


PERÍODO: 2018.2		
CENTRO UNIVERSITÁRIO CTC DEPARTAMENTO DE INFORMÁTICA		
INF2791	Prof <sup>a</sup> . Clarisse Sieckenius de Souza	
Tópicos em IHC II	CARGA HORÁRIA TOTAL: 45	CRÉDITOS: 3
Explainable AI IA Explicável	Pré-requisito: Formal Nenhum (Recomenda-se algum conhecimento sobre Inteligência Artificial e Redes Neurais)	

<b>OBJETIVOS</b>	<ol style="list-style-type: none"> <li>1) Introduzir a questão dos ‘significados humanos’ embutidos em – e atribuídos a – aplicações contemporâneas de Inteligência Artificial</li> <li>2) Promover uma reflexão especificamente voltada para responsabilidade social e conduta ética no desenvolvimento de sistemas com aprendizado automático (Machine Learning), redes neurais (Neural Nets) e certos tipos de sistemas de apoio a decisão baseados em grandes volumes de dados (Data-Driven Decision Support Systems).</li> </ol>
<b>EMENTA</b>	<ul style="list-style-type: none"> <li>• <i>General Data Protection Regulation</i> (GDPR), a regulamentação europeia de Maio de 2018 sobre direitos dos donos de dados (Data Owners), em particular o <i>direito à explicação</i> (Right to Explanation);</li> <li>• Diferentes perspectivas sobre o que é uma “explicação”; explicabilidade e interpretabilidade de modelos de aprendizado automático;</li> <li>• O problema da explicação automática nas décadas de 1980 e 1990; o que mudou por volta de 2015;</li> <li>• Valores e vieses embutidos em algoritmos (<i>algorithmic accountability</i>);</li> <li>• Itens para uma agenda de pesquisa sobre responsabilidade social no desenvolvimento de modelos e aplicações de aprendizado de máquina.</li> </ul>
<b>PROGRAMA</b>	<p><b>Aula 1:</b> Introdução ao assunto; Panorâmica do Estado da Arte; a <i>GDPR</i>;</p> <p><b>Aula 2:</b> Aspectos do direito a explicação;</p> <p><b>Aula 3:</b> Exercitando interpretações do ‘direito a explicação’ em cenários reais e plausíveis;</p>

	<p><b>Aula 4:</b> Seminário sobre o que constitui um <i>explicação</i>;</p> <p><b>Aula 5:</b> Seminário sobre o que constitui um <i>explicação</i> (continuação);</p> <p><b>Aula 6:</b> <i>Explicações</i> em sistemas especialistas (IA simbólica);</p> <p><b>Aula 7:</b> Desafios de sistemas que utilizam redes neurais ou aprendizado profundo;</p> <p><b>Aula 8:</b> Desafios de sistemas que utilizam redes neurais ou aprendizado profundo (continuação);</p> <p><b>Aula 9:</b> Seminário sobre <i>Explainable AI</i>;</p> <p><b>Aula 10:</b> Seminário sobre <i>Explainable AI</i> (continuação);</p> <p><b>Aula 11:</b> Valores e vieses embutidos em algoritmos (<i>algorithmic accountability</i>);</p> <p><b>Aula 12:</b> Valores e vieses embutidos em algoritmos (continuação);</p> <p><b>Aula 13:</b> Itens para uma agenda de pesquisa sobre responsabilidade social no desenvolvimento de modelos e aplicações de aprendizado de máquina;</p> <p><b>Aula 14:</b> Itens para uma agenda de pesquisa sobre responsabilidade social no desenvolvimento de modelos e aplicações de aprendizado de máquina (continuação);</p> <p><b>Aula 15:</b> Resumo e encerramento da disciplina.</p>
<p><b>AVALIAÇÃO</b></p>	<p>Os alunos serão avaliados pela participação e apresentação em seminários da disciplina (aulas 4, 5, 9, 10), além de um relatório final (de até 6 páginas de texto, imagens, tabelas e gráficos, excetuadas as referências bibliográficas), contendo sua proposta de itens para uma agenda de pesquisa, não necessariamente pessoal, sobre <i>Explainable AI</i>. A disciplina inclui 1h por semana de estudos individuais sem horário fixo para os alunos se dedicarem ao estudo da bibliografia indicada.</p>
<p><b>BIBLIOGRAFIA PRINCIPAL</b></p>	<ol style="list-style-type: none"> <li>1) Alan Dix, <b>Neural Networks and Pattern Recognition in Human-computer Interaction</b>, pp. 429-451, Upper Saddle River, NJ, USA. Ellis Horwood. 1992.</li> <li>2) Victoria Bellotti and Keith Edwards, <b>Intelligibility and Accountability: Human Considerations in Context-Aware Systems</b>, <i>Human-Computer Interaction</i>, vol. 16, no. 2-4, pp. 193-212, Taylor &amp; Francis. 2001.</li> <li>3) Federica Russo and Jon Williamson, <b>Interpreting Causality in the Health Sciences</b>, <i>International Studies in the Philosophy of Science</i>, vol. 21, no. 2, pp. 157-170, 2007.</li> <li>4) Simone Stumpf, Vidya Rajaram, Lida Li, Weng-Keen Wong, Margaret Burnett, Thomas Dietterich, Erin Sullivan and Jonathan Herlocker, <b>Interacting meaningfully with machine learning systems: Three experiments</b>, <i>International Journal of Human-Computer Studies</i>, vol. 67, no. 8, pp. 639-662, 2009.</li> <li>5) Jonathan Grudin, <b>AI and HCI: Two fields divided by a common focus</b>, <i>AI Magazine</i>, vol. 30, no. 4, pp. 48, 2009.</li> </ol>

