

On the Multilingual and Genre Robustness of EmoGraphs for Author Profiling in Social Media

Francisco Rangel and Paolo Rosso



UNIVERSITAT
POLITÈCNICA
DE VALÈNCIA

autoritas®

7 Features

Given a graph and its adjacency matrix $A = a_{i,j}$, where $a_{i,j} = 1$ if node i is linked to a node j and 0 otherwise:	
$x_n = \lambda \sum_{i \in M(n)} x_i = \lambda \sum_{i \in G} a_{n,i} x_i$ where λ is a constant representing the greatest eigenvalue associated with the centrality measure.	
It is the ratio of all shortest paths from one node to another node in the graph that pass through x :	
$BC(x) = \sum_{i,j \in N - \{x\}} \frac{\sigma_{i,j}(x)}{\sigma_{i,j}}$ Where $\sigma_{i,j}$ is the total number of shortest paths from node i to j and $\sigma_{i,j}(x)$ is the total number of those paths that pass through x .	
Nodes-edges ratio $\max(E) = N * (N - 1)$	
Average degree $\text{Averaging all nodes degrees. Scaling it to [0,1]}$	
Diameter $d = \max_{i \in N} D(i)$ where $D(N)$ is the eccentricity	
Density $D = \frac{2 E }{(N+1)(N-1)}$	
Modularity $Bonacich V.D., Guilleman J.L., Lambiotte R. Left brain, right brain: Fast unfolding of communities in large networks. Journal of Statistical Mechanics: Theory and Experiment, vol. 2008 (10), pp. 10008 (2008)$	
Clustering coefficient $\text{Watts-Strogatz} cc_1 = \frac{\sum_{i \in N} C(i)}{N}$	
Average path length $\text{Brandes, U. A Faster Algorithm for Betweenness Centrality. In: Journal of Mathematical Sociology 28(2). pp. 169-177 (2001)}$	

7 Features

Author profiling aims at identifying different traits such as age and gender of an author on the basis of her writings. We propose the novel EmoGraph graph-based approach where morphosyntactic categories are enriched with semantic and affective information following the next steps.

1 Morpho-syntactic annotation with Freeling

He estado tomando cursos en línea

VAIP1S0 VAP00SM VMG0000 NCMP000 RG

sobre temas valiosos que disfruto estudiando

SPS00 NCMP000 AQOMP0 PROCN000 VMIP1S0 VMG0000

y que podrían ayudarme a hablar en público

CC PROCN000 VMIC3P0 VMN0000 SPS00 VMN0000 SPS00

público .

NCMS000 Fp

Resources

Freeling	http://nlip.lsi.upc.edu/freeling/
WordNet Domains (+EuroWordnet)	http://wndomains.fbk.eu/ http://www.ilc.uva.nl/EuroWordNet/
Semantic Classification of Verbs	Lovin, B. English Verb Classes and Alternants. University of Chicago Press, Chicago. (1993)
Polarity Lexicon	Hu, M., Liu, B. Mining and Summarizing Customer Reviews. In: Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, Seattle, Washington, USA, pp. 168-177 (2004)
Spanish Emotion Lexicon	Sidorov, G., Miranda, S., Viveros, F., Gelbukh, A., Castro, N., Velásquez, F., Díaz, I., Suárez, S., Treviño, A., Gordon, J.: Empirical Study of Opinion Mining in Spanish Tweets. 11th Mexican International Conference on Artificial Intelligence, MICAI, pp. 1-14 (2012)



PAN-API4 dataset

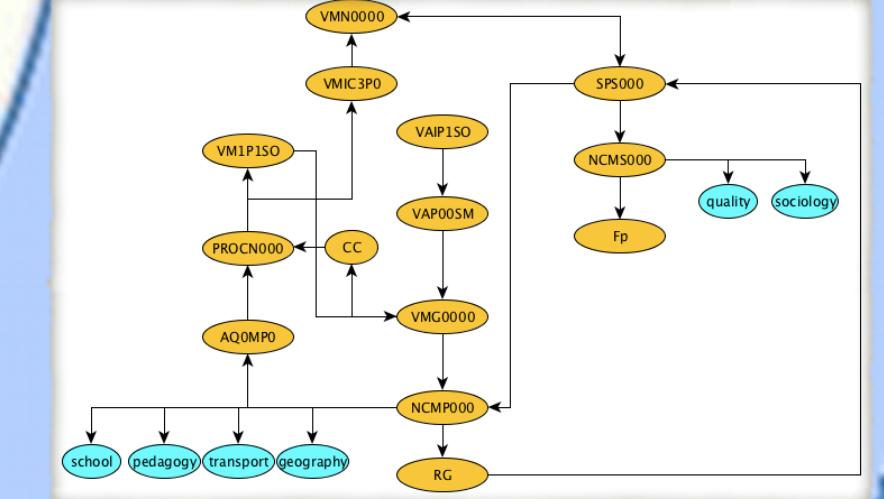
Social Media	Blogs		Twitter		Reviews		
	English	Spanish	English	Spanish	English	Spanish	English
18-24	680	150	10	4	12	4	74
25-34	900	180	24	12	56	26	200
35-49	980	138	32	26	58	46	200
50-64	790	70	10	10	26	12	200
65+	26	28	2	2	2	2	147
Σ	3376	566	78	54	154	90	821

