

# BlogSet-BR: A Brazilian Portuguese Blog Corpus

<sup>1</sup>Henrique D.P. dos Santos, <sup>2</sup>Vinicius Woloszyn, <sup>1</sup>Renata Vieira

<sup>1</sup>Pontificia Universidade Catolica do Rio Grande do Sul  
Av. Ipiranga, 6681, Building 32 - Porto Alegre - Brazil

<sup>2</sup>Instituto de Informática – Universidade Federal do Rio Grande do Sul  
Caixa Postal 15.064 – 91.501-970 – Porto Alegre – RS – Brazil

henrique.santos.003@acad.pucrs.br, vwoloszyn@inf.ufrgs.br, renata.vieira@pucrs.br

## Abstract

The rich user-generated content found on internet blogs have always attracted the interest of scientific communities for many different purposes, such as from opinion and sentiment mining, information extraction or topic discovery. Nonetheless, an extensive corpora is essential to perform most of Natural Language Processing involved in these tasks. This paper presents BlogSet-BR, an extensive Brazilian Portuguese corpus containing 2.1 billions words extracted from 7.4 millions posts over 808 thousand different Brazilian blogs. Additionally, a survey was conducted with authors to draw a profile of Brazilian bloggers.

**Keywords:** Blog as Corpus, BlogSet-BR, Brazilian Portuguese Corpus

## 1. Introduction

Several efforts have been made to build a large corpora based on user-generated content, since they are crucial for many different Natural Language Processing tasks, such as opinion mining, sentiment analysis, topic detection and Age/Gender detection (Burton et al., 2009; Buck et al., 2014; Santos et al., 2017). From all content available on the internet, blogs have been often employed as a main source of user-generated content (Agarwal and Liu, 2008; Agarwal and Liu, 2009; Santos et al., 2012). Nonetheless, there is still a lack of a large semi-structured corpus that also contains author profiles in Brazilian Portuguese. To illustrate this situation, Table 1<sup>1</sup> gives an overview of the initiatives to create a Portuguese corpora based on general content.

Besides Linguateca resources in Table 1, we added other four corpora: the brWac corpus for Brazilian Portuguese built by downloading text from the web (Boos et al., 2014), Buscapé, a Portuguese product reviews corpus extracted from a collaborative review site (Hartmann et al., 2014), the Portuguese Wikipedia dump and BlogSet-BR corpus. The first three corpora listed in the table contains the authorship of the content.

This paper describes BlogSet-BR, a semi-structured large corpus containing author information extracted from Brazilian Portuguese blogs. It contains more than 7.4 millions posts resulting in 2.1 billion words. It is the first Brazilian Portuguese corpus about blogs. In this paper, we also conducted a survey with the authors to create a profile of Brazilian bloggers. The main contributions of this work are:

- The first semi-structured large corpus of posts from Brazilian Portuguese blogs;
- A corpus with authorship, date-time and label information attached to the text;

- Information retrieved from the author profiles, such as gender, age, and educational level.

This work describes the building process of the BlogSet-BR collection and also descriptive statistics about the data. The rest of this paper is organized as follow: first, Section 2. presents related work; in Section 3., we discuss the three phases of creating the BlogSet-BR collection; section 4. provides an overview of the statistics and content on the collected data; section 5. describes the profile of Brazilian bloggers obtained through a survey conducted with authors; section 6. presents the conclusions and suggestions for further work.

## 2. Background

Extract content from the web is a constant effort by academic and industry researchers. Such data sets allow the accomplishment of many different tasks, such as relevance ranking for online documents and several other machine learning tasks (Woloszyn et al., 2016; Woloszyn et al., 2017). For instance, the Common Crawl project maintains an open repository of web crawl data that can be accessed and analyzed by any research group<sup>2</sup>. This corpus has been used to build language models (Roziński and Stokowiec, 2016), and to analyze word frequencies (Buck et al., 2014). However, this project does not contain structured information to allow a social network analyses about the authors. This gap is filled by datasets built with blog content, enabling analyses of web data together with user relation and temporal characteristics. The TREC Conference organizers built a blog corpus for research purposes (Macdonald and Ounis, 2006) making 100,649 blogs in the English language available for TREC shared tasks.

The ICWSM 2009 Spinn3r dataset (Burton et al., 2009) is another example of corpora extracted from blogs. It has more than 44 million blog posts for the English language

<sup>1</sup>Adapted from <http://www.linguateca.pt/ACDC/>

<sup>2</sup><http://commoncrawl.org/>



the distribution of the dates of the collected posts after the year 2003. The posts before 2003 represent only 0.005% of the corpus.

