

V Workshop de Dissertações de Mestrado do PPgSI (2018)

FICHA DE PESQUISA

DADOS GERAIS				
Título do projeto de pesquisa	Atribuição autoral de textos digitais utilizando embeddings			
Orientando	José Eleandro Custódio			
Orientador(es)	Ivandr� Paraboni			
Semestre no curso, na data do workshop	<input type="checkbox"/> 2º semestre	<input checked="" type="checkbox"/> 3º semestre	<input type="checkbox"/> 4º semestre	<input type="checkbox"/> 5º semestre
Qualificação	[<input type="checkbox"/>] Qualificação já realizada em: dd/mm/aaaa [<input checked="" type="checkbox"/>] Realização da qualificação planejada para: 29/10/2018			
Defesa	Prazo máximo para depósito: 29/10/2018		Realização da defesa planejada para: 01/05/2019	
Linha e Área de pesquisa	Gestão e Desenvolvimento de Sistemas: [<input type="checkbox"/>] BD [<input type="checkbox"/>] Engenharia de Software [<input type="checkbox"/>] Gestão de TI [<input type="checkbox"/>] IHC		Inteligência de Sistemas: [<input checked="" type="checkbox"/>] IA [<input type="checkbox"/>] Processamento Gráfico [<input type="checkbox"/>] Reconhecimento de Padrões	
Área de aplicação	[<input type="checkbox"/>] Ambientes corporativos [<input type="checkbox"/>] Bioinformática [<input type="checkbox"/>] Biometria [<input type="checkbox"/>] Dispositivos móveis [<input type="checkbox"/>] Educação	[<input type="checkbox"/>] Educação a distância [<input type="checkbox"/>] Governo eletrônico [<input type="checkbox"/>] Internet [<input type="checkbox"/>] Jogos [<input type="checkbox"/>] Jogos sérios	[<input checked="" type="checkbox"/>] Língua Natural [<input checked="" type="checkbox"/>] Linguística [<input type="checkbox"/>] Processos de Negócio [<input type="checkbox"/>] Quimioinformática	[<input type="checkbox"/>] Redes Sociais [<input type="checkbox"/>] Robótica [<input type="checkbox"/>] Saúde [<input type="checkbox"/>] [outro – escrever]
Publicações associadas ao projeto de mestrado	Foram publicados dois trabalhos: - “Similaridade de Textos aplicada à Verificação Autoral” apresentado em poster no HDRio2018 - EACH-USP Ensemble Cross-domain Authorship Attribution: Notebook for PAN at CLEF 2018 – Um artigo descrevendo o método utilizado no sistema de atribuição autoral utilizando na competição PAN-CLEF e que atingiu o melhor resultado geral da edição de 2018.			

DESCRIÇÃO DO PROJETO DE PESQUISA	
Contextualização / motivação	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	Métodos tradicionais de processamento de língua natural foram amplamente estudados para AA. No entanto, técnicas recentes de embeddings foram pouco estudadas e não há uma análise definitiva de suas eficácias para a AA.
Objetivo geral da pesquisa	Comparar os diferentes métodos de embeddings com métodos baselines estabelecidos para a atribuição autoral.
Trabalhos relacionados	Sobre trabalhos que analisam um tipo de embeddings para AA: Document embeddings learned on various types of n-grams for cross-topic authorship attribution. Gómez-Adorno, Helena et. al (2018); Continuous N-gram Representations for Authorship Attribution. Sari et al (2017); Sobre trabalhos que analisam técnicas reconhecidas para AA: Authorship attribution in portuguese using character N-grams Markov, Ila et al (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).
Proposta para solução	Estender o trabalho realizado para o PAN-CLEF-2018 através de: - Aplicação de embeddings aplicados nos níveis caracter, palavra e POS tag. - Aplicar um comitê de máquina para combinar os diferentes métodos.
Dados	O projeto utilizará conjuntos de dados disponibilizados pelos organizadores do PAN-CLEF para edição de 2018 e demais edições.
Validação	Será feito através da comparação de medidas estatísticas entre resultados obtidos e o baseline PAN-CLEF 2018.

