

Ruprecht-Karls-Universität Heidelberg
Institut für Informatik
Lehrstuhl für Parallele und Verteilte Systeme

Masterarbeit
Automated Checking of Clean Code
Guidelines for Python

Name: Enrico Kaack
Matrikelnummer: 3534472
Betreuer: Prof. Dr. Artur Andrzejak
Datum der Abgabe: 31.10.2020

Ich versichere, dass ich diese Master-Arbeit selbstständig verfasst und nur die angegebenen Quellen und Hilfsmittel verwendet habe.

Heidelberg, 31.10.2020

Zusammenfassung

Abstract auf Deutsch

Abstract

This is the abstract.

Contents

List of Figures	1
List of Tables	2
1. Introduction	5
1.1. Objectives and Contributions	6
1.2. Structure	8
2. Background and Related Work	9
2.1. Code Quality	9
2.2. Clean Code	11
2.2.1. Naming	11
2.2.2. Functions	12
2.2.3. Comments	14
2.2.4. Data Structures and Objects	15
2.2.5. Classes	16
2.2.6. Exception Handling	17
2.2.7. Additional Rules	18
2.3. Source Code Analysis	19
2.4. Clean Code Complexity Levels	20
2.5. Quantitative Metrics for Code Quality	24
2.5.1. Cyclomatic Complexitiy	25
2.5.2. Halstead Complexity Measures	26
2.5.3. Software Maintainability Index	27
2.6. Tools for Code Quality Analysis	28
2.6.1. PyLint	29
2.6.2. PMD Source Code Analyzer Project	29
2.6.3. Codacy	31
2.6.4. Sonarqube	31

3. Approach	34
3.1. Clean Code Analysis Platform	34
3.1.1. Architecture	35
3.1.2. Return None Plugin	39
3.1.3. Improvements	45
3.2. Clean Code Classification	46
3.2.1. Challenges	48
3.2.2. Dataset	51
3.2.3. Processing	53
3.2.4. Models	60
4. Quantitative Evaluation	62
4.1. Research Questions	62
4.1.1. RQ1: What Is the Utility of the CCAP Besides Existing Tools? .	62
4.1.2. RQ2: How Do Different Models Compare on the Task of Detecting Non-Clean Code?	65
4.1.3. RQ3: Do Machine Learning-Based Models Cover a Larger Variety of Cases Than Rule-Based Checker?	75
5. Conclusion	82
5.1. Summary	82
5.2. Future Work	83
Bibliography	84
Appendices	90
A. RQ2: Model Performances	90
B. RQ3: Model Performances on Manipulated Data	95
C. Tokenizer Types	99
D. Dataset	101

List of Figures

1.1. Motivating example for clean code vs non-clean code	6
2.1. Illustration Law of Demeter	16
2.2. Control-flow graph visualisation for a code sample	25
3.1. High-level schematic overview of the CCAP.	36
3.2. Output of the HTML Output Plugin, displayed in a browser.	45
3.3. Example commit messages that underline the inconsistency in commit message naming.	49
3.4. Split of the dataset. First, 20% of all files are seperated into a holdout set for later testing. The remaining 80% of the files are split into 90% training and 10% test data.	54
3.5. Schematic representation of the pipeline to train and test models.	55
3.6. Sliding window approach on source code containing a problematic <code>return</code> <code>None</code>	58
3.7. Pipeline for code manipulation and evaluation.	58
4.1. F1 performance of the models, sorted by the baseline f1 score.	78
C.1. Possible token types of the tokenizer for our type encoding. Extracted from the <code>tokenize.tok_name</code> attribute.	100

List of Tables

2.1. Overview of our categorization of automated approaches to check for compliance with the clean code rules.	21
3.1. Class distribution for the dataset.	50
3.2. Open-source repositories we used in our dataset and their corresponding domain.	52
3.3. General and problem specific metrics for the train, test and holdout set. .	53
3.4. Class frequencies for all problem types on the train, test and holdout set. We observe a comparable class distribution among the datasets for each problem type.	53
4.1. Dataset sizes for different resampling strategies with f1 performance on the test set. The f1 performance for the speciifc resampling is the random forest classifier with 100 trees, no class weighting and type encoding. . .	68
4.2. Dataset sizes for different resampling strategies with f1 performance on the test set. The f1 performance for the speciifc resampling is the gradient boosting classifier with 200 stages, a learning rate of 0.2 and type encoding.	70
4.3. Dataset sizes for different resampling strategies with f1 performance on the test set. The f1 performance for the speciifc resampling is the LSTM with an embedding size of 32, 10 lstm cells, a training on 3 epochs with a batch size of 256.	72
4.4. Best performing classifier for RQ2 based on f1 score of the holdout set. For the SVM, the test score of the configuration with available data is used.	74
A1. random forest	91
A2. gradient boosting classifier	92
A3. lstm	93
A4. svm. NaN value for configuration that we did not test due to time constraints.	94

A1. Manipulated code for random forest	96
A2. Manipulated code gradient boosting classifier	97
A3. Manipulated code lstm	98
A1. Open-source repositories we used for our dataset.	101

1. Introduction

The rise of information technology with computers over the last decades allowed humanity to process data and solve problems in unprecedented performance. Computers require instruction in order to fulfil its purpose. The common language between humans and computers are programming languages. Software developers use programming languages to write source code that instructs computers to process data and solve problems. While the machine only follows the instructions described by the source code, the humans have to understand the purpose of the instructions.

In today's world, software becomes a vast market with a revenue of over 450 billion U.S. dollars in 2019 (TODO source <https://www.statista.com/forecasts/963597/software-revenue-in-the-world>). Creating new code often requires to read and understand preexisting code. As a consequence, it is essential to make code as readable and understandable as possible.

The objectives of the clean code principles are readability and understandability of source code. A set of rules and recommendations about writing code to be easy to read and understand. For the computer, it is irrelevant if the human can understand the code since it only executes the instructions. However, any developer who reads clean code profits from it. In Figure 1.1, the upper listing shows an implementation of an algorithm without paying attention to clean code rules. In the bottom listing, the same algorithm is implemented following the clean code rules such as having one purpose per function or method calls in if-conditions. Executing the code will yield the same result. However, it is far simpler to understand the bottom than the top code. Reading the clean `get_students_eligible_for_exam` function let the developer immediately understand the two eligibility criterias for the exam. The comments like in the upper example are not even necessary, since the code is self-describing. Consequently, using or modifying the right code is much easier and time-saving.

An automated tool for checking violations of the clean code rule would preserve the readability and understandability of the code for future modifications. Integrated into a build pipeline, it could reject code changes if necessary and act as a quality gate.

1. Introduction

```
class Student():
    def __init__(self, name):
        self.name = name
        self.exercises = []

def get_students_eligible_for_exam(students):
    #filter for students with enough points
    enough_points = []
    for student in students:
        points = 0
        for e in student.exercises:
            points += e
        if points >= 60:
            enough_points.append(student)

    #only students with at least 8 completed exercise
    eligible = []
    for student in enough_points:
        if len(student.exercises) >= 8:
            eligible.append(student)
    return eligible
```

```
class Student():
    def __init__(self, name):
        self.name = name
        self.exercises = []

class ExamEligibilityChecker():
    MIN_NUMBER_OF_POINTS = 60
    NUMBER_OF_SOLVED_EXERCISES = 8

    def has_student_enough_points(student):
        total_points = 0
        for exercise in student.exercises:
            total_points += exercise
        return total_points >= ExamEligibilityChecker.MIN_NUMBER_OF_POINTS

    def has_student_enough_completed_exercises(student):
        return len(student.exercises) >= ExamEligibilityChecker.NUMBER_OF_SOLVED_EXERCISES

def get_students_eligible_for_exam(students):
    eligible_students = []
    for student in students:
        if ExamEligibilityChecker.has_student_enough_points(student)
        and ExamEligibilityChecker.has_student_enough_completed_exercises(student):
            eligible_students.append(student)
    return eligible_students
```

Figure 1.1.: Motivating example code that checks, if a student is eligible for an exam. The first listing shows a smaller, more chaotic version of an implementation. The second example applies some clean code rules like a single purpose per function or method calls in if-conditions. Further clean code improvements are possible.

Furthermore, it can teach students about clean code rules. In practical works, an automated checker could assess the students' code and recommend improvements. The students can change the code and get immediate feedback if they bettered or worsen the code.

1.1. Objectives and Contributions

This thesis focus on the automated detection of clean code violations on Python source code. We have the following main objectives:

- Design and implement a platform that can be extended with different automated checkers.
- Compare machine learning models on detecting code patterns that violate the clean code rules.
- Evaluate the generalisation capabilities of unseen pattern variations.

To fulfil the objectives, we made the following contributions:

- We provide an overview of code quality, clean code rules, quantitative metrics and existing tools.
- We introduce a taxonomy for implementation complexity of automated checker for clean code violations.
- We contribute a design and implementation of the clean code analysis platform (CCAP) with focus on extensibility, useability and integration capabilities.
- We extend the CCAP with checkers for returning `None` and using direct comparisons in conditionals.
- We compare our CCAP with preexisting tools.
- We collect a dataset and evaluate a support-vector classifier, random forest classifier, gradient boosting classifier and an LSTM-based neural network on detecting problematic code patterns for that we may not find a deterministic checker.
- We introduce an automated way to manipulate the code to contain unseen problem variations.
- We simulate the impact of having a hand-collected dataset on the machine-learning models if the samples do not contain all possible variations of a problem pattern. We give ideas on how to handle the situation in an time-optimising way.

For evaluation, we will discuss the following research questions:

RQ1 What is the utility of the CCAP besides existing tools?

RQ2 How do different models compare on the task of detecting non-clean code?

RQ3 Do machine-learning-based models cover a larger variety of cases than rule-based checker?

1.2. Structure

The remainder of this thesis is organised as follows: In Chapter 2, we introduce code quality and clean code principles. We then classify the clean code principles depending on the implementation complexity of an automated checker. We introduce several quantitative metrics to measure code quality and provide an overview of existing tools that can determine the code quality and clean code violations to an extent. In Chapter 3, we describe our approach for the CCAP and the machine-learning-based clean code classification. We evaluate the CCAP and machine-learning models and answer the research questions in Chapter 4. Finally, we conclude our work with a summary and outlook in Chapter 5.

2. Background and Related Work

2.1. Code Quality

The term code quality defines an evaluation of source code concerning the readability and understandability of the source code. If it is well-written and can be understood by developers easily, it is considered to be high quality. TODO rewrite to reliability, robustness, maintainability, readability siehe Softwarequalität wikipedia seite, gibt besseren eindruck

One of the main reasons why companies and developers aim for high code quality is the maintainability of source code. With a better understandability of the source code, changing it becomes more efficient and less error-prone. Software maintenance consists mainly of those activities that involve changing source code:

1. Adding code for new features
2. Implementing changed requirements by modifying code
3. Fixing bugs to ensure correct functionality of the software.

The importance of good maintainability of source code is based on the associated cost. For project planning, different sources suggest planning 66 to 75% of the total software costs for maintenance [45, 14]. A recent study from 2018 between Stripe and Harris Poll found that developers spent 42% of their time maintaining code [38]. Increasing the effectiveness of maintaining code is consequently cost-relevant. The application of agile software development methodologies like Scrum is another important contributor for code to be changeable with the least possible effort. Scrum is a agile software development practice of building software in increments while having a running software at the end of each increment. One of the advantages is the adaptability of this model to changing requirements, since, with the end of each increment, a requirement change can be implemented [35]. The easiness of code changeability due to code quality becomes a business-critical requirement and can positively impact the maintainability in the following ways [6]:

1. Well-written code makes it easy to determine the location and the way source code has to be changed.
2. A developer can implement changes more efficient in good code.
3. Easy to understand code can prevent unexpected side-effects and bugs when applying a change.
4. Changes can be validated easier.

Besides the maintainability, a higher code quality makes it easier for new developers to understand the code base and be productive since it is easier to read and understand the code.

The understandability and readability of high-quality source code have implications going further than maintainability. The International Organization for Standardization provides the standard ISO/IEC 25010:2011 that measures code quality by the contribution of the following characteristic [36]:

1. Reliability
2. Performance efficiency
3. Security
4. Maintainability

TODO: additionally, the standard mentions functional suitability, compatibility, useability and portability Reliability is the ability of a system to run without critical exception that would cause the system to crash, become unavailable or result in an inconsistent state. The most important factors for reliability are correct multi-threading with access to shared resources and correct memory allocation and deallocation. Performance efficiency describes an optimal use of the resources like CPU-time, memory or network bandwidth. The security aspect focus on code that minimise possible attack vectors. Especially common vulnerabilities like SQL injections shall be prevented.

To sum up, code quality is a concept to describe source code from the view of a developers ability to understand and change code. This view influence characteristics like maintainability or reliability that can be used as metrics to measure code quality.

2.2. Clean Code

In chapter 2.1, we defined the terminology of code quality. Yet the question remains, how to achieve a high code quality. One approach is the clean code concept, coined by Robert C. Martin in his book “Clean Code: A Handbook of Agile Software Craftsmanship” [25]. Basically, source code can be seen as clean code or as non-clean chaotic code. The author describes clean code principles in the form of practical rules. Those practical rules shall result in intuitive code that is easy to understand and read, i.e. the code quality increases.

Chaotic code, on the other hand, is harder to read and understand. The author argues, developers produce chaotic code in a conflict between deadline pressure and the necessary effort to make code more intuitive and follow clean code principles [25]. The former pressure is created by deadlines that focus on the visible output like the functionality of the program, whereas the latter is not directly visible in the final product. If the success criteria are solely based on the visible output, the code quality may suffer by fast-to-write chaotic code. This behaviour is shortsighted since an accumulation of chaotic code reduces productivity over time. A larger system with chaotic code will slow down later modifications or additions of code. With the costs associated with maintainability as described in chapter 2.1, reducing the maintenance effort is a logical step. With the clean code principles, developers can read and understand code in a more intuitive way and the described costs of chaotic code will decrease [25].

The following sections explain the clean code guidelines following the book by Robert C. Martin [25]. TODO give a more detailed overview of the next chapters

2.2.1. Naming

The naming section covers several rules for naming variables, functions, types and classes. Good naming can make it easier to read and understand code. “Good” naming is an opinionated topic; For Robert C. Martin, the key factor of good naming is the descriptiveness of names [25]. Descriptive names provide the reader with enough information to understand what a variable stores or what a function computes. Abbreviations or symbolic names like *a1*, *a2* are not descriptive and do not provide information about the meaning. Using symbolic names for mathematical expressions could be an exception to the rule if the naming mirrors the expression.

To come up with a descriptive name, it is helpful to include a verb or verb expression for function names, because functions express an action. For class names, a noun

emphasises the object character of the class.

Another helpful property of a name is an easy pronunciation since it is easier to read and to talk about the code with other developers. For reading, the pronunciation is necessary for the inner subvocalization¹. Non-dictionary words or unusual abbreviations are harder to pronounce. Consequently, it is useful to make longer but descriptive and pronounceable names, especially since the autocomplete feature of IDEs will free the developer from typing the long name. Additionally, searching for long names works better than for short names, because long names are more likely to be unambiguous compared to shorter names. An acceptable exception to this rule are short variables in a small scope (e.g. variable *i* in a loop), but using *i* in a large scope could be ambiguous and troublesome for searching and understanding.

An old naming convention is the Hungarian naming convention. In Hungarian naming, the type is encoded as a prefix of a variable name. TODO sample and source. Nowadays, the type encoding is redundant to the automatic type inference of IDEs. The automatic type inference also prevents confusion for the reader if the type in the notation and the actual type are inconsistent. The same logic applies to a prefix for member variables and methods, since an IDE can automatically highlight those different tokens.

2.2.2. Functions

This chapter describes several rules regarding functions, that Robert C. Martin characterized [25].

First, functions should be small in length. Exceeding 20 lines should not be necessary in most cases. If a function is small, it can be read more easily and without scrolling.

Second, inside a function, if-, else- and while-statements should contain a function call for the condition and one function call for the body. With a descriptive naming (see 2.2.1), a function call documents the condition and body in a concise and readable way. In sample TODO, reading a condition with function calls does not require the reader to decrypt the logical expression. This can also save additional comments that explain the meaning of the condition. Furthermore, if the body also contains a function call, it is fast to understand the cause-effect relation of the if expression.

Third, functions should fulfil one purpose. Since functions may have to call several functions subsequently to perform the required computation, Robert C. Martin expands the rule that functions should only operate on one abstraction level [25]. This would

¹the internal speech when reading

result in the following structure: Functions with a low abstraction level handle data access and manipulation, e.g. string manipulation. Functions on a middle abstraction layer orchestrate the low-level operations. On the next higher level, the mid-level functions are orchestrated etc. Following this structure, a function has one purpose on the low level and one orchestration function on a higher abstraction layer.

Fourth, the number of function arguments should be three or lower. With many function arguments, it becomes harder to call a function. Simultaneously, testing becomes harder, especially if all argument combinations should be tested. Functions with none or one argument simplify testing a lot. Since testing can reveal bugs, testable code is crucial for reducing the number of bugs. If a function requires more arguments, it may be advisable to bundle those in a config object to pass multiple arguments as one. With a config object, arguments can be logically grouped that would help the reader to understand arguments.

Fifth, side-effects in a function body should be avoided. Side-effects happen if a function modifies a variable outside its scope without explicitly mentioning this in the function name. This leads to dependencies between the functions that are not obvious to a developer who checks the signature, i.e. the name, input arguments and return type, to use the function. To spot a side-effect, the developer would have to read through the function, although reading the signature should be enough to use the function. With side-effects, the function behaves unexpected and time-consuming mistakes are the consequence. Especially side-effects that initialise other objects result in a time-dependency that is hard to identify and could be overseen at all.

The next rule forbids the returning of error codes and suggests raising an exception if the programming language supports it. Raising an exception allows better separation between application logic and error handling code since the error code checking interrupts the reading flow of code. On the other hand, catching a raised exception is separated from the logic for a non-error execution and can be separated into an additional function for an even cleaner structure. TODO: see example In this case, error handling counts as the one purpose of this function.

Last, the important principle for software engineering applies directly to functions: Don't repeat yourself [18]. Repeated code is dangerous and chaotic since it requires changes in multiple locations if it has to be modified. This makes duplicated code very prone to copy and paste errors and small mistakes that may not be obvious during development but will lead to a fatal crash at runtime. Code clone detection is an active research field, especially for finding semantic clones without syntactic similarity. For

instance, Büch and Andrzejak suggest an AST-based recursive neural network [9] and Saini et al. present an approach that can even detect duplications with just minor syntactic similarities [34].

2.2.3. Comments

Comments are part of every programming language and can play an important role in code quality. Subjectively good comments clarify the meaning of the code and help to understand the code. However, comments can become outdated and wrong or provide useless information [25]. In a perfect world, the programming language and the programmer would be expressive enough that commenting is not required for clarity. Following the clean code guidelines for naming and functions makes many descriptive comments obsolete, since the function description is encoded in the function name.

Comments will not help to turn chaotic code into clean code. If a developer explains a line of code by a comment, it is often more helpful to call a function with a descriptive name to replace the comment. Since comments are not part of the program logic, if the explaining comment is not updated with the code, the comment gives misinformation to the reader like described before.

In brief, before writing a comment, the developer should think about expressing the same in code. Robert C. Martin gives some exceptions for situations, where comments can be helpful or necessary [25]:

Legal notes: Legal notes like copyright or license information and author mentions may be necessary. Although they should be short and link to an external licensing document in full extent.

Explaining comments: Some explaining comments can be helpful and are not easy to encode in normal code. For instance, an explanation of a complex regular expression may be too complicated for a function name and a comment provides the space for a sufficient explanation. This increase the readability since the developer does not need to interpret the regular expression manually.

Intention: Explaining an intention which does not become obvious by reading the source code might also be a valid use-case for a comment.

Warning for consequences: Some parts of the code can have special consequences like not being thread-safe or using many system resources. A warning can spare a developer from having trouble when using the function.

Emphasise : A comment to emphasise the importance of a seemingly unimportant part of the code prevents breaking modifications of the code.

2.2.4. Data Structures and Objects

This chapter describes the rules of data structures and objects. Both store data but the way this data get exposed differ. Objects hide data behind abstractions and have functions that work with these data. Data structures expose the data and do not provide functionality.

For objects, I. Holland defined the Law of Demeter (LoD) [24]. In object-oriented programming, a method m of an object O may only call methods of the following components:

- the object O itself
- the parameters of the method
- objects that are created within m 's scope
- instance variables of O

In other words, the object O only calls methods on direct “neighbours” or on itself. The object does not call a method of its neighbour to subsequently call a method on the returned object from the neighbour. By restricting the allowed method calls outside the object itself to direct “neighbours”, the dependencies between objects are reduced and the modularity increased. As a consequence, changing methods only affects the direct neighbours and not objects in potentially different parts of the software. Figure 2.1 shows an example for the LoD. In the method `test`, a function call on `B` returns a reference to `C`. Since `C` is not a direct neighbor of `A`, the function call violates the LoD. The modularity of objects suffer, since a change of the method `do_something()` would require changes in object `A`. This close coupling of an `A` and `C` is circumvented with a method `do_something_on_C` that wraps the functionality of `C`, so a change in `C` would only require a change in `B`. With that common layer, `A` and `C` are decoupled. Object `A` in this case does not know how the certain functionality is performed and it does not matter. Bad code practice like the violation of the LoD in Figure 2.1 is called a train wreck, since subsequent method calls look like multiple train carriages. Such train wrecks can be improved by the described wrapping function.

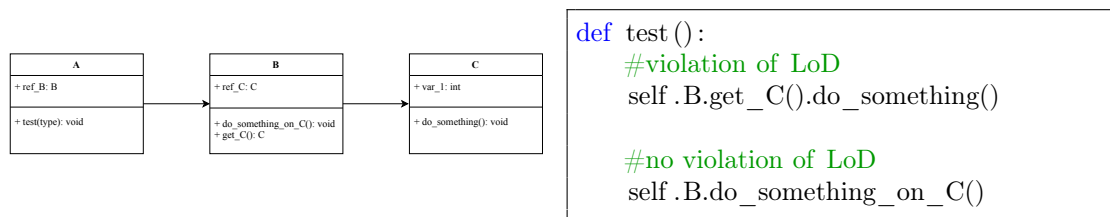


Figure 2.1.: A sample illustration of the Law of Demeter. The figure on the left shows a class diagram for the classes *A*, *B* and *C*. The listing on the right shows the implementation of the *test* method of class *A*.

Data structures are often used as Data Transfer Objects (DTO) to communicate with other processes or services. These objects do not contain any functions; instead, they only have accessible member variables. A specialisation of DTOs are Active Records that contain additional methods for data storage. They are used to represent a data source like a database. For DTOs and Active Records, it is bad practice to insert business logic into these objects [25]. With business logic, the data structure becomes a hybrid between an object and a data structure. It becomes unclear, if the purpose of a hybrid is to store data or to expose methods to modify those data in a controlled way. Without the clear purpose of the object or data structure, it becomes harder to understand for developers.

