

Semantic-based Features for Author Profiling Identification: First insights

Delia-Irazú Hernández¹, Rafael Guzmán-Cabrera², Antonio Reyes³
and Martha-Alicia Rocha^{1,4}

¹Universidad Politécnica de Valencia, España

²Universidad de Guanajuato, México

³Instituto Superior de Intérpretes y Traductores, México

⁴Instituto Tecnológico de León, México

INTRODUCTION

Nowadays, social interaction through Internet is becoming a major problem due to the insufficient control regarding the authenticity of user profiles. Author profiling is the task of identifying personal characteristics of Internet users (such as age, gender, native language) based on analysing their interactions, mainly, considering their textual patterns. Author profiling has various applications such as security, forensics, marketing, among others.

FEATURES DESCRIPTION

- **Signature** \rightarrow explicit linguistic markers within a text.
- **Chatslang** \rightarrow words and expressions commonly used in internet forums.
- **Context** \rightarrow the presence of discriminating clusters across the classes.
- **Emotionality** \rightarrow the use of words to communicate emotions, feelings, moods, etc.
- **Semantic similarity** \rightarrow the semantic relatedness among the words of a text.
- In addition a list of **Bag of Words** was used. Finally, the **Jaccard similarity** coefficient was applied in order to focus on informative words rather than only on frequent ones.

EXPERIMENTS

- Each conversation is represented as a numerical vector in which each entry represent a feature.
- We make different combinations of the features proposed and we classified the conversations using various learning algorithms.

Experiment	Description
<i>SBF</i>	Semantic similarity + Signatures + ChatSlang + Emotionality
<i>SBF_M</i>	Semantic similarity + Emotionality
<i>SBF+BOW</i>	<i>SBF</i> + Bag of Words
<i>SBF+Jaccard</i>	<i>SBF</i> + Jaccard similarity coefficient
<i>SBF+Jaccard+BOW</i>	<i>SBF+BOW</i> + Jaccard similarity coefficient
<i>Jaccard</i>	Jaccard similarity coefficient
<i>Jaccard+BOW</i>	Jaccard similarity coefficient + Bag of Words
<i>Jaccard+Context</i>	Jaccard similarity coefficient + Context

Table: Features Combination

RESULTS

The learning algorithms applied to this combinations were Naive Bayes (NB), Support Vector Machines (SVM), Multilayer Perceptron (MP), Decision Tree (J48), and a bagging of classifiers (NB + SVM +J48). The table below introduces the results obtained from different experiments.

Experiments	NB	SVM	MP	J48	Bagging	Average
<i>SBF</i>	16.66	19.66	18.3	17.66	18.67	17.99
<i>SBF+BOW</i>	22	15.66	-	21.66	20.67	19.77
<i>SBF+Jaccard</i>	19	20.60	18.33	15.33	20	18.31
<i>SBF+Jaccard+BOW</i>	23	15.66	-	22	21	20.22
<i>Jaccard</i>	22.33	17.33	22.33	14.67	17.66	18.11
<i>Jaccard+BOW</i>	23	15.33	-	19.33	21.33	19.22
<i>Jaccard+Context</i>	17.66	21.6	-	17.33	21.33	18.86

Table: Results of Author Profiling classifiers

PAN RESULTS

- We defined a final model which was integrated with the features: ***SBF_M*+Jaccard+BOW**. According to our best results, the NB classifier was used to participate in the author profiling task in PAN 2013 competition. The results are shown in the next table.

Task	Accuracy		
	total	Gender	Age
English	0.2816	0.5671	0.5061
Spanish	0.1757	0.4982	0.3554

Table: Author Profiling Evaluation PAN 2013

CONCLUTIONS

- This research is based on semantic features.
- After analyzing the results, we could realize that the author profiling task has a high lever of overlap between classes; hence, the dificulty of correctly identifying the classes increases subtientially.
- The future work consists of developing an algorithm for principal components analysis (PCA) in order to obtain highly discriminating features.