- 6) Jon Kolko, **Abductive Thinking and Sensemaking: The Drivers of Design Synthesis**, *Design Issues*, vol. 26, no. 1, pp. 15-28, 2010.
- 7) Lorenzo Casini, Phyllis McKay Illari, Federica Russo and Jon Williamson, **Models for Prediction, Explanation and Control: Recursive Bayesian Networks**, *Theoria. Revista de Teoría, Historia y Fundamentos de la Ciencia*, vol. 26, no. 1, pp. 5-33, January. 2011.
- 8) Illari, Phyllis and Russo, Federica. *Causality: Philosophical Theory meets Scientific Practice*. Oxford University Press, Oxfprd. Oxford University Press. 2014
- 9) Cathy O'Neil, *Weapons of math destruction: How big data increases inequality and threatens democracy*. Broadway Books, New York, NY. Broadway Books. 2016.
- 10) Zachary C. Lipton, **The mythos of model interpretability**, *arXiv:1606.03490*, 2016.
- 11) Lucian Leahu, **Ontological Surprises: A Relational Perspective on Machine Learning**, in *Proceedings of the 2016 ACM Conference on Designing Interactive Systems*, pp. 182-186, New York, NY, USA. ACM. 2016.
- 12) C. Hill, R. Bellamy, T. Erickson and M. Burnett, **Trials and tribulations of developers of intelligent systems: A field study**, in *2016 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*, pp. 162-170, September. 2016.
- 13) Sandra Wachter, Brent Mittelstadt and Luciano Floridi, **Why a right to explanation of automated decision-making does not exist in the general data protection regulation**, *International Data Privacy Law*, vol. 7, no. 2, pp. 76-99, 2017.
- 14) Andrew D. Selbst and Julia Powles, **Meaningful information and the right to explanation**, *International Data Privacy Law*, vol. 7, no. 4, pp. 233-242, 2017.
- 15) Tim Miller, **Explanation in Artificial Intelligence: Insights from the Social Sciences**, *arXiv preprint arXiv:1706.07269*, 2017.
- 16) Lilian Edwards and Michael Veale, **Slave to the Algorithm: Why a Right to an Explanation Is Probably Not the Remedy You Are Looking for**, *Duke Law and Technology Review*, vol. 16, pp. 1-65, 2017.
- 17) Maja Brkan, **AI-supported Decision-making Under the General Data Protection Regulation**, in *Proceedings of the 16th Edition of the International Conference on Artificial Intelligence and Law*, pp. 3-8, New York, NY, USA. ACM. 2017.
- 18) Reuben Binns, **Algorithmic Accountability and Public Reason**, *Philosophy & Technology*, May. 2017.
- 19) Sean Penney, Jonathan Dodge, Claudia Hilderbrand, Andrew Anderson, Logan Simpson and Margaret Burnett, **Toward Foraging for Understanding of StarCraft Agents: An Empirical Study**, in *23rd International Conference on Intelligent User Interfaces*, pp. 225-237, New York, NY, USA. ACM. 2018.
- 20) Judea Pearl and Dana Mackenzie, *The Book of Why: The New Science of Cause and Effect*. Basic Books., New York, NY. Basic Books.. 2018.