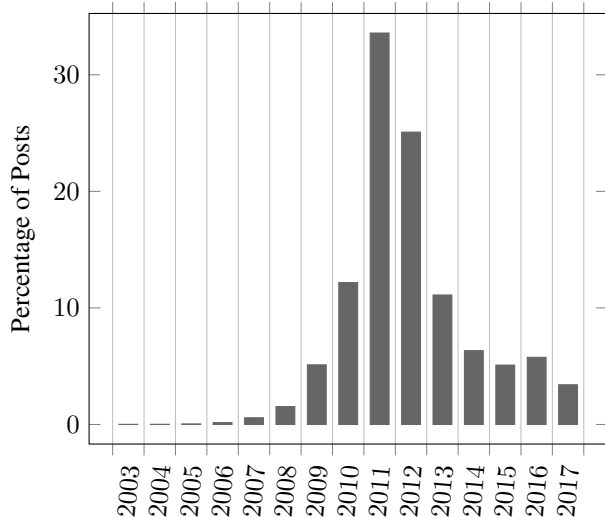


Figure 2: Posts Year Distribution

The posts start decreasing in Blogspot platform after 2011. This behavior could represent the migration of users to Twitter and Facebook, both grow exponentially in Brazil in the same period (Yokoyama and Sekiguchi, 2014).

### 5. Bloggers Profile

A survey<sup>4</sup> was conducted with 4,332 Brazilian authors from the Blogspot platform to establish a profile of bloggers who use the platform. The requested pieces of information were:

- Age
- Gender
- Topic of interest
- Pageviews
- Update frequency
- Educational level

The majority of authors are male (61%). In addition, almost every user finished high school (92%) and most of them add new a new post at least once a week (85%). Figure 3 shows the age distribution on the collected survey.

Regarding frequency, 32% create new content every day and 53% publish every week. Most users create their own content (86%), giving their own opinion about mundane facts.

Table 2 shows the distribution of topics, which vary: books, fashion, technology, education, politics, movies, and arts. In this question, users were allowed to select multiple choices.

In Figure 4 we show an overview of answers by Brazilian bloggers regarding their educational level. This was not a mandatory question, but 97% of all those who answered

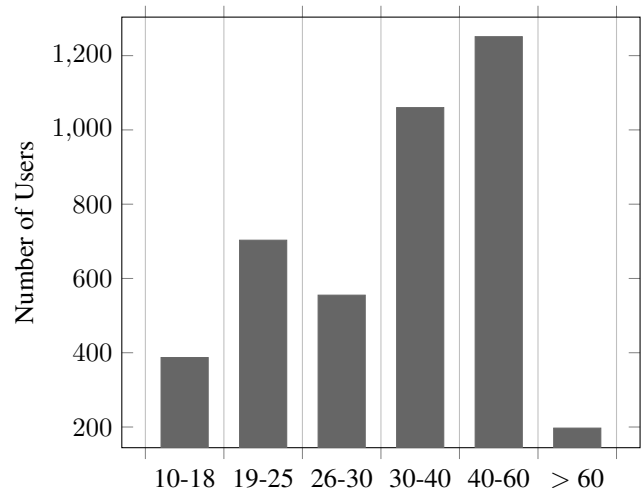


Figure 3: User's Age Distribution

Topic	% of Users
Arts	29%
Education	28%
Literature / Books	28%
Music	27%
Politics	24%
Philosophy	19%
Friendship	17%
Movies	17%
Health	17%
Technology	16%

Table 2: Distribution of the Top 10 Topics selected by users

finished high school and the majority of them are at the university.

Considering social media, most interviewees use Facebook (94%) and Twitter (67%). This rich information about age, gender, and educational level of Brazilian bloggers could be employed in different tasks, such as writing style de-

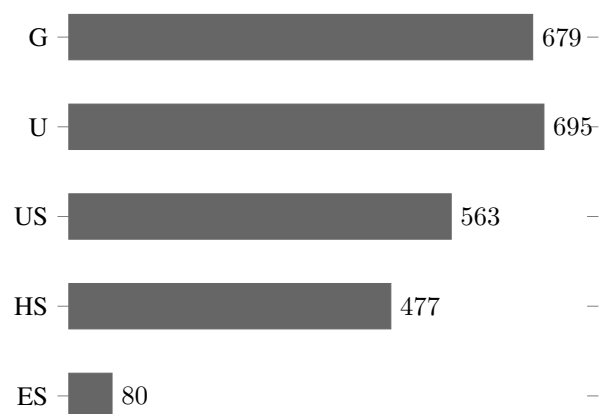


Figure 4: User's Educational Level: G (Graduated), U(Undergraduated), US (Undergraduated Student), HS (High School), ES (Elementary School)

<sup>4</sup>Survey conducted in 2012. Not every author in the survey is in the corpus BlogSet-BR, collected in 2017.

tection and readability assessment (Herring and Paolillo, 2006; Argamon et al., 2007).

