

II Workshop de Dissertações de Mestrado do PPgSI

FICHA DA PESQUISAⁱ

DADOS GERAIS	
Título do projeto de pesquisa	Tratamento do desbalanceamento de classes em problemas de química medicinal.
Orientando	Suzana Gomes Claudino
Orientador(es)	Profa. Dra. Káthia Maria Honório e Profa. Dra. Patrícia Rufino Oliveira
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: 18/09/2015 <input type="checkbox"/> Qualificação planejada para: ____/____/____
Defesa	Prazo máximo para depósito: 31/01/2016 Depósito planejado para: 20/12/2015
Linha e Área de pesquisa	<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	<input type="checkbox"/> Ambientes Corporativos <input type="checkbox"/> Educação <input type="checkbox"/> Linguagem Natural <input type="checkbox"/> Redes Sociais <input 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 checked="" type="checkbox"/> Química <input type="checkbox"/> [outro – escrever]

DESCRIÇÃO DO PROJETO DE PESQUISA	
Contextualização / motivação	A busca de novos medicamentos é uma área importante da química medicinal, porém não é uma tarefa simples. Técnicas de classificação e aprendizagem de máquina podem ser utilizadas para auxiliar algumas tarefas neste processo.
Problema de pesquisa	A classificação de conjuntos de dados de natureza química sofre constantemente com o problema de desbalanceamento de dados, quando se trata de classificação binária entre compostos ativos e inativos frente a um alvo biológico.
Objetivo geral da pesquisa	Aplicar a técnica HC-kNN a fim de tratar o problema de desbalanceamento de classes na classificação de conjuntos de dados de natureza química.
Trabalhos relacionados	<p><i>A Hybrid Coupled k-Nearest Neighbor Algorithm on Imbalance Data</i> – Técnica proposta para classificação de conjuntos mistos com características numéricas e categóricas em conjuntos com dados desbalanceados.</p> <p><i>Influence relevance voting: an accurate and interpretable virtual high throughput screening method.</i> – Técnica proposta para tratamento de desbalanceamento na classificação de dados químicos.</p> <p><i>Machine learning methods for property prediction in chemoinformatics: Quo Vadis?</i> – Apresentar o estado da arte em pesquisas relacionadas a técnicas de aprendizagem de máquina na busca por novos medicamentos.</p>
Justificativa e relevância	O desbalanceamento de classes é um problema de classificação, pois muitas técnicas levam em consideração a informação a priori de que todos os compostos têm a mesma importância, o que faz com que a classificação seja tendenciosa para a classe majoritária e no contexto de quimioinformática, a classe de maior interesse é justamente a classe minoritária (compostos que apresentam rotulo “ativo”, ou seja, que apresentam atividade biológica frente a um alvo biológico).
Proposta para Solução	Está sendo avaliado o desempenho da técnica HC-kNN na classificação de compostos químicos e realizada a comparação com outras técnicas como o k-NN, variando as medidas de similaridade utilizadas.
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	Serão utilizados a validação cruzada e análise da curva ROC para a validação dos dados.
Limitações	Aquisição de base de dados com quantidade de elementos satisfatórios e definição confiável do limiar de atividade biológica que apresentem o comportamento de desbalanceamento severo.

