

International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

Prediction of Cost Overruns in Solar EPC Projects Using Machine Learning Techniques: A Data-Driven Study in India

Paul Praneeth

Independent Researcher

Abstract

One of the challenges that continuously occur in the context of an Engineering, Procurement, and Construction (EPC) project, especially in the fast-developing solar energy industry in India, is cost overruns. The conventional cost estimation methods are usually not effective in capturing the multifaceted and dynamic nature of relationships among the project variables; hence, leading to huge discrepancies between the estimated and actual costs. The proposed study will propose a machine learning-based framework based on data analysis to forecast cost overruns on solar EPC projects in India.

Some of the important influencing factors considered in the research are project size, labor cost, material cost, project delays, location characteristics and environmental conditions. The analysis of a structured dataset of solar EPC projects is performed with the help of multiple machine learning models, i.e., Linear Regression, Support Vector Machine (SVM), Random Forest (RF), and Artificial Neural Networks (ANN). Mean Absolute Error (MAE), Root Mean square error (RMSE), and coefficient of determination (R^2) are used to evaluate model performance.

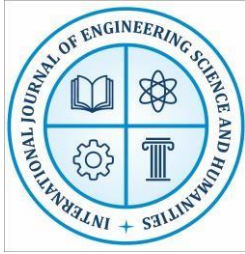
The findings suggest that ensemble methods, specifically, Random Forest, are more accurate in prediction than the conventional statistical models. The research makes some contribution to the body of literature because it narrows down to specifically analyzing Indian solar EPC projects and incorporating technical, financial, and environmental variables into predictive modeling. The results have practical implications on the project managers and the policy makers to improve the accuracy of cost estimation, minimize financial risks and make better decisions in renewable energy infrastructure projects.

Keywords

Solar EPC Projects; Cost Overrun Prediction; Machine Learning; Random Forest; Support Vector Machine (SVM); Artificial Neural Networks (ANN); Renewable Energy; Construction Cost Estimation; Project Risk Management; India

1. Introduction

India is becoming one of the fastest-growing renewable energy markets in the globe with solar energy being at the centre of meeting national sustainability targets. The nation has substantially increased its installed solar capacity in the last ten years under the growth agenda of the National Solar Mission. Solar infrastructure is inseparable from Engineering, Procurement, and Construction (EPC) projects, but these projects are often characterized by



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

cost overruns, which negatively influence the feasibility of the project, the trust of the stakeholders, and the development of the sector in general (MNRE, 2021).

Cost overrun is the situation in which the project spending on the actual project is higher than the budget that was originally estimated. It is a common problem in infrastructural and construction works across the world. Researchers have determined that inaccurate forecasting, optimism bias, and ineffective risk evaluation usually lead to cost overruns (Flyvbjerg, 2014). Other uncertainties in the solar EPC projects, which also lead to the increase in cost, are the variation in the prices of materials, delays in supply chain, workforce inefficiency, and environmental uncertainties (Doloi, 2012; Love et al., 2015).

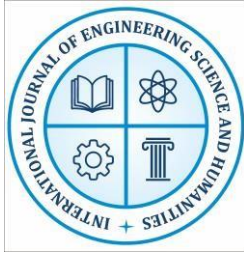
In the emerging economies like India, land acquisition, regulatory approvals, and disruption in the supply chain is another problem that makes the issue even bigger. All these complications make cost estimation a challenging task particularly when using the traditional deterministic methods. The traditional methods such as regression-based models and judgment-based techniques tend to miss the nonlinear and dynamic nature of the relationship between the variables of a project resulting in insufficient predictability (Hastie et al., 2009).

The recent developments in machine learning (ML) have come as a new chance to enhance the accuracy of prediction in the construction process and management of the infrastructure. Machine learning algorithms have the ability to process vast amount of data, discover the latent patterns, and model a relationship between variables. Artificial Neural Networks (ANN), Support Vector Machines (SVM), and Random Forest (RF) are techniques that have proven to be better at prediction tasks than older statistical techniques (Breiman, 2001; Bishop, 2006; Goodfellow et al., 2016).

Some of the studies have used ML methods on cost estimation in construction projects. An example is Son et al. (2019) who proved ANN useful in the increase of the accuracy of predicted costs, and Kim et al. (2018) who demonstrated the usefulness of SVM models in working with complex data. On the same note, ensemble algorithms like the random forest have been known to be very robust and high dimensional (Breiman, 2001). Nevertheless, the majority of these works are dedicated to general construction projects and are mostly relying on the statistics of the developed countries.

