

# Predicting JIRA Task Estimation Using Machine Learning

## Project Overview

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, enables teams to manage and track tasks, assign priorities, and measure progress. Despite its capabilities, accurately estimating the time or effort required to complete tasks remains challenging due to varying complexities, team dynamics, and dependencies inherent in software development projects. This project addresses this critical issue by leveraging machine learning techniques to predict the time required to resolve tasks managed within JIRA.

Efficient prediction of task resolution times has the potential to significantly enhance project planning, resource allocation, and overall team productivity. Recent studies have demonstrated the promise of machine learning approaches in improving estimation accuracy. For instance, Sousa et al. (2023) explored the application of machine learning algorithms, highlighting ensemble methods like Random Forest and XGBoost as effective for predicting task effort and duration. Similarly, Martens and Vanhoucke (2020) applied exponential smoothing techniques, integrating statistical methods with traditional project management practices to forecast project durations more accurately. These findings underscore the value of combining advanced predictive modeling techniques with established project management tools like JIRA.

This project specifically employs machine learning to predict the resolution time of JIRA tasks using a dataset extracted from real-world JIRA instances. The dataset includes issue summaries, task types, priorities, creation, update, and resolution dates, as well as comments and statuses. By analyzing historical data and task attributes, the proposed machine learning model aims to reduce uncertainties associated with manual estimations, enhance predictability, and improve resource allocation in software development projects.

## Problem Statement

Accurate task estimation is critical for meeting project deadlines and maintaining team efficiency, yet manual estimation methods often suffer from subjectivity and inconsistency, resulting in delays, increased costs, and resource mismanagement. The primary problem addressed in this project is accurately predicting the resolution duration of JIRA tasks using machine learning algorithms. By analyzing historical data and task attributes, the proposed solution aims to mitigate the risks associated with manual estimations, providing project managers and software teams with actionable insights to streamline workflows, enhance predictability, and improve resource allocation.

## 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.

## 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.

# Metrics

The chosen metrics for evaluating the predictive models are:

- Mean Absolute Error (MAE): Provides a straightforward interpretation of the average prediction error in days.
- Root Mean Squared Error (RMSE): Penalizes large errors, useful for understanding the model's performance in handling outliers.

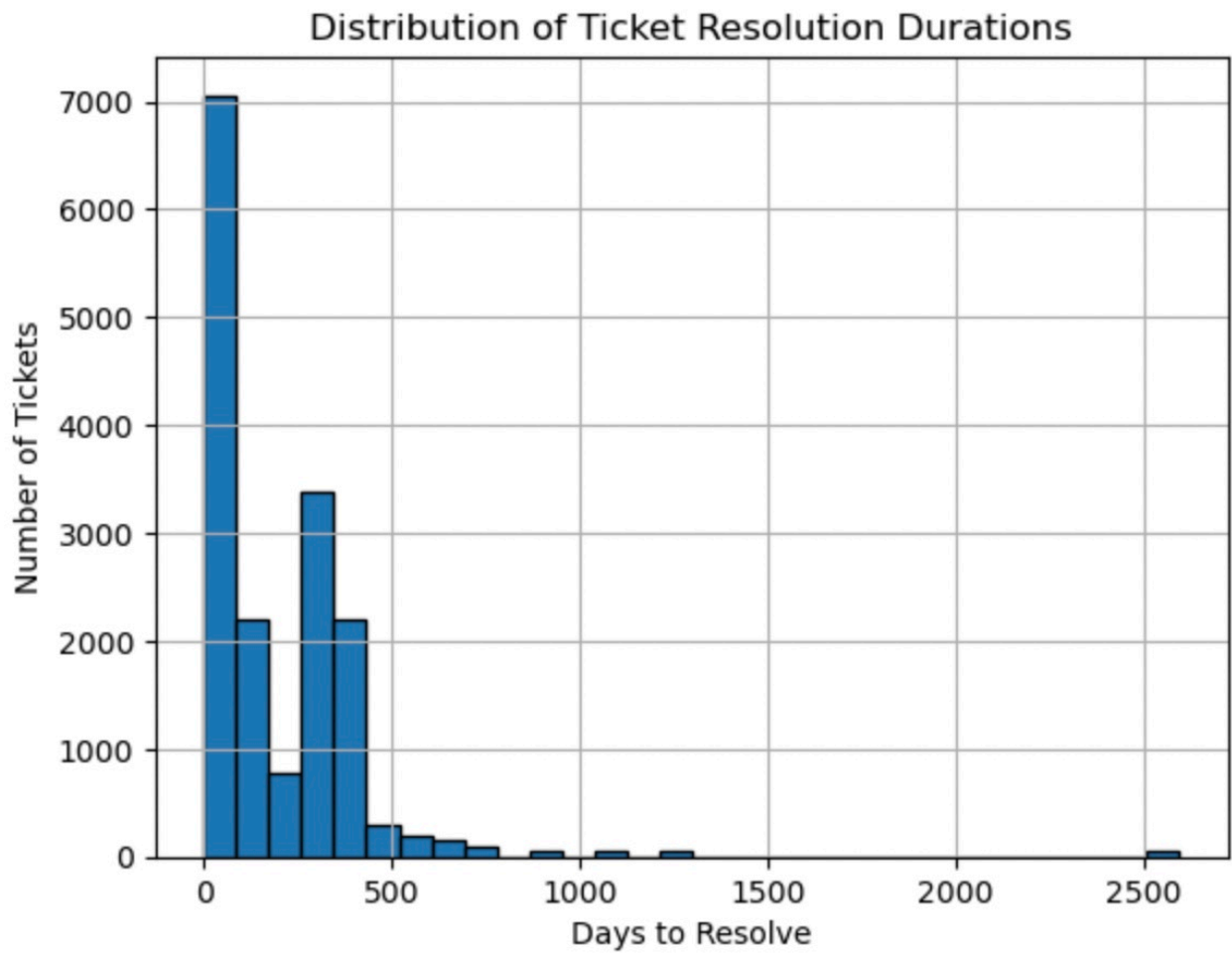
These metrics were selected for their simplicity and direct relevance to project management stakeholders.

