

LEWIS UNIVERSITY

Structural Quality & Software Evolution

A Thesis

By

Alison Major

Department of Computer and Mathematical Sciences

Submitted in partial fulfillment of the requirements

for the degree of
Master of Science in Computer Science,

Concentration in Software Engineering

April 30, 2022

The undersigned have examined the thesis entitled ‘**Structural Quality & Software Evolution**’ presented by **Alison Major**, a candidate for the degree of **Master of Science in Computer Science (Concentration in Software Engineering)** and hereby certify that it is worthy of acceptance.

Advisor

Date

Program Director

Date

Department Chair

Date

Abstract

Some software engineering projects fail to evolve, making them obsolete. When we build new software solutions, we consider several categories during planning: cost, time-to-deliver, and sometimes to a smaller extent, quality. Though quality can be a factor we consider when planning a project, it is a complex attribute to understand and measure. We care about quality, as it can impact the longevity of a project. Low-quality software can be complicated to enhance and evolve. Software that fails to evolve will fail to generate user engagement, leading to revenue loss. Many tools are available for enforcing standards, some built into integrated development environments (IDEs) (like Visual Studio Code) and others as third-party linter tools. Two such tools are Pylint and Radon, which can analyze Python projects. In particular, we will focus on the refactor violations noted in Pylint and Radon's code reports, as these warnings can lead us to code smells. We review some projects and resources to understand the correlation between software structure quality and its impacts on a system's evolving ability. For example, we find that our reviewed projects have a low count of refactoring messages and middle to high maintainability indexes, indicating that these measurements may help estimate the quality of a system. With this understanding, we explore ways to improve the evolution of a software system through tools and suggestions.

Acknowledgements

This research and working towards my degree has been an endeavor and one that I would not dream of tackling on my own.

First, my thanks go to God, the ultimate creator, who formed us all to be creators fitting our skills. And to His son, Jesus, who gave his life so that the sins of the world may be forgiven, and rose from the dead and now sits at the right hand of God.

Second, much gratitude goes to my very creative and patient family. I am thankful for my husband, Chris, who was foundational in finding the time and focus on working on my studies. My amazing kids, Ewan and Gwynnie, have provided patience, encouragement, and numerous trips to the library. I could not have gotten through this without all of your support.

Finally, I am grateful for my coworkers and my professors. The Robotic Process Automation (RPA) team at Sysco has been cheering me on from the sidelines. The professors at Lewis University have provided mentorship and guidance. They have all encouraged me as I've worked to apply my knowledge from the last decade of my experience as a software developer and gather new information that has helped us build better software solutions.

Contents

Abstract	ii
Acknowledgements	iii
List of Tables	vi
List of Figures	viii
1 Introduction	1
1.1 Maintainability Index & Refactor Scores	4
1.2 Paper Structure	6
2 Background and Literature Review	8
2.1 Keeping Users Engaged Long Term	9
2.1.1 Why does software evolution matter?	9
2.1.2 How do we ensure software evolution?	11
2.2 The Impact of Structural Quality	12
2.2.1 Software Maintenance	12
2.2.2 Software Evolution	14
2.2.3 Measuring Maintainability	18
2.2.4 Maintainability Scores	20
2.2.5 Other Maintainability Characteristics	22
2.2.6 Documentation and Maintainability	25
2.3 Related Work	27
2.3.1 Considering Data Sets	27
2.3.2 Design Patterns and Software Quality	28
2.3.3 Software Architecture and Maintainability	29
3 Methodology	31
3.1 Initial Repository Set	33
3.2 Filtered Repository Set	35

4	Results	37
5	Conclusions and Recommendations	44
	References	48
	Appendix	53
7.1	Pylint: Refactor Scores	53
7.2	Radon: Maintainability Index	54
7.3	Data: Repository Refactor Messages	55
7.4	Data: Repository Refactor-to-SLOC Ratio	57
7.5	Data: Repository Refactor By Message	60
7.5.1	Raven-Python: Most Common Refactor Messages . . .	61
7.5.2	Scrapy: Most Common Refactor Messages	62
7.5.3	Numba: Most Common Refactor Messages	63
7.5.4	Cython: Most Common Refactor Messages	64
7.5.5	YouTube-DL: Most Common Refactor Messages	65
7.5.6	Electrum: Most Common Refactor Messages	66
7.5.7	Sympy: Most Common Refactor Messages	67
7.5.8	Ansible: Most Common Refactor Messages	68
7.5.9	Salt: Most Common Refactor Messages	69
7.6	Data: Pylint Scores for Best & Worst	70
7.6.1	Cython Pylint Scores	71
7.6.2	Raven-Python Pylint Scores	72

List of Tables

3.1	Criteria used to filter down the initial set of repositories. . . .	36
3.2	List of the 46 repositories for research focus.	36
4.1	Total refactor messages per repository for the 3 best and 3 worst offenders. Also provided with the refactor message that had the most warnings and its total count (and the percentage of that message from the total refactor warnings for that repository).	38
4.2	Total refactor messages per repository, total project source line of code (SLOC), and the ratio of refactor messages to SLOC. .	39
4.3	Total refactor messages per repository, total project source line of code (SLOC), the ratio of refactor messages to SLOC, average Maintainability Index (MI).	40
7.1	Total refactor messages per repository. Also provided with the refactor message that had the most warnings and its total count (and the percentage of that message from the total refactor warnings for that repository).	56
7.2	Total refactor messages per repository, total project source line of code (SLOC), the ratio of refactor messages to SLOC, average Maintainability Index (MI) with standard deviation, and code comment-to-SLOC ratio.	59
7.3	The 6 most common refactor messages for Raven-Python. . . .	61
7.4	The 10 most common refactor messages for Scrappy.	62
7.5	The 10 most common refactor messages for Numba.	63
7.6	The 10 most common refactor messages for Cython.	64
7.7	The 10 most common refactor messages for YouTube-DL. . . .	65
7.8	The 10 most common refactor messages for Electrum.	66
7.9	The 10 most common refactor messages for SymPy.	67

7.10	The 10 most common refactor messages for Ansible.	68
7.11	The 10 most common refactor messages for Salt.	69
7.12	Pylint score from various releases of Cython, the “best” ranked repository for maintainability.	71
7.13	Pylint score from various releases of Raven-Python, the “worst” ranked repository for maintainability.	72
7.14	The latest Pylint score from Sentry-Python, which replaced Raven-Python.	72

List of Figures

1.1	Frequency distribution of ISO 25010 product quality characteristics in OSS quality models, as found by Adewumi et. al in 2016 [1].	3
2.1	A simplified visual of Lehman’s fifth law, “Conservation of Familiarity.”	17
2.2	The eight product quality characteristics and sub-characteristics.	30
4.1	Diligent development communities can keep refactor warnings low, regardless of system size (lines of code).	42
4.2	Many repositories average in the mid-score to high-score. Radon considers 20 points and up to be very maintainable.	43
5.1	A snapshot of the badges from <i>SymPy</i> ’s repository.	45

Chapter 1

Introduction

We have several areas of concern when building software systems: cost, delivery timeline, and quality. The cost and time-to-market are often the two problems given the highest priority in a project. However, engineers must consider the software quality to preserve the system's longevity. Despite their importance, the code and architecture quality can be challenging to understand and measure.

The art of programming began as something very tedious and error-prone, with developers writing every step of a process in a machine language. As writing machine code was not a great way to grow the industry, developers created programming languages. With human-friendly languages, programmers could reduce the manual effort in writing code and improve the quality of their work by allowing the machine to transform the high-level code into something executable and efficient. [2]

With the creation of programming languages, the ability to build systems and applications has become a skill that almost anyone can learn. Open source projects create communities of learning and advancement. However, all systems need structure and an investment of time.

When we think about programming projects, we can assume that as time goes on and changes and additions occur within a system's source code, the complexity of that system will grow. However, when we manage the code structure, we can keep the complexity in check, allowing systems to evolve. For example, developers can maintain this structure through simple steps like having readable code. They can also consider more complex aspects, like how coupled and cohesive a system is.

One way to understand the quality of a system is to discuss its “maintainability,” the ease of receiving new features or resolving bugs. For example, developers may find that adjusting one area to add a new feature requires touching several other code areas in tightly coupled systems. Some code measuring systems provide a Maintainability Index (MI), a well-known quality measure. While the effectiveness of quantifying software quality is debatable, many quality models include MI measurements [3], [1].

Adewumi et al. found that 55% of the existing Open Source Software (OSS) quality models measure maintainability, making it the most common quality characteristic measured [1]. From “Fig. 1.1,” we can infer that maintainability is more important than functional stability. If the code is maintainable and accessible, missing features can be incorporated. However, it

would be hard to implement if the code is not well documented and is difficult to read and understand [1].

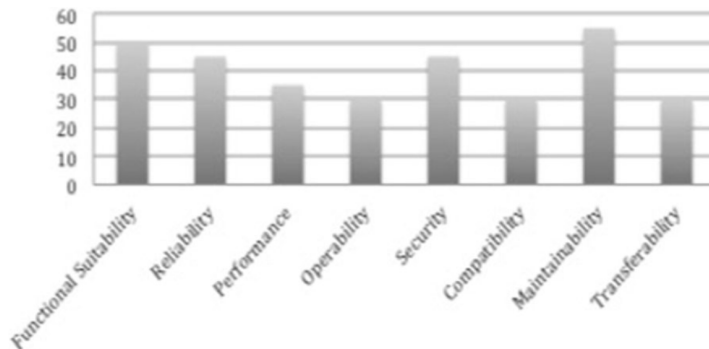


Figure 1.1: Frequency distribution of ISO 25010 product quality characteristics in OSS quality models, as found by Adewumi et. al in 2016 [1].

Code smells are used extensively by practitioners to identify low-quality spots in the software system. These areas would need the teams' attention and are good candidates for refactoring. Rather than building comprehensive models with all possible software characteristics, developers should instead focus on the essential characteristics: maintainability, usability, and maintenance capacity of a software community [1].

