

# Test Task for Project

## Logical coupling analysis in Git repositories

Roman Kovalev

May 2024

### Introduction

The test task for this project involves developing a tool to calculate pairs of developers who contribute most frequently to the same files/modules in a GitHub repository. This report describes what has been implemented for the test task, but its main value lies in discussing what else can be done for this problem.

### Problem Analysis

So, we have a GitHub repository. This means we have a list of commits in the repository with all the contained information: commit date, author, and modified files.

To begin, it makes sense to limit the number of commits we will analyze. This is necessary because the idea of "grouping developers by similarity" implies up-to-date data based on the set of active developers and the interests of current project participants. Additionally, it's worth mentioning that repositories of large projects include dozens or even hundreds of daily commits, which can provide us with a lot of irrelevant information. Since we are interested in current information, it is reasonable to limit ourselves to commits from the past few days.

Now that we have defined our input data, let's introduce a numerical estimation of the "similarity" between two developers based on the files/modules they have modified. This can be done in several ways, but we will choose the simplest: the number of files where changes were made by both the first and second developers.

We can then consider a complete weighted graph where vertices represent developers and edge weights represent the number of shared files for each pair.

Then, we can rephrase the original problem as follows: *Find the maximum weight matching in a complete weighted graph with no negative weights.*

### Implemented Solution

In this task, I implemented a greedy algorithm that, at each iteration, selects a pair of developers who are most similar among all available developers.

The time complexity of this algorithm is  $O(n^3)$ .

Unfortunately, the greedy algorithm does not guarantee that the found matching will have the maximum weight, but we will obtain some pseudo-optimal division.

## Exact Solution

Let's reduce our problem to the Assignment Problem.

### Assignment Problem:

Given a bipartite weighted complete graph. We need to find a maximum weight matching.

Let the original graph be represented as:  $G = \{V, E = V \times V, f : E \rightarrow \mathbb{Z}_+\}$  Let's duplicate each vertex - set  $V'$  - and consider graph  $H$ :

$$H = \{V \cup V', E', g : E' \rightarrow \mathbb{Z}_+\}$$

The weight of an edge in graph  $H$  between two vertices is zero if they both belong to the same set  $V$  or  $V'$ , and is equal to the weight of the edge between their preimages in the original graph otherwise. The weight of an edge between copies of the same vertex is also zero.

Now we have a bipartite weighted complete graph  $H_{n,n}$  vertices. If we find a solution to the Assignment Problem on this graph, we will straightforwardly restore a maximum weight matching on the original graph.

It is worth mentioning that the Assignment Problem is solvable in  $O(n^4)$ . The reduction of the original problem to the Assignment Problem works in  $O(n^2)$ , so the best division into pairs can be achieved in  $O(n^4)$ .

In the [paper](#), it is shown that there are randomized solutions for the Assignment Problem with worse asymptotics but better practical performance. In our case, it would be reasonable to use one of these solutions.

## Algorithm Demonstration

A public repository of the library [pytorch](#) was used as the test repository.

1. Run the file *Main.py*.
2. Specify the full path to the cloned repository (*/home/roman-not-hehe/Desktop/sandboxes/pytorch*).
3. Specify the number of days to analyze the repository.
4. Obtain information about the number of commits by authors who made changes to the repository during this time and the total number of commits.
5. Since the pairs obtained might likely be "leftover" pairs of developers, meaning their "similarity" is zero, we need to choose whether to display these pairs on the screen or not.
6. Next, the program displays pairs of developers with the number of files they both modified.
7. Finally, the program outputs the average number of shared files per pair, as well as their median value.

```
Run - Intershps_summer_2024
Run: main x
Please enter the path to your git-repo:
/home/roman-not-hehe/Desktop/sandboxes/pytorch
For how many days back do you want the information?
10
You have requested information for the last 10 days
Number of commits: 96
Number of contributors: 62
Do you want to output zero pairs? yes/no
no
Printing pairs of authors and the number of shared files:
eellison -- Sam Larsen : 1
Edward Z. Yang -- Avik Chaudhuri : 5
PyTorch MergeBot -- Daniel Javady : 5
Nikita Shulga -- aaitzhan : 1
Jiang, Yanbing -- Pearu Peterson : 1
Stefan-Alin Pahontu -- Mikayla Gawarecki : 1
Denis Vieriu -- Roy Hvaara : 1
David Chiu -- Michael Lazos : 1
Jez Ng -- David Berard : 1
Chien-Chin Huang -- Brian Hirsh : 1
The average 'similarity' of coupling: 0.5806451612903226
The median 'similarity' of coupling: 0
Process finished with exit code 0
```

The first example

```
Run - Intershps_summer_2024
Run: main x
/home/roman-not-hehe/Desktop/NUP/04_spring_2024/Intershps_summer_2024/.venv/bin/python /home/roman-not-hehe/Desktop/NUP/04_spring_2024
Please enter the path to your git-repo:
/home/roman-not-hehe/Desktop/sandboxes/pytorch
For how many days back do you want the information?
365
You have requested information for the last 365 days
Number of commits: 12786
Number of contributors: 1043
Do you want to output zero pairs? yes/no
no
Printing pairs of authors and the number of shared files:
JenDL -- Zheng, Zhaoqiong : 1
Aleksai Nikiforov -- vinithakv : 12
Apuva Jain -- y-sq : 1
IvanKobzarev -- Roger Lam : 2
Connor Baker -- TachikakaMin : 1
Alexandre Ghelfi, PhD -- Milan Straka : 1
Noam Siegel -- Dmitry Ulyanov : 1
katotaisei -- Amr Elshennawy : 1
Lucy Qiu -- Justin Yip : 11
dshi7 -- Ronan Gautier : 1
Lengyue -- Oren Leung : 2
Nikita Karetnikov -- Mwiza Kunda : 10
Flavio Sales Truzzi -- blzheng : 4
```

```
Run - Intershps_summer_2024
Run: main x
Please enter the path to your git-repo:
/home/roman-not-hehe/Desktop/sandboxes/pytorch
For how many days back do you want the information?
365
You have requested information for the last 365 days
Number of commits: 12786
Number of contributors: 1043
Do you want to output zero pairs? yes/no
no
Printing pairs of authors and the number of shared files:
Aarni Koskela -- Aiden Nibali : 1
Louis Feng -- shaoyf42 : 4
Rohan -- Scott Roy : 1
Will Constable -- Rodrigo Kumpera : 21
Kunal Bhatta -- nidefawl : 1
Xiaoya Xiang -- Leon Gao : 1
Aaron Orenstein -- Valentine233 : 22
Scott Wolchok -- Benson Ma : 8
Niklas Nolte -- milesial : 1
Fadi Botros -- SherlockNoMad : 1
Pian Pawakapan -- Michael Suo : 18
aaitzhan -- min-jean-cho : 2
rzou -- Richard Zou : 55
Matt Tinnel -- Kamil Dzieniszewski : 1
Fabrice Pont -- Alan Ji : 1
cyyever -- lkct : 2
Daniel Javady -- Lu Fang : 1
mashaobin -- SandishKumarHN : 1
The average 'similarity' of coupling: 8.360153256704981
The median 'similarity' of coupling: 1
Process finished with exit code 0
```

The second example