

Advanced Detection of Source Code Clones via an Ensemble of Unsupervised Similarity Measures

Jorge Martinez-Gil

*Software Competence Center Hagenberg GmbH
Softwarepark 32a, 4232 Hagenberg, Austria
jorge.martinez-gil@scch.at*

Abstract

The capability of accurately determining code similarity is crucial in many tasks related to software development. For example, it might be essential to identify code duplicates for performing software maintenance. This research introduces a novel ensemble learning approach for code similarity assessment, combining the strengths of multiple unsupervised similarity measures. The key idea is that the strengths of a diverse set of similarity measures can complement each other and mitigate individual weaknesses, leading to improved performance. Preliminary results show that while Transformers-based CodeBERT and its variant GraphCodeBERT are undoubtedly the best option in the presence of abundant training data, in the case of specific small datasets (up to 500 samples), our ensemble achieves similar results, without prejudice to the interpretability of the resulting solution, and with a much lower associated carbon footprint due to training. The source code of this novel approach can be downloaded from <https://github.com/jorge-martinez-gil/ensemble-codesim>.

Keywords: Knowledge Engineering, Similarity Measures, Code Similarity

1. Introduction

Automatically assessing the degree of similarity between source code fragments presents several challenges ranging from syntactic variability or the use of complex algorithmic structures to the employment of diverse abstraction levels across programming languages [18]. While established techniques and tools exist to address this problem, significant limitations must be addressed. For example, the best-performing approaches need vast training data, their insights are practically intelligible to the human operator, and their associated carbon footprint during the training phase does not bring positive societal outcomes [30].

To overcome these issues, this research proposes a novel ensemble approach for improved code similarity assessment. The method combines the outputs of multiple similarity measures

to produce better similarity scores. This approach is based on the principle that the strengths of diverse similarity measures can complement each other and mitigate individual weaknesses, leading to improved systemic performance.

Extensive empirical evaluations on several datasets have demonstrated the effectiveness of our ensemble learning method. Our method outperformed traditional single-model approaches in all cases (both in scenarios with small and large datasets). Furthermore, our method achieved comparable performance to state-of-the-art methods when solving small datasets without reliance on deep learning techniques, making it a solution to consider under circumstances such as the absence of significant amounts of training data, the need for interpretability and the wish for a reduced carbon footprint when training.

Therefore, our contribution is not only to guide the aggregation of appropriate unsupervised similarity measures for code similarity but also to pave the way for applications such as source clone detection [35] or code plagiarism assessment [31] in realistic scenarios where there is not too much data for training. In this way, the major contributions of this study are:

- We build upon previous work [28] to curate a comprehensive collection of unsupervised similarity measures to compare the likeness of different code fragments. This collection facilitates a systematic exploration of the most effective measures for assessing source code similarity at low computational cost.
- Furthermore, we investigate the learning of ensembles through bagging and boosting techniques. We aim to develop models with enhanced code similarity assessment capabilities using the abovementioned similarity measures.
- Empirical evaluation of the proposed ensembles focuses on key metrics such as precision, recall, and f-measure. Our findings suggest that ensembles outperform individual similarity measures and have a strong potential to rival even the most sophisticated state-of-the-art techniques in the absence of significant training data.

The remainder of this paper is structured as follows: Section 2 provides the background on the code similarity assessment challenge and establishes its importance in the software development field. Section 3 explains the unsupervised similarity measures to address this challenge and how we can rely on them to build ensembles with superior performance. Section 4 empirically evaluates the ensembles introduced before using popular benchmark datasets. Section 5 offers a detailed analysis and discussion of the experimental results. Finally, the paper concludes with lessons learned and outlines future lines of work.

2. State-of-the-art

The existence of many programming languages and individual coding styles poses a challenge for determining the degree of likeness between source code fragments [1]. Despite that, unsupervised similarity measures demonstrate reasonable adaptability when working with code written in multiple languages or according to various styles, so it can be accepted that they are good indicators of code similarity [28]. This feature is important in many scenarios, such as maintenance, whereby increased complexity is commonplace [15].

The research community widely acknowledges that the emerging strategies based on language models offer superior performance in this context [19]. However, our hypothesis is that unsupervised methods for code similarity evaluation offer several advantages for addressing key aspects that should not be overlooked:

- **Code understanding:** Capturing the meaning of code is a complex task, as functionally equivalent code fragments may exhibit superficial differences or, on the contrary, share surface-level similarities while having entirely different purposes. Some unsupervised techniques specifically focus on the semantic intent rather than syntactic details.
- **Independence from labeled data:** Since unsupervised techniques eliminate the need for labeled training data, they reduce the overhead of establishing a ground truth. While labeled examples remain valuable for validating the effectiveness of these methods, these examples are not a prerequisite for deployment.
- **Noise filtering:** Codebases often contain comments and other non-functional elements embedded within the core logic. Specific unsupervised methods can distinguish between meaningful code patterns and just incidental material.
- **Stylistic diversity:** Unsupervised methods can accommodate the diversity of programming languages and coding styles. This mitigates the need for a complex strategy that attempts to account for all possible variations.

