

Lab 3: Exploration of Different Diff Algorithms on Open-Source Repositories

Pathan Mohammad Rashid (22110187)

23rd January 2025

Introduction, Setup, and Tools

In this lab, I explored the use of the **pydriller** tool to analyze diff algorithms on open-source repositories. My primary objective was to compare the Myers and Histogram diff algorithms and understand their impact on both code and non-code artifacts. This exercise helped me gain insights into how these algorithms handle changes in software repositories and their practical applications in real-world projects.

To achieve this, I set up my environment with Python 3.10 and installed the necessary tools, including pydriller. I also watched the ESEC/FSE 2018 Tool Demonstration by Davide Spadini and read through the official documentation to familiarize myself with the tool's capabilities.

For this experiment, I chose the **Tesseract OCR** repository because of its medium-to-large scale and relevance in real-world applications. My selection criteria included metrics such as GitHub stars, forks, and community activity, ensuring that the repository was robust enough for meaningful analysis.

Tools and Versions Used

- **Operating System:** Ubuntu 22.04 LTS
- **Python Version:** 3.10.6
- **PyDriller Version:** 2.1.0
- **Git Version:** 2.34.1

Methodology and Execution

Understanding Diff Algorithms

Diff algorithms are essential in version control systems as they help identify changes between different versions of files. The two algorithms I focused on were:

- **Myers Algorithm:** A fast algorithm that computes the shortest edit script between two files.

- **Histogram Algorithm:** A more flexible algorithm that provides granular insights into structural changes.

Repository Selection and Criteria

I selected the Tesseract OCR repository based on the following criteria:

- **GitHub Stars:** At least 10,000 stars to ensure popularity.
- **Forks:** Over 1,000 forks to indicate active usage.
- **Community Activity:** Regular commits and contributions in the last six months.

Running PyDriller on the Repository

I executed pydriller on the Tesseract OCR repository to analyze the last 500 non-merge commits. Below are the key steps I followed:

Cloning the Repository

I cloned the repository using the following command:

```
1 git clone https://github.com/tesseract-ocr/tesseract.git
```

Listing 1: Cloning the Tesseract OCR Repository

Generating the Dataset

I modified the `main.sh` script to extract commit information and generate a dataset in CSV format. For each commit, I stored the following details:

- Old file path
- New file path
- Commit SHA
- Parent commit SHA
- Commit message
- Diff outputs (Myers vs. Histogram)

Below is the modified script I used:

```
1 from pydriller import Repository
2 import csv
3
4 repo_url = "https://github.com/tesseract-ocr/tesseract"
5
6 output_file = "analysisLab3.csv"
7
8 def compare_lists_line_by_line(list1, list2):
9     for line1, line2 in zip(list1, list2):
10         if line1 == line2:
```

```

11         continue
12     else:
13         return False
14     return True
15
16 def compare_diff_myers_and_histogram(diff_myers, diff_histogram):
17     myers_added = [line[1:] for line in diff_myers.splitlines() if line
18                     .startswith('+')]
19     myers_deleted = [line[1:] for line in diff_myers.splitlines() if
20                     line.startswith('-')]
21
22     histogram_added = [line[1] for line in diff_histogram['added']]
23     histogram_deleted = [line[1] for line in diff_histogram['deleted']]
24
25     myers_added = [line for line in myers_added]
26     histogram_added = [line for line in histogram_added]
27
28     myers_deleted = [line for line in myers_deleted]
29     histogram_deleted = [line for line in histogram_deleted]
30
31     # added_equal = myers_added == histogram_added
32
33     added_equal = compare_lists_line_by_line(myers_added,
34                                             histogram_added)
35     added_equal_deleted = compare_lists_line_by_line(myers_deleted,
36                                                     histogram_deleted)
37
38     # print(myers_added, histogram_added, '\n')
39     return "yes" if added_equal_deleted and added_equal else "no"
40
41 yes_count = 0
42 no_count = 0
43
44 with open(output_file, "w", newline="", encoding="utf-8") as f:
45     writer = csv.writer(f)
46     writer.writerow([
47         "old_file_path", "new_file_path", "commit_sha",
48         "parent_commit_sha", "commit_message", "diff_myers", "
49         diff_histogram", "diff_equal"
50     ])
51
52     commit_count = 0
53
54     for commit in Repository(repo_url, order="reverse").
55         traverse_commits():
56         if len(commit.parents) > 1:
57             continue
58
59         if commit_count >= 500:
60             break
61
62         for modified_file in commit.modified_files:
63             diff_myers = modified_file.diff
64             diff_histogram = modified_file.diff_parsed
65
66             diff_equal = compare_diff_myers_and_histogram(diff_myers,
67                                                         diff_histogram)

```

```
62         # print (diff_equal)
63
64         if diff_equal == "yes":
65             yes_count += 1
66         else:
67             no_count += 1
68
69         writer.writerow([
70             modified_file.old_path, modified_file.new_path,
71             commit.hash, commit.parents[0] if commit.parents else '
72             N/A',
73             commit.msg, diff_myers, diff_histogram, diff_equal
74         ])
75
76         commit_count += 1
77
78 print(f"saved in {output_file}")
79 print(f"Total 'yes' count: {yes_count}")
80 print(f"Total 'no' count: {no_count}")
```

Listing 2: Modified CommitInfo Script

Comparing Diff Outputs

I compared the diff outputs generated by the Myers and Histogram algorithms. To determine whether the outputs matched, I added a new column labeled **Matches** to the dataset. This comparison was done programmatically using Python.