2.2.5. Classes

For classes, the clean code principles describe multiple rules that improve the readability and understandability of classes.

The first rule specifies the recommended size of a class. Unlike counting lines as for functions, the size of a class is the number of responsibilities. A responsibility of a class can be manipulating data or writing data to a file. A single responsibility per class is described as Single-Responsibility-Principle and offers several advantages: First, the naming can reflect this responsibility, so the class functionality is easy to understand. Second, changing one functionality should only require changes in the one class that is responsible for this specific function. And last, we can reuse the class in different projects if we need that functionality. If the class fulfilled two responsibilities, we would have to remove the unnecessary functionality before reusing the class in a different project. Reaching back to the examples of one responsibility being data manipulation and one is writing data to a file, reusing the latter functionality would require us to strip out the data manipulation. The Single-Responsibility-Principle leads to a system of many small classes with one, clearly defined responsibility.

The second rule recommends classes to have high cohesion. Cohesion is high if a method manipulates many instance variables. If all methods of a class manipulate all instance variables, the class has peak cohesion. With a high cohesion, the methods and the class are a logical unit. A low cohesion indicates to split the class into multiple classes since the methods are not a logical unit. As a result of splitting classes into smaller logical units, will more likely comply with the Single-Responsibility-Principle. (TODO add metric for cohesion)

The last rule provides a guideline on how to handle internal or third-party class dependencies. If a dependency changes, the class depending on it should not change. To accomplish this decoupling from the implementation of a dependency, classes should not depend on the dependency directly but instead on an abstract class (or interface). That abstract class describes the concept of the dependency that is unlikely to change. For example, for a database dependency, the abstract class defines the concepts of storing and retrieving an entry. For the specific database, a specific class implements the abstract class or interface and provides the functionality. Changing the database only changes the implementation in the specific class but not the concept. Therefore, there is no need to modify a class that depends on the concepts from the abstract class. This isolation of dependencies is especially useful if the dependencies are not under the developer's control and may be changed any time by a third party. The abstraction is also valuable for testing, since the dependency can be replaced by a so-called mock class that only simulates the behaviour. This allows testing to be independent of the dependency and to locate potential errors in the own system or the dependency.

2.2.6. Exception Handling

This chapter covers clean code rules for error handling. Basically, the following ways are methods for handling errors: First, a function can return an error code if some exception occurred. The function caller has to check for the error code and act accordingly. Second, some programming languages support the concept of throwing and catching an exception. And last, a function can return a null value to indicate the execution failure.

In section 2.2.2, we described the clean code rule that prefers throwing an exception over returning an error code. The supporting argument of Robert C. Martin is the better separation between application logic and error handling [25]. Error codes have to be checked immediately, so the error checking interrupts the code for the application logic. By throwing and catching errors, the error handling can be completely extracted into a sperate function and can be removed from the main application logic. Furthermore,

the try-catch block enforces transactional behaviour. At the end of either the try or the catch block, the application should be in a consistent state. Error codes hide this explicit transactional behaviour.

Notwithstanding the clean code rules, Go as more recent programming language² returns explicit error types instead of throwing an exception. This was designed to “encourage you to explicitly check for errors where they occur”, instead of “throwing exceptions and sometimes catching them” [15]. So the practical application of the clean code rule for throwing exception depends on the design of the programming language.

The third method of error handling is returning a null value. Instead of checking for a specific error code, the caller would check for a null return value. If a developer misses the check, the program will terminate with an unhandled `NullPointerException` at runtime. Tony Hoare introduced the null reference in 1965 and later called it his “billion-dollar mistake” [20], since it leads to many bugs and security vulnerabilities throughout the decades. Languages like Kotlin are designed with null safety enforced. Kotlin distinguishes between nullable and non-nullable references, with a compiler enforcing null checking [2]. If a language does not enforce null safety with a compiler, a special case like an empty collection or an optional type can be returned. Returning null introduce the same problems like error code handling such as mixing application and error handling logic and is additionally not explicitly marked as an error code. From a function signature, it may not be obvious that the return value may be null if the language does not support compile-time checking or optional types. Robert C. Martin describes returning null as a violation of the clean code rules [25].

Independent of the error handling method is the next rule for error handling with a third-party library. In section 2.2.5, we described wrapping dependencies like third-party libraries to decrease the coupling to the dependency. The wrapping concept also applies to error handling, since the wrapper can unify the exceptions, so the application does not check for dependency-specific error types. Changing the dependency requires only a change in the wrapper implementation but not everywhere this wrapper is used.

2.2.7. Additional Rules

Robert C. Martin describes more rules for system design, multi-threading, testing and third-party code [25]. Additionally, the author describes the concept of code smells. In theory, following these rules seems to lead to clean code. In practice, however, developers

²designed in 2007 as a response to problems with C++, Java and Python [1]

tend to violate rules in some situations. A code smell is a code characteristic that could indicate such a violation.

This work will focus on the described rules and implement a detection for some rule violations.

2.3. Source Code Analysis

In order to define and measure code quality, methods for analysing a program and its source code are necessary.

To analyse source code, different techniques are bundled as static code analysis. A static code analysis gathers information about the structure of the source code and program [31].

The following provides a non-comprehensive overview of the most important methods for code analysis and representation. Source code itself is stored as encoded text files. This representation offers no structural information.

A first processing step is a lexical analysis, also known as tokenization. The tokenization transforms the character stream from the text file into a stream of tokens. This is possible due to the syntactic definition of a programming language, that allows to split the raw text into typed tokens [27]. A token type can be for instance a keyword, operator or an identifier.

With the token stream, a parser can build an Abstract Syntax Tree (AST) that represents the abstract structure of the source code. Since the AST represents the code structure, it is possible to traverse the tree-like data structure to analyse the code structure.

Another method is the control flow analysis. For the control flow analysis, the program has to be represented as a control flow graph. The control flow graph represents possible execution paths in a program and consists of nodes and edges. A node represents a basic block, i.e. a code sequence without branching. A directed edge represent jumps in the control flow; e.g. an if-statement would have two outgoing edges, one for a true condition and one for a condition evaluated as false [4]

Besides the control flow, it is possible to analyse the data flow in a data flow graph. With the data-flow analysis it is possible to collect information about the read and write operations to variables. In combination with the control flow graph, the data flow for different control flow can be calculated [27].

Last, a call graph can be computed. A call graph maps all method or function calls

starting from the main method. It allows to spot dependencies between modules [31].

2.4. Clean Code Complexity Levels

We described the clean code rules in Section 2.2. This work focus on the automated checking of clean code rules. Some rules can be checked straightforward with an algorithm, whereas others may be too abstract to be checked in an automated way. In this Section, we want to categorise the clean code rules into different complexity levels. The complexity levels are based on the complexity of implementing an automated rule checker.

We define the taxonomy of complexity in three levels, based on the different source code analysis techniques and effort needed:

Basic Level The level of basic complexity covers all rules that can be checked by analysing the text, token stream or AST of the program.

Advanced Level Rules with medium complexity require a control flow, data flow, call graph analysis or more advanced, deterministic analysis techniques. Furthermore, a combination of different analysis methods from the basic complexity level will be categorised as medium complexity.

High Level The high level of complexity covers the clean code rules we think can not be checked with traditional, deterministic analysis techniques such as those described in Section 2.3. For this level, it may be indispensable to use statistical methods like machine learning approaches to create an automated checker. It is important to note that we do not prove that no deterministic algorithm exists.

Based on these levels, we categorise the clean code rules described in Section 2.2. Table 2.1 gives an overview over all rules and our assigned complexity level.

For naming, all rules have in common, that an algorithm would have to extract the names of various entities like variables, classes, functions from the source code first. The AST provides the names and entity types they describe. Extracting all names is a matter of walking over the AST. We see this in the basic complexity level.

The first important rule is the descriptiveness of names. Assessing the descriptiveness of a name is a problem with a high level of complexity. An algorithm would have to understand the purpose of the source code to determine if the name describes the source code. It is an ongoing research effort to use machine learning on code for different tasks.

Category	Rule	Level		
		Basic	Advanced	High
Naming	Descriptiveness			X
Naming	Pronunciation	X	X	
Naming	Names in Small Scope		X	
Naming	Hungarian Notation	X		
Functions	Length	X		
Functions	Conditional Function Call	X		
Functions	Purpose			X
Functions	Levels of Abstraction			X
Functions	Number of Arguments	X		
Functions	Side-Effects		X	X
Functions	Error Code Detection			X
Functions	Duplicated Code			X
Comments	Detecting Necessary Comments			X
Objects	Law Of Demeter		X	
Objects	Mix with Data Structures	X		
Classes	Single Responsibility		X	X
Classes	Decouple Dependencies		X	
Errors	Return None	X		

Table 2.1.: Overview of our categorization of automated approaches to check for compliance with the clean code rules.

The tasks of code summarisation explore the use of machine learning to generate a natural language description of source code. In a paper from 2018, Alon et al. describe a neural network to represent code as vectors and “predict a method’s name from the vector representation of its body” [5]. Based on that and following research, it may be possible to compare the summarisation result with the chosen name to detect non-descriptive names. Nevertheless, it may be harder to apply such concepts to variable names, since a variable name has less context information compared to the method name with its method body.

The next rule covers the easy pronunciation of names. A rule violation would be the use of abbreviations or non-dictionary words in variable names that are hard to pronounce. An approach with a basic level of complexity could be a minimum number of characters for every name. A more advanced approach could utilise naming conventions such as camel case or underscore separated words to extract the single words of a name and compare it with a dictionary. While it would help to prevent spelling mistakes and

should ensure pronounceability, a dictionary may not contain all words. Especially with domain-specific words, this could lead to false alarms.

Variable names in a small, local scope are excepted from the rule of pronounceability and thus longer names. To determine the scope of a variable, the data flow analysis could be used. With the analysis, the previous rule could include this exception. Due to the data-flow analysis, we categorise this rule into the advanced level of complexity.

Detecting Hungarian notation could be accomplished with regular expressions on the variable names since the Hungarian notation follows predefined notation strategies. Therefore, such a rule checker would be a basic level of complexity.

For functions, the clean code principles define multiple rules. First, functions should be smaller than 20 lines of code. A checker for this rule would be at the basic level of complexity since the AST contains all the necessary information about the function body. If the metadata such as line number is retained during AST construction, a checker can calculate the lines of code for the function body.

The next rule covers the function call in the condition and body part of a conditional such as an if-statement. An if-statement is represented as a node with a condition, body and an `else` descendant. An algorithm can scan the AST to check the condition and body nodes of the conditional statement. Such an checker would have a basic level of complexity.

Regarding the purpose of a function, the clean code rules suggest a function to have one purpose and to operate on one abstraction level. Counting the number of purposes requires an understanding of the semantic of the code. Similar understanding would be necessary for the level of abstraction. We argue such a checker would be in the high level of complexity category since the understanding of code could only be achieved with machine learning. The aforementioned research field of code summarisation could be beneficial for this task, too. A natural language description of a function may require the model to “understand” the purpose and task in order to articulate it in natural language. It may be possible to use work in this field to detect functions with multiple purposes or abstraction levels.

The rule that limits the number of function arguments is on a basic level of complexity. The AST contains the information about the code structure and therefore, all function arguments. Counting those requires the algorithm to search for function nodes in the AST and count the number of arguments.

The next rule forbids side-effects of a function. Detecting side-effects can be broken down into a two-step process. First, tracing the data manipulation of functions is at the

advanced level of complexity. With a data-flow analysis and call graph, it is possible to check what variables a function manipulates. Second, the classification whether an effect of a function is a side-effect requires an understanding of the intended effect of a function. The intended effect is described in the name. We argue that the understanding of the intended effect is at the high level of complexity.

The clean code rules discourage the use of error code to indicate exceptions. If the error codes are not part of the programming language specification (as they are not for Python), it could be challenging for an automated checker to distinguish between returned error codes and non-error values. There may be indicators such as an early return of a constant value or comparing a returned value with a constant value. This indicators as a heuristic may be good enough for practical use. We would classify such a checker at a high level of complexity.

The automated checking for repeated code is an active research topic. As described in Section 2.2.2, efforts are made with several tools to detect repeated code. Especially the detection of duplicated semantic is hard with just a few syntactic similarities. Therefore, we categorise a code duplication checker as a high level of complexity.

The clean code rules suggest to not use comments except for legal notes, explaining comments, comments stating intention, warnings for consequences and emphasising comments. An automated checker would flag all comments, except they contain such an exception. Finding all comments in code is possible on the token stream. While this has a basic level of complexity, detecting the exceptions is harder. The first comment at the start of a file could be seen as a legal note and a warning for consequences may include the term “warning”. However, distinguishing necessary explanations or intentions from useless comments is a task in the high level of complexity category.

For objects, the clean code principles require compliance with the Law of Demeter that restrict methods allowed to be called inside an object’s method. An automated checker may find the violations by analysing the ASTs of all classes. It can then assess if a method call is violating the LoD. Due to the analysis of multiple AST to determine allowed “neighbours” of an object, we see this checker on the advanced level of complexity.

The next rule discourages the mixing of data structures with objects by introducing logic. Such a hybrid would have accessible member variables and methods. An algorithm can detect those hybrids based on an analysis of the AST. We see such a checker in the basic level of complexity.

For classes, the clean code rule enforces a single responsibility per class. As discussed when detecting a single purpose of a function, detecting a single responsibility of a class

requires an understanding of the code. Similarly, a checker would fall into the high level of complexity. On the other hand, it may be possible to use the cohesion metric as an indicator for too many responsibilities of a class. A high cohesion will result in smaller classes that form a logical unit. Smaller logical units can be expected to have fewer responsibilities. Cohesion can be measured as *Lack of Cohesion in Methods* (LCOM) as described by Hitz and Montazeri [19]. An analysis of the AST and call graph should provide all information to calculate the LCOM number. Therefore, we see an automated checker for cohesion on the advanced level. Using the cohesion as a heuristic for the responsibilities of a class, we could detect potential violations against the Single-Responsibility-Principle with an advanced level of complexity.

An automated check for the decoupling of dependencies can be performed with an advanced level of complexity. All required information is in the code. Finding all class dependencies is possible by combining several ASTs. Following the dependencies, it is possible to determine whether they inherit from an abstract class or implement an interface. In that case, the checker would assume a dependency decoupling. Based on our classification schema, this checker would be in the advanced level of complexity.

Last, the clean code rule prevents the return of the null value as an error code or as an indication for no result. A checker could prevent this with a basic level of complexity. It could search for all `return` statements in the AST and check for the `null` constant.

2.5. Quantitative Metrics for Code Quality

To define the quality of a software product, it is necessary to have methods to express the code quality as quantitative units. This has multiple advantages:

- A metric sums up the quality of a project in a single unit.
- A quantitative approach tracks the changes in code quality over time. Therefore, it is obvious if code quality improves or not. In case the quality undercuts a threshold, special measures like mandatory refactoring can be undertaken.
- A developers performance can be evaluated based on the code quality. Since the maintainability and reliability of the software depends on the code quality, this is a good incentive to enforce high-quality work.
- A certain level of code quality may be required by a contract. As a result, a customer may have fewer bugs and a smaller maintenance effort.

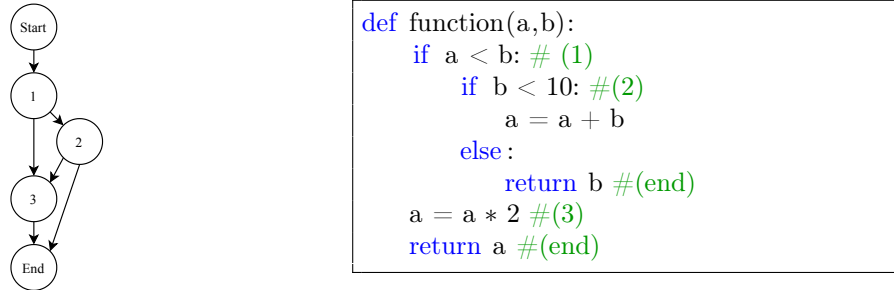


Figure 2.2.: Control-flow graph visualisation for a code sample. The graph has 6 edges, 5 nodes and one connected component. The cyclomatic complexity is 3.

Metrics for code quality are a subset of general software metrics. The following sections describe common software metrics that express code quality.

2.5.1. Cyclomatic Complexitiy

Cyclomatic Complexity is a metric for the complexity of a section of code. It was introduced by Thomas McCabe. It measures the complexity by counting the number of linearly independent execution paths [26].

An execution path is the subsequent execution of instructions. With control flow structures such as if statements, two execution paths are possible, depending on the evaluation of the if condition. An execution path is linearly independent if it includes one subpath that is not part of any other path. A control flow graph represents all possible control flows as described in 2.3. Figure 2.2 shows a control flow graph for a sample function.

The cyclomatic complexity on a control-flow graph as the number of linearly independent execution paths is defined as:

$$M = E - N + 2P \quad (2.1)$$

E is the number of edges, N the number of nodes and P the number of connected components. A connected graph is a subgraph, in which all nodes are connected to each other by a path. Following this definition, the cyclomatic complexity can be calculated. Figure 2.2 shows a visualisation of a sample python function. The control-flow graph has six edges, five nodes and is one connected component. Using equation 2.1, the cyclomatic complexity is 3.

MacCabe recommends limiting the cyclomatic complexity to 10 [26]. A lower cyclomatic complexity improves testability since the complexity represents the number of

execution paths that need to be tested. Therefore, M is the upper bound for the number of test cases for full branch coverage. Software that has to comply to safety standards like ISO 26262 (for electronics in automobiles) or IEC 62304 (for medical devices) are mandated to have a low cyclomatic complexity [22, 21].

Although cyclomatic complexity is used throughout the industry, several shortcomings are under critique. First, complexity from data flow is ignored. Working with a larger number of variables and operations, the code becomes more complex, but the cyclomatic complexity does not take this into account. Second, nested code structures are not considered by the metric, although it adds additional difficulty for understanding [46].

2.5.2. Halstead Complexity Measures

Maurice Halstead introduced the Halstead complexity measures (HCM) in 1977 [17]. Halstead approached the complexity measure with an empirical approach by defining observable and measurable properties and put them into relations.

The measurable properties are operators and operands. Operators are symbols in expressions like an addition, subtraction or multiplication symbol. Operands are the values and variables that are manipulated by the operators.

The sum and number of distinct operators and operands are counted as:

- η_1 as the number of distinct operators
- η_2 as the number of distinct operands
- N_1 is the sum of all operators
- N_2 is the sum of all operands

From these base properties, Halstead derived additional properties:

- Program vocabulary size:

$$\eta = \eta_1 + \eta_2$$

- Program length as the sum of all operators and operands:

$$N = N_1 + N_2$$

- The volume of a program in terms of program length and program vocabulary is defined as:

$$V = N * \log_2 \eta$$

Based on those properties, code metrics can be calculated. First, the difficulty of understanding a program is calculated as:

$$D = \frac{\eta_1}{2} * \frac{N_2}{\eta_2}$$

The major contributors to difficulty are the number of distinct operators η_1 and the sum of all operands N_2 .

Second, the combination of difficulty and volume is the required effort for understanding or changing code follows:

$$E = D * V$$

Third, the effort for changes translates into real time following the relation:

$$T = \frac{E}{18} s.$$

Last, the number of bugs correlates with the following relation [46]:

$$B = \frac{V}{3000}$$

HCM focus on the complexity introduced by data manipulation but ignores complexity from the control-flow. Since the cyclomatic complexity focus on the control-flow complexity but not on the data-flow complexity, it is practical to use both together to supplement each drawback [46].). Additionally, the correlation with the number of bugs is based on the programmer's skill estimation with a fixed value of 3000. Since this varies between projects, such an exponential and fixed number is doubtful [46].

2.5.3. Software Maintainability Index

The software maintainability index was developed by Dan Coleman and Paul Oman in 1994. 16 HP engineers evaluated 16 software systems and scored it in a range from 0 to 100, with 100 representing best maintainability [12]. Following a regression analysis, they identified the following equation to match the maintainability of the evaluated systems:

$$MI = 171 - 5.2 * \ln \bar{V} - 0.23 * \bar{M} - 16.2 * \ln \overline{LOC} + 50 * \sin \sqrt{2.4 * C}$$

\bar{V} is the average Halstead Volume, \bar{M} the average cyclic complexity, LOC the lines of code and C as fraction of comments.

The software maintainability index was defined many years ago with a limited sample size of developers and projects. Additionally, programming languages have changed significantly over time. A study by Sjøberg et al. suggest no correlation between the software maintainability index and actual maintenance effort in a controlled environment [37]. Although the study only analysis four software projects and lacks generalisation, they find a strong anti-correlation with different maintainability metrics. Only the code size seems to correlate with the actual maintainability. The former seems to be consistent with a systematic review on software maintainability predictions and metrics by Riaz et. al [33].

2.6. Tools for Code Quality Analysis

It is good practice to use static code analysis tools to improve code quality. A static code analysis examines a program by analysing an abstract syntax tree, the control or data flow, a pointer analysis or an abstract (approximated) execution. In comparison to dynamic analysis, the code is not executed. The result of static analysis can be based on approximations, so there might be false-positive results. An expert has to examine the result and correct results if necessary [31].

Besides the use of static code analysis for compiler and optimisations, it is also helpful for analysing the code quality. Different categories of code quality-related principles can be analysed:

Code Guidelines: A static code analysis can ensure compliance to structural, naming and formatting code guidelines.

Standard Compliance: A requirement for a software may be compliance to an industry-standard like IEC 61508 or ISO 26262 [39, 22]. Static code analysis can ensure or at least assist with the compliance.

Code Smells: Some code smells follow a known pattern and a static code analysis can detect those smells. Since code smells can be an indication for chaotic code, it is best if these smells are detected and removed early on.

Bug Detection: Although bug detection with static analysis can not detect all bugs, every detected bug before a software release is important. Examples for detected

bugs are unhandled expressions or concurrency issues [13].

Security Vulnerability Detection: Static code analysis tools can detect some security vulnerabilities like SQL Injection Flaws, buffer overflows and missing input validation with high confidence. Since these are some common, easy to exploit vulnerabilities, an analysis for security vulnerabilities can increase the security of the overall software [44].

The following sections provide an overview of a few static code analysis tools with a focus on code quality and maintainability. We selected them based on popularity and if the projects are still maintained.

2.6.1. PyLint

PyLint is a code analysis tool for the Python programming language. It is open-source and licensed under the GNU General Public License v2.0 and available for all common platforms. PyLint can be executed as a standalone program or can be integrated into common IDEs like Eclipse. A continuous integration pipeline may include PyLint as well to ensure an analysis on every build.

With PyLint, the developer can make sure the code complies to the PEP 8 style guide for python coding [43]. This includes name formatting, line length and more. It does not calculate a metric for the code but instead warns about violated principles. Additionally, PyLint can detect common errors like missing import statements that may cause the program to crash at the start or later at runtime. To support refactoring, PyLint can detect duplicated code and will suggest refactoring the code.

PyLint can be configured to ignore some checks and to disable specific rules. To expand the ruleset, a developer can write a "checker", an algorithm to check for a specific rule. The algorithm can analyse the raw file content, a token stream of the source code or an AST representation. The checker can raise a rule violation by providing the location information and the problem type to the PyLint framework and it will be included in the PyLint analysis report.

