

6th PPgSI's Dissertations Workshop 2019

Authorship attribution with authorship characterization techniques

Author:	Caio Deutsch			
Advisor:	Dr. Ivandré Paraboni			
Co-advisor:	NA			
Research lines:	<input type="checkbox"/> Systems Development and Management		<input checked="" type="checkbox"/> Systems Intelligence	
Research areas:	<input type="checkbox"/> Database <input type="checkbox"/> Software engineering <input type="checkbox"/> Information technology management <input type="checkbox"/> Human-Computer Interaction		<input checked="" type="checkbox"/> Artificial intelligence <input type="checkbox"/> Graphics processing <input checked="" type="checkbox"/> Pattern recognition <input type="checkbox"/> Optimization	
Application areas:	<input type="checkbox"/> Enterprise environments / Business processes <input type="checkbox"/> Bioinformatics <input type="checkbox"/> Biometrics <input type="checkbox"/> Mobile devices <input type="checkbox"/> Economy <input type="checkbox"/> Education / Distance learning <input type="checkbox"/> E-government <input checked="" type="checkbox"/> Internet / Social Networks <input type="checkbox"/> Games / Serious games <input checked="" type="checkbox"/> Linguistics / Natural Language <input type="checkbox"/> Cheminformatics <input type="checkbox"/> Robotics <input type="checkbox"/> Health <input type="checkbox"/> Other Which? _____ <input type="checkbox"/> General*			
Period in the program (at the workshop date):	<input type="checkbox"/> 2 nd semester	<input checked="" type="checkbox"/> 3 rd semester	<input type="checkbox"/> 4 th semester	<input type="checkbox"/> 5 th semester
Qualifying:	<input type="checkbox"/> Qualifying held in: dd/mm/yyyy		<input checked="" type="checkbox"/> Plan for qualifying in: 15/09/2019	
Defense:	Deadline for deposit: dd/mm/yyyy		Plan for defending in: 15/06/2020	
Publications associated with the master's project:	No publications to date.			

The research project summary

Context:

Authorship attribution (AA), a sub-area of Natural Language Processing (NLP), seeks to identify the author of a given text in a set of possible authors.

This task of identifying authors has been studied by experts in the field for at least six decades, where the analyzed texts were still handwritten and not digital as in the present day.

Technological advancement has led to the creation of large data sources for NLP studies, specifically for the AA problem. These studies have been addressing areas such as news, where so-called fake news cause controversy on common themes in our current society, such as politics and state elections.

Research problem:

Author attribution is increasingly showing importance for various social activities, especially forence analysis. Applications can help, for example, unravel mysteries of fake news authors, authoring source code, or even identifying aliases.

Studies involving AA demonstrate modest results and motivate the exploration of different techniques to improve the accuracy of current models. This is the case of (UCELAY et al., 2016), where the authors used authoritative characterization (CA) techniques to predict gender and age. These models could be used to solve the AA problem.

In addition, the problem of AA can be approached for several languages, however the vast majority of studies are focused on the English language, and studies in Brazilian Portuguese are rare.

Research objective:

Develop models capable of identifying authors of a particular text using characterization authorship (CA) techniques, in order to obtain superior results to traditional AA models.

Characteristics of the proposed solution:

The classification model is composed of a set of three other classifiers: n-grams of characters patterns, n-grams of non-diacritically distorted characters, and n-gram words and a binary variable which simulates a four classifier which will be a CA classifier in the next experiment. The four variables are ensemble with a logistic regression and used to predict a AA problem.

Theoretical foundations:

- Part of speech
- Text distortion
- Word embeddings

Correlated works:

- CUSTODIO, J. E.; PARABONI, I. Each-usp ensemble cross-domain authorship attribution: Notebook for pan at clef 2018
- UCELAY, M.; VILLEGAS, M.; FUNEZ, D.; CAGNINA, L.; ERRECALDE, M.; RAMIREZ-DE-LA-ROSA, G.; VILLATORO-TELLO, E. Profile-based approach for age and gender identification
- ROCHA, A.; SCHEIRER, W. J.; FORSTALL, C. W.; CAVALCANTE, T.; THEOPHILO, A.; SHEN, B.; CARVALHO, A. R. B.; STAMATATOS, E. Authorship attribution for social media forensics.

Validation

The experiments will be validated using conventional machine learning measures such as recall, accuracy, precision, and F measurement.

Limitations, risks, and threats:

The scope of this project is limited to AA problems using CA techniques with machine learning methods. Computational methods based on complex networks, graphs and compression models will not be considered. The mentioned corpus are limited to domains of literature, Social networks such as, facebook, twitter, whatsapp and others. The languages considered are Brazilian Portuguese and English.

Scientific contribution:

This research aims to advance the frontier of knowledge about the problem of AA using CA techniques using language and domain independent models. By developing and studying AA models using CA techniques, it is expected to advance knowledge of the relationship between language and the factors that determine authorship in both areas. In particular, it is expected to advance studies to the Portuguese language.

The research method

Genre (choose ONE)	<input type="checkbox"/> Theoretical research	<input checked="" type="checkbox"/> Practical research	<input type="checkbox"/> Empirical research	<input type="checkbox"/> Methodological research
Nature (choose ONE)	<input type="checkbox"/> Basic research		<input checked="" type="checkbox"/> Applied research	
Approach (choose ONE)	<input type="checkbox"/> Quantitative research	<input type="checkbox"/> Qualitative research	<input checked="" type="checkbox"/> Quali-quantitative research	
Literature review* (you can choose more than one)	<input type="checkbox"/> Narrative review	<input type="checkbox"/> Meta-analysis	<input type="checkbox"/> Theoretical review	
	<input checked="" type="checkbox"/> Descriptive review	<input type="checkbox"/> Qualitative systematic review	<input type="checkbox"/> Realistic review	
	<input type="checkbox"/> Scoping review	<input type="checkbox"/> Umbrella review	<input type="checkbox"/> Critical review	
Main technical procedure (choose ONE)	<input type="checkbox"/> Experimental research	<input type="checkbox"/> <i>Survey</i>	<input type="checkbox"/> Ethnographic research	
	<input type="checkbox"/> Bibliographic research	<input checked="" type="checkbox"/> Case study	<input type="checkbox"/> Grounded theory	
	<input type="checkbox"/> Documental research	<input type="checkbox"/> Participatory research	<input type="checkbox"/> Design science	
	<input type="checkbox"/> <i>Ex-post-facto</i> research	<input type="checkbox"/> Research-action	<input type="checkbox"/> Other Which? _____	
Data analysis (you can choose more than one)	<input checked="" type="checkbox"/> Descriptive statistics	<input type="checkbox"/> Statistical test	<input type="checkbox"/> Discourse analysis	
	<input checked="" type="checkbox"/> Inferential statistics	<input checked="" type="checkbox"/> Content analysis	<input type="checkbox"/> Others: _____	

* Definition of types of literature reviews established by Paré, G., Trudel M-C., Jaana M., Kitsiou, S. Synthesizing Information systems knowledge: A typology of literature reviews. In: Information & Management 52, p. 183-199, 2015. DOI: 10.1016/j.im.2014.08.008

Next steps:

- Finish the text for the qualification exam
- Develop a CA classifier to predict the binary variable used in experiment 1
- Write down the results of this new experiment
- Schedule the defense