

Pós-Graduação da Faculdade de Saúde Pública da USP (FSP/USP)

Disciplina EPI 5717 Machine Learning para Predições em Saúde

2º Semestre – 2022

Créditos: 6 – 90 horas

Local: FSP/USP

Horário: quartas-feiras das 14 às 18 horas

Professor Responsável

Alexandre Chiavegatto Filho

Justificativa:

O rápido aumento na quantidade de dados tem aberto novas oportunidades para a saúde brasileira. Entre as várias novidades proporcionadas pelo big data em saúde, a mais promissora é o uso de modelos preditivos de inteligência artificial, conhecidos como machine learning. A disciplina tem como objetivo apresentar essa área em rápido crescimento com foco nas suas aplicações práticas, além de discutir seus benefícios, limitações e possíveis uso na área da saúde. O foco do curso será no tipo de dado mais coletado em saúde, i.e. dados estruturados/tabulares, e será utilizada a linguagem Python.

Programa

- 1 – Perspectivas do uso de inteligência artificial em saúde.
- 2 – Pré-processamento dos dados (padronização, one-hot encoding, imputação, outliers, rebalanceamento, vazamento de informação).
- 3 – Sobreajuste e divisão da amostra em treino, validação e teste.
- 4 – Mensuração da performance de algoritmos preditivos (área abaixo da curva ROC, precisão, recall, especificidade, valor predito negativo e raiz quadrada do erro quadrático médio).
- 5 – Algoritmos para predição de variável dependente contínua e categórica (regressões penalizadas com lasso e ridge, redes neurais, random forests, XGBoost, lightGBM e catboost).
- 6 – Técnicas de otimização de hiperparâmetros.
- 7 – Estratégias para a seleção de variáveis preditoras (Boruta).
- 8 – Aprendizado federado e aprendizado online (contínuo).
- 9 – Estratégias para a identificação da importância de variáveis preditoras (Shapley values).
- 10 – Desafios éticos do uso de machine learning em saúde.

Avaliação

A avaliação será realizada por meio de um trabalho final (60%) e exercícios realizados ao longo da disciplina (40%), sendo destes 20% referente aos entregáveis e 20% à discussão de artigos.

Discussão de artigos: todos os alunos devem entregar antes do início da aula uma revisão de menos de uma página para cada um dos quatro artigos, com um parágrafo de resumo e o resto uma avaliação sobre a importância e a qualidade do artigo em questão.

Diferentes alunos serão escolhidos para cada artigo, em que um irá apresentar o artigo e os outros irão liderar o debate do artigo em relação à sua importância e qualidade.

Observação

Para realizar o curso é necessário ter conhecimentos pelo menos básicos de estatística e programação.

Bibliografia

Batista AFM, Chiavegatto Filho ADP. Machine Learning aplicado à Saúde. In: Artur Ziviani; Natalia Castro Fernandes; Débora Christina Muchaluat Saade. (Org.). Livro de Minicursos. Niterói, RJ: Sociedade Brasileira de Computação, 2019.

Chiavegatto Filho ADP, Batista AFM, dos Santos HG. Data leakage in health outcomes prediction with machine learning. *Journal of Medical Internet Research* 2021; 23(1).

Fernandes TF, de Oliveira TA, Teixeira CE, Batista AFM, Costa GD, Chiavegatto Filho ADP. A multipurpose machine learning approach to predict COVID-19 negative prognosis in Sao Paulo, Brazil. *Scientific Reports* 2021; 3343(11).

Geron A. Mãos à obra: aprendizado de máquina com Scikit-Learn & TensorFlow. Alta Books; 2019.

Raschka S, Mirjalili V. Python Machine Learning - Third Edition: Machine Learning and Deep Learning with Python, scikit-learn, and TensorFlow 2. Packt Publishing; 2020.

Topol E. Deep Medicine: How artificial intelligence can make healthcare human again. Basic Books; 2019.

Extra:

Metz C. Genius Makers: The Mavericks Who Brought A.I. to Google, Facebook, and the World. Cornerstone; 2021.

Cronograma – 2023

Data	Tópico	Referência
08/08	Aula 1 - Discussão do conteúdo programático e apresentação da área.	Lones MA. How to avoid machine learning pitfalls: a guide for academic researchers. arXiv:2108.02497. 2021.
10/08	Aula 2 – Monitoria.	
15/08	Aula 3 – Pré-processamento dos dados.	“A Comprehensive Guide to Data Preprocessing” https://neptune.ai/blog/data-preprocessing-guide
17/08	Aula 4 – Monitoria.	
22/08	Aula 5 – Sobreajuste, viés e variância; divisão da amostra em treino, validação e teste.	Raschka S. Model evaluation, model selection, and algorithm selection in machine learning. arXiv:1811.12808. 2020.
24/08	Aula 6 - Monitoria	
29/08	Entregável 1: pré-processamento e divisão da amostra em treino e teste. Aula 7 - Mensuração da performance de algoritmos preditivos e otimização de hiperparâmetros.	“Tour of Evaluation Metrics for Imbalanced Classification” https://machinelearningmastery.com/tour-of-evaluation-metrics-for-imbalanced-classification/
31/08	Aula 8 – Monitoria.	
12/09	Aula 9 – Principais algoritmos de machine learning para dados estruturados: regressões penalizadas, redes neurais e algoritmos de árvore (árvores de decisão, random forests e gradient boosting: XGBoost, lightGBM e catboost).	Al-Shari H, Saleh YA, Odabas A. Comparison of Gradient Boosting Decision Tree Algorithms for CPU Performance. 2021.
14/09	Aula 10 – Monitoria.	
19/09	Aula 11 - Estratégias para a seleção de variáveis preditoras.	Degenhardt F. Evaluation of variable selection methods for random forests and omics data sets. Brief Bioinform. 2019;20(2):492-503.

21/09	Aula 12 – Monitoria.	
26/09	Aula 13 – Preparação das apresentações.	
28/09	Aula 14 - Discussão de artigos científicos de machine learning em saúde.	<p>Clift et al. Predicting 10-year breast cancer mortality risk in the general female population in England: a model development and validation study. <i>Lancet Digit Health</i> 2023;5(9):e571-e581.</p> <p>Yang et al. A gender specific risk assessment of coronary heart disease based on physical examination data. <i>NPJ Digit Med</i> 2023;6(1):136</p> <p>Finlayson SG et al. The clinician and dataset shift in artificial intelligence. <i>N Engl J Med</i> 2021; 385(3):283-286.</p> <p>Futoma J et al. The myth of generalisability in clinical research and machine learning in health care. <i>Lancet Digit Health</i> 2020;2(9):e489-e492.</p>
03/10	Entregável 2: predição. Aula 15 - Estratégias para a identificação da importância de variáveis preditoras.	Molnar CM et al. Interpretable machine learning -- a brief history, state-of-the-art and challenges. arXiv:2010.09337. 2020
05/10	Aula 16 – Monitoria.	
10/10	Aula 17 - Aprendizado federado (federated learning) e aprendizado online (contínuo).	<p>T. Li et al. Federated Learning: Challenges, Methods, and Future Directions. <i>IEEE Signal Processing</i> 2020;37(3):50-60.</p> <p>Hoi SCH et al. Online Learning: A Comprehensive Survey. arXiv:1802.02871. 2018</p>
19/10	Aula 18 – Monitoria.	
24/10	Aula 19 - Desafios éticos do uso de machine learning em saúde.	<p>Char DS et al. Identifying Ethical Considerations for Machine Learning Healthcare Applications. <i>Am J Bioeth</i> 2020;20