

Author Profiling using Complementary Second Order Attributes and Stylometric Features



Konstantinos Bougiatiotis and Anastasia Krithara {bogas.ko, akrithara}@iit.demokritos.gr

National Center for Scientific Research "Demokritos", Athens, Greece

Introduction

Overview

- Author Profiling: Predict author traits based solely on text
- Novelty: PAN'16 features cross-genre evaluation(train on Twitter texts and test on other genres)

Data

Users: 1070 | Tweets: 562812

Second Order Attributes(SOA)

1. Calculate word-profile vectors \rightarrow Find **descriptive terms per class**, exploiting the per-class frequency of the words

$$t_{i,j} = \sum_{k:d_k \notin P_j} \log(1 + \frac{tf_{i,k}}{\operatorname{len}(d_k)} * w_k)$$

2. Map documents in profile space, using the word-profile vectors, from step 1, of the containing terms for each document

- **(**) Tasks: Age and Gender
- Languages: English, Spanish and Dutch(gender only)

$d_{k,j} = \sum_{i:t_i \in d_k} \frac{tf_{i,k}}{len(d_k)} \times \vec{t_i}$

System Workflow



Sample representation of the SOA method



Weighted SOAComplementary

- **?** Use the **complementary classes** for each word-class relation \rightarrow More even amount of data for each class \rightarrow **Robust estimates** and lesser bias
- **?** Weights inversely proportional to class frequency \rightarrow Terms re-

^aLópez-Monroy et al.: INAOE's participation at PAN'13: Author Profiling task-Notebook for PAN at CLEF 2013. In: CLEF 2013 Evaluation Labs and Workshop lated with rare profiles, aggregate more weight $\rightarrow Prior \ knowledge$ will help sparsely populated classes

Train Results(4-fold CV)

Models	English		Spanish		Dutch
	Age	Gender	Age	Gender	Gender
N-grams(PAN'15)	47.0	74.8	49.6	68.8	76.8
SOA	47.5	76.2	54.0	72.8	76.0
SOAC	49.1	76.8	50.4	71.6	76.8
W-SOAC	49.1	76.8	50.4	72.8	76.8
N-grams + W-SOAC	50.0	$\boldsymbol{77.5}$	52.0	73.2	78.1

Stylometric-Structural Features(PAN15')



Test Results						
Dataset	Language	Subtask	Accuracy%			
Test-1	Dutch	Gender	44.00			
	Fnelich	Age	30.46			
		Gender	53.45			
	Spanich	Age	29.69			
	spanisn -	Gender	60.94			
Test-2	Dutch	Gender	41.60			
	Fngligh	Age	55.13			
		Gender	69.23			
	Snanich	Age	32.14			
	- spanisn	Gender	67.86			

Number of	Number of	Number of	Tf-idf of	Bag of	Ngram	Word length	Number of	
Hashtags	Links	Mentions	Ngrams	Smileys	Graphs		Uppercase	

Finally, selected top **3000** frequent **3-grams** of chars(age) and **unigrams** of chars(gender)





Grant Agreements No. FP7-610928

Conclusions

✓ Stylometry and Discriminative features both capture gender information well enough. Also boosted performance through fusion
✓ Age considerably more difficult than gender to predict, across all languages and regardless of the methodology
✓ Different performance in the two test datasets, highlight the added

difficulty of the cross-genre task