We acknowledge the superior capabilities of language models when confronted with this challenge. However, we also recognize and aim to demonstrate that alternative models, which excel in specific areas, could be employed under certain plausible conditions across a wide range of everyday scenarios. For instance, these models may be more suitable in situations where only limited datasets are available or where the working logic of the proposed solution needs to be comprehended.

2.1. Unsupervised Code Similarity

There are various techniques to assess the similarity between source code fragments, each with its unique perspective based on certain features of the code fragments being compared. The

following is a list that we have already proposed previously [28], consisting of twenty-one distinct unsupervised approaches. Table 1 presents all these similarity measures.

Table 1: Comparison of Different Unsupervised Similarity Approaches

Approach	Description	Ref.
Abstract Syntax Trees (ASTs) Similarity	Examines the hierarchical structures of code to identify structural similarities between different ASTs.	[22]
Bag-of-Words Similarity	Counts the occurrence of words in texts, comparing these frequencies while ignoring word order and structural context.	[7]
Code Embeddings Similarity	Uses vector representations of code for comparison, evaluating similarity based on these embeddings. This technique is used without recalibration (zero-shot) here.	[2]
Comments Similarity	Focuses on comparing code comments to facilitate documentation and understanding, using traditional text similarity methods.	[29]
Fuzzy Matching Similarity	Adapts string comparison techniques for minor syntactical variations to code similarity assessment.	[37]
Function Calls Similarity	Compares source code fragments based on their functions names.	[42]
Graph-based Similarity	Relies on the relationships in graphs, representing data structures or code dependencies, to determine similarity.	[44]
Jaccard Similarity	Uses set intersection and union to evaluate similarity in text analysis and information retrieval contexts.	[13]
Levenshtein Similarity	Measures the minimum number of edits needed to transform one string into another, reflecting similarity.	[25]
Longest Common Subsequence (LCS) Similarity	Identifies the longest sequence of items that two series share, indicating similarity.	[5]
Metrics Similarity	Involves computing various code-related metrics and comparing these values to estimate similarity.	[32]
N-grams Similarity	Compares texts based on the occurrence of contiguous sequences of n items, evaluating shared n-grams.	[8]
Output Analysis Similarity	Assesses the similarity of program outputs by applying traditional text similarity metrics.	[29]

Continued on next page

Table 1 continued from previous page

Approach	Description	Ref.
Perceptual Hashing Similarity	Adapts techniques from image similarity, using hash codes from visual representations of code to measure similarity.	[33]
Program Dependence Graph Similarity	Assesses the similarity between codes by analyzing their program dependence graphs, which represent the dependencies between their elements.	[24]
Rolling Hash Similarity	Employs a dynamically updating hash to compare substrings in large texts, adapted for code comparison.	[14]
RKR-GST Similarity	Targets code similarity detection by identifying maximal sequences of contiguous matching tokens in texts.	[41]
Semdiff Similarity	Detects semantic differences between program versions, evaluating the impact of code changes on semantics.	[16]
Semantic Clone Similarity	Assesses code fragment similarity based on the semantic meanings of program element names.	[11]
TF-IDF Similarity	Common in text analysis, compares texts based on the relative importance of terms, using weighted term frequencies.	[20]
Winnow Similarity	Identifies similar texts by hashing and comparing text fingerprints.	[36]

According to existing research [28], it is unrealistic to think that these techniques alone can deliver performance equivalent to the state-of-the-art. However, we aim to show how their aggregation can do so under certain conditions.

2.2. Supervised Code Similarity

Supervised code similarity is a well-studied field and has yielded the best results to date [38, 39, 40, 45, 43]. Although among the existing approaches, CodeBERT [10] and its variant GraphCodeBERT [12] stand out. CodeBERT is a variant of BERT [9], and empirical results show that it is the leading technology to understand source code expressed in different general-purpose programming languages. It excels at tasks like generating documentation, creating code, and comparing similar pieces of code. The idea behind CodeBERT is to combine techniques from language processing with an understanding of code and training on a wide range of programming languages and their associated comments. Results show CodeBERT is very good at detecting similarities between source code fragments, as demonstrated by its performance in the resolution of benchmark datasets for clone detection [10]. However, an additional fine-tuning phase on the pre-trained model over the actual data is usually required to obtain optimal results.

2.3. Contribution over the state-of-the-art

This research introduces a novel ensemble learning approach for code similarity assessment. The rationale behind combining multiple unsupervised similarity measures is to leverage the strengths of individual approaches while mitigating their weaknesses. This strategy is helpful in scenarios where deep learning is not optimal due to the absence of significant training data, when prioritizing model interpretability without significantly compromising performance, or even when seeking to minimize the environmental impact.

