# UTILIZING PREDICTIVE ANALYTICS AND MACHINE LEARNING FOR ENHANCED PROJECT RISK MANAGEMENT AND RESOURCE OPTIMIZATION

## SAKILA AKTER JAHAN[*]
*Illinois State University*[*]

**\*Corresponding Author:** SAKILA AKTER JAHAN

## ABSTRACT

*Effective risk management is critical for minimizing unforeseen costs and ensuring project success through proactive strategies. This study introduces an innovative real-time risk management framework leveraging predictive analytics and machine learning (ML). By analyzing historical project data, this approach identifies potential risks, emphasizing parameters such as task durations, resource allocation, and project outcomes. A t-distributed Stochastic Neighbor Embedding (t-SNE) technique optimizes feature selection, reducing dimensionality while retaining essential data properties. Model evaluation metrics include accuracy, precision, recall, and F1-score. The results indicate that the Gradient Boosting Machine (GBM) outperforms previous models, achieving 85% accuracy, 82% precision, 85% recall, and an 80% F1-score. Furthermore, predictive analytics significantly improves resource utilization efficiency (85%) and reduces project costs by 10%, compared to 70% and 5%, respectively, achieved by traditional methods. While GBM demonstrates superior overall performance, Logistic Regression (LR) offers favorable precision-recall trade-offs, underscoring the importance of tailored model selection in project risk management.*

**Keywords:** *Predictive Analytics, Project Risk Management, Machine Learning, Risk Prediction, Data-Driven Strategies, Historical Data*

## 1. INTRODUCTION

Despite technological advancements and innovative methodologies, the success rate of IT projects remains suboptimal, with failed projects incurring substantial financial losses. The inherent risks and complexities associated with IT initiatives necessitate robust risk identification and management strategies. Effective collaboration between predictive analytics and risk management is pivotal to addressing these challenges, as traditional methods often fail to predict and mitigate risks dynamically.

Risk management in IT projects involves identifying potential hazards related to requirements, design, resource allocation, technical integration, and feasibility. As each project is unique, risk factors evolve, requiring continuous monitoring and adaptation. Failure to detect and address such risks can lead to complete project collapse. Leveraging historical data for predictive analytics allows for a comprehensive analysis of potential outcomes, their likelihood, and their impact on organizational objectives.

Emerging methodologies, particularly machine learning (ML) and data mining, have demonstrated significant potential in identifying patterns and making accurate predictions from large datasets. This research seeks to explore the integration of predictive analytics and ML to enhance risk identification and management, addressing the persistent high failure rates of IT projects.

### 1.1 Motivation and Novelty

The increasing complexity of modern projects poses significant challenges for traditional risk assessment methods, which struggle to adapt to dynamic environments and large-scale data. This study addresses these limitations by introducing an ML-driven framework for real-time risk prediction. By utilizing historical project data, the approach enhances decision-making and promotes timely risk responses. Feature engineering using t-SNE optimizes data pre-processing, reducing dimensionality while preserving critical predictive factors. This methodology, combined with advanced model selection techniques, differentiates the research from prior studies that relied on static or less sophisticated models.

### 1.2 Contribution

This study makes notable contributions to the domain of project risk management through the following:

- **Data-Driven Risk Identification:** Introduces an innovative ML-based framework for analyzing historical project data to enhance risk management precision and decision-making.
- **Feature Engineering with t-SNE:** Demonstrates the application of t-SNE to manage large datasets, ensuring efficient risk modeling by retaining critical features while reducing dimensionality.
- **Model Selection and Performance:** Highlights the effectiveness of Gradient Boosting Machine (GBM) in predicting project risks and contrasts its performance with Logistic Regression (LR), contributing to advancements in risk modeling techniques.
- **Comprehensive Evaluation Metrics:** Establishes a robust assessment framework using accuracy, precision, recall, and F1-score to ensure reliable and actionable risk predictions.

### 1.3 Justification

The proposed methodology is validated by the unpredictable nature of project interactions, which linear techniques such as Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA) fail to address adequately. The use of t-SNE ensures the preservation of local data structures during feature engineering. Rigorous pre-processing, including handling missing values, outlier removal, and feature scaling, enhances the quality of data inputs for the GBM model. The comprehensive evaluation framework ensures the reliability and practical applicability of the model, contributing to improved risk identification, resource optimization, and cost efficiency in real-world scenarios.

