

Новый ли репозиторий?

2

3

4

5

6

Проверка, есть ли этот репозиторий в базе данных

git -C {repo_path} log -p --unified=0 -w --ignore-blank-lines

Запись информации в датасет

Очистка от ненужных строк

1 commit c1f3475f3348faa6e2c6115eeec01e11c0f10b38

Создание новых столбцов

Для каждой строки определяется ее тип:

путь к файлу до и после изменений;

текст или код, который поменялся

дополнительно обрабатываются

номер строки, на которой начался блок изменений и количество затронутых строк

Перед записью в новое поле данные очищаются от лишних символов и при необходимости

Добавляются дополнительные поля для удобства при дальнейшей работе (например, segment_id)

a/ETL.R

Создается второй датасет, в котором весь измененный текст в рамках каждого блока объединяется в

одну строку. При этом из нее удаляются все нечитаемые символы, в т.ч. пробелы, табуляции

filter(!grepl('deleted',lines))%%filter(!grepl('new',lin..

if(!(grepl("commit",dflines[i])grepl("---",dflines[i])gre...

write_parquet(changesDf,paste(gsub("/","_",strsplit(...

renamefromgithub.com_tidyverse_ggplot2.git.parqu...

git_diff_cmd-glue('git-C{repo_dir}log-p--unified0-w--.

library(arrow)repo_url-"https://github.com/tidyverse..

dir repo-getRepo("https://github.com/tidyverse/ggpl..

repo_dir-getRepo(repo_url)git_log_cmd-glue('git-C{r...

write_parquet(git_commit_history,paste(gsub("/","_".

Поля, где были пропущены значения, заполняются данными из поля выше

src_file

a/ETL.R

/dev/null

b/ETL.R.

b/ETL.R

b/ETL.R

b/README.md

a/README.md

dst_file

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

b/FTI R

b/ETL R

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R.

b/ETL.R

dst_file

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R.

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

b/ETL.R

Данные из двух датасетов (см. шаг 6 и 8) объединяются в один датасет через segment id

b/README.md

count_del

47

28

0

13

Заполнение пропущенных значений

Удаляются все записи, где значение src_code было NA

b/ETL.R

NA TRUE

NA TRUE

NA TRUE

2 TRUE

NA FALSE

NA FALSE

3 TRUE

NA FALSE

NA FALSE

NA FALSE

0 TRUE

count_del *

0

0

1

1

0

0

1

0

1

1

0

0

0

0

0

0

0

0

0

0

0

0

0

0

count_del = start_add

0

1

0

1

0

1

1

1

0

0

1

3

0

0

1

1 while(in){

0 print(j)

0 print("pull")

0 print("clone")

1 while(i50000){

start_del *

47

47

54

54

86

86

88

90

15

18

21

39

39

39

39

39

39

39

39

39

39

39

39

39

39

47

54

86

88

90

15

18

21

39

2

11

25

28

0

13

13

start_add

count_add

48

56

89

16

37

1

start del

start_add *

48

48

56

89

92

95

16

18

37

37

37

37

37

37

37

37

37

37

37

37

48

56

89

92

14

16

18

37

13

27

39

1

13

2 filter(!grepl('deleted',lines))%% filter(!grepl('new',lin...

3 if(!(grepl("commit",dflines[i])grepl("---",dflines[i])gre.

1 write_parquet(changesDf,paste(gsub("/","_",strsplit(...

54 renamefromgithub.com_tidyverse_ggplot2.git.parqu...

54 git_diff_cmd-glue('git-C{repo_dir}log-p--unified0-w--..

3 library(arrow)repo_url-"https://github.com/tidyverse...

11 dir_repo-getRepo("https://github.com/tidyverse/ggpl..

11 repo_dir-getRepo(repo_url)git_log_cmd-glue('git-C{r...

2 write_parquet(git_commit_history,paste(gsub("/","_".

28 library(glue)#Длявызовасистемныхкомандlibrary(...