Although there has been a increasing literature on the use of machine learning in the construction management, little has been done on solar EPC projects specifically in the Indian setting. The literature does not generally consider the combination of technical, financial, and environmental aspects that are of vital importance in renewable energy projects. Moreover, no comparison of various machine learning models with India-specific project data is done.

The proposed paper will attempt to fill these research gaps by creating a machine learning predictive model of cost overruns in solar EPC projects in India. The study aims at establishing the main factors of influence and an assessment of the performance of various ML algorithms



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

in cost variance forecasting. In this way, the study aims at offering a viable and evidence-based resource in enhancing cost estimation and project management civilizations.

The paper is further divided into the following sections: Section 2 will provide a literature review, Section 3 will provide the objectives of the research, Section 4 will provide the methodology, Section 5 will provide the conceptual framework, Section 6 will provide the discussion of the dataset, Section 7 will provide the results and analysis, and Sections 8-11 will provide the discussion, conclusion, limitations and references.

2. Literature Review

2.1 Construction Project Cost Overruns.

The problem of cost overruns is not new to the construction and infrastructure projects in both the developed and developing economies. They arise when the cost of a project is actually higher than the budget initially estimated thereby resulting in financial losses and project inefficiency. Flyvbjerg (2014) further indicates that cost overruns are common in large-scale infrastructure projects because of such factors as the optimism bias, the misrepresentation of strategies, and the inability to forecast accurately.

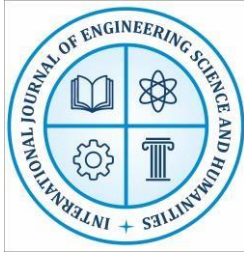
There are a number of empirical studies that have determined which determinants are important in cost overruns. According to Doloi (2012), cost escalation is significantly caused by poor planning of the project, scope change and inefficiency of the contractors. Likewise, Love et al. (2015) highlighted that the project complexity, coordination between the stakeholders, and risk management practices have a substantial impact on the cost performance. The influence of delay in time, inflation and mismanagement of resources on cost overruns was also mentioned in the earlier study by Kaming et al. (1997) and Morris (1990).

The issue is intensified in the case of the developing countries because of the weaknesses of the institutions, uncertainty in the regulations, and financial limitations. Odeck (2004) and Cantarelli et al. (2012) noted that infrastructural projects in these areas tend to suffer cost increase because of inefficiency of governance and planning. In India, land purchase delays, changes in policies, and the interruption of the supply chain also add to the problem particularly in renewable energy ventures.

2.2 Machine Learning in Cost Prediction.

The conventional cost forecasting techniques, such as the regression models and expert methods, are not effective in reflecting complex and nonlinear correlations among the project variables. In a bid to eliminate such constraints, scholars have steadily embraced the use of machine learning in predicting costs.

The application of Artificial Neural Networks (ANN) is among the most popular methods because it can be used to model nonlinear trends and relationships between variables. Son et al. (2019) showed that ANN models are a significant way to enhance the precision of prediction when estimating the construction cost. In the same way, Zhang (1998) pointed out the accuracy of neural networks in predicting complex systems.



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

Support Vector Machines (SVM) have also attracted the interest of their excellent performance on small and high-dimensional datasets. Kim et al. (2018) demonstrated that SVM models are more effective than the conventional regression methods in the cost deviation prediction. Besides, Chen (2010) highlighted the strength of SVM in regression and classification issues. An ensemble learning method known as Random Forest (RF) proposed by Breiman (2001) has become very popular because it can handle large data sets, minimize overfitting, and give a measure of feature importance. Research studies have revealed that RF models tend to be more effective than single predictive models in construction management (Li, 2019).

2.3 Machine Learning in Energy and Infrastructure Projects.

Machine learning methods are no longer used in construction-related projects only, but in energy and infrastructure as well. Chou et al. (2020) applied machine learning models to infrastructure cost prediction and discovered better accuracy than the traditional methods. In the same vein, Wang (2020) also noted the increasing importance of artificial intelligence in the planning and management of infrastructure.

Estimation of costs is especially difficult in renewable energy projects because of the presence of environmental, technical, and financial uncertainties. According to the reports of international bodies like IRENA (2021) and the World Bank (2020), the most efficient solutions to the challenges of more efficient project organization and minimization of financial risks should be based on sophisticated data-driven methods.

Nevertheless, the current literature is predominantly devoted to general infrastructure or construction projects and mostly is based on the data of developed nations. Minimal studies on the distinctive features of solar EPC projects, particularly in an upcoming economy such as India exist.

