

II Workshop de Dissertações de Mestrado do PPgSI

FICHA DA PESQUISAⁱ

DADOS GERAIS																	
Título do projeto de pesquisa	Aprendizado de dados positivos e não rotulados para o tratamento de incerteza na rotulação de dados de química medicinal																
Orientando	João Carlos Silva de Souza																
Orientador(es)	Profa. Dra. Patrícia Rufino Oliveira e Profa. Dra. Káthia Maria Honório																
Momento atual	<input type="checkbox"/> 3º semestre <input checked="" type="checkbox"/> 4º semestre <input type="checkbox"/> 5º semestre <input type="checkbox"/> 6º semestre																
Qualificação	<input checked="" type="checkbox"/> Qualificação já realizada em: 17/09/2015 <input type="checkbox"/> Qualificação planejada para: ____/____/____																
Defesa	Prazo máximo para depósito: _30/01/2017 Depósito planejado para: 19/12/2015																
Linha e Área de pesquisa	<table style="width: 100%; border: none;"> <tr> <td style="width: 50%; border: none;"> <input type="checkbox"/> Gestão e desenvolvimento de Sistemas: <input type="checkbox"/> BD <input type="checkbox"/> Gestão de SI <input type="checkbox"/> Eng. de Software <input type="checkbox"/> IHC </td> <td style="width: 50%; border: none;"> <input checked="" type="checkbox"/> Inteligência de Sistemas: <input checked="" type="checkbox"/> IA <input type="checkbox"/> Rec. de Padrões <input type="checkbox"/> Proc. Gráfico </td> </tr> </table>	<input type="checkbox"/> Gestão e desenvolvimento de Sistemas: <input type="checkbox"/> BD <input type="checkbox"/> Gestão de SI <input type="checkbox"/> Eng. de Software <input type="checkbox"/> IHC	<input checked="" type="checkbox"/> Inteligência de Sistemas: <input checked="" type="checkbox"/> IA <input type="checkbox"/> Rec. de Padrões <input type="checkbox"/> Proc. Gráfico														
<input type="checkbox"/> Gestão e desenvolvimento de Sistemas: <input type="checkbox"/> BD <input type="checkbox"/> Gestão de SI <input type="checkbox"/> Eng. de Software <input type="checkbox"/> IHC	<input checked="" type="checkbox"/> Inteligência de Sistemas: <input checked="" type="checkbox"/> IA <input type="checkbox"/> Rec. de Padrões <input type="checkbox"/> Proc. Gráfico																
Área de aplicação	<table style="width: 100%; border: none;"> <tr> <td style="width: 25%; border: none;"><input type="checkbox"/> Ambientes Corporativos</td> <td style="width: 25%; border: none;"><input type="checkbox"/> Educação</td> <td style="width: 25%; border: none;"><input type="checkbox"/> Linguagem Natural</td> <td style="width: 25%; border: none;"><input type="checkbox"/> Redes Sociais</td> </tr> <tr> <td style="border: none;"><input checked="" type="checkbox"/> Bioinformática</td> <td style="border: none;"><input type="checkbox"/> Educação a Distância</td> <td style="border: none;"><input type="checkbox"/> Linguística</td> <td style="border: none;"><input type="checkbox"/> Robótica</td> </tr> <tr> <td style="border: none;"><input type="checkbox"/> Biometria</td> <td style="border: none;"><input type="checkbox"/> Internet</td> <td style="border: none;"><input type="checkbox"/> Processos de Negócio</td> <td style="border: none;"><input type="checkbox"/> Saúde</td> </tr> <tr> <td style="border: none;"><input type="checkbox"/> Economia</td> <td style="border: none;"><input type="checkbox"/> Jogos</td> <td style="border: none;"><input type="checkbox"/> Química</td> <td></td> </tr> </table>	<input type="checkbox"/> Ambientes Corporativos	<input type="checkbox"/> Educação	<input type="checkbox"/> Linguagem Natural	<input type="checkbox"/> Redes Sociais	<input checked="" type="checkbox"/> Bioinformática	<input type="checkbox"/> Educação a Distância	<input type="checkbox"/> Linguística	<input type="checkbox"/> Robótica	<input type="checkbox"/> Biometria	<input type="checkbox"/> Internet	<input type="checkbox"/> Processos de Negócio	<input type="checkbox"/> Saúde	<input type="checkbox"/> Economia	<input type="checkbox"/> Jogos	<input type="checkbox"/> Química	
<input type="checkbox"/> Ambientes Corporativos	<input type="checkbox"/> Educação	<input type="checkbox"/> Linguagem Natural	<input type="checkbox"/> Redes Sociais														
<input checked="" type="checkbox"/> Bioinformática	<input type="checkbox"/> Educação a Distância	<input type="checkbox"/> Linguística	<input type="checkbox"/> Robótica														
<input type="checkbox"/> Biometria	<input type="checkbox"/> Internet	<input type="checkbox"/> Processos de Negócio	<input type="checkbox"/> Saúde														
<input type="checkbox"/> Economia	<input type="checkbox"/> Jogos	<input type="checkbox"/> Química															

