

VI Workshop de Dissertações de Mestrado do PPgSI 2019

Atribuição autoral multicanais de textos digitais

Autoria de:	José Eleandro Custódio			
Orientação de:	Prof. Dr. Ivandrê Paraboni			
Coorientação de:				
Linha de pesquisa:	<input type="checkbox"/> Gestão e Desenvolvimento de Sistemas		<input checked="" type="checkbox"/> Inteligência de Sistemas	
Área de pesquisa:	<input type="checkbox"/> Banco de dados <input type="checkbox"/> Engenharia de software		<input checked="" type="checkbox"/> Inteligência artificial <input type="checkbox"/> Processamento gráfico	
	<input type="checkbox"/> Gestão de tecnologia da informação <input type="checkbox"/> Interação humano computador		<input type="checkbox"/> Reconhecimento de padrões <input type="checkbox"/> Otimização	
Área de aplicação:	<input type="checkbox"/> Ambientes corporativos / Processos de negócio		<input type="checkbox"/> Bioinformática <input type="checkbox"/> Biometria <input type="checkbox"/> Dispositivos móveis	
	<input type="checkbox"/> Economia <input type="checkbox"/> Educação / Educação a distância		<input type="checkbox"/> Governo eletrônico <input type="checkbox"/> Internet / Redes sociais	
	<input type="checkbox"/> Jogos / Jogos sérios		<input checked="" type="checkbox"/> Linguística / Língua natural <input type="checkbox"/> Químioinformática <input type="checkbox"/> Robótica	
	<input type="checkbox"/> Saúde <input type="checkbox"/> Outra Qual? _____		<input type="checkbox"/> Geral*	
Semestre no curso (na data do workshop):	<input type="checkbox"/> 2º semestre	<input type="checkbox"/> 3º semestre	<input type="checkbox"/> 4º semestre	<input checked="" type="checkbox"/> 5º semestre
Qualificação:	<input checked="" type="checkbox"/> Qualificação já realizada em: 30/11/2018		<input type="checkbox"/> Realização da qualificação planejada para: 30/11/2019	
Defesa:	Prazo máximo para depósito: 27/01/2020		Realização da defesa planejada para: 30/11/2019	
Publicações associadas ao projeto de mestrado:	<ul style="list-style-type: none"> • CUSTODIO, J. E. ; PARABONI, I. . Multichannel Open-set Cross-Domain Authorship Attribution. In: Conference and Labs of the Evaluation Forum (CLEF-2019), 2019, Lugano, Switzerland. Working Notes of the CLEF-2019 Evaluation Labs, 2019. v. 2380. • CUSTODIO, J. E. ; PARABONI, I. . Similaridade de textos aplicada à verificação autoral. In: I Congresso Internacional em Humanidades Digitais (HDRio-2018), 2018, Rio de Janeiro. Proceedings of HDRio-2018, 2018. p. 403-407. • CUSTODIO, J. E. ; PARABONI, I. . EACH-USP Ensemble Cross-Domain Authorship Attribution. In: Conference and Labs of the Evaluation Forum (CLEF-2018), 2018, Avignon. Working Notes Papers of the CLEF 2018 Evaluation Lab, 2018. v. 2125. 			

Resumo do projeto de pesquisa

Contexto:

A atribuição autoral(AA) de textos digitais busca identificar o autor de um texto baseando-se nos rastros deixados pelo mesmo ao escrever. Os rastros, ou estilo de escrita, podem ser examinados através da análise da utilização das palavras, seqüências de caracteres, pontuações e preferências gramaticais. No entanto, essas características são impactadas pelo domínio do texto, língua e quantidade de autores.

Problema de pesquisa:

Os métodos computacionais aplicados à atribuição autoral tendem a utilizar apenas um tipo de extração de características. No entanto, os rastros de autoria podem ser encontrados em diversos níveis linguísticos, como por exemplo, nível sintático, semântico ou lexical.

Objetivo de pesquisa:

Construir um método computacional que combine diversos tipos de conhecimento e de extração de características.

Caracterização da solução em desenvolvimento:

O método proposto irá construir classificadores que usem como extração de características n-gramas de palavras, n-gramas de caracteres, n-gramas de anotações POS e *embeddings*. Os modelos de POS e *embeddings* usarão bases pré-treinadas bibliotecas prontas como *embeddings* Google e a biblioteca Python Spacy. Os modelos construídos serão então combinados em um método único que irá fornecer a classificação dos documentos gerando a marcação do autor do texto.