Our work’s focus is eminently practical. In fact, we evaluate our ensembles on benchmark datasets to confirm the approach’s superiority over traditional techniques. Furthermore, we offer our code as open-source since this availability might promote transparency and collaborative advancement.

3. Ensemble Approach

We propose using an ensemble of unsupervised similarity measures to improve the performance when assessing code similarity. The rationale behind combining several unsupervised similarity measures is to offer a more accurate evaluation than strategies using just one measure [27]. Therefore, we first define the problem, introduce the concept of an ensemble, and discuss bagging and boosting as the two most prominent techniques for building ensembles. These techniques for building ensembles help aggregate different unsupervised similarity measures, leading usually to an enhanced strategy.

3.1. Problem Statement

We can define code similarity as follows: Given a set of code fragments $S = \{C_1, C_2, \dots, C_n\}$, we aim to find a function $f : S \times S \rightarrow [0, 1]$ that measures the similarity between any two fragments C_i and C_j .

The function f should map a pair (C_i, C_j) to a continuous value within the real interval $[0, 1]$:

1. $f(C_i, C_j) = 0$ indicates absolute code dissimilarity.
2. $f(C_i, C_j) = 1$ indicates equivalent code.
3. The value of $f(C_i, C_j)$ increases with increasing similarity between C_i and C_j .

Some applications can be built from this definition [34]. For example, code clone detection can be performed by applying a threshold value to differentiate between clones and non-clones [23].

$$\text{Clone}(C_i, C_j) = \begin{cases} \text{True}, & \text{if } f(C_i, C_j) > \theta \\ \text{False}, & \text{otherwise} \end{cases}$$

In addition, something similar can also be done to determine whether or not the code was plagiarized in an academic environment [3], for example. These two challenges of a very similar nature will be studied in more detail later in the paper.

3.2. Ensemble Definition

Let $S = \{s_1, s_2, \dots, s_n\}$ be a set of code fragments and $M = \{m_1, m_2, \dots, m_{21}\}$ be a set of unsupervised similarity measures. An ensemble E aims to aggregate similarity measures such as:

$$E(s_a, s_b) = \alpha_1 \cdot m_1(s_a, s_b) + \dots + \alpha_{21} \cdot m_{21}(s_a, s_b)$$

where α_i represents the weight of each similarity measure, and $\sum_{i=1}^{21} \alpha_i = 1$.

Each $m_i(s_a, s_b)$ returns a similarity score in the real range $[0, 1]$, where 0 indicates no similarity, and 1 indicates equivalent code fragments. The ensemble E , therefore, also returns a score in the real range $[0, 1]$, effectively aggregating the individual similarity assessments into a measure.

There are several techniques for learning ensembles, almost all based on a purely data-driven strategy. In the context of this work, and for practical reasons, we focus on two: bagging and boosting. Other techniques, such as stacking, will not be considered here because they have already been addressed in similar problems in the past, such as [26]. Although an in-depth study of each of these techniques is beyond the scope of this study, we offer below a brief description of how these ensemble-building strategies operate.

3.3. Bagging

The idea behind bagging is to train multiple instances of the same model with different data subsets and aggregate their predictions [6]. This strategy takes a list of code fragments, the number of subsets to create, and unsupervised similarity measures as input. It then averages the results of applying all measures to each subset.

3.4. Boosting

The idea behind boosting is to train multiple models, with each model instance focuses on the errors of the previous one [4]. It is necessary to operate by sequentially applying a list of code fragments, the number of rounds to run, and a list of unsupervised measures as input. This sequential application allows for iterative refinement of the code similarity scores, with the weights being updated based on performance, leading to the final scores.

3.5. Differences with the Transformer architecture

Using an ensemble of unsupervised semantic similarity measures for code similarity assessment offers several advantages over relying solely on a model like CodeBERT [9], or its more powerful variant GraphCodeBERT [12]. The ensemble approach combines different unsupervised similarity

measures, each capturing unique aspects of code, thus providing a diverse analysis that can reduce individual biases and improve performance.

Ensembles are also more flexible and scalable and do not require massive labeled data, making them suitable for diverse programming languages and code domains. Therefore, our approach benefits from the strengths of multiple unsupervised similarity measures, potentially outperforming all of them. Transformers-inspired solutions can face challenges with interpretability due to their complex attention mechanisms and might struggle with out-of-distribution or noisy code fragments. However, if sufficient data is available and there are no requirements beyond performance, results from transformer architectures should be much superior.

4. Empirical Evaluation

In previous work [28], we show that unsupervised similarity measures exhibit strengths in different domains. Some excel in textual comparisons, making them ideal for identifying cloned text. Others prioritize the detection of functional similarity, while specific measures specialize in uncovering structural resemblances between code and textual elements. The optimal choice of measure ultimately depends on the composition of the benchmark dataset. However, it is not reasonable to ask a human expert to do that since it would require many hours of work, and the result would be prone to errors since the task is far from trivial. The good news is that bagging and boosting algorithms, with their data-driven strategies, will decide the inclusion and importance of each of the simple unsupervised similarity measures.