DESCRIÇÃO DO PROJETO DE PESQUISA	
Contextualização / motivação	O projeto aborda o problema da incerteza na rotulação de dados e seus efeitos em processos de classificação, mais especificamente, este trabalho se destina a tratar o problema de incerteza em conjuntos de dados medicinais com o principal intuito de melhorar os modelos de classificação criados a partir de dados incertos ou não rotulados.
Problema de pesquisa	Reduzir o impacto na classificação de dados de química medicinal em cenários onde existe incerteza ou ausência na rotulação.
Objetivo geral da pesquisa	Utilizar uma técnica semi-supervisionada de aprendizagem de máquina para estimar uma rotulação confiável para os compostos candidatos a novos fármacos.
Trabalhos relacionados	Machine learning methods for property prediction in chemoinformatics: Quo Vadis? Graph kernels for chemical informatics Genome-wide sequence-based prediction of peripheral proteins using a novel semi-supervised learning technique
Justificativa e relevância	A rotulação incorreta das instâncias utilizadas na fase de treinamento do modelo prejudica a eficiência da classificação de novos compostos. Desse modo, encontrar uma forma de estimar corretamente o rótulo dos candidatos a fármacos, otimiza o processo de descoberta de medicamentos. Apesar de alguns trabalhos terem obtido sucesso na análise da estrutura dos compostos, este trabalho propõe uma análise que não depende de conversão de dados, ao contrário do que acontece principalmente nas análises baseadas em grafos.
Proposta para Solução	Será avaliada uma técnica de aprendizado semi-supervisionado que além dos dados rotulados, utiliza os dados rotulados com incerteza, na tentativa de estimar uma rotulação confiável. Em seguida, a nova rotulação será utilizada para construir um classificador baseado no aprendizado extremo de máquinas.
Dados	Será utilizado a base de dados pública PubChem Bioassay Data Set, na qual existe a possibilidade de extrair o valor do índice de atividade e a rotulação, e uma base de dados contendo dados sobre compostos para o tratamento de diabetes <i>Mellitus</i> e síndrome metabólica.
Forma de validação	Após o processamento dos dados pretende-se analisar classificador final comparando casos onde os dados foram estimados no aprendizado PU e sem aprendizado PU.
Limitações	Aquisição de base de dados com quantidade de elementos satisfatórios e definição confiável do limiar de atividade biológica.
Resultados esperados	<p>Contribuições científicas: Espera-se provar que as rotulações estimadas pelo aprendizado PU contribuam de forma positiva na construção de um classificador baseado em aprendizado extremo de máquinas.</p> <p>Contribuições tecnológicas: Uso de máquinas de aprendizado extremo na classificação de dados de química medicinal.</p>

MÉTODO DE PESQUISA

Gênero	<input type="checkbox"/> Pesquisa teórica	<input checked="" type="checkbox"/> Pesquisa prática	<input type="checkbox"/> Pesquisa empírica	<input type="checkbox"/> Pesquisa metodológica
Natureza	<input type="checkbox"/> Pesquisa básica/pura	<input checked="" type="checkbox"/> Pesquisa aplicada		
Objetivo	<input type="checkbox"/> Pesquisa descritiva	<input checked="" type="checkbox"/> Pesquisa exploratória	<input type="checkbox"/> Pesquisa explicativa	
Abordagem	<input checked="" type="checkbox"/> Pesquisa quantitativa	<input type="checkbox"/> Pesquisa qualitativa	<input type="checkbox"/> Pesquisa mista (quali-quant)	
Procedimento(s) técnico(s)	<input type="checkbox"/> Pesquisa experimental <input type="checkbox"/> Pesquisa bibliográfica <input type="checkbox"/> Pesquisa documental <input type="checkbox"/> Pesquisa <i>ex-post-facto</i> <input type="checkbox"/> Pesquisa de levantamento	<input type="checkbox"/> Pesquisa com <i>survey</i> <input type="checkbox"/> Estudo de caso <input type="checkbox"/> Pesquisa participante <input type="checkbox"/> Pesquisa-ação <input type="checkbox"/> Pesquisa etnográfica	<input type="checkbox"/> Pesquisa netnográfica <input checked="" type="checkbox"/> Teoria fundamentada em dados (<i>grounded theory</i>) <input type="checkbox"/> Ciência do projeto (<i>Design science research</i>)	
Fonte(s) de dados	<input type="checkbox"/> pesquisa de laboratório	<input type="checkbox"/> pesquisa de campo	<input checked="" type="checkbox"/> pesquisa bibliográfica	
Técnica(s) / Instrumento(s) de coleta de dados	<input type="checkbox"/> medição <input type="checkbox"/> questionário <input type="checkbox"/> entrevista <input type="checkbox"/> grupos focais	<input type="checkbox"/> formulário <input type="checkbox"/> <i>benchmark</i>	<input type="checkbox"/> observação (direta / participante) <input type="checkbox"/> diário de campo/notas de campo <input type="checkbox"/> análise documental (ou de artefatos) <input checked="" type="checkbox"/> Análise de banco de dados públicos de química medicinal	
Técnica(s) de análise de dados	<input checked="" type="checkbox"/> Análise quantitativa: <input type="checkbox"/> Estatística descritiva <input type="checkbox"/> Estatística inferencial		<input type="checkbox"/> Análise qualitativa: <input type="checkbox"/> Análise de conteúdo <input type="checkbox"/> Análise do discurso	

CRONOGRAMA

	2014												2015												2016											
	1	2	3	4	5	6	7	8	9	10	11	12	1	2	3	4	5	6	7	8	9	10	11	12	1	2	3	4	5	6	7	8	9	10	11	12
Levantamento Bibliográfico																																				
Revisão sistemática																																				
Discretização																																				
Aprendizado PU																																				
Algoritmos Suporte																																				
Experimentos																																				
Resultados																																				
Produção de Artigo																																				
Escrita da dissertação																																				
Depósito da dissertação																																				
Defesa																																				
Ajustes Finais																																				

ⁱ Esta ficha é uma adaptação da usada no “VIII Workshop de Teses e Dissertações em Sistemas de Informação (WTDSI 2015)” realizado como parte do “XI Simpósio Brasileiro de Sistemas de Informação (SBSI 2015)”