# Design And Development Of A Machine Learning Framework To Analyze Software Project Complexity

Sheikh Kashif Ahmed<sup>1</sup>, Lalchandra Gaund<sup>2</sup> <sup>1</sup>Thakur Shayamnarayan Engineering College,Mumbai <sup>2</sup>Thakur Polytechnic, Mumbai (sheikhkashif@tsecmumbai.in, lalchandra.co@tpoly.in )

Abstract—This report explores the use of machine learning algorithms to analyze and predict the complexity of software projects, addressing the limitations of traditional methods such as COCOMO II and Function Point Analysis[1][2]. These methods, while historically significant, struggle to adapt to modern development paradigms like Agile and DevOps due to their reliance on static parameters. The proposed approach leverages machine learning techniques, specifically Random Forest and Principal Component Analysis (PCA), to provide a dynamic and data-driven framework for complexity analysis. The study identifies critical research gaps, including the lack of integration between feature importance ranking and predictive modeling, as well as insufficient real-world validation across diverse datasets. The methodology involves comprehensive data preprocessing, feature selection, and model validation using robust metrics such as precision, recall, and Mean Absolute Error (MAE). By applying these advanced algorithms, this research aims to enhance project planning, risk mitigation, and resource allocation. The findings are expected to contribute significantly to the field of software engineering, providing a scalable solution for analyzing project complexity and bridging the gap between traditional estimation techniques and modern machine learning methodologies[3]. The study also highlights limitations, including dependency on data quality and challenges in generalizing the model across various domains, while proposing future directions for extending this research.

Keywords—COCOMO II, Function Point Analysis, machine learning; software project complexity, Random Forest, Principal Component Analysis(PCA), traditional estimation technique, Mean Absolute Error(MAE)

#### I. INTRODUCTION

Software project complexity significantly impacts resource allocation, cost estimation, and project timelines. Traditional methods often rely on static parameters, making them inadequate for dynamic development environments. With tha.1 A. Traditional Approaches advent of Agile, DevOps, and other iterative development methodologies, there is an urgent need for adaptive and datadriven approaches to assess project complexity effectively[4]. Machine learning provides a promising avenue for addressing these challenges, offering tools capable of handling large datasets, identifying patterns, and making predictions based on real-time inputs. This report explores how machine learning algorithms, such as Random Forest and PCA, can be utilized to analyze software project complexity, bridging the gap

between traditional estimation methods and modern technological demands[5].

#### II. **RATIONALE OF STUDY**

FiThe increasing complexity of software projects stems from factors like rapid technological advancements, shorter development cycles, and the need to accommodate diverse user requirements. Traditional approaches, including expert judgment and algorithmic models like COCOMO II, are often too rigid or subjective, failing to capture the nuances of modern projects. By integrating machine learning algorithms, this study aims to develop a framework that adapts to evolving requirements, enabling accurate complexity analysis and aiding decision-making processes in project management[6]. The rationale also includes addressing inefficiencies in resource allocation, risk assessment, and delivery timelines through predictive insights.

#### III. NATIONAL AND INTERNATIONAL STATUS

BeGlobally, machine learning in software engineering is gaining traction, with significant advancements in areas such defect prediction, cost estimation, and resource as optimization. Countries like the United States, Germany, and China are at the forefront, leveraging AI technologies in software development. Nationally, institutions such as the IITs and NITs in India have contributed to research in predictive modeling and software metrics. Despite these efforts, there is a noticeable gap in applying machine learning to analyze software project complexity3]. Many studies remain limited to academic settings, with insufficient emphasis on real-world validation or integration into industry practices.

#### LITERATURE REVIEW IV.

COCOMO II: Relies on static parameters and predefined equations, making it unsuitable for modern iterative development practices.

Function Point Analysis: Useful for estimating project size but struggles to adapt to evolving requirements and flexible methodologies like Agile.

Expert Judgment: Relies on subjective assessments, often leading to inconsistent results and biases.

#### B. Machine Learning for Complexity

Promising studies have explored defect prediction, cost estimation, and risk assessment using machine learning.

