

# **AutoJudge: Programming Problem Difficulty and Score Prediction Using Machine Learning**

Divyansh Bansal  
IIT Roorkee  
B.Tech (Second Year)

2024-28

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Problem Statement</b>	<b>3</b>
<b>3</b>	<b>Dataset Description</b>	<b>3</b>
<b>4</b>	<b>Data Preprocessing</b>	<b>3</b>
<b>5</b>	<b>Feature Engineering</b>	<b>4</b>
5.1	Text Features . . . . .	4
5.2	Handcrafted Features . . . . .	4
<b>6</b>	<b>Model Architecture</b>	<b>4</b>
6.1	Classification . . . . .	4
6.2	Regression . . . . .	4
<b>7</b>	<b>Experimental Setup</b>	<b>4</b>
<b>8</b>	<b>Results and Evaluation</b>	<b>5</b>
8.1	Classification Results . . . . .	5
8.2	Regression Results . . . . .	5
<b>9</b>	<b>Web Interface</b>	<b>5</b>
<b>10</b>	<b>Sample Predictions</b>	<b>7</b>
<b>11</b>	<b>Conclusion</b>	<b>7</b>
<b>12</b>	<b>References</b>	<b>7</b>

# 1 Introduction

Competitive programming platforms contain a large number of problems with varying difficulty levels. These difficulty labels are typically assigned manually, making them subjective and inconsistent. Automating difficulty estimation can help learners, instructors, and problem setters.

This project, **AutoJudge**, uses Natural Language Processing (NLP) and Machine Learning techniques to predict the difficulty of programming problems using only their textual descriptions. The system outputs both a categorical difficulty label and a continuous difficulty score.

## 2 Problem Statement

The objective of this project is to build a machine learning system that:

- Predicts the difficulty class (Easy / Medium / Hard)
- Predicts a numerical difficulty score (1–5)
- Provides interpretable insights for predictions

The main challenges include class imbalance, overlapping difficulty characteristics, and stable regression.

## 3 Dataset Description

The dataset consists of 4,112 competitive programming problems stored in JSON Lines format. Each entry contains the problem title, description, input/output details, a difficulty label, and a numeric difficulty score.

The dataset is imbalanced, with Hard problems appearing more frequently than Easy ones.

## 4 Data Preprocessing

The following preprocessing steps were applied:

- Lowercasing and text normalization
- Removal of special characters and extra spaces
- Concatenation of all textual fields
- Title weighted twice to emphasize importance

These steps ensure clean and consistent input for feature extraction.

## 5 Feature Engineering

### 5.1 Text Features

TF-IDF vectorization was applied using unigrams, bigrams, and trigrams with stop-word removal and sublinear term frequency scaling.

### 5.2 Handcrafted Features

Additional features include:

- Text length and word count
- Sentence statistics
- Mathematical symbol density
- Numeric constraint detection (e.g.,  $n \geq 10^5$ )
- Algorithmic pattern detection (DP, graphs, BFS/DFS, binary search)

All numeric features were standardized using `StandardScaler`.

## 6 Model Architecture

### 6.1 Classification

Three classifiers were trained: Gradient Boosting, Random Forest, and Support Vector Machine. SMOTE was used to handle class imbalance. An ensemble approach combined the probabilistic outputs of these models.

### 6.2 Regression

To avoid unstable predictions, class-conditioned regression was implemented. Separate Random Forest regressors were trained for Easy, Medium, and Hard problems.

## 7 Experimental Setup

- Train-test split: 80% training, 20% testing
- Classification metrics: Accuracy, Precision, Recall, F1-score
- Regression metrics: MAE and RMSE

All experiments were conducted using Python and Scikit-learn.

## 8 Results and Evaluation

### 8.1 Classification Results

The ensemble classifier achieved an accuracy of **54.92%**.

Class	Precision	Recall	F1-score
Easy	0.52	0.41	0.45
Medium	0.45	0.32	0.38
Hard	0.60	0.77	0.67

Table 1: Classification Performance

<b>Actual \ Predicted</b>	<b>Easy</b>	<b>Hard</b>	<b>Medium</b>
<b>Easy</b>	62	46	45
<b>Hard</b>	24	299	66
<b>Medium</b>	34	156	91

Table 2: Confusion Matrix for Difficulty Classification

The confusion matrix shows that most misclassifications occur between neighboring difficulty levels (Easy–Medium and Medium–Hard), which reflects the inherent ambiguity in difficulty labeling rather than random errors.