## 6. Conclusion and Further Work

In this paper, we described the motivations, details, and building process of BlogSet-BR. Additionally, we conducted a survey with authors to create their profile, which resulted in a rich description of Brazilian bloggers. This corpus is the biggest corpus with authorship information for the Brazilian Portuguese language. These types of dataset are useful for topic detection and other NLP tasks. All the content of this work is available on the web<sup>5</sup>. It is possible to download the raw data in the JSON format, only the Brazilian blog posts in the CSV format, and the survey with 4 thousand users in the XLS format.

## Acknowledgments

This work was partially supported by the CAPES (Coordenação de Aperfeiçoamento de Pessoal de Nível Superior) Foundation (Brazil), PUCRS (Pontifícia Universidade Católica do Rio Grande do Sul), and UFRGS (Universidade Federal do Rio Grande do Sul).

## 7. Bibliographical References

- Adar, E. and Adamic, L. A. (2005). Tracking information epidemics in blogspace. In *Proceedings of the 2005 IEEE/WIC/ACM international conference on web intelligence*, pages 207–214. IEEE Computer Society.
- Agarwal, N. and Liu, H. (2008). Blogosphere: research issues, tools, and applications. *ACM SIGKDD Explorations Newsletter*, 10(1):18–31.
- Agarwal, N. and Liu, H. (2009). Modeling and data mining in blogosphere. *Synthesis lectures on data mining and knowledge discovery*, 1(1):1–109.
- Argamon, S., Koppel, M., Pennebaker, J. W., and Schler, J. (2007). Mining the blogosphere: Age, gender and the varieties of self-expression. *First Monday*, 12(9).
- Buck, C., Heafield, K., and Van Ooyen, B. (2014). N-gram counts and language models from the common crawl. In *LREC*, volume 2, page 4.
- Burton, K., Java, A., and Soboroff, I. (2009). The icwsm 2009 spinn3r dataset. In *Third Annual Conference on Weblogs and Social Media (ICWSM 2009)*. AAAI.
- Herring, S. C. and Paolillo, J. C. (2006). Gender and genre variation in weblogs. *Journal of Sociolinguistics*, 10(4):439–459.
- Lerman, K. and Ghosh, R. (2010). Information contagion: An empirical study of the spread of news on digg and twitter social networks. *ICWSM*, 10:90–97.
- Macdonald, C. and Ounis, I. (2006). The trec blogs06 collection: Creating and analysing a blog test collection. *Department of Computer Science, University of Glasgow Tech Report TR-2006-224*, 1:3–1.
- Mishne, G. et al. (2005). Experiments with mood classification in blog posts. In *Proceedings of ACM SIGIR 2005 workshop on stylistic analysis of text for information access*, volume 19, pages 321–327. Citeseer.

- Quan, C. and Ren, F. (2009). Construction of a blog emotion corpus for chinese emotional expression analysis. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 3-Volume 3*, pages 1446–1454. Association for Computational Linguistics.
- Roziewski, S. and Stokowiec, W. (2016). Language-crawl: A generic tool for building language models upon common-crawl. In *LREC*.
- Santos, R. L., Macdonald, C., McCreadie, R., Ounis, I., Soboroff, I., et al. (2012). Information retrieval on the blogosphere. *Foundations and Trends® in Information Retrieval*, 6(1):1–125.
- Santos, H. D. P., Woloszyn, V., and Vieira, R. (2017). Portuguese personal story analysis and detection in blogs. In *Proceedings of the International Conference on Web Intelligence*, pages 709–715. ACM.
- Santos, H. D. P. (2013). Identificação de autoridades em tópicos na blogosfera brasileira usando comentários como relacionamento. Master’s thesis, Federal University of Rio Grande do Sul.
- Woloszyn, V., Santos, H. D. P., and Wives, L. K. (2016). The influence of readability aspects on the user’s perception of helpfulness of online reviews. *Revista de Sistemas de Informação da FSMA*.
- Woloszyn, V., Santos, H. D. P., Wives, L. K., and Becker, K. (2017). Mrr: an unsupervised algorithm to rank reviews by relevance. In *Proceedings of the International Conference on Web Intelligence*, pages 877–883. ACM.
- Yokoyama, M. H. and Sekiguchi, T. (2014). The use of social network sites in the workplace: A case study in brazilian companies. *Brazilian Business Review*, 11(2):87.

## 8. Language Resource References

- Boos, R., Prestes, K., Villavicencio, A., and Padró, M. (2014). brwac: a wacky corpus for brazilian portuguese. In *International Conference on Computational Processing of the Portuguese Language*, pages 201–206. Springer.
- Hartmann, N., Avanço, L., Balage Filho, P. P., Duran, M. S., Nunes, M. d. G. V., Pardo, T. A. S., Aluísio, S. M., et al. (2014). A large corpus of product reviews in portuguese: Tackling out-of-vocabulary words. In *LREC*, pages 3865–3871.

<sup>5</sup><http://www.inf.pucrs.br/linatural/blogset-br>