of 953 related to Self-Consolidated Concrete. SCC is the acronym for Self-Consolidated Concrete, high-flow concrete, supposed to self-consolidate. This type of concrete really demonstrates its superb applicability especially in the dense and complex structures with reinforcement. Several design variables, namely the water-to-cement ratio, aggregate properties, and admixtures affect this characteristic; therefore, making it rather difficult to predict precisely its slump flow behavior. GEP was applied to analyze a dataset obtained from a sequence of experiments and yielded a predictive model for slump flow. Descriptive analysis tools, regression techniques, as well as error metrics, MSE, RMSE, and R², have been used to test the robustness and reliability of the model. The results have also indicated that, based solely on the mixed design parameters, it is possible to predict slump flow by GEP. Significant relationships have also been found between values of slump flow and input factors.

Keywords: Slump Flow Prediction, Gene Expression Programming (GEP), Self-Consolidated Concrete (SCC), Statistical Analysis, Predictive Modeling

1. Introduction

Modern construction has made Self-Consolidating Concrete important, as it flows under its own weight and fills intricate formwork without the application of mechanical vibration. This aspect of SCC makes it more valuable for difficult constructions, including bridges, high-rise structures, and precast units, where conventional vibrated concrete cannot be applied. Prediction of slump flow is considered the main challenge in SCC application. It is also one of the critical parameters controlling workability and fluidity in the mix. The predicted slump flow is very useful to achieve the optimization of mix design to meet the required specifications on performance [1]. This means the flowability of the concrete, when it is allowed to spread freely in a conical mold, is a function of many parameters ranging from water-to-cement ratios and aggregate gradations to admixture contents and environmental conditions like temperature and humidity. Along with inherent complexities of behavior, these make advanced modeling techniques vital to predict slump flow exactly. Traditionally, empirical methods and very basic regression models have been used for slump flow prediction, but these models often cannot capture the complex nonlinear relationships between variables involved [2]. It was evident from the latest research studies that more precise

predictive models, capable of grasping dynamic interactions between mix design parameters and environmental factors, were needed.

Gene Expression Programming (GEP) is one of the high-performance techniques in machine learning, which has attracted the interest of civil engineering scientists for the modeling of complex systems. GEP essentially is a genetic algorithm approach to which mathematical expressions or models evolve as computer programs and, thus is particularly suited to jobs such as those related to the prediction of properties of concrete. Unlike the traditional regression models, GEP can accommodate nonlinearities and interactions between multiple input variables, providing more accurate and robust predictions. Some successful applications of GEP have been found in predicting different properties of concrete, including compressive strength, durability, and shrinkage [3], [4]. There has been increasing concern in recent years to focus specifically on GEP when predicting the slump flow of SCC. While statistical regression models or ANNs have been used within the traditional approach to predict workability of concrete, the inherent advantage of GEP lies within its ability to evolve interpretable and flexible symbolic models. Previous studies have shown that GEP can even surpass other modeling techniques in both predictive accuracy and interpretability, which makes it potentially a more attractive tool for engineers and researchers [5]. For example, Hamidi et al. proved that the slump flow and other SCC properties can be precisely predicted by machine learning models like GEP rather than traditional methods [6].

This study develops a predictive model for SCC slump flow through GEP and a statistical analysis of the outcome. In this research, a total of 953 experimental data points were applied in training and validation wherein the source of research was considered. All mix design parameters known to influence significantly the slump flow of SCC were considered in the dataset, including water-cement ratio, aggregate types, admixture content, and temperature. These are used to evaluate the performance of the GEP model using performance metrics such as Mean Squared Error, Root Mean Squared Error, and R-squared values. This study demonstrates the potential of GEP as a reliable tool to predict the slump flow of SCC by comparing the model with other conventional prediction methods.

2. Significance of Study

The main aim of this statistical analysis is the evaluation and comparison of predictive accuracy of the Gene Expression Programming (GEP) model and multiple regression analysis for predicting slump flow in self-consolidating concrete (SCC). Therefore, using multiple statistical tools correlated with the analysis like errors-the MAE, RMSE, and R²-permits testing the model. The primary goal here will be to compare whether a more realistic relationship among different mix parameters and slump flow can be achieved from models. Thus, there arises a more accurate, feasible model in optimizing mixes along with ensuring quality control processes, while the performance also gets better in its related construction application.

3. Methodology

This flowchart fig 1 represents a well-structured workflow for machine learning (ML), data-driven analysis. It involves the collection of relevant information through various sources for data collection followed by the actual development of ML models using data collected. In the next stages, the analysis of data can be done when the model has been built using the performance analysis of the built model and derived insights. Finally, the workflow ends with results in which the findings are given, with meaningful conclusions derived from analyzing the data. This step-by-step process ensures a systematic and efficient methodology about predictive modeling and decision-making.

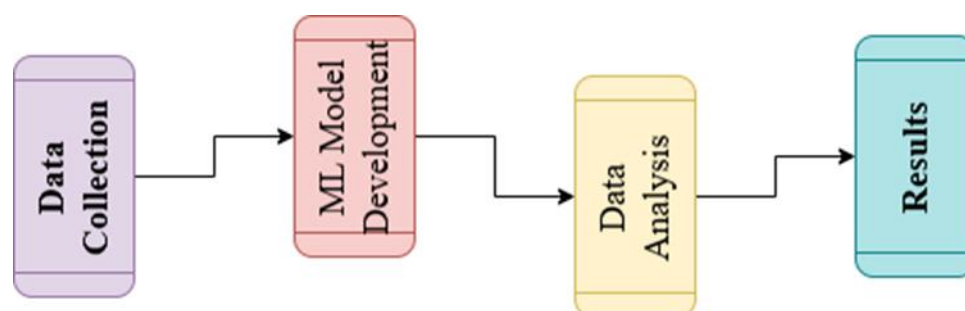


Figure 1. Flow Chart of Methodology

3.1 Data collection

Between 2014 to 2024, 76 published publications yielded a comprehensive dataset of more than 953 SCC mixes [7]. Slump flow obtained only through empirical testing were not included. The specific gravity of the ingredients, which was used to calculate derived properties like paste volume, was one of the crucial mixture design parameters that were noted. Important new and rheological characteristics were collected.

