

A Single Author Style Representation for the Author Verification Task

Notebook for PAN at CLEF 2014

Cristhian Mayor^{1,3}, Josue Gutierrez³, Angel Toledo³, Rodrigo Martinez³, Paola Ledesma^{2,3}, Gibran Fuentes³, and Ivan Meza³

¹Institut National des Sciences Appliquees de Lyon (INSA Lyon)
<http://www.insa-lyon.fr/>

²Escuela Nacional de Antropologia e Historia (ENAH)
<http://www.enah.edu.mx>

³Instituto de Investigaciones en Matematicas Aplicadas y en Sistemas (IIMAS)
Universidad Nacional Autonoma de Mexico (UNAM)

Abstract This paper presents our experience implementing three approaches for the ‘PAN 2014 Author Identification’ [3,1] task using the same representation for the author’s style. Two of our approaches extend previous successful approaches: naive Bayes [4] and impostor [8] methods. The third approach is based on original research on sparse representation for text documents. We present results with the official development and test corpora.

1 Introduction

Author verification has multiple applications on several areas including information retrieval and computational linguistics, and has an impact in fields such as law and journalism [2,5,9]. In this edition of the *PAN 2014 Author Identification*, the task was formally defined as follows¹:

Given a small set (no more than 5, possibly as few as one) of “known” documents by a single person and a “questioned” document, the task is to determine whether the questioned document was written by the same person who wrote the known document set.

This year the documents were in four languages and four genres with the following combinations: essays and reviews in Dutch, essays and novels in English and articles in Greek and Spanish.

In this work we present the implementation of three approaches to perform the authorship verification task based on the same document representation. In particular, we repeatedly use a vector space representation of documents as described in our approach last edition [6], but we agglomerate them to obtain a single author’s style representation. We applied two well-known methods for verifying the author: naive Bayes [4] and impostor [8]. Additionally, we implemented a novel approach using sparse representation [10].

¹ As described in the official website of the competition <http://pan.webis.de/> (2014).

2 Author's Style representation

The representation for an author's style is generated in two stages. First, we represent the documents from an author using the vector space model [7]:

$$d = (w_1, w_2, \dots, w_m)$$

where w_i is a frequency or weight of a *word* of the vocabulary of size m for the word i . A common representation of a vector space model is the bag of words model, in which the words represent actual words of the document and the frequency count occurrences of such words in the represented document.

Additionally to the bag of words we use the following feature frequencies to represent the documents:

Bag of words Frequencies of words in the document.

Bigram Frequencies of two consecutive words.

Trigram Frequencies of three consecutive words.

Prefix Frequencies of prefixes of words.

Suffix Frequencies of suffixes words.

Prefix bigram Frequencies of two consecutive prefixes of words.

Suffix bigram Frequencies of two consecutive suffixes words.

Stop words Frequencies of stop words.

Stop words bigram Frequencies of two consecutive stop words.

Punctuation Frequencies of punctuations.

Words per sentence Frequencies of words per sentence.

In the second stage an author is represented by the sum of the vectors representing the documents written by her or him. This cumulative vector is normalized by the number of documents by the author and the resulting normalized vector represents the style of the author. This representation is not novel, however many approaches on author verification optimize the representation on the domain, in our experiments we keep the same representation independently of the domain. Some of the implemented approaches require an instance of a document by a certain author, to accomplish this it we sample a document from the author's style representation.

3 Approaches

The approaches implemented in the task are described next.

3.1 Impostors

This method consists on iteratively compare the vector distance between the author's document to the questioned document versus the distances between several impostor documents to the questioned document. With these distances a score is built up based on how many times the author and questioned documents are closer than the impostor and questioned documents. For this approach we follow the description of the method

by *Seidman* (2013) [8]. We modify the method to work on more than one set of features and instead to use impostors from the web we used the training corpus as source of impostors. Additionally, we extended the approach to produce a probability as output based on repetition of the algorithm since the document instances were randomly sampled.

3.2 Naive Bayes

This method consists of sampling from the author and the impostor style representations two document instances for each. A probability score is then calculated using the common term between the questioned document and the author's documents. On the other hand, an alternative score is calculated between the questioned document and the impostor documents. These scores are derived using Bayes². The purpose of the score is to capture the probability that the document was created by the same author, if the score for the author is higher than the impostor, we consider it as evidence of authorship. We iterate n times over this method to calculate the probability of authorship.