1 imgsrc"images/План% 20работы-02.png"data-fig-ali... FALSE

count_add

2

1

3

0

0

0

54

0

11

28

count_add = segment_id

2

1

3

3

0

0

0

0

54

54

54

54

54

54

54

54

1

1

2

3

3

5

7

8

9

9

9

9

9

9

9

9

9

9

9

9

9

9

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

15

repo

TRUE

git-analysis

FALSE

FALSE

TRUE

FALSE

FALSE

FALSE

FALSE

TRUE

1 filter(!grepl('^deleted', line

3 if (!(grepl("^commit", df\$

5 write_parquet(changesDf, pa

4 @@ -47,0 +48,2 @@ df <- git_diff_df %>%

5 + filter(!grepl('^deleted', lines)) %>%

6 + filter(!grepl('^new', lines)) %>%

10 @@ -86,0 +89,3 @@ while (i < n) {

7 @@ -54 +56 @@ j = 1

9 + while (i < 50000) {

• хэш коммита

10 @@ -86,0 +89,3 @@ while (i < n) {

14 @@ -88 +92,0 @@ while (i < n) {

3

5

6

7

8

9

10

11

12

13

14

15

16

commit

11 + if (!(grepl("^commit", df\$lines[i]) | grepl("^---",

17 + write_parquet(changesDf, paste(gsub("/", "_", strs... FALSE

segment_id † is_add †

1 TRUE

2 FALSE

2 TRUE 3 TRUE

4 FALSE

5 TRUE

6 FALSE 7 FALSE

8 FALSE

9 FALSE

11 FALSE

12 FALSE

12 TRUE

13 TRUE

Удаляются ненужные столбцы

1 clf3475f3348faa6e2c6115eeec01e11c0f10b38

2 c1f3475f3348faa6e2c6115eeec01e11c0f10b38

3 c1f3475f3348faa6e2c6115eeec01e11c0f10b38

4 clf3475f3348faa6e2c6l15eeec0lel1c0f10b38

5 c1f3475f3348faa6e2c6115eeec01e11c0f10b38

12 2d49d9c9dd7d77f91c906869548ae799938cc99d

13 2d49d9c9dd7d77f91c906869548ae799938cc99d

14 2d49d9c9dd7d77f91c906869548ae799938cc99d

15 2d49d9c9dd7d77f91c906869548ae799938cc99d

Удаление дубликатов

8

9

10

2

21 2d49d9c9dd7d77f91c906869548ae799938cc99d

22 2d49d9c9dd7d77f91c906869548ae799938cc99d

Объединение кода

Sun Mar 30 00:45:05 2025 +0300

-47,0 +48,2 @@ df <- git_diff_df %>%

+ filter(!grepl('^deleted', lines)) %>% + filter(!grepl('^new', lines)) %>% @@ -54 +56 @@ j = 1 -while (i < n) { +while (i < 50000) {

1 commit c1f3475f3348faa6e2c6115eeec01e11c0f10b38

2 Author: Maria <maria.ut.005@yandex.ru> 3 Date: Sun Mar 30 00:45:05 2025 +0300

11 @@ -47,0 +48,2 @@ df <- git_diff_df % >% 12 + filter(!arepl('^deleted', lines)) %>% 13 + filter(!grepl('^new', lines)) %>%

17 @@ -86,0 +89,3 @@ while (i < n) {

Changes in files

7 diff -- git a/ETL.R b/ETL.R 8 index 19bdfd0..481335a 100644

14 @@ -54 +56 @@ j = 1 15 -while (i < n) { 16 + while (i < 50000) {

10 +++ b/ETL.R

lines

2 --- a/ETL.R

3 +++ b/ETL.R

8 -while (i < n) {

@@ -88 +92,0 @@ while (i < n) {

@@ -90,0 +95 @@ while (i < n) {

Changes in files

--- a/ETL.R +++ b/ETL.R

diff --git a/ETL.R b/ETL.R index 19bdfd0..481335a 100644

Если да, то получаем номер последнего коммита last_commit

git -C {repo_path} log {last_commit}..HEAD -p --unified=0 -w --ignore-blank-lines

commit clf3475f3348faa6e2c6115eeec01e11c0f10b38 (HEAD -> master, origin/master, origin/HEAD)
Author: Maria <maria.ut.005@yandex.ru>