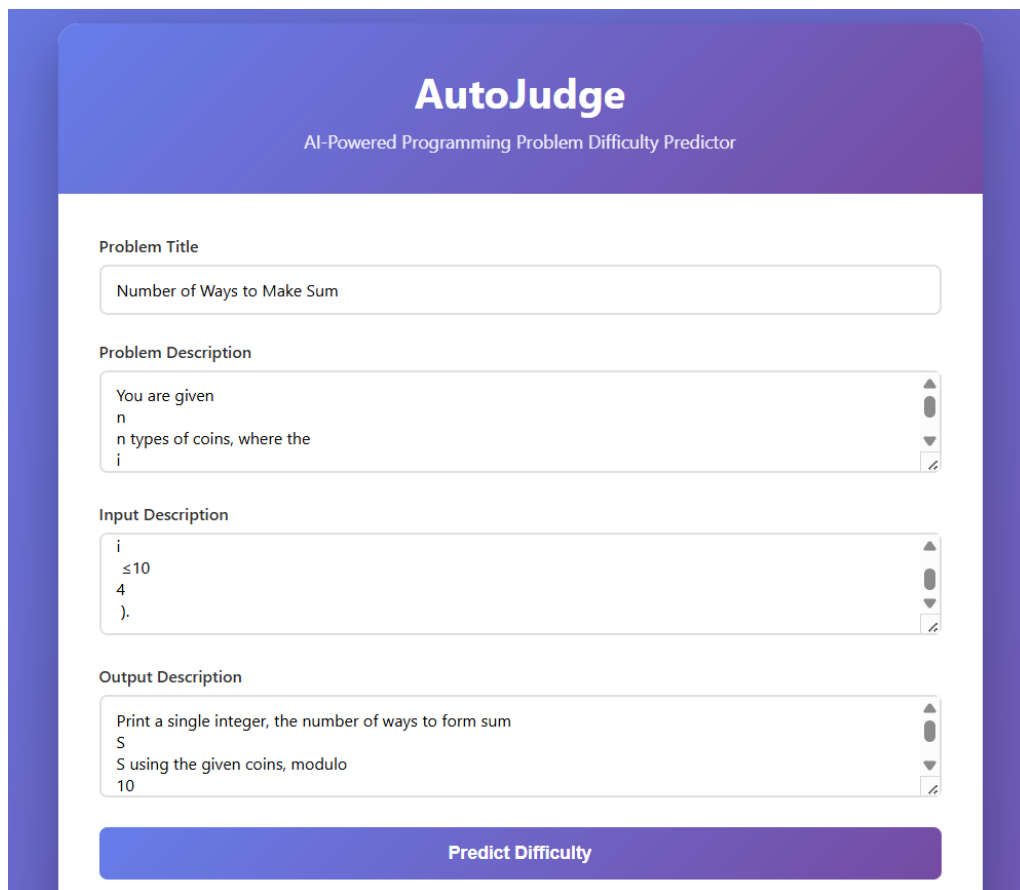
### 8.2 Regression Results

- MAE: 1.366
- RMSE: 1.775

These results indicate stable and realistic score predictions.

## 9 Web Interface

A Flask-based web application was developed to demonstrate real-time predictions.



**AutoJudge**  
AI-Powered Programming Problem Difficulty Predictor

**Problem Title**  
Number of Ways to Make Sum

**Problem Description**  
You are given  
n  
n types of coins, where the  
i

**Input Description**  
i  
≤10  
4  
).

**Output Description**  
Print a single integer, the number of ways to form sum  
S  
S using the given coins, modulo  
10

**Predict Difficulty**

Figure 1: AutoJudge Web Interface

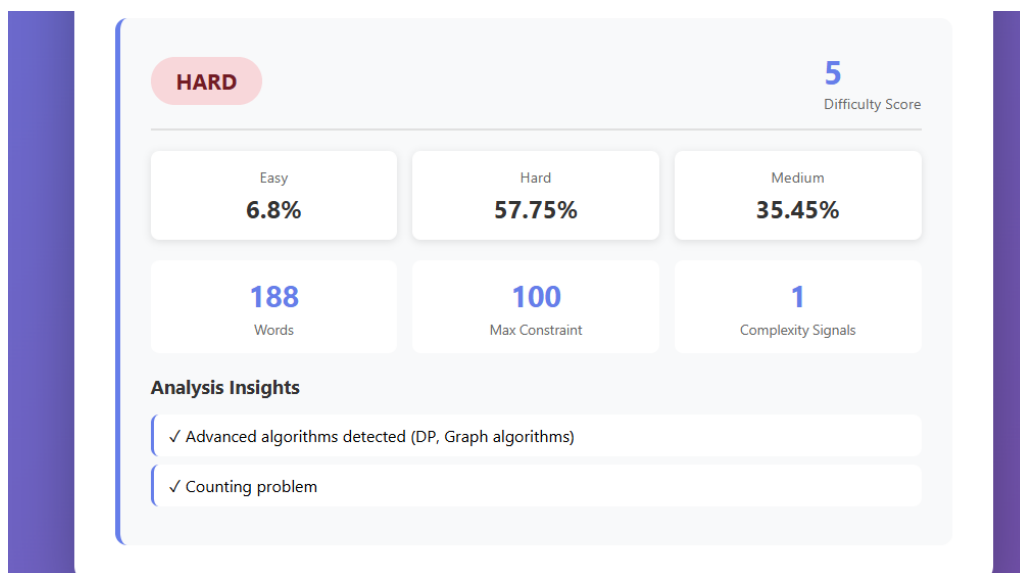


Figure 2: Predicted Difficulty and Score Output

The interface displays predicted difficulty, score, and algorithmic insights.

## 10 Sample Predictions

- Simple array problem: Easy, score  $\approx 1.6$
- Graph shortest path problem: Hard, score  $\approx 4.8$

## 11 Conclusion

AutoJudge demonstrates how NLP and machine learning can be used to estimate programming problem difficulty using only textual data. The ensemble classification and class-conditioned regression approach ensures stable and interpretable predictions.

Future improvements may include transformer-based models and the use of solution code.

## 12 References

- Scikit-learn Documentation
- NLP Feature Engineering Techniques
- Competitive Programming Platforms