

Is Stack Overflow in Portuguese more attractive for Brazilian users?

Angela Lozano*, Bogdan Vasilescu[†], Miguel Botto Tobar^{‡§}, Wesley Torres[§], and Alexander Serebrenik[§],

*Vrije Universiteit Brussel, Brussels, Belgium

[†]Carnegie Mellon University, Pittsburgh, PA, USA

[‡]University of Guayaquil, Guayaquil, Ecuador

[§]Eindhoven University of Technology, Eindhoven, The Netherlands

Abstract—The abstract goes here.

I. INTRODUCTION

Software development and maintenance are activities that often involves many concepts and reference documents (citar). Many software aspects may be changed over time. In order to work with them and details involved in a software project, developers often need helps from one another. Nowadays a widely used way is for developers to ask questions and/or answer them in various online forums. StackOverflow (SO)¹ is a Q&A site with more than six million registered users. SO also has the version on Portuguese, Russian and Spanish.

Alexander ►What is the goal of your study?◄

Alexander ►Here you need to explain why did you decide to focus on Brazil as opposed to any other country in the world.◄

Alexander ►Separate the methodological discussion (what have I done? for example, downloaded the data, run a tool) from the results (what has it delivered, for instance, 5.954 user profiles). Make a separate section "Methodology" and a separate section "Results".◄

II. METHODOLOGY

A. Data Extraction

The data extraction has been performed on March 6, 2016, and included data from November 2013 to February 2016 from the Stack Exchange (SE) data dump² and for this study, we considered the Portuguese version (SOPT). The XML files corresponding to the tags, users, and posts were transferred to a MySQL database, through a R function per type of file (i.e., posts, users, and tags).

B. Data Preprocessing

We cleansed the tables eliminating few users were due to lack of data. None of these users had AccountId (i.e., user identifier for all stackExchange websites), LastAccessDate, WebsiteUrl, Location, UpVotes, DownVotes or Age. All of these users have the same display name (i.e., "a25bedc5-3d09-41b8-82fb-ea6c353d75ae"), and whenever they have a ProfileImageUrl, it is the same³. These accounts were created

at different times from November 2015 to February 2016. We could not come up with a plausible reason for these anonymous users having the same display name but no other data, they do not seem to have anything in common. In total 3 SOPT users have been eliminated.

We focused on Brazilian users thus to identify their location we used `countryNameManager`⁴.

C. Searching Process

Consequently, the locations were identified, and before starting the search a group of students were selected based on whether they had experience in use SO (in Spanish version) or GitHub, as detailed below:

The search started doing a manual inspection by each user profile based on `userId`, i.e. `http://pt.stackoverflow.com/users/1919/`, where **1919** is the `userId`. On the user profile, we looked the email address, if it was not available, we checked whether the user has a GitHub account or a personal web page.

- With Github account we found out the email address below `userName` if it was not on, we used a browser extension `gitDiscovered`⁵ to discover the email address or we checked his/her public activity looking for at he/she did a git command⁶, and then we searched the email address using a Github API⁷ by `userName`.
- With the personal web page, we searched the email address on the section "about me" or *sobre me* in the Portuguese language.
- If none of the above, we used GitHub and searched by `userName` from SOPT, and compared profile picture, location, skills between SOPT and search results on GitHub, and then we followed the steps above mentioned with GitHub account, in order to find the user and get the email address.

In order to ensure the accuracy of results, we selected a random group of 15 users each 500 profiles and searched using the steps above mentioned. Whether we found new email addresses or missing information, we chose another one random group and applied the manual inspection again.

¹<https://stackoverflow.com/>

²<https://archive.org/details/stackexchange>

³<https://www.gravatar.com/avatar/?s=128&d=identicon&r=PG&f=1>

⁴<https://github.com/tue-mdse/countryNameManager>

⁵<https://gitdiscovered.com/>

⁶<https://git-scm.com/docs/git-push>

⁷<https://api.github.com/users/username/events/public>

III. INTERVIEW

In order to understand how Brazilians use the Portuguese version of StackOverflow, we decided to conduct a semi-structured interview because **Wesley** ► *I WILL EXPLAIN THE REASON ABOUT WE CHOOSE THIS KIND OF INTERVIEW*◄ . We interviewed 4 Brazilians developers who work in different regions of Brazil. One of these developers never used the Portuguese version of StackOverflow, but we interviewed him just to get his point of view about the Portuguese version of StackOverflow. All of these interviews were conducted in Portuguese then they were translated to English. Both the Portuguese and English versions of the interviews can be downloaded in **Wesley** ► *Add the address*◄ .

Brazil is a big country and it might have different **Wesley** ► *I will add some word here that I dont know it yet =>*◄ for software development. To try to cover this diversity, we used the social media to call for developers who would like to participate of the interview. We select developers from Santa Catarina, São Paulo, Brasília and Pernambuco; south, southeast, center-west and northeast of Brazil, respectively.

Wesley ► *I will check the international names of these States*◄

We recorded the audio of our first interview, however during the transcription process we realized that this was not the best approach to follow because it took too long **Wesley** ► *I will think in another sentence/word*◄ , this process can take up to eight hours per hour of audio as described by Hove et al.[1]. Thus, we decided to conduct the other interviews using the Skype chat.