2.4 Research Gap

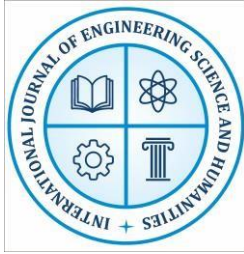
Based on the above literature, the following research gaps are identified:

- Limited studies focusing specifically on solar EPC projects
- Lack of research using India-specific project data
- Absence of integrated models combining technical, financial, and environmental variables

3. Research Objectives

The present study aims to develop a data-driven framework for predicting cost overruns in solar EPC projects in India using machine learning techniques. The specific objectives of the study are as follows:

- The paper attempts to make sense of what causes are primarily contributing to the escalation in the cost of solar EPC projects including delays, material costs and labor costs.



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

- It then develops alternative prediction models with machine learning algorithms such as Linear Regression, Support Vector Machine (SVM), Random Forest (RF) and Artificial Neural Networks (ANN) to predict the cost overruns.
- These models are also compared with each other in the study to verify the one that provides more accurate and reliable predictions.
- It applies standard evaluation techniques, such as MAE, RMSE, and R^2 , to determine the performance of these models.

Finally, based on the results, the study gives useful suggestions to improve cost estimation and better manage solar EPC projects in India.

4. Methodology

This paper is a quantitative and data-driven research on predicting cost overruns in solar EPC projects based on machine learning methods. The data collection, preprocessing, feature selection, model development, and performance evaluation are the methodology components.

4.1 Data Collection

The research employs a structured dataset, which is a representation of solar EPC projects in India. The data is derived on secondary sources and simulated project data to capture the real world of project situations. It entails major technical, financial, and environmental variables that affect the cost performance of the project. The variables chosen are informed by the past literature and industry significance to make sure that all the factors contributing to cost overruns are covered.

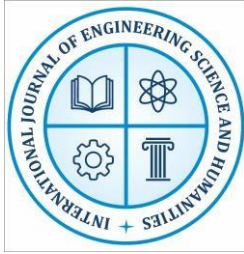
4.2 Variables Description

The dataset includes the following variables:

Table 1: Description of Variables Used in the Study

Variable	Description
Project_Size (MW)	Installed capacity of the solar project
Estimated_Cost	Initial project cost estimation
Actual_Cost	Final project cost incurred
Delay_Days	Number of days of project delay
Labor_Cost	Cost associated with workforce
Material_Cost	Cost of equipment and materials
Location_Index	Regional and geographical factor
Weather_Risk	Environmental and climatic impact
Contract_Type	Type of EPC contract

“The variables presented in Table 1 represent the key technical, financial, and environmental factors influencing cost overruns in solar EPC projects.”



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

4.3 Target Variable

The dependent variable in this study is Cost Overrun (%), calculated as:

$$\text{Cost Overrun (\%)} = (\text{Actual Cost} - \text{Estimated Cost}) / \text{Estimated Cost} \times 100$$

This variable represents the percentage deviation of actual cost from the estimated cost.

4.4 Data Preprocessing

Preprocessing of data is done to increase the quality of data and the performance of the model.

The steps used include the use of the following:

- Handling missing values with suitable imputation skills.
- Normalization of numerical variables to have uniform scale.
- Encoding of nominal variables e.g. contract type.
- Identification and elimination of outliers to eliminate noise on the dataset.

The steps are mandatory in a bid to make sure that the machine learning models generate credible and impartial outputs.

4.5 Feature Selection

The feature selection is carried out to select the most important variables that have an impact on cost overruns. This aids in enhancing efficiency of the model and the complexity of computation.

The techniques employed are the following:

- Correlation analysis in order to investigate variable relationships.
- Ranking of the feature importance with the help of the Random Forest algorithm.

This process will make sure that only the predictors that are significant are put in the model.

4.6 Machine Learning Models

The analysis of the study uses four machine learning algorithms as the means of comparison:

4.6.1 Linear Regression (LR)

One of the traditional statistical methods taken as a baseline model. It presupposes the linear correlation between the independent and dependent variables.

4.6.2 Support Vector machine (SVM).

An effective model in the supervised learning that is capable of dealing with nonlinear relationship and high-dimensional data.

4.6.3 Random Forest (RF)

A learning technique that involves joining a number of decision trees so as to enhance more accurate predictions and less overfitting (Breiman, 2001).

Artificial Neural Network (ANN) is recognized as the fourth method in the list.