### 1.4 Paper Structure

This paper is structured as follows:

- **Section II:** Literature review highlighting existing risk management approaches.
- **Section III:** Methodology detailing the ML models and evaluation techniques.
- **Section IV:** Results and comparative analysis of ML models.
- **Section V:** Conclusion and future research directions.

Table 1. Summary of previous study on project risk management using machine learning.

| Authors | Focus | Methodologies | Key Findings | Limitations | Future Work |
|---|---|---|---|---|---|
| Roy [11] | Risk assessment in construction. | Risk matrix, analysis of unstructured text and image data. | Highlights ML's role in digitalising construction, emphasises need for expertise, data security issues. | Limited generalizability due to project-specific datasets. | Explore other area for broader insights. |
| Unnamed Authors [12] | Predicting software project failure. | Logistic Regression (LR), Naive Bayes (NB), Support Vector Machine (SVM), Decision Trees (DT), Neural Networks, Adaptive Neuro-Fuzzy Inference Systems. | Develops a reliable risk assessment model applicable to any software project at any lifecycle stage. | Models may not account for all project variables and complexities. | Test models on diverse projects to enhance robustness and adaptability. |
| Elokby et al. [13] | IT project success in telecom and IT industries. | Risk management procedures, qualitative and quantitative evaluations. | Successful project metrics identified; emphasizes comprehensive risk management practices. | Limited to the telecom and IT sectors; results may not be generalised. | Expand the study to include other industries and sectors. |
| Owolabi et al. [14] | Predicting completion risk using Big Data. | Linear Regression, Regression Tree, Random Forest (RF), SVM, Deep Neural Networks (DNNs). | Random Forest found effective in predicting delays in PPP projects with lower average prediction error. | Dependence on historical data may not capture future project dynamics. | Incorporate real-time data and feedback loops for adaptive modelling. |
| Mahdi et al. [15] | Literature review on ML in software risk analysis. | Methodological innovations in ML, literature analysis. | Identifies patterns in ML methodologies; provides foundation for future research in software project risk assessment. | The review scope may overlook recent developments in ML techniques. | Update the review periodically to include emerging methodologies. |
| Burkov et al. [16] | Qualitative risk assessment in projects. | Qualitative risk assessments, three-point risk scale. | Critiques reliance on qualitative evaluations suggest the need for improved risk management strategies. | The subjective nature of qualitative assessments may lead to bias. | Develop quantitative framework to complement assessments. |

## 2. Literature Review

This section presents a comprehensive review of existing methodologies and strategies for risk management in projects. Key studies are summarized in Table 1, emphasizing advancements and gaps in the field.

- **Roy (2023):** Explores ML applications in construction risk management, focusing on data consistency challenges and security concerns in unstructured datasets.
- **2022 Study:** Investigates ML models for early prediction of software project failures, leveraging techniques such as Logistic Regression, Neural Networks, and Adaptive Neuro-Fuzzy Systems.
- **Elokby et al. (2021):** Examines risk management practices in Egypt's IT and telecom industries, highlighting the role of qualitative and quantitative risk evaluations in project success.
- **Owolabi et al. (2020):** Proposes predictive models using Big Data Analytics to forecast delays in PPP projects, demonstrating the efficacy of Random Forest in reducing prediction errors.
- **Mahdi et al. (2020):** Reviews ML methodologies in software risk analysis, identifying trends and offering insights into improving software project outcomes.
- **Burkov et al. (2020):** Discusses qualitative risk assessments and risk aversion techniques, advocating for more advanced methodologies to enhance risk management.
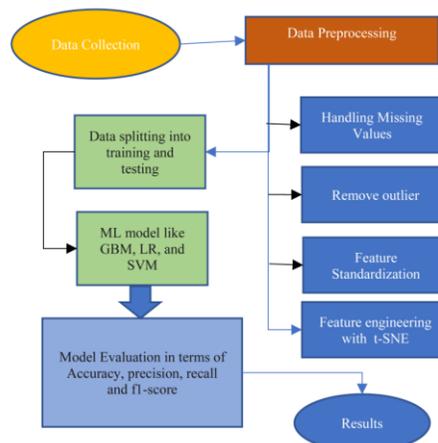- 



Figure 1. Flowchart for project Risk management.

