

Анализ качества ведения репозиториях GitHub для языка C++

Томилин Илья Сергеевич
Высшая школа программной инженерии
Санкт-Петербургский политехнический
университет Петра Великого
Санкт-Петербург, Россия
tomilin.is@edu.spbstu.ru

Фёдоров Иван Алексеевич
Высшая школа программной инженерии
Санкт-Петербургский политехнический
университет Петра Великого
Санкт-Петербург, Россия
fedorov.ia@edu.spbstu.ru

Никифоров Игорь Валерьевич
Высшая школа программной инженерии
Санкт-Петербургский политехнический
университет Петра Великого
Санкт-Петербург, Россия
nikiforov_iv@spbstu.ru

Аннотация—В данной статье представлен один из возможных подходов, позволяющий создавать программное решение для сбора набора данных с github, а также производить анализ на основе полученного набора данных и визуализировать результаты. В качестве объекта исследования был выбран анализ репозиториях C++ на качество их ведения. Данная статья может помочь упростить принятие решения в выборе средств и подхода для достижения цели анализа репозиториях.

Ключевые слова—*data, analysis, github, bigdata*

I. ВВЕДЕНИЕ

GitHub — крупнейший веб-сервис для хостинга IT-проектов и их совместной разработки. На 2021 год GitHub используют больше 73 миллионов разработчиков, больше 4 миллионов организаций и содержится более 200 миллионов различных репозиториях. Одни репозитории ведутся с хорошим стилем (наличие README, описания, манифеста системы сборки и.т.п.), а другие наоборот не следуют этому правилу. Особенно актуально соблюдение этих правил для репозиториях языка C++, поскольку они требуют повышенного качества ведения. Анализ качества репозитория необходим, чтобы выявить насколько качественно разработчики подходят к ведению своих репозиториях, стоит ли использовать структуру случайно выбранного репозитория или более тщательно искать пример качественного репозитория с точки зрения его ведения при создании своего проекта. В статье решаются задачи связанные с анализом:

- лицензий репозиториях, наличием README файла;
- наличием Dockerfile, для определения того имеет ли проект CI;
- используемого стиля отступов в файлах с исходным кодом;
- наиболее популярных языков, которые используют в одном репозитории совместно с C++;
- систем сборки;
- названий веток, наличием gitflow.

II. ОБЗОР АНАЛОГИЧНЫХ РЕШЕНИЙ

Схожие задачи анализа репозиториях на различные метрики рассматривались в статьях, краткое содержание которых будет изложено ниже.

A. Анализ качества свободно разрабатываемого программного обеспечения

В статье [2] рассматривались какие-либо существенные взаимосвязи между качеством проекта и характеристиками членов команды, которая работает над проектом. За основу показателя качества проекта были выбраны два таких показателя, первый из них отображает популярность проекта, второй качество поддержки, которое предоставляют члены команды. Первый показатель определялся на основе количества “звезд” репозитория, второй на основе анализа “выживаемости” пользователей проекта.

B. Анализ настроений комментариев к репозиториям

В статье [3] рассматривались эмоции комментариев, которые по мнению авторов могут оказывать большое влияние на производительность, качество выполнения задач. Авторы использовали лексический анализ настроений для изучения эмоций выраженных в комментариях к репозиторию, анализировали взаимосвязь с различными факторами такими как: язык программирования, время, дата.

C. Анализ важности и влияния репозиториях

В следующей статье [4] исследовалась важность репозиториях, поскольку тысячи разработчиков, которые размещают свои проекты на github превозносят новые нерешенные проблемы тем самым мотивируя других людей на участие в решение этих проблем. В статье описываются достаточно сложные алгоритмы исследования этой проблемы на основе обращений пользователей к репозиториям.

D. Анализ наличия системы непрерывной интеграции

В статье [5] авторы затронули такую тему как наличие систем непрерывной интеграции в репозиториях. Производилось исследование наличие системы Travis CI в репозиториях с различными языками программирования.

Е. Анализ репозиторий связанных с биоинформатикой

В статье [6] проводилось исследование репозиторий связанных с биоинформатикой. Анализировался исходный код этих репозиторий для изучения взаимосвязей между свойствами кода, деятельностью разработчиков, программным обеспечением. Было проанализировано 1720 репозиторий и получено большое количество необходимой информации. Данное исследование стало первым в данном направлении.

Ф. Анализ качества ведения Dockerfile в репозиториях

В статье [7] проводилось исследование содержимого Dockerfile с целью определения его качества. Для авторов качество Dockerfile определялось в написании конкретной версии образов, наличие строки с описанием разработчика и соблюдение хорошего стиля написания Dockerfile согласно рекомендациям описанным на сайте с документацией.

III. ОПИСАНИЕ МОДЕЛИ ДАННЫХ И ХАРАКТЕРИСТИКИ НАБОРА ДАННЫХ

А. Описание модели

Данные для анализа формируются в json документы со структурой показанной на рис. 1. В структуре находятся поля, которые заполняются на основе API github и разбора директории загруженного репозитория.

```
1  {
2    "all_lines": 0,
3    "another_files": 0,
4    "another_languages": [
5      ],
6    "branches": [
7      ],
8    "build_system": [
9      ],
10   "contributors_female": 0,
11   "contributors_male": 0,
12   "contributors_total_count": 0,
13   "contributors_unrecognized": 0,
14   "count_spaces": 0,
15   "count_tabs": 0,
16   "cpp_src_files": 0,
17   "cpp_src_lines": 0,
18   "created_at": "",
19   "description": "",
20   "dockerfile": false,
21   "forks": 0,
22   "full_name": "",
23   "indentation_method": "",
24   "indentation_method_dominated": "",
25   "issue_titles": [],
26   "issues_close": 0,
27   "issues_open": 0,
28   "license": "Other",
29   "readme": false,
30   "stars": 0,
31   "total_issues": 0,
32   "updated_at": ""
33 }
```

Рисунок 1 – Структура json документа

В. Характеристики набора данных

За все время сбора данных, были проанализированы репозитории за следующие года: 2010, 2011, 2012, 2013, 2014 и частично 2015. Общее число проанализированных github репозиторий составило – 178469 штук, fork не рассматривались. Также следует отметить, что это приблизительно 6% от общего числа всех созданных репозиторий для языка C++ в период с 2010-01-01 до

2021-12-01 годы. Общий объем загруженных репозиторий составил 4.4 ТБ. Количество часов, которое потребовалось для загрузки репозиторий составило 387 часов. Общий объем json документов с данными составил 758 МБ. Для загрузки документов с данными в индекс elasticsearch, потребовалось 40 минут.

IV. ПРЕДЛАГАЕМОЕ ПРОГРАММНОЕ РЕШЕНИЕ И ОПИСАНИЕ ЕГО ОСОБЕННОСТЕЙ

Для упрощения взаимодействия с программным решением и автоматизацией работы для конечного пользователя был создан программный дистрибутив содержащий 9 директорий и 30 файлов в которые входят: автоматизационные сценарии, приложение для разбора API github, приложение для разбора загруженных репозиторий, Dockerfile и docker-compose файлы для возможности сбора и анализа набора данных в специальных, изолированных, от основной операционной системы окружениях, дополнительные файлы и сценарии необходимые для корректной работы приложения. Архитектура приложения показана на рис. 2.

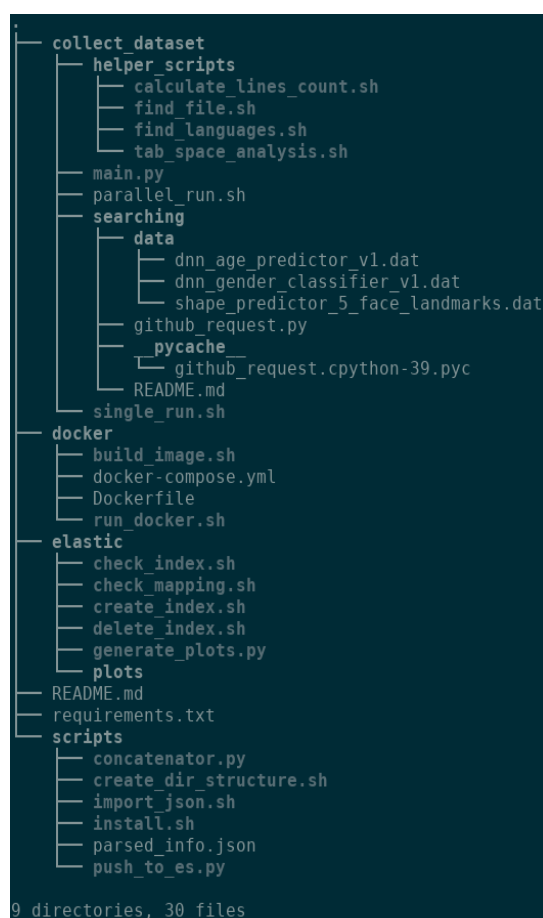


Рисунок 2 – Структура поставляемого дистрибутива

Дистрибутив был отлажен и протестирован в ОС Linux, дистрибутив Debian 11 «bullseye».

V. ОПИСАНИЕ РЕАЛИЗАЦИИ И ТЕХНОЛОГИИ

А. Архитектура программного решения

Программное решение было реализовано с использованием скриптового языка программирования Python, в качестве автоматизационных скриптов используются shell сценарии написанные на bash. Архитектура приложения показана на рис. 3.

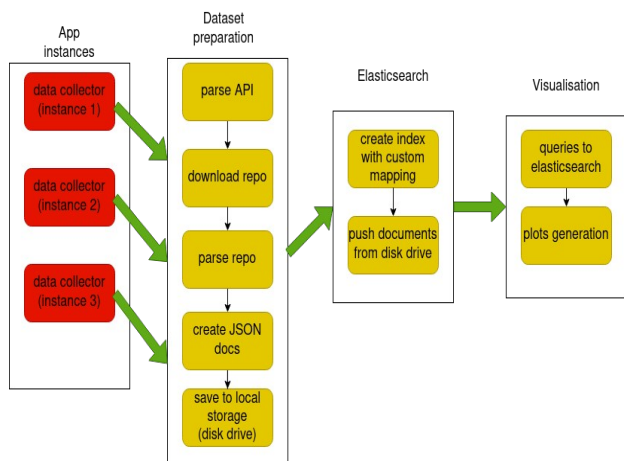


Рисунок 3 – Архитектура приложения

Архитектура разбита на три основные части:

- создания набора данных на основе информации API github и разбора репозиторийев;
- запуск экземпляра elasticsearch, создание нового индекса и загрузка полученных данных из первого пункта в индекс;
- создание запросов к индексу для получения необходимых метрик на основе набора данных и визуализация полученных метрик.

В. Описание библиотек

Для запросов к API github использовались библиотеки requests, curl. Для разбора содержимого репозитория:

- анализ репозитория для поиска различных расширений файлов и их количества;
- частота вхождение того или иного файла по его типу.

использовалась библиотека cloc. Пример выходного результата работы библиотеки cloc показан на рис. 4.

```

33 text files.
33 unique files.
4 files ignored.

github.com/AlDanial/cloc v 1.86 T=0.02 s (1215.5 files/s, 90916.3 lines/s)
-----
Language             files  blank  comment  code
-----
Python                5       220      320      932
Bourne Shell          15        69       15      348
XML                   5         0        0      207
Dockerfile            1         7         8       31
Markdown              2         3         0       30
YAML                  1         4        24       20
JSON                  1         0         0        6
-----
SUM:                 30       303      367     1574
  
```

Рисунок 4 – Результат работы cloc

Для поиска Dockerfile, файла с README использовались самописные сценарии на bash. Также все сценарии отвечающие за автоматизацию различных действий, написаны на bash.

За модуль распознавания пола по изображению отвечает библиотека dlib-models. Библиотека представляет сеть ResNet с 29 слоями, для создания этой сети была взята сеть из статьи [1], которая подверглась небольшим изменениям по слоям. Данная нейронная сеть тренировалась на наборе данных из 3 миллионов лиц. Пример распознавания пола моделью показана на рис. 5.



Рисунок 5 – Распознавание пола dlib-models с выходной визуализацией

Для визуализации данных использована библиотека matplotlib. Данная библиотека позволяет создавать различные типы графиков. API, которое предоставляет библиотека является достаточно низкоуровневым, что позволяет очень гибко настраивать различные параметры графиков, однако усложняет написание кода. Пример графика полученный с помощью matplotlib показан на рис. 6.

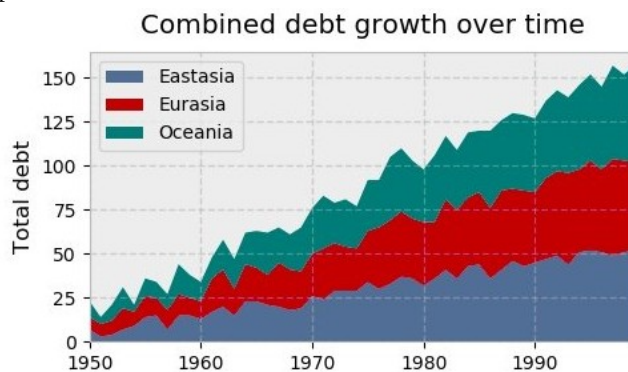


Рисунок 6 – Пример графика созданный с помощью matplotlib

С. Потребление ресурсов

Расход ресурсов процессорного времени и оперативной памяти не является узким местом для приложения, потребление ресурсов с тремя запущенными экземплярами приложения показано на рис. 7.

```

3512019 tomilin 20 0 9104 2956 2148 S 0.0 0.0 0:00.83 /usr/bin/bash
3930902 tomilin 20 0 6892 3128 2792 S 0.0 0.0 0:00.00 /usr/bin/bash/bin/bash -e /
4006379 tomilin 20 0 387M 95592 9916 S 5.1 0.6 9:07.94 /usr/bin/python3.9 ./main.p
544500 tomilin 20 0 15028 3972 3704 S 0.1 0.0 0:00.02 /usr/bin/git clone -v -
544505 tomilin 20 0 19024 6832 3368 S 0.4 0.0 0:00.06 /usr/lib/git-core/git
544501 tomilin 20 0 6692 3340 3132 S 0.0 0.0 0:00.00 /usr/lib/git-core/git
544502 tomilin 20 0 97156 14332 7556 S 0.9 0.1 0:00.12 /usr/lib/git-core/git
3648561 tomilin 20 0 65568 0 0 T 0.0 0.0 0:00.09 /usr/bin/vim.gtk3|vim elastic
1333538 tomilin 20 0 8124 8 4 S 0.0 0.0 0:00.01 /usr/bin/bash
1641760 tomilin 20 0 6892 68 24 S 0.0 0.0 0:00.00 /usr/bin/bash/bin/bash -e /
2920249 tomilin 20 0 380M 84640 8076 R 6.3 0.5 23:45.17 /usr/bin/python3.9 ./main.p
3745319 tomilin 20 0 8844 8 4 S 0.0 0.0 0:00.37 /usr/bin/bash
2535421 tomilin 20 0 6892 0 0 S 0.0 0.0 0:00.00 /usr/bin/bash/bin/bash -e /
711837 tomilin 20 0 398M 107M 8096 S 5.7 0.6 41:46.67 /usr/bin/python3.9 ./main.p
  
```

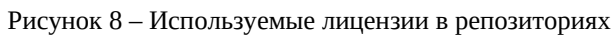
Рисунок 7 – Потребление ресурсов приложения запущенного в трех экземплярах

Как было выяснено, основная проблема связана со временем анализа репозиторийев, которое тратится на загрузку репозитория, анализ изображений разработчиков, разбор содержимого репозитория.

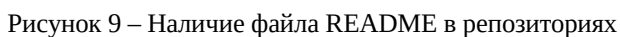
VI. РЕЗУЛЬТАТЫ АНАЛИЗА РЕПОЗИТОРИЕВ

В результате анализа набора данных были получены искомые результаты по рассматриваемым критериям.

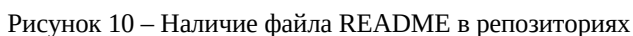
На основе анализа лицензий рис. 8, можно сделать вывод, что большинство репозиторий не используют какую-либо лицензию. Самая популярная из известных лицензий стала MIT License.



Выполнив анализ репозитория на наличие README файла рис. 9, стало понятно, что примерно половина репозитория не имеют README файла, однако с небольшим отрывом в ~1.5% лидируют репозитории в которых README присутствует.



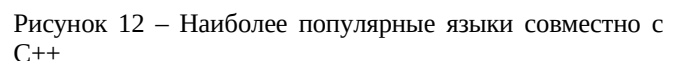
Выполнив анализ репозитория на наличие Dockerfile манифеста рис. 10, стало понятно, что меньше 1% из общего числа репозитория имеют Dockerfile.



Выполнив анализ отступов в репозиториях рис. 11, удалось выяснить, что преобладает смешанный стиль, затем идут пробелы и на последнем месте табуляция.



Выполнив анализ расширений файлов с исходным кодом рис. 12, был определен список языков, если исключить языки сборочных манифестов, то самыми популярными стали shell, python и C.



Проанализировав содержимое репозитория на наличие самых популярных систем сборки cmake, make, qmake, meson, ms build system рис. 13, было получено процентное соотношение в котором на первом месте оказалась система сборки cmake, которая считается достаточно неудобной и сложной.

СПИСОК ЛИТЕРАТУРЫ

- [1] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," In CVPR, 771–777, 2016.
- [2] O. Jarczyk, B. Gruszka, S. Jaroszewicz, "Quality analysis of open-source software", SocInfo 2014, LNCS 8851, pp. 80–94, 2014.
- [3] E. Guzman, D. Azocar, Y. Li, "Sentiment Analysis of Commit Comments in GitHub: An Empirical Study," In MSR 2014.
- [4] Y. Hu, J. Zhang, X. Bai, S. Yo and Z. Yang, "Influence analysis of Github repositories," SpringerPlus 5, 2016.
- [5] M. Beller, G. Gousios, A. Zaidman, "Oops, My Tests Broke the Build: An Explorative Analysis of Travis CI with GitHub," In MSR 2017.
- [6] P. Russell, R. Johnson, S. Ananthan, B. Harnke and N. Carlson, "A large-scale analysis of bioinformatics code on GitHub," PLoS ONE 13, 2018.
- [7] J. Cito, G. Schermann, J. Wittern, P. Leitner, S. Zumberi, H. Gall, "An Empirical Analysis of the Docker Container Ecosystem on GitHub," In MSR, 2017.

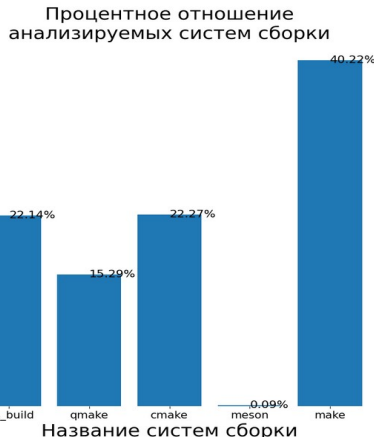


Рисунок 13 – Процентное отношение рассматриваемых систем сборки

G. Самые популярные названия веток и наличие gitflow

Выполнив анализ репозитория на наименования веток рис. 14, название master оказалось на первом месте.



Рисунок 14 – Самые популярные названия веток

Наличие комбинации веток, которое используется в gitflow рис. 15 оказалось у относительно малого числа репозиториях.

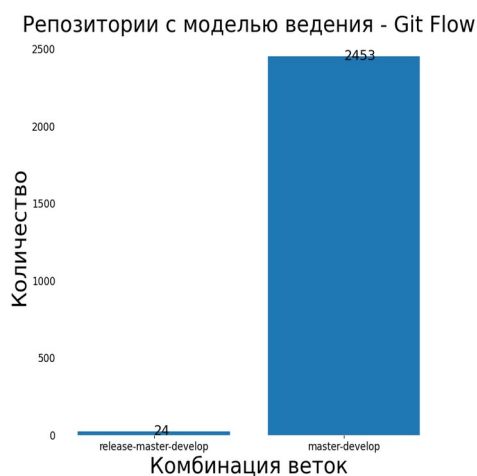


Рисунок 15 – Модель ветвления gitflow в репозиториях