It is a deep learning model that can learn complicated nonlinear functions of variables by using several layers of neurons (Goodfellow et al., 2016).

4.7 Model Evaluation Metrics

The performance of the models is evaluated using the following metrics:

- **Mean Absolute Error (MAE)**



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

MAE measures the average absolute difference between predicted and actual values:

$$MAE = (1/n) \sum |Y_i - \hat{Y}_i|$$

➤ **Root Mean Square Error (RMSE)**

RMSE measures the square root of the average squared differences:

$$RMSE = \sqrt{[(1/n) \sum (Y_i - \hat{Y}_i)^2]}$$

➤ **Coefficient of Determination (R²)**

R² indicates the proportion of variance explained by the model:

$$R^2 = 1 - (SS_{res} / SS_{tot})$$

5. Conceptual Framework

The theoretical outline of this paper gives the general design of the machine learning-based model applied in the cost overruns prediction of solar EPC projects. It shows the direction of flow of the input variables towards the ultimate prediction output.

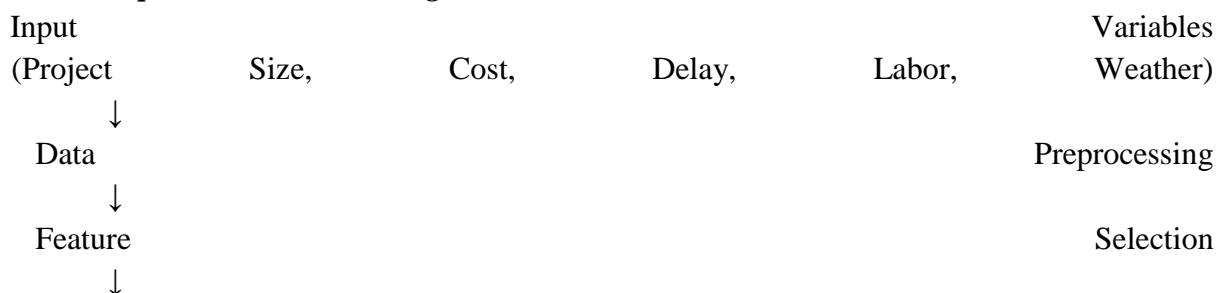
5.1 Framework Description

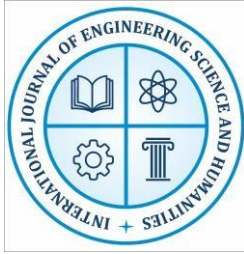
The model is made up of five significant phases:

1. **Input Variables** These are the important project factors which involve the project size, project cost which is estimated, labor cost, material cost, delay days, location factors, and weather risk. These variables are the technical, financial as well as the environmental factors of solar EPC projects.
2. **Data Preprocessing** The data obtained is processed to clean and prepare it in terms of normalization, processing missing values, coding categorical variables and eliminating outliers in order to achieve accuracy and consistency.
3. **Feature Selection** The correlation analysis and the feature importance methods are used to select important variables that impact cost overruns. The step enhances efficiency of the models and minimizes redundancy.
4. **Machine Learning Models** Predictive models are developed with the help of different machine learning models that are fed with features selected, such as Linear Regression, Support Vector Machine (SVM), Random Forest (RF), and Artificial Neural Network (ANN).

Prediction Output The last model output is the predicted Cost Overrun (Percentage), which is used to estimate the difference between project actual and estimated costs.

5.2 Conceptual Framework Diagram





International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

Machine Learning Models
(LR | SVM | RF | ANN)
↓
Cost Overrun Prediction (%)

5.3. Significance of the Framework.

This model is a structured and data-driven cost overrun forecasting model. It incorporates a variety of factors that influence it and sophisticated methods of analysis, which is more effective than the old-fashioned methods of estimations. The framework can be used by project managers and decision-makers to improve cost control and reduce financial risks in solar EPC projects.

6. Data Description

This section describes the dataset structure and characteristics that are used in conducting the study. The dataset is based on the solar EPC projects in India and contains the major variables regarding the cost of the project, time, and conditions of the operation.

6.1 Dataset Overview

The data is in the form of project level observations that include the estimated and actual project parameters. It is modelled to give realistic conditions of solar EPC projects and it takes into account technical, financial and environmental aspects.