Techniques such as decision trees, support vector machines, and neural networks have been employed, demonstrating improved prediction accuracy compared to traditional models[6].

However, these studies often focus narrowly on specific metrics rather than providing a holistic approach to complexity analysis.

#### C. Research Gaps

Insufficient integration of feature importance ranking with advanced predictive algorithms[7].

Limited exploration of dimensionality reduction techniques like PCA in software project complexity analysis.

Lack of real-world validations using diverse and dynamic datasets.

Minimal focus on addressing multicollinearity and overfitting in predictive modeling.

## V. PROPOSED RESEARCH METHODOLOGY

### A. Algorithm Selection

Random Forest is an ensemble learning technique that builds multiple decision trees and combines their outputs for a more accurate and stable prediction[10]. It works by leveraging two kev ideas:

Bootstrap Sampling (Bagging): Training each tree on a random subset of data with replacement.

Feature Randomness: Considering only a random subset of features at each split in the tree.

By aggregating the predictions of multiple trees, Random Forest minimizes overfitting (common in individual decision trees) and improves generalization.

Steps in Random Forest

**1.Data Preparation** 

Split the data into training and testing datasets.

2.Bootstrap Sampling

For each tree, a random subset of the training data is selected with replacement (some samples may appear multiple times). **3.Tree Building** 

At each split in the tree, only a random subset of features is considered to find the best split. This introduces feature diversity among the trees.

4.Aggregation (Ensemble)

Classification: Combine predictions of all trees using majority voting.

Regression: Take the average of predictions across all trees.

#### Example

Suppose you are predicting whether a loan applicant is Low Risk or High Risk using the features: Income, Age, Credit Score, and Loan Amount.

**Bootstrap Sampling:** 

From the training dataset of 1,000 applicants, multiple subsets are created with replacement.

Tree 1 gets Dataset A, Tree 2 gets Dataset B, and so on.

Tree Building:

Tree 1 might split based on Age and Loan Amount.

Tree 2 might split based on Credit Score and Income.

Each tree learns patterns based on its sampled data and selected features.

Prediction by Each Tree:

After training, each tree predicts for a new applicant:

Tree 1: Low Risk.

Tree 2: High Risk.

Tree 3: Low Risk.

Tree 4: Low Risk.

Tree 5: High Risk.

Aggregation (Ensemble Voting):

Classification Aggregation: Take a majority vote from all trees. Predictions: [Low Risk, High Risk, Low Risk, Low Risk, High Risk].

Majority vote  $\rightarrow$  Low Risk (3 votes for Low Risk vs. 2 for High Risk).

Regression Aggregation: Take the average prediction from all trees.

Final Prediction:

The Random Forest outputs Low Risk as the aggregated result.

PCA is a linear transformation technique used for dimensionality reduction. It identifies new axes (principal components) that capture the most variance in the data, enabling you to represent data in fewer dimensions while preserving as much information as possible[11].

Steps in PCA:

Standardize the Data:

Convert features to have zero mean and unit variance.

Compute Covariance Matrix:

The covariance matrix measures the relationships between features. For n features, it's an n×n matrix.

Find Eigenvalues and Eigenvectors:

Eigenvalues indicate the amount of variance captured by each principal component.

Eigenvectors define the direction of the principal components.

Select Principal Components:

Sort eigenvalues in descending order and select the top k eigenvalues and their corresponding eigenvectors.

Transform Data:

Project the original data onto the selected principal components.

### **Example of PCA:**

Suppose you have a dataset with three features:

Height (cm)

Weight (kg) Age (years)

Standardize the Data:

Mean-center and normalize Height, Weight, and Age.

Compute Covariance Matrix:

A 3x3	matrix	sho	wing	how	features	vary	together:
Covariance Matrix =			Var(Height)			Cov(Height, Weight)	
			Cov(Weight, Height)			$\operatorname{Var}(\operatorname{Weight})$	
			Cov	v(Age,	$\operatorname{Height})$	Cov	(Age, Weight)

**Eigenvalues and Eigenvectors:** 

Eigenvalues:  $[2.8, 1.1, 0.1] \rightarrow$  Variance captured by each component.

