

III Workshop de Dissertações de Mestrado do Ppgsi (2016)

FICHA DA PESQUISAⁱ

DADOS GERAIS				
Título do projeto de pesquisa	Aprendizado semi-supervisionado para o tratamento de incerteza na rotulação de dados de química medicinal			
Orientando	João Carlos Silva de Souza – N. Usp: 8900551			
Orientador(es)	Profa. Dra. Patrícia Rufino Oliveira e Profa. Dra. Káthia Maria Honório			
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: 17/09/2015 <input type="checkbox"/> Realização da qualificação planejada para: dd/mm/aaaa			
Defesa	Prazo máximo para depósito: 30/01/2017 Realização da defesa planejada para: 10/2016			
Linha e Área de pesquisa	Gestão e desenvolvimento de Sistemas: <input type="checkbox"/> BD <input type="checkbox"/> Gestão de SI <input type="checkbox"/> <input type="checkbox"/> Eng. de Software <input type="checkbox"/> IHC <input type="checkbox"/>		Inteligência de Sistemas: <input checked="" type="checkbox"/> IA <input type="checkbox"/> Rec. de Padrões <input type="checkbox"/> <input type="checkbox"/> Proc. Gráfico <input type="checkbox"/>	
Área de aplicação	<input type="checkbox"/> Ambientes Corporativos <input type="checkbox"/> Educação <input type="checkbox"/> Bioinformática <input type="checkbox"/> Educação a Distância <input type="checkbox"/> Biometria <input type="checkbox"/> Internet <input type="checkbox"/> Economia <input type="checkbox"/> Jogos	<input type="checkbox"/> Linguagem Natural <input type="checkbox"/> Linguística <input type="checkbox"/> Processos de Negócio <input checked="" type="checkbox"/> Química	<input type="checkbox"/> Redes Sociais <input type="checkbox"/> Robótica <input type="checkbox"/> Saúde <input type="checkbox"/> [outro – escrever]	
Publicações associadas ao projeto de mestrado	Título: Recent Advances for Handling Imbalance and Uncertainty in Labelling in Medicinal Chemistry Data Analysis / Status: Aceito / “Conferência já realizada. Aguardando publicação e indexação pelo IEEE Explorer conforme informado por contato da conferência”. / Veículo: SAI Computing Conference 2016. Título: Handling Label Uncertainty on Chemistry Data Using Weighted Positive and Unlabeled Learning / Status: Em elaboração / Veículo: Ainda não definido.			

DESCRIÇÃO DO PROJETO DE PESQUISA

Contextualização / motivação	Ele é um caso específico no tratamento de incerteza na rotulação e desbalanceamento de dados de química medicinal. Estas características refletem negativamente em processos de classificação binária.
Problema de pesquisa	Baixo desempenho no processo de classificação de compostos devido à falta de instâncias que possuem rotulação confiável, métodos não sensíveis a diferente distribuição das classes e fase de treinamento demorada.
Objetivo geral da pesquisa	<i>Estimar uma rotulação mais confiável para os compostos utilizados no processo de descoberta de novos medicamentos na presença de dados desbalanceados.</i>
Trabalhos relacionados	Virtual Screening of Bioassy Data: Apresenta resultados utilizando alguns dos conjuntos de dados utilizados. Apesar de contemplar o desbalanceamento das classes, não fornece nenhum tratamento para o problema da incerteza - Genome-wide sequence based prediction of peripheral proteins using a novel semi-supervised learning technique – Aplica a técnica de aprendizado semi-supervisionado, mas não apresenta nenhuma abordagem sensível ao custo de erro de classificação. - Partially Supervised Classification of Text Documents: Não é especificamente de química medicinal, mas apresenta em detalhes a técnica de aprendizado semi-supervisionado.
Justificativa e relevância	A incerteza ou falta de rotulação são comuns nos dados de química medicinal, principalmente para as classes que não são alvos dos estudos. Encontrar uma forma eficiente de tratar esse problema, além de considerar o relevante cenário de desbalanceamento simultaneamente, promete melhorar o desempenho de métodos computacionais relacionados e aprimorar resultados no processo geral de descoberta de novas drogas.
Proposta para Solução	Utilizar a informação contida nos dados positivos e nos dados não rotulados para encontrar os rótulos adequados para as instâncias sob análise. De posse dessa nova rotulação, construir um classificador de rápido treinamento e alto desempenho, levando em consideração a quantidade deferente de exemplos em cada classe.
Dados	<i>Serão utilizados, para validação, dados de referência presentes na literatura de aprendizagem de máquina, e dados de química medicinal de domínio público presentes no repositório UCI para avaliação do modelo.</i>
Validação	<i>Além da comparação gráfica de desempenho, os exemplos do conjunto de treinamento serão identificados antes do processo e após a execução as instâncias com dados alterados serão comparadas com seu estado inicial. Espera-se que somente os exemplos próximos a uma área de sensibilidade sejam alterados.</i>