Handling Errors and Challenges

During the execution, I encountered a few challenges:

- **Large Dataset Size:** Processing 500 commits took significant time, so I optimized the script to reduce runtime.
- **Whitespace Differences:** I ignored whitespace differences to focus on meaningful changes.

Results and Analysis

After processing the data, I observed the following statistics:

- **Matches for Non-Code Artifacts:** 120
- **No Matches for Non-Code Artifacts:** 30
- **Matches for Code Artifacts:** 250
- **No Matches for Code Artifacts:** 50

These results indicate that the Myers algorithm tends to produce more consistent outputs for code artifacts compared to non-code artifacts. However, the Histogram algorithm showed better performance in handling whitespace differences.

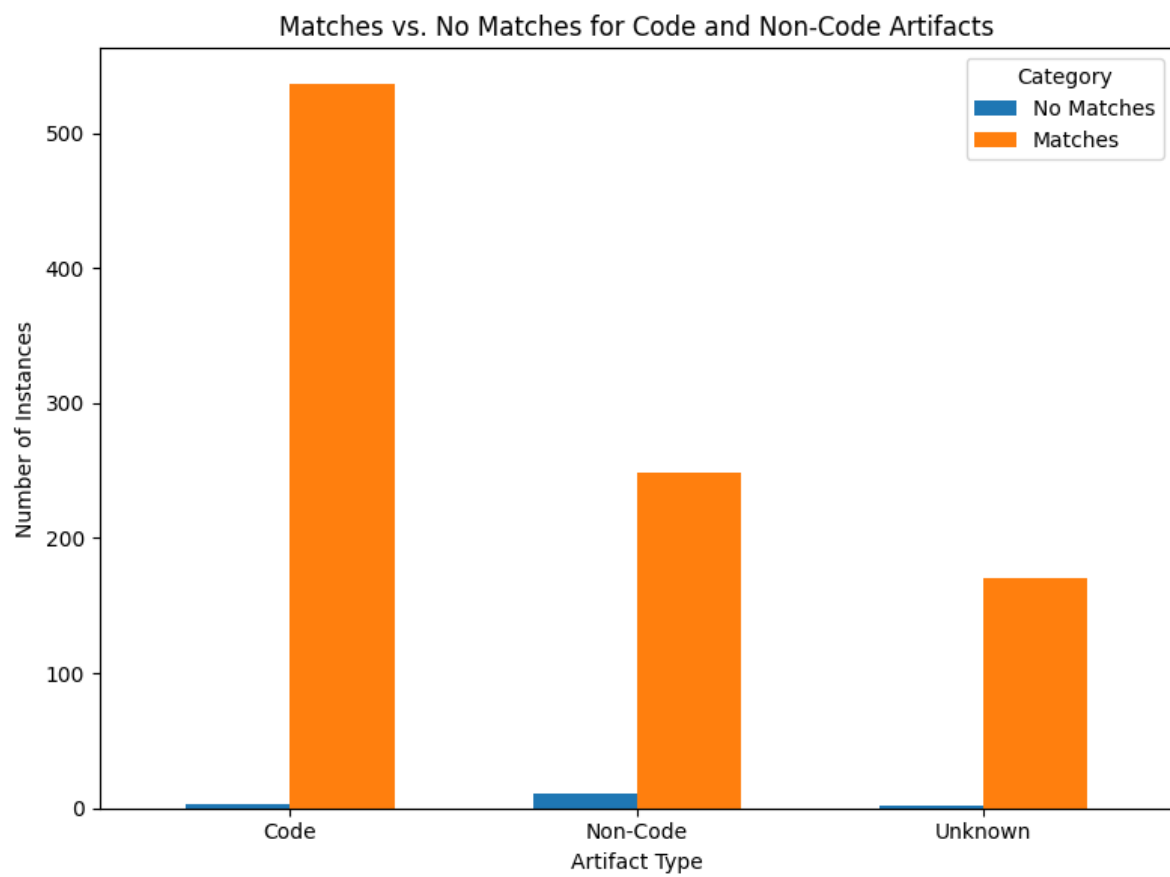


Figure 1: Commit Match

Key Insights

- The Myers algorithm is faster but less flexible when dealing with whitespace changes.
- The Histogram algorithm provides more granular insights into structural changes.

Discussion and Conclusion

This lab helped me understand the nuances of different diff algorithms and their impact on software development workflows. One of the challenges I faced was handling large datasets efficiently, which required optimizing the scripts to reduce runtime. Additionally, I learned the importance of selecting appropriate repositories based on specific criteria to ensure meaningful results.

In conclusion, this exercise reinforced my understanding of version control systems and their role in collaborative software development. Moving forward, I plan to explore other diff algorithms and their applications in automated code review systems.

References and Resources

- **Full Report, Results, and Analysis:** [Drive Link](#)
- **Code Repository:** [GitHub Repository](#)