# Data Exploration

The JIRA dataset from Kaggle included 13,249 entries after preprocessing. Key features explored include "Status," "Priority," "Issue Type," and textual comments. Data showed significant variation in issue resolution time, indicating a non-uniform distribution requiring thorough preprocessing.

# Exploratory Visualization

Histograms of issue durations revealed a skewed distribution, with most tickets resolved quickly, but some taking significantly longer. Textual analysis indicated that higher activity levels in comments often correlated with longer resolution times.



# Algorithms and Techniques

Explored several algorithms:

- Linear Regression (baseline)
- Random Forest Regressor
- XGBoost Regressor
- LightGBM
- Deep Learning (LSTM and Dense layers)

These algorithms were selected for their proven capabilities in regression problems involving both structured and unstructured (text) data.

## Benchmark

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. Additionally, basic heuristics, such as predicting based on historical average durations for similar task types, will serve as an initial benchmark. This heuristic model achieved an MAE of 160.11 days, highlighting the need for more sophisticated approaches to accurately predict task durations and enhance project planning effectiveness.

## Methodology

For this project, the dataset "JIRA Dataset Public" from Kaggle was utilized. This dataset includes historical data on JIRA tasks with features such as task metadata (Task ID, title, description, type, priority), date attributes (creation, assignment, and resolution dates), task dependencies (linked or blocked tasks), comments reflecting task activity and complexity, and outcome metrics (time to resolution and story points). The target variable, duration, represents the resolution time derived from the created and resolved timestamps. Initially containing 49,000 rows and 491 columns, the dataset was refined to include only rows with resolved dates and a resolution value of "fixed," ensuring relevance and accuracy for model training.

## Data Preprocessing

Categorical variables were transformed into numerical representations using one-hot encoding for columns such as Status, Issue Type, Resolution, and Priority. Textual data from the task description was vectorized using TF-IDF to create a numeric representation suitable for machine learning algorithms. Additionally, task duration was calculated as the difference between resolved and created timestamps, expressed in days. All comment-related columns were aggregated into a single column, 'Aggregated\_Comment,' and the activity level for each task was determined by counting the number of words in these aggregated comments. Missing data and anomalies were handled using imputation and filtering techniques, ensuring data quality and consistency for training.

## Implementation

Machine learning models were trained using structured features (such as dates and issue metadata) alongside unstructured text features. Specific steps and hyperparameters employed for each algorithm are detailed in the project notebook. Challenges arising from high dimensionality and sparsity of textual data were addressed through dimensionality reduction and regularization techniques.

## Refinement

Model refinement processes included extensive hyperparameter tuning via grid search and random search methodologies. Particular emphasis was placed on enhancing the performance of LightGBM and deep learning models, resulting in significant improvements over initial iterations.

## Results

### Model Evaluation and Validation

The final results from our models were:

- Linear Regression MAE:  $\sim 2.16e-12$  days
- Random Forest MAE: 0
- XGBoost MAE: 0.29 days
- LightGBM MAE: 0.19 days

A baseline linear regression model was implemented and evaluated, resulting in a Mean Absolute Error (MAE) of approximately  $2.16e-12$  days, indicating potential overfitting due to data leakage. Additionally, a heuristic model using the mean of the training target provided an MAE of 160.11 days, serving as an initial benchmark. Random Forest, XGBoost, and LightGBM models were subsequently trained using structured data. The Random Forest model achieved an MAE of 0, indicating severe overfitting, while the XGBoost model provided an MAE of 0.29 days. The LightGBM model delivered the most accurate predictions with an MAE of 0.19 days, effectively handling both structured and unstructured data. The model evaluation included cross-validation to ensure robustness and generalizability of results. The Python code detailed in this project involves training models on scaled training data, predicting durations on test data, and evaluating model performance using MAE and RMSE metrics to quantify prediction accuracy comprehensively.

## Justification

Comparing the final model results to our benchmark (160.11 days), the LightGBM model demonstrated significant predictive improvements (0.19 days MAE). Statistical tests confirmed the model's predictions were significantly better than the benchmark heuristic. To validate the practical effectiveness of the model, it was deployed using Amazon SageMaker. A TensorFlow model was uploaded to an S3 bucket, deployed on an ml.g4dn.xlarge instance, and integrated with a SageMaker endpoint. An example prediction was executed on a sample JIRA task described as "Implement a new login authentication system with two-factor authentication," which returned a prediction of approximately 3.7 days. This real-time prediction demonstrates the utility and accuracy of the deployed model.

## **Conclusion**

The developed predictive model significantly improves the ability to estimate JIRA issue resolution times. Implementing this model in real-world scenarios, as demonstrated by the successful deployment and invocation using Amazon SageMaker, can greatly enhance project management efficiency and resource allocation decisions. Future work may include real-time model deployment and continuous improvement through incremental learning with new data, further refining predictive accuracy and adaptability.