Получение изменений в каждом файле для каждого коммита

@@ -86,0 +89,3 @@ while (i < n) {
+ if (!(grep]("^commit", df\$lines[i]) | grepl("^---", df\$lines[i]) | grepl("^\\+\\+", df\$lines[i]) | grepl("^@@", df\$lines[i]))) {

+write_parquet(changesDf, paste(gsub("/", "_", strsplit(repo_url, "//")[[1]][2]),"_changes.parquet", sep=""))

*не повторяется в программе, представлен для полноты понимания алгоритма

git-analysis

git-analysis

git-analysis

git-analysis

*на скриншот не поместился столбец is_add, который указывает на характер изменений (+ или -) 1 commit clf3475f3348faa6e2c6l15eeec0lel1c0f10... TRUE 1 c1f3475f3348faa6e2c6115eeec01e11c0f10b38 2 --- a/ETL.R FALSE 3 +++ b/ETL.R FALSE 4 @@ 47,0 +48,2 @@ df <- git_diff_df % >% FALSE 5 + filter(!grepl('^deleted', lines)) %>% FALSE 6 + filter(!grepl('^new', lines)) %>%

FALSE

FALSE

FALSE

FALSE

FALSE

while(in){ while(i50000){

print("pull")

print("clone")

#gsub("/","\\\\",dir)

- 6 c1f3475f3348faa6e2c6115eeec01e11c0f10b38 a/ETL.R 7 c1f3475f3348faa6e2c6115eeec01e11c0f10b38 8 clf3475f3348faa6e2c6115eeec01e11c0f10b38 a/ETL.R c1f3475f3348faa6e2c6115eeec01e11c0f10b38 a/ETL.R 2d49d9c9dd7d77f91c906869548ae799938cc99d 2d49d9c9dd7d77f91c906869548ae799938cc99d a/ETL.R
- src_file 1 clf3475f3348faa6e2c6l15eeec0lellc0f10b38 a/ETL.R 3 c1f3475f3348faa6e2c6115eeec01e11c0f10b38 a/ETL.R 5 clf3475f3348faa6e2c6115eeec0le11c0f10b38 a/ETL.R 8 c1f3475f3348faa6e2c6115eeec01e11c0f10b38 a/ETL.R c1f3475f3348faa6e2c6115eeec01e11c0f10b38