3.2 GEP Model Development

By using Gene Expression Programming to develop a predictive model for slump flow in self-consolidating concrete, the mix parameters used include water-cement ratio, superplasticizer content, and proportions of fine aggregate. GEP is indeed an evolutionary algorithm that produces mathematical models through the systematic evolution of the structure and the parameters of expressions, so this makes GEP particularly effective in capturing non-linear relationships of complex interactions between the input variables and the output. Start generating a population of random mathematical expressions as chromosomes. These mathematical expressions are then evaluated in terms of prediction accuracy and the best-performing ones undergo selection, crossover, and mutation to create generations. This process continues in an iterative manner until it identifies the most accurate model for prediction. GEP is a powerful tool for modeling nonlinear interactions between variables and, therefore, can predict slump flow in SCC more effectively than linear regression methods. The flexibility and accuracy of the GEP model have been demonstrated in several studies, where it has been successfully applied to predict concrete properties like compressive strength and workability [8].

3.3 Data Analysis

3.3.1 Descriptive Analysis

One of the most advanced, adaptable, and popular methods in the field of sensory analysis is descriptive analysis [9]. Data exploration basically consists of descriptive analysis to summarize and interpret major attributes of a dataset toward in-depth understanding of hidden underlying patterns. In the research work, descriptive analysis is conducted on the various mixes of parameters like the water-cement ratio, percentage of superplasticizer, and fine aggregates on SCC with slump flow value. The objective was to understand the distribution, central tendency, and variability of

these variables. Table 1 shows the descriptive statistics of different concrete mix parameters, namely cement, total powder (TP), fine aggregate (FA), coarse aggregate (CA), water, admixture (Adm), maximum grain size (MGS), and slump. This table provides mean, median, mode, standard deviation, variance, skewness, and kurtosis and minimum, maximum, and range values. These figures help in understanding distribution, variability, and trends in the dataset. They can, therefore, be helpful for optimizing concrete mix design and analyzing its properties.

Table 1. Descriptive Analysis

	Cement	TP	FA	CA	Water	Adm	MGS	Slump
Mean	377.91	517.24	853.93	795.90	183.93	1.76	16.07	679.63
Standard Error	3.35	2.91	3.37	3.25	0.78	0.06	0.10	2.08
Median	382.25	520.00	861.80	804.00	180.00	1.35	16.00	689.89
Mode	500.00	550.00	875.00	750.00	176.00	0.37	16.00	700.00
Standard Deviation	103.25	89.69	104.08	100.20	24.20	1.76	3.16	64.05
Sample Variance	10659.54	8043.63	10833.47	10040.58	585.72	3.10	9.97	4101.81
Kurtosis	-0.36	3.34	0.63	0.41	2.90	12.23	-0.58	2.04
Skewness	0.06	0.91	-0.05	-0.43	0.55	3.21	-0.35	-0.90
Range	598.70	672.00	831.00	782.00	229.90	12.74	17.00	500.00
Minimum	108.30	250.00	369.00	400.00	101.60	0.10	8.00	380.00
Maximum	707.00	922.00	1200.00	1182.00	331.50	12.84	25.00	880.00
Sum	359773.70	492407.86	812937.30	757692.67	175105.92	1671.39	15302.10	647005.79
Count	952.00	952.00	952.00	952.00	952.00	952.00	952.00	952.00
Largest (1)	707.00	922.00	1200.00	1182.00	331.50	12.84	25.00	880.00
Smallest (1)	108.30	250.00	369.00	400.00	101.60	0.10	8.00	380.00
Confidence Level(95.0%)	6.57	5.70	6.62	6.37	1.54	0.11	0.20	4.07

Figure 2 is a violin plot with error bars, displaying the distribution of various construction material parameters such as Cement TP, FA, CA, Water, Adm, MGS and Slump. The violin plots illustrate the density of the data, with color gradients representing different concentration levels. The error bars denote the minimum and maximum values (Min~Max), while the black dots indicate the mean values of each parameter. The spread of each violin plot suggests variability in data distribution while some parameters showing a wider range and some being more concentrated. This visualization effectively highlights the central tendency and dispersion of the dataset.

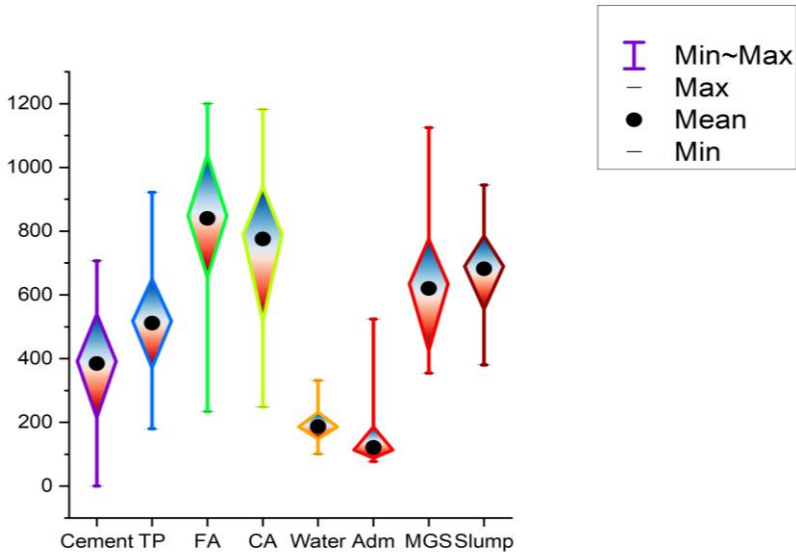


Figure 2. Box Plot of Descriptive Analysis

3.3.2 Correlation

A statistical technique for assessing the degree of association between two quantitative variables is correlation analysis. Whereas a weak correlation indicates that the variables are barely

related, a high correlation indicates that two or more variables have a substantial association [10]. A correlation matrix and a heat map shown in table 2 and figure 3 are useful tools in data analysis to understand relationships between variables. The correlation matrix provides numerical values indicating the strength and direction of linear relationships between pairs of variables, often ranging from -1 to 1, where -1 represents a perfect negative correlation, 0 shows no correlation, and 1 represents a perfect positive correlation. This matrix is visually represented using a gradient of colors, so patterns and significant relationships can be seen briefly. Together, they help in finding interdependencies among variables, highly correlated features, and guide decisions on feature selection or data preprocessing for predictive modeling.