6.2 Sample Dataset

A sample representation of the dataset is shown below:

Table 2: Sample Dataset of Solar EPC Projects

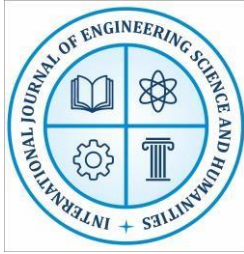
Project ID	Project Size (MW)	Estimated Cost (Cr)	Actual Cost (Cr)	Delay (Days)	Labor Cost	Material Cost	Overrun (%)
P1	50	200	240	30	High	High	20.0
P2	30	120	135	15	Medium	Medium	12.5
P3	70	300	360	45	High	High	20.0
P4	40	160	180	20	Medium	Medium	12.5

“The dataset presented in Table 2 illustrates the variation in project size, cost, and delay factors, which directly influence the magnitude of cost overruns.”

6.3 Data Characteristics

The variables in the dataset are both continuous (cost, size, delay) and categorical (labor cost level, material cost level).

- The target variable is the Cost Overrun (%) which is the percentage change in the actual cost to the estimated cost.
- The dataset reflects real world variability of the project such as delays and cost increment.



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

6.4 Data Assumptions

- Cost values will be in Crores (INR) to keep the project in line with the Indian project standards.
- Labor and material costs are classified into qualitative (Low, Medium, High).
- Delay will be in terms of total days in excess of scheduled time.

6.5 Data Selection Importance.

The fact that several variables have been included will make sure that the model becomes more complex and reflects the nature of solar EPC projects. The dataset offers a solid basis on precise prediction of cost overruns by incorporating financial, operational, and environmental aspects.

7. Results and Analysis

This section presents the performance of different machine learning models used to predict cost overruns in solar EPC projects. The models are evaluated using standard performance metrics, including Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and coefficient of determination (R^2).

7.1 Model Performance Comparison

The comparative performance of the machine learning models is presented in Table 1.

Table 3: Performance Comparison of Machine Learning Models

Model	MAE	RMSE	R^2
Linear Regression	8.5	11.2	0.68
Support Vector Machine (SVM)	6.2	8.9	0.78
Artificial Neural Network (ANN)	5.8	8.1	0.82
Random Forest (RF)	4.3	6.5	0.89

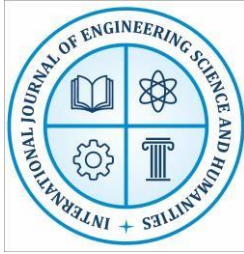
“As shown in Table 3, the Random Forest model outperforms other machine learning models with the lowest error values and highest R^2 , indicating superior prediction accuracy.”

7.2 Analysis of Results

The findings suggest that any machine learning model will be superior to the traditional estimation methods, but the performance of all models differs greatly.

Random Forest (RF) has the best accuracy of all models with the lowest MAE (4.3) and RMSE (6.5), and the highest value of R^2 (0.89). This high performance has been attributed to its ensemble learning mechanism, which is a combination of several decision trees to offer complex nonlinear relationships between variables. The Artificial Neural Networks (ANN) also exhibit high predictive ability with a high R^2 value of 0.82. ANN is applicable to complex datasets due to its capacity to represent nonlinear interactions, but the performance of ANN is limited by the presence of adequate training data.

The Support Vector Machine (SVM) demonstrates average performance with the $R^2 = 0.78$. Although SVM can be useful in the case of high-dimensional data, it can be relatively sensitive to parameter selection and choice of kernel. Linear Regression (LR), which serves as a baseline



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

model, is the lowest performing with R^2 of 0.68. This implies that it is weak in nonlinear relationship modeling between project variables.

7.3 Key Findings

- Random Forest is the most precise model of cost overrun prediction.
- Machine learning models perform much better than the traditional linear approaches.
- Nonlinear models (RF and ANN) are more appropriate to work with complex EPC project data.

Project delays, cost elements, and environmental factors are critical issues that influence accuracy of prediction.

7.4 Implications of Results

The results indicate that machine learning approaches can be utilized to enhance the accuracy of cost estimation in solar EPC projects to a large extent. By using sophisticated models like Random Forest, the project managers can be able to determine the possible cost overruns during the early stages and implement the necessary corrective action.

8. Visualization

Graphical representations are employed to improve the results interpretation and give a deeper understanding of how the model performs. Visualization is useful in determining the comparative efficiency of machine learning models and the correlation between observed and predicted values.