1.1 Maintainability Index & Refactor Scores

“Software maintainability is one of the fundamental external quality attributes and is recognised as a research area of primary concern in software engineering [4].” As such, we need to find ways to measure the maintainability of a software system. As mentioned earlier, many OSS quality models already include maintainability as a characteristic measurement.

Pylint is a static analysis tool that identifies several classes of code quality concerns. Refactor violations are relevant to our study, which reports various code smells. However, we can assume there must be some correlation between the Maintainability Index and the type and number of code smells in a software system, quantified by the Pylint refactor score.

To compute various code metrics, Radon is another Python tool available [5]. One such score that Radon provides is the Maintainability Index (MI), which tells us how easy it should be to support and change the program. Radon calculates this value using a formula that consists of SLOC (Source Lines Of Code), Cyclomatic Complexity, and Halstead volume.

This study explores such assumptions and systematically investigates any correlation between Radon’s Maintainability Index metric and the Pylint Refactor message count. Furthermore, we perform analysis on specific refactor violations to reveal and shed light on the relative effectiveness of the different refactor violations and their relationship to Maintainability Index.

The structural quality of a software system will impact the software evolu-

tion. If the project has poor structural quality, the architecture will minimize its ability to evolve, and the software system will eventually “die off” so to speak.

We will look at many open-source Python systems using Pylint and Radon and attempt to correlate the data from the scores to the level of ease in adding new features to the system. The Pylint and Radon results will determine if a system is more maintainable with better scores.

1.2 Paper Structure

In Chapter 2, we will dig into a deeper background of the topic, exploring ideas of why software systems need to keep users engaged long term (Section 2.1). We will explore automated measurements that provide evaluation scores of software systems. By using some of these quality and maintainability scores, we can see how structure impacts evolution (Section 2.2). Additionally, we will explain how quality is measured and the different attributes that can factor into maintainability. We will also review related works (Section 2.3).

With more background on the problem, we can then review the methodology for our research in Chapter 3. Here we will review where we found our initial data set (Section 3.1) and what criteria we used to filter it to a manageable size for our tests (Section 3.2).

Chapter 4 will review the results of our research using the methodology previously explained. We will explore the number of refactoring messages for each repository that we study relative to the size of the lines of code. Additionally, we will review the average Maintainability Index (MI) for each project reviewed. From there, we will compare the scores among the group of projects.

We will then provide conclusions and recommendations in Chapter 5. In addition, the practices commonly considered the best to follow in computer programming are reinforced by much of our research, leading us to recom-

mend well-known methods in automated delivery pipelines.

Chapter 2

Background and Literature Review

In this chapter, we will gain a better understanding of the topic at hand, first by understanding long-term user engagement. Then, we will explore the impact of structural quality within a software system. Once we understand the problem thoroughly, we will determine a valuable dataset to that we can apply our theories, as well as explore related works.

2.1 Keeping Users Engaged Long Term

When developing a new system or a new software idea, getting the project off the ground and front of users is one thing. However, keeping that project alive with a thriving community of engaged users is another.

The systems we create could be customer-facing web applications, games, or internal applications used to carry out tasks. Regardless of the system, the product will no longer provide usefulness without evolving with the user's needs. As users become familiar with a system that is designed to depend, at least in part, on their attitudes and their practices, the users will modify their behavior to minimize their effort or maximize the effectiveness of the tool [2]. Even in a corporate setting with internal business systems, users will need to change; how a system can adapt to those needs requires a level of flexibility.

2.1.1 Why does software evolution matter?

“Software evolution is the continual development of software after its initial release to address changing stakeholder and/or market requirements.” [6]

When a system cannot evolve, the users primarily feel the impact. However, this impact will eventually get back to those who created and continue to support the system. With users that are either unsatisfied or unable to

use the system any longer, the engagement levels will drop. The decline in users will ultimately result in a loss of income, as the system can no longer deliver to the needs of its audience.

Because organizations invest large amounts of money in the software systems they create, they depend on their continued success. Software evolution will allow the system to adapt to new or changing business requirements, fix bugs and defects, and integrate with other systems that have changed and evolved that may share the same software environment.

While using a software system, inevitably, users will stumble into situations that even the best quality assurance testers will miss. Defects will require fixing once discovered. **THOUGHT: This is a maintenance thing; maybe inappropriate to include here**

To keep a system up-to-date, we must add new features. For example, there may be a need to significantly improve a system's performance or reliability if the user base expands. For a business or group to maintain user satisfaction, the software must continue to evolve; this will result in a system's increased size and complexity, as well as quality decline (unless very closely monitored over time) [7].

Security can also impact the need for a system to be maintained. For example, nefarious people can uncover new ways to infiltrate a system, so it is vital to stay on top of the newest versions of dependencies and technologies to avoid potential breaches of data and experience.

2.1.2 How do we ensure software evolution?

Because the maintainability of a system can ultimately influence the ability to generate revenue, we must find ways to ensure that a project will evolve. One of these ways could be to ensure that a project continues to be considered “maintainable” throughout its lifetime. This system characteristic will ensure that we can fix bugs quickly, but new features should be easy to add as the users’ needs evolve.

Not only is the user engagement necessary for the activation of the evolution process, but the engagement of the development community itself, especially when considering open source projects, is critical. As Adewumi et al. point out, a critical feature that distinguishes open source software is built and maintained by a community [8]. Additionally, the quality of the community will determine the quality of the open-source software [9].

2.2 The Impact of Structural Quality

2.2.1 Software Maintenance

The structural quality of a software system will impact the software evolution. Similarly, as a project evolves, it is likely to degrade, as developers often make changes quickly and in ways that the original design did not anticipate [10]. However, if the project has poor structural quality, its ability to evolve will be minimized, and the software system will eventually “die off.”

There is much planning involved in all software creation projects in what the product will be, will do, who it is for, and so on. One of the things that should also be on the planning list is long-term maintenance and growth. So how do we build a thing that will be easier to add features to down the road?

Let us define maintainability in the context of software. For example, it would be easy to maintain if a system is easy to debug and easy to add new features. These new features are generally considered minor features and may often be reported as bugs by users when, in reality, they are looking for functionality enhancements [11].

“Software maintenance in software engineering is the modification of a software product after delivery to correct faults, to improve performance or other attributes.” [11]

It may be easier to understand what characteristics define a system with poor maintainability. These types of systems will have poor code quality,

leading to defects. For example, there could be undetected vulnerabilities or vulnerabilities. However, it may be that the system is overly complex. In addition to the complexity, it could be hard to read due to the poor naming of dead (unused) code throughout the source code.

A project is known to have good maintainability when there is an enforced set of clean and consistent standards for the code. These standards often involve having human-readable names for functions, methods, and variables. Minimize any complex code and methods should be small and focus on a single thing. Parts of the system are decoupled and organized, making it easy to work on different parts with low impact on unrelated parts. For example, the code is DRY (there is limited redundancy in the code), unused code has been removed, and there is a level of documentation that supports an easy understanding of the system.

Why should we care about whether the code is maintainable? There is an assumption that the majority of cost over the lifetime of a project falls under maintainability. Fred Brooks, in his book “The Mythical Man-Month” (1975), even claimed that over 90% of the costs for a typical software system come up in the maintenance phase [12]. In 1977, Meir M. Lehman noted that maintenance accounts for 70% of a program budget, with 30% spent on development. In 1993, there was another observation that the development of a typical project uses only 20-40% of the resources (money, time, effort), with 60-70% used for maintenance activities [13].

Once the bulk of the system is off the ground and live worldwide, how

well the team can improve the system with new features and fix bugs, even working on different parts in parallel, can be impacted by its maintainability. Any successful piece of software will inevitably need to be maintained.

TODO: Anything about SOLID and when too many files are being touched?

[10]

2.2.2 Software Evolution

There is a distinction to be made between **software maintenance** and **software evolution**. We will refer to software maintenance as bug resolution and minor functional improvements. For example, we can consider this routine maintenance when we must fix a broken route in the application or provide a subtle enhancement to the user experience. However, when we look at upgrades, adaptations to the changing and growing needs of the user, or migrating the system to new technology, we can refer to this as the evolution of the software.

The evolution of software can result from new laws that have come into being. As technology changes, governing bodies must continually revisit data collection and information sharing policies. As a result, changes in technology and laws may lead to adaptations in the software systems.

It is also fair to say that systems will change because we can never fully determine a user's needs at the start of a project. So, for example, it would be safe to say that users' needs will change over time. However, the changing user needs can lead to a never-ending project that will always need some

form of enhancement.

Meir “Manny” Lehman and László “Les” Bélády contributed to a list of laws involving software evolution known as Lehman’s Laws that describe a balance between forces that drive new developments while also slowing progress. These laws apply to programs written to perform some real-world activity. The software environment and the software’s behavior are linked; additionally, this program category assumes that the program needs to adapt to varying requirements and circumstances in that environment. Eight laws were created and are listed below. [14]

TODO: Read and reference “An Empirical Study of Lehman’s Law on Software Quality Evolution” [7].

1. **Continuing Change** *(1974)*
2. **Increasing Complexity** *(1974)*
3. **Self Regulation** *(1974)*
4. **Conservation of Organisational Stability** *(1978)*
5. **Conservation of Familiarity** *(1978)*
6. **Continuing Growth** *(1991)*
7. **Declining Quality** *(1996)*
8. **Feedback System** *(1996)*

The first law, “Continuing Change,” tells us that it will become progressively less satisfactory if a system does not adapt. The second, “Increasing Complexity,” explains that as a system evolves, the complexity will increase unless there is work done to maintain or reduce complexity. Increased complexity can be due to the added code volume from new features or an increasing number of developers editing the code. Unless we actively address this phenomenon of increased complexity during changes, it can impact the maintainability (and the ability of a project to continue evolving) in the future.