Table 2. Correlation Matrix

	Cement	TP	FA	CA	Water	Adm	MGS	Slump
Cement	1	0.41777	-0.11831	-0.01473	0.12457	0.20519	0.03903	0.01321
TP	0.41777	1	-0.29438	-0.13573	-0.08167	0.29617	0.06029	0.07632
FA	-0.11831	-0.29438	1	-0.51598	-0.19853	-0.17397	-0.0378	-0.09974
CA	-0.01473	-0.13573	-0.51598	1	-0.17238	-0.15023	0.06748	-0.02747
Water	0.12457	-0.08167	-0.19853	-0.17238	1	-0.01874	-0.11882	-0.01864
Adm	0.20519	0.29617	-0.17397	-0.15023	-0.01874	1	-0.14521	-0.07665
MGS	0.03903	0.06029	-0.0378	0.06748	-0.11882	-0.14521	1	0.07617
Slump	0.01321	0.07632	-0.09974	-0.02747	-0.01864	-0.07665	0.07617	1

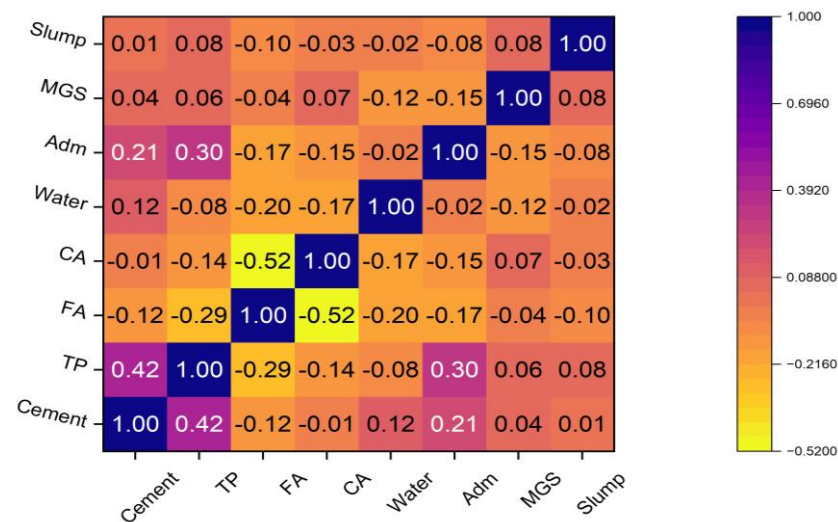


Figure 3. Heat Map

3.3.3 Scatter Metrix and Scatter Plot

One of the most popular and well-understood visual representations of bivariate data is a scatter plot. They can also be used with high-dimensional data by using the scatter plot matrix representation or dimensionality reduction [11]. A scatter plot shown in figure 4 are graphical tools to visualize the relationship between the variables in a data set. A scatter plot is a two-dimensional graph where individual data points are plotted along two axes, which represent two variables. It enables the discovery of some trends, patterns, or outliers and lets us understand the nature of the relation (linear, non-linear, or none) between variables.

A scatter matrix is an array of scatter plots that shows all pairwise relationships between all combinations of variables in a dataset. Each cell in the matrix displays the scatter plot of one variable against another, giving a comprehensive overview of the relationships in the dataset. Scatter matrices are especially useful in finding correlations, clusters, or trends across multiple variables at once.