4.1. Baseline

In this study, we compare our ensembles using two strategies. First, a weak baseline consists of obtaining better results than each similarity measure considered individually. Secondly, a strong baseline that consists of the state-of-the-art conformed by CodeBERT [9], and its variant GraphCodeBERT [12].

In this regard, the state-of-the-art models can have a previous fine-tuning phase that begins with loading and randomly splitting a dataset of code fragments into training, validation, and test sets. Then, the approach’s core involves training the model to discern clone pairs, guided by a trainer configured with specific training arguments like epoch count, batch size, and learning rate adjustments. Finally, performance metrics are calculated to evaluate the model’s effectiveness as the average value of a given number of executions. In our study, we considered up to 10 independent executions.

4.2. *Small-scale Dataset*

Firstly, we use the IR-Plag dataset¹ [21]. This dataset was designed for benchmarking code similarity techniques. It contains code files deliberately crafted to simulate patterns of academic plagiarism. Although primarily intended for plagiarism detection, its characteristics align with our goals due to the overlap between plagiarized and cloned code, which involves replication (though the intent may differ). The IR-Plag dataset encompasses a variety of complexities without strictly classifying the clone types.

The dataset consists of seven original code files. A substantial portion, 355 files (77%), are labeled as plagiarized, indicating a high prevalence of duplication. 105 files are considered non-plagiarized, potentially representing modified or derivative works. This yields a total of 467 code files. Collectively, these files contain 59,201 tokens, with 540 unique tokens underlining the dataset’s lexical diversity. File sizes exhibit variation, ranging from 286 tokens (largest) to 40 tokens (smallest), with an average file size of approximately 126 tokens. Therefore, although small in size, it is a dataset with a high compositional diversity of programming elements.

4.3. *Large-scale Dataset*

Secondly, we use the BigCloneBench dataset, which is a benchmark collection created to evaluate how well different code clone detection strategies perform². It is helpful to tackle the challenge of finding duplicate code. It includes a variety of programming languages, and its data comes from actual software and student projects. This mix makes it an interesting resource for evaluating clone detection in real-life situations. This dataset is usually used in training solutions to make code more maintainable.

The dataset is divided into three subsets: training, validation, and testing. The training set contains 901,028 items, while the validation and testing sets contain 415,416 items each. This distribution is outlined to ensure comprehensive training, validation, and testing of the model. This helps compare different strategies and determine which ones are best at finding code clones.

4.4. *Evaluation Criteria*

In principle, measuring the accuracy of the different techniques is appropriate. That is the percentage of occasions in which the strategy under study succeeded or failed in predicting the similarity of code fragments. However, this metric is strongly discouraged for unbalanced datasets because the mere fact of constantly repeating the label of the most numerous classes in the test set will yield a high value that may be misleading.

For that reason, the community prefers other metrics to assess the precision and coverage of a clone detection strategy, such as:

¹<https://github.com/oscar-karnalim/sourcecodeplagiarismdataset>

²<https://github.com/microsoft/CodeXGLUE/tree/main/Code-Code/Clone-detection-BigCloneBench>

- Precision: This measures the accuracy of an approach by determining the fraction of true positive identifications among all detected code fragments. Higher precision indicates better reliability. It is defined as:

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

- Recall: It evaluates the completeness of an approach by measuring the fraction of true positive clones in the dataset that were correctly identified. Higher recall implies that fewer actual clones were missed.

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

- F-measure: This metric assesses the performance of an approach by harmonizing precision and recall, ensuring a balanced evaluation of both metrics.

$$F\text{-measure} = 2 \cdot \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

This evaluation method is popular because it emphasizes the importance of positive classes by separately assessing false positives (precision) and false negatives (recall). It imposes penalties on the model for not identifying positive cases and for incorrect positive predictions. Subsequently, one can compute a harmonic mean of these two metrics, known as the F-measure, which serves as a basis for ranking different techniques.

4.5. Hardware requirements

Please note that all the computations of similarity measures and the learning of their ensembles has been done on a CPU 11th Gen Intel(R) Core(TM) i7-1185G7 a 3.00GHz with 32GB for RAM (TDP of 50W approx.), while the configuration and execution of the baselines with CodeBERT and GraphCodeBERT has required a GPU Tesla V100-PCIE-16GB (TDP of 300W approx.).

4.6. Results

Table 2 shows us the results obtained for the first benchmark dataset, i.e., the small-scale IR-Plag. It can be seen that some unsupervised similarity measures are pretty good on their own if one looks at the f-measure. However, it is striking that ensembles obtain the best results, so we far exceed our weak baseline. Please note that we do not study time consumption issues here as they have been studied previously [28].