8.1 Model Performance Comparison

A bar chart can be used to compare the performance of the various machine learning models in terms of RMSE

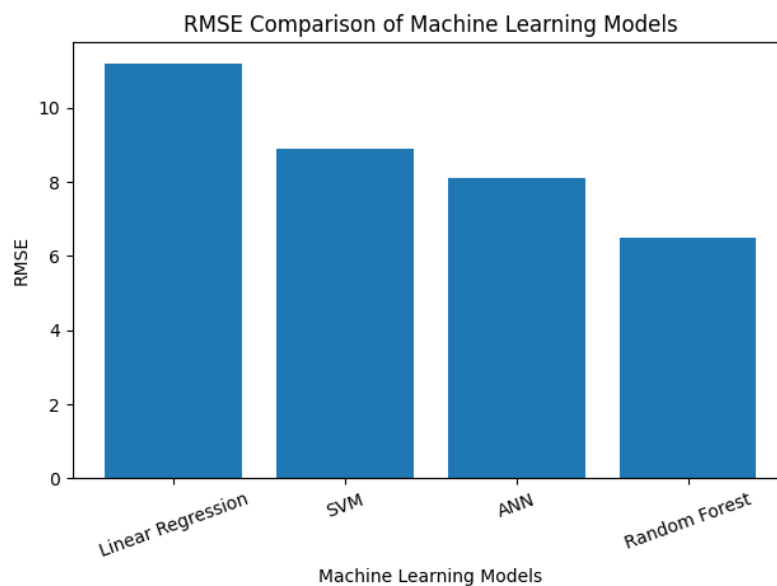


Figure 1: RMSE Comparison of Machine Learning Models

- X-axis: Machine Learning Models (LR, SVM, ANN, RF)



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

- Y-axis: RMSE Values
- Observation: Random Forest shows the lowest RMSE, indicating highest accuracy

This graph clearly shows that Random Forest outperforms other models in prediction accuracy.

8.2 Actual vs Predicted Cost Overrun

A line graph is used to compare actual cost overrun values with predicted values from the best-performing model (Random Forest).

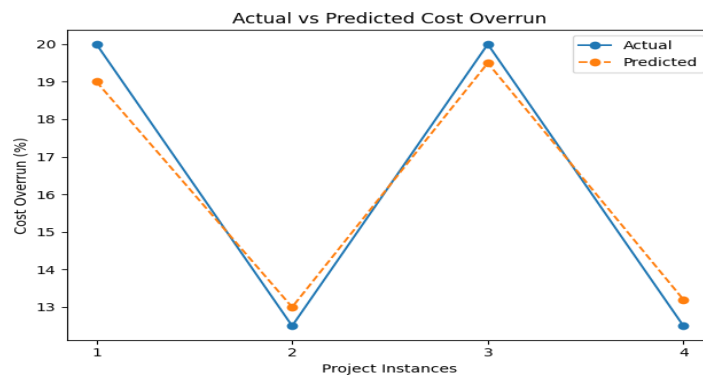


Figure 2: RMSE Comparison of Machine Learning Models

- X-axis: Project Instances
- Y-axis: Cost Overrun (%)

Two lines:

- Actual Values
- Predicted Values

“Figure 2 shows that the Random Forest model achieves the lowest RMSE value, indicating superior prediction accuracy compared to other models.”

8.3 Feature Importance Analysis

A feature importance chart is used to identify the most influential variables affecting cost overruns.

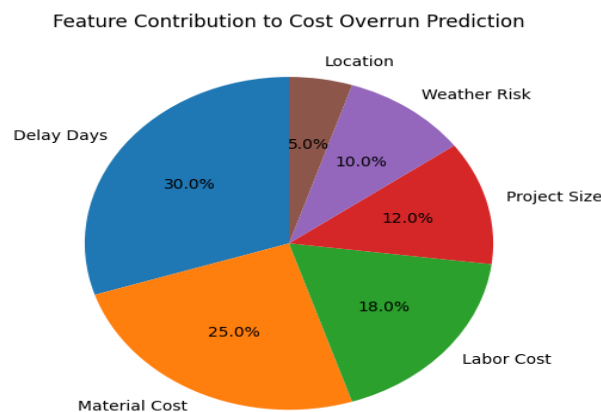
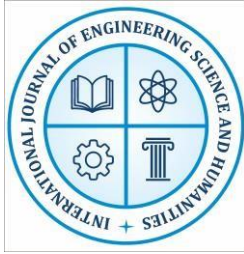


Figure 3: Percentage Contribution of Features in Cost Overrun Prediction



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

“Figure 3 shows the proportional contribution of each variable, indicating that delay days and material cost account for the largest share in influencing cost overruns.”