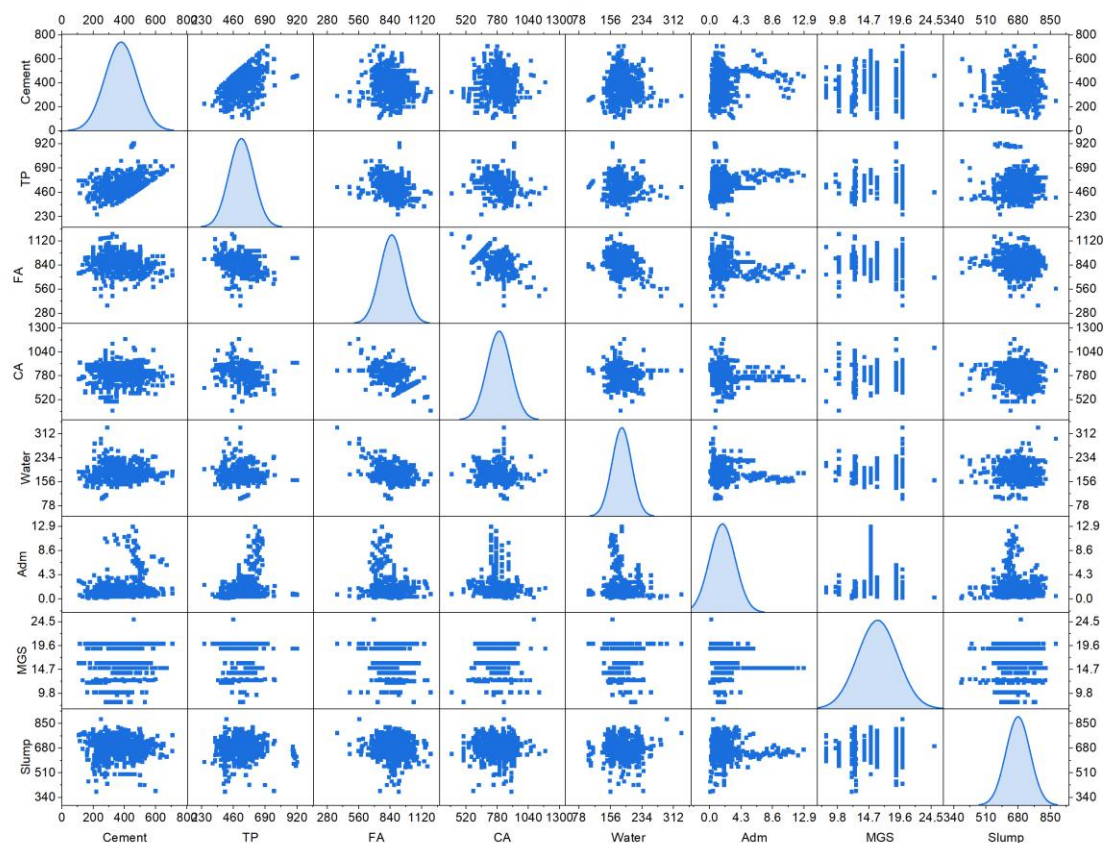


Figure 4. Scatter Plot

4. Results

Regression analysis between predicted values and actual values helps us determine how well the predictions made by a predictive model reflect the real relationship among the input variables and their outcome. In this process, the predicted values can be considered as the modeled outcomes, and actual values are the observed or actual outcomes from the data used. This analysis would show how good the performance is and how accurate the given model is.

4.1 Linear Regression

Linear regression is the statistical method used in establishing the relationship between predicted and actual value by fitting the best straight line through data points. It predicts a dependent variable based on some linear combination of predictors, with parameters representing the interception, slope, and error term. Homoscedasticity as well as normally distributed residuals are considered to have a linear relationship. It is widely used in any field to understand the variation between variables, trends as well as forecasting, though it offers interpretable results for decision-making. Linear Regression between GEP and slump shown in figure 5.

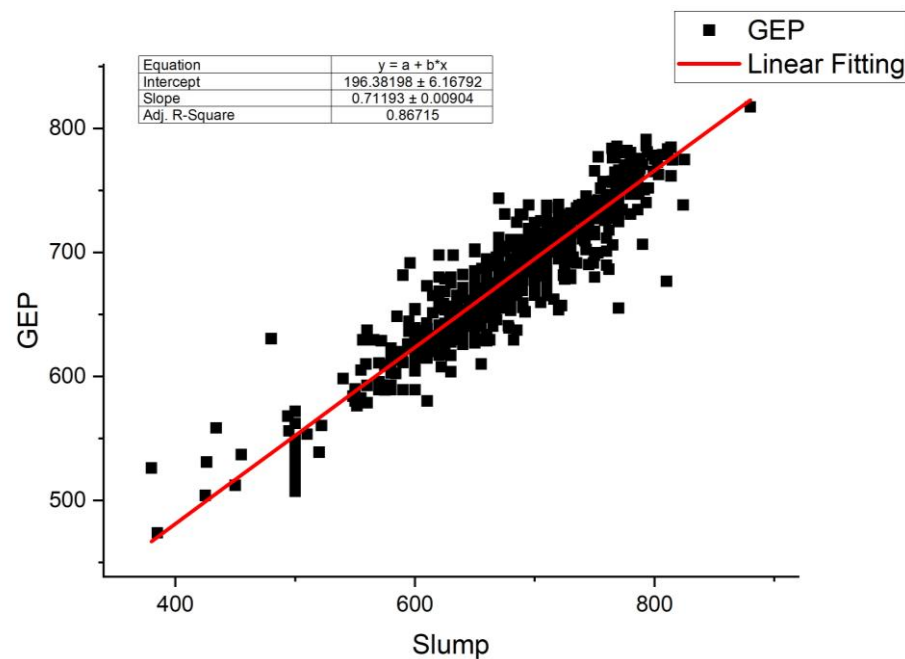


Figure 5. Linear Fitting

Top Plot (Actual vs. Predicted Plot)

Here is in figure 6 a plot of the actual values of the dependent variable. This could be slumping flow; the predicted values by a Gene Expression Programming model are also shown. The idea behind this plot is to visually assess whether the predicted values of the model tend to be in line with the actual observations. If the lines are closely overlapped, it would mean a good fit and accurate predictions by the model.

Bottom Plot (Residual Plot):

This scatter plot plots the residual (the difference between an actual and a predicted) against the order or index of observations. The motivation is to check the randomness in the distribution of residuals. Ideally residuals should be randomly scattered along zero, indicating that models do not have systematic error and assumptions like linearity and homoscedasticity are satisfied.

