

Путьков Дмитрий Александрович

Заславский Марк Маркович

Каф. МО ЭВМ, Программная инженерия,

15 декабря, 2017г.

Сбор статистики с репозиториев Github

Аннотация

В данной статье описаны исследования возможных решений для сбора информации с репозиториев по заданным критериям и обработка и сохранение данных результатов в гугл таблице. Также в исследовании описываются аналоги выбранного инструмента, описаны их сильные и слабые стороны, обосновано, показано, почему они не могут быть использованы в качестве решения поставленной задачи.

Было предложено создать инструмент, который бы собирал статистики по всем заданным критериям с репозиториев и строил бы по ним графики. Преимуществом подхода, описанного в данной статье, является покрытие всех поставленных требований к инструменту подобного рода. Приложение запускается из командной строки, интерфейсом визуализации является: терминал - при запуске и работе приложения, Google Spreadsheet - как итог работы программы.

Введение

Причиной изучения данной темы стала необходимость проверки большого количества работ студентов и просмотр преподавателям активности студентов и сравнения результатов групп между собой. Большая часть подобного рода работ представляется для проверки на Github, который также широко используется для учебных проектов на кафедре МОЭВМ <https://github.com/moevm>. Поэтому требовалось найти решение, которое бы помогло преподавателю в отслеживании активности не только групп, но и студентов, что позволило бы повысить уровень и качество проверки работ студентов и выставление оценок по более широкому количеству критериев.

Решение должно собирать статистику с репозиториях, в которых располагаются проекты: количество pull requests , commits , issues.

Цель исследования - создание инструмента, позволяющего иметь визуальное представление об изменении статистики репозиториях с течением времени, сравнения этих репозиториях и их участников между собой. Инструмент подобного рода - один из инструментов, которые необходимы на проектах с большим количеством людей, для отображения активности членов проектов. Также актуальность данного исследования подкрепляется относительно небольшим и, в большинстве своем, узконаправленными инструментами для сбора подобного рода статистики.

Обзор предметной области

Для решения задачи необходимо исследовать существующие решения. Для начала были рассмотрены уже существующие способы сбора статистики, проведено сравнения их между собой по выделенными нами критериями:

сбор статистики в Github

Github имеет свой собственный инструмент для сбора статистики с репозитория широкому спектру критериев. Также он автоматически строит по этим критериям графики.

Стандартные свойства Git [5]

У Git есть команда `git shortlog -s -- author=author@gmail.com` , которая получает количество commits по данному пользователю. Также можно получить более конкретную информацию(а именно : сколько строчек добавлено\удалено).

Caelum-git- reports [1]

Это генератор статистики для репозиториях. Он может получать количество commits по всему репозиторию и по отдельному человеку, при этом учитывая промежуток времени, за который требуется пользователю.

"Ручной" сбор информации

Вероятно, самый распространенный, самый полный и самый долгий способ сбора статистики: Просматривать информацию ин-

дивидуально по каждому репозиторию и фиксируя все данные, которые интересуют пользователя.

Критерии для сравнения:

Скорость получения данных

Этот критерий определяет, сколько понадобится времени для получения интересующих данных.

Ассортимент получаемых данных.

Данный критерий определяет, какие именно данные : commits, issues или pull requests могут быть получены при использовании одного из способов.

Визуализация данных

Этот критерий оценивает, полученную информацию. Насколько удобно будет воспринимать в виде таблиц или графиков.

Удобство сравнения статистик разных репозиториев

Определяет, насколько удобно использовать способ для сравнения большого числа разных репозиториев и/или участников.

Таблица 1. Сравнение существующих аналогов.

Аналоги	Скорость	Ассортимент	Удобство для чтения	Удобство сравнения
Инструмент Github	Очень быстро/автоматически	Широкий	Имеются графики	Неудобно
Caelum	Быстро	Только commits	Без графиков	Неудобно
Ручной способ	Медленно	Широкий	Без графиков	Неудобно

Таким образом, можно сделать вывод, что некоторые из представленных в таблице способ эффективны в большинстве случаев, однако ни один из них

не в состоянии выполнить все поставленные задачи полностью. Самым удобным можно выделить базовый инструмент GitHub, единственный минус которого в том, что графики строятся только по одному репозиторию. В остальном же, и количество информации, и широта рассматриваемых критериев полностью, и даже больше, покрывает нужды. Что касается остальных критериев, то они либо узконаправленные, либо слишком медленные, как ручной способ.