Approach	Precision	Recall	F-Measure
Abstract Syntax Tree	0.95	0.30	0.45
Bag-of-Words	0.93	0.19	0.31
CodeBERT	0.91	0.03	0.05
Comment Sim.	0.88	0.40	0.55
Output Analysis	0.88	0.93	0.90
Function Calls	0.85	0.67	0.75
Fuzzy Matching	0.82	0.50	0.62
Graph Matching	0.82	0.91	0.86
Rolling Hash	0.94	0.20	0.33
Perceptual Hash	0.89	0.41	0.57
Jaccard	0.99	0.26	0.42
Longest Common Subsequence	0.78	0.97	0.87
Levenshtein	0.82	0.76	0.79
Metrics comparison	0.81	0.57	0.67
N-Grams	0.84	0.83	0.84
Program Dependence Graph	0.85	0.39	0.53
Rabin-Karp	0.87	0.47	0.61
Semantic Clone	0.94	0.24	0.38
Semdiff method	0.87	0.25	0.39
TDF-IDF	0.97	0.17	0.29
Winnow	0.89	0.65	0.75
Boosting	0.88	0.99	0.93
Bagging	0.95	0.97	0.96

Table 2: Summary of the results obtained for the IR-Plag dataset

Table 3 establishes a comparative for the state-of-the-art published to this dataset. On the one hand, it can be seen that the ensembles can perform better than CodeBERT. The reason is that they perform better without large data for training. This is even though several combinations of different partitions have been evaluated: 60-20-20, 70-15-15, and 80-10-10 on training, validation and test respectively. We are reporting here the best average obtained. On the other hand, GraphCodeBERT matches our best results, although with the disadvantages already discussed throughout this study, so it might make sense to use our solution when working with small datasets.

Approach	Precision	Recall	F-Measure
CodeBERT	0.72	1.00	0.84
Output Analysis	0.88	0.93	0.90
Boosting	0.88	0.99	0.93
GraphCodeBERT	0.98	0.95	0.96
Bagging	0.95	0.97	0.96

Table 3: State-of-the-art for the IR-Plag dataset. Due to the small size of the dataset, several datasets partitions have been tested for the CodeBERT and GraphCodeBERT fine-tuning process

Table 4 shows us the results of the different unsupervised similarity measures over the BigCloneBench dataset. These results estimate over 25,000 samples (6% of the dataset) because we do not have enough computational power to analyze the complete dataset. We also discard similarity calculations in scenarios that require more than 20 seconds of computation, otherwise the process may require many days of computation. Furthermore, similarity by output analysis cannot be performed in this context due to the high computational requirements of packaging, compiling, executing and comparing several million code samples. Anyway, if desired and if sufficient computational power is available, our solution can be used on the total set.

Approach	Precision	Recall	F-Measure
Abstract Syntax Tree	0.15	0.47	0.23
Bag-of-Words	0.27	0.62	0.38
CodeBERT	0.18	0.29	0.23
Comment Sim.	0.20	0.31	0.25
Output Analysis	0.00	0.00	0.00
Function Calls	0.40	0.62	0.48
Fuzzy Matching	0.38	0.43	0.40
Graph Matching	0.16	0.46	0.23
Rolling Hash	0.21	0.23	0.22
Perceptual Hash	0.15	0.78	0.24
Jaccard	0.27	0.62	0.38
Longest Common Subsequence	0.22	0.47	0.30
Levenshtein	0.33	0.23	0.27
Metrics comparison	0.15	0.72	0.25
N-Grams	0.32	0.45	0.37
Program Dependence Graph	0.18	0.44	0.26
Rabin-Karp	0.15	0.22	0.18
Semantic Clone	0.23	0.43	0.30
Semdiff method	0.18	0.26	0.22
TDF-IDF	0.27	0.63	0.38
Winnow	0.20	0.71	0.31
Bagging	0.85	0.42	0.56
Boosting	0.79	0.45	0.57

Table 4: Summary of the results obtained for the BigCloneBench dataset

As can be seen, each technique individually shows a performance far below what would be required in a real environment. Therefore, there are better candidates for integration into real-world systems. However, their systematic aggregation through bagging or boosting techniques yields superior results, so we far exceeded our weak baseline.

Table 5 establishes a comparative for the state-of-the-art published to this dataset. For the fine-tuning of CodeBERT and GraphCodeBERT, we rely on the configuration outlined by the authors of the original papers. At the same time, the numbers relative to our strategies are

estimated. The reason is that the bagging and boosting strategies have been tested using the seven most promising unsupervised similarity measures from the 415,416 samples of the test partition. These promising measures have been obtained from a preliminary study, please refer to the source code for more information. The reason is, again, not having enough computational resources to operate over the whole search space. Considering the twenty-one approaches should yield better results, albeit only slightly.