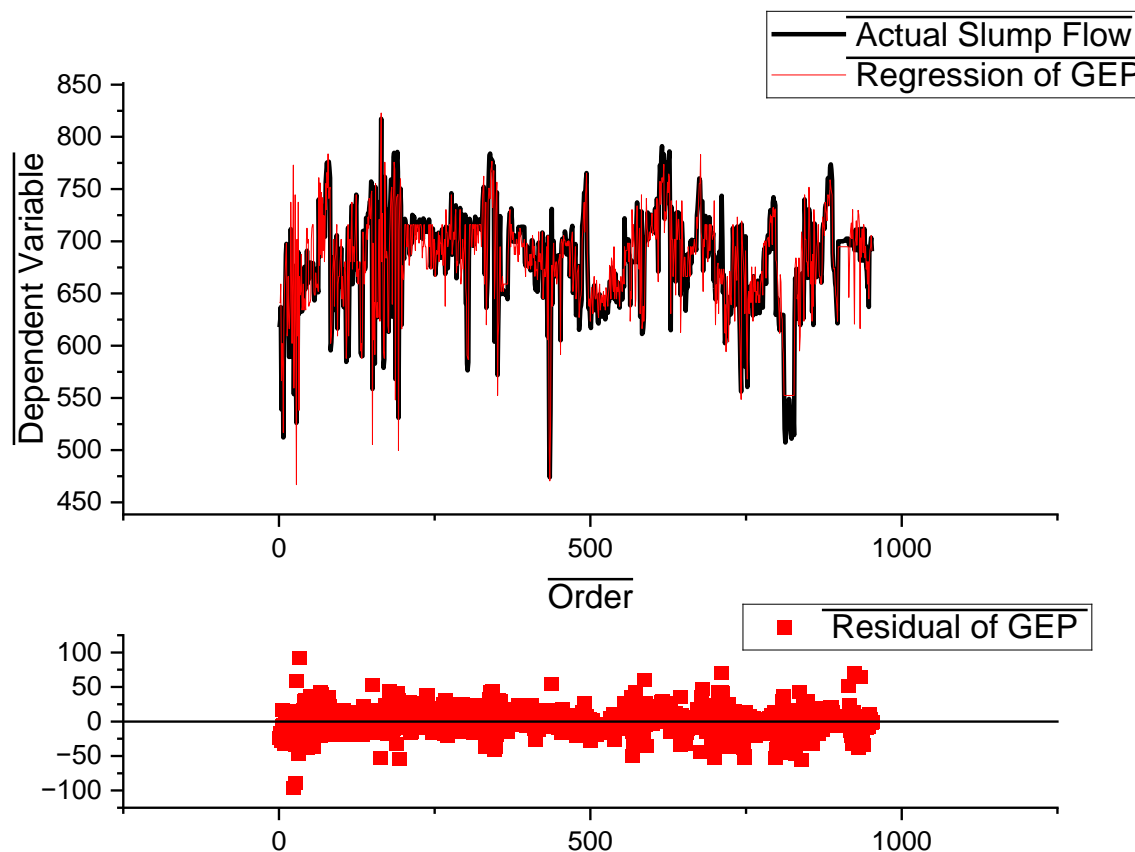


Figure 6. GEP and Slump Comparison and Residual Plot

4.2 Comparison of Models

A Taylor diagram shown in figure 7 the standard deviation SD and correlation coefficient R of multiple machine learning models in relation to observed data at a given time, which gives an ability to compare graphically their accuracy. How accurate each model is from the reference data, or experimental data, is plotted in the graphic while The R -value is the distance from the origin. The accuracy of models with points closer to the reference point is higher since they both show greater alignment with observed data.

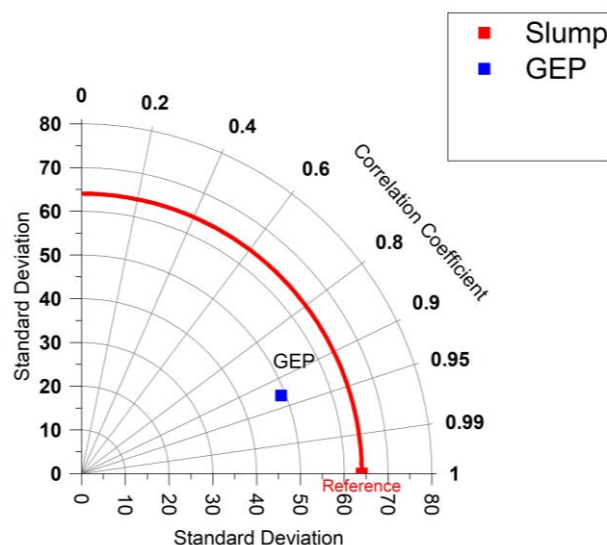


Figure 7. Taylor Diagram

4.3 Assessment of Model

4.3.1 Error Metrix

Table 3 show the results of evaluating the established models' prediction accuracy using various performance measures. The GEP model of slump flow prediction of Self-Consolidated Concrete (SCC) was found to be reliable with an R2 value of 0.87, which shows that the model explains 87% of the variability in slump flow data. The Root Mean Square Error (RMSE) value of 25.66 and Mean Absolute Error (MAE) value of 17.63 reflect the precision of the model in predicting slump flow values with comparatively minimal margins of error. The findings affirm the precision and reliability of the model in replicating the complex relationships among mix design parameters and slump flow behavior of SCC.

Table 3. Error Metrix

Metric	Value
R2	0.87
RMSE	25.66
MAE	17.63

4.3.2 T-Test

The t-test and S-test are both statistical methods used for hypothesis testing, but they serve different purposes and are based on different assumptions. The T-test is one type of statistical test used to compare the means of two groups. It is one of the most widely used statistical hypothesis tests in pain research. [12]. Table 4 presents the results of a T-test comparing the predicted slump flow values generated by Gene Expression Programming (GEP) to actual slump flow values. The key findings from the t-test are that the mean slump flow for actual data (Slump) is 679.6279, while the mean predicted slump flow using GEP is 680.2271. The difference between the means is very small, suggesting that the GEP model produces predictions close to the actual values. Moreover, the variance in the actual slump flow (4101.808) is higher than the variance in the GEP predictions (2397.083), indicating that the predicted values are less dispersed compared to the actual measurements. Additionally, the calculated t-statistic is -0.22933, which is relatively small in magnitude, indicating a minimal difference between the two groups. The p-value for the two-tailed test is 0.818639, which is much higher than a typical significance threshold (e.g., 0.05). This suggests that there is no statistically significant difference between the actual and predicted values. The critical t-values (1.645655 for one-tail and 1.961212 for two-tail) further confirm that the observed t-statistic falls well within the acceptance region, meaning we fail to reject the null hypothesis.