	<p>21) James Larus, Chris Hankin, Siri Granum Carson, Markus Christen, Silvia Crafa, Oliver Grau, Claude Kirchner, Bran Knowles, Andrew McGettrick, Damian Andrew Tamburri and Hannes Werthner, <b>When Computers Decide: European Recommendations on Machine-Learned Automated Decision Making</b>, New York, NY, USA. ACM. 2018. Technical Report: Informatics Europe &amp; EUACM 2018 (20 p.).</p> <p>22) Ashraf Abdul, Jo Vermeulen, Danding Wang, Brian Y. Lim and Mohan Kankanhalli, <b>Trends and Trajectories for Explainable, Accountable and Intelligent Systems: An HCI Research Agenda</b>, in <i>Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems</i>, pp. 582:1-582:18, New York, NY, USA. ACM. 2018.</p>
<p><b>BIBLIOGRAFIA COMPLEMENTAR</b></p>	<p>1) Jaime Nubiola, <b>Abduction or the Logic of Surprise</b>, <i>Semiotica</i>, vol. 2005, no. 153-1/4, pp. 117-130, 2005.</p> <p>2) Lorenzo Magnani, <b>An Abductive Theory of Scientific Reasoning</b>, <i>Semiotica</i>, vol. 153, no. 1/4, pp. 261-286, 2005.</p> <p>3) Maria Eunice Quilici Gonzalez and Willem Ferdinand Gerardus Haselager, <b>Creativity: Surprise and Abductive Reasoning</b>, <i>Semiotica</i>, vol. 153, no. 1/4, pp. 325-342, January. 2005</p> <p>4) T. Kulesza, S. Stumpf, M. Burnett, W. K. Wong, Y. Riche, T. Moore, I. Oberst, A. Shinsel and K. McIntosh, <b>Explanatory Debugging: Supporting End-User Debugging of Machine-Learned Programs</b>, in 2010 IEEE Symposium on Visual Languages and Human-Centric Computing, pp. 41-48, September. 2010</p> <p>5) P. J. G. Lisboa, <b>Interpretability in Machine Learning – Principles and Practice</b>, in Fuzzy Logic and Applications: 10th International Workshop, WILF 2013, Genoa, Italy, November 19-22, 2013. Proceedings, pp. 15-21, Cham. Springer International Publishing. 2013.</p> <p>6) Todd Kulesza, Margaret Burnett, Weng-Keen Wong and Simone Stumpf, <b>Principles of Explanatory Debugging to Personalize Interactive Machine Learning</b>, in Proceedings of the 20th International Conference on Intelligent User Interfaces, pp. 126-137, New York, NY, USA. ACM. 2015.</p> <p>7) Robert Folger and Christopher Stein, <b>Abduction 101: Reasoning processes to aid discovery</b>, <i>Human Resource Management Review</i>, vol. 27, no. 2, pp. 306-315, 2017.</p> <p>8) Simone Stumpf, Simonas Skrebe, Graeme Aymer and Julie Hobson, <b>Explaining smart heating systems to discourage fiddling with optimized behavior</b>, in IUI 2018 Workshop on Explainable Smart Systems (ExSS) - March 7-11, 2018, pp. 1-4, Tokyo, Japan. Online Publication at: <a href="http://explainablesystems.comp.nus.edu.sg/">http://explainablesystems.comp.nus.edu.sg/</a>. 2018</p> <p>9) Jonathan Dodge, Sean Penney, Andrew Anderson and Margaret Burnett, <b>What should be in an XAI explanation? What IFT reveals</b>, in IUI 2018 Workshop on Explainable Smart Systems (ExSS) - March 7-11, 2018, pp. 1-4, Tokyo, Japan. Online Publication at: <a href="http://explainablesystems.comp.nus.edu.sg/">http://explainablesystems.comp.nus.edu.sg/</a>. 2018.</p>

	<p>10) Jonathan Dodge, Sean Penney, Claudia Hilderbrand, Andrew Anderson and Margaret Burnett, <b>How the Experts Do It: Assessing and Explaining Agent Behaviors in Real-Time Strategy Games</b>, in Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems, pp. 562:1-562:12, New York, NY, USA. ACM. 2018.</p>
--	---