Eigenvectors: Directions of principal components.

Select Top-k Components:

Retain the first 2 components (with the largest eigenvalues) to reduce dimensionality.

Project Data:

Transform the data onto the 2 principal components.

Result: A 2D representation of your data, combining Height and Weight into one component while retaining most of the variance.

Combined Example: Random Forest + PCA

Suppose you have a dataset with 50 features to predict credit default risk.

Apply PCA:

Reduce 50 features to the top 10 principal components that capture 95% of the variance.

This step reduces redundancy and noise.

Train Random Forest:

Use the reduced 10-component dataset to train a Random Forest model.

Random Forest will build diverse trees based on these components and output predictions.

Make Predictions:

Use the Random Forest ensemble to predict credit risk with high accuracy and less computational cost.

#### VI. IMPLEMENTATION PLAN

Phase 1: Literature Review and Dataset Preparation

Conduct an extensive review of existing literature. Identify and preprocess datasets for machine learning.

Phase 2: Model Development

Implement Random Forest for feature selection. Apply PCA for dimensionality reduction.

Phase 3: Validation and Testing

Validate the model using real-world datasets. Compare results with traditional methods.

Phase 4: Documentation and Publication

Document findings and submit to peer-reviewed journals.

### VII. EXPECTED RESEARCH OUTCOMES

1.Enhanced Predictions Improved accuracy in predicting software project complexity compared to traditional methods.

2.Better Project Management Effective resource allocation, risk assessment, and timeline management.

3.Contributions to the Field Development of a replicable framework for complexity analysis in software engineering, bridging the gap between traditional estimation techniques and modern machine learning methodologies.

#### REFERENCES

Mow (Henight & Aggensen, M. (2003). On the sensitivity of COCOMO Coof (Wreight exige) on model. International Conference on Software Engineering, Var(Age) Dubey, S. S., & Rana, A. (2018). Limitations of Function Point for Agile

- [2] Software Environment: A Case Study. International Journal of Computer Science and Technology.
- [3] Dubey, S. K., Rana, A., & Singh, R. K. (2023). Software Complexity Prediction Model: A Combined Machine Learning Approach. Advances in Intelligent Systems and Computing.
- [4] Khan, M. A., Khan, S. U., & Anjum, A. (2024). Hybrid Feature Selection Method for Predicting Software Defect. Journal of Engineering and Applied Science.
- Chen, H., Xu, B., & Zhong, K. (2024). Enhancing Software Effort [5] Estimation through Reinforcement Learning-Based Feature Selection. arXiv preprint arXiv:2403.16749.
- [6] Osman, H., Ghafari, M., & Nierstrasz, O. (2018). The Impact of Feature Selection on Predicting the Number of Bugs. arXiv preprint arXiv-1807 04486
- [7] Tran, N., Tran, T., & Nguyen, N. (2024). Leveraging AI for Enhanced Software Effort Estimation: A Comprehensive Study and Framework Proposal. arXiv preprint arXiv:2402.05484.
- [8] Uc-Cetina, V. (2023). Recent Advances in Software Effort Estimation using Machine Learning. arXiv preprint arXiv:2303.03482.
- Singh, A., Singh, R., & Agrawal, S. (2022). Machine Learning for [9] Predicting Software Project Complexity. Journal of Software Engineering Research and Development.
- [10] Gupta, M., & Sharma, P. (2021). Evaluating the Effectiveness of PCA in Software Cost Estimation Models. International Journal of Advanced Computer Science and Applications.
- Kumar, S., & Mehta, R. (2023). AI-Driven Solutions for Modern [11] Software Development Challenges. IEEE Transactions on Software Engineering.



Mtech(CSE), Assistant Professor at TCET, Mumbai with 18+ year of teaching experience.



ME(pursuing), Lecturer at TPoly, Mumbai with 5+ years of experience.