Table 4. T-Test of the Data

	<i>Slump</i>	<i>GEP</i>
Mean	679.6279	680.2271
Variance	4101.808	2397.083
Observations	952	952
Pooled Variance	3249.445	
df	1902	
t Stat	-0.22933	
P(T<=t) one-tail	0.409319	
t Critical one-tail	1.645655	
P(T<=t) two-tail	0.818639	
t Critical two-tail	1.961212	

4.3.3 Z-Test

Utilized Z scores to evaluate several distinct approaches, including as fold changes, Z ratios, and Z and t statistical tests, for forecasting noteworthy alterations in gene expression [13]. Table 5

presents the results of a z-test comparing the actual slump flow values (Slump) and the predicted values from Gene Expression Programming (GEP). The test aims to determine whether there is a statistically significant difference between the two distributions. The mean actual slump flow is 679.6279, and the mean predicted slump flow from GEP is 680.2271. The difference between these means is minimal, suggesting that the GEP model closely predicts the actual slump flow; endorsement of the t-test. The Z-test assumes a null hypothesis (H_0) that the difference between the actual and predicted values is zero. The z-score is -0.22934, which is small in magnitude, indicating that the observed difference is not substantial. The p-value for the two-tailed test is 0.818602, which is much greater than the standard significance level 0.05 or 0.01. This means we fail to reject the null hypothesis, implying that there is no statistically significant difference between actual and predicted slump flow values. The critical z-values further confirm that the observed z-score falls within the acceptance region.

The Z-test results indicate that there is no statistically significant difference between the actual and GEP-predicted slump flow values. This suggests that Gene Expression Programming (GEP) is an effective predictive tool for estimating slump flow in self-consolidating concrete. However, while the predictions are statistically close to the actual values, further evaluation using different datasets and performance metrics (e.g., RMSE, R^2) could enhance the robustness of the model.

Table 5. Z-Test of the Data

	<i>Slump</i>	<i>GEP</i>
Mean	679.6279	680.2271
Observations	952	952
Hypothesized Mean Difference	0	
z	-0.22934	
P(Z<=z) one-tail	0.409301	
z Critical one-tail	1.644854	
P(Z<=z) two-tail	0.818602	
z Critical two-tail	1.959964	

4.4 Shape Analysis

Explanation/SHAP analysis is about using methods from mathematical explanations to explain machine-learning models, especially complex ones, such as deep neural networks, ensemble methods, or black-box models. Therefore, the importance of the explanation in SHAP can be described in the following. The significant nonlinear correlations between independent and dependent variables are shown by the SHAP value graphs [14].

4.5 Feature Importance Plots

SHAP values help understand which features have the greatest impact on the model's predictions. It helps in feature identification and ranking based on their contribution so that you can focus on the most influential factors.

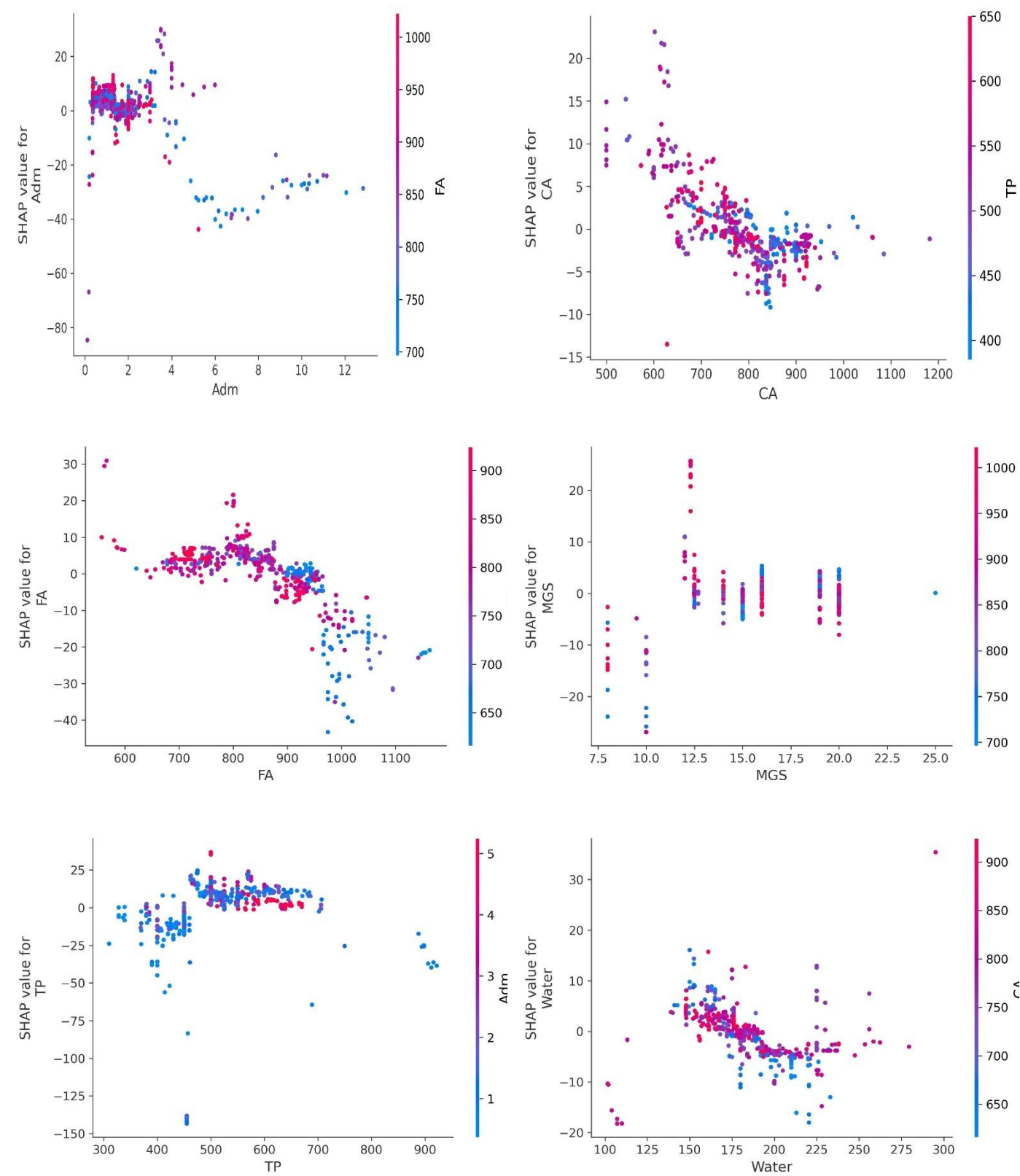


Figure 8. Parameters Dependence 270

4.6 Force Plot 271

The purpose of a force plot in SHAP (SHapley Additive exPlanations) analysis is to visually 272
explain the contributions of individual features to a model's prediction for a specific instance (data 273

