

**VI Workshop de Dissertações de Mestrado do PPgSI
2019**

Caracterização autoral a partir de textos utilizando redes neurais artificiais

| | | | | |
|--|---|--------------------------------------|--|---|
| Autoria de: | Rafael Felipe Sandroni Dias | | | |
| Orientação de: | Profa. Dr. Ivandré 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 type="checkbox"/> Gestão de tecnologia da informação <input type="checkbox"/> Interação humano computador | | <input checked="" type="checkbox"/> Inteligência artificial <input type="checkbox"/> Processamento gráfico <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"/> Quimioinformática <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"/> 2ª semestre | <input type="checkbox"/> 3ª semestre | <input type="checkbox"/> 4ª semestre | <input checked="" type="checkbox"/> 5ª semestre |
| Qualificação: | <input checked="" type="checkbox"/> Qualificação já realizada em: 29/06/2018 | | <input type="checkbox"/> Realização da qualificação planejada para: | |
| Defesa: | Prazo máximo para depósito: 03/09/2019 | | Realização da defesa planejada para: 03/09/2019 | |
| Publicações associadas ao projeto de mestrado: | <ul style="list-style-type: none"> • 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. | | | |

Resumo do projeto de pesquisa

Contexto:

A caracterização autoral (do inglês, Author Profiling) é uma tarefa computacional de reconhecimento de características sociais (i.e., caracterização de gênero, faixa etária, grau de educação etc) de autores a partir de textos. Na literatura, modelos de caracterização autoral utilizam, em sua grande maioria, conhecimentos linguísticos especializados para cada tipo de tarefa, idioma e domínio de texto. Contudo, estudos recentes de processamento de língua natural (PLN) têm obtido resultados satisfatórios com uso de redes neurais artificiais para extração de informações sintáticas e semântica em problemas de classificação de textos.

Problema de pesquisa:

Os modelos computacionais de caracterização autoral tendem a utilizar diversos métodos baseados em conhecimento linguístico prévio. No entanto, em diversas tarefas de PLN, os métodos de redes neurais artificiais apresentam resultados satisfatórios para extração de informações sintáticas e semântica em problemas de classificação de textos.

Objetivo de pesquisa:

Pesquisa para elaboração de modelos independentes de idioma e domínio textual baseados em redes neurais artificiais 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 devidamente anotados com informações de gênero, faixa etária, grau de religiosidade, grau de escolaridade, área de formação e visão política. Os corpúscos são: PAN-CLEF 2013, The Blog Authorship, B5-post, BRMoral e BlogSet-BR. Tais corpúscos serão utilizados para avaliação dos modelos de caracterização autoral propostos fazendo uso da medida F para comparação com sistemas de baseline.

Caracterização da solução em desenvolvimento:

Este estudo faz uso de técnicas de redes neurais artificiais para problemas de classificação de textos, tais como redes neurais de convolução (CNN) e redes neurais recorrentes (i.e., Long-Short Term Memory - LSTM). Além de métodos de representação distribuída de palavras (i.e., word embeddings) e caracteres (i.e., char embeddings), tais como word2vec e fasttext.

Fundamentos:

Os modelos de redes neurais do tipo CNN e LSTM, e representação distribuída de palavras, são conhecidos na literatura por capturar elementos semânticos e sintáticos da escrita.

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.

Validação:

As respostas dos modelos computacionais propostos e sistemas de baseline são comparadas, utilizando medida F e acurácia. Os corpúscos são separados entre conjuntos de treinamento e testes para tal avaliação.

Limitações, riscos e ameaças:

O estudo será limitado aos conjuntos de dados selecionados e às tarefas de caracterização autoral por eles suportadas.

Contribuição científica:

As contribuições científicas previstas para este trabalho são a definição de tarefas de caracterização autoral, assim como avançar a fronteira de conhecimento de abordagens baseadas em redes neurais artificiais para as tarefas de caracterização autoral de modo que gerem resultados de referência para futuros estudos desta área.

Contribuição tecnológica (se pertinente):

As contribuições tecnológica previstas para este trabalho são a organização de diversos corpúsculos e tarefas de caracterização autoral, modelos baseados em redes neurais artificiais para estas tarefas e resultados de referência para futuros estudos desta área.

Método de pesquisa

| | | | | |
|--|--|--|--|--|
| Gênero (escolha UM) | <input type="checkbox"/> Pesquisa teórica | <input checked="" type="checkbox"/> Pesquisa prática | <input 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-quantitativa) | |
| 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 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 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:

Atividades para agosto/2019.

1. Revisão do texto (até 15 de agosto)
2. Depósito da dissertação (até 03 de setembro)
3. Defesa da dissertação (até 15 de setembro)

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.