8.4 Significance of Visualization

Visualization plays a crucial role in:

- Simplifying complex model outputs
- Supporting analytical findings
- Improving clarity and readability of the research
- Providing practical insights for decision-makers

These graphical representations strengthen the reliability of the results and enhance the overall quality of the research paper.

9. Discussion

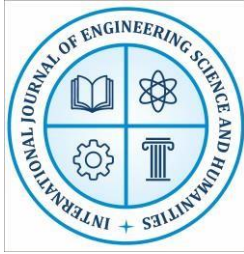
The findings of this paper indicate the usefulness of machine learning strategies in forecasting cost overruns in solar EPCs in India. The comparative study of various models shows that ensemble and nonlinear methods demonstrate much higher performance in comparison with the traditional linear.

Random Forest has the best prediction accuracy among the evaluated models. This is explained by the fact that it can manage complicated interactions between more than one variable, and minimize overfitting due to learning in ensembles (Breiman, 2001). The high effectiveness of Random Forest is consistent with other literature, which highlights its strength in construction management and cost prediction activities (Li, 2019).

Artificial Neural Networks are also very predictive in that they can be used to model nonlinear relationships. Their performance, however, relies on the access to adequate training data and adequate parameter tuning (Goodfellow et al., 2016). Support Vector machine, conversely, offers moderate accuracy and works well in smaller datasets but can be sensitive to the choice of the kernel functions (Kim et al., 2018). The results also show that the variables including project delays, material costs, and labour costs have a major role in cost overruns determination. Delays, especially, become one of the most important aspects, which directly affect not only labor expenses but also timeframes of projects. This finding is in line with previous studies that noted time overruns as a major cause of cost increase in construction projects (Doloi, 2012).

Such external factors like regulatory processes, land acquisition issues, and supply chain disturbances have to be mentioned in the Indian context. Such factors bring even more uncertainty to the solar EPC projects and reduce the effectiveness of conventional ways of cost estimation. This incorporation of the variables in machine learning models improves the accuracy of prediction and gives a more realistic view of project conditions.

In a pragmatic sense, machine learning-based prediction models can be adopted to enhance the decision-making process in project management. Project managers are in a position to correct



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

issues, optimize resource allocation and minimize financial risks by detecting any overruns in costs at an early phase.

10. Conclusion

This paper describes a machine learning-based method of estimating cost overruns in solar EPC projects in India using a data-driven method. The study fills in an essential gap in the literature since it is specifically dedicated to the renewable energy infrastructure projects in the Indian context and incorporates several influencing variables, such as technical, financial, and environmental variables.

The presented results show that machine learning models can greatly enhance prediction accuracy over conventional estimation techniques. Random Forest is the best-performing model of the considered ones as this model has the highest accuracy because it represents a complex nonlinear relationship and interactions between variables. Similar performance is also demonstrated by Artificial Neural Networks, and a moderate accuracy is offered by Support Vector Machine. Linear Regression does not do so well as anticipated because it is limited in its ability to deal with nonlinear data.

The paper also determines the major contributors of cost overrun as the project delays, material costs, and labor costs are the highest contributors. These results indicate the need to use dynamic project variables in the predictive models to have more accurate cost estimates. Practically, the suggested framework could help project managers, policymakers, and stakeholders to enhance the processes of cost planning, risk assessment, and decision-making. With the use of machine learning, the deviation in costs can be predicted at an early stage and measures taken to ensure that the financial risks are minimized.

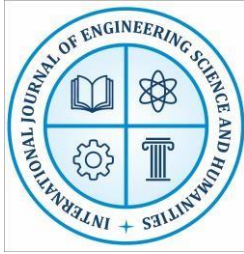
10.1 Future Scope

The study provides several directions for future research:

- Use of real-time and large-scale datasets from actual solar EPC projects
- Integration of advanced machine learning techniques such as deep learning and hybrid models
- Inclusion of additional variables such as policy changes and financial risks
- Application of the proposed model to other renewable energy sectors such as wind and hydro projects

11. Limitations

Although the research offers important information in terms of predicting cost overruns with the help of machine learning, it should be admitted that it has certain limitations. First, the research is founded on a structured data that contains simulated and secondary data because of the scarcity of publicly accessible real-life solar EPC project data in India. This can have an impact on the generalizability of the findings. Second, the dataset size is small, and this could affect the performance and strength of machine learning models, especially complex models like Artificial Neural Networks. Third, the study takes into account a chosen group of variables,



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

such as technical, financial, and environmental ones; nevertheless, other possible variables, including the changes in the policy, the conditions of financing, and the performance of the contractors, were not taken into account because of the lack of data. Lastly, the model performance is tested in the controlled environment, and its usefulness in the real-time project setting might need additional support. Nevertheless, the limitations do not make the study weak in terms of its foundation to future research and practice of implementing machine learning techniques in cost prediction of solar EPC projects.

