

# Project Proposal: Predicting JIRA Task Estimation Using Machine Learning

## Domain Background

The field of software engineering and project management heavily relies on accurate task estimation for effective planning and resource allocation. JIRA, a widely used task management tool, allows teams to manage and track tasks, assign priorities, and measure progress. However, estimating the time or effort required to complete a task, particularly in software development, remains challenging due to varying complexities, team dynamics, and dependencies. Leveraging machine learning to predict task duration or story points can revolutionize project planning by reducing uncertainties and improving productivity.

Recent research has explored the application of machine learning techniques to improve estimation accuracy. For instance, Sousa, André O., et al. "Applying Machine Learning to Estimate the Effort and Duration of Individual Tasks in Software Projects." IEEE Access (2023) investigated the application of machine learning algorithms to predict effort and duration for individual tasks in software projects, demonstrating the potential of ensemble methods like Random Forest and XGBoost in improving estimation accuracy (IEEE link). Additionally, Martens, Annelies, and Mario Vanhoucke. "Integrating corrective actions in project time forecasting using exponential smoothing." Journal of Management in Engineering 36.5 (2020) applied exponential smoothing techniques to project duration forecasting, highlighting the potential of integrating statistical methods with traditional project management practices.

These studies underscore the potential of combining machine learning with project management tools like JIRA to enhance task duration estimation, thereby improving planning and resource allocation in software development projects.

## Problem Statement

Accurate task estimation is critical for meeting project deadlines and maintaining team efficiency. Manual estimation methods often suffer from subjectivity and inconsistency, leading to missed deadlines and resource mismanagement. The primary problem is: How can machine learning algorithms be used to predict the time or story points required for JIRA tasks based on historical data and task attributes? By solving this problem, project managers and software teams can achieve better predictability and resource allocation.

## Datasets and Inputs

For this project, the dataset "JIRA Dataset Public" from Kaggle will be utilized. This dataset contains historical data on JIRA tasks and includes features such as:

- Task Metadata: Task ID, title, description, type (e.g., Bug, Feature, Improvement), and priority.
- Date Attributes: Creation date, assignment date, and resolution date.
- Task Dependencies: Information on linked or blocked tasks.
- Comments: Number of comments, reflecting task activity and complexity.
- Outcome Metrics: Time to resolution (duration) and story points.

These inputs will serve as the foundation for training a machine learning model to predict story points and completion time.

The target variable, duration, represents the time taken to resolve a task. This field is derived by subtracting the created column (timestamp indicating when the task was created) from the resolved column (timestamp indicating when the task was resolved). The original dataset contains 49,000 rows and 491 columns.

However, after cleaning the data to include only rows with a resolved date and a resolution value of "fixed," the final dataset size will be smaller, ensuring relevance and accuracy for model training. This refined dataset provides a quantitative measure of task completion time, which is essential for building predictive models.

## Solution Statement

The proposed solution involves developing a machine learning model that leverages historical JIRA task data to predict the estimated time to complete a task or its story points. This solution will:

1. Preprocess the dataset by cleaning and encoding categorical and textual data.
2. Use natural language processing (NLP) techniques to analyze textual fields such as the title and description.
3. Train a supervised learning model using features such as task type, priority, number of dependencies, and historical metrics.
4. Provide predictions for story points and completion time for new tasks.

The solution will also include interpretability features to explain the predictions and allow project managers to make informed decisions.

## Benchmark Model

As a baseline, a simple linear regression model will be implemented to predict task duration based on numerical features such as the number of dependencies and priority. This benchmark will provide a point of comparison for more advanced models. Additionally, basic heuristics, such as calculating the average task duration or story points for similar task types, will serve as an initial comparison metric.

## Evaluation Metrics

To measure the performance of the proposed model, Mean Absolute Error (MAE) will be used as the primary evaluation metric. MAE provides a clear and interpretable measure of the average absolute difference between the predicted and actual task durations or story points. By focusing on MAE, the project ensures a consistent and comparable metric across different models. Additionally, RMSE and  $R^2$  may be reported as supplementary metrics for a deeper analysis of model performance.

## Outline of Project Design

1. Data Preprocessing:
  - Load and clean the dataset, handling missing values and inconsistencies.
  - Perform feature engineering, including:
    - Encoding categorical variables (e.g., priority, task type).
    - Extracting insights from textual fields using NLP techniques like BERT or TF-IDF.
    - Calculating derived metrics such as task dependencies and activity levels.
2. Model Development:
  - Implement baseline models, such as linear regression and basic heuristics.
  - Train advanced machine learning models like Random Forest, XGBoost, or LightGBM for structured data.
  - Experiment with deep learning models for textual data, such as LSTM or Transformers.
  - Combine structured and unstructured data using multi-modal approaches.

### 3. Model Evaluation:

- Split the dataset into training, validation, and test sets.
- Evaluate models using the defined metrics and compare against the benchmark model.

### 4. Deployment:

- Deploy the final model as a REST API for integration with JIRA or other task management systems.
- Provide an interface for users to input new task details and receive predictions.

### 5. Testing and Iteration:

- Test the model on new data to ensure generalizability and accuracy.
- Incorporate feedback and refine the model as needed.

### 6. Deliverables:

- A trained machine learning model capable of predicting story points and task duration.
- A detailed report and visualizations explaining model predictions.
- An API or tool for integration with project management systems.

## Conclusion

This project aims to leverage historical JIRA data to build a machine learning model that can predict task duration and story points. By automating and improving estimation accuracy, this tool will provide project managers and developers with actionable insights, leading to better resource planning and execution. The project's success will be measured by its ability to outperform benchmark models and its practical utility in real-world scenarios.