2.6.2. PMD Source Code Analyzer Project

PMD is an "extensible cross-language static code analyser" for Java, JavaScript and more. As an open-source project, it is licensed under a BSD-style license and is available for macOS, Linux and Windows-based systems. It integrates into build systems like

Maven and Gradle as well as into common IDEs like Eclipse and IntelliJ. PMD can run as part of a continuous integration pipeline and is included in the automated code review tool Codacy (see TODO).

Depending on the target language, PMD supports different rules. For Java, PMD has rules in the following categories:

Best practices: Best practices include rules like one declaration per line or using a logger instead of a *System.out.print()* in production systems.

Coding Style and Design Several rules to improve the readability of the code like naming conventions, ordering of declarations and design problems like a violation of the Law of Demeter, Cyclomatic complexity calculation and detection of god classes.

Multithreading: Rules to mainly prevent the use of not thread-safe objects or methods. Due to the nature of multi-threading and unpredictable scheduling of threads, problems like unpredictable values of variables or deadlocks may occur. Since they may not occur in every execution, they are hard to spot and to debug. Consequently, warnings of using non-thread-safe code may save hours of debugging.

Performance: These rules flag known operations that have hidden performance implications. For example, string concatenation with the plus operator in Java causes the Java Virtual Machine to create an internal String buffer. This can slow down the program execution if numerous String concatenations are performed.

Security: Security-wise, PMD only checks for hardcoded values for cryptographic operations like keys and initialisation vectors.

Error Prone Code: PMD checks for several known code structures that will or may cause a bug at execution. Some rules are part of the clean code principles described in 2.2 and some rules are specific to the target language.

A user can expand the PMD ruleset in two ways: A XPath expression can be specified and will be validated against the AST by the PMD framework. For more control, it is possible to write a Java-based plugin that implements a custom AST traverser. The latter allows for more sophisticated rules and checks.

2.6.3. Codacy

Codacy is a software to "automate your code quality" (TODO source) for more than 30 programming languages. It is available as a cloud-based subscription service with a free tier for open-source projects and a self-host option for enterprise customers. Codacy runs as cloud software and a user can connect a GitHub, GitLab or Bitbucket repository that will be scanned automatically on every push or trigger a scan of local files with a command-line program. Additionally, a badge can be added to the readme page to show off the analysis results.

Codacy can be seen as a platform that runs multiple different "analysis engines". These analysis engines combine multiple tools depending on the language. For Python, Codacy uses PyLint, Bandit (a scanner for security issues) and a metric plugin. Due to the licensing of those analysis engines, the engines are open-sourced, whereas the Codacy platform is proprietary.

Depending on the language, Codacy supports scanning for common security and maintainability issues. The later is faced with scanning for Code Standardization, Test Coverage and Code Smells to reduce technical debt.

Codacy allows customisation by disabling certain rules and changing rule parameters like the pattern to fit the rules to the project.

2.6.4. Sonarqube

Sonarqube is an analysis tool to maintain high code quality and security. It supports 15 programming languages in the open-source version licensed under the GNU Lesser General Public License, Version 3.0. In the Developer Edition, Sonarqube supports 22 languages and 27 languages in the enterprise edition. Sonarqube is a self-hostable server application with a SonarScanner client module that can be integrated into build pipelines like Gradle and Maven as well as a command-line tool for other build pipelines. The SonarScanner reports the result to the server, that serves a website to review the results. With the Developer Edition, branch and pull request analysis are possible and a pull request will be annotated with the analysis results. TODO source

For Python, Sonarqube offers more than 170 rules. These rules are in the following categories:

Code Smells: Code Smells like duplicated String literals are flagged with four different levels of severity (from higher to lower severity): Blocker, Critical, Major, Minor. Explanations and examples are accessible for the developer to understand and

fix the issue. The analysis is more in-depth than AST-based solutions since it additionally uses control and data flow analysis.

Bug: Multiple rules cover bugs that would definitively result in a runtime exception and program termination. As an example, calling a function with the wrong number of arguments will be flagged with the highest severity, since it will raise a `TypeError` during runtime. Although programmers may notice issues like this during coding and testing, the issue can remain hidden if it only triggers a special execution path of the system.

Security Hotspot: The Security Hotspot analysis unfolds pieces of code that may be a real vulnerability and requires a human review. The scanning has a hotspot detection for seven out of the OWASP Top Ten Web Application security risks (TODO src). A flagged hotspot contains a detailed explanation of the reason for being flagged and a guide on how to review this hotspot. The recommendation to double-check a piece of code can help to prevent a security vulnerability from being deployed to production.

Security Vulnerabilities: A security vulnerability analysis reveals code that is at risk of being exploited and has to be fixed immediately. For example, misconfigurations of cryptographic libraries can be revealed easily and possible damage can be proactively prevented.

Additionally, Sonarqube offers several metrics like test coverage and a custom derivate of Cyclomatic Complexity named Cognitive Complexity [10]. The authors see several shortcomings in the original Cyclomatic Complexity model:

- Parts of code with the same Cyclomatic Complexity do not necessarily represent an equal difficulty for maintenance.
- Cyclomatic Complexity does not include modern language features like try/catch and lambdas. Therefore, the score can not be accurate anymore.
- Cyclomatic Complexity lacks useability on a class or module level. Small class with complex methods would have the same Cyclomatic Complexity as large classes with low complexity per method (such as getter or setter methods). The informative value of the metric is low.

Cognitive Complexity addresses these points, especially the incorporation of modern language features and meaningfulness on class and module level.

The calculation is based on three rules [10]:

Ignore Shorthand Structures : Method calls condense multiple statements into one easy to understand statement. Therefore, method calls as shorthand are ignored for the score. Similarly, shorthand structures like the null-coalescing operator reduce the Cognitive Complexity compared to an extensive null check and are therefore ignored for score calculation.

Break in the linear control flow : A break in the expected linear control flow by loop structures and conditionals adds to Cognitive Complexity as well as to Cyclomatic Complexity. Additionally, catches, switches, sequences of logical operators, recursion, and jumps also adds to Cognitive Complexity, since these structures break the linear control flow.

Nested flow-breaking structures : Flow-breaking structures that are nested additionally increase the Cognitive Complexity, since it is harder to understand than non-nested structures.

The method-level Cognitive Complexity score represents a relative complexity difference between methods that would have the same Cyclomatic Complexity. Furthermore, the aggregated Cognitive Complexity on a class or module level differentiates between classes with many, simple methods and classes with a few, but complex methods. As a result, domain classes (containing mainly getters/setters) have a lower score compared to classes with complex code behaviour [10].

Developers can extend Sonarqube similarly to PMD. They can write a Java plugin that can access the AST but also a semantic model of the code. Additionally, extensions can be created that provide the functionality to other extensions. The Java plugin is compiled into a .jar file and placed into the plugin dir of the Sonarqube installation. For simpler rules, XPath expressions are possible through the web interface and allow a quick extension. For instance, if developers observe bad code style during a pull request review, they can quickly write a rule to enforce this rule in all subsequent pull requests automatically.

TODO Integration in continuous integration, manual review

Tools review for all tools that check different stuff

Single Responsible principle, Open-Closed principle

3. Approach

In this chapter, we first describe our approach in designing and implementing the Clean Code Analysis Platform. Afterwards, we depict our approach to using machine learning models to detect violations of clean code rules.

3.1. Clean Code Analysis Platform

The goal of the design and implementation of the Clean Code Analysis Platform (CCAP) is a tool for software developers to improve the code quality of existing and new code. The tool accepts a directory containing source code files as input and analyses the input for snippets of improvable code quality. If the analysis classifies a code snippet as problematic, it should help the developer to improve the snippet with information about the problem. Ultimately, this should train the developer to spot problematic code by its own and to write clean code by default, so the number of alerts should decrease over time. At the same time, the overall software quality of a project increases immediately at rewriting a marked snippet and in the long term at training the developer to write code with higher quality.

In order to use the tool effectively, the design and implementation should cover the following requirements:

Useability: The CCAP should be an easy-to-use tool. Developers shall be able to install and run the tool without difficulties. The extra effort of using this tool should be small, and the developer should use the tool in his day-to-day workflow without additional friction. The developers can interpret the results and localise the issue immediately.

Extensibility: The extension of the detectable clean code violations should be easy. A clear defined interface for extensions is required. An extension developer would not need specific knowledge about the internal architecture of the tool. The extensibility allows the desired workflow of a developer finding problematic code in

an, e.g. peer-review, formalising it into an extension and sharing this extension with the team. With each iteration, the code quality of all team members would increase.

Integration: The tool should be easy to integrate into different systems. These systems include local workflows like git pre-commits or build systems and remote continuous integration/delivery/deployment pipelines.

A more specific requirement is Python as an input and extension language. After JavaScript, Python is the second most popular programming language 2019 according to the Github statistics [16] and the third most popular language according to the TIOBE index (TODO <https://www.tiobe.com/tiobe-index/>). Besides the general popularity, Python is heavily used in the scientific community for machine learning and at universities for teaching programming. These groups are part of the potential target audience, and students in particular can benefit from automated reporting of low-quality source code.

3.1.1. Architecture

The CCAP architecture is divided into the main part and two extension possibilities: An extension with analysis plugins adds automated checkers for more rules that are validated by the system. Adding an output plugin allows specifying the output format to fit custom workflow needs. The main part consists of four components: A core component to act as an orchestration unit, a component for handling the source code input and a component for handling the analysis plugins as well as output plugins. The design follows the requirements and goals for the platform. Figure 3.1 provides a high-level overview of the components.

The core component contains the main function and handles argument validation and parsing. Furthermore, it orchestrates other components by initialising and executing those. This process is divided into the argument parsing, initialisation and execution phase: In the argument parsing phase, the command line arguments are parsed and validated. It validates the existence of the required input directory argument and the optional plugin path configuration for the analysis plugins and the output plugins. Additionally, the logging level and the output format can be defined. The latter determines, which output plugin will be used, although the existence of the specified plugin is not validated in this phase. A parsing or validation error will cause a program termination without further processing. The initialisation phase instantiates all components and the

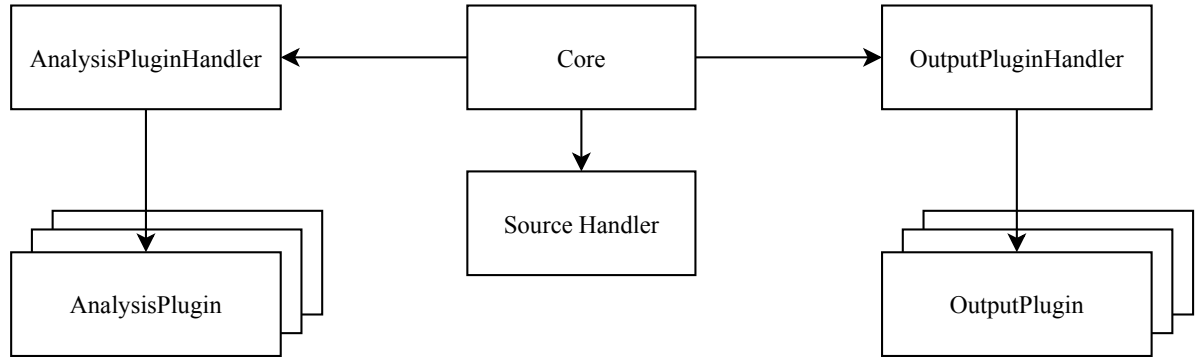


Figure 3.1.: High-level schematic overview of the CCAP. The core component is the orchestration unit, the source handler parse the source code files and the two plugin handler search for and execute the analysis and output plugins.

analysis plugin handler will scan the specified directories for plugins and keeps an index of all found plugins in memory. The output plugin manager scans for an output plugin, that satisfies the specified argument. If no plugin matches the output argument, the program will terminate and display a failure message to indicate the problem. In the execution phase, the input handling component scans the input directory for files ending with `.py` and parses the source code into an Abstract Syntax Tree (AST) per file. In the next step, the core passes the parsed data to the analysis plugin handler. The latter will execute all plugins for all files and collects the results. Afterwards, the core component calls the output plugin component to output the results. If no exceptions occur during the execution phase, the program will be terminated successfully.

The input component scans the given input directory for all Python source code files and parse the source code into an AST. For scanning the input directory, an algorithm will walk recursively over all folders and files. The detection of Python source code files is based on the file ending `.py`. The algorithm will return a list of file paths and the corresponding file content. Next, the AST parser is called and will add a parsed AST object to the list besides the file path and content. This list will be passed to the analysis plugins by the analysis plugin handler. An alternative approach would be to not read and parse the code in the input component, but instead, let the plugins read and parse the file content if needed. With many files to scan, the latter approach would have a lower memory footprint since the file content and the AST will not be held in memory. However, every analysis plugin has to perform an expensive read operation from disk and the performance scales with the number of files and the number of analysis plugins. If the input component reads all files, the information is held in the main memory and the performance only scales with the number of files. Since the files are text-based, we

expect the number of files needed to exceed the main memory to be high enough to fit most projects.

The analysis plugin handler manages all analysis plugins. This component finds all plugins, executes the plugins and collects the reported results. During the initialisation phase of the core component, the analysis plugin handler will scan the plugin directory for all Python files. It imports all Python files and scans those for classes, that inherit from an abstract **AbstractAnalysisPlugin** class. The abstract class defines all methods that need to be implemented in the specific plugin subclass. The analysis plugin handler instantiates all found classes. During the core execution phase, the core receives a list with the file name, file content and the parsed AST. All plugins are called on a specific entry point method that is defined in **AbstractAnalysisPlugin** with the aforementioned list. The plugin will return an instance of **AnalysisReport** with the plugins metadata and a collection of problems. The report is collected for every plugin into a **FullReport**. Additionally, the report contains information about the overall plugin execution time and run arguments. After all plugins have been executed for all files, the analysis plugin handler returns the **FullReport** to the core component.

The last component in the execution chain is the output plugin handler that will pass the **FullReport** to the specified output plugin. It implements the same algorithm as the analysis plugin handler to find all plugins that inherit from **AbstractOutputPlugin**. Instead of keeping track of all plugins, only the plugin that corresponds to the output format argument is instantiated. The output plugin has an entry point as defined in **AbstractOutputPlugin** that is called with all the collected results.

Analysis Plugins

Analysis plugins provide the easy extensibility of the platform by developers. All users of the tool can extend the set of problems it can detect by implementing a plugin in Python and placing it into the plugin directory of the tool. In order to be compatible with the core components, a plugin has to inherit from the **AbstractAnalysisPlugin** class. First, the abstract class introduces a class member variable **metadata** of the class **PluginMetaData**. A specific plugin class sets this class member in the constructor to provide a plugin name, author and optional website. This metadata is used in the output to show more information about the plugin that reported a problem. Second, **AbstractAnalysisPlugin** specifies a **do_analysis** method as an entry point that the platform will call. It accepts a list of **ParsedSourceFile** objects that contains the file path, the file content and the corresponding AST for all input files. With this

3. Approach

information, the plugin can implement any logic to detect problems. For example, the plugin can traverse the AST to look for specific node types, it can use a tokenizer on the content of the files and analyse the token stream or it can run sophisticated machine learning algorithms on the source code.

The plugin can even import third-party libraries, although the user has to install those on the system. After detecting all problems, the plugin returns an **AnalysisReport**. The report contains the plugins metadata and a list of found problems. These problems are instances of a problem class inheriting from **AbstractAnalysisProblem**. The abstract class expects a file path and line number as constructor arguments and requires the plugin developer to override the problem name and description. Listing 3.1 shows an example implementation. The problem name and description will be shown in the final output and should adhere to the following guidelines:

- Have a proper name that allows the experienced developer to recognise the problem quickly.
- Explain what code construct is problematic.
- Give reasons why this code is seen as problematic.
- Show guidance and examples on how to fix the problem and improve the code.

Although it is possible to use one plugin for multiple, different problem types, having one plugin for one problem type helps to reuse and share the plugin. Additionally, it can be disabled easily by removing the plugin from the plugin folder and complies to the Single-Responsibility Rule.

```
class ReturnNoneProblem(AbstractAnalysisProblem):
    def __init__(self, file_path, line_number):
        self.name = "Returned None"
        self.description = "Returning None is dangerous since the caller has to check for None
        . Otherwise, a runtime exception may occur."
        super().__init__(file_path, line_number)
```

Listing 3.1: Example for overwriting the **AbstractAnalysisProblem** with a specific implementation. The problem name and description are overwritten.

Steps to create an analysis plugin In order to extend the CCAP with an analysis plugin, the following steps are required for a developer:

1. Create a `.py` file with a class inheriting from `AbstractAnalysisPlugin`.
2. Instantiate `PluginMetaData` and assign it to the `metadata` member.
3. Define a problem class inheriting from `AbstractAnalysisProblem` and set the problem name and description following the guidelines above.
4. Implement the `do_analysis` method with a `ParsedSourceFile` parameter. Return an `AnalysisReport` instance with all found problems.
5. Place the `.py` file into the analysis plugin directory of the tool.

While these are the minimum required steps to implement the plugin, the developer is free to add additional methods, classes or import libraries as necessary. Furthermore, it is advisable to implement several tests to ensure the plugins correctness. CCAP uses the `pytest`¹ library for testing. Implementing a test is as easy as writing a function with a `“test_”` prefix and running the `pytest` command. The plugin can be tested by simulating a call from the CCAP core with a mocked `ParsedSourceFile`.

3.1.2. Return None Plugin

A rather simple analysis plugin demonstrates the capabilities of the CCAP: The Return None Plugin scans the source code for functions that return the `None` value. Returning `textttNone` may result in runtime exceptions if the function caller does not expect and handle a potential `None` return. Although `None` can be returned directly or as a value of the returning variable, this plugin only focus on the explicit `return None` statement to showcase the options for the developer to analyse the source code. We will later refer to this problem type as `RETURN_NONE` (RN).

Detecting a `return None` statement is possible in multiple ways. The following shows the possibilities for a developer to implement such a detector:

Regular Expression: In the `do_analysis` method, a developer has access to the source code as a string. Therefore, it is straightforward to use a regular expression to detect a `return None` statement:

¹<https://docs.pytest.org/en/stable/>

3. Approach

```
import re
matches = re.finditer(r"return None", source_file.content, re.MULTILINE | re.DOTALL)
```

Since the regex library only returns the start and end index of matches, these have to be converted to line numbers. Afterwards a `ReturnNoneProblem` instance can be created for every match and added to the `AnalysisReport` for this plugin.

This approach uses regular expressions to match patterns in a string without utilising the structures of the code. It is a simple but powerful way and most developers are familiar with regular expressions. On the flip-side, source code may have various syntactic ways to express a semantic. Consequently, the regular expressions have to be designed carefully to cover all variations. For instance, it is possible to encounter a doubled whitespace, that would break the aforementioned regular expression. Additionally, regular expressions do not operate on the structure of the code; therefore, they can not detect high-level patterns on the code structure.

Tokenization: The process of dividing a character stream into a sequence of tokens is called tokenization (also known as lexical analysis). With the Python tokenizer, a token can contain multiple characters and has a token type like name, operator or indentation. A token sequence provides more information about the code structure that can be used to detect problematic patterns. A token sequence for a `return None` statement would be the following:

```
[...( type: "NAME", value: "return"), (type: "NAME", value: "None")...]
```

A simple algorithm would scan the sequence for two subsequent name tokens with the values `return` and `None`. The abstraction level of a token sequence is higher than of a character sequence. Since the whitespace between `return` and `None` provides no semantic meaning, it is removed on this abstraction level. The regular expression-based approach would have to deal with problems as the doubled whitespace. In contrast, the token-based approach profits from the higher abstraction level and can access meaningful tokens directly.

Abstract Syntax Tree : After tokenization, a parser takes the token stream and parses it into a hierarchical data structure like an Abstract Syntax Tree (AST). With an AST, the code is represented structurally and it is possible to traverse the tree following the structure of the code. The AST consists of nodes with different types and children. For instance, a node represents an `if` expression and has a `test`,

```
for node in ast.walk(a.ast):
    if isinstance(node, ast.Return):
        return_value = node.value
        if isinstance(return_value, ast.Constant) or isinstance(return_value, ast.NameConstant):
            if return_value.value is None:
                problems.append(ReturnNullProblem(a.file_path, return_value.lineno))
```

Listing 3.2: Detecting a `return None` by analysing the AST data structure.

body and `orelse` children. It is possible to access the condition (`test` children) and the body of the if statement as descendant nodes in the AST. Analysing an AST allows a higher abstraction level than analysing the token string since the AST represents the code structure. Therefore, automated checkers for more abstract rules that analyse the code structure are possible.

To detect a `return None` statement in the abstract syntax tree, an algorithm would traverse all AST nodes looking for a return-typed node. If found, the `value` descendant contains the expression after the `return` keyword. A `None` value would be represented as a node of type `Constant` or `NameConstant`, depending on the parser version. The `value` descendant of the node may then be checked for equality to `None`. All AST nodes contain the corresponding line number that is necessary for creating a problem message. See Listing 3.2 for a Python implementation.

Analysing the AST works best to find problems in the code structure since the AST represents the code structure in a well-defined, traversable data structure. Although simpler patterns like the `return None` could be detected using regular expressions, a detection on AST level allows detecting the semantic meaning instead of the syntactic representation of a problem. For more complex problems, it is inevitable to use the AST most of the time.

This implementation covers only the simple case, in which `return None` is written explicitly. Of course, several variations are possible, that would not be detected by this plugin and would require more sophisticated algorithms. Evaluating, if machine learning models trained on this simple rule can also detect such modifications is part of RQ3.

Since detection alone is not a great help for the user, a useful description of the problem is necessary. As described in Section 2.2.6, there are some alternatives to returning `none`. If returning `None` indicates an error, it is better to raise an exception and provoke an explicit error handling in order to prevent runtime errors. If the standard return type is a collection, an empty list is more appropriate, since the function caller will most likely write the logic for an unknown amount of list items. However, if the standard return

type is a single object, it is not apparent what to return. With PEP484, an optional type was introduced as a type hint in Python 3.5 [42]. A static type checker for Python like mypy outputs a warning if a developer forgets to check an optional for not being `None` [23].

Condition with Comparison The second plugin detects a direct comparison in a conditional statement. For better readability and understandability, the clean code rules suggest to use a function call with an appropriate name that evaluates to `true` or `false`. This function call allows a natural reading of the conditional statement without deciphering the meaning of the boolean logic.

An if-statement consists of a condition, a body and an else part. If the condition evaluates to `true`, the body is executed. Otherwise, the else part is executed. Any logical expression or a function call will be evaluated in the condition. A simple algorithm to detect a comparison in the condition would take the AST, search for `if` nodes, and check if the condition part is a compare node. However, the simple algorithm would not detect a negation and logic AND/OR conjunction. Consequently, a recursive algorithm has to follow logic operators and checks the expressions for comparison nodes. The simple approach does not provide value for the CCAP, but we use it for the machine learning evaluation. We later refer to the code detected with this algorithm as `CONDITION_COMPARISON_SIMPLE`.