point). It provides a way to understand how each feature, along with its corresponding value, influ- 274
ences the final prediction. 275

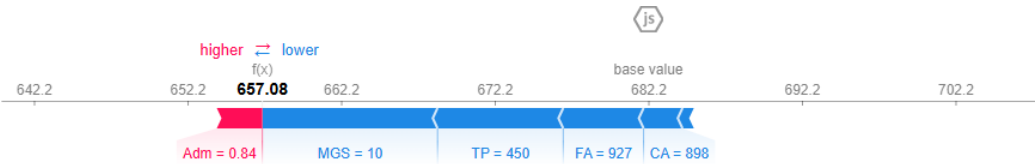


Figure 9. Force Plot

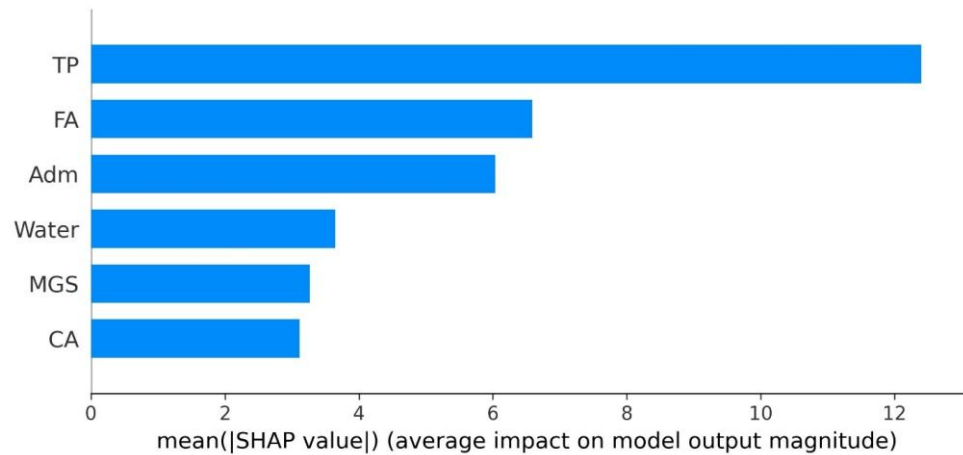


Figure 10. Parameters Importance

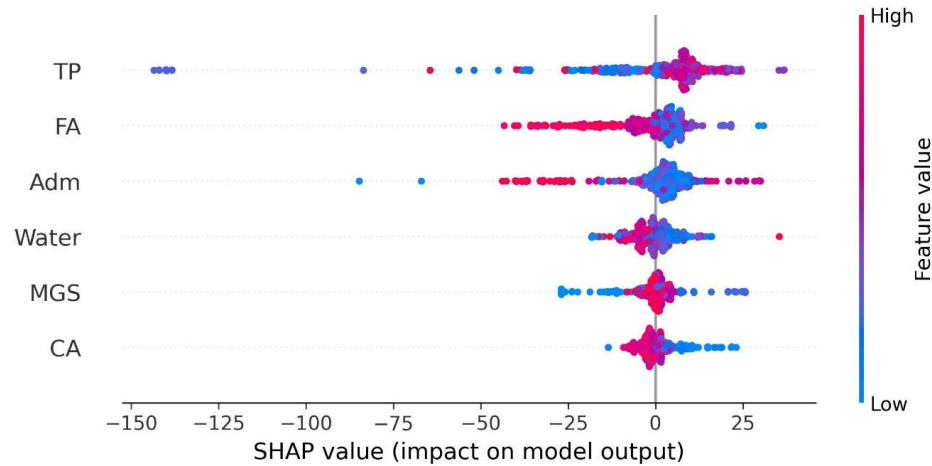
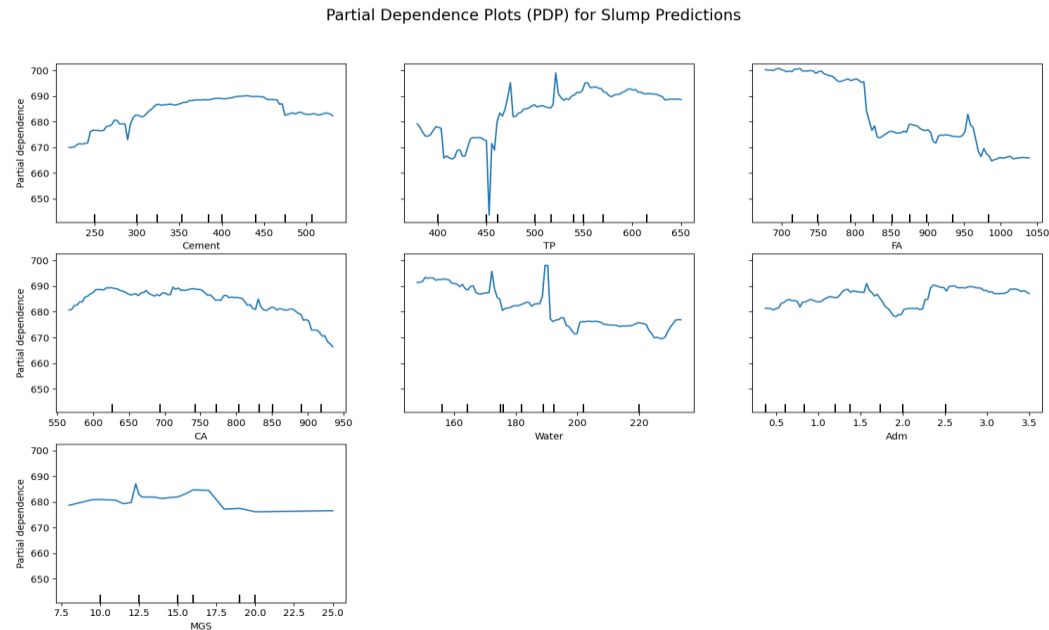