2d49d9c9dd7d77f91c906869548ae799938cc99d

2d49d9c9dd7d77f91c906869548ae799938cc99d

fa3a78ff139856b542759417fadad2cd76b05461

fa3a78ff139856b542759417fadad2cd76b05461

fa3a78ff139856b542759417fadad2cd76b05461

fa3a78ff139856b542759417fadad2cd76b05461

a4fee3e2f8ebd87aeb31d86a962360166aea9b15

Объединение данных

1 c1f3475f3348faa6e2c6115eeec01e11c0f10b38

16 fa3a78ff139856b542759417fadad2cd76b05461

17 f1b13932a6e5dd9d8f3873514f53880f56547511

18 a4fee3e2f8ebd87aeb31d86a962360166aea9b15

Создание датасета

repo

[▲] repo

2 dotfiles

6 learnGo

7 learnWeb

8 magma

9 MersOff

10 mqwerty 11 mqwerty_bot

14 reports_i2z1

15 Students_db

17 tic-tac-toe 18 TMP

2 adff96e6b01eef3e3d4dcc3b62457a6b51365522

3 100b53956053f4d30bb9bae201c249a49c6688ce

5 d6b76d02fb5262c264000e3931e2e1989aaa5460

6 5a0119593db5ec60dff88c3f4cd719fda3d324f6

7 de87cef63cld9ea6854efa7dc97b06ld9c9144la

9 b12375c0489e41017d5eb86eed046d1ef89fcec7

LO 878b140ba0163722ffef8b6a0030c58954b8f05a

L2 b326babcbf2514393a7bb08feb3410608b52ef68

L3 d4ca1080fe0e14bf51d5412b10c159bcc2d2ebbd

Мария

@ MariUt005

L1 c8e479642315a7353998c706c838debd1dd8e222

8 a68e9787ff8f67cb6c53f70eda3e2adfe6460fb1

4 c7c02c6db07f809edbefac7b141ff458fbb867f3

16 tetris

12 qq 13 regex

91 f1b13932a6e5dd9d8f3873514f53880f56547511

22 f418704449370c5ba16131037f6eed177f018247

11 2d49d9c9dd7d77f91c906869548ae799938cc99d

12 2d49d9c9dd7d77f91c906869548ae799938cc99d

Удаляются все дублирующиеся записи

2 c1f3475f3348faa6e2c6115eeec01e11c0f10b38 a/ETL.R 3 c1f3475f3348faa6e2c6115eeec01e11c0f10b38 a/ETL.R b/ETL.R 54 4 c1f3475f3348faa6e2c6115eeec01e11c0f10b38 86 a/ETL.R b/ETL.R 5 c1f3475f3348faa6e2c6115eeec01e11c0f10b38 a/ETL.R b/ETL.R. 6 c1f3475f3348faa6e2c6115eeec01e11c0f10b38 7 2d49d9c9dd7d77f91c906869548ae799938cc99d 15 a/ETL.R b/ETL.R. 8 2d49d9c9dd7d77f91c906869548ae799938cc99d b/ETL.R. 18 9 2d49d9c9dd7d77f91c906869548ae799938cc99d 21 b/ETL.R. 10 2d49d9c9dd7d77f91c906869548ae799938cc99d 39 b/ETL.R 11 2d49d9c9dd7d77f91c906869548ae799938cc99d a/ETL.R b/ETL.R 39 12 fa3a78ff139856b542759417fadad2cd76b05461 b/ETL.R 13 fa3a78ff139856b542759417fadad2cd76b05461 b/ETL.R 14 fa3a78ff139856b542759417fadad2cd76b05461 25 a/ETL.R b/ETL.R 15 fa3a78ff139856b542759417fadad2cd76b05461 a/ETL.R b/ETL.R. 25

a/ETL.R

/dev/null

a/README.md

Добавление датасета в базу данных

Этап 3. Сохранение пути к репозиторию

Полученный датасет записывается в git_diff в DuckDB

1 git-analysis D:/Творчество/git-analysis/MariUt005/git-analysis

D:/Творчество/git-analysis/MariUt005/learnWeb D:/Творчество/git-analysis/MariUt005/magma

D:/TBopчecтвo/git-analysis/MariUt005/MersOff

D:/Творчество/git-analysis/MariUt005/mqwerty

D:/Творчество/git-analysis/MariUt005/qq

D:/Творчество/git-analysis/MariUt005/regex

D:/Творчество/git-analysis/MariUt005/tetris

D:/Творчество/git-analysis/MariUt005/mqwerty_bot

D:/TBopчecтвo/git-analysis/MariUt005/reports_i2z1

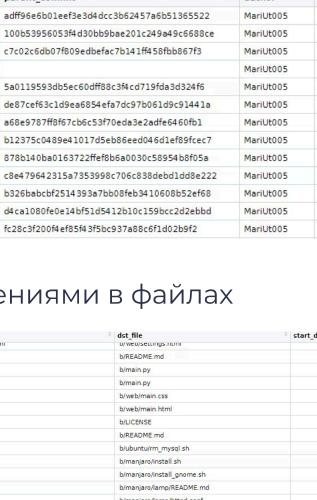
D:/Творчество/git-analysis/MariUt005/Students_db