Approach	Precision	Recall	F-Measure
Deckard [17]	0.93	0.02	0.03
RtvNN [40]	0.95	0.01	0.01
Bagging	0.85	0.42	0.56
Boosting	0.79	0.45	0.57
CDLH [39]	0.92	0.74	0.82
ASTNN [45]	0.92	0.94	0.93
CodeBERT [10]	0.95	0.93	0.94
FA-AST-GMN [38]	0.96	0.94	0.95
TBBCD [43]	0.94	0.96	0.95
GraphCodeBERT[12]	0.95	0.95	0.95

Table 5: State-of-the-art for the BigCloneBench dataset

5. Discussion

Our experiments demonstrate the effectiveness of ensemble learning for assessing code similarity. Although the results are questionable when the training data are abundant, equivalent performance is observed when working with small datasets (467 samples in the case we have studied).

In any case, integrating diverse unsupervised similarity measures seems effective since our approach surpasses the performance of single-measure strategies regarding code clone detection. Our ensemble approach mitigates the inherent limitations of individual similarity measures. Aggregating results from multiple measures compensate for the shortcomings of any single similarity measure. For instance, a structure-focused measure may overlook semantically equivalent code with superficial syntactic differences.

Furthermore, the adaptable nature of our ensemble approach offers other significant benefits. We can incorporate a variety of similarity measures, enabling us to tailor the approach for specific applications. If plagiarism detection is the primary objective, structure-focused measures could be given higher weight. Semantic-oriented measures might be prioritized for code reuse identification, etc. In general, our research yielded some lessons learned, detailed below.

5.1. Lessons Learned

Our experiments show that some unsupervised code similarity methods work reasonably well for finding source code clones. Therefore, it could make sense to integrate these methods to help improve software development under certain circumstances. For example,

1. Our ensembles can find duplicate code to rewrite as reusable parts at a low cost. This is important since code reuse is vital to many software-related tasks.
2. Real-world code can be unstructured. Our ensembles of unsupervised methods can manage this noise and variation, particularly relevant in scenarios with little labeled data.
3. A human operator cannot easily interpret the Transformers-like architectures on which CodeBERT is based. Moreover, they have a huge carbon footprint when trained [30]. Our approach may be reasonable in situations where these two aspects are critical.

However, there are some limitations and areas for improvement where the performance of our ensembles is not optimal, which are discussed below.

5.2. Limitations

Our proposed approach has yielded promising results; however, some aspects still need improvement. First, when dealing with abundant training data, the performance of our ensembles falls substantially below the benchmarks set by state-of-the-art models.

Second, scalability presents a challenge; as codebases expand in size and complexity, the computational burden of implementing multiple similarity measures increases considerably. This escalation demands integrating more efficient optimization techniques to manage and mitigate the increased processing requirements.

Lastly, while our ensemble method offers theoretical interpretability, its practical use demands a more pragmatic approach. Further work in this direction is crucial to demonstrate how the model identifies similar code fragments, thereby enhancing transparency and usability in practice.

6. Conclusions and Future Work

Evaluating code similarity is a crucial component of many software development tasks. In this study, we have explored unsupervised similarity measures, leveraging their independence from labeled data and adaptability to diverse coding styles. Our findings show the efficacy of ensemble learning approaches in code similarity assessment. We have developed a strategy that surpasses traditional techniques by aggregating multiple measures, achieving performance comparable to state-of-the-art standards in the absence of significant training data.

As code becomes more complex, there is a growing need for accurate and efficient ways to measure code similarity. The software industry’s ever-changing nature shows us the vital role that unsupervised similarity measures could play in this context. While supervised approaches are superior in performance, unsupervised techniques often offer advantages of another kind.

Future work must further explore the potential of unsupervised code clone detection methods. Future research could also evaluate the efficacy of other ensemble techniques and examine the applicability of transfer learning for improved performance. The ultimate goal is to advance code analysis automation. This development could significantly reduce the need for manual intervention across all phases of the software development process, thereby increasing efficiency and productivity in the field.

Acknowledgments

The research reported in this paper has been funded by the Federal Ministry for Climate Action, Environment, Energy, Mobility, Innovation, and Technology (BMK), the Federal Ministry for Digital and Economic Affairs (BMDW), and the State of Upper Austria in the frame of SCCH, a center in the COMET - Competence Centers for Excellent Technologies Programme managed by Austrian Research Promotion Agency FFG.

References

- [1] Ain, Q. U., Butt, W. H., Anwar, M. W., Azam, F., & Maqbool, B. (2019). A systematic review on code clone detection. *IEEE access*, 7, 86121–86144.
- [2] Alon, U., Zilberstein, M., Levy, O., & Yahav, E. (2019). code2vec: Learning distributed representations of code. *Proceedings of the ACM on Programming Languages*, 3, 1–29.
- [3] Aniceto, R. C., Holanda, M., Castanho, C., & Da Silva, D. (2021). Source code plagiarism detection in an educational context: A literature mapping. In *2021 IEEE Frontiers in Education Conference (FIE)* (pp. 1–9). IEEE.
- [4] Bentéjac, C., Csörgő, A., & Martínez-Muñoz, G. (2021). A comparative analysis of gradient boosting algorithms. *Artificial Intelligence Review*, 54, 1937–1967.
- [5] Bergroth, L., Hakonen, H., & Raita, T. (2000). A survey of longest common subsequence algorithms. In *Proceedings Seventh International Symposium on String Processing and Information Retrieval. SPIRE 2000* (pp. 39–48). IEEE.
- [6] Breiman, L. (1996). Bagging predictors. *Machine learning*, 24, 123–140.

