# Code Smell Detection in Python Code Using Machine Learning

Mayur Jain
*Department of Computer Science*
*California State University, Long Beach*
California, USA
mayurjain333@gmail.com

*Abstract*—This paper presents a machine learning approach for detecting code smells in Python code snippets. Leveraging the Code4ML dataset, which contains annotated code snippets from Kaggle notebooks, we develop a model that uses TF-IDF vectorization and a Random Forest classifier to classify code blocks based on structural issues. Hyperparameter tuning and model optimization are incorporated to achieve a high level of accuracy. The final model demonstrates satisfactory precision, recall, and accuracy, making it suitable for automated code review and quality assurance tasks.

## I. INTRODUCTION

Code smells are indicators of potential structural weaknesses in code that could lead to maintainability and scalability issues. While not bugs, they make the code more complex and harder to modify. Detecting code smells manually is both time-consuming and subjective. Therefore, automated solutions are essential, especially in large-scale projects. This research aims to create a machine learning-based code smell detection model for Python code using the Code4ML dataset, a large corpus of annotated Python code extracted from Kaggle notebooks.

## II. DATASET

The Code4ML dataset used in this study is an enriched, large-scale corpus of annotated Python code snippets and metadata. The dataset, sourced from Kaggle, contains approximately 2.5 million code snippets from around 100,000 Jupyter notebooks. Key components of the dataset include:

- **code_blocks.csv**: Contains the raw code snippets.
- **kernels_meta.csv**: Metadata for each notebook, including Kaggle scores, comments, and upvotes.
- **competitions_meta.csv**: Describes Kaggle competitions, including data type and competition task.
- **markup_data.csv**: Annotated code blocks with semantic types, enabling deeper structural analysis.
- **vertices.csv**: Maps semantic types and subclasses to numeric IDs, facilitating interpretability of annotated code blocks.

These files are interconnected, allowing for comprehensive metadata-driven insights. For instance, `code_blocks.csv` is linked to `kernels_meta.csv` via the `kernel_id` column, while `kernels_meta.csv` is linked to `competitions_meta.csv` through `comp_name`. The dataset provides a robust foundation for detecting code smells based on patterns in large-scale Python code.

## III. METHODOLOGY

Our approach leverages a Random Forest Classifier trained on code snippets represented by TF-IDF vectorized features. The model is fine-tuned using hyperparameter tuning to maximize classification performance.

### A. Data Preprocessing

The `code_block` column in `code_blocks.csv` contains raw Python code snippets, which are transformed into TF-IDF vectors to serve as model features. TF-IDF (Term Frequency-Inverse Document Frequency) captures the relative importance of words across different code blocks, facilitating better feature representation for machine learning models.

### B. Model Training

We use a Random Forest Classifier due to its robustness with high-dimensional data. Additionally, RandomizedSearchCV is employed for hyperparameter tuning, optimizing parameters such as the number of estimators, max depth, and minimum samples split to achieve high accuracy.

### C. Model Comparison

To evaluate our model, we compared its performance with other standard ML algorithms, including:

- **Logistic Regression**: A linear model that provided lower precision and recall due to its limitations with complex, non-linear relationships.
- **Support Vector Machine (SVM)**: Achieved reasonable accuracy but required more computational power, making it less suitable for large datasets.
- **Decision Tree Classifier**: Performed moderately well but lacked the ensemble power of the Random Forest model, leading to overfitting on complex code patterns.

The Random Forest model outperformed these alternatives in terms of accuracy and computational efficiency, particularly due to its ensemble nature and robustness to overfitting.

## IV. ENHANCEMENTS IN CODE SMELL DETECTION AND VISUALIZATION

This section describes the updates made to the 'detect_code_smells.py' and 'visualization_code_smells.py' scripts as part of the project enhancements.

## A. Modifications in Code Smell Detection

The 'detect_code_smells.py' script was updated to include additional types of code smells for more comprehensive detection. These include:

- Long Method: Detects methods with lines exceeding 20 tokens.
- Duplicate Code: Identifies repetitive code patterns.
- Magic Numbers: Flags hardcoded numeric values not part of loops or ranges.
- Too Many Arguments: Checks for functions with more than three arguments.
- Dead Code: Detects unused variables or placeholder code (e.g., 'pass', 'TODO').
- Large Classes: Identifies classes with more than five methods or attributes.
- Complex Conditions: Flags overly complex conditional statements.

The 'additional_smells' column in the results now captures all these smell types, enabling better categorization and visualization.

## B. Enhancements in Visualization

The 'visualization_code_smells.py' script was enhanced to provide additional insights. Key updates include:

- Added a stacked bar chart showcasing the top 10 kernels with the most code smells.
- Introduced a bar chart visualizing the distribution of each code smell type using data from the 'additional_smells' column.
- Improved data cleaning to handle and accurately count smell types, ensuring no extraneous symbols (e.g., brackets, quotes) affect the analysis.

## C. Visual Results

The following visual results were generated as part of the analysis and visualization:



Fig. 1. Stacked bar chart showing the top 10 kernels with the most code smells.

## V. RESULTS

Our final model achieved high classification performance, as shown in Table I, with an accuracy of 96.6%, precision of 97%, and recall of 84% for detecting code smells.
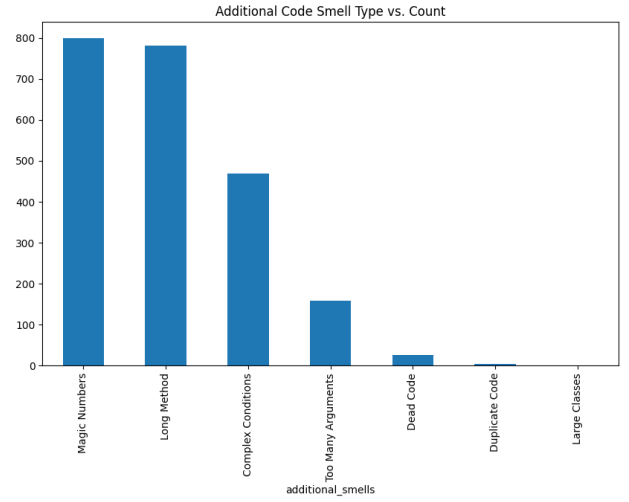


Fig. 2. Bar chart showing the distribution of additional code smell types.

### TABLE I
### MODEL PERFORMANCE METRICS

| Metric | Precision | Recall | F1-Score |
|---|---|---|---|
| Class 0 | 0.97 | 0.99 | 0.98 |
| Class 1 | 0.93 | 0.84 | 0.88 |
| **Accuracy** | | 0.966 | |

## A. Execution Screenshots

The execution of the Python scripts was captured as follows:



Fig. 3. Model output showing accuracy and performance metrics.



## VI. DISCUSSION

Our model demonstrates the potential of machine learning in automating the detection of code smells in Python. The

Random Forest Classifier, combined with TF-IDF vectorization, effectively identifies structural issues in code. However, the model's performance may vary based on the type of code smell and code block complexity. Future improvements could involve experimenting with deep learning techniques such as transformers or CodeBERT embeddings to capture semantic meaning more effectively.

## VII. FUTURE ENHANCEMENTS

While our model shows promising results, there are several areas for potential improvement:

- **Advanced Embeddings**: Future work could explore deep learning-based embeddings, such as CodeBERT, to capture semantic context beyond the word frequency captured by TF-IDF.
- **Multi-language Support**: Extending the dataset to support multiple programming languages would increase the model's applicability.

- **Integration with CI/CD**: Incorporating this model into CI/CD pipelines would enable real-time code quality checks during development.

## VIII. CONCLUSION

This paper presents a machine learning approach for detecting code smells in Python code using the Code4ML dataset. The model leverages TF-IDF vectorization and Random Forest classification to identify potential structural issues in code, achieving a high level of accuracy. The results indicate that this approach is viable for automated code review and quality assurance processes.

## ACKNOWLEDGMENT

## REFERENCES

[1] M. Fowler, K. Beck, J. Brant, W. Opdyke, and D. Roberts, *Refactoring: Improving the Design of Existing Code*. Boston, MA, USA: Addison-Wesley, 1999.

[2] S. Bavota, R. Oliveto, M. Lanza, and D. Poshyvanyk, "Code smell detection: Towards a machine learning-based approach," in *Proceedings of the 27th IEEE International Conference on Software Maintenance (ICSM)*, 2011, pp. 440–443.

[3] J. A. Camargo Cruz and W. F. Ericson, "Machine learning for code smell detection: An empirical study," in *IEEE Transactions on Software Engineering*, vol. 47, no. 3, pp. 754-767, 2021.

[4] Code4ML Dataset is available at: https://zenodo.org/records/13918465