Таким образом, можно сделать вывод, что существует достаточное количество аналогов, но ни один из них полностью не подходит для достижения поставленной цели. Именно поэтому было решено создать инструмент, который бы показывал хорошие результаты по всем данным критериям.

Выбор метода решения

Видом решения задачи является скрипт, запускаемый из терминала командой.

Требования к решению:

- Сравнение несколько репозиториев между собой.
- Сравнение участников всех репозиториев между собой
- Визуализация представления статистики в виде графиков
- Графический показ изменения статистики с течением времени.

На основании сравнения аналогов было решено, что:

- Решение должно представлять собой инструмент, который бы требовал минимальных затрат времени пользователя.
- Этот инструмент должен быть удобен в использовании, то есть следует свести действия пользователя к минимуму.
- Инструмент должен содержать в себе и возможность сравнения репозиториев и их участников между собой, а также сохранять статистику за период предыдущих использований скрипта, для визуального представления развития репозитория.
- Кроме этого, было решено внедрить такое свойство, как построение графиков по этим репозиториям для более удобного просмотра статистики.
- Далее, следовало учесть возможность того, что некоторые участники проектов не важны для статистики. Пример: преподаватели в репозитории групп. Поэтому инструмент должен был быть снабжён возможностью игнорировать статистику указанных пользователей.

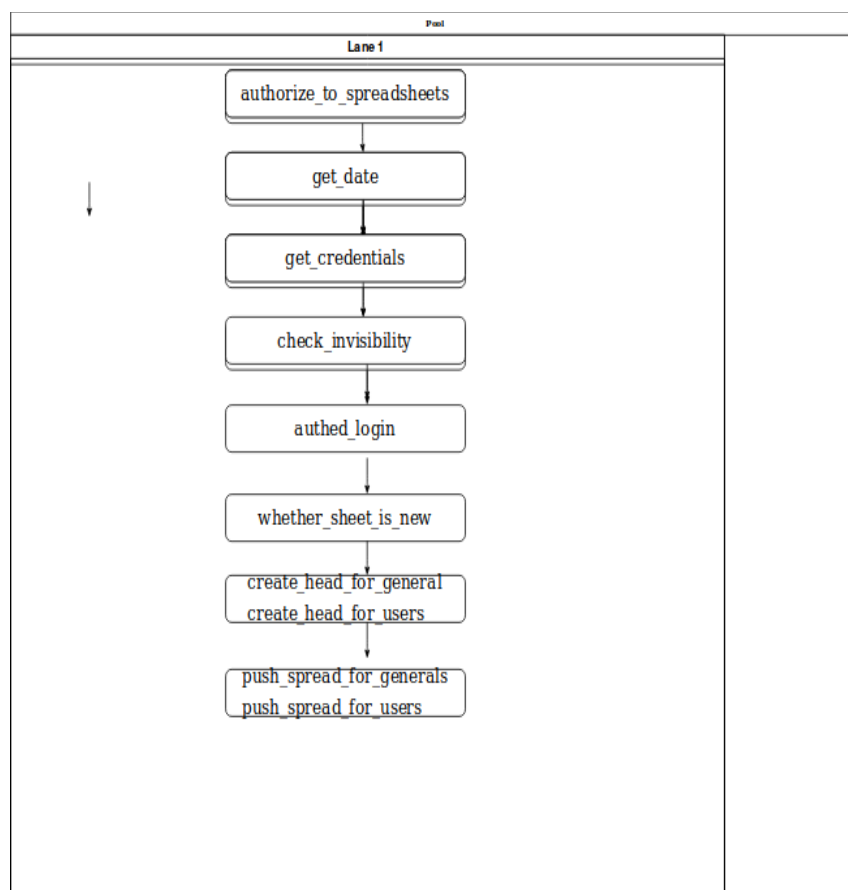
- Также следовало учесть, что может появиться необходимость сравнения участников из разных репозиторий между собой, поэтому должны быть построены отдельные графики только по людям.

Архитектура

* Сценарий использования:

1. Запустить скрипт из терминала командой.
2. По завершении работы скрипта, запустить Google apps скрипт, если это первый запуск скрипта для заданных репозиторий.
3. Просмотр результатов работы приложения в Google Spreadsheet.