Figure 11. Parameters influence

4.7 Partial Dependence Plot

Partial Dependence Plot (PDP) is a versatile tool for visualizing what it means for one feature 283
to be related to predictions in a machine learning model and keeping all other variables fixed. This 284
is handy especially when trying to analyze specific impacts of features on prediction output from 285
complex models, like Random Forests and Gradient Boosting Machines or a black-box model. 286



5. Conclusions

The study employed GEP on the SCC with regression and error analysis as basis. From results, a far more superior fit was reported when GEP predicted slump flow due to nonlinearity capturing superiority compared with regression models, using various tools to check up statistics, which had RMSE, MAE, and R² for determining a model performance of a prediction method.

- Model Performance: The GEP model achieved high predictive accuracy with metrics like $R^2=0.87$, $RMSE = 25.66$, and $MAE = 17.63$, highlighting its ability to model SCC slump flow effectively.
- Correlation Insights: High correlations were discovered from a comprehensive correlation matrix as far as relationships of input parameters water-cement ratio, superplasticizer, and aggregates percentages are concerned with moderate to high influence on the slump flow.
- Descriptive Analysis: Statistical investigation revealed normal distribution in the slump flow data, and the mean and standard deviations obtained were 679.63 mm and 64.05, respectively, which ensured homogeneity of the data set.
- The plot for actual vs. predicted slump flow values in case of GEP showed a closer alignment compared to the other plot, visualizing the plots as well as the residual randomness that outperform linear regression models.
- Model Validation: T-tests and Z-tests affirmed the slight differences in mean actual and predicted values, hence confirmed the soundness of GEP.
- Feature Importance: SHAP ranked features such as water-cement ratio and admixtures as key contributors, enabling interpretability, and thus aid in mix optimization strategies.

Abbreviations

The following abbreviations are used in this manuscript:

TP	Total Powder
FA	Fine aggregate
CA	Coarse Aggregate
Adm	Admixture
MGS	Maximum Grain Size
SSC	Self-Consolidated Concrete
GEP	Gene Expression Programming

References

1. Amine el Mahdi Safhi ^a, Hamed Dabiri ^b, Ahmed Soliman ^a, Kamal H. Khayat ^c "Prediction of self-consolidating concrete properties using XGBoost machine learning algorithm: Part 1–Workability" doi: <https://doi.org/10.1016/j.conbuildmat.2023.133560>
2. Alireza Mohebbi, Mohammad Shekarchi¹, Mehrdad Mahoutian and Shima Mohebbi, "Modeling the effects of additives on rheological properties of fresh self-consolidating cement paste using artificial neural network"
3. Muhammad Raheel, Mudassir Iqbal, Rawid Khan, Muhammad Alam, Marc Azab & Sayed M. Eldin "Application of gene expression programming to predict the compressive strength of quaternary-blended concrete" <https://link.springer.com/article/10.1007/s42107-023-00573-w>
4. Yaser Gamil* "Machine learning in concrete technology: A review of current research, trends, and applications" <https://doi.org/10.3389/fbuil.2023.1145591>
5. Soo-Duck Hwang, Kamal H. Khayat, and Olivier Bonneau "Performance-Based Specifications of Self-Consolidating Concrete Used in Structural Applications"
6. Mosbeh R. Kaloop ^{1,2,3}ORCID, Pijush Samui ⁴, Mohamed Shafeek ⁵ and Jong Wan Hu ^{1,2}, "Estimating Slump Flow and Compressive Strength of Self-Compacting Concrete Using Emotional Neural Networks" <https://doi.org/10.3390/app10238543>
7. A.M. el Safhi, A Comprehensive Self-Consolidating Concrete Dataset for Advanced Construction Practices, 2024, <https://doi.org/10.5281/zenodo.10569517>.
8. Ferreira, C. (2001). "Gene expression programming: A new adaptive algorithm for solving problems."
9. Sarah E. Kemp, May Ng, Tracey Hollowood, Joanne Hort "Introduction to Descriptive Analysis" <https://doi.org/10.1002/9781118991657.ch1>
10. "Correlation analysis" Julian A. <https://hdl.handle.net/11563/153064>
11. Lin Shao ^a, Timo Schleicher ^b, Michael Behrisch ^b, Tobias Schreck ^a, Ivan Sipiran ^c, Daniel A. Keim ^b "Guiding the exploration of scatter plot data using motif-based interest measures", <https://doi.org/10.1016/j.jvlc.2016.07.003>
12. Tae Kyun Kim "T test as a parametric statistic" Tae Kyun Kim DOI: <https://doi.org/10.4097/kjae.2015.68.6.540>
13. P. Vawter ^{† ‡}, William J. Freed [†], Kevin G. Becker, "Analysis of Microarray Data Using Z Score Transformation, Chris Cheadle ^{*}, [https://doi.org/10.1016/S1525-1578\(10\)60455-2](https://doi.org/10.1016/S1525-1578(10)60455-2)
14. Chao Yang, Mingyang Chen, Quan Yuan, "The application of XGBoost and SHAP to examining the factors in freight truck-related crashes: An exploratory analysis", <https://doi.org/10.1016/j.aap.2021.106153>