It is not the first time that a instant message tool was used to conduct an interview [2]. This approach has been discussed in the social sciences [3], [4] **Wesley** ► *I will improve this paragraph*◄

talk about barriers [5].. Igor Thesis is about " Supporting newcomers to overcome the barriers to contribute to open source software projects "

metodology by [6]

The guidelines; there is no right/wrong answer; better record the interview. [7]

In our first interview, we interviewed a developer from Sao Paulo (Southeast of Brazil). The interviewed was made in , Subject 1

Subject 2

Subject 3

Subject 4

summary of out finds: **Wesley** ► *I will remove the names...*◄

- All of them complained about the Portuguese content. They think the English version is more complete.
- Marcio are not interesting in help the community.
- Alex and Karina said that if they had an account they would help others. Alex have an account but he did not have it when he had the chance to help.
- They do not make the search on SO, first they use google, then the google send them to SO.
- Marcio thinks that on-line translation tool is good enough, so he can use the English version without any problems.

- Karina and Marcio always find a solution for their problems, so they never had to make any question on SO.
- Alex thinks that the Portuguese version will soon not be necessary anymore, because he thinks that English is essential for those who work in the IT.
- Giovanni and Marcio thinks that some people do not use the PT version because they dont know that there is a PT version. Karina thinks that some people do not use it, because of the poor PT content.
- Giovanni prefer be more active in the PT version because it is new and needs more help. And Alex prefer be more active in the EN version because (according to him) English is the official language for software development.

Text from Igor's Theses.. remove it

"All the interviews followed a semi-structured script and were conducted using textual based chat tools, like Google Talk. We chose this mean once the participants are used to this kind of tool for their professional and personal activities. The interviews were conducted following three different scripts, used according to the participant's profiles. The scripts were validated during the pilot interviews and by one specialist in qualitative studies, and one specialist in Open Source Software"

"We understand that the use of textual chat as the interview means can be considered a threat. The possibility of context change and the execution of parallel activities that distract the interviewees can be a negative aspect of using this mean. The use of Instant Messengers has been discussed in the social sciences (Opdenakker, 2006; Hinchcliffe and Gavin, 2009), and they point out that there is a set of positive effects of using these tools. In our case, we chose to use this means once the participants are used to the environment (they could choose the IM that they were more used to), and electronic means are the default (and preferred) way of communication in OSS projects."

Things to write:

- The type of interview, how it was conducted, why we did this way

The Criteria we used to select people - Brazilians from industry that use Stack Overflow in Portuguese.

- show the guideline - The guidelines; Make it clear that we said this to the interviewed: there is no right/wrong answer and if we could record the interview...

- write that all everyone that was interviewed agreed about use chat - skype -

IV. RESULTS

We identify the location of 7.264 users of SOPT which corresponds to 27% of its users. As we foresaw, most of the users of SOPT are located at Portuguese-speaking countries, in particular in Brazil (see Table I). Although there is a wide range of non-Portuguese speaking countries users, when looking at percentages these countries only represent 2

For half of the users whose location was identified, we could identify their gender (65%). Females are an overwhelming minority (4% SOPT users).

Country	Total
Brazil*	5954
Portugal*	599
United States	220
United Kingdom	80
Canada	44
France	25
Germany	42
India	30
The Netherlands	20
Mozambique*	14
Angola*	8
Cape Verde*	4
Other non Portuguese countries	224
None	19415

Table I

USER'S LOCATION IN SOPT. PORTUGUESE SPEAKING COUNTRIES ARE MARKED WITH AN ASTERISK.

REFERENCES

- [1] S. E. Hove and B. Anda, "Experiences from conducting semi-structured interviews in empirical software engineering research," in *Proceedings of the 11th IEEE International Software Metrics Symposium*, ser. METRICS '05. Washington, DC, USA: IEEE Computer Society, 2005, pp. 23–. [Online]. Available: <http://dx.doi.org/10.1109/METRICS.2005.24>
- [2] I. F. Steinmacher, "Supporting newcomers to overcome the barriers to contribute to open source software projects," Ph.D. dissertation, University of São Paulo, 2 2015.
- [3] V. Hinchcliffe and H. Gavin, "Social and Virtual Networks: Evaluating Synchronous Online Interviewing Using Instant Messenger," *The Qualitative Report*, vol. 14, no. 2, 2009. [Online]. Available: <http://www.nova.edu/ssss/QR/QR14-2/hinchcliffe.pdf>
- [4] R. Opdenakker, "Advantages and disadvantages of four interview techniques in qualitative research," *Forum Qualitative Sozialforschung / Forum: Qualitative Social Research*, vol. 7, no. 4, 2006. [Online]. Available: <http://www.qualitative-research.net/index.php/fqs/article/view/175>
- [5] P. J. G. C. P. Denae Ford, Justin Smith, "Paradise unplugged: Identifying barriers for female participation on stack overflow," *International Symposium on the Foundations of Software Engineering (FSE)*, 2016.
- [6] C. M. Gerpheide, R. R. Schiffelers, and A. Serebrenik, "Assessing and improving quality of qvto model transformations," *Software Quality Journal*, vol. 24, no. 3, pp. 797–834, Sep. 2016. [Online]. Available: <http://dx.doi.org/10.1007/s11219-015-9280-8>
- [7] C. B. Seaman, "Qualitative methods in empirical studies of software engineering," *IEEE Trans. Softw. Eng.*, vol. 25, no. 4, pp. 557–572, Jul. 1999. [Online]. Available: <http://dx.doi.org/10.1109/32.799955>

ACKNOWLEDGMENT

The authors would like to thank...