Following these considerations, an advanced algorithm would scan for `if` nodes in the AST. The conditional part as a boolean expression is evaluated in a recursive. It can be a boolean operation (like AND/OR) or a unary operator (like NOT). In these cases, the function will be called recursively with the respective expressions. If the conditional part is neither a boolean operation nor a unary operator, the function returns `true` if a single expression is not a function or method call. The clean code rules recommend a function or method call instead of a comparison to maintain readability. The last case would be the base case in the recursion. Listing 3.3 shows the Python implementation of this algorithm and the detected code patterns. We will later refer to the problem type of this advanced algorithm as `CONDITION_COMPARISON (CC)`.

Output Plugins

With output plugins, CCAP adds additional flexibility towards the output format. Depending on environmental requirements, the output can be adapted with custom logic by defining an output plugin. For instance, running the tool locally could write the


```
def _check_if_direct_comparison(self, node):
    if isinstance(node, ast.BoolOp):
        violated = False
        #check all expressions of the boolean operator
        for value in node.values:
            if self._check_if_direct_comparison(value):
                violated = True
        return violated
    elif isinstance(node, ast.UnaryOp):
        return self._check_if_direct_comparison(node.operand)

    return not isinstance(node, ast.Call)
```

Listing 3.3: Recursive function to analyse an if statement for direct comparisons. Since a condition should contain a method call, the function returns False if this is not the case.

results to the standard output in a human-readable way or create a formatted HTML file to be displayed in the browser. A machine-readable JSON output may be preferred when running inside an automated workflow.

Output plugins follow similar concepts as analysis plugins. A plugin class inherits from `AbstractOutputPlugin`. The abstract class introduces the metadata member of the `PluginMetaData` class to encapsulate plugin name and author information. Additionally, a second member variable `output_format` has to be defined with a short abbreviation for the output format. This field is used by the output plugin handler to select the output plugin by the run arguments. It should be globally unique, or the first found plugin will precede other plugins with equal naming. We suggest using a custom prefix to the format name to maintain global uniqueness.

As an entry point, the method `handle_report` has to be implemented. The method provides an instance of *FullReport* as its argument. The `report` field holds a collection of `AnalysisReport` for every analysis plugin that was executed. With metadata and problem information, the output plugin has access to all required information to produce the desired output.

Steps to create an output plugin To extend the output capabilities of the CCAP, the following steps are, similar to analysis plugins, required:

1. Create a `.py` file with a class inheriting from `AbstractOutputPlugin`.
2. Instantiate `PluginMetaData` and assign it to the metadata member.

3. Approach

```
Analysis Report on /Users/enrico/MA_CleanCodeAnalyser/test_programs.
Analyse Plugins: Condition Method Call Plugin, Simple Condition Method Call Plugin, Return None
  (Null) Plugin, Sample Plugin.
Summary: Found 16 problem(s)
-----
PROBLEMS:
  PLUGIN NAME: Condition Method Call Plugin by Enrico Kaack <e.kaack@live.de>
    Found: Explicit comparison in condition in /Users/enrico/MA_CleanCodeAnalyser/
    test_programs/return_none.py:18
    Explicit comparisons in conditions should be replaced by method call for better
    readability
  ...
  PLUGIN NAME: Return None (Null) Plugin by Enrico Kaack <e.kaack@live.de>
    Found: Returned None in /Users/enrico/MA_CleanCodeAnalyser/test_programs/
    return_none.py:4
    Returning None is dangerous since the caller has to check for None. Otherwise, a
    runtime exception may occur.
```

Listing 3.4: Example output to `stdout` of the Standard Output Plugin. Besides a list of problems, it also outputs additional

3. Set the `output_format` field with a unique abbreviation for this output format. Since it should be unique, a custom prefix prevents a name collision with pre-existing plugins.
4. Implement the `handle_report` method with a `ParsedSourceFile` parameter.
5. Place the `.py` file into the output plugin directory of the tool.

Standard Output Plugin The Standard Output Plugin writes the formatted output to the `stdout` stream. It is enabled by default if no output plugin is explicitly specified.

The output is divided into a general part and the problem list. The former contains the input path, all executed plugins, and a summary field with the number of problems in total. The latter displays the list with problems, grouped by analysis plugin. For each problem, the problem name, file path, line number and description is printed. A colon separates the file path and line number. Some terminals parse the path and line number so that a user can click the path and it opens the default editor with the cursor in the correct line. A sample output is printed in Listing 3.4.

HTML Output Plugin For larger projects, the Standard Output Plugin may not be sufficient to get an overview of the project or to generate a report. Therefore, the HTML Output Plugin creates a `report.html` file that can be opened in every browser.

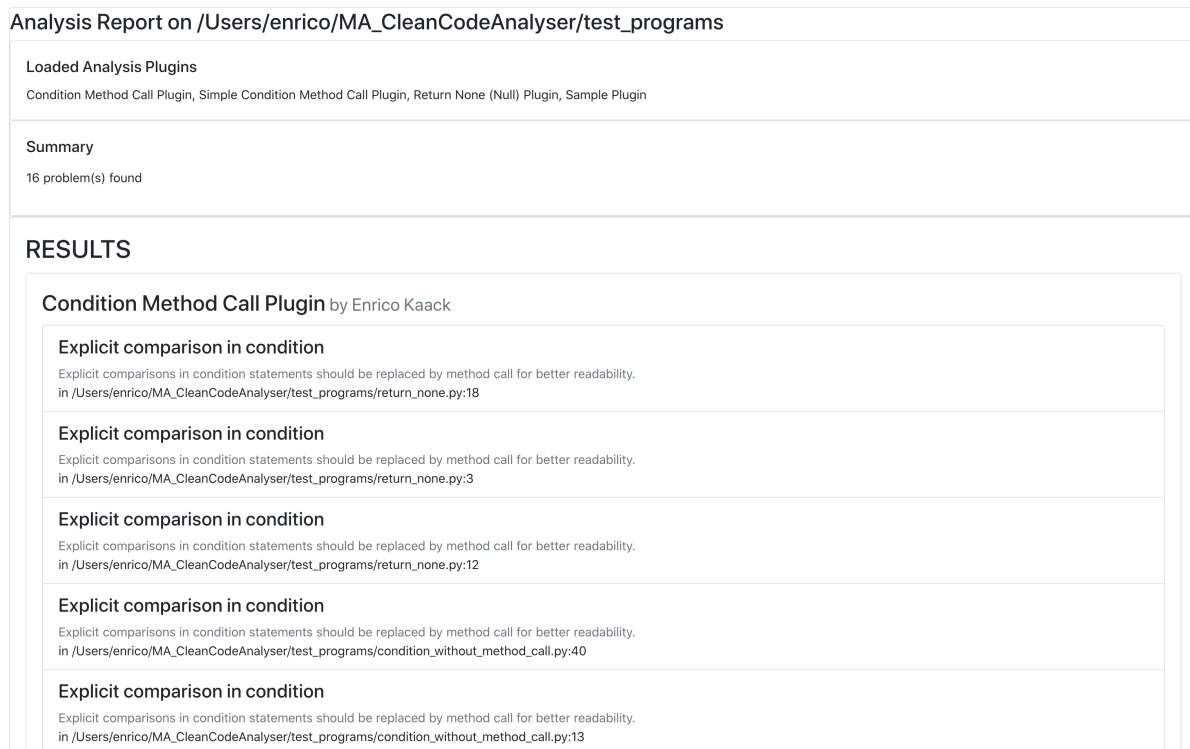


Figure 3.2.: Output of the HTML Output Plugin, displayed in a browser.

The report provides a cleaner overview and the scrolling capabilities of a browser. It has similar general data at the top as the Standard Output Plugin, but it moves the problem description into a tooltip to offer a more compact overview. See Figure 3.2 for an example screenshot.

3.1.3. Improvements

While the core functionality exists, extensions for the CCAP would aim for improved useability. One crucial improvement in useability would be an IDE integration for common IDEs like Visual Studio Code. A good starting point would be the Language Server Protocol². This IDE integration requires a language client as a Visual Studio Code extension and a protocol-compliant language server. The latter would wrap the CCAP and modify it slightly to process a single, changed document and return the results. Since the CCAP architecture is modular, this should be possible without larger modifications.

Another feature would be a configuration to disable specific analysis plugins for a

²<https://microsoft.github.io/language-server-protocol/>

project. With the current architecture, this would be possible by moving the analysis plugin file outside the plugin directory. A configuration with command-line parameters or with a hidden config file per project would have a better user experience. The latter could also be committed to the version control system, so every team member has the same configuration.

Continuing on configuration possibilities, having configuration possibilities for plugins would increase the flexibility for a plugin developer. Introducing configurable warn levels like **warning** or **error** could help in continuous integration pipelines to decide if a build succeeds with warnings or if a build should fail because of error-level problems. Some rule violation could be seen as recommendations to improve (warning level), whereas other violations may be unacceptable (error level).

Last, a feature to disable problem reporting for a specific code location could bring a boost in user acceptance. While some clean code rules are objective and can be measured precisely, some rules may not apply to every occurrence of the situation. For instance, the clean code guidelines generally suggest not to have more than three function arguments; it may be necessary or even inevitable to have four arguments. The CCAP would report this as a problem, but the user should be able to decide if it is acceptable. In this case, the problem type on this specific location should not be reported in the future. In the current version, the user has no way to ignore a specific problem. Consequently, the number of problems the user has to ignore will accumulate over time until the user abandons the tool since it does not provide an added value over the frustration of manually ignoring problems.

3.2. Clean Code Classification

In the previous chapter, we presented a platform to check for clean code rules. This platform requires analysis plugins to detect specific problems. These handwritten rules work well for rules with a basic level of complexity. However, certain rules of an advanced level of complexity may require a complicated analysis of the AST, the data-flow or the inter-module dependencies. Furthermore, rules with a high level of complexity can be subjective, and we do not see a way of detecting those rule violations with a deterministic algorithm.

Therefore, this chapter introduces an approach to evaluate different supervised machine learning models to classify code into clean code or problematic code. Furthermore, we describe how we modify the code to test the generalisation capabilities of the models.

The mentioned machine learning models include random forest classifier, support vector machines, gradient boosting classifier and a recurrent neural network model based on Long Short Term Memory units.

The objectives are twofold: For RQ2, we train, evaluate and compare the different models on the detection of different code violations. Then we manipulate the source code to test in RQ3, how well the models generalise to non-seen variations of the rule violations.

To achieve our objectives, we will train our models on the detection of three problem types:

1. The `RETURN_NONE` (RN) problem type covers a `return None` statement in Python code (see Section 3.1.2). A detector would report such a statement as this problem type.
2. A detector for the `CONDITION_COMPARISON_SIMPLE` (CCS) problem types detects a comparison operator in the condition of an if-statement. It does not detect the comparison if it is nested with logical expressions. The detection algorithm is further described in Section 3.1.2.
3. With the `CONDITION_COMPARISON` (CC) type, we describe whether the condition inside an if-statement contains other expressions than boolean logic and function calls. The CC type is a superset of the CCS problems.

From the three problem types, we think the RN type is the easiest type to spot since the model only has to spot the same statement in every sample. The CC type is more complex to spot and the model would have to learn the connection between the if-statement and the possible logic expressions. The same applies to the CCS type, that should be easier to detect than the CC type due to fewer variations. Additionally, we use the trained model on the CCS subset and test it with a manipulated CC superset to assess the generalisation to unseen variations in RQ3.

In the following, we describe the approach in detail. First, we describe the challenges and our solutions. Second, we describe our dataset and our split in train, test and holdout set. Afterwards, we describe our pipeline of processing the raw files into samples with a fixed size, our encoding approach and our code manipulation for RQ3. In the final section, we describe the machine learning models and parameters we use for all research questions.

3.2.1. Challenges

Lack of datasets The most crucial preconditions of supervised machine learning experiments are labelled datasets. For clean code detection, we did not find a suitable dataset with labelled clean code violations in our research. As a result, we evaluated different approaches to create such a dataset:

Solution 1: Related research fields like code completion often use the py150 dataset to train models [32]. This dataset contains 150.000 python files, sourced from GitHub. For our supervised classification tasks, this dataset is inadequate due to its missing labels.

Solution 2: Another possible data source could be git commit histories. We could scan the history for commits that fixed unclean code. The previous code could then be labelled “unclean” and the committed code would be “clean”. Same applies to issues and referenced pull requests on GitHub. We dismissed this approach for several reasons:

1. There is no standardised annotation in neither issues descriptions nor git commit messages that would indicate reliably if the code change is fixing a violation of the clean code rules.
2. Even if we would be able to find commits that improve chaotic code, we could not ensure automatically, that the commit message is correct and only the improvement is included in the commit.
3. Searching for commits, we found several clean code improvements and refactorings of large portions of the codebase in one commit. Additionally, it is often not explicitly mentioned, which rule applied for the improvement. Consequently, training with such data could only lead to binary classification into clean code or unclean code. From a practical standpoint, it is advantageous to name the explicit rule violation and to explain how to improve.

Figure 3.3 offers an illustrative example of commit messages from different projects. First, there is no consistent annotation schema to display clean code optimisation. Second, a commit may include multiple changes like the second example. Finally, the first commit shows a large refactoring that most likely fixes several problems and not only clean code violations.

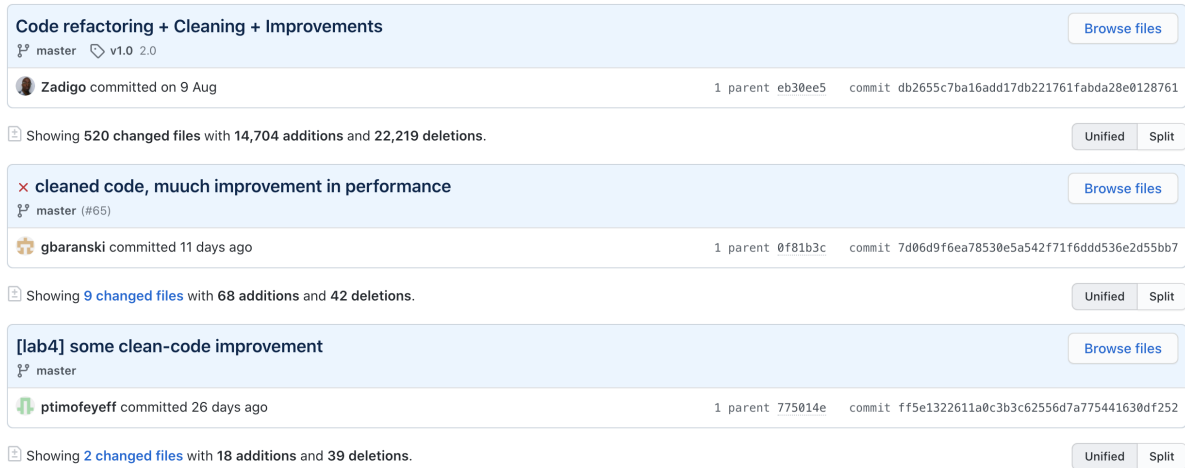


Figure 3.3.: Example commit messages from different repositories. Those examples highlight (1) the missing naming convention, (2) inclusion of additional performance improvements, and (3) several edits in multiple files [30, 7, 41]. Those commits were found by searching GitHub for “clean code improvements” and filtering for results of type commits and language Python.

Solution 3: Since Python code is available in abundance in the form of open-source projects, it is possible to manually label source code that violates clean code rules. For the following reasons, we dismissed this approach:

1. Curating a hand-labelled dataset is a time-consuming task.
2. Manual labelling does not ensure correct labelling. A human can make mistakes or can subjectively misinterpret chaotic code as clean code. The quality of the dataset could suffer and limit the machine learning models performance.

Solution 4: Given any amount of Python files, we could use the analysis plugins from Section 3.1.1 to reliably detect violations of the corresponding rule. Based on this labelled data, we could train the classifier to detect those rule violations. The analysis plugins only cover rule violations of basic level of complexity, so the trained classifier could only detect rule violations that could also be detected by a basic checker. However, if there would be a labelled dataset for more complex clean code rules, we could transfer the findings to these more complex rules.

We choose this solution to solve the dataset problem and train classifier. We can automatically generate labelled data with correct and consistent labels. Furthermore, this solution allows for scaling the dataset size if necessary. We will use the analysis plugins to label the three problem types: `RETURN NONE (RN)`, `CONDITION COMPARISON`

SIMPLE (CCS) and **CONDITION COMPARISON** (CC). However, this implies that we train the classifier to detect rule violations that could be detected by simple analysis plugins. If a labelled dataset for rules with a high level of complexity becomes available, we hope to transfer our findings to those rule violations as well. In RQ3, we will then test those classifiers on unseen variations of the rule violations to test the generalisation of the models. This gives an outlook on the transfer the training on rules with a high-level complexity by simulating a hand-labelled training set that is limited in problem variations.

Imbalanced Dataset A second challenge is a major imbalance between the samples labelled as problematic code and clean code. For the three different problem types, RN, CCS, and CC, the class distribution is shown in Table 3.1. The imbalance is especially severe for the RN type since only two out of 1000 samples contain the rule violation. For the CCS and CC type, the class are slightly more balanced with 26 and 49 samples with rule violations out of 1000.

We combined the following solutions:

Solution 1: Due to the imbalance, we do not choose accuracy as a suitable evaluation metric, since classifying all samples as clean code would result in an accuracy of 99.79% for the RN problem type. Instead, we use precision and recall as metrics, with an f1 score as a single, weighted metric. More details about the metrics are in Section 4.1.2.

Solution 2: We will apply undersampling and oversampling techniques on the data and evaluate the difference in model performance. Undersampling removes random samples from the majority class and thus balance class labels. Oversampling replicates random samples from the minority class to balance class labels while also increasing the overall amount of data. This technique is only used for the training data since it would distort the evaluation results when applied to test data as well.

Label	RN	CCS	CC
0	99.79%	97.36%	95.14%
1	0.21%	2.64%	4.86%

Table 3.1.: Class distribution for the dataset. Label [0] represents clean code and label [1] marks problematic code of the corresponding category.

Solution 3: Our approach of training different classifier is a solution to the data imbalance since different classifiers have a different sensitivity to data balance.

SVM quadratic complexity in training The implementation of the support vector machine has a quadratic time complexity on the number of training data [3]. Consequently, the training takes much time and the oversampling strategy worsen the training time further since it introduces more samples.

Solution: We reduce the number of experiments with SVM and just evaluate some simple cases for baseline performance.

3.2.2. Dataset

The source of our dataset are open-source Python projects on GitHub. We queried the top starred Python repositories and handselected 18 repositories. Although most top starred repositories are data science frameworks, we explicitly choose projects from other domains such as web server, automation and containerisation. In Table 3.2, we list all projects with the corresponding domain.

A script downloaded all projects in their main branch with the current head. See Table A1 for the corresponding git hashes. Afterwards, we removed all non-python files. Next, we uniformly sampled 20% of all files and separated those into a holdout set for testing and the code manipulation in RQ3. Additionally, we perform our train/test split on the file level and not on the sample level. We describe the reason for the file level split later in Section 3.2.3.

We ensured to have similar data and label distribution on all data sets. Table 3.3 shows general and problem-specific metrics on file level for the train, test and holdout set. The problem-specific metrics are collected after processing the code files as described in Section 3.2.3. We define the metrics as follows:

We calculate the *average lines of code per file* for n files with l_i as the lines of code for file i with the following equation:

$$\text{average LOC per File} = \frac{\sum_{i=0}^n l_i}{n}$$

For the *proportion of lines of code containing a problem*, we assume one line contains the problem. Therefore, we define this metric for n files with l_i lines of code for the i th

3. Approach

Organisation	Project	Domain
scikit-learn	scikit-learn	Data Science
pytorch	pytorch	Data Science
explosion	spaCy	Data Science
bokeh	bokeh	Data Science
pandas-dev	pandas	Data Science
keras-team	keras	Data Science
django	django	Web Server
pallets	flask	Web Server
tiangolo	fastapi	Web Server
tornadoweb	tornado	Web Server
encode	django-rest-framework	Web Server
jakubroztocil	httpie	Web Requests
home-assistant	core	Home Automation
certbot	certbot	Automation
google	python-fire	Automation
apache	airflow	Automation
ansible	ansible	Automation
docker	compose	Containerization

Table 3.2.: Open-source repositories we used in our dataset and their corresponding domain.

file and p_i problems in the i th file as follows:

$$LOCs \text{ containing problem} = \frac{\sum_{i=0}^n p_i}{\sum_{i=0}^n l_i}$$

Last, we calculate the number of problems per file for n files and p_i problems in the i th file as:

$$Problems \text{ per File} = \frac{\sum_{i=0}^n p_i}{n}$$

The most important metrics for the data distribution is the proportion of lines of code containing the problem type and the average number of problems per file. For lines of code containing the problem type, we observe a maximal difference of 0.02 percentage points between the datasets for the RN problem type, 0.11 percentage points for CCS and 0.13 percentage points for the CC problem type. This directly translates into the class distribution as shown in Table 3.4 with a maximal difference 0.06, 0.26 and 0.31 percentage points for the three problem types. The number of problems per file is lower in the test set for the CCS and CC problem type. Nevertheless, the size of a file may

	Metric	Train	Test	Holdout
	Lines of Code	2,530,455	246,813	671,554
	Number of Files	13,330	1,481	3,702
	average LOC per File	189.83	166.65	181.40
Return None	LOCs containing problem	0.07%	0.09%	0.08%
	Problems per File	0.14	0.15	0.14
Condition Comparison Simple	LOCs containing problem	1.29%	1.21%	1.32%
	Problems per File	2.44	2.02	2.4
Condition Comparison	LOCs containing problem	2.34%	2.31%	2.44%
	Problems per File	4.44	3.85	4.42

Table 3.3.: General and problem specific metrics for the train, test and holdout set.

Problem Type	Label	Train	Test	Holdout
Return None	[0]	99.79%	99.73%	99.78%
	[1]	0.21%	0.27%	0.22%
Condition Comparison Simple	[0]	97.34%	97.49%	97.23%
	[1]	2.66%	2.51%	2.77%
Condition Comparison	[0]	95.13%	95.21%	94.9%
	[1]	4.87%	4.79%	5.1%

Table 3.4.: Class frequencies for all problem types on the train, test and holdout set.

We observe a comparable class distribution among the datasets for each problem type.

vary, and thus the number of problems per file may vary too. Since the lines of code containing the problem and the class frequencies are comparable, our datasets have a comparable data distribution.

TODO per-project analysis table in attachments?