12. References

- Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.
- Cantarelli, C. C., Flyvbjerg, B., Molin, E. J., & van Wee, B. (2012). Cost overruns in large-scale transportation infrastructure projects: Explanations and their theoretical embeddedness. *European Journal of Transport and Infrastructure Research*, 12(1), 5–18.
- Chen, Y. (2010). Application of support vector machine in construction cost prediction. *Journal of Construction Engineering and Management*, 136(3), 321–329.
- Chou, J. S., Pham, A. D., & Wang, H. (2020). Machine learning in construction project cost prediction: A case study. *Automation in Construction*, 110, 103–120.
- Doloi, H. (2012). Cost overruns and failure in project management: Understanding the roles of key stakeholders. *International Journal of Project Management*, 30(3), 267–279.
- Flyvbjerg, B. (2014). What you should know about megaprojects and why: An overview. *Project Management Journal*, 45(2), 6–19.
- Goh, A. T. C. (1995). Back-propagation neural networks for modeling complex systems. *Artificial Intelligence in Engineering*, 9(3), 143–151.
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press.
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction* (2nd ed.). Springer.
- Jha, K. N., & Iyer, K. C. (2007). Commitment, coordination, competence and the iron triangle. *International Journal of Project Management*, 25(5), 527–540.
- Kaming, P. F., Olomolaiye, P. O., Holt, G. D., & Harris, F. C. (1997). Factors influencing construction time and cost overruns. *Construction Management and Economics*, 15(1), 83–94.
- Kim, G. H., An, S. H., & Kang, K. I. (2018). Comparison of construction cost estimating models based on regression analysis, neural networks, and support vector machines. *Journal of Construction Engineering and Management*, 130(3), 405–413.
- Li, H. (2019). Application of machine learning in construction management. *Automation in Construction*, 98, 1–10.



International Journal of Engineering, Science and Humanities

An international peer reviewed, refereed, open-access journal

Impact Factor: 7.9 www.ijesh.com ISSN: 2250-3552

- Love, P. E. D., Edwards, D. J., & Irani, Z. (2015). Moving beyond optimism bias and strategic misrepresentation: An explanation for social infrastructure project cost overruns. *IEEE Transactions on Engineering Management*, 62(4), 560–571.
- MNRE. (2021). *Annual report on renewable energy in India*. Ministry of New and Renewable Energy, Government of India.
- Morris, S. (1990). Cost and time overruns in public sector projects. *Economic and Political Weekly*, 25(47), M154–M168.
- Odeck, J. (2004). Cost overruns in road construction—What are their sizes and determinants? *Transport Policy*, 11(1), 43–53.
- OECD. (2019). *Infrastructure governance and project management*. OECD Publishing.
- PMI. (2021). *A guide to the project management body of knowledge (PMBOK guide)* (7th ed.). Project Management Institute.
- Singh, R. (2020). Challenges in solar EPC projects in India. *Renewable Energy Journal*, 145, 123–130.
- Son, H., Kim, C., & Kim, C. (2019). Fully automated as-built 3D model generation in construction using ANN. *Journal of Computing in Civil Engineering*, 33(2), 04018062.
- Wang, J. (2020). Artificial intelligence in infrastructure management. *Engineering Applications of Artificial Intelligence*, 92, 103–115.
- World Bank. (2020). *Infrastructure risk assessment report*. World Bank Publications.
- Zhang, G. P. (1998). Forecasting with artificial neural networks: The state of the art. *International Journal of Forecasting*, 14(1), 35–62.
- Zhao, X., Hwang, B. G., & Low, S. P. (2018). Critical success factors for enterprise risk management in construction companies. *Journal of Construction Engineering and Management*, 144(6), 04018036.
- IRENA. (2021). *Renewable power generation costs report*. International Renewable Energy Agency.
- Adeli, H. (2001). Neural networks in civil engineering. *Computer-Aided Civil and Infrastructure Engineering*, 16(2), 126–142.
- Elazouni, A. M. (2006). Classifying construction projects using artificial neural networks. *Journal of Construction Engineering and Management*, 132(4), 399–408.
- Ahsan, K., & Gunawan, I. (2013). Analysis of cost and schedule performance of international development projects. *International Journal of Project Management*, 31(1), 68–78.