Добавление датасета в базу данных

Полученный датасет записывается в repo_path в DuckDB

* src_file

- 1 book_reader D:/Творчество/git-analysis/MariUt005/book_reader D:/Творчество/git-analysis/MariUt005/dotfiles 3 Enlightsion D:/Творчество/git-analysis/MariUt005/Enlightsion 4 exam mqwerty D:/Творчество/git-analysis/MariUt005/exam mqwerty 5 git-analysis D:/Творчество/git-analysis/MariUt005/git-analysis D:/Творчество/git-analysis/MariUt005/learnGo
- D:/Творчество/git-analysis/MariUt005/tic-tac-toe D:/Творчество/git-analysis/MariUt005/TMP езультат Датасет с краткой историей коммитов parent_commit author 1 75f5da303f0baaf391961148af1690c5ac776ec2 adff96e6b01eef3e3d4dcc3b62457a6b51365522 2021-07-24 09:22:06 +0300 | Some latest changes



Кристина

@ gigwrld

2021-07-09 07:07:52 +0300 Update README.md

2021-03-14 10:03:29 +0300 Add files via upload

2021-03-14 10:02:27 +0300 Add files via upload

2020-09-11 18:58:21 +0300 Update install.sh

2020-09-11 18:49:22 +0300 Create readme

2020-05-08 17:18:01 +0300 upd2

2021-03-14 10:03:17 +0300 Delete install_gnome.sh

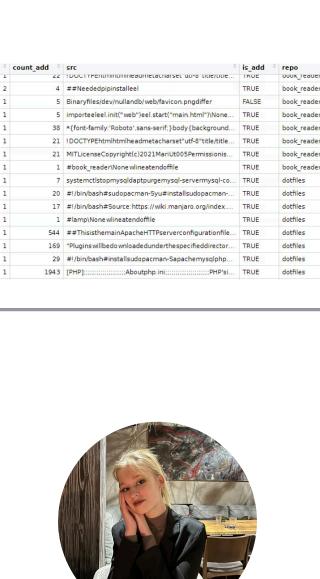
2021-07-09 07:00:39 +0300 Initial commit

2021-07-09 07:05:33 +0300 | Created main menu page

2021-06-27 07:09:22 +0300 | Added rm mysql.sh in ubuntu/

2021-04-14 23:41:27 +0300 Added LAMP and reorganized some files

2020-05-08 17:52:39 +0300 | Change the branch of the repository "lamp"



Александра

@ gusarovaal

Made with **GAMMA**

book reader

book_reader

book_reader

book_reader

dotfiles

dotfiles

dotfiles

dotfiles

dotfiles

dotfiles

dotfiles

dotfiles

Датасет с изменениями в файлах 31 /31304303100441391901140411090C34C//0eCZ 32 adff96e6b01eef3e3d4dcc3b62457a6b51365522 a/README.md b/README.md 34 100b53956053f4d30bb9bae201c249a49c6688ce /dev/null b/main.py 35 100b53956053f4d30bb9bae201c249a49c6688ce /dev/null b/web/main.css 37 c7c02c6db07f809edbefac7b141ff458fbb867f3 b/LICENSE 38 c7c02c6db07f809edbefac7b141ff458fbb867f3 b/README.md 39 d6b76d02fb5262c264000e3931e2e1989aaa5460 40 5a0119593db5ec60dff88c3f4cd719fda3d324f6 b/manjaro/install.sh 41 5a0119593db5ec60dff88c3f4cd719fda3d324f6 b/manjaro/install_gnome.sh 42 5a0119593db5ec60dff88c3f4cd719fda3d324f6 b/maniaro/lamp/README.md 43 5a0119593db5ec60dff88c3f4cd719fda3d324f6 b/manjaro/lamp/httpd.conf 45 5a0119593db5ec60dff88c3f4cd719fda3d324f6 b/manjaro/lamp/install.sh b/manjaro/lamp/php.ini 46 5a0119593db5ec60dff88c3f4cd719fda3d324f6 Авторы