# Mental lexicon for personality identification in texts

**Anonymous Author(s)**
Affiliation
Address
email

## 1 Introduction

Personality identification from texts is a relative new area of interest in the natural language processing (NLP) community. The benefits of helping to identify the personality of a subject solely on the text they write are manifolds. For one, it can help directly to the authors of such texts to understand their social interactions, and their behaviour in general [5, 12]. Beyond that, personality identification is useful for many other research areas. For instance, in human computer interactions (HCI), interactive systems may be able to adapt to user's personality, providing a better experience [2]. In education, building intelligent tutors compatible with the student's personality can improve, not only the experience of the student with the system, but also the system could provide more adequate material from a educative program in accordance to the particular student's preferences [10, 6].

From the NLP perspective, personality identification from texts can be treated as an author profiling problem. Author profiling consists on, given a text, determine some demographics characteristics of the author of such text. In this context, the representation of a given text such that the model can extract relevant information according to the specific demographic interest [7, 1] is of a relevant importance.

In the mental health context, the main interest is not only to build accurate systems, but to provide interpretable results that in turn, would serve as additional and reliable elements to a therapist. Accordingly, we focused on developing an automatic method for personality identification, able to provide valuable information regarding the language usage of subjects being analyzed.

Specifically, we use the linguistic theory behind lexical availability to first compute a set of relevant mental lexicon from groups of subjects (e.g. *introverts* vs *extroverts* for the Extroversion trait) and then we use this mental lexicon in a representation stage. For our experiments, we use two data sets: English essays and Spanish essays; these datasets use the Big Five Model of Personality [9].

## 2 Lexical Availability as language descriptor

Lexical availability methods were developed to provide useful vocabulary to immigrants in early 60's in France [13]; where word's frequencies do not necessary means importance of such a word in a given context. Traditionally, the lexical availability elicitation approach consists on ask to a group of subjects to write, in a small period of time (usually 2 to 5 minutes), a set of terms given a specific center of interest [4, 13].

We use a linguistically motivated approach aiming to identify those lexical markers that represent the words springing to mind in response to a specific topic. Lexical Availability score (LA) measures the ease with which word is generated in a given communicative situation [4], and allows to obtain the *mental lexicon* which represents the vocabulary flow usable of a group of people [3].

In general, the terms with greater LA score can be seen as the most important ones for a group of people with the same personality trait. Thus, we computed the mental lexicon for each pole in a trait,
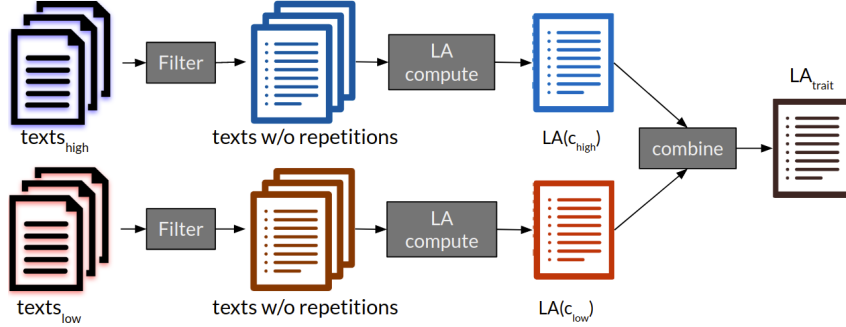
Figure 1: Schema to generate a mental lexicon given a set of written texts.

Table 1: Results with the best configuration from our proposed method and traditional baseline. In bold are mark results of our method when outperform the baseline.

| Trait | RxPI Spanish[11] | | English essays[8] | |
| --- | --- | --- | --- | --- |
| | F-macro (Ours) | F-macro (Baseline) | F-macro (Ours) | F-macro (Baseline) |
| EXT | **0.6018** | 0.5640 | 0.5753 | 0.5788 |
| AGR | 0.5697 | 0.5711 | **0.5615** | 0.5530 |
| CON | **0.5857** | .5800 | 0.5795 | 0.5806 |
| STA | **0.6026** | 0.5828 | **0.5918** | 0.5785 |
| OPE | 0.5704 | 0.5722 | **0.6414** | 0.6237 |

36  and then a general list ($LA_{trait}$) was generated to be use in a vectorial representation model with
37  dimension equal to $|LA_{trait}|$.

## 3   Proposed framework and evaluation

39  We proposed the method in Figure 1 to use lexical availability for texts representation. Our method
40  has three main processes: The *filter process* generates a list of terms without repetitions given any
41  instance text. The *LA compute process* computes the lexical availability score of a list of terms as
42  $LA(t_j) = \sum_{i=1}^{n} e^{(-2.3*\frac{i-1}{n-1})} * \frac{f_{ij}}{I}$, where $t_j$ is the term $j$ in a list; $n$ is the lowest position of a term
43  $j$ in some list; $i$ is the position of term $j$ in a list; $f_{ij}$ is the number of lists in where term $j$ appears
44  in position $i$, and $I$ is the total number of lists. Finally, the *combine process*, takes as input the lists
45  generated for each class and using set operations combine them into a single general list that we
46  called $LA_{trait}$.

47  Once we have the mental lexicon of a trait (a.k.a. $LA_{trait}$), we use the scores and terms in this
48  list to generate a vector representation of a given instance text. In order to weight each term in
49  our vector, we use three approaches as follows. If $w_k$ is the weight of a term $k$ and $LA(w_k)$
50  is the score of lexical aviability of word $k$ in the list $LA$ then: 1) $w_k^{global} = LA_{trait}(w_k)$, 2)
51  $w_k^{comb} = LA_{trait}(w_k) * LA_{instance}(w_k)$ where $LA_{instance}$ is the score of a term in the unseen
52  instance, and 3) $w_k^{tfla} = tf * LA_{trait}(w_k)$, where $tf$ is the frequency of the term ($w_k$) in the unseen
53  instance.

54  To compare our performance in classification, we used three representation baselines: n-grams
55  of words and characters, and a dictionary based representations such as LIWC. For each of these
56  baselines we experimented with several configuration parameters (e.g. the number of $n$). To train a
57  model we used traditional learning algorithms such as probabilistic, decision trees, support vector
58  machine, and instance based. Table 1 shows the results with the best parameters for our method as
59  well as for the baselines.

60  Our ongoing work in this project is to analyze the semantic categories in each lists that are relevant
61  to the expert when identify the personality of a subject. At the same time we want to use more
62  sophisticated methods that take advantage of our proposed representation to improve the classification
63  performance.

# References

[1] S. Argamon, M. Koppel, J. Fine, and A. R. Shimoni. Gender, genre, and writing style in formal written texts. *TEXT*, 23:321–346, 2003.

[2] T. W. Bickmore and R. W. Picard. Establishing and maintaining long-term human-computer relationships. *ACM Trans. Comput.-Hum. Interact.*, 12(2):293–327, June 2005.

[3] R. M. J. Catalán. *Lexical availability in English and Spanish as a second language*, volume 17. Springer, 2013.

[4] N. R. Dimitrijević. A comparative study of the lexical availability of monolingual and bilingual schoolchildren. 1981.

[5] K. A. Holder. M. D. Temperament and happiness in children. *Journal of Happiness Studies*, 2009.

[6] M. Komarraju and S. J. Karau. The relationship between the big five personality traits and academic motivation. *Personality and Individual Differences*, 39(3):557 – 567, 2005.

[7] M. Koppel, S. Argamon, and A. R. Shimoni. Automatically categorizing written texts by author gender. *Literary and Linguistic Computing*, 17(4):401–412, 2002.

[8] F. Mairesse, M. A. Walker, M. R. Mehl, and R. K. Moore. Using linguistic cues for the automatic recognition of personality in conversation and text. *Journal of Artificial Intelligence Research (JAIR*, pages 457–500, 2007.

[9] R. R. McCrae and P. T. Costa Jr. Personality trait structure as a human universal. *American psychologist*, 52(5):509, 1997.

[10] M. Pavalache-Ilie and S. Cocorada. Interactions of students' personality in the online learning environment. *Procedia - Social and Behavioral Sciences*, 128:117 – 122, 2014.

[11] G. Ramírez-de-la Rosa, E. Villatoro-Tello, and H. Jiménez-Salazar. Txpi-u: A resource for personality identification of undergraduates. *Journal of Intelligent & Fuzzy Systems*, 34(5):2991–3001, May 2018.

[12] G. B. Svendsen, J.-A. K. Johnsen, L. Almås-Sørensen, and J. Vittersø. Personality and technology acceptance: the influence of personality factors on the core constructs of the technology acceptance model. *Behaviour & Information Technology*, 32(4):323–334, 2013.

[13] A. Ávila Muñoz and J. Sanchez Saez. Fuzzy sets and prototype theory: Representational model of cognitive community structures based on lexical availability trials. *Review of Cognitive Linguistics*, 12, 01 2014.