These studies underscore the growing significance of data-driven approaches and ML in modern risk management practices, providing a foundation for the proposed research framework.

### 3. Methods and Materials

This section outlines a systematic methodology for applying predictive analytics and machine learning (ML) to enhance project risk management. The approach begins with comprehensive data collection, followed by rigorous preprocessing and feature engineering to prepare the data for analysis. The Gradient Boosting Machine (GBM) algorithm is employed as the core predictive model, with evaluation based on metrics such as accuracy, precision, recall, and F1-score. This workflow ensures robust risk prediction and supports informed decision-making in project management.

### 3.1 Data Collection

The foundation of this study is a robust dataset encompassing historical project information. This dataset includes:

- Project timelines and task schedules.
- Resource allocations across tasks and projects.
- Project outcomes, such as completion status and success metrics.

By capturing diverse variables influencing project success, the study ensures a comprehensive representation of the factors contributing to project risks.

### 3.2 Data Preprocessing

Preprocessing is essential for ensuring the data's quality and consistency. Key steps include:

- **Handling Missing Values:** Missing data points are imputed using statistical techniques such as the median or mode to preserve dataset integrity.
- **Outlier Removal:** Outliers are identified and excluded to prevent skewed results that could bias the model's learning process.
- **Data Normalization:** Numerical features are standardized to ensure consistent scaling, preventing features with larger magnitudes from disproportionately influencing the model.

These steps ensure that the dataset is clean, reliable, and ready for effective analysis.

### 3.3 Feature Engineering with t-SNE

Feature engineering is a pivotal step in constructing an effective predictive model. This study employs **t-Distributed Stochastic Neighbor Embedding (t-SNE)** to reduce data dimensionality while retaining critical structural relationships.

Unlike linear techniques such as Principal Component Analysis (PCA) or Linear Discriminant Analysis (LDA), t-SNE is adept at capturing non-linear relationships, which are prevalent in project data (e.g., dependencies between tasks and outcomes). This ensures that the predictive model remains focused on the most influential features without sacrificing data integrity.

$$a_{ij}^* = \frac{a_{ij} - \mu_{aj}}{\sigma_{aj}} \tag{1}$$

### 3.4 Data Partitioning

The dataset is divided into two subsets:

- **Training Set (70%):** Used to train the predictive models.
- **Testing Set (30%):** Used to evaluate model performance and generalizability.

This partitioning strategy aligns with standard practices in machine learning, ensuring robust model validation.

### 3.5 Classification Using Gradient Boosting Machine (GBM)

The **Gradient Boosting Machine (GBM)** serves as the primary classification algorithm. GBM leverages ensemble learning to iteratively improve prediction accuracy. Key attributes of GBM include:

- Integration of **decision trees** as base learners.
- Optimization of the loss function through gradient descent, minimizing error iteratively.

This approach ensures accurate and reliable project risk predictions.

$$\hat{F}(x_i) - \min_{f(x_i)} \sum_{i=1}^{n} \mathcal{L}(y_i, F(x_i)) \tag{2}$$

### 3.6 Performance Evaluation

The model's effectiveness is assessed using four key metrics:

1. **Accuracy:** Measures the proportion of correct predictions, reflecting the model's reliability in identifying project risks.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \tag{3}$$

2. **Precision:** Indicates the proportion of true positives among all predicted positives, crucial for minimizing false alarms.

$$Precision = \frac{TP}{TP + FP} \qquad (4)$$

3. **Recall:** Assesses the model's ability to identify actual risks, critical for minimizing undetected risks.

$$Recall = \frac{TP}{TP + FN} \qquad (5)$$

4. **F1-Score:** Provides a harmonic mean of precision and recall, ensuring a balanced evaluation of the model's performance.

$$F1\text{-}Score = \frac{2 * Precision * Recall}{Precision + Recall} \qquad (6)$$

These metrics collectively offer a comprehensive evaluation framework, enabling the identification of optimal models for practical risk management applications.

This methodology demonstrates a rigorous, data-driven approach to project risk prediction, combining advanced ML techniques with thorough evaluation to improve project management outcomes.