3.2.3. Processing

For the processing steps, we create a pipeline with the d6tflow framework³. Each processing step is a task that stores its results depending on the parameter configuration in a pickle file. This allows to define a pipeline once and run it with different parameter combinations. The scheduler automatically determines what tasks to run and what task

³<https://github.com/d6t/d6tflow>

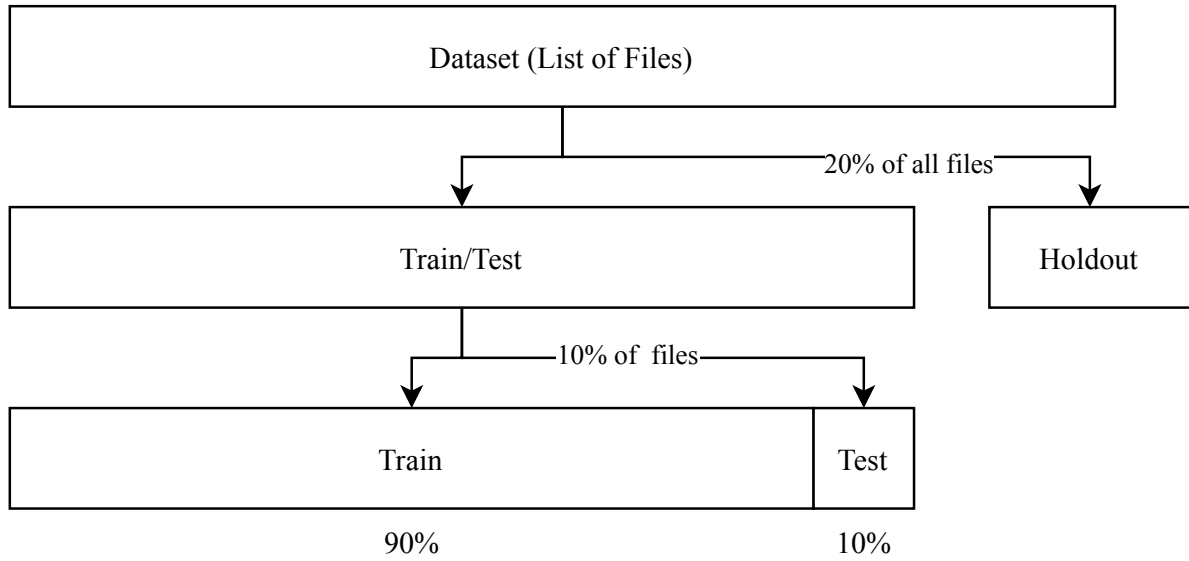


Figure 3.4.: Split of the dataset. First, 20% of all files are separated into a holdout set for later testing. The remaining 80% of the files are split into 90% training and 10% test data.

outputs are already stored. Consequently, running the pipeline is more efficient and repeatable.

The processing pipeline consists of tasks to read the files into a data structure, to process source code and find the problems, to create a vocab dictionary, to convert the data into labelled samples. In Figure 3.5, we provide a schematic representation of the pipeline.

Files into Internal Data Structure

First, a scanner walks recursively over all subdirectories of the input folder to find all files. If the file has a `.py` extension, its file path and the file content will be stored in a dictionary. The dictionaries for all files are collected into a list. Reading all files into main memory increases the performance for downstream tasks since main memory access is faster than disk access. Although the size of the system's main memory limits the dataset size, a file contains text and is just several kilobytes in size. The train dataset with 13,330 files is 86.47MB in accumulated file size.

Problem Detection

As a next step, the analysis plugins from the CCAP (see Section 3.1.1) process every file and store the line number along with the problem type (RN, CCS, or CC). As a result,

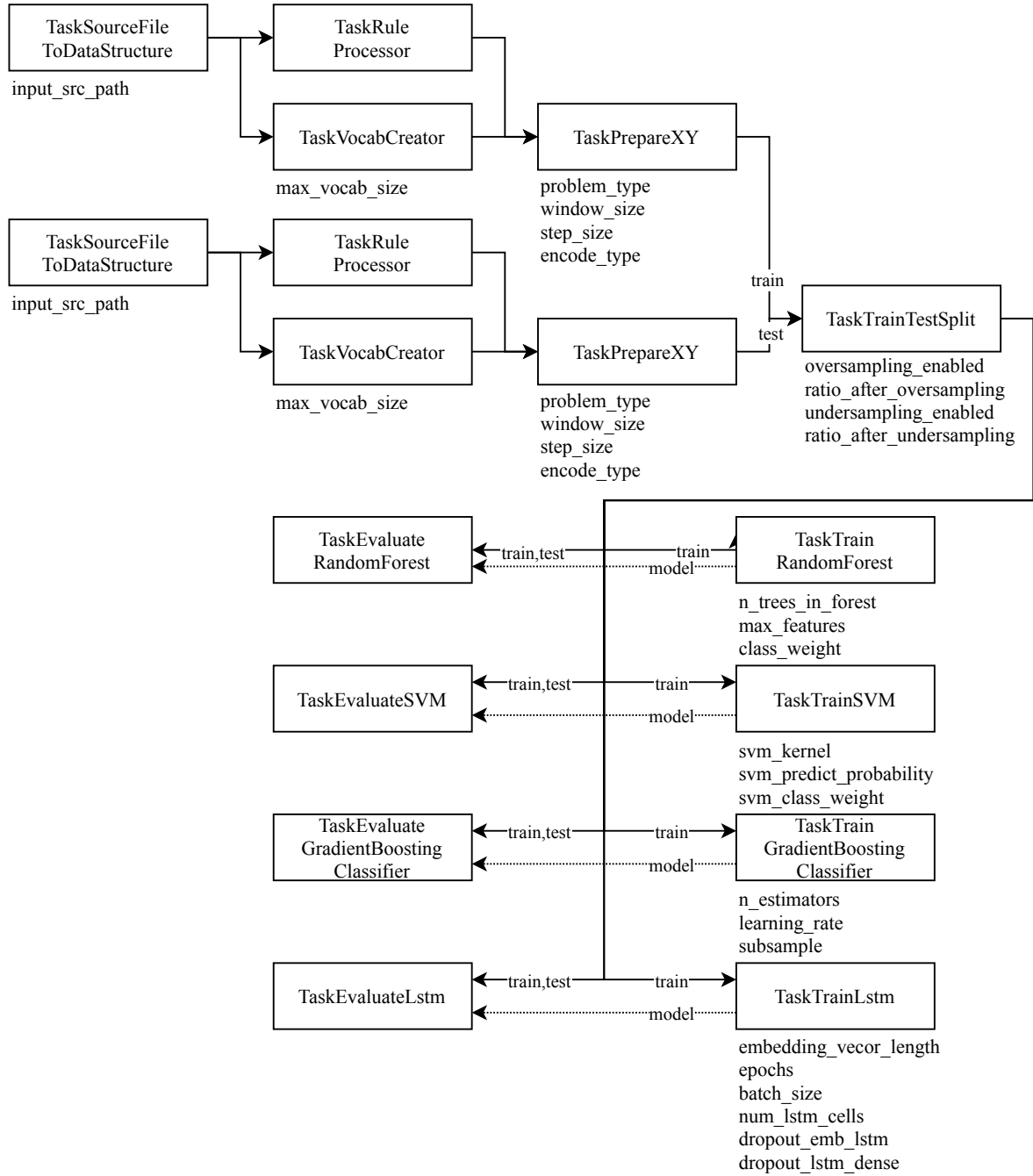


Figure 3.5.: Schematic representation of the pipeline for RQ2 to train and test models with configurable parameters. First, the pipeline reads the source code into memory, extract the vocabulary (`TaskVocabCreator`) and identify all problem types (`TaskRuleProcessor`). It then converts the code of dynamic length into fixed-size samples (`TaskPrepareXY`). This sub-pipeline runs separately for the train and test set. In `TaskTrainTestSplit`, both datasets are combined and over- or undersampling is applied to the train set. The train set is used to train different classifiers whereas the train and test set are used for evaluation.

for every file, a list of problematic line numbers and the corresponding type is available for further processing.

Data Encoding

The data encoding step transforms the internal data representation into the input vector x and output vector y . One (x, y) pair describes a sample with input and ground-truth label. The transformation happens in the following multi-stage process:

Fixed Length Sequence The length of the source code is dynamic, whereas the input size of our models is fixed. Therefore, the character stream of variable length has to be transformed into a token stream of fixed size. To extract meaningful tokens from the character stream, we use the Python tokenizer from its standard library. The tokenizer separates the character stream based on the python syntax definition into tokens. All tokens contain a token type (like a name or operator token), the corresponding characters in the source code and a start and end position (line and column number). Afterwards, the token stream is divided into a fixed size token sequence using a sliding window approach. The step and window size are configurable, and we use a window size of 20 and a step size of 3 in our experiments. The sliding window operation is part of the `TaskPrepareXY` step in the pipeline.

Vocabulary Creation In order to feed the textual token values into a model, we need to transform it into a numeric representation. Since our token values can be the entire value of a string literal, we can not encode every single token value into a number. We create a vocabulary with a limited size that contains the most common token values and a numerical representation for all other values. Therefore, the occurrence of each token value is counted and sorted based on their frequency. To ignore potential capitalisation mismatches, all token values are lower-case. The overall size of the vocabulary is configurable. If the size is smaller than the size of the distinct token value set, the least common token values will be replaced by an unknown token. The unknown token has a numeric representation one bigger than the vocabulary size. We found that using a vocabulary size of 100,000 results in an acceptably low 3.9% of token values being unknown. In our pipeline, the task `TaskVocabCreator` creates the vocabulary.

Index-Based Encoding The textual token values are encoded with an Index-Based encoding. The most common token values have a numerical representation based on

their index in a vocabulary (created with all token values of the train data). Optional, the type of each token is encoded by its numerical representation in the standard library (see Figure C.1). The type and value encodings are then combined alternating in a final encoding vector x_i for every sample i .

Label Extraction With the internal representation, the ground truth is encoded as a line number. This representation has to be transformed in a label y_i for the given input vector x_i for sample i . Label 1 is assigned if the sample contains the given problem type. A sample contains a problem type if it contains all tokens from the problematic line. If the sample is missing a token of the problematic line at the beginning or end, the label is 0 for non-problematic code. After the label extraction in the pipeline step **TaskPrepareXY**, we have a list of input vectors x and a list of ground truth labels y . With this, we can train and test our models.

Train/Test Split

In Section 3.2.2, we described to create a train, test and holdout dataset by file-level separation and not by splitting the samples. The reason for this approach is our sliding window concept. Since we use the sliding window to convert source files with dynamic length into fixed-length token sample samples, we may have one problematic code pattern in multiple samples. The difference would be in the location inside the sample and the surrounding tokens, whereas the pattern remains the same. If we would do a train/test split on a sample level, we may introduce samples covering one problematic pattern in both datasets. Figure 3.6 illustrates multiple windows covering one sample.

This unclean separation would diminish the validity of metrics calculated from the test and holdout dataset. Therefore, we split on file level and perform the previous preprocessing separately for training and test dataset. Only the vocabulary generation happens on the training file for both processing steps. In **TaskPrepareXY** in the pipeline, we then merge those two pipelines for training and testing (see Figure 3.5).

Code Manipulation

In research question RQ3, we evaluate the models' generalisation to detect similar pattern the model was not trained on. Therefore, we manipulate the code that it is still a rule violation, but the original analysis plugins would not spot them. The pipeline for code manipulation differs from the pipeline for RQ2 by an additional code manipulation task and removal of the training step (see Figure 3.7).

3. Approach

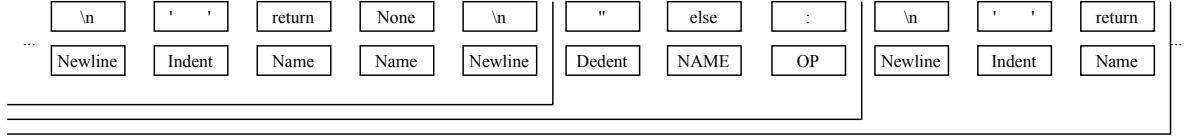


Figure 3.6.: Sliding window approach on source code containing a problematic `return None`. Multiple windows will contain the problematic sample. Consequently, we can not split on a sample level, since the same problem would be in the train and test set. Instead, we split on file level to have a clean separation between the datasets.

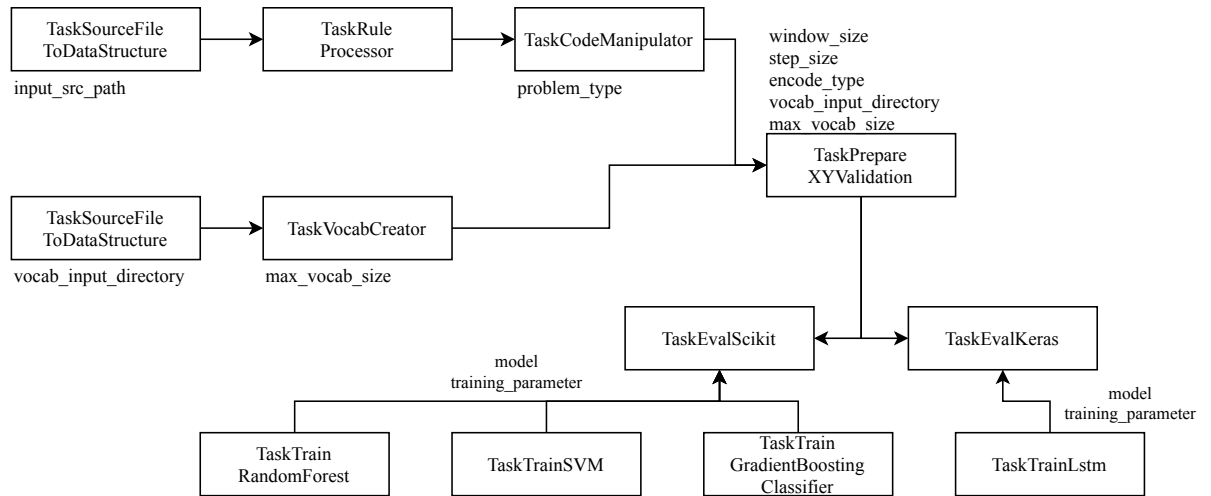


Figure 3.7.: Pipeline for code manipulation and evaluation. The preprocessing part is like in Figure 3.5 but with an additional *TaskCodeManipulator* to manipulate the code as described in Section 3.2.3. The random forest classifier, svm and gradient boosting classifier require different evaluation code as the LSTM, so two *TaskEval** are used.

Return None For the RN problem type, we manipulate the code to have an inline if statement that only returns `None` in one branch. Therefore, a function may still return `None`, but the analysis plugin would not detect it. Listing 3.5 shows an example for this case. Although a modification of the analysis plugin could cover the variations as well, we want to see how well machine learning models could perform on this task.

To modify the original code, we use the preprocessed code with problems detected as in Section 3.2.3. For every source code line containing the RN problem, we use regular expressions to replace the `return None` with a variation as seen in Listing 3.5.

The subsequent data encoding step is similar to the one for model training, although this modified samples will only be used during evaluation.

```
def f(a,b):
    #detected by the analysis plugin
    return None

    #not detected by the analysis plugin
    return None if a < b else b

    #not detected by the analysis plugin
    return a if a < b else None
```

Listing 3.5: Samples for returning `None`. The analysis plugin would flag the first return; the second and third return are modified variations that would be ignored by the analysis plugin. The performance of the machine learning models on detecting the latter will be evaluated.

Condition Comparison We want to test the generalisation of the model trained on the CCS types by evaluating on the labelled data for the CC type. Since the dataset for the CC type includes the CCS subset, we manipulate the CC dataset to only contain samples the CCS was not trained on. The sample in Listing 3.6 shows the manipulation to ensure that all CC problems were not part of the CCS training set. The implementation uses a similar regular expression approach as described in Section 3.2.3

```
def f(a,b):
    #detected by CONDITION COMPARISON SIMPLE
    #detected by CONDITION COMPARISON
    if a < b:
        pass

    #manipulated
```

```
#not detected by CONDITION COMPARISON SIMPLE
#detected by CONDITION COMPARISON
if not a < b:
    pass

#manipulated
#not detected by CONDITION COMPARISON SIMPLE
#detected by CONDITION COMPARISON
if is_smaller(a,b) and b < c::
    pass
```

Listing 3.6: Sample statements for the difference between the two analysis plugins CC and CCS.

3.2.4. Models

Classifying code samples is a binary classification problem since we consider each problem type as a separate classification task. We train and evaluate a random forest classifier, a support-vector machine-based classifier, a gradient boosting classifier and a neural network with LSTM cells.

Random Forest We use a random forest classifier with 100 decision trees to perform the binary classification. Additional to over- and undersampling, we experiment with an automatic class weighting inverse proportional to class frequency. The other hyperparameter follows the default implementation of the scikit-learn library: The quality measure of a split follows the Gini function, the tree depth is unlimited and the number of features to consider for each split is the square-root of the number of features overall. We use the implementation in the scikit-learn library in version 0.23.2 [29].

Support Vector Machine We use the support-vector classification from the scikit-learn library [29]. As a kernel function, we choose the radial basis function for a performance baseline. Additionally, we try the class weighting inverse proportional to the class frequency to combat the imbalanced dataset.

Gradient Boosting Classifier For gradient boosting classifier, we mainly vary the number of boosting steps and the learning rate. Furthermore, we experiment with stochastic gradient boosting by setting the subsample parameter to 0.4 and 0.7. Other hyperparameter settings follow the default value of the scikit-learn library [29].

Neural Networks with LSTM Cells Our neural network consists of five layers (see Listing 3.7). First, we use an embedding layer with an embedding size of 32. The input is the Index-Based encoded samples without type encoding. This layer is trained end-to-end with the complete network on the trainings-data. Furthermore, this layer covers most trainable parameters of the model. As a hidden layer, we use 10 LSTM cells with a prior and posterior dropout layer with a default 0.2 dropout. To perform binary classification, our model ends with a dense mapping onto a single neuron. The binary result is determined based on a 0.5 threshold for the output of the last neuron.

We vary the hyperparameters for batch size, number of epochs, embedding size and number of LSTM cells. For implementation, we use Keras in version 2.4.3 [11].

Layer (type)	Output Shape	Param #
embedding (Embedding)	(None, 20, 32)	3200064
dropout (Dropout)	(None, 20, 32)	0
lstm (LSTM)	(None, 10)	1720
dropout_1 (Dropout)	(None, 10)	0
dense (Dense)	(None, 1)	11
Total params: 3,201,795		
Trainable params: 3,201,795		
Non-trainable params: 0		

Listing 3.7: Summary of our LSTM network. We use an embedding layer of size 32, 10 LSTM cells, a dropout layer before and after the LSTM layer with a default dropout ratio of 0.2 and a dense mapping to a single output neuron for binary classification with a threshold of 0.5.

4. Quantitative Evaluation

In this chapter, we will evaluate the approach described in Chapter 3. First, we will define the research questions. Afterwards, we will describe the motivation, details of the approach, and the results per research question. We explain the results we observe and give possible improvements.

4.1. Research Questions

RQ1 What is the utility of the CCAP besides existing tools?

RQ2 How do different models compare on the task of detecting non-clean code?

RQ3 Do machine learning-based models cover a larger variety of cases than rule-based checker?

4.1.1. RQ1: What Is the Utility of the CCAP Besides Existing Tools?

Motivation Several tools are established that aim to improve code quality and detect unclean code. Every new tool with similar claims has to prove its usefulness beside the preexisting tools. We compare the CCAP with existing tools based on the design goals expandability, useability and integration. Additionally, we will mention and compare useful features unique to existing tools.

Approach To evaluate, if the CCAP is a useful addition, we compare different features to see, if it can supplement or replace existing tools. As described in section 2.6, the compared tools are Sonarcube, PMD Source Code Analyser Project, Codacy and PyLint.

Results Table TODO shows a summary of the features.

Finding 1: Expandability For expandability of analysis plugins, CCAP is comparable to PyLint, that also offer plugins to analyse the raw string, the token stream or the AST. PMD allows plugins to analyse the AST or define XPath rules. Since XPath rules can be used in the CCAP as well (using a third-party library in the plugin), PMD offers fewer expansion possibilities. Sonarqube, on the other hand, provides the most possibilities for extension. Not only can developers analyse the AST, specify XPath rules, but they have access to an additional semantic model of the code that provides direct access to structures such as methods, parameters, return values. This semantic model simplifies analysis in comparison to the AST analysis. Additionally, Sonarqube allows plugins to expose an API for other plugins to use. This different type of plugin further increases the possibilities for rule checking plugins. The least expandability offers Codacy, that allow customising or disabling existing rules, but does not offer to add rules to the existing ruleset.

Regarding output plugins, no compared tool offers the output customisation of CCAP. PMD and PyLint have different predefined output formats like JSON, HTML, CSV or text. Although PyLint allows customising the message format using a formatting string as a command-line argument, they do not offer extensions for the output. Sonarqube displays the analysis reports in its WebUI. The scanner tool, SonarScanner, sends the reports to the server component, that renders the results in the browser. Codacy's local scanner produces text output; the cloud version also has a WebUI.

Finding 2: Integration shortcomings The CCAP has significant shortcomings in its integrations into IDEs, build processes and CI pipelines. All contestant tools provide plugins for build systems like Gradle or Maven. Except for Codacy, all tools have IDE integrations into the most common editors. Sonarqube and Codacy have integrations into source control systems. Same applies to PMD and PyLint since they are included in Codacy and therefore have the same integrations.

Finding 3: Useability For a plugin developer, the useability of the plugin system is simple for CCAP and PyLint. Both tools have the same plugin mechanism and a simple plugin interface. PMD has a similar, straightforward plugin interface, but its Java plugin has to be bundled into a jar file before adding to the classpath. The latter applies to Sonarqube as well; additionally, the powerful and rich plugin API makes it harder to use.

For the user, installing CCAP is like installing every other python package. Running

from the command-line is fast and can be automated, e.g. using git pre-commit hooks. The same workflow would be achieved with PyLint. Sonarqube is more complex to install due to its multiple components, but after the setup, the WebUI and integrations into most developers workflow have the best useability. Codacy as a cloud service only requires permission to access the source code repository to start scanning the code.

Finding 4: Multi-Language scanning Sonarqube, PMD and Codacy offer multi-language scanning. This multi-language ability is advantageous for code repositories with multiple languages and for reusing the same tools and infrastructure for multiple projects with different programming languages. Supporting multi-language scanning is not possible with CCAP and would require a modification of the architecture: The source handler would have to scan the files with respect to their programming language and transform those into a language-independent representation. Other parts of the architecture would remain similar. Although multi-language support was not a design goal, it would be a significant advantage for adoption.

Finding 5: Maturity and Community-Support The maturity and community-support of tools impact the integration of the tools into the developer workflow. Some integrations are community-made and shared, so everybody has an advantage. The community support would also be necessary for CCAP to write and share additional analysis plugins and to integrate into different workflows. Since the CCAP is at the beginning, it does not have a community that promote the integration and extension by other developers.

Summary In summary, CCAP offers a good, but not best in class expandability with analysis plugins. It has a high useability for plugin developer and users. The lack of integrations for different workflows reflects the lack of community support and maturity. We see the CCAP as an addition to the powerful Sonarqube. The simple, yet robust expandability of the CCAP could supplement the shortcomings of Sonarqube in its powerful, but complex expandability. For instance, a code reviewer could code problematic code into a simple python plugin, distributes it to the team and would hopefully not reencounter the same problem. Additionally, the CCAP could be used to teach about Clean Code and to enforce specific coding rules to students, since the local setup is straightforward. Clean Code rules could be turned into analysis plugins by students or teachers without having to understand a complicated interface.

