

Final Project

An Empirical Study of the Maintainability of ChatGPT's answers

Thomas Trépanier

Montréal, Canada

thomas.trépanier@polytml.ca

Abstract—

Index Terms—llm, maintainability, ChatGPT, software repository mining

I. INTROCUCTION

II. RELATED WORK

III. METHODOLOGY

The methodology to answer *RQ1* and *RQ2* is mostly the same and can be divided into 5 steps: (1) data collection, (2) code snippet extraction, (3) code snippets filtering, (4) code smell detection, and (5) code smell analysis. The main difference between the two research questions is the data collection step. For *RQ1*, the data used is the DevGPT dataset provided by the MSR2024 challenge [1]. For *RQ2*, a dataset was created using StackOverflow's API following the methodology described in section III-A.

All of the code and algorithms to extract and process data was written in Python. The code is available on GitHub ¹.

A. Data Collection

1) *DevGPT Dataset*: All of the provided snapshots were used in the study, as well as all of the JSON files in each of the snapshots. The JSON files contain information about the conversations between the developers and the ChatGPT bot.

2) *StackOverflow Dataset*: In order to build a dataset of StackOverflow's code snippets, the following approach was used:

- 1) **Fetching Questions from StackOverflow's API**: The first step was to query StackOverflow's API to get a list of Python related questions asked between November 23rd 2018 and November 23rd 2023. These dates were chosen as they represent a recent period of time that would correspond to the kind of code snippets found in the DevGPT dataset.

Furthermore, the questions had to have at least five answers with one of them being accepted, and having been viewed 1000 times. These criteria were chosen to ensure that the questions were relevant and that

the code snippets found in the answers would be of relatively good quality.

- 2) **Retrieving the answers to the questions**: The second step was to fetch the answers to the questions retrieved in step 1. To do so, the *question_id* of each question was used to query StackOverflow's API and retrieve the answers associated with the question. The answers were then stored in a JSON file.

B. Code Snippet Extraction

1) *DevGPT Dataset*: The DevGPT dataset contains conversations between developers and the ChatGPT bot. Each JSON file contains a data structure called `ChatgptSharing`, itself containing a list of `Conversation`. Each `Conversation` contains the prompt given by the developer, the answer given by ChatGPT, and a data structure called `ListOfCode` containing the code snippets included in ChatGPT's answer.

This `ListOfCode` itself provides the code snippet emitted by ChatGPT as well as the programming language used in the snippet.

In order to build a dataset of all of the code snippets generated by ChatGPT, a program was written to parse all of the JSON files of all of the snapshots and extract all the `ListOfCode` data structures.

2) *StackOverflow Dataset*: The answers retrieved from StackOverflow's API consisted of HTML bodies containing HTML elements to delimit the code snippet. Therefore, the HTML bodies were parsed using the following Regexp to extract the code snippet itself: `<pre><code>(.*?)</code></pre>`.

The snippet was then parsed from its HTML format to a Unicode format and stored in a JSON file.

C. Code Snippets Filtering

Not all code snippets extracted were used for the code smell detection, since not all of them were suited for the analysis. The definition of a *suitable* code snippet is as follows:

¹<https://github.com/ThomasTrepanier/log6307-final-project>

- 1) The code snippet must be in Python. This is because the code smell detection tool used, *PyScent*, is a Python tool.
- 2) The code snippet must be at least 5 lines long, and contain either a function (defined with the `def` keyword) or a class (defined with the `class` keyword). This is to filter out short answers or code snippets that were not actually Python code, such as JSON objects or library import statements.
- 3) The code snippet must be written entirely in english. This is because non-unicode characters are not supported by *PyScent*.
- 4) The code snippet must not contain any syntax errors. This is because *PyScent* is not able to detect code smells in code that does not compile. To remove the snippets with syntax errors, a program was written to delete export each code snippet to a `.py` file and try to compile it. If a `SyntaxError` or `TypeError` was thrown by Python, the file was assumed to be erroneous and deleted from the dataset.

D. Code Smell Detection

The last step was to detect code smells in the code snippets. To do so, the *PyScent* tool was used. *PyScent* is a Python tool that detects code smells in Python code and the source code can be found on Github ².

PyScent detects the 11 following code smells:

- **Long Method:** A method with over 50 statements.
- **Long Parameter List:** A method that takes more than 5 arguments.
- **Long Branches:** An if statement with more than 12 branches.
- **Too Many Attributes:** A class that has more than 7 attributes.
- **Too Many Methods:** A class that has more than 20 public methods.
- **Useless Try/Except Clauses:** A try/except clause that catches too many exception (using the `Exception` or `Standard Error`), or whose catch clauses are empty.
- **Shotgun Surgery:** A Class that calls more than 5 methods from other classes.
- **Class Cohesion:** A class that has cohesion score lower 30%. The cohesion score is computed using the Cohesion project on Gitbut ³.
- **Code Complexity:** A block that has a cyclomatic complexity score lower than 'C'. The cyclomatic complexity score is computed using the Radon project ⁴.
- **Long Lambda:** A lambda function that has more than 60 characters.

- **Long List Comprehension:** A list comprehension that has more than 72 characters.

To analyse the code snippets, *PyScent* needs them to be in `.py` files under one directory. Therefore, a program was written to export all of the code snippets to `.py` files and run *PyScent* on them.

The results of the code smell detection are presented in section IV.

IV. RESULTS

V. THREATS TO VALIDITY

VI. THREATS TO VALIDITY

VII. CONCLUSION

REFERENCES

- [1] NAIST-SE. "DevGPT: Studying Developer-ChatGPT Conversations". MSR-2024 Challenge. <https://github.com/NAIST-SE/DevGPT>

²<https://github.com/whyjay17/Pyscent>

³<https://github.com/mschwager/cohesion>

⁴<https://pypi.org/project/radon/>