3.3 Sparse

This methodology has been successfully applied to the face recognition task, in which the identity of a face image has to be determined from a set of known faces [10]. We adapted this methodology to the authorship verification task. The method consists on identifying the components that contribute to the questioned document from samples of documents from a set of authors. The rationale is that the biggest contribution of components should be elements from a single author. In order to identify the components the method proposes the following l^1 -minimization:

$$\begin{aligned} \text{minimize} \quad & x = \operatorname{argmin} \|x\|_1 \\ \text{subject to} \quad & Ax = y \end{aligned} \tag{1}$$

Where y is the questioned document, A is the matrix of n samples from different m candidate authors (impostors), and x is the variable to minimize which represent the contribution from each candidate. So that multiplying the samples by the contribution we could generate the questioned document. From the resulting variable x_0 we can quantify the residuals given by Ax versus y and decide which author contributes with more components. We adapt this method to produce a probability as result by iterating k times over the full method.

4 Results

Table 4 presents the overall development results obtained with each approach. The results we obtained using the same document representation presented in section 2 and the same parameters for all languages and genres.

The results for each language and genre for our best system (i.e., sparse) are presented in Table 4.

² This method is inspired in the following document: <http://cs229.stanford.edu/proj2009/Leahy.pdf>.

Table 1. Overall results for all approaches.

Approach	AUC	C@1	Score
<i>impostor</i>	62%	56%	35%
<i>n-gram</i>	64%	57%	36%
<i>Sparse</i>	72%	68%	48%

Table 2. Detailed final scores for language and genre for the *sparse* approach on the testing corpus.

Approach	AUC	C@1	Score
Dutch reviews	93%	88%	82%
Dutch essays	57%	52%	30%
English essays	57%	56%	32%
English novels	66%	61%	41%
Greeks articles	82%	75%	62%
Spanish news	75%	71%	54%
<i>Overall</i>	70%	65%	50%

5 Discussion

In the preparation of the authorship verification systems we implemented three approaches: *impostor*, *n-gram* and *sparse* methodologies. During development we tested all of them on the same representation for the author’s style and the same parameter for the four languages and three genres of the task. During development and testing the best results were achieved using the *sparse* methodology which is interesting to us since it is the first time such method is applied to the task of authorship verification.

References

1. Gollub, T., Potthast, M., Beyer, A., Busse, M., Rangel, F., Rosso, P., Stamatatos, E., Stein, B.: irecent trends in digital text forensics and its evaluation. in pamela forner. In: Information Access Evaluation meets Multilinguality, Multimodality, and Visualization. 4th International Conference of the CLEF Initiative (CLEF 13) (2013)
2. Juola, P.: Authorship attribution. Found. Trends Inf. Retr. 1(3), 233–334 (Dec 2006)
3. Juola, P., Stamatatos, E.: Overview of the author identification task at pan 2013. In: CLEF 2013 Evaluation Labs and Workshop - Online Working Notes (2013)
4. Kešelj, V., Peng, F., Cercone, N., Thomas, C.: N-gram-based author profiles for authorship attribution. In: Proceedings of the conference pacific association for computational linguistics, PACLING. vol. 3, pp. 255–264 (2003)
5. Koppel, M., Schler, J., Argamon, S.: Computational methods in authorship attribution. Journal of the American Society for Information Science and Technology 60(1), 9–26 (2009)
6. Ledesma, P., Fuentes, G., Jasso, G., Toledo, A., Meza, I.: Distance learning for author verification. In: Proceedings of the conference pacific association for computational linguistics, PACLING. vol. 3, pp. 255–264 (2003)

7. Salton, G., Wong, A., Yang, C.S.: A vector space model for automatic indexing. *Commun. ACM* 18(11), 613–620 (Nov 1975)
8. Seidman, S.: Authorship verification using the impostors method. In: *CLEF 2013 Evaluation Labs and Workshop - Online Working Notes* (2013)
9. Stamatatos, E.: A survey of modern authorship attribution methods. *Journal of the American Society for Information Science and Technology* 60(3), 538–556 (2009)
10. Wright, J., Yang, A.Y., Ganesh, A., Sastry, S.S., Ma, Y.: Robust face recognition via sparse representation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 31(2), 210–227 (2009)