## 4. Result Analysis and Discussion

This section evaluates the outcomes of employing predictive analytics for project risk management, focusing on the performance of the Gradient Boosting Machine (GBM) model. Metrics such as accuracy, precision, recall, and F1-score are analyzed to validate the approach. Comparative analysis with other machine learning models (e.g., Logistic Regression and Support Vector Machines) highlights the advantages and trade-offs of the proposed methodology.

Table 2. Historical project data based GBM model performance.

| Measures | Gradient boosting machine |
| --- | --- |
| Accuracy | 85 |
| Precision | 82 |
| Recall | 85 |
| F1-Score | 80 |

### 4.1 Result Analysis

**Performance Metrics:**

The GBM model achieved strong results across the key metrics:

- **Accuracy:** 85%
- **Precision:** 82%
- **Recall:** 85%
- **F1-Score:** 80%

These metrics demonstrate a well-balanced performance, with the GBM model excelling in classification accuracy and risk prediction.
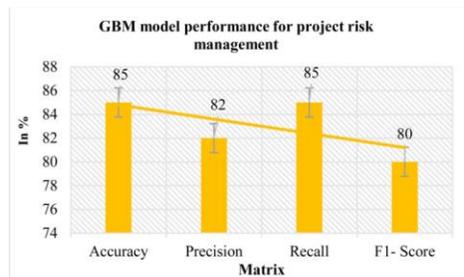


Figure 2. GBM model performance.

**Model Optimization:**

- Figure 3 shows the improvement in model accuracy through hyperparameter tuning.
  - **X-Axis:** Number of iterations.
  - **Y-Axis:** Accuracy percentage.
- The upward trend on the graph underscores the model's adaptability and refinement capabilities, enhancing predictions with each iteration.
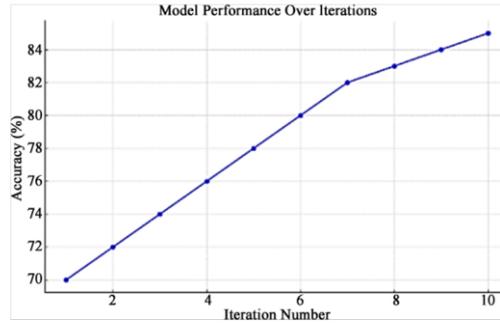
Figure 3. Line graph for model performance over iterations.

**Resource Optimization:**
- Figure 4 compares **Predictive Analytics** with **Traditional Allocation**:
  - **Resource Utilization Efficiency:** Predictive Analytics achieved 85% compared to 70% for traditional methods.
  - **Project Cost Reduction:** Predictive Analytics achieved 10%, double the 5% of traditional methods.
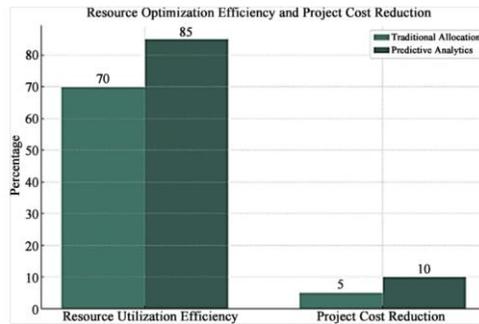  -



Figure 4. Resource optimization efficiency.

This showcases predictive analytics' superior ability to optimize resources and reduce costs, leading to timely and cost-effective project completions.


**Prediction Accuracy:**
- Figure 5 presents a scatter plot comparing predicted and actual project outcomes.
  - Most data points closely align with the ideal prediction line, represented by red dashes.
  - The proximity of predicted values to actual outcomes confirms the GBM model's reliability in forecasting project risks and returns.
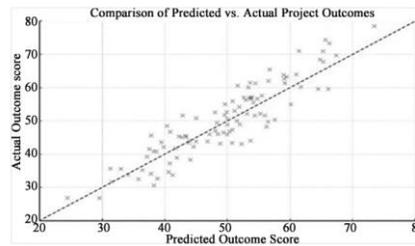


Figure 5. Comparison of predicted vs. actual project outcomes.

