

V Workshop de Dissertações de Mestrado do PPgSI (2018)

FICHA DE PESQUISA

DADOS GERAIS				
Título do projeto de pesquisa	Caracterização autoral a partir de textos			
Orientando	Rafael Felipe Sandroni Dias			
Orientador(es)	Prof. Dr. Ivandré Paraboni			
Semestre no curso, na data do workshop	<input type="checkbox"/> 2º semestre	<input type="checkbox"/> 3º semestre	<input checked="" type="checkbox"/> 4º semestre	<input type="checkbox"/> 5º semestre
Qualificação	<input checked="" type="checkbox"/> Qualificação já realizada em: 29/06/2018 <input type="checkbox"/> Realização da qualificação planejada para: dd/mm/aaaa			
Defesa	Prazo máximo para depósito: 03/09/2019 Realização da defesa planejada para: 04/2019			
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 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	Publicado: Rafael Felipe Sandroni Dias; Ivandré Paraboni (2018) Author profiling using word embeddings with subword information. 6th Author Profiling Task at PAN 2018: Multimodal Gender Identification in Twitter (PAN-CLEF 2018) Publicado: Rafael Felipe Sandroni Dias; Ivandré Paraboni (2018) Caracterização autoral de usuários do facebook: gênero, idade, religiosidade e área de formação. I Congresso Internacional em Humanidades Digitais (HDRio-2018). Publicado: Hsieh, Fernando Chiu; Rafael Felipe Sandroni Dias; Ivandré Paraboni (2018) Author Profiling from Facebook Corpora. 11th International Conference on Language Resources and Evaluation (LREC-2018) pp. 2566-2570. Miyazaki, Japan.			

DESCRIÇÃO DO PROJETO DE PESQUISA	
Contextualização / motivação	Na literatura, modelos de caracterização autoral (CA) utilizam, em sua grande maioria, conhecimentos linguísticos especializados para cada tipo de tarefa, idioma e domínio de texto. Estudos recentes de processamento de língua natural (PLN) têm obtido resultados promissores em modelos independentes de idioma e domínio com uso de técnicas baseadas em redes neurais artificiais.
Problema de pesquisa	Gerar uma abordagem independente do uso de conhecimento linguístico especializado para cada tarefa, idioma e domínio, para obter as características autorais de um indivíduo, e que apresentem maior eficiência que modelos tradicionais utilizados na área de CA, que fazem uso de conhecimento especializado para computar características (feature engineering) em textos.
Objetivo geral da pesquisa	Proposta de pesquisa para elaboração de um modelo independente de idioma e domínio aplicado ao reconhecimento de características autorais à partir de documentos escritos nos idiomas português, inglês e espanhol, nos domínios de redes sociais, blogs, entrevistas e formulários do governo, e devidamente anotados com informações de gênero, idade, grau de religiosidade, grau de escolaridade, área de formação, visão política e localidade.
Trabalhos relacionados	Fatima apresenta em seu trabalho uma abordagem multilíngue para predição de gênero e idade, considerando características de estilo de escrita e conteúdo, além de informações de n-gramas de caracteres. Sap, desenvolve um modelo léxico independente de domínio de texto, considerando a predição de idade como um problema de regressão. Sierra, apresenta resultados iniciais com modelos baseados em word embeddings e redes neurais de convolução (CNNs) para predição de gênero e idade.
Proposta para solução	Abordagem que permita o reconhecimento de características autorais de forma independente de conhecimento linguístico especializado para cada tarefa, idioma e domínio. Para isso, será desenvolvido uma estratégia combinando modelos de word embeddings e redes neurais profundas para extração de características e classificação de textos.
Dados	O projeto utiliza cópulas públicas devidamente anotados com informações de gênero, idade, grau de religiosidade, grau de escolaridade, área de formação, visão política e localidade, nos idiomas português, inglês e espanhol. Tais cópulas serão utilizados para avaliação do modelo de caracterização autoral independente de idioma e domínio. Os cópulas são: PAN-CLEF 2013, PAN-CLEF 2014, The Blog Authorship, B5-Post, BRMoral, BlogSet-BR e e-sic.