4.1.2. RQ2: How Do Different Models Compare on the Task of Detecting Non-Clean Code?

Motivation As we have shown with the CCAP and the analysis plugins, it is possible to write an algorithm that can detect specific code patterns like non-clean code patterns. If we can write a well-tested detection algorithm, we do not need error-prone machine learning models to detect code patterns. Nonetheless, if we can not design an algorithm that detects the desired code pattern reliably, we may perform better using a machine learning model. Additionally, if we encode a subjective pattern, it may not be possible to design an algorithm to detect such a pattern objectively. As a solution, we could label code parts as our pattern and use machine learning approaches to extract rules to detect such a pattern implicitly.

For this research question, we will evaluate different machine learning models on different code patterns and compare the results.

Approach For answering the research question, we will use the approach detailed in chapter 3.2. The dataset, consisting of 18 projects from GitHub, is prepared as described in chapter 3.2.2. The data is split into 90% train and 10% test data. We use the training dataset to train Random Forest, Gradient Boosting Classifier, Support Vector Machines and an LSTM-based neural network. Due to the class imbalance, we will additionally train with an oversampled dataset (oversampling rate of 0.5) and an undersampled dataset (undersampling rate of 0.5). Since our test set contains real-world code samples, we expect it to represent a real-world label distribution, and we, therefore, accept the potential effects of the data imbalance. Nevertheless, we will try to train with an over- and undersampled dataset to observe the effect on the performance.

Our evaluation metrics are recall, precision and the combined f1 score. Due to the data imbalance, accuracy is not a meaningful metric. A high recall means a high detection rate of non-clean code. The costs of a false negative prediction are potential costs in unclean code (maintainability, understandability). With high precision, a detected non-clean code sample is likely to be a non-clean code sample and the system reports less false positives. The immediate cost of false-positive predictions (resulting in a lower precision) is the lost time of the developer. Additionally, the cry wolf effect comes into play [8]. In our case, a developer who has seen a false alarm will take subsequent alarms less serious. Additionally, a study has shown that a higher false alarm rate in the advisory warning system in cars results in a lower, subjective evaluation of the system [28]. A low subjective evaluation of developers would lead to a decrease in adoption rate, which

consequently leads to more unfixed problematic code. Since both recall and precision are essential to the success of our models, we decided to use the equally weighted f1 score as a single-value evaluation metric. Additionally, to compensate for the cascading effect of a low precision due to the cry wolf effect, we require a precision of 0.8 as a satisficing metric. Accordingly, we accept all models that return two false positive and eight true positive samples. On the other hand, we define the recall as our optimising metric.

We evaluate all models for the three different problem types RETURN_NONE (RN), CONDITION_COMPARISON_SIMPLE (CCS) and CONDITION_COMPARISON (CC). We expect the CC type to be the most difficult to learn since it is the most complex rule. CCS and RN should be easier to learn due to less possible variations in the code structure.

Results We list all hyperparameter configurations and the corresponding training, test and holdout performance for all problem types in the appendix in Tables A1 to A4. We report a detailed analysis in a separate finding for each model. Additionally, we describe our findings when comparing the models for the given classification tasks.

Finding 1: Support Vector Classifier *The support vector classifier in our configurations is not able to detect any problem type.* Due to the quadratic time complexity during training described in Section 3.2.1, we were only able to test a few configurations. We tested configurations with a general subsampling and undersampling to reduce the input size and to achieve a manageable training time. Therefore, we do not have data for non-resampled datasets or oversampled datasets to draw reasonable conclusions.

We can only observe some effects: First, the class balancing increases the recall for all problem types up to 0.79 for RN, 0.5 for CCS and 0.55 for CC. Second, undersampling with a ratio of 0.1 result in a f1 performance of zero for the tested RN and CCS problem types. Third, the performance for the CC type is double the performance for the CCS type (0.11 vs 0.06) and the performance for the CCS type is ten times the performance of the RN type. The zero performance with an 0.1 undersampling ratio may be an indicator of a strong sensitivity for class imbalance. But we have no supporting data for this hypothesis. Jiang Tian et. al describes a poor performance on imbalanced datasets with conventional support vecotr machines [40].

In conclusion, the support vector classifier does not fulfil our acceptance criteria of precision higher than 0.8. Especially the sensitivity for class imbalance makes the use for our problem domain with severely imbalanced datasets pointless. Furthermore, train time limits our experiments. It is important to note that we did no further investigation in improving the train time or performance. Maybe we could have increased classification and runtime performance by choosing a different kernel function such as a linear kernel function or by following strategies addressing this problem as proposed by Tian et al. [40]. However, we decided to invest our efforts into the other, more promising models.

Finding 2: Random Forest *Random forest classifier performs well and fulfils our satisficing condition on precision.* It achieves f1 scores of over 0.8 for the RN type and over 0.9 for the CCS and CC type. The best performing random forest classifier uses a 0.5 oversampling ration, 100 trees, type encoding and no additional class weighting.

Our observation indicates that the random forest classifier profit from more data, even with data duplication. We base this thesis on our observations with over- and undersampling the training data. All oversampled configurations perform better on the test set than their not resampled counterparts. This performance improvement is especially large for the RN type, with an f1 score increase of 0.03-0.05. For the CC and CCS problem type, the performance growth is smaller, with only around 0.01 in f1 score. We explain the difference between the problem types in the difference in class frequency. For the RN problem type, only 0.21% of the samples are labelled positive, whereas 2.66% of the samples for CCS and 4.87% of the samples for CC are labelled positive. Oversampling the data with a ratio of 0.5 results in more duplicated samples from the minority class for the RN problem type than for CCS and CC. Table 4.1 contains detailed measures of the train set size after oversampling and the performance metrics.

For undersampling, we observe a decrease in performance for all problem types. For the RN problem type, the f1 performance drops up to 0.51, whereas the performance decrease for CCS and CC is lower with up to 0.12 and 0.07, respectively (undersampling ratio of 0.5). We explain this difference again with the difference in class frequency and the resulting change in training size. With an undersampling ratio of just 0.1, fewer samples from the majority class get deleted and the performance for the RN type only deteriorates by 0.08. The different performance decrease between the CCS and CC type can also result from the difference in class frequency. The undersampling with an ratio

4. Quantitative Evaluation

	Test F1 Score			Dataset		
	RN	CCS	CC	RN	CCS	CC
No resampling	0.7599	0.9208	0.9211	4,950,329	4,950,329	4,950,329
0.5 oversampling	0.8379	0.9329	0.929	7,409,925	7,228,141	7,064,008
1.0 oversampling	0.8358	0.9287	0.9273	9,879,900	9,637,522	9,418,678
0.1 undersampling	0.6857	0.9268	0.929	114,169	1,447,248	2,650,890
0.5 undersampling	0.2583	0.806	0.8632	31,137	394,704	722,970

Table 4.1.: Dataset sizes for different resampling strategies with f1 performance on the test set. The f1 performance for the specific resampling is the random forest classifier with 100 trees, no class weighting and type encoding.

of 0.5 reduces the training size for CCS to 394,704 whereas for CC just to 722,970. With an undersampling ratio of just 0.1, the effect is smaller and only visible for the RN type (see Table 4.1).

Our experiments show that the random forest classifier is not very sensitive to class imbalance. Only an extreme class imbalance as for the RN type with a class distribution 0.21% vs 99.79% limits the classifier. We conclude this thesis from two main observations. First, all performance metrics for the CCS and CC problem type are comparable for all configuration without over- or undersampling. If the model would be sensitive to class imbalance, the higher class imbalance of the CCS type (97.34% vs 2.66%) should reflect in a worse performance than the CC type. We do not observe such performance differences. By contrast, the class frequency of positive samples for the RN type is by a factor twelve smaller than for the CCS type. The resulting severe class imbalance seems to affect the performance negatively. Second, the class weighting has no positive influence on the performance of the classifier. With balanced class weighting, the weighting for each class is adjusted inversely proportional to the class frequency. If the random forest classifier were sensitive to class imbalance, we would expect a performance improvement. Instead, we see a slight deterioration in f1 performance up to 0.04.

Type encoding increases the performance of the random forest classifier slightly for the CCS and CC type of up to 0.02. For the RN problem type, we observe a similar improvement for oversampling with type encoding. For undersampling with type encoding, we observe a significant decrease of $\tilde{0.09}$ in f1 performance. We expected type encoding to have a more significant impact. The types should work as an additional feature the classifier can utilize. With the index-based encoding of the token values, 100,001 possible values can be encoded. The number of different token types is signifi-

cantly smaller with only 61 different types (see the complete list of possible token types in Figure C.1). Especially for function and method names, it would allow the model to treat two different function calls as one. However, the improvement is lower than expected. At least for the CCC type, the NAME type of a method call should be easier to learn on than the changing value of the method name.

Overall, the random forest classifier handles class imbalance fairly well and can profit from more data, even if it contains duplicates from oversampling. Undersampling decrease the performance of the classifier. Type encoding has slight positive impact on the performance.

Finding 3: Gradient Boosting Classifier *The gradient boosting classifier performs well on all problem types and fulfils our satisfying metric.* We achieve an f1 score of over 0.9 for the RN and CC type and 0.84 for the CCS type. This best performing configuration is trained with 300 boosting stages, a learning rate of 0.2 and type encoding.

The gradient boosting classifier is sensitive to less training samples or samples with less variety. Two main observations allow this conclusion. First, we observe a decrease in performance for under- and oversampling (see Table 4.2). For undersampling the RN type with a ratio of 0.5, the f1 score drops from 0.8306 by roughly 45% and by approximate 41% for oversampling with a ratio of 0.5. For the CCS type, the f1 score drops from 0.773 by approximately 9% for over- and undersampling. The smallest drop is for the CC type from 0.8567 by nearly 7% for over- and undersampling. With undersampling, we reduce the number of examples and with oversampling, we increase the minority class without increasing the variety of samples. Since the performance characteristics for over- and undersampling are comparable, we conclude that more samples with the same variety further decrease the model’s performance. The decrease in f1 score is caused by a drop in precision, i.e. an increase of false positives. With less data, the classifier tends to be biased towards the positive class. Second, the larger performance drop for the RN problem type and the smaller decrease for the CC type correlates with the number of training samples. A larger decrease in training data results in a larger drop in performance.

The random forest classifier is insensitive towards class imbalance. One supporting point is the aforementioned performance decrease with over- or undersampling. Additionally, the performance of the RN problem type is comparable to the CC problem type despite their difference in class imbalance (0.21% class frequency for RN vs 4.87% for CC).

4. Quantitative Evaluation

	Test F1 Score			Dataset		
	RN	CCS	CC	RN	CCS	CC
No resampling	0.8306	0.773	0.8567	4,950,329	4,950,329	4,950,329
0.5 oversampling	0.4847	0.6948	0.7933	7,409,925	7,228,141	7,064,008
0.5 undersampling	0.4545	0.6836	0.7925	31,137	394,704	722,970

Table 4.2.: Dataset sizes for different resampling strategies with f1 performance on the test set. The f1 performance for the speciifc resampling is the gradient boosting classifier with 200 stages, a learning rate of 0.2 and type encoding.

The random forest classifier has problems learning the CCS type. Contradictory to our assumptions, the CC problem type is easier to learn than the CCS type. The precision is lower by up to 0.05, whereas the recall is lower by up to 0.11 (without resampling and ≥ 200 boosting stages). We assume, the gradient boosting classifier has trouble learning the CCS subset and instead learns the more general superset of CC problems. To recap, the CCS problem type only labels a direct comparison in the if-condition as problem whereas the CC type consider nested boolean expressions. With the CCS subset, we restrict the classifier to only detect a direct comparison. We analysed the classification of the test samples, especially the false-positive ones. If the classifier learn the structure of the CC superset, the classifier would label samples as positive that are only positive for the CC type. We randomly sampled 20 false-positive samples and found 11 samples like in Listing 4.1. Those 11 samples contain the a direct comparison nested inside boolean operators. Such nesting is only labeled as a CC problem, not as CCS problem. We think our restriction in problem types collide with the classifier’s ability to detect the more gneral pattern of a comparison somewhere inside the if-condition.

Subsampling the train data for the individual base learner during training has no impact. The performance for all subsampling ratios are identical, so the individual base learners are trained with the same data as before.

The gradient boosting classifier is robust to overfitting the training data. We observe a comparable train and test performance for all non-resampled configurations. Only with under- or oversampling, the classifier seems to overfit the train data, since the test performance is worse the train performance. Another indicator for the robustness is the performance increase when using more boosting stages. Increasing the boosting stages increase the recall and precision for all problem type. Our best model with 300 boosting stages and a 0.2 learning rate can probably be more fine-tuned by increasing the boosting stages and adapting the learning rate. Again, being restraint by the increase in training

```

    if input.dim() == 4 and mode == 'bicubic' :
        assert align_corners is not
DEDENT,NAME,NAME,OP,NAME,OP,OP,OP,NUMBER,
NAME,NAME,OP,STRING,OP,NEWLINE,INDENT,NAME,
NAME,NAME,NAME
#####

input_data = [ indices , values , default ]
if <UNKNOWN> and n > 1 :

NEWLINE,NAME,OP,OP,NAME,OP,NAME,OP,NAME,OP,
NEWLINE,NAME,NAME,NAME,NAME,OP,NUMBER,OP,
NEWLINE,INDENT
#####
:
<UNKNOWN>
if target is none or source is none :
    return false

OP,NEWLINE,INDENT,STRING,NEWLINE,NAME,NAME,
NAME,NAME,NAME,NAME,NAME,NAME,OP,NEWLINE,
INDENT,NAME,NAME,NEWLINE,DEDENT

```

Listing 4.1: False-positive samples for the best performing random forest classifier with 300 boosting stages, a learning rate of 0.2, type encoding and without resampling for the CCS type. We removed spaces between tokens if necessary, but leave the indentation as in the sample.

time, we did not fine-tune the classifier further.

In conclusion, the gradient boosting classifier is robust to overfitting and somewhat improves its performance with more boosting steps. We can also observe robustness for class imbalance and on the other hand sensitivity towards less training data. Subsampling has no effect and we observe an unexpected high difficulty of learning the CCS type for this classifier.

Finding 4: LSTM-Based Neural Network *Our implementation of an LSTM-based neural network performs well and fulfils our satisficing condition.* We achieve an f1 score of over 0.97 for all problem types without resampling, 10 LSTM cells, an embedding size of 32 and training with 2 epochs and 256 batch size.

Since we have a remarkably high performance, we do not vary all parameters. Instead, we use a 0.5 under- and oversampling as for other models, vary the batch size and number of epochs and try reducing the size of the embedding layer since this layer contains the majority of the trainable parameter (see listing 3.7).

4. Quantitative Evaluation

	Test F1 Score			Dataset		
	RN	CCS	CC	RN	CCS	CC
No resampling	0.987	0.9808	0.9782	4,950,329	4,950,329	4,950,329
0.5 oversampling	0.9859	0.9781	0.9745	7,409,925	7,228,141	7,064,008
0.5 undersampling	0.6294	0.9564	0.9561	31,137	394,704	722,970

Table 4.3.: Dataset sizes for different resampling strategies with f1 performance on the test set. The f1 performance for the speciifc resampling is the LSTM with an embedding size of 32, 10 lstm cells, a training on 3 epochs with a batch size of 256.

Oversampling does not have a measurable impact on our LSTM-based neural network. The f1 performance with a 0.5 oversampling is not significantly different from the performance without oversampling. We observe overfitting on the train data with a near-perfect recall and precision. The test performance remains higher than 0.97 f1 score. Although the model overfits to the train data, it still generalises well enough with comparable performance to the models trained on the original class distribution.

The LSTM-based neural network is sensitive to fewer train data. We draw this conclusion based on the performance deterioration for undersampled train sets (see Table 4.3). For the RN type, undersampling results in a significant drop in performance from 0.987 by roughly 36%. For the CCS and CC type, the performance decrease is smaller at less than 3%. The train set of the RN type is significantly smaller after undersampling than for the CCS and CC type, due to the lower class frequency of the minority class of only 0.21% before undersampling. Therefore, more samples from the majority class are deleted during undersampling and the variety decreases. In combination with many trainable parameters in our model, the model fails to fit the training data correctly. Responsible for the performance drop is a decrease in precision, whereas the recall remains high. In other words, the false-positive rate increases by factor 100 to 0.3%. By removing negative samples and variety in training, we think it becomes more challenging for the model to detect negative samples correctly. On the other hand, the high recall indicates a good detection capability for the positive class. Since we did not decrease the number of positive samples, the model can learn on the full variety of those models.

Our model configuration has too many trainable parameters for the problem complexity. We base this assumption on two observations. First, we halved the size of the embedding layer to 16 and consequently reducing the 3,201,795 trainable parameters of the model to roughly half the size. If the model needed all parameters to fit the training

data accurately, we would expect a performance drop. Instead, we do not observe any significant change in performance. Due to this observation, we assume a further reduction of the embedding size may be possible without performance loss. However, since this configuration trains in a reasonable time of under an hour without a GPU and has a test performance, we do not investigate this too powerful model further. Second, a variation of the epochs and batch size for the training yields no significant difference. We think, the problem detection is too simple for the model and it quickly reaches its optimum. Further training does not increase the model's performance further.

To sum up, we achieve near-perfect performance with our neural network model. We acknowledge, that our model has too many trainable parameters that we do not need for the same performance as shown with the reduction of the embedding layer size. The model does not benefit from resampling. Oversampling has no impact on the test performance and undersampling impairs the training of the model's weight since it reduces the overall number of samples.

Finding 6 We found a weakness in the way we labeled samples with the rule checkers. We decided to label a sample as rule violation, if all tokens of the problematic line are present. With the sliding window approach, it is possible to have all tokens for the detection in the sample, but missing the first or last newline or indentation tokens. The sample would therefore be labeled as clean code, although the problematic pattern is there. In the analysis of the false-positives for the gradient boosting classifier, we found samples for both cases in the false-positives as seen in Listing 4.2. In the first sample, not all preceding DEDENT types are within the window, whereas the trailing NEWLINE token is missing in the second sample.

To improve this weakness, we would strip out all DEDENT and possibly NEWLINE tokens from the data before applying the sliding window. Due to time constraints, we are not able to retrain all models with the improved approach. However, the problematic line will occur in at least another sample, if the line is shorter than 16 tokens. It will most likely not impact the variety of problem patterns, but it will slightly reduce our performance score due to false-positives.

Summary Table 4.4 shows the performance score of the best performing model of each classifier for all problem types. The SVM performs poorly and does not meet our mandatory satisficing condition of precision higher than 0.8. Therefore, we will not evaluate it further. On the other hand, the LSTM-based neural network model

```

if 'add_host' in <UNKNOWN> :
<UNKNOWN>
    <UNKNOWN> = <UNKNOWN> . get ( 'add_host' ,
DEDENT,DEDENT,DEDENT,NAME,STRING,NAME,
NAME,OP,NEWLINE,COMMENT,NL,INDENT,NAME,
OP,NAME,OP,NAME,OP,STRING,OP
#####
if vhost.ssl :
points += 3

if points > <UNKNOWN> :
NL,DEDENT,NAME,NAME,OP,NAME,OP,NEWLINE,
INDENT,NAME,OP,NUMBER,NEWLINE,NL,DEDENT,
NAME,NAME,OP,NAME,OP

```

Listing 4.2: False-positives of the best gradient boosting classifier with 300 boosting stages, a learning rate of 0.2, type encoding and without resampling for the CCS type. In the first sample, not all preceding DEDENT tokens are in this window. For the second example, the trailing NEWLINE token is missing in the window. Without this encoding issue, both samples would have been detected correctly.

		RN	CCS	CC
Random Forest				
	F1	0.8486	0.9352	0.9351
	Recall	0.7404	0.8936	0.8976
	Precision	0.9939	0.9809	0.9759
SVM				
	F1	0.0066	0.0622	0.1115
	Recall	0.0134	0.5004	0.5535
	Precision	0.0044	0.0332	0.062
Gradient Boosting Classifier				
	F1	0.9328	0.8489	0.9088
	Recall	0.9005	0.7953	0.8766
	Precision	0.9675	0.9102	0.9434
LSTM				
	F1	0.9933	0.9805	0.9783
	Recall	0.9933	0.9779	0.9868
	Precision	0.9933	0.9831	0.97

Table 4.4.: Best performing classifier for RQ2 based on f1 score of the holdout set. For the SVM, the test score of the configuration with available data is used.

outperforms all other models on this classification task for all problem types.

Both random forest and gradient boosting classifier, fulfil the satisficing condition for all problem types. Interestingly, they supplement each other: For the RN problem type, the gradient boosting classifier performs better with an f1 score of 0.9328, whereas the random forest classifier surpasses the gradient boosting classifier for the CCS and CC problem type with 0.9352 and 0.9351 f1. As discussed, we explain the former result with the better insensitivity to a class imbalance of the gradient boosting classifier. In contrast, the latter seems to be harder to learn for the gradient boosting classifier. Especially the difficulties for the CCS type are unexpected and do not affect the random forest classifier or neural network.

4.1.3. RQ3: Do Machine Learning-Based Models Cover a Larger Variety of Cases Than Rule-Based Checker?

Motivation For the previous research question, we trained models on a dataset. We created the dataset by downloading open-source repositories and applying our analysis plugins from Section 3.1.1 to label samples as clean or non-clean code. With this automated approach, it is simple to generate a dataset and to scale the dataset to any size needed for training. Such algorithms (checkers) that detect non-clean code are the prerequisite for the dataset and successful training of machine learning models. This research question wants to put the machine learning approach into a scenario, in which this prerequisite is not met and we lack a checker to generate a labelled dataset. This scenario transfers to clean code rules with a high level of complexity (as defined in Section 2.4). For these rules, we may not be able to analyse the source code or another representation with a checker to generate a representative dataset. We could (1) hand-label samples and train a machine learning model on those samples. Alternatively, we may (2) only find a deterministic checker that detects a subset of the rule violations. We hope that model trained on either of those datasets generalise well to all patterns of the rule violation. Both scenarios boil down to a lack of samples for all or many structural variations of the rule violation and a discrepancy of the test set and the real world with more variety. We will evaluate both scenarios for our models from the previous experiment in this research question.

Approach To answer this research question, we use all models for all three problem types from the previous research question: We use the models trained on the RN problem

type to simulate scenarios (1). First, we manipulate the holdout dataset as described in 3.2.3. The code still contains returns of `None`, although it is not explicitly written as two consecutive code tokens. We want to explore if a model generalised well enough to detect unseen variations of the problem type. With the manipulation, we ensure that all samples contain the unseen variations and our performance metrics only refer to those samples.

Second, we simulate scenario (2) of having trained on a subset of all variations with the CCS models. The CCS model was trained on a subset of the CC problem patterns. We manipulate all occurrences of the CC problem type to contain only patterns that are not in the subset for the CCS training. It is important to note that the class frequency of positive samples changes from the trained 2.51% to the 5.1% positive samples for the manipulated CC type in the holdout set.

Last, we evaluate the models trained on the CC type with our manipulated code to provide a baseline for the CCS trained models. We expect the CC model to show similar performance than in the previous experiments since it was trained on the problem structure of the manipulated CC dataset.