2 Graph creation POS sequence

(ES) He estado tomando cursos en línea sobre temas valiosos que disfruto estudiando y que podrían ayudarme a hablar en público.

(EN) I have been taking online courses about valuable subjects that I enjoy studying and might help me to speak in public.

3 Topics with WordNet Domains



cursos NCMP000 transport geography
público NCMS000 pedagogy school
hablar VMN0000 sociology quality
language VMN0000

Conclusions

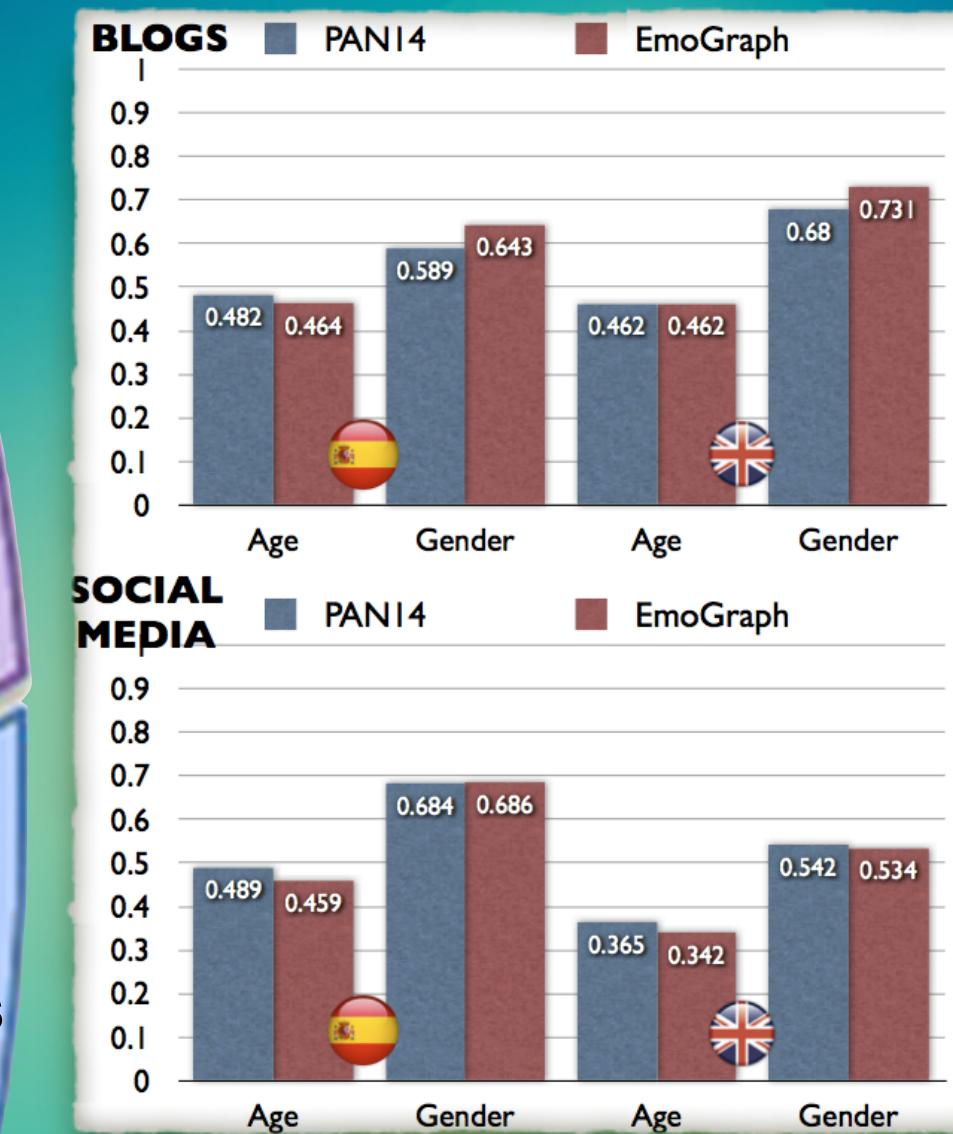
EmoGraph remains robust against different:

*Genres: Blogs, Twitter, Social Media and Reviews

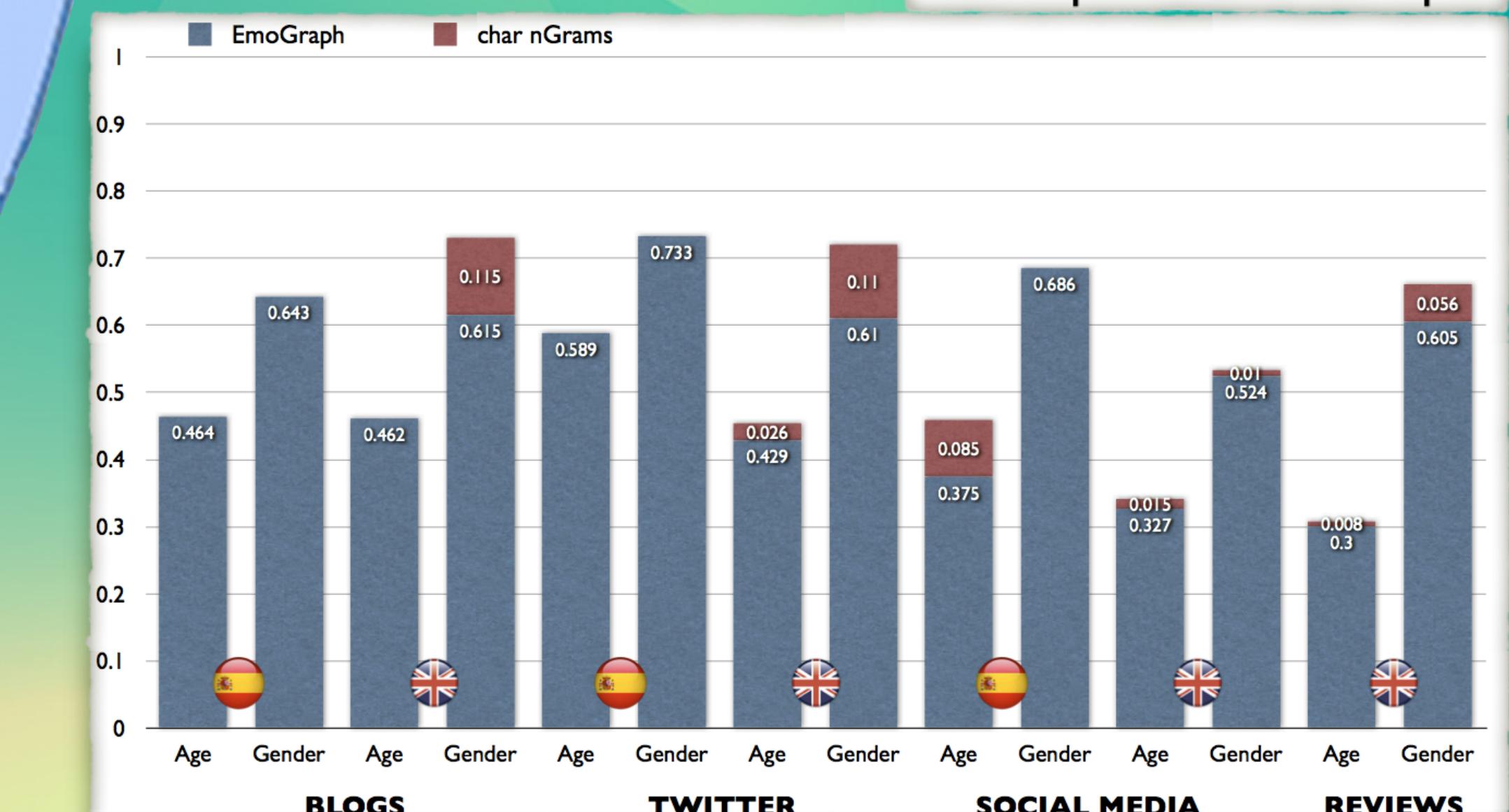
*Languages: Spanish, English

More competitive in Spanish

Age and Gender Identification



The Impact of EmoGraphs



Future Work

As future work we plan to investigate further the application of:

cost-sensitive machine learning techniques



autoritas®