Lehman’s fifth law, “Conservation of Familiarity,” explains how the average incremental growth does not change over time as a system evolves. The people interacting with the system, such as the developers, business persons, or users, must continue using and working within the system at the same “level of mastery.” If the system grows and changes excessively, the mastery will drop, slowing down the next set of changes. Slower deployment of additional changes could be because the source code or architecture has become more complex (impacting the developers’ ability to adapt and enhance the system) or because the user features have changed. For example, the system audience needs time to master the new interfaces or tools. However, the average incremental growth will remain steady because of this natural “slow-down” for excessive change. For example, we can see a simplified visual in “Fig. 2.1” showing that when the number of changes spikes (excessive growth in a system), it will be followed by an iteration of fewer changes. Over

time, excessive growth followed by fewer changes leads to a nearly consistent average of incremental growth (the thick, horizontal line).

Number of Changes vs. Time

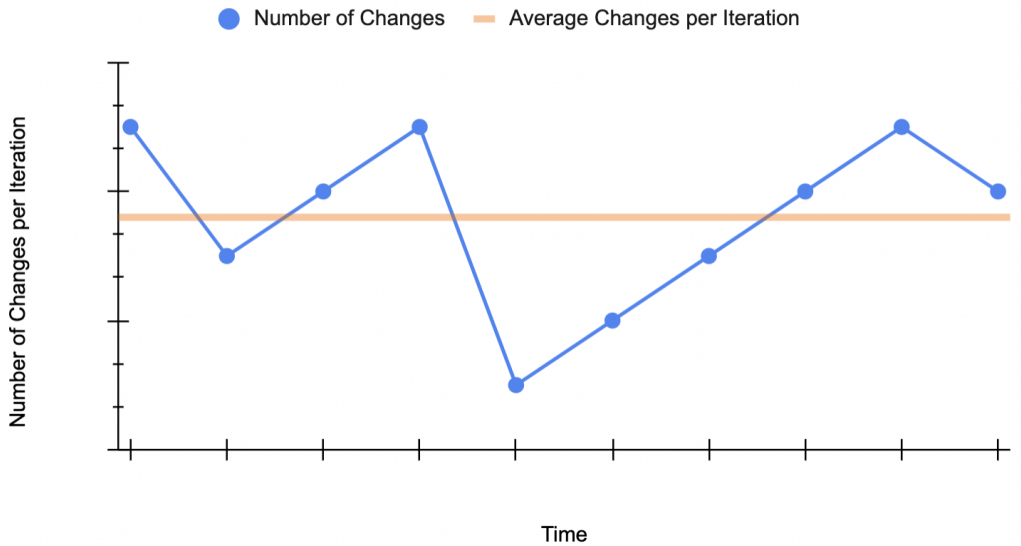


Figure 2.1: A simplified visual of Lehman’s fifth law, “Conservation of Familiarity.”

In Lehman’s sixth law, “Continuing Growth,” we see that the system user’s satisfaction will drop without continually increasing the functional content. Along with a similar idea, the law about “Declining Quality” states that if the operational environment for the system does not change, the system’s quality will appear to decline. Yu and Mishra performed focused research on supporting Lehman’s Laws, especially in the case of the seventh law about declining quality. They did this by defining a metric for software quality and reviewing the bug history, growth of the size, complexity, and quality of two large open-source systems [7]. Therefore, we must continue

adapting for even the appearance of the maintained quality of a system.

With many characteristics surrounding software evolution, we benefit from the Internet, which has positively improved the experience. Two common resources currently available to developers have impacted software evolution [6]:

1. The rapid growth of the World Wide Web and Internet Resources make it easier for users and engineers to find related information.
2. Open source development where anybody could download the source codes and modify it has enabled fast and parallel evolution (through forks).

These two suggestions are very evident in modern development. For example, a developer may regularly use resources like StackOverflow to find solutions to problems and use open-source tools that the developer and their team can contribute to or adjust to their specific needs.

2.2.3 Measuring Maintainability

Despite the nuanced differences between *maintainability* and *evolution*, the two characteristics run parallel to each other. It will also be easier to evolve if a system is easy to maintain. If we can measure our system's maintainability, we can also determine if our system is in a good position to continue evolving to meet our future needs.

“The unit cost of change must initially be made as low as possible and its growth, as the system ages, minimized. Programs must be made more alterable, and the alterability maintained throughout their lifetime. The change process itself must be planned and controlled.” [2]

Several tools attempt to provide some value around these ideas. In this paper, we will focus on the metrics that Pylint provides, specifically looking into their Refactor score of Pylint.

We will look at many open-source Python systems using Radon and Pylint, then attempt to correlate the data from the scores to the level of ease in adding new features to the system. This review will determine if a more maintainable system correlates with better Pylint scores.

2.2.4 Maintainability Scores

First, let us consider our original understanding of software maintainability. While this definition focuses primarily on bug fixes and minor enhancements, maintainable projects should also have ease in their ability to evolve. Therefore, we can study the impact maintainability (structural quality) has on software evolution by reviewing the scores provided by automated code review tools.

In this study, we will be using Pylint and will focus on the values of the Refactor score regarding a set of open-source Python systems. We must understand what Pylint itself is doing to understand the scores we will be working with,

Through the documentation of Pylint, we can understand how to use it and the scores it will provide [15]. The Pylint score itself is calculated by the following equation [16]:

$$10.0 - ((\text{float}(5 * e + w + r + c) / s) * 10)$$

Numbers closer to 10 reflect systems that have fewer errors, fewer warnings, and overall better structure and consistency. In the above equation, we are using the following values [17]:

- **statement** (**s**): the total number of statements analyzed
- **error** (**e**): the total number of errors, which are likely bugs in the code

- **warning** (**w**): the total number of warnings, which are python specific problems
- **refactor** (**r**): the total number of refactor warnings for bad code smells
- **convention** (**c**): the total number of convention warnings for programming standard violations

The Refactor score is of particular interest to us and considers many features meticulously outlined on the Pylint site [18]. These types of warnings include many checks, like prompting when to simplify a boolean condition, a useless **return**, and more. This score, in particular, will be part of our focus.

Pylint will check the code for code smells based on the definitions for documented checks to calculate the Refactor score. For every infraction, the score increases by one count.

“In computer programming, a **code smell** is any characteristic in the source code of a program that possibly indicates a deeper problem.” [19]

We can use these Refactor scores to help us spot architecture smells. After all, code smells can point the way to deeper problems in our system. However, there are fundamental established design principles that we should consider when creating software; code smells alert us to areas that have deviated from these principles. These smells are drivers for refactoring and when addressed,

can help us maintain the integrity of our architecture rather than creating a patchwork construction. Because of the relation of refactoring scores to the code structure itself, we will be spending much of our focus on this particular value.

Finally, concerning Python, it is also helpful to be familiar with PEP 8, as this is the default set of standards that Pylint uses to judge Python code [20]. This standard leads to making code more readable and more consistent, which may contribute to the code being more maintainable than without the standards. These standards cover things like indentation spacing, maximum line length, where to insert new lines, how to handle imports, and more. By defining a set of standards, teams can ensure they have a defined set of rules so that any contributors to the code understand the expectations (and so that automated systems like Pylint can enforce those standards to maintain readability and consistency).

2.2.5 Other Maintainability Characteristics

The authors of “Measurement and refactoring for package structure based on complex network” recently reviewed a similar idea focusing on cohesion and coupling over time for a project [21]. In a software system, we desire low coupling (allowing for changes to one area to remain independent of changes to another area) and high cohesion (indicating reduced complexity in modules, which improves maintainability). Through a few experiments on open-source software systems, the authors determined that their algorithm that

calculated metrics could improve package structures to have high cohesion and low coupling. Their study gives us confidence that metrics around the software’s structure can provide value in keeping systems in a maintainable state, which allows for software evolution.

Another variable that may impact the maintainability of code is readability. For example, in the article “How does code readability change during software evolution?” the authors have addressed this concern and found that most source codes were readable within the sample they reviewed. Additionally, a minority of commits changed the readability; if a file is less readable, it was likely that it remained that way and did not improve [22]. This variable (readability) in the maintainability of a software system can influence how easy or difficult it is to make a change. The authors also found that big commits, usually associated with adaptive changes (a form of software evolution), were the most prone to reduce code readability [22]. We assume that smaller commits are almost always better and can lead to more readable code.

Piantadosi et al. found that changes in readability, whether improvements or disintegrations, often occurred unintentionally [22]. By enforcing the PEP 8 standard, we know that Pylint is encouraging systems to remain readable. Therefore, projects that use some form of automated system in their pipeline benefit from keeping their project on track, limiting the effects of readability on a software’s potential for evolution.

The paper “Standardized code quality benchmarking for improving soft-

ware maintainability” provides insights into how the code’s maintainability is impacted by the technical quality of source code [23]. Within their paper, the authors seek to show four key points: (1) how easy it is to determine where and how the change is made, (2) how easy it is to implement the change, (3) how easy it is to avoid unexpected effects, and (4) how easy it is to validate the changes. Their approach has shown that some tools and methods can improve and maintain technical quality within their projects, allowing systems to continue to evolve at a reasonable pace.

2.2.6 Documentation and Maintainability

We assume that the Refactor score in projects should correlate to the system’s evolution. The first pass through the data is not conclusive in this detail, as the projects reviewed have many other factors contributing to the evolution of the project (number of contributors, size of the code system, and more). We assume that the correlation between software quality and software evolution would indicate that the better-scoring code systems are readable. In addition, it would be helpful to understand whether there are any similarities in the system’s documentation that could contribute to improved software evolution of a system.

In the textbook “Software Architecture in Practice,” chapter 18 provides some insight into documentation around architecture [24]:

“If you go to the trouble of creating a strong architecture, one that you expect to stand the test of time, then you *must* go to the trouble of describing it in enough detail, without ambiguity, and organizing it so that others can quickly find and update the needed information.”

The book describes how documentation holds the results of significant design decisions, providing valuable insights into decisions down the road. While not directly related to the Pylint Refactor score and not within the source code itself, it is still helpful to remind ourselves that documentation can also influence the ability of a software system to evolve.

“Our study has shown that the primary studies provide empirical evidence on the positive effect of documentation of designs pattern instances on programme comprehension, and therefore, maintainability.”

“...developers should pay more effort to add such documentation, even if in the form of simple comments in the source code.”

In research done by Wedyan and Abufakher (quoted above), documenting design patterns helped enhance code understanding [25]. In turn, comprehensibility impacts the maintainability of the code positively, reinforcing the impact that documentation can have and how it ties nicely into considerations for software structure.

Borrego et al. discussed that software development teams must have access to architecture knowledge, as that is the base on which they will build and understand the technical part of any software development cycle [26]. However, there is recognition that an absence of knowledge management is a factor in the failure of a development project [26]. Bjørnson and Dingsøy noted that software engineering knowledge is managed chiefly through repositories that will support knowledge flows [27].

2.3 Related Work

In our research, we are running with the assumption that Maintainability Index (MI) is our primary indicator. Therefore, we will look for the correlation between the MI and other Pylint scores. We hope to find which correlations align and which are the most important.

2.3.1 Considering Data Sets

When exploring the correlations between maintainability and refactoring, many sources are available for research. For example, some researchers have looked at proprietary systems as they evolve, while others have chosen open-source code available to the general public.

A study conducted by Baishakhi Ray, Daryl Posnett, Premkumar Devanbu, and Vladimir Filkov begins by programmatically collecting a sample set of projects on GitHub that vary in languages. Then the group of projects is appropriately thinned out, resulting in a final set used for the review. The results are then studied for the impact different programming languages may have on the code quality [28]. Finally, their research determined which languages were more prone to defects and which individual languages were more related to individual bugs rather than bugs overall.

The authors of “Predicting Maintainability with Object-Oriented Metrics - An Empirical Comparison” performed a similar study to what we are doing here. Their study focuses on object-oriented software (specifically C/C++

and Java) and correlation analysis between object-oriented metrics and software maintainability. Janke et al. looked for the best metrics to predict maintainability [29]. This study focuses on a few hand-picked software systems with an analysis of the changelogs. Our study, however, will be of a larger scale (about 50 software systems) and focused solely on Python-heavy projects.

2.3.2 Design Patterns and Software Quality

In the paper “Impact of design patterns on software quality: a systematic literature review” the authors compared the use of design patterns to software evolution and maintainability. They found that design patterns provided flexibility when reviewing changes that extended (evolved) software [25].

“Changes performed in a class can be corrective, adaptive, perfective, or preventive. These changes can occur due to new requirements, debugging, changes that propagate from changes in other classes and refactoring.”

Wedyan and Abufakher found that there were two reasons that a class had more frequent changes [25]:

1. The class was easy to extend.
2. The class correlated to other classes (raising alarms about class modularity).

With these findings in mind, we intentionally aim to focus our research on changes for system extensions and adaptations rather than bug fixes that appeared to be more considerable change due to high coupling. This paper focused on Refactor scores (code smells) rather than Error scores (bugs) within the system.

2.3.3 Software Architecture and Maintainability

In the research done in “Software Architecture Metrics: A Literature Review”, the authors discuss how early detection of issues within the software’s architecture is key to mitigating the risk of poor performance and can lower the cost of repairing faults [30]. While most developers have had access to these metrics for several decades, the industry and open-source community have not latched onto their use for keeping code in easy-to-work-with conditions.

The review done by Coulin et al. called out five essential qualities of software architecture [30]:

1. Maintainability
2. Extensibility
3. Simplicity, Understandability
4. Re-usability
5. Performance

Focusing on these qualities can narrow down the choice between different design options to an ideal solution. Keeping these five qualities in top-of-mind for new (and changed) code allows for easier future development and evolution of the software system.

ISO/IEC 25010:2011 is a detailed standard for software quality that contains eight product quality characteristics [31]. Each characteristic is further comprised of various sub-characteristics. “Fig. 2.2” lists these characteristics, with maintainability being one of the most significant characteristics in some studies [32], [1].

Product Quality Model			
FUNCTIONAL SUITABILITY	PERFORMANCE EFFICIENCY	PORTABILITY	SECURITY
Completeness	Time Behavior	Adaptability	Confidentiality
Correctness	Resource Utilization	Installability	Integrity
Appropriateness	Capacity	Replaceability	Non-Repudiation
USABILITY	RELIABILITY	MAINTAINABILITY	Accountability
Appropriateness	Maturity	Modularity	Authenticity
Recognizability	Availability	Reusability	COMPATIBILITY
Learnability	Fault Tolerance	Analysability	
Operability	Recoverability	Modifiability	
User Error Protection		Testability	
UI Aesthetics			Co-Existence
Accessibility			Interoperability

Figure 2.2: The eight product quality characteristics and sub-characteristics.

Chapter 3

Methodology

There are several factors to consider when reviewing data and considering which software systems to consider. A prominent starting place is with open source software (OSS), which is freely available to study.

To find projects of a caliber worth studying, we may also consider the maintenance capacity.

“Maintenance capacity refers to the number of contributors to an OSS project and the amount of time they are willing and able to contribute to the development effort as observed from versioning logs, mailing lists, discussion forums and bug report systems. Furthermore, sustainability refers to the ability of the community to grow in terms of new contributors and to regenerate by attracting and engaging new members to take the place of those leaving the community.” [1]

Open source projects, as stated above, have many factors and variables that contribute to the successful ability to maintain and evolve a project. The architecture must be built in a fashion that can be enhanced, but the community of developers must be adequately engaged and interested in continuous improvement.

3.1 Initial Repository Set

We have established that we have a problem with projects that fail to evolve, resulting in a loss of revenue. We also understand that evolving software is essential to keep users engaged; without it, there is an appearance in the decline of quality, the program becomes less satisfactory to the user, and the potential for competitors to outpace us with features available. We must now understand how we can ensure that our systems evolve. For this, we will look to understand how the system’s structural quality impacts software evolution.

The work done by Dr. Omari and Dr. Martinez involves collecting a sub-set of Python projects that we can use for further research. The bulk of their effort is to determine which classifiers to use to pare down the public set of Python systems into a good collection for further analysis [33]. In addition, we have used the work they have provided to select appropriate Python systems for review by collecting meta-data on these code systems.

Selection criteria employed by Omari and Martinez included “popular projects with long development history and multiple release cycles [33].” All projects are open source, enabling other researchers to access the defined set. Additionally, they captured meta-data used to define our subset of repositories for our particular focus.

From their subset of repositories, we reviewed current Pylint scores from each of the 132 systems. This set gives us a sampling of data that we can now

dig deeper into, comparing similar systems (similar size, a similar number of contributors, and more) and their evolution process by reviewing past commits rather than merely the system's current state here.

3.2 Filtered Respository Set

Once we had a narrowed set of projects from GitHub to ones primarily written in Python, we culled the set more using several criteria:

1. Projects that are at least 80% Python
2. Projects with a long history of commits
3. Projects with large development teams (community of contributors)
4. Projects with many releases
5. Projects of a substantial age

Armed with this list, we were able to use the metadata from GitHub for each of our repositories already collected and determine a cross-section of these criteria that would result in about 50 repositories for further study.

Beginning with the languages field from GitHub, we could quickly narrow down projects that had at least 80% of the code in Python. In our set, 103 repositories contained 80% or more Python code.

With this narrowed set, we then looked to see which percentile all the remaining criteria would yield the desired number of repositories. We determined that using the value at the 20th percentile in each of the above categories would yield the size set we would need.

Criteria	20th Percentile Value
Number of Commits	2,968
Number of Contributors	90
Number of Releases	44
Age (in months)	66.4

Table 3.1: Criteria used to filter down the initial set of repositories.

Table 3.1 shows the values found for each of our criteria. Using these values as our minimum requirements, we can narrow our repository set to 46 repositories (see Table 3.2).

Repository Names
ansible astropy autobahn-python aws-cli beets biopython boto buildbot celery cobbler conda cython django django-rest-framework electrum fail2ban gensim luigi matplotlib mongoengine mitmproxy mongo-python-driver mopidy networkx paramiko nova numba pandas peewee pelican pip pyramid ranger raven-python salt scikit-image scikit-learn scrapy sentry sqlalchemy swift sympy tornado web2py werkzeug youtube-dl

Table 3.2: List of the 46 repositories for research focus.

Chapter 4

Results

With our narrowed focus on repositories, we can now spend more time looking through the data collected. Specifically, we gathered data on the Refactor Warnings provided by Pylint. In Table 4.1 we have a shortened listing of the repositories from our set with their total count of refactoring warning messages, in addition to the most commonly appearing to refactor warning message. For the entire set, see Table 7.1 in the Appendix.

Regarding the number of refactoring warning messages, we can see that the “worst” project in our set was the repository Sympy [34], with 14,206 refactor messages. However, project size may also weigh how many warnings and errors are present, as projects with more lines of code are presumed to have more errors.

If we instead look at the ratio of refactoring warnings to the total lines of source code, we get a very different picture, provided in Table 4.2. We

Repository Name	Msg Count	Top Msg	Top Msg Count
sympy	14,206	too-many-arguments	6,601 (46%)
ansible	10,431	no-else-return	1,711 (16%)
salt	7,814	too-many-arguments	1,640 (21%)
ranger	109	no-else-return	29 (27%)
sentry	66	no-self-use	32 (48%)
raven-python	20	too-few-public-methods	7 (35%)

Table 4.1: Total refactor messages per repository for the 3 best and 3 worst offenders. Also provided with the refactor message that had the most warnings and its total count (and the percentage of that message from the total refactor warnings for that repository).

calculated these values using the following equation, r is the total refactor message count for the project, c is the total number of source lines of code in the project, and 100 is a multiplier to make the numbers easier to wrap our heads around. So the numbers closer to 0 are better ratios, as there are fewer refactor messages in the project relative to its size.

$$RefactorRatio = r/c * 100$$

In this case, the repository Raven-Python [35] (which was the “best” in regards to total refactor message count) is now the “worst” with the highest ratio of refactoring warnings compared to the lines of source code in the project. Raven-Python also has the least count of SLOC in the entire

repository set!

Repository	Total Refactor Msgs	Total Project SLOC	Ratio
raven-python	20	1,474	1.35
scrapy	381	58,768	0.64
numba	126	20,192	0.62
electrum	722	425,576	0.16
youtube-dl	1,003	667,075	0.15
cython	1,718	1,183,863	0.14

Table 4.2: Total refactor messages per repository, total project source line of code (SLOC), and the ratio of refactor messages to SLOC.

Given that we now have an idea of which projects have the most refactor messages per line of source code (“worst” offenders) and the projects with the least refactor messages (“best” projects for maintainability), we can look a little closer at the most frequent messages. For example, common to all six repositories at a high frequency is the message “no-self-use” and “no-else-return”. The message “no-self-use” means that ‘self’ is used as an argument but not in the method; this should be handled differently. For example, the message “no-else-return” highlights when an unnecessary code block follows an if-conditional.

Some of the other common messages call out the use of “too many” of different types of objects. For example, “too-many-branches”, “too-many-arguments”, “too-many-locals”, “too-many-statements”...

“Software maintainability these days has become one of the essential external attributes of software, which further forms a basis

of research for many researchers working in the fields related to software engineering. Software maintainability can be described as how a particular software system can be changed concerning the number of Lines of Code (LOC).” [32]

Now, it would be helpful to understand where we stand on the Maintainability Index (MI) for these projects, as this is our form of measurement that will help us understand how the code might be able to evolve. The MI is calculated on a per-module basis, but for the sake of conversation (and acknowledging that we will lose some nuances by reducing it this way), let us find the average MI at the project level and add this to our table (Table 4.3). Also included is the standard deviation for each average, with the lowest standard deviation being 13.75 and the highest at 31.14. Most averages have a standard deviation between 20 and 30 points.

Repo	Refactor Msgs	Project SLOC	Ratio	Avg Project MI
raven-python	20	1474	1.35	87.02 \pm 13.75
scrapy	381	58768	0.64	64.47 \pm 17.72
numba	126	20192	0.62	62.55 \pm 21.08
electrum	722	425576	0.16	39.41 \pm 28.06
youtube-dl	1003	667075	0.15	54.16 \pm 19.95
cython	1718	1183863	0.14	31.02 \pm 29.19

Table 4.3: Total refactor messages per repository, total project source line of code (SLOC), the ratio of refactor messages to SLOC, average Maintainability Index (MI).

With our more minor data set on hand, we ran Radon against each repository to gain a Maintainability Index. We averaged the Maintainability Index

of all files within each project for a very simplified view. This gave us a single number for general comparison with each repository.

The “Avg MI” column in the table contains the calculated average value for each file in the project. In Radon’s grading system, any value over 20 is considered a “very maintainable” code. We also used a count of the Pylint REFACTOR messages and calculated the ratio of messages to source-line-of-code for how many code smells each project contained relative to its size. This calculated ratio is shown in the “Ratio” column. Listed in the table are the “best 3” (cython, youtube-dl, and electrum) and the “worst 3” (raven-python, scrapy, and numba) regarding their refactor ratios.

When reviewing the Maintainability Inde for our repositories at either end of the spectrum and in context with the entire set (see Appendix 7.2), we will find that while Raven-Python has the highest number of refactoring warnings when compared to the SLOC count. Raven-Python also has the second-highest average MI across all of its project files.

Of the entire set, our highest average MI is 87.38 ± 19.47 (belonging to Sentry), and our lowest average MI is 28.85 ± 27.37 (belonging to Matplotlib). Regardless, these average values are above a score of 20, which is Radon’s lowest score provided for an “A” rank, that is, what would be considered a project “very high” maintainability.

With this information, we can get an idea of where these projects stand regarding their refactor message ratios compared to each other. For example, the histogram in Fig. 4.1 shows that the majority of our repository set has a

meager ratio of refactoring messages. However, the mass projects on the low-end could indicate that our projects are highly maintainable based on our assumptions that refactor messages are related to maintainability. Having few refactor messages (relative to the number of lines of code) means we have a smaller number of code smells. Therefore, infrequent refactor warnings are helpful for our architecture and quality!

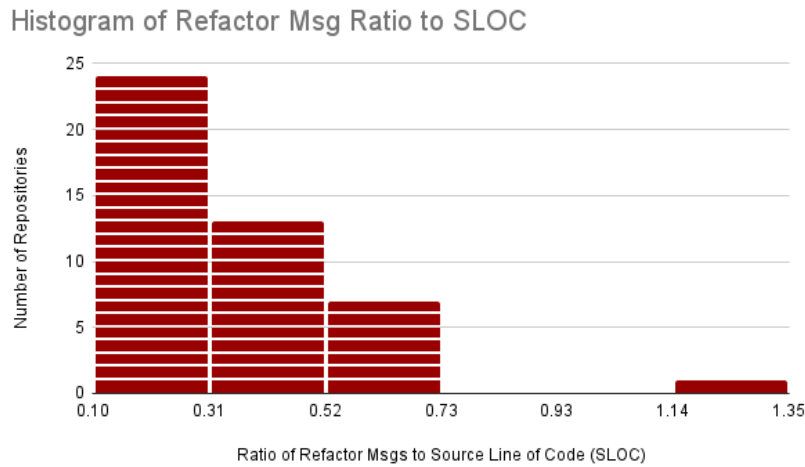


Figure 4.1: Diligent development communities can keep refactor warnings low, regardless of system size (lines of code).

The data set we chose may impact the ratios, lending to them all being on the low end. Perhaps open-source projects with highly engaged communities tend to keep their code in a maintainable state because this would be the only way for significant involvement (over 90 contributors per project). It has been interesting to see how low the ratio of refactoring messages has been in our data set.

To view the same repository set with their Maintainability Index averages,

we have another histogram in Fig. 4.2. This Radon score can range from 0 (awful) to 100 (perfect). Radon considers any value above 20 to be an “A” or “very maintainable”. All of our projects receive an “A” grade, with the majority in the middle range (around 40 to 50). All projects in our data set were all mid-to high-range scores.

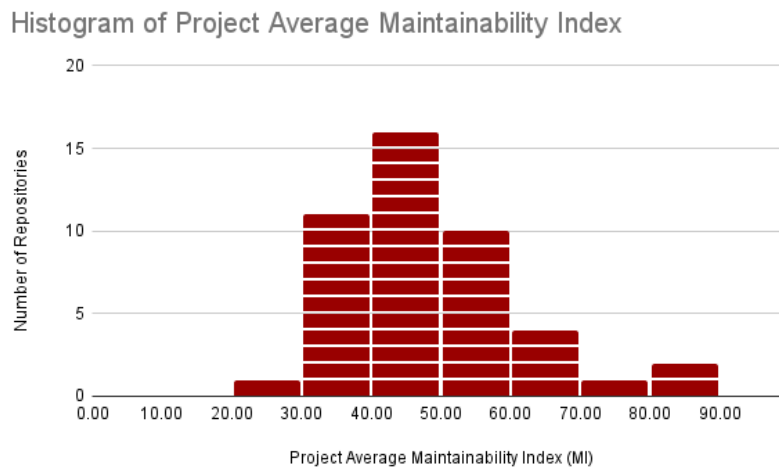


Figure 4.2: Many repositories average in the mid-score to high-score. Radon considers 20 points and up to be very maintainable.

Chapter 5

Conclusions and Recommendations

By collecting data and drawing our conclusions from it, with help from the insights from the studies done before ours, we may better understand metrics that can be useful regarding maintainability. Good projects will inevitably continue to grow and evolve. For example, understanding methods to keep code refactor on a certain level makes code easy to change. We may also find that projects with worsening scores slow down with updates and have reduced engagement.

When reviewing our data with current project Refactor message ratios, our three worst offenders were *Raven-Python*, *Scrapy*, and *Numba* [35], [36], [37]. *Raven-Python* has been deprecated in favor of a newer system, *Sentry-Python*. Despite their poor scores, the other two projects are still active in

development. The two active projects are included in thousands of projects each, which may be the reason for continued development despite the potential difficulty in maintenance.

Projects that may be open source or have many contributors are especially vulnerable to maintainability degrading over the evolution of a project. Having a reliable metric can be very useful in programmatically avoiding code smells and keeping code in a state that is easy to manage through simple metric checks in deployment pipelines.

While not one of our worst refactor-to-SLOC ratios, *SymPy* is a repository with a large number of refactoring warnings. We can see an example of code smell issues in reviewing some current symptoms that *SymPy* is experiencing, with only 72% code coverage and a failing build (see "Fig. 5.1"). Despite the engagement and continued development, we suspect that the actual adaptations and evolution of the software may be complex with this code.



Figure 5.1: A snapshot of the badges from *SymPy*'s repository.

We have further work to do in this study to gain a better understanding. With several "best" and "worst" Python software systems, we will look into the history of the projects' commits. It would be helpful to see how the Refactor scores have changed over time and if the rate at which changes were pushed correlated to the increase or decrease in that Refactor score.

Additional data can be gathered from this set that may provide more insights than this first brush of the data provides us. Understanding the impact of structural quality on the evolution of a project can provide compelling perspectives.

From our set, the very best and worst repositories are still actively improving and evolving, even today. An exception to this is our "worst" repository, raven-python, which was deprecated (this means the code is no longer supported), but in exchange for an improved system. The raven-python community did not die but instead recognized that a better architecture was needed for the system to continue to evolve. So they built a new code system (Sentry-Python) and shifted support there, where they have a better architecture and ability to evolve.

Open-source systems and any project with many contributors or have many changes are vulnerable to degrading quality. It is just the nature of change over time. However, when we review a set of popular repositories that already have a long history of commits (these repositories are still active and over five years old!), we find that they have good maintainability.

Where do we go from here? Regarding the data, we would like to do more studies to understand the correlation between the Pylint metrics and Maintainability Index to gain better insights into our assumptions.

So far, we have found a reinforcement of what many of us already know. Good architecture is vital for software to be able to evolve. Reliable quality metrics can help keep an architecture ready for enhancements. Regardless

of the codebase, teams should understand what standards will work best for their software solution. Then auto-enforce these standards in the development pipeline. Automated pipelines keep everyone honest and will improve the project's ability to evolve.

References

- [1] A. Adewumi, S. Misra, N. Omoregbe, B. Crawford, and R. Soto, “A systematic literature review of open source software quality assessment models,” *SpringerPlus*, vol. 5, no. 1936, 2016.
- [2] M. M. Lehman, “Programs, life cycles, and laws of software evolution,” *Proceedings of the IEEE*, vol. 68, no. 9, pp. 1060–1076, September 1980.
- [3] A. V. Deursen, “Think twice before using the “maintainability index”,” 2014, <https://avandeursen.com/2014/08/29/think-twice-before-using-the-maintainability-index/> [Online; accessed 23-Jan-2022]. [Online]. Available: <https://avandeursen.com/2014/08/29/think-twice-before-using-the-maintainability-index/>
- [4] H. Alsolai and M. Roper, “A systematic literature review of machine learning techniques for software maintainability prediction,” *Information and Software Technology*, 2019.
- [5] M. Lacchia, “Radon 4.1.0 documentation,” 2020, <https://radon.readthedocs.io/en/latest/commandline.html> [Online; accessed 19-April-2022]. [Online]. Available: <https://radon.readthedocs.io/en/latest/commandline.html>
- [6] W. contributors, “Software evolution — Wikipedia, the free encyclopedia,” 2021, https://en.wikipedia.org/wiki/Software_evolution [Online; accessed 12-December-2021]. [Online]. Available: https://en.wikipedia.org/wiki/Software_evolution
- [7] L. Yu and A. Mishra, “An empirical study of lehman’s law on software quality evolution,” *Int J Software Informatics*, vol. 7, no. 3, pp. 469–481, 2013.

- [8] K. Haaland and A.-K. Groven, “Free/libre open source quality models-a comparison between two approaches,” 2010.
- [9] I. Samoladas, G. Gousios, D. Spinellis, and I. Stamelos, “The sqo-oss quality model: Measurement based open source software evaluation,” *IFIP International Federation for Information Processing*, vol. 275, pp. 237–248, 2008.
- [10] R. C. Martin, “Design principles and design patterns,” *Object Mentor*, 2000.
- [11] W. contributors, “Software maintenance — Wikipedia, the free encyclopedia,” 2021, https://en.wikipedia.org/wiki/Software_maintenance [Online; accessed 17-December-2021]. [Online]. Available: https://en.wikipedia.org/wiki/Software_maintenance
- [12] F. Brooks, *The Mythical Man-Month*. Addison-Wesley, 1975.
- [13] “Ieee standard for software maintenance,” *IEEE Std 1219-1993*, pp. 1–45, 1993.
- [14] W. contributors, “Lehman’s laws of software evolution,” 2021, https://en.wikipedia.org/wiki/Lehman%27s_laws_of_software_evolution [Online; accessed 12-December-2021]. [Online]. Available: https://en.wikipedia.org/wiki/Lehman%27s_laws_of_software_evolution
- [15] Logilab and contributors, “Pylint,” Logilab, 2020, <https://pylint.org/> [Online; accessed 14-December-2021]. [Online]. Available: <https://pylint.org/>
- [16] P. Logilab and contributors, “Pylint features,” Logilab and PyCQA, 2021, https://pylint.pycqa.org/en/latest/technical_reference/features.html#reports-options [Online; accessed 14-December-2021]. [Online]. Available: https://pylint.pycqa.org/en/latest/technical_reference/features.html#reports-options
- [17] R. Kirkpatrick, “A beginner’s guide to code standards in python - pylint tutorial,” 2016, <https://docs.pylint.org/en/1.6.0/tutorial.html> [Online; accessed 18-December-2021]. [Online]. Available: <https://docs.pylint.org/en/1.6.0/tutorial.html>

- [18] P. Logilab and contributors, “Pylint features,” Logilab and PyCQA, 2021, https://pylint.pycqa.org/en/latest/technical_reference/features.html#refactoring-checker [Online; accessed 14-December-2021]. [Online]. Available: https://pylint.pycqa.org/en/latest/technical_reference/features.html#refactoring-checker
- [19] W. contributors, “Code smell — Wikipedia, the free encyclopedia,” 2021, https://en.wikipedia.org/wiki/Code_smell [Online; accessed 18-December-2021]. [Online]. Available: https://en.wikipedia.org/wiki/Code_smell
- [20] P. S. Foundation and contributors, “Pep 8 – style guide for python code,” Heroku Application, 2021, <https://www.python.org/dev/peps/pep-0008/> [Online; accessed 14-December-2021]. [Online]. Available: <https://www.python.org/dev/peps/pep-0008/>
- [21] Y. Zhou, Y. Mi, Y. Zhu, and L. Chen, “Measurement and refactoring for package structure based on complex network,” *Applied Network Science*, vol. 5, no. 50, 2020.
- [22] V. Piantadosi, F. Fierro, S. Scalabrino, A. Serebrenik, and R. Oliveto, “How does code readability change during software evolution?” *Software Qual J*, vol. 25, pp. 5374–5412, 2020.
- [23] R. Baggen, José, P. Correia, K. Schill, and J. Visser, “Standardized code quality benchmarking for improving software maintainability,” *Software Qual J*, vol. 20, pp. 287–307, 2012.
- [24] L. Bass, P. Clements, and R. Kazman, *Software Architecture in Practice: Third Edition*. Addison-Wesley Professional, 2012.
- [25] F. Wedyan and S. Abufakher, “Impact of design patterns on software quality: a systematic literature review,” *IET Software*, vol. 14, no. 1, 2020.
- [26] G. Borrego, A. L. Morán, R. R. P. Cinco, O. M. Rodríguez-Elias, and E. García-Canseco, “Review of approaches to manage architectural knowledge in agile global software development,” *IET Software*, vol. 11, no. 3, pp. 77–88, 2017.

- [27] F. O. Bjørnson and T. Dingsøy, “Knowledge management in software engineering: A systematic review of studied concepts, findings and research methods used,” *Information and Software Technology*, vol. 50, no. 11, pp. 1055–1068, 2008. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0950584908000487>
- [28] B. Ray, D. Posnett, P. Devanbu, and V. Filkov, “A large-scale study of programming languages and code quality in github,” *COMMUNICATIONS OF THE ACM*, vol. 60, no. 10, pp. 91–100, 2017.
- [29] M. Dagpinar and J. H. Janke, “Predicting maintainability with object-oriented metrics - an empirical comparison,” *Reverse Engineering - Working Conference Proceedings*, pp. 155–164, 12 2003.
- [30] T. Coulin, M. Detante, W. Mouchère, and F. Petrillo, “Software architecture metrics: A literature review,” 2019.
- [31] T. C. I. J. S. 7, Software, and systems engineering, “Iso/iec 25010:2011: Systems and software engineering — systems and software quality requirements and evaluation (square) — system and software quality models,” ISO, 2011, <https://www.iso.org/obp/ui/#iso:std:iso-iec:25010:ed-1:v1:en> [Online; accessed 14-December-2021]. [Online]. Available: <https://www.iso.org/obp/ui/#iso:std:iso-iec:25010:ed-1:v1:en>
- [32] S. Gupta and A. Chug, “An extensive analysis of machine learning based boosting algorithms for software maintainability prediction,” *International Journal of Interactive Multimedia and Artificial Intelligence*, vol. 7, no. 2, pp. 89–109, October 2021.
- [33] S. Omari and G. Martinez, “Enabling empirical research: A corpus of large-scale python systems,” 2018.
- [34] S. contributors, “SymPy,” 2021, <https://github.com/sympy/sympy> [Online; accessed 18-December-2021]. [Online]. Available: <https://github.com/sympy/sympy>
- [35] —, “Sentry,” 2021, <https://github.com/getsentry/raven-python> [Online; accessed 18-December-2021]. [Online]. Available: <https://github.com/getsentry/raven-python>

- [36] —, “Scrapy,” 2021, <https://github.com/scrapy/scrapy> [Online; accessed 23-April-2022]. [Online]. Available: <https://github.com/scrapy/scrapy>
- [37] N. contributors, “Numba,” 2022, <https://github.com/numba/numba> [Online; accessed 23-April-2022]. [Online]. Available: <https://github.com/numba/numba>
- [38] C. contributors, “Cython,” 2022, <https://github.com/cython/cython> [Online; accessed 23-April-2022]. [Online]. Available: <https://github.com/cython/cython>
- [39] Y.-D. contributors, “Youtube-dl,” 2022, <https://github.com/rg3/youtube-dl> [Online; accessed 23-April-2022]. [Online]. Available: <https://github.com/rg3/youtube-dl>
- [40] E. contributors, “Electrum - lightweight bitcoin client,” 2022, <https://github.com/spesmilo/electrum> [Online; accessed 23-April-2022]. [Online]. Available: <https://github.com/spesmilo/electrum>
- [41] A. contributors, “Ansible,” 2021, <https://github.com/ansible/ansible> [Online; accessed 17-December-2021]. [Online]. Available: <https://github.com/ansible/ansible>
- [42] S. contributors, “Salt,” 2021, <https://github.com/saltstack/salt> [Online; accessed 18-December-2021]. [Online]. Available: <https://github.com/saltstack/salt>

Appendix

7.1 Pylint: Refactor Scores

The score in Pylint is a value out of 10, with 10 being the best.

Pylint has a number of types of messages:

1. Convention
2. Error
3. Fatal
4. Information
5. Refactor
6. Warning

There is a very large list of all messages at [Overview of All Pylint Messages](#).

We spend most of our time reviewing the Refactor messages, found at [Pylint Refactor Messages](#).

Additionally, the Convention messages are also helpful for review, found at Pylint Convention Messages.

7.2 Radon: Maintainability Index

In this paper, we use Radon as one method for finding the Maintainability Index (MI) for projects, at the module (file) level.

MI Score	Rank	Maintainability
100 - 20	A	Very high
19 - 10	B	Medium
9 - 0	C	Extremely Low

7.3 Data: Repository Refactor Messages

TODO: fill this table with ALL the repositories

```
-- Find total number of refactor messages
SELECT
    t1.name AS "Repository Name",
    SUM(t1.n::numeric) AS "Total Refactor"
FROM new_pylint_metrics_by_msg AS t1
    WHERE t1.code = 'R'
GROUP BY t1.name
ORDER BY SUM(t1.n::numeric) DESC;

-- Find "worst offender" messages for each repository
SELECT
    t2.name AS "Repository Name",
    t2.msg AS "Most Common Refactor Message",
    SUM(t2.n::numeric) AS "Most Common Refactor Message Count"
FROM new_pylint_metrics_by_msg AS t2
    WHERE name = 'sympy' -- set repository name here
    AND code = 'R'
GROUP BY t2.name, t2.msg
ORDER BY SUM(t2.n::numeric) DESC
LIMIT 1;
```


Repository Name	Msg Count	Top Msg	Top Msg Count
sympy	14,206	too-many-arguments	6,601 (46%)
ansible	10,431	no-else-return	1,711 (16%)
salt	7,814	too-many-arguments	1,640 (21%)
ranger	109	no-else-return	29 (27%)
sentry	66	no-self-use	32 (48%)
raven-python	20	too-few-public-methods	7 (35%)

Table 7.1: Total refactor messages per repository. Also provided with the refactor message that had the most warnings and its total count (and the percentage of that message from the total refactor warnings for that repository).

7.4 Data: Repository Refactor-to-SLOC Ratio

It is good to note that the value for the Source Line of Code (SLOC), when added to multi-line comments, single-line comments, and blank lines, should add up to the total lines of code in the project [5]. For this table, we are looking only at the SLOC, which are lines that are only actual code (no comments or blank lines).

```
-- https://radon.readthedocs.io/en/latest/commandline.html#the-raw-command
-- The equation sloc+multi+singlecomments+blank=loc should always hold.
```

```
SELECT
  t1_radon_raw.name AS "Repository",
  SUM(
    CASE
      WHEN t3_python_metrics.code = 'R' THEN n::numeric
      ELSE 0
    END
  ) AS "Total Refactor Msgs",
  SUM(t1_radon_raw.sloc::numeric) AS "Total Project SLOC",
  CONCAT(TRUNC(
    SUM(
      CASE
        WHEN t3_python_metrics.code = 'R' THEN n::numeric
        ELSE 0
      END
    ) / SUM(t1_radon_raw.sloc::numeric) * 100,
    2 -- Number of decimal places
  ), ' ')
  ) AS "Ratio",
  CONCAT(
```

```

        TRUNC(AVG(t2_radon_mi.mi::numeric), 2),
        ' ±',
        TRUNC(STDDEV_POP(t2_radon_mi.mi::numeric), 2)
    ) AS "Avg Project MI",
    CASE
        WHEN AVG(t2_radon_mi.mi::numeric) >= 20 THEN 'A'
        WHEN AVG(t2_radon_mi.mi::numeric) >= 10 THEN 'B'
        ELSE 'C'
    END AS "MI Rank",
    CASE
        WHEN AVG(t2_radon_mi.mi::numeric) >= 20 THEN 'Very high'
        WHEN AVG(t2_radon_mi.mi::numeric) >= 10 THEN 'Medium'
        ELSE 'Extremely low'
    END AS "Maintainability",
    CONCAT(TRUNC(
        SUM(t1_radon_raw.comments::numeric)
        / SUM(t1_radon_raw.sloc::numeric)
        * 100,
        2 -- Number of decimal places
    ), '%') AS "Total Project Comment-to-SLOC Ratio"
FROM new_radon_raw_metrics AS t1_radon_raw
INNER JOIN new_radon_mi_metric AS t2_radon_mi
    ON t1_radon_raw.module_name = t2_radon_mi.module_name
    INNER JOIN new_pylint_metrics_by_msg AS t3_python_metrics
        ON t1_radon_raw.module_name = t3_python_metrics.module_name
        INNER JOIN new_repo_details AS t4_repo_details
            ON t1_radon_raw.name = t4_repo_details.repo_name_2
GROUP BY t1_radon_raw.name
ORDER BY (
    SUM(
        CASE
            WHEN t3_python_metrics.code = 'R' THEN n::numeric
            ELSE 0
        END
    ) / SUM(t1_radon_raw.sloc::numeric)
) DESC;

```

Repository	Refactor Msgs	Project SLOC	Ratio	Avg Project MI	MI Rank	Maintainability	Comment-to-SLOC Ratio
raven-python	20	1474	1.35	87.02 \pm 13.75	A	Very high	8.41%
scrapy	381	58768	0.64	64.47 \pm 17.72	A	Very high	5.59%
numba	126	20192	0.62	62.55 \pm 21.08	A	Very high	13.11%
sentry	66	10770	0.61	87.38 \pm 19.47	A	Very high	8.89%
boto	1398	237761	0.58	58.00 \pm 22.88	A	Very high	14.91%
celery	1626	302080	0.53	41.73 \pm 24.26	A	Very high	6.00%
pyramid	452	85302	0.52	70.25 \pm 24.18	A	Very high	10.29%
mopidy	211	40381	0.52	59.15 \pm 20.46	A	Very high	7.19%
sympy	14206	2873930	0.49	31.16 \pm 26.81	A	Very high	8.08%
aws-cli	585	126303	0.46	61.73 \pm 20.44	A	Very high	15.89%
luigi	299	68973	0.43	53.56 \pm 22.51	A	Very high	14.81%
networkx	568	147166	0.38	44.66 \pm 20.48	A	Very high	24.00%
scikit-learn	1485	396109	0.37	41.74 \pm 19.67	A	Very high	13.88%
peewee	163	44595	0.36	31.26 \pm 20.32	A	Very high	4.76%
mongo-python-driver	418	117448	0.35	43.62 \pm 23.74	A	Very high	12.44%
buildbot	1129	318794	0.35	56.60 \pm 23.29	A	Very high	17.88%
django	1720	485681	0.35	57.28 \pm 28.51	A	Very high	12.97%
pandas	1197	357537	0.33	42.07 \pm 25.86	A	Very high	9.24%
mitmproxy	488	148319	0.32	56.21 \pm 21.53	A	Very high	5.60%
biopython	1930	615445	0.31	43.54 \pm 23.97	A	Very high	17.19%
django-rest-framework	241	77488	0.31	50.68 \pm 28.56	A	Very high	9.35%
werkzeug	348	116937	0.29	34.06 \pm 20.89	A	Very high	6.20%
autobahn-python	541	182842	0.29	46.86 \pm 27.75	A	Very high	19.61%
sqlalchemy	2040	703340	0.29	32.98 \pm 29.94	A	Very high	7.17%
cobbler	221	77348	0.28	58.76 \pm 22.87	A	Very high	12.01%
scikit-image	499	187819	0.26	53.87 \pm 23.64	A	Very high	9.64%
nova	2593	981563	0.26	64.72 \pm 31.14	A	Very high	16.51%
paramiko	257	97387	0.26	43.79 \pm 22.00	A	Very high	13.78%
tornado	326	123535	0.26	37.14 \pm 22.48	A	Very high	17.28%
conda	660	260363	0.25	44.82 \pm 27.94	A	Very high	16.42%
beets	471	186165	0.25	43.66 \pm 26.32	A	Very high	15.64%
astropy	1965	780617	0.25	40.20 \pm 28.61	A	Very high	19.38%
salt	7814	3131332	0.24	45.69 \pm 24.78	A	Very high	8.43%
ranger	109	44487	0.24	45.74 \pm 25.48	A	Very high	9.49%
pelican	117	48096	0.24	35.95 \pm 19.42	A	Very high	7.82%
swift	1159	523880	0.22	37.13 \pm 24.09	A	Very high	14.30%
pip	1215	590140	0.20	47.52 \pm 28.78	A	Very high	12.61%
matplotlib	1926	943406	0.20	28.85 \pm 27.37	A	Very high	11.10%
fail2ban	157	84340	0.18	48.41 \pm 20.03	A	Very high	30.77%
mongoengine	145	80473	0.18	30.75 \pm 29.57	A	Very high	9.24%
ansible	10429	5818380	0.17	46.82 \pm 22.87	A	Very high	5.98%
web2py	1555	903496	0.17	37.78 \pm 29.88	A	Very high	10.16%
electrum	722	425576	0.16	39.41 \pm 28.06	A	Very high	9.14%
youtube-dl	1003	667075	0.15	54.16 \pm 19.95	A	Very high	5.04%
cython	1718	1183863	0.14	31.02 \pm 29.19	A	Very high	12.74%

Table 7.2: Total refactor messages per repository, total project source line of code (SLOC), the ratio of refactor messages to SLOC, average Maintainability Index (MI) with standard deviation, and code comment-to-SLOC ratio.

7.5 Data: Repository Refactor By Message

This section contains information from the repositories that were on the edges of refactor message counts.

The sections for the projects with the highest ratio of refactor messages compared to the source lines of code are contained in Subsections 7.5.1 (Raven-Python), 7.5.2 (Scrapy), and 7.5.3 (Numba).

The sections for the projects with the lowest ratio of refactor messages compared to the source lines of code are contained in Subsections 7.5.4 (Cython), 7.5.5 (YouTube-DL), and 7.5.6 (Electrum).

The sections for the projects with the most total refactor messages are contained in Subsections 7.5.7 (SymPy), 7.5.8 (Ansible), and 7.5.9 (Salt).

```
SELECT
    t2.name AS "Repository Name",
    t2.msg AS "Most Common Refactor Message",
    SUM(t2.n::numeric) AS "Most Common Refactor Message Count"
FROM new_pylint_metrics_by_msg AS t2
WHERE name = 'sympy' -- Change repository name for each one
AND code = 'R'
GROUP BY t2.name, t2.msg
ORDER BY SUM(t2.n::numeric) DESC
LIMIT 10;
```

7.5.1 Raven-Python: Most Common Refactor Messages

This repository, Raven-Python [35], had the highest ratio of refactor messages when compared to the number of source lines of code.

Only 6 were returned because there is such a low number of refactor warning messages for this project.

Repository Name	Most Common Refactor Message	Most Common Refactor Message Count
raven-python	too-few-public-methods	7
raven-python	useless-object-inheritance	5
raven-python	no-self-use	4
raven-python	no-else-return	2
raven-python	inconsistent-return-statements	1
raven-python	too-many-branches	1

Table 7.3: The 6 most common refactor messages for Raven-Python.

7.5.2 Scrapy: Most Common Refactor Messages

This repository, Scrapy [36], had the second highest ratio of refactor messages when compared to the number of source lines of code.

Repository Name	Most Common Refactor Message	Most Common Refactor Message Count
scrapy	useless-object-inheritance	102
scrapy	no-self-use	79
scrapy	inconsistent-return-statements	49
scrapy	no-else-return	44
scrapy	too-few-public-methods	42
scrapy	too-many-arguments	29
scrapy	too-many-instance-attributes	17
scrapy	too-many-locals	5
scrapy	no-else-raise	4
scrapy	too-many-return-statements	4

Table 7.4: The 10 most common refactor messages for Scrapy.

7.5.3 Numba: Most Common Refactor Messages

This repository, Numba [37], had the third highest ratio of refactor messages when compared to the number of source lines of code.

Repository Name	Most Common Refactor Message	Most Common Refactor Message Count
numba	too-few-public-methods	54
numba	no-else-return	25
numba	useless-object-inheritance	13
numba	no-self-use	8
numba	too-many-branches	5
numba	too-many-public-methods	4
numba	too-many-arguments	4
numba	too-many-locals	3
numba	too-many-return-statements	2
numba	too-many-instance-attributes	2

Table 7.5: The 10 most common refactor messages for Numba.

7.5.4 Cython: Most Common Refactor Messages

This repository, Cython [38], had the lowest ratio of refactor messages when compared to the number of source lines of code.

Repository Name	Most Common Refactor Message	Most Common Refactor Message Count
cython	no-else-return	417
cython	no-self-use	293
cython	too-many-branches	182
cython	too-many-arguments	162
cython	useless-object-inheritance	132
cython	too-many-locals	107
cython	too-few-public-methods	79
cython	too-many-statements	79
cython	inconsistent-return-statements	61
cython	too-many-instance-attributes	53

Table 7.6: The 10 most common refactor messages for Cython.

7.5.5 YouTube-DL: Most Common Refactor Messages

This repository, YouTube-DL [39], had the second lowest ratio of refactor messages when compared to the number of source lines of code.

Repository Name	Most Common Refactor Message	Most Common Refactor Message Count
youtube-dl	too-many-locals	413
youtube-dl	too-many-branches	116
youtube-dl	inconsistent-return-statements	85
youtube-dl	too-many-statements	84
youtube-dl	no-else-return	83
youtube-dl	too-many-arguments	49
youtube-dl	no-self-use	49
youtube-dl	no-else-raise	31
youtube-dl	useless-object-inheritance	26
youtube-dl	too-many-nested-blocks	19

Table 7.7: The 10 most common refactor messages for YouTube-DL.

7.5.6 Electrum: Most Common Refactor Messages

This repository, Electrum [40], had the third lowest ratio of refactor messages when compared to the numbe

Repository Name	Most Common Refactor Message	Most Common Refactor Message Count
electrum	no-self-use	163
electrum	too-many-locals	87
electrum	inconsistent-return-statements	77
electrum	no-else-return	75
electrum	too-few-public-methods	69
electrum	too-many-arguments	51
electrum	too-many-instance-attributes	38
electrum	too-many-statements	38
electrum	too-many-branches	32
electrum	too-many-public-methods	26

Table 7.8: The 10 most common refactor messages for Electrum.

7.5.7 Sympy: Most Common Refactor Messages

This repository is SymPy [34], one of the repositories with the highest count of refactor warning messages.

Repository Name	Most Common Refactor Message	Most Common Refactor Message Count
sympy	too-many-arguments	6,601
sympy	no-else-return	2,813
sympy	no-self-use	831
sympy	inconsistent-return-statements	809
sympy	too-many-locals	741
sympy	too-many-branches	608
sympy	too-many-return-statements	376
sympy	too-many-statements	294
sympy	too-many-nested-blocks	159
sympy	too-many-ancestors	148

Table 7.9: The 10 most common refactor messages for SymPy.

7.5.8 Ansible: Most Common Refactor Messages

This repository is Ansible [41], one of the repositories with the highest count of refactor warning messages.

Repository Name	Most Common Refactor Message	Most Common Refactor Message Count
ansible	no-else-return	1,711
ansible	too-many-branches	1,265
ansible	inconsistent-return-statements	1,138
ansible	useless-object-inheritance	1,060
ansible	no-self-use	869
ansible	too-many-locals	816
ansible	too-many-statements	633
ansible	too-many-arguments	582
ansible	too-few-public-methods	512
ansible	too-many-nested-blocks	446

Table 7.10: The 10 most common refactor messages for Ansible.

7.5.9 Salt: Most Common Refactor Messages

This repository is Salt [42], one of the repositories with the highest count of refactor warning messages.

Repository Name	Most Common Refactor Message	Most Common Refactor Message Count
salt	too-many-arguments	1,640
salt	no-else-return	1,565
salt	too-many-branches	1,024
salt	too-many-locals	891
salt	too-many-statements	495
salt	too-many-nested-blocks	340
salt	inconsistent-return-statements	278
salt	too-many-return-statements	242
salt	too-few-public-methods	206
salt	useless-object-inheritance	206

Table 7.11: The 10 most common refactor messages for Salt.

7.6 Data: Pylint Scores for Best & Worst

This data was generated by cloning the existing repository and then deciding which folder contained most python files. Then check out a specific commit hash (beginning at the stage where our data was generated and moving towards current) and run pylint. The Pylint rating is a score out of 10.

```
##### Cython #####
# Clone the repository
git clone https://github.com/cython/cython.git

# Navigate into the repository folder
cd cython

##### Raven-Python #####
# Clone the repository
git clone https://github.com/getsentry/raven-python

# Navigate into the repository folder
cd raven-python

##### Generic Steps #####
# Checkout the commit hash
git checkout <hash>

# Find the date of the last commit at this stage
git log

# Run pylint
pylint ./*
```

7.6.1 Cython Pylint Scores

Release	Commit Hash	Commit Date	Pylint
(start)	433e6992ca89e0c7059b87cbae0e9536f11aa58f	14-Dec-2018	8.11
0.29.25	488e21a34259258210f0be92c58618e1ea8a928f	6-Dec-2021	8.14
0.29.26	3028e8c7ac296bc848d996e397c3354b3dbbd431	16-Dec-2021	8.13
3.0.0a10	ddaaa7b8bfe9885b7bed432cd0a5ab8191d112cd	6-Jan-2022	7.95
0.29.27	229a4531780863c8a5c311d6b3c70a545988f85f	28-Jan-2022	8.13
0.29.28	27b6709241461f620fb25756ef9f1192cc4f589a	17-Feb-2022	8.13
(latest)	d48d0a038e2838d3bd2981e2687557a86936076b	21-Apr-2022	7.91

Table 7.12: Pylint score from various releases of Cython, the “best” ranked repository for maintainability.

7.6.2 Raven-Python Pylint Scores

Release	Commit Hash	Commit Date	Pylint
(start)	2d55add71f3a4d4cbb86ba56770b0e5038cb57cf	24-Oct-2011	2.68
5.3.0	0c9b082a5f6dc70e7526e9107688a19d75a9da45	30-Apr-2015	4.93
5.4.0	3b0f6a0374b122d0b26cbdd641a70c4e01316598	6-Jul-2015	5.01
6.2.1	12b388a76dba4d1de82314f0184efcfcd1e2070a	22-Sep-2017	5.92
6.3.0	1a0d697dd11459bfa8ce933b3b53264098dad60d	29-Oct-2017	5.96
6.4.0	15ccf495f070a9dd8b3a0501cd1e7225429ca29c	11-Dec-2017	5.93
6.5.0	a2923ea42aeb0bfc8180dbec64dd037dfe381fc1	17-Jan-2018	5.95
6.6.0	f579e6809b01d27da5fe515d8572b497c98b4b43	14-Feb-2018	5.95
6.7.0	2c4ed64beecf9af4b45d4028eae7d540fa8df767	18-Apr-2018	5.85
6.8.0	595d69558d8a093b72423b07404c25863b1d0f31	12-May-2018	5.86
6.10.0	d7d14f61b7fb425bcb15512f659626648c494f98	19-Dec-2018	5.89
(latest)	5b3f48c66269993a0202cfc988750e5fe66e0c00	18-Jan-2021	5.89

Table 7.13: Pylint score from various releases of Raven-Python, the “worst” ranked repository for maintainability.

Release	Commit Hash	Commit Date	Pylint
(latest)	4cce4b5d9f5b34379879a332b320e870ce0ce1ad	20-Apr-2022	7.03

Table 7.14: The latest Pylint score from Sentry-Python, which replaced Raven-Python.