We expect the best performing CC model from the previous experiments to outperform all CCS models and to show similar performance like in the previous research question.

Results We evaluate the models’ performance for both scenarios: The RN problem type simulates scenario (1) whereas the CCS problem type represents scenario (2). We list all evaluation results in the appendix in Tables A1 to A3.

Finding 1 *The detection of the manipulated RN problem type works poorly with all models.* The random forest classifier has an f1 score of less than 0.1, the gradient boosting classifier of 0.24 and the LSTM-based neural network of 0.39. The precision of all models is below 0.26 and would therefore not comply with the satisficing condition from the previous experiment of 0.8. Our simulation for scenario (1) fails.

It is obvious that our models do not generalise on the problem structure. Since all models perform significantly worse than in the previous experiment, it is feasible to search for the reason in the problem type. We think the model did not learn the problem structure during training. Our reasoning goes as follows. The problem type appears in the training set only in one variation: A `return` keyword is followed by a subsequent `None` value. The surrounding tokens varied, but only the subsequent occurrence of these two tokens had to be learnt to perform to detect the problem type. For the previous

experiment, this behaviour resulted in the best performance. To perform well on this experiment, the model should have learned a more sophisticated approach of a `None` value after the `return` keyword before the end of the line. Since this was not part of the training data, it is obvious that the model learnt the simple subsequent detection of this pattern, since it yields the best performance in the previous experiment.

Finding 2 *The classification of the modified CC samples with the CCS models yields better results than for the RN type.* For the random forest classifier, the performance on the modified CC dataset is at 0.72 f1 score. The best gradient boosting classifier in this experiment achieves 0.83 and the best LSTM-based neural network reaches 0.66 f1 score. Compared to the performance for the RN types, the performance values are at least double as high, so the model has learnt the pattern from the subset better. We attribute this to the wider variety in problem pattern in the source code: Since the CCS checker looks for comparison in the condition, it will label different comparison operators such as `==`, `<` or `>` as rule violations, whereas the only variation for the RN type is `return None` as subsequent tokens.

Finding 3 *The more fine-tuning during training, the less the model generalise for the modified problems.*

All models on the modified problem types behave contrary to the unmodified problem types in the previous experiments. The hyperparameter combinations with a worse performance in the previous experiment score better in this evaluation. By contrast, the best model configurations from the previous experiment are among the worst-performing models for this evaluation. In Figure 4.1, we plotted the f1 performance of the baseline, manipulated RN and CCS. For each classifier, we sorted the models descending on the baseline f1 score. The baseline performance is the performance of the CC model from the previous experiment on the manipulated CC data. It correlates with the performance of the CC model on the original holdout data in the previous experiment. The overall picture shows, that the models with the best baseline performance have a low performance on the manipulated datasets. With decreasing baseline performance, the performance of the manipulated CCS increases for the random forest and gradient boosting classifier. The manipulated RN type follows this pattern for the random forest classifier and the LSTM.

We observe this behaviour on different hyperparameters we tuned for better performance in the previous experiment. It applies to both problem types. For the random



Figure 4.1.: F1 performance of the models, sorted by the baseline f1 score. The baseline performance correlates with the performance of the CC type in the previous experiment. With a worse performance in the previous experiment, the performance for the manipulated types increases.

forest classifier, undersampling the training set lead inevitably to worse performance in the previous experiment than without undersampling. In this experiment, the under-sampled random forest classifier leads to the overall best scores. For gradient boosting classifier, the models without resampling scored best in the previous experiment and are performing worst in this experiment. Additionally, increases in boosting stages that lead to an improvement in the previous experiment now lead to a deterioration in performance. The LSTM-based neural network results draw a similar picture. The undersampling harmed the performance, whereas, in this experiment, a 0.5 undersampling achieves the only f1 score over zero for the RN type and the best for the manipulated CC type.

We would explain this observation separate for the RN and CCS type: For the RN type, the models probably only learnt the direct token sequence (as described in Finding 1). The better it learnt this sequence, the better its performance on the previous experiment. For this experiment, the models that only learnt to look for a `return` and a `None` somewhere in the token stream achieve a higher result. In contrast, it performed poorly on the previous experiment. Consequently, we assume the poor performing models from the previous experiment followed this pattern. (TODO can i proof this point with FP from the evaluation? maybe find samples in which there are just these two tokens)

The CCS problem type, simulating scenario (2), offered more variety to learn from, as discussed in Section 4.1.3. If the model learnt from this variety a general problem structure, fine-tuning the hyperparameters for the previous experiment would restrict this generality to only CCS patterns. We assumed in Section 4.1.2 that the gradient boosting classifier get restricted by the CCS type and is more likely to learn the CC pattern. In this experiment, the best gradient boosting classifier configuration surpasses all other models. Additionally, it surpasses the baseline performance measure of the model trained on the unmodified CC type that we assumed to perform better. This solidifies our assumptions about the learning of the CCS type. We think this explanation applies to other models (TODO false positives from previous experiment anschauen und sehen ob da CC type die indikatoren waren).

The explanations differ for both scenarios, but their common ground is an insufficient amount of variety. While the RN type has only one variety in the problem pattern, the CCS rules restricted the variety to the subset.

Improvements As identified in the findings (4.1.3 and 4.1.3), we see the root cause for the bad performance in the learning of the specific code pattern and not the general structure of a problem type. We identified the lack of variety as the main reason, manifesting as the restriction of the learning on the CCS type. For the RN type, the general structure on a token level would be something like a `None` value somewhere after a `return` but before a line break. A slightly more advanced structure could be extracted from the CCS type: A comparison inside a condition (starting with an `if` statement and closing with the colon) should be marked. With the chosen approach, the model did not train on the general problem structure of the clean code violation but explicit code patterns.

We see several ways to improve: First, we did not provide the model with the full variety of problem patterns. Therefore, we suggest hand-label or hand-generated pat-

terns with explicit variations. In the motivation for this research question, we state that it may not be practically possible to collect many samples with the full variety of problem patterns for an adequate training corpus. For the improvement, we need a base collection of samples. This base collection can be obtained by an algorithm that can only detect a subset of the problem patterns (like in our approach with the CCS type) or a similar enough code pattern (like we tried with the RN type). On top of that, we explicitly create samples with more variety by either hand-labelling existing code or writing new code. With the hand-crafted samples, we fine-tune the model to increase the variety of code patterns the model can generalise on. Simultaneously, we loosen the indirect restriction of the dataset covering only a subset of the data. By combining a scalable, automatic generated dataset with hand-labelled samples, we may be able to improve the models' performance on the samples that we could not label automatically. Additionally, we should use some hand-labelled samples for more variety in the test corpus to tune the hyperparameters to more possible problem variations in real datasets.

Second, we may compensate for less variety in training samples with a different encoding. Using the encoding to focus on the code structure may require fewer variations to extract the problematic code structure. We use an indexed-based encoding of the source codes token stream. With this approach, we do not exploit the code structure in our encoding. Other machine learning tasks on code profit from other encoding types that exploit more structure from the code. Alon et al. introduced code2vec to represent code as paths in the AST [5], capturing the structure of the code. Similarly, Zhang et al. split the AST into smaller statement trees and apply a recurrent neural network approach to generate a code representation [47]. Since we identified the lack of learning the problem structure from the samples, an encoding approaches that explicitly take advantage of the structure of code may improve the performance of this experiment. It may be even possible to train the model on a subset of possible problem variations like we tried with the manipulated CCS type without the restricting effect, since the problem structure may be inferred from the subset alone. Additionally, this may be combined with the fine-tuning on hand-labelled samples to increase the performance to a practical level. The satisficing condition of precision higher than 0.8 we defined in the previous experiment would also apply to this experiment. A lower precision would lead to the cry wolf effect and the practical use would therefore be small.

Summary In summary, the detection of the modified RN type performs poorly with all models and configurations. We see the reason in the low variety of different problem patterns. The experiment with the RN type should simulate the scenario in which we could only obtain a dataset with hand-labelled samples that do not cover all possible patterns. Based on the analysis of the results, we admit that a hand-labelled sample set will probably contain a greater variety of samples. With more variety, the result would most likely be more comparable to the performance of the CCS type with more variation. If we define the subset of the CCS type as a hand-labelled set, we can interpret the results for this problem type in the context of scenario (1).

The performance on the modified CC type is more promising than for the RN type. We explain this by more variety in the training samples, although the models seem to be restricted to this subset (especially the gradient boosting classifier). If we interpret the performance with the CCS type model on the modified CC type data as scenario (1), the variations of our hand-labelled samples will restrict the models' performance. Furthermore, it would lead to a misleading parameter tuning, since our test set would not contain all real-world variations.

For improvement, we suggest to hand-craft additional samples with a different variation of the problem pattern. This can be used for increasing the variety of the training set with fine-tuning or to better approximate the test set with the real world to tune the models' parameter towards real-life datasets. Additionally, encoding the data by utilising the source code structure, less variety may be necessary to learn the variations. If this is the case, it would have positive practical implications, since it would require less manual code labelling or creating.

5. Conclusion

5.1. Summary

Software maintenance is a considerable cost factor for software development. With high code quality, the maintenance efforts can be optimised. With the clean code principles, developers have a ruleset to write understandable and readable code that will positively impact code quality. The focus of this thesis lies on the automated checking of such rules.

We introduced a taxonomy of complexity level for the clean code rules, depending on the level of complexity for an automated checker. On the highest level, we grouped clean code rules that would potentially require machine-learning approaches to learn the detection of a rule violation.

The contribution of this work is twofold: First, we designed and implemented the clean code analysis platform (CCAP). It is a tool that checks for clean code violations. One of the essential design aspects is the easiness to extend the tool to cover additional clean code rules. In comparison with preexisting tools, we see a high value for users and extension developers, that write automated checker for new rules. On the downside, it lacks the integration into different workflows. The CCAP is an addition to Sonarqube. The strength in extensibility supplement the rather complicated extensibility of Sonarqube. Furthermore, CCAP is a useful tool to teach about clean code principles.

Second, we compared different machine learning models on detecting violating code patterns. A support vector classifier was not able to complete the task with acceptable performance. The LSTM-based neural network outperforms all other models on the task, but it is sensitive towards less training data due to many trainable parameters. A random forest classifier fulfilled our acceptance criteria. It can handle moderate class imbalance and profit from oversampling the training samples. The gradient boosting classifier is insensitive towards class imbalance which is a useful property for the task. When classifying for a single clean code violation, most lines of code do not contain the violating pattern, so the training dataset is imbalanced. The gradient boosting classifier

also fulfils our acceptance criteria and performs better than the random forest on two out of three problematic code patterns.

Furthermore, we tried to simulate a situation in which the dataset may not cover the full variety of possible problem patterns. We tested the trained models on unseen problem variations. No model performed on an acceptable level. We assume, our models did not learn the structure of the problematic pattern but just the possible variations. We can think of two main reasons: First, the variation in the training set is too limited for the model to extract the general structure. We suggest to explicitly introduce variety by hand-generating different patterns to fine-tune the model and to use in the test set for hyperparameter tuning. Second, our encoding approach does not utilise the structure of the source code. Since our models did not extract the structure from our encoding, encoding the structure as input features would ease the learning of the problem structure.

5.2. Future Work

TODO CCAP language server protocol IDE integration, different warning levels, disabling plugins, ignore a specific rule violation line

Bibliography

- [1] Go at Google: Language Design in the Service of Software Engineering. <https://talks.golang.org/2012/splash.article>.
- [2] Null Safety - Kotlin Programming Language. <https://kotlinlang.org/docs/reference/null-safety.html>.
- [3] Abdiansah Abdiansah and Retantyo Wardoyo. Time Complexity Analysis of Support Vector Machines (SVM) in LibSVM. *International Journal of Computer Applications*, 128(3):28–34, October 2015.
- [4] Frances E. Allen. Control flow analysis. *ACM SIGPLAN Notices*, 5(7):1–19, July 1970.
- [5] Uri Alon, Meital Zilberstein, Omer Levy, and Eran Yahav. Code2vec: Learning Distributed Representations of Code. *arXiv:1803.09473 [cs, stat]*, October 2018.
- [6] Robert Baggen, José Pedro Correia, Katrin Schill, and Joost Visser. Standardized code quality benchmarking for improving software maintainability. *Software Quality Journal*, 20(2):287–307, 2012.
- [7] Grzegorz Baranski. Cleaned code, much improvement in performance. [gbaranski/houseflow@7d06d9f](https://github.com/gbaranski/houseflow/commit/7d06d9f6ea78530e5a542f71f6ddd536e2d55bb7). <https://github.com/gbaranski/houseflow/commit/7d06d9f6ea78530e5a542f71f6ddd536e2d55bb7>, September 2020.
- [8] Shlomo Breznitz. *Cry Wolf: The Psychology of False Alarms*. Lawrence Erlbaum Associates, Hillsdale, N.J, 1984.
- [9] Lutz Buch and Artur Andrzejak. Learning-Based Recursive Aggregation of Abstract Syntax Trees for Code Clone Detection. In *2019 IEEE 26th International Conference on Software Analysis, Evolution and Reengineering (SANER)*, pages 95–104, Hangzhou, China, February 2019. IEEE.

- [10] G Ann Campbell. Cognitive Complexity-A new way of measuring understandability. *SonarSource SA*, 2018.
- [11] François Chollet et al. Keras. 2015.
- [12] D. Coleman, D. Ash, B. Lowther, and P. Oman. Using metrics to evaluate software system maintainability. *Computer*, 27(8):44–49, August 1994.
- [13] Aurelien Delaitre, Bertrand Stivalet, Elizabeth Fong, and Vadim Okun. Evaluating Bug Finders – Test and Measurement of Static Code Analyzers. In *2015 IEEE/ACM 1st International Workshop on Complex Faults and Failures in Large Software Systems (COUFLESS)*, pages 14–20, Florence, Italy, May 2015. IEEE.
- [14] Galorath. Accurately Estimate Your Software Maintenance Costs, January 2019.
- [15] Andrew Gerrand. Error handling and Go. <https://blog.golang.org/error-handling-and-go>, December 2011.
- [16] GitHub Inc. The State of the Octoverse. <https://octoverse.github.com/>, 2019.
- [17] Maurice Howard Halstead et al. *Elements of Software Science*, volume 7. Elsevier New York, 1977.
- [18] Wang Haoyu and Zhou Haili. Basic Design Principles in Software Engineering. In *2012 Fourth International Conference on Computational and Information Sciences*, pages 1251–1254, Chongqing, China, August 2012. IEEE.
- [19] Martin Hitz and Behzad Montazeri. Measuring coupling and cohesion in object-oriented systems. In *Proceedings of International Symposium on Applied Corporate Computing*, pages 25–27, 1995.
- [20] Tony Hoare. Null References: The Billion Dollar Mistake, March 2009.
- [21] ISO/TC 210. IEC 62304:2006 Medical device software — Software life cycle processes. Standard, International Organization for Standardization, Geneva, CH, May 2006.
- [22] ISO/TC 22/SC 32. ISO 26262:2018 Road vehicles — Functional safety. Technical report, International Organization for Standardization, December 2018.
- [23] Jukka Lehtosalo et al. Mypy-Optional Static Typing for Python. <http://mypy-lang.org>.

- [24] K.J. Lieberherr and I.M. Holland. Assuring good style for object-oriented programs. *IEEE Software*, 6(5):38–48, September 1989.
- [25] Robert C Martin. *Clean Code: A Handbook of Agile Software Craftsmanship*. Pearson Education, 2009.
- [26] T.J. McCabe. A Complexity Measure. *IEEE Transactions on Software Engineering*, SE-2(4):308–320, December 1976.
- [27] Torben Ægidius Mogensen. *Introduction to Compiler Design*. Undergraduate Topics in Computer Science. Springer International Publishing : Imprint: Springer, Cham, 2nd ed. 2017 edition, 2017.
- [28] Frederik Naujoks, Andrea Kiesel, and Alexandra Neukum. Cooperative warning systems: The impact of false and unnecessary alarms on drivers’ compliance. *Accident Analysis & Prevention*, 97:162–175, December 2016.
- [29] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [30] John Pendenque. Code refactoring + Cleaning + Improvements . Zadigo/ecommerce_template@db2655c. https://github.com/Zadigo/ecommerce_template/commit/db2655c7ba16add17db221761fabda2 August 2020.
- [31] Herbert Prahofer, Florian Angerer, Rudolf Ramler, and Friedrich Grillenberger. Static Code Analysis of IEC 61131-3 Programs: Comprehensive Tool Support and Experiences from Large-Scale Industrial Application. *IEEE Transactions on Industrial Informatics*, 13(1):37–47, February 2017.
- [32] Veselin Raychev, Pavol Bielik, and Martin Vechev. Probabilistic model for code with decision trees. *ACM SIGPLAN Notices*, 51(10):731–747, 2016.
- [33] Mehwish Riaz, Emilia Mendes, and Ewan Tempero. A systematic review of software maintainability prediction and metrics. In *2009 3rd International Symposium on Empirical Software Engineering and Measurement*, pages 367–377, Lake Buena Vista, FL, USA, October 2009. IEEE.

- [34] Vaibhav Saini, Farima Farmahinifarahani, Yadong Lu, Pierre Baldi, and Cristina V. Lopes. Oreo: Detection of clones in the twilight zone. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering - ESEC/FSE 2018*, pages 354–365, Lake Buena Vista, FL, USA, 2018. ACM Press.
- [35] Ken Schwaber and Mike Beedle. *Agile Software Development with Scrum*. Series in Agile Software Development. Prentice Hall, Upper Saddle River, NJ, 2002.
- [36] ISO Central Secretary. ISO/IEC 25010:2011: Systems and software engineering — Systems and software Quality Requirements and Evaluation (SQuaRE) — System and software quality models. Standard ISO/IEC 25000:2014, International Organization for Standardization, Geneva, CH, 2011.
- [37] Dag I K Sjøberg, Dag Sjøberg, Bente Anda, and Audris Mockus. Questioning Software Maintenance Metrics: A Comparative Case Study. page 4.
- [38] Stripe and Harris Poll. The Developer Coefficient - Software engineering efficiency and its \$3 trillion impact on global GDP. Technical report, September 2018.
- [39] TC 65/SC 65A. IEC 61508:2010 Parts 1-7. Standard, International Electrotechnical Commission, April 2010.
- [40] Jiang Tian, Hong Gu, and Wenqi Liu. Imbalanced classification using support vector machine ensemble. *Neural Computing and Applications*, 20(2):203–209, March 2011.
- [41] Pavel Timofeev. [lab4] some clean-code improvement · ptimofeyeff/distributed-computing@ff5e132. <https://github.com/ptimofeyeff/distributed-computing/commit/ff5e1322611a0c3b3c62556d7a775441630df252>, August 2020.
- [42] Guido Van Rossum, Jukka Lehtosalo, and Lukasz Langa. PEP 484—type hints. <https://www.python.org/dev/peps/pep-0484/>, 2014.
- [43] Guido Van Rossum, Barry Warsaw, and Nick Coghlan. PEP 8: Style guide for Python code, 2013.
- [44] Dave Wichers and Eitan Worcel. Source Code Analysis Tools | OWASP. https://owasp.org/www-community/Source_Code_Analysis_Tools.

- [45] S.W.L. Yip and T. Lam. A software maintenance survey. In *Proceedings of 1st Asia-Pacific Software Engineering Conference*, pages 70–79, Tokyo, Japan, 1994. IEEE Comput. Soc. Press.
- [46] Sheng Yu and Shijie Zhou. A survey on metric of software complexity. In *2010 2nd IEEE International Conference on Information Management and Engineering*, pages 352–356, Chengdu, China, 2010. IEEE.
- [47] Jian Zhang, Xu Wang, Hongyu Zhang, Hailong Sun, Kaixuan Wang, and Xudong Liu. A Novel Neural Source Code Representation Based on Abstract Syntax Tree. In *2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE)*, pages 783–794, Montreal, QC, Canada, May 2019. IEEE.