- [7] Corley, C. D., & Mihalcea, R. (2005). Measuring the semantic similarity of texts. In *Proceedings of the ACL workshop on empirical modeling of semantic equivalence and entailment* (pp. 13–18).
- [8] Damashek, M. (1995). Gauging similarity with n-grams: Language-independent categorization of text. *Science*, 267, 843–848.
- [9] Devlin, J., Chang, M., Lee, K., & Toutanova, K. (2019). BERT: pre-training of deep bidirectional transformers for language understanding. In J. Burstein, C. Doran, & T. Solorio (Eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)* (pp. 4171–4186). Association for Computational Linguistics. doi:10.18653/v1/n19-1423.
- [10] Feng, Z., Guo, D., Tang, D., Duan, N., Feng, X., Gong, M., Shou, L., Qin, B., Liu, T., Jiang, D., & Zhou, M. (2020). Codebert: A pre-trained model for programming and natural languages. In T. Cohn, Y. He, & Y. Liu (Eds.), *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020* (pp. 1536–1547). Association for Computational Linguistics volume EMNLP 2020 of *Findings of ACL*. URL: <https://doi.org/10.18653/v1/2020.findings-emnlp.139>. doi:10.18653/V1/2020.FINDINGS-EMNLP.139.
- [11] Gabel, M., Jiang, L., & Su, Z. (2008). Scalable detection of semantic clones. In *Proceedings of the 30th international conference on Software engineering* (pp. 321–330).
- [12] Guo, D., Ren, S., Lu, S., Feng, Z., Tang, D., Liu, S., Zhou, L., Duan, N., Svyatkovskiy, A., Fu, S. et al. (2020). Graphcodebert: Pre-training code representations with data flow. *arXiv preprint arXiv:2009.08366*, .
- [13] Haque, S., Eberhart, Z., Bansal, A., & McMillan, C. (2022). Semantic similarity metrics for evaluating source code summarization. In *Proceedings of the 30th IEEE/ACM International Conference on Program Comprehension* (pp. 36–47).
- [14] Hartanto, A. D., Syaputra, A., & Pristyanto, Y. (2019). Best parameter selection of rabin-karp algorithm in detecting document similarity. In *2019 International Conference on Information and Communications Technology (ICOIACT)* (pp. 457–461). IEEE.
- [15] Higo, Y., Ueda, Y., Kamiya, T., Kusumoto, S., & Inoue, K. (2002). On software maintenance process improvement based on code clone analysis. In *Product Focused Software Process Improvement: 4th International Conference, PROFES 2002 Rovaniemi, Finland, December 9–11, 2002 Proceedings 4* (pp. 185–197). Springer.

- [16] Horwitz, S. (1990). Identifying the semantic and textual differences between two versions of a program. In *Proceedings of the ACM SIGPLAN 1990 conference on Programming language design and implementation* (pp. 234–245).
- [17] Jiang, L., Mishergbi, G., Su, Z., & Glondu, S. (2007). Deckard: Scalable and accurate tree-based detection of code clones. In *29th International Conference on Software Engineering (ICSE'07)* (pp. 96–105). IEEE.
- [18] Juergens, E., Deissenboeck, F., Hummel, B., & Wagner, S. (2009). Do code clones matter? In *2009 IEEE 31st International Conference on Software Engineering* (pp. 485–495). IEEE.
- [19] Karmakar, A., & Robbes, R. (2021). What do pre-trained code models know about code? In *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)* (pp. 1332–1336). IEEE.
- [20] Karnalim, O. (2020). Tf-idf inspired detection for cross-language source code plagiarism and collusion. *Computer Science*, 21.
- [21] Karnalim, O., Budi, S., Toba, H., & Joy, M. (2019). Source code plagiarism detection in academia with information retrieval: Dataset and the observation. *Informatics in Education*, 18, 321–344.
- [22] Karnalim, O., & Simon (2020). Syntax trees and information retrieval to improve code similarity detection. In *Proceedings of the Twenty-Second Australasian Computing Education Conference* (pp. 48–55).
- [23] Karnalim, O. et al. (2021). Explanation in code similarity investigation. *IEEE Access*, 9, 59935–59948.
- [24] Krinke, J. (2001). Identifying similar code with program dependence graphs. In *Proceedings Eighth Working Conference on Reverse Engineering* (pp. 301–309). IEEE.
- [25] Levenshtein, V. I. (1966). Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet physics doklady* (pp. 707–710). volume 10.
- [26] Martinez-Gil, J. (2022). A comprehensive review of stacking methods for semantic similarity measurement. *Machine Learning with Applications*, 10, 100423. doi:10.1016/j.mlwa.2022.100423.
- [27] Martinez-Gil, J. (2023). A comparative study of ensemble techniques based on genetic programming: A case study in semantic similarity assessment. *Int. J. Softw. Eng. Knowl. Eng.*, 33, 289–312. doi:10.1142/S0218194022500772.