**4.2 Discussion**
The results validate the integration of predictive analytics into project management:
- The GBM model exhibited robust performance with accuracy (85%), precision (82%), recall (85%), and F1-score (80%).
- The steady improvement in model accuracy with hyperparameter tuning suggests future iterations could yield even better results.
- Predictive analytics demonstrated significant advantages in resource utilization (85%) and cost reduction (10%), outperforming traditional methods.
- The alignment of predicted and actual project outcomes highlights the model's reliability and effectiveness in real-world scenarios.

Table 3. ML model comparison on historical dataset.

| Model | Accuracy | Precision | Recall | F1-Score |
|-------|----------|-----------|--------|----------|
| GBM | 85 | 82 | 85 | 80 |
| LR | 71 | 83 | 77 | 87 |
| SVM | 83 | 84 | 77 | 83 |

These findings underscore the transformative potential of predictive analytics in project management, especially in improving decision-making and optimizing resource allocation.
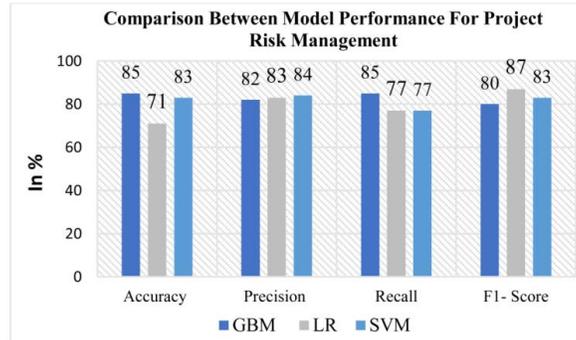


Figure 6. Comparison between model performance.

### 4.3 Comparative Study

A comparative analysis of GBM, Logistic Regression (LR), and Support Vector Machines (SVM) was conducted:

| Metric | GBM | LR | SVM |
|--------|-----|-----|-----|
| Accuracy | 85% | 71% | 83% |
| Precision | 82% | **83%** | **84%** |
| Recall | **85%** | 72% | 77% |
| F1-Score | 80% | **87%** | **83%** |

- **GBM:** Achieves the highest accuracy and recall, making it ideal for overall risk prediction.
- **LR:** Offers the best precision and F1-score, suitable when minimizing false positives is critical.
- **SVM:** Balances accuracy and precision, offering competitive performance across metrics.

These findings highlight the trade-offs between models, with GBM being the most accurate but LR offering the best precision-recall balance.

### 5. Conclusion and Future Scope

**Conclusion:**

This study highlights the effectiveness of predictive analytics in project risk management. The GBM model, with an accuracy of 85% and strong overall performance, emerges as a reliable tool for identifying project risks and optimizing resources. The comparative analysis further demonstrates how different machine learning models cater to varied project management needs.

**Future Scope:**

- **Enhanced Risk Variables:** Future research should include a broader range of risk-related variables to improve prediction accuracy.
- **Advanced Feature Selection:** Exploring more sophisticated feature engineering methods may enhance model performance.
- **Expanded Datasets:** Incorporating larger and more diverse datasets could improve the generalizability of findings.
- **Additional ML Techniques:** Testing alternative machine learning algorithms could provide further insights into effective risk management tools.

This study lays the foundation for advancing machine learning-driven approaches in project risk management, enabling more informed and efficient decision-making.

## References

1. Yang, K., Lin, Y. and Chen, L. (2023) Discovering Critical Factors in the Content of Crowdfunding Projects. Algorithms, 16, Article 51. https://doi.org/10.3390/a16010051

2. Rekha, J.H. and Parvathi, R. (2015) Survey on Software Project Risks and Big Data Analytics. Procedia Computer Science, 50, 295-300. https://doi.org/10.1016/j.procs.2015.04.045

3. Cruz, M.T., Ganapathy, S.A. and Yasin, N.Z.B.M. (2018) Knowledge Management S. R. Bauskar et al. DOI: 10.4236/jdaip.2024.124030579

4. Journal of Data Analysis and Information Processing and Predictive Analytics in IT Project Risks. International Journal of Trend in Scientific Research and Development, 8, 209-216. https://doi.org/10.31142/ijtsrd19142

5. Thomas, J. (2024) Optimizing Bio-Energy Supply Chain to Achieve Alternative Energy Targets. Journal of Electrical Systems, 20, 2260-2273. https://doi.org/10.52783/jes.3176

6. Alotaibi, E.M. (2023) Risk Assessment Using Predictive Analytics. International Journal of Professional Business Review, 8, e01723. https://doi.org/10.26668/businessreview/2023.v8i5.1723

7. Anumandla, S.K.R., Yarlagadda, V.K., Vennapusa, S.C.R. and Kothapalli, K.R.V. (2020) Unveiling the Influence of Artificial Intelligence on Resource Management and Sustainable Development: A Comprehensive Investigation. Technology & Management Review, 5, 45-65.

8. Brandtner, P. (2022) Predictive Analytics and Intelligent Decision Support Systems in Supply Chain Risk Management—Research Directions for Future Studies. In: Yang, X.S., Sherratt, S., Dey, N. and Joshi, A., Eds., Proceedings of Seventh International Congress on Information and Communication Technology, Springer, 549-558. https://doi.org/10.1007/978-981-19-2394-4_50

9. de Langhe, B. and Puntoni, S. (2020) Leading with Decision-Driven Data Analytics. MIT Sloan Management Review.

10. Araz, O.M., Choi, T., Olson, D.L. and Salman, F.S. (2020) Role of Analytics for Operational Risk Management in the Era of Big Data. Decision Sciences, 51, 1320-1346. https://doi.org/10.1111/deci.12451

11. Dimitriadou, A. and Gregoriou, A. (2023) Predicting Bitcoin Prices Using Machine Learning. Entropy, 25, Article 777. https://doi.org/10.3390/e25050777

12. Roy, A. (2023) Risk Analysis of Implementing Machine Learning in ConstructionProjects. https://www.diva-portal.org/smash/record.jsf?pid=diva2:1845289

13. Ibraigheeth, M. and Abu Eid, A.I. (2022) Software Project Risk Assessment Using Machine Learning Approaches. American Journal of Multidisciplinary Research & Development, 4, 35-41.

14. Elokby, E.A., Alawi, N.A., Abdelgayed, A.T.A. and Al-hodiany, Z.M. (2021) Does Project Risk Managemet Matter for the Success of Information Technology Projects in Egypt. 2021 2nd International Conference on Smart Computing and Electronic Enterprise (ICSCEE), Cameron Highlands, 15-17 June 2021, 243-250. https://doi.org/10.1109/icscee50312.2021.9498167

15. Owolabi, H.A., Bilal, M., Oyedele, L.O., Alaka, H.A., Ajayi, S.O. and Akinade, O.O. (2020) Predicting Completion Risk in PPP Projects Using Big Data Analytics. IEEE Transactions on Engineering Management, 67, 430-453. https://doi.org/10.1109/tem.2018.2876321

16. Mahdi, M.N., M.H, M.Z., Yusof, A., Cheng, L.K., Mohd Azmi, M.S. and Ahmad, A.R. (2020) Design and Development of Machine Learning Technique for Software Project Risk Assessment—A Review. 2020 8th International Conference on Information Technology and Multimedia (ICIMU), Selangor, 24-26 August 2020, 354-362, https://doi.org/10.1109/icimu49871.2020.9243459

17. Burkov, V., Burkova, I., Barkalov, S. and Averina, T. (2020) Project Risk Management. 2020 2nd International Conference on Control Systems, Mathematical Modeling, Automation and Energy Efficiency (SUMMA), Lipetsk, 11-13 November 2020, 145-148. https://doi.org/10.1109/summa50634.2020.9280817

18. Thomas, J. (2021) Enhancing Supply Chain Resilience through Cloud-Based SCM S. R. Bauskar et al. DOI: 10.4236/jdaip.2024.124030 580

19. Journal of Data Analysis and Information Processing and Advanced Machine Learning: A Case Study of Logistics. Journal of Emerging Technologies and Innovative Research, 8, e357-e364.

20. Zeng, H., Yang, C., Zhang, H., Wu, Z., Zhang, J., Dai, G., et al. (2019) A LightGBM-Based EEG Analysis Method for Driver Mental States Classification. Computational Intelligence and Neuroscience, 2019, Article 3761203. https://doi.org/10.1155/2019/3761203