Appendices

A. RQ2: Model Performances

Parameter					RN			CCS			CC			
over-samp.	under-samp.	#trees	class weight	encode type		Train	Test	Holdout	Train	Test	Holdout	Train	Test	Holdout
0.5	-	100.0	None	True	<i>F1</i>	1.0	0.8379	0.8486	1.0	0.9329	0.9352	1.0	0.929	0.9351
					<i>Rec</i>	1.0	0.7221	0.7404	1.0	0.8898	0.8936	1.0	0.8888	0.8976
					<i>Prec</i>	1.0	0.9978	0.9939	1.0	0.9803	0.9809	1.0	0.9731	0.9759
1.0	-	100.0	None	True	<i>F1</i>	1.0	0.8358	0.8367	1.0	0.9287	0.9322	1.0	0.9273	0.9324
					<i>Rec</i>	1.0	0.719	0.7225	1.0	0.8824	0.8876	1.0	0.8853	0.8925
					<i>Prec</i>	1.0	0.9978	0.9937	1.0	0.9801	0.9815	1.0	0.9736	0.9761
0.5	-	100.0	balanced	True	<i>F1</i>	1.0	0.8352	0.8385	1.0	0.9287	0.9321	1.0	0.9267	0.9329
					<i>Rec</i>	1.0	0.7182	0.7249	1.0	0.8826	0.8882	1.0	0.885	0.893
					<i>Prec</i>	1.0	0.9978	0.9942	1.0	0.9799	0.9806	1.0	0.9727	0.9765
0.5	-	100.0	None	False	<i>F1</i>	1.0	0.8287	0.8519	1.0	0.9191	0.9217	1.0	0.921	0.9248
					<i>Rec</i>	1.0	0.708	0.7453	1.0	0.869	0.8723	1.0	0.8736	0.8785
					<i>Prec</i>	1.0	0.9989	0.9939	1.0	0.9752	0.977	1.0	0.9739	0.9764
1.0	-	100.0	None	False	<i>F1</i>	1.0	0.8162	0.8268	1.0	0.9189	0.9186	1.0	0.9161	0.9214
					<i>Rec</i>	1.0	0.69	0.7084	1.0	0.8692	0.8674	1.0	0.8656	0.872
					<i>Prec</i>	1.0	0.9989	0.9926	1.0	0.9747	0.9763	1.0	0.9729	0.9767
-	-	100.0	None	False	<i>F1</i>	0.9999	0.7814	0.7987	0.9998	0.9057	0.9106	0.9998	0.9116	0.9173
					<i>Rec</i>	0.9998	0.6413	0.6658	0.9999	0.8406	0.8489	0.9998	0.8495	0.8571
					<i>Prec</i>	1.0	1.0	0.9979	0.9998	0.9818	0.982	0.9998	0.9835	0.9866
-	-	100.0	None	True	<i>F1</i>	0.9999	0.7699	0.7666	0.9998	0.9208	0.9243	0.9998	0.9211	0.9269
					<i>Rec</i>	0.9998	0.6264	0.6226	0.9999	0.8637	0.8694	0.9998	0.8675	0.8759
					<i>Prec</i>	1.0	0.9987	0.9972	0.9998	0.9859	0.9865	0.9998	0.9817	0.9841
-	-	100.0	balanced	True	<i>F1</i>	0.9999	0.7506	0.7503	0.9998	0.9064	0.9087	0.9998	0.9071	0.9121
					<i>Rec</i>	1.0	0.6013	0.6018	1.0	0.8385	0.8404	1.0	0.8414	0.8484
					<i>Prec</i>	0.9997	0.9987	0.9959	0.9996	0.9863	0.9891	0.9997	0.9839	0.9863
-	-	100.0	balanced	False	<i>F1</i>	0.9999	0.7475	0.7405	0.9998	0.8894	0.8924	0.9998	0.8913	0.898
					<i>Rec</i>	1.0	0.5973	0.5888	1.0	0.8112	0.8149	1.0	0.8136	0.8227
					<i>Prec</i>	0.9997	0.9987	0.9976	0.9996	0.9843	0.9863	0.9997	0.9853	0.9885
-	0.1	100.0	None	True	<i>F1</i>	1.0	0.6857	0.654	0.9999	0.9268	0.928	0.9999	0.929	0.9345
					<i>Rec</i>	1.0	0.9027	0.9275	1.0	0.92	0.9243	0.9999	0.8981	0.9064
					<i>Prec</i>	1.0	0.5529	0.5051	0.9999	0.9336	0.9317	0.9999	0.9621	0.9644
-	0.5	100.0	None	False	<i>F1</i>	1.0	0.3469	0.3024	1.0	0.8087	0.8157	1.0	0.8581	0.8624
					<i>Rec</i>	1.0	0.9364	0.9712	1.0	0.966	0.9684	1.0	0.9551	0.9566
					<i>Prec</i>	1.0	0.2128	0.1791	1.0	0.6954	0.7046	1.0	0.779	0.785
-	0.5	100.0	balanced	True	<i>F1</i>	1.0	0.2709	0.2357	1.0	0.8111	0.8181	1.0	0.8671	0.871
					<i>Rec</i>	1.0	0.9513	0.9775	1.0	0.967	0.9672	1.0	0.9537	0.9562
					<i>Prec</i>	1.0	0.158	0.134	1.0	0.6985	0.7088	0.9999	0.795	0.7998
-	0.5	100.0	None	True	<i>F1</i>	1.0	0.2583	0.2227	1.0	0.806	0.8151	1.0	0.8632	0.8693
					<i>Rec</i>	1.0	0.9568	0.9782	1.0	0.9692	0.9709	1.0	0.9562	0.9602
					<i>Prec</i>	1.0	0.1493	0.1257	1.0	0.6899	0.7024	1.0	0.7867	0.7941

Table A1.: random forest

A. RQ2: Model Performances

Parameter						RN			CCS			CC			
over-samp.	under-samp.	stages	learn-rate	sub-sample	encode type		Train	Test	Holdout	Train	Test	Holdout	Train	Test	Holdout
-	-	300	0.2	1.0	True	<i>F1</i>	0.9364	0.9196	0.9328	0.8477	0.8419	0.8489	0.9087	0.9008	0.9088
						<i>Rec</i>	0.9015	0.8666	0.9005	0.7926	0.7902	0.7953	0.8767	0.8673	0.8766
						<i>Prec</i>	0.9742	0.9796	0.9675	0.9111	0.9009	0.9102	0.9431	0.937	0.9434
-	-	200	0.2	1.0	True	<i>F1</i>	0.8554	0.8306	0.8499	0.7801	0.773	0.7837	0.8644	0.8567	0.8667
						<i>Rec</i>	0.7695	0.7276	0.765	0.6934	0.6886	0.6992	0.8057	0.7951	0.8081
						<i>Prec</i>	0.9629	0.9676	0.956	0.8916	0.881	0.8915	0.9324	0.9286	0.9344
-	-	200	0.2	0.4	True	<i>F1</i>	0.8554	0.8306	0.8499	0.7801	0.773	0.7837	0.8644	0.8567	0.8667
						<i>Rec</i>	0.7695	0.7276	0.765	0.6934	0.6886	0.6992	0.8057	0.7951	0.8081
						<i>Prec</i>	0.9629	0.9676	0.956	0.8916	0.881	0.8915	0.9324	0.9286	0.9344
-	-	200	0.2	0.7	True	<i>F1</i>	0.8554	0.8306	0.8499	0.7801	0.773	0.7837	0.8644	0.8567	0.8667
						<i>Rec</i>	0.7695	0.7276	0.765	0.6934	0.6886	0.6992	0.8057	0.7951	0.8081
						<i>Prec</i>	0.9629	0.9676	0.956	0.8916	0.881	0.8915	0.9324	0.9286	0.9344
-	-	100	0.2	1.0	True	<i>F1</i>	0.6438	0.6177	0.6441	0.6055	0.5974	0.6021	0.7417	0.7314	0.7437
						<i>Rec</i>	0.4908	0.4584	0.4896	0.4757	0.4699	0.472	0.6241	0.6114	0.6252
						<i>Prec</i>	0.9354	0.9465	0.9412	0.8329	0.8198	0.8314	0.9138	0.9099	0.9175
0.5	-	200	0.2	1.0	True	<i>F1</i>	0.9936	0.4847	0.4165	0.9704	0.6948	0.7007	0.9624	0.7933	0.791
						<i>Rec</i>	0.9984	0.9882	0.9989	0.9843	0.9826	0.9842	0.9723	0.9693	0.9711
						<i>Prec</i>	0.9889	0.3211	0.2631	0.9568	0.5374	0.544	0.9528	0.6714	0.6672
-	0.5	200	0.2	1.0	True	<i>F1</i>	0.9944	0.4545	0.3817	0.9684	0.6836	0.6917	0.9632	0.7925	0.7905
						<i>Rec</i>	0.9977	0.9874	0.9986	0.9816	0.9793	0.9817	0.9733	0.9714	0.972
						<i>Prec</i>	0.9912	0.2952	0.2359	0.9556	0.525	0.534	0.9533	0.6693	0.6661
0.5	-	100	0.2	1.0	True	<i>F1</i>	0.9831	0.2906	0.244	0.945	0.5727	0.5866	0.9363	0.7233	0.7266
						<i>Rec</i>	0.9914	0.9765	0.9954	0.9602	0.9588	0.962	0.938	0.9338	0.9383
						<i>Prec</i>	0.975	0.1707	0.1391	0.9302	0.4083	0.422	0.9347	0.5902	0.5928
-	0.5	100	0.2	1.0	True	<i>F1</i>	0.9843	0.2767	0.2278	0.943	0.5641	0.579	0.9381	0.7234	0.727
						<i>Rec</i>	0.9935	0.9827	0.9958	0.957	0.9525	0.9588	0.9416	0.9373	0.9409
						<i>Prec</i>	0.9751	0.161	0.1286	0.9294	0.4007	0.4147	0.9346	0.589	0.5924
0.5	-	100	0.1	1.0	True	<i>F1</i>	0.9662	0.2047	0.1716	0.9111	0.4684	0.4863	0.8831	0.6059	0.613
						<i>Rec</i>	0.9714	0.9529	0.9775	0.9217	0.9186	0.9242	0.868	0.8636	0.869
						<i>Prec</i>	0.9611	0.1146	0.0941	0.9007	0.3143	0.33	0.8988	0.4667	0.4735
-	0.5	100	0.1	1.0	True	<i>F1</i>	0.9662	0.1901	0.1586	0.9115	0.478	0.4944	0.8853	0.6127	0.6216
						<i>Rec</i>	0.9737	0.9584	0.9814	0.9194	0.9151	0.9215	0.8694	0.8636	0.872
						<i>Prec</i>	0.9587	0.1055	0.0863	0.9038	0.3235	0.3379	0.9019	0.4747	0.4829

Table A2.: gradient boosting classifier

Parameter								RN			CCS			CC		
over-samp.	under-samp.	emb-size	epochs	batch-size	#lstm-cells	dropout		Train	Test	Holdout	Train	Test	Holdout	Train	Test	Holdout
-	-	32	2	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.9936 0.9912 0.996	0.9894 0.9851 0.9937	0.9933 0.9933 0.9933	0.9905 0.9952 0.9858	0.9782 0.9864 0.97	0.981 0.9865 0.9755	0.9876 0.9941 0.9812	0.9774 0.9857 0.9692	0.9778 0.9852 0.9706
-	-	16	2	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.9925 0.9891 0.9958	0.989 0.9843 0.9937	0.9933 0.9933 0.9933	0.9864 0.9959 0.9771	0.9783 0.9906 0.9664	0.9791 0.9886 0.9698	0.9841 0.9763 0.992	0.9724 0.9594 0.9858	0.9734 0.9599 0.9872
-	-	32	3	64	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.9919 0.985 0.9988	0.9881 0.9804 0.996	0.991 0.987 0.995	0.9897 0.9964 0.9831	0.9806 0.9903 0.9712	0.9808 0.9889 0.9728	0.9877 0.9919 0.9835	0.9783 0.9837 0.9728	0.9781 0.9834 0.9728
-	-	32	2	512	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.9933 0.9916 0.9951	0.9878 0.9835 0.9921	0.9916 0.9919 0.9912	0.9905 0.9939 0.9872	0.9806 0.9842 0.977	0.9797 0.9818 0.9776	0.9868 0.996 0.9777	0.976 0.9894 0.963	0.977 0.99 0.9644
-	-	32	3	512	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.9935 0.9889 0.9982	0.9873 0.9796 0.9952	0.9915 0.9887 0.9943	0.992 0.9972 0.9869	0.9813 0.9886 0.9742	0.9809 0.988 0.974	0.9886 0.9967 0.9807	0.9769 0.9887 0.9655	0.9774 0.9885 0.9665
-	-	32	3	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.9945 0.9936 0.9953	0.987 0.9859 0.9882	0.9914 0.9947 0.9881	0.9928 0.992 0.9936	0.9808 0.9795 0.982	0.9805 0.9779 0.9831	0.9906 0.9956 0.9857	0.9782 0.9868 0.9698	0.9783 0.9868 0.97
0.5	-	32	3	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	1.0 1.0 0.9999	0.9859 0.989 0.9828	0.9864 0.9951 0.9779	0.9995 1.0 0.9991	0.9781 0.9814 0.9747	0.9795 0.9826 0.9764	0.9986 0.9999 0.9973	0.9745 0.9853 0.964	0.9746 0.9865 0.9631
0.5	-	32	3	256	100.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	1.0 1.0 1.0	0.9858 0.9796 0.9921	0.9905 0.9905 0.9905	0.9997 1.0 0.9995	0.9785 0.9755 0.9816	0.9777 0.9732 0.9823	0.9993 1.0 0.9986	0.9761 0.9756 0.9765	0.9756 0.9766 0.9747
-	-	16	3	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.9899 0.998 0.982	0.9844 0.9906 0.9783	0.9838 0.9961 0.9719	0.9908 0.9923 0.9893	0.9797 0.9807 0.9787	0.9812 0.9808 0.9815	0.9876 0.9879 0.9874	0.9759 0.9751 0.9767	0.976 0.975 0.9771
-	0.5	32	3	256	100.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.9994 1.0 0.9988	0.8594 1.0 0.7534	0.818 1.0 0.6921	0.9986 0.9998 0.9974	0.9555 0.9985 0.916	0.9548 0.9981 0.9152	0.998 0.9994 0.9966	0.9707 0.9961 0.9465	0.9722 0.9965 0.9491
-	0.5	32	3	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.9978 1.0 0.9956	0.6294 1.0 0.4593	0.5796 1.0 0.4081	0.9985 0.999 0.9979	0.9564 0.9958 0.9201	0.957 0.9965 0.9206	0.9974 0.9992 0.9955	0.9561 0.9963 0.919	0.9588 0.9957 0.9245

Table A3.: lstm

A. RQ2: Model Performances

Parameter						RN			CCS			CC			
over-samp.	under-samp.	kernel	sub-sample	class weight	encode type		Train	Test	Holdout	Train	Test	Holdout	Train	Test	Holdout
-	0.5	rbf	0.5	None	True	<i>F1</i>	0.0364	0.0066	nan	nan	nan	nan	nan	nan	nan
						<i>Rec</i>	0.0187	0.0134	nan	nan	nan	nan	nan	nan	nan
						<i>Prec</i>	0.6552	0.0044	nan	nan	nan	nan	nan	nan	nan
-	0.5	rbf	0.5	balanced	True	<i>F1</i>	0.5257	0.0052	nan	nan	nan	nan	nan	nan	nan
						<i>Rec</i>	0.8069	0.766	nan	nan	nan	nan	nan	nan	nan
						<i>Prec</i>	0.3898	0.0026	nan	nan	nan	nan	nan	nan	nan
-	0.3	rbf	0.5	balanced	True	<i>F1</i>	0.4141	0.0052	nan	nan	nan	nan	nan	nan	nan
						<i>Rec</i>	0.8353	0.7877	nan	nan	nan	nan	nan	nan	nan
						<i>Prec</i>	0.2753	0.0026	nan	nan	nan	nan	nan	nan	nan
-	0.1	rbf	0.5	None	True	<i>F1</i>	0.0	0.0	nan	nan	nan	nan	nan	nan	nan
						<i>Rec</i>	0.0	0.0	nan	nan	nan	nan	nan	nan	nan
						<i>Prec</i>	0.0	0.0	nan	nan	nan	nan	nan	nan	nan
-	0.5	rbf	0.3	balanced	True	<i>F1</i>	nan	nan	nan	0.4476	0.0622	nan	0.4685	0.1115	nan
						<i>Rec</i>	nan	nan	nan	0.5129	0.5004	nan	0.5647	0.5535	nan
						<i>Prec</i>	nan	nan	nan	0.3971	0.0332	nan	0.4004	0.062	nan
-	0.3	rbf	0.3	balanced	True	<i>F1</i>	nan	nan	nan	0.3676	0.0629	nan	0.3807	0.1124	nan
						<i>Rec</i>	nan	nan	nan	0.5151	0.501	nan	0.5622	0.5497	nan
						<i>Prec</i>	nan	nan	nan	0.2857	0.0336	nan	0.2878	0.0626	nan
-	0.5	rbf	0.3	None	True	<i>F1</i>	nan	nan	nan	0.0318	0.0208	nan	0.0267	0.0189	nan
						<i>Rec</i>	nan	nan	nan	0.0164	0.0134	nan	0.0137	0.0106	nan
						<i>Prec</i>	nan	nan	nan	0.5553	0.0461	nan	0.5741	0.0829	nan
-	0.1	rbf	0.3	None	True	<i>F1</i>	nan	nan	nan	0.0	0.0	nan	nan	nan	nan
						<i>Rec</i>	nan	nan	nan	0.0	0.0	nan	nan	nan	nan
						<i>Prec</i>	nan	nan	nan	1.0	0.0	nan	nan	nan	nan

Table A4.: svm. NaN value for configuration that we did not test due to time constrains.

B. RQ3: Model Performances on Manipulated Data

B. RQ3: Model Performances on Manipulated Data

Parameter					RN	CC	Base CC	
over-samp.	under-samp.	#trees	class weight	encode type				
-	0.5	100.0	balanced	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0539 0.3022 0.0296	0.7158 0.758 0.678	0.8186 0.9782 0.7038
-	0.5	100.0	None	False	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0509 0.2099 0.0289	0.6829 0.6899 0.6761	0.8221 0.9783 0.709
-	0.5	100.0	None	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0507 0.3033 0.0277	0.7148 0.7626 0.6726	0.812 0.9796 0.6933
-	0.1	100.0	None	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0031 0.0038 0.0027	0.3411 0.21 0.9091	0.9419 0.9207 0.9641
0.5	-	100.0	balanced	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0295 0.015 0.9019	0.9415 0.9072 0.9784
-	-	100.0	balanced	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0122 0.0062 0.8855	0.8954 0.8165 0.9911
1.0	-	100.0	None	False	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0157 0.0079 0.8327	0.8855 0.8059 0.9825
-	-	100.0	None	False	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0036 0.0018 0.8095	0.8709 0.7754 0.9931
1.0	-	100.0	None	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0509 0.0261 0.9295	0.9355 0.8947 0.9802
0.5	-	100.0	None	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0317 0.0161 0.895	0.9453 0.9124 0.9806
-	-	100.0	None	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0103 0.0052 0.8567	0.9341 0.8858 0.9879
-	-	100.0	balanced	False	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0039 0.0019 0.8527	0.7969 0.665 0.9942
0.5	-	100.0	None	False	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0141 0.0071 0.883	0.8984 0.8266 0.9838

Table A1.: Manipulated code for random forest

Parameter							RN	CC	Base CC
over-samp.	under-samp.	stages	learn-rate	sub-sample	encode type				
0.5	-	200	0.2	1.0	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.24 0.7676 0.1422	0.83 0.978 0.7209	0.7519 0.9854 0.6078
-	0.5	200	0.2	1.0	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.2309 0.8196 0.1344	0.8291 0.9773 0.7199	0.7466 0.9853 0.601
-	-	100	0.2	1.0	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.15 0.0891 0.4729	0.5254 0.3629 0.9513	0.8648 0.8225 0.9117
0.5	-	100	0.2	1.0	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.1488 0.8706 0.0814	0.7314 0.9781 0.5841	0.6883 0.9789 0.5307
-	0.5	100	0.2	1.0	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.1435 0.9055 0.0779	0.7303 0.9783 0.5827	0.6874 0.9817 0.5288
0.5	-	100	0.1	1.0	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.1063 0.891 0.0565	0.6377 0.9758 0.4736	0.5957 0.977 0.4285
-	0.5	100	0.1	1.0	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0993 0.9077 0.0525	0.6454 0.9786 0.4815	0.607 0.9769 0.4402
-	-	200	0.2	0.4	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0963 0.0558 0.349	0.3292 0.1995 0.9406	0.8948 0.8805 0.9096
-	-	200	0.2	1.0	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0963 0.0558 0.349	0.3292 0.1995 0.9406	0.8948 0.8805 0.9096
-	-	200	0.2	0.7	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0963 0.0558 0.349	0.3292 0.1995 0.9406	0.8948 0.8805 0.9096
-	-	300	0.2	1.0	True	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0563 0.0322 0.2222	0.2333 0.1331 0.942	0.9174 0.9105 0.9243

Table A2.: Manipulated code gradient boosting classifier

B. RQ3: Model Performances on Manipulated Data

Parameter							RN	CC	Base CC	
over-samp.	under-samp.	emb-size	epochs	batch-size	#lstm-cells	dropout				
-	0.5	32	3	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.3895 0.847 0.2529	0.6568 0.5009 0.9534	0.9502 0.996 0.9085
-	-	32	2	512	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.3195 0.1902 0.9968	0.9694 0.9931 0.9467
-	-	16	3	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.2831 0.165 0.9956	0.9698 0.9737 0.9659
0.5	-	32	3	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.119 0.0633 0.9851	0.9692 0.9947 0.945
-	-	32	2	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.1898 0.1049 0.9938	0.9728 0.988 0.958
-	-	32	3	512	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0037 0.0019 0.7047	0.9714 0.9929 0.9508
-	-	32	3	64	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0857 0.0448 0.9742	0.9721 0.9788 0.9655
-	-	32	3	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.0451 0.0231 0.9645	0.9711 0.9937 0.9494
-	-	16	2	256	10.0	0.2 0.2	<i>F1</i> <i>Rec</i> <i>Prec</i>	0.0 0.0 0.0	0.5738 0.4036 0.9925	0.972 0.9676 0.9765

Table A3.: Manipulated code lstm

C. Tokenizer Types

0: 'ENDMARKER'	31: 'TILDE'
1: 'NAME'	32: 'CIRCUMFLEX'
2: 'NUMBER'	33: 'LEFTSHIFT'
3: 'STRING'	34: 'RIGHTSHIFT'
4: 'NEWLINE'	35: 'DOUBLESTAR'
5: 'INDENT'	36: 'PLUSEQUAL'
6: 'DEDENT'	37: 'MINEQUAL'
7: 'LPAR'	38: 'STAREQUAL'
8: 'RPAR'	39: 'SLASHEQUAL'
9: 'LSQB'	40: 'PERCENTEQUAL'
10: 'RSQB'	41: 'AMPEREQUAL'
11: 'COLON'	42: 'VBAREQUAL'
12: 'COMMA'	43: 'CIRCUMFLEXEQUAL'
13: 'SEMI'	44: 'LEFTSHIFTEQUAL'
14: 'PLUS'	45: 'RIGHTSHIFTEQUAL'
15: 'MINUS'	46: 'DOUBLESTAREQUAL'
16: 'STAR'	47: 'DOUBLESASH'
17: 'SLASH'	48: 'DOUBLESASHEQUAL'
18: 'VBAR'	49: 'AT'
19: 'AMPER'	50: 'ATEQUAL'
20: 'LESS'	51: 'RARROW'
21: 'GREATER'	52: 'ELLIPSIS'
22: 'EQUAL'	53: 'OP'
23: 'DOT'	54: 'AWAIT'
24: 'PERCENT'	55: 'ASYNC'
25: 'LBRACE'	56: 'ERRORTOKEN'
26: 'RBRACE'	57: 'COMMENT'
27: 'EQEQUAL'	58: 'NL'
28: 'NOTEQUAL'	59: 'ENCODING'
29: 'LESSEQUAL'	256: 'NT_OFFSET'
30: 'GREATEREQUAL'	

Figure C.1.: Possible token types of the tokenizer for our type encoding. Extracted from the `tokenize.tok_name` attribute.

D. Dataset

Organisation	Repository	Branch	Git Hash
scikit-learn	scikit-learn	master	8ba49f628092fe3fb1deea101910006aa6c76d49
pytorch	pytorch	master	2b70f827373fe5c90e15ce7ba0a1d331dcd57379
ansible	ansible	devel	9792d631b1b92356d41e9a792de4b2479a1bce44
django	django	master	9c92924cd5d164701e2514e1c2d6574126bd7cc2
pallets	flask	master	6d3f87ee071ea061f0af2ba7d29b148841090b59
keras-team	keras	master	1a3ee8441933fc007be6b2beb47af67998d50737
jakubroztocil	httpie	master	100872b5cf493b355453a9dff89f39b0c124d524
home-assistant	core	dev	dc84196202a4683881ddac9942fe1c9af6de3b5b
certbot	certbot	master	2a047eb526ee02d4bb3bb3668260ac6d481b99d3
pandas-dev	pandas	master	e666f5cc59a8631a5ee3e38c20ad8e08531cbec3
docker	compose	master	76963e44add9810c1d906b5fbd052dc3fb480479
tiangolo	fastapi	master	a6897963d5ff2c836313c3b69fc6062051c07a63
tornadoweb	tornado	master	4a4cf9ce5867d29b4557a0b28eca7a56e77207de
encode	django-rest-framework	master	7f3a3557a050147dd2420aa41d1bf7ddd7f9818e
google	python-fire	master	79d85df2706b2dab5dd5075f2a76953743ee9bf3
explosion	spaCy	master	7d7b65ffd42c56ce3a0aa73b18196eb20a1dcc24
bokeh	bokeh	master	66231e91c3a806de73d77f7c3617e028d28e7c6c
apache	airflow	master	3a349624a20d3432dc75e337d6ffb1109a50e451

Table A1.: Open-source repositories we used for our dataset.