- [28] Martinez-Gil, J. (2024). Source code clone detection using unsupervised similarity measures. In P. Bludau, R. Ramler, D. Winkler, & J. Bergsmann (Eds.), *Software Quality as a Foundation for Security - 16th International Conference on Software Quality, SWQD 2024, Vienna, Austria, April 23-25, 2024, Proceedings* (pp. 21–37). Springer volume 505 of *Lecture Notes in Business Information Processing*. URL: https://doi.org/10.1007/978-3-031-56281-5_2. doi:10.1007/978-3-031-56281-5_2.
- [29] Martinez-Gil, J., & Chaves-Gonzalez, J. M. (2020). A novel method based on symbolic regression for interpretable semantic similarity measurement. *Expert Syst. Appl.*, *160*, 113663. doi:10.1016/j.eswa.2020.113663.
- [30] Martinez-Gil, J., & Chaves-Gonzalez, J. M. (2022). Sustainable semantic similarity assessment. *J. Intell. Fuzzy Syst.*, *43*, 6163–6174. URL: <https://doi.org/10.3233/JIFS-220137>. doi:10.3233/JIFS-220137.
- [31] Novak, M., Joy, M., & Kermek, D. (2019). Source-code similarity detection and detection tools used in academia: a systematic review. *ACM Transactions on Computing Education (TOCE)*, *19*, 1–37.
- [32] Nuñez-Varela, A. S., Pérez-Gonzalez, H. G., Martínez-Perez, F. E., & Soubervielle-Montalvo, C. (2017). Source code metrics: A systematic mapping study. *Journal of Systems and Software*, *128*, 164–197.
- [33] Ragkhitwetsagul, C., Krinke, J., & Marnette, B. (2018). A picture is worth a thousand words: Code clone detection based on image similarity. In *12th IEEE International Workshop on Software Clones, IWSC 2018, Campobasso, Italy, March 20, 2018* (pp. 44–50). IEEE Computer Society. URL: <https://doi.org/10.1109/IWSC.2018.8327318>. doi:10.1109/IWSC.2018.8327318.
- [34] Roy, C. K., Cordy, J. R., & Koschke, R. (2009). Comparison and evaluation of code clone detection techniques and tools: A qualitative approach. *Science of computer programming*, *74*, 470–495.
- [35] Saini, N., Singh, S. et al. (2018). Code clones: Detection and management. *Procedia computer science*, *132*, 718–727.
- [36] Schleimer, S., Wilkerson, D. S., & Aiken, A. (2003). Winnowing: local algorithms for document fingerprinting. In *Proceedings of the 2003 ACM SIGMOD international conference on Management of data* (pp. 76–85).

- [37] Singla, N., & Garg, D. (2012). String matching algorithms and their applicability in various applications. *International journal of soft computing and engineering*, 1, 218–222.
- [38] Wang, W., Li, G., Ma, B., Xia, X., & Jin, Z. (2020). Detecting code clones with graph neural network and flow-augmented abstract syntax tree. In *2020 IEEE 27th International Conference on Software Analysis, Evolution and Reengineering (SANER)* (pp. 261–271). IEEE.
- [39] Wei, H., & Li, M. (2017). Supervised deep features for software functional clone detection by exploiting lexical and syntactical information in source code. In *IJCAI* (pp. 3034–3040).
- [40] White, M., Tufano, M., Vendome, C., & Poshyvanyk, D. (2016). Deep learning code fragments for code clone detection. In *Proceedings of the 31st IEEE/ACM international conference on automated software engineering* (pp. 87–98).
- [41] Wise, M. J. (1993). String similarity via greedy string tiling and running karp-rabin matching. *Online Preprint, Dec*, 119, 1–17.
- [42] Xu, M., Wu, L., Qi, S., Xu, J., Zhang, H., Ren, Y., & Zheng, N. (2013). A similarity metric method of obfuscated malware using function-call graph. *Journal of Computer Virology and Hacking Techniques*, 9, 35–47.
- [43] Yu, H., Lam, W., Chen, L., Li, G., Xie, T., & Wang, Q. (2019). Neural detection of semantic code clones via tree-based convolution. In *2019 IEEE/ACM 27th International Conference on Program Comprehension (ICPC)* (pp. 70–80). IEEE.
- [44] Zager, L. A., & Verghese, G. C. (2008). Graph similarity scoring and matching. *Applied mathematics letters*, 21, 86–94.
- [45] Zhang, J., Wang, X., Zhang, H., Sun, H., Wang, K., & Liu, X. (2019). A novel neural source code representation based on abstract syntax tree. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)* (pp. 783–794). IEEE.