Fundamentos:

Os n-gramas de caracteres e palavras mais frequentes são reconhecidos na literatura por capturar elementos inconscientes da escrita.

Trabalhos relacionados:

- Sobre trabalhos que analisam um tipo de embeddings para AA: Document embeddings learned on various types of authorship attribution. Gómez-Adorno, Helena et. al (2018); Continuous N-gram Representations for Authorship Attribution (2017); Authorship attribution in portuguese using character N-grams (2017); Authorship attribution using text distortion Stamatatos, Efethathios (2017).
- Modelos matemáticos que descrevem o conceito de embeddings:
- Distributed Representations of Words and Phrases and their Compositionality. Mikolov, Thomas et al (2013); A Neural Probabilistic Language Model. Bengio, Yoshua et al (2003).

Validação:

São usados conjuntos de dados públicos que possuam desempenho conhecido, e os experimentos serão desenvolvidos usando validação cruzada, exploração de parâmetros com grid search.

Limitações, riscos e ameaças:

Esse trabalho irá lidar primordialmente com o problema de atribuição autoral de conjunto fechado, onde os autores dos textos desconhecidos são necessariamente um dos autores do conjunto de treinamento. Demais variações poderão ser exploradas apenas para expansão do conhecimento sobre o tema

Os métodos propostos serão limitados a um conjunto de até 30 autores candidatos (30 classes), estando fora do escopo, o cenário com milhares de autores.

Contribuição científica:

O trabalho pretende aumentar o conhecimento sobre o funcionamento dos n-gramas de caracteres e outros para atribuição autoral, pretende expandir o conhecimento sobre as formas de combinação de modelos computacionais.

Contribuição tecnológica (se pertinente):

Método de pesquisa				
Gênero (escolha UM)	<input type="checkbox"/> Pesquisa teórica	<input type="checkbox"/> Pesquisa prática	<input checked="" type="checkbox"/> Pesquisa empírica	<input type="checkbox"/> Pesquisa metodológica
Natureza (escolha UMA)	<input type="checkbox"/> Pesquisa básica		<input checked="" type="checkbox"/> Pesquisa aplicada	
Abordagem (escolha UMA)	<input checked="" type="checkbox"/> Pesquisa quantitativa	<input type="checkbox"/> Pesquisa qualitativa	<input type="checkbox"/> Pesquisa mista (quali-quant)	
Revisão de literatura* (você pode escolher mais de uma)	<input type="checkbox"/> Revisão narrativa	<input type="checkbox"/> Meta-análise	<input type="checkbox"/> Revisão teórica	
	<input type="checkbox"/> Revisão descritiva	<input checked="" type="checkbox"/> Revisão sistemática qualitativa	<input type="checkbox"/> Revisão realística	
	<input type="checkbox"/> Revisão de escopo	<input type="checkbox"/> Revisão <i>guarda-chuva</i>	<input type="checkbox"/> Revisão crítica	
Procedimento técnico principal (escolha UM)	<input checked="" type="checkbox"/> Pesquisa experimental	<input type="checkbox"/> Pesquisa com <i>survey</i>	<input type="checkbox"/> Pesquisa etnográfica	
	<input type="checkbox"/> Pesquisa bibliográfica	<input type="checkbox"/> Estudo de caso	<input type="checkbox"/> Teoria fundamentada em dados	
	<input type="checkbox"/> Pesquisa documental	<input type="checkbox"/> Pesquisa participante	<input type="checkbox"/> Ciência do projeto	
	<input type="checkbox"/> Pesquisa <i>ex-post-facto</i>	<input type="checkbox"/> Pesquisa-ação	<input type="checkbox"/> Outra Qual? _____	
Análise de dados (você pode escolher mais de uma)	<input checked="" type="checkbox"/> Estatística descritiva	<input type="checkbox"/> Teste estatístico	<input type="checkbox"/> Análise do discurso	
	<input checked="" type="checkbox"/> Estatística inferencial	<input type="checkbox"/> Análise de conteúdo	<input type="checkbox"/> Outros: _____	

* Definição de tipos de revisões de literatura estabelecida por 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

Próximas atividades:

[Descrição dos experimentos e confecção da dissertação.](#)

Opcional: Forneça um esquema gráfico que mostre aspectos de sua pesquisa. Por exemplo: um fluxograma para construção da sua solução ou um infográfico para sua proposta de pesquisa. Se necessário, use a quarta página.