Алгоритм работы в общем виде можно представить следующим образом:



1. Подключение к Google

Spreadsheet - функция `authorize_to_spreadsheets`. Для доступа к Google Spreadsheet требуется документ в формате `.json`, генерируемый прямо при создании проекта Google.

Затем, происходит подключение к таблице, передавая путь к этому файлу и ссылке <https://spreadsheets.google.com/feeds>

2. Генерация страниц и запись текущей даты - функция `get_date`. Теперь, имея доступ к таблице, следует создать список листов, чтобы в дальнейшем записывать добавлять туда информацию. В переменную записывается текущая дата в формате дд-мм-гг.

3. Получение и парсинг аргументов - функция `get_credentials`. Используя библиотеку `argparse`, функция считывает следующие аргументы:

- CREDENTIALS название файла, хранящие `client_id` `client_secret` для Github.
- REPO список репозитория, информацию о которых требуется получить.
- INVISIBILITY_ACCESS файл, содержащий список участников, информацию о которых не требуется учитывать.
- LINK Ссылка на Google Spreadsheet, в которую будет записываться информация.

4. Учёт исключений, переданных как аргументы - функция `check_invisibility`. Создание списка участников, полученный через `INVISIBILITY`.

5. Запись данных в таблицу - функция `creating_worksheet`, включает следующие подпункты :

- Авторизация на Github - функция `authed_login`. [3]
- Проверка, создаем ли мы новую таблицу или обновляем уже существующую - функция `whether_sheet_is_new`.
- Создаем оформление : на одном листе будет 4 таблицы, каждый содержит название, столбцы - дату обновления, строки - репозитории или участники - функция `create_head_for_general` или `create_head_for_users`.
- Для каждого репозитория из списка вызывается функция, общая сложность которой $O(\text{countOfAssignees}) + O(\text{countOfIssues}) + O(\text{countOfCommits}) + O(\text{countOfPull})$ - функция `Get_info`.
- Запись полученных данных в таблицы: `push_spread_for_generals` или `push_spread_for_users`.

6. Данный пункт является самым долгим по времени выполнения. Это связано, прежде всего, с подключением к Github data и записи их в Google spreadsheet.

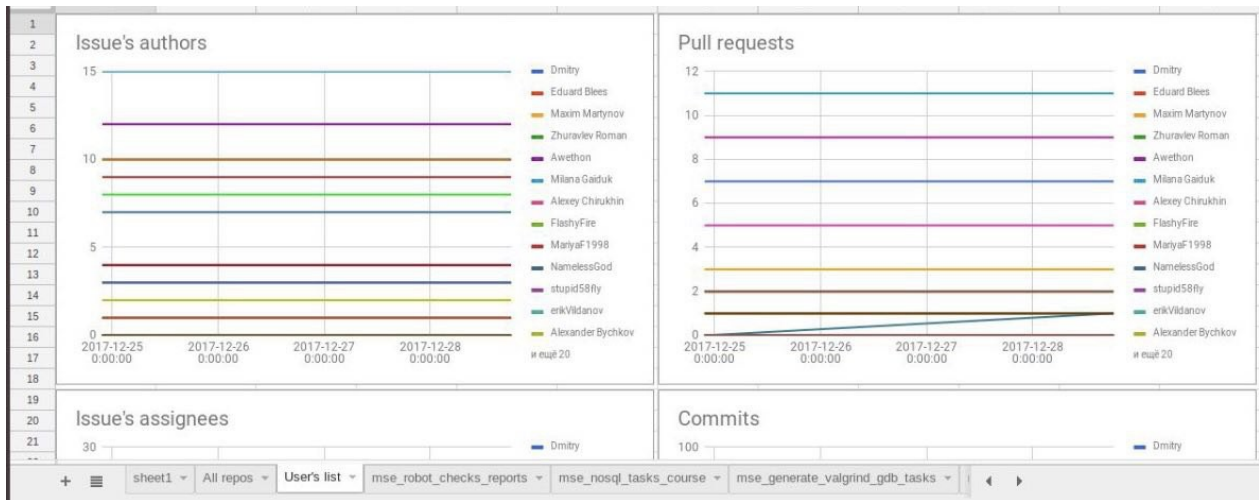
7. Скрипт, написанный на GoogleAPPScript [2], создающий графики:

- При первой записи данных в новую таблицу следует запустить скрипт, передав в него ссылку на таблицу. В дальнейшем это не требуется, т.к. скрипт создает графики по изначально очень большому рэнджу из клеток, поэтому графики будут изменяться самостоятельно.
- Входные данные:
 - Ссылка на таблицу
 - Файл с client_id, client_secret github
 - Файл со списком репозиторийев
 - Файл со списком людей-исключений(необязательный параметр)
- Выходные данные:
 - Ссылка на созданную таблицу
 - Ссылка на гугл-скрипт

Пользовательский интерфейс:

```
./get_info.py --credentials credFile --repos reposFile --invisibility  
peopleFile --link link
```

Результат в виде графиков :



- Ссылка на таблицу
 - --credentials - файл с client_id, client_secret github
 - --repos - файл со списком репозиторийев
 - --invisibility - файл со списком людей-исключений(необязательный параметр)
- Выходные данные:
 - Ссылка на созданную таблицу
 - Ссылка на Google App script

Техническая реализация:

- Библиотека PyGitHub для Python (2 and 3) позволяющая получать доступ к GitHub API v3. Эта библиотека позволяет управлять ресурсами GitHub: repositories, user profiles, and organizations в приложении, написанном на Python.[4]
- Модуль argparse упрощает создание удобных интерфейсов командной строки. Программа определяет, какие аргументы она требует, а argparse будет определять, как разбирать их из sys.argv. Модуль argparse также автоматически генерирует сообщения справки и использования и выдает ошибки, когда пользователи дают программе недопустимые аргументы.[6]
- Библиотека gspread позволяет управлять Google Spreadsheet.[7]
- Запись в таблицу по названию или ссылке.

- `oauth2client.service_account` – Это клиентская библиотека для доступа к ресурсам, защищенным OAuth 2.0. [8]
- Модуль `datetime` предоставляет классы для управления датами и временем как простым, так и сложным способом. Хотя арифметика даты и времени поддерживается, фокус реализации заключается в эффективном извлечении атрибутов для форматирования и обработки вывода.[9]
- Интерфейсом пользователя является командная строка.
- Данные для участников хранятся в словаре, которая создается при итерации по списку репозиторий. Она обнуляется каждый раз при записи данных нового репозитория. Также существует общий список данных, в который заносится имя репозитория и информация по нему. Такой словарь имеет жизненный цикл, равный жизненному циклу работы программы.
- Тесты, используемые при проверки работы :
Интеграционный тест.

Заключение

Таким образом, был создан инструмент, который оптимально подходит для решения поставленных задач. Скорость сбора информации с Github невысокая, но объем полученной информации высокий. Также полученная информация отображается в виде графиков- такая визуализация информации является удобной для восприятия. Можно выделить следующие недостатки данного решения :

1. Необходимость самим создавать инструмент
2. Невысокая скорость сбора информации с репозиторий
3. Привязка к Google Spreadsheet и гугл-скриптам. В этом случае необходимо хранить файл-ключ в папке со скриптом для корректной работы приложения.
4. Не автоматизированная работа гугл-скрипта: при первом сборе информации по репозиторию, требуется вручную запустить гугл-скрипт для создания графиков. В дальнейшем, при обновлении информации по репозиторию, этого делать не требуется.

В итоге, цель, ради которой создавалось приложение, была достигнута. Были решены все поставленные задачи. Однако, в дальнейшем данный инструмент может быть доработан в следующих направлениях : полная автоматизация работы и внедрение графического интерфейса для более комфортной работы с приложением, увеличение скорости работы скрипта.

Список литературы

1. github.com/caelum/git-reports/wiki [1]
2. <https://developers.google.com/apps-script/> [2]
3. <https://www.twilio.com/blog/2017/02/an-easy-way-to-read-and-write-to-a-google> [3]
4. <https://github.com/PyGithub/PyGithub> [4]
5. <https://help.github.com/articles/viewing-a-repository-s-network/> [5]
6. <https://pypi.python.org/pypi/argparse> [6]
7. <https://gspread.readthedocs.io/en/latest/> [7]
8. <https://github.com/google/oauth2client> [8]
9. <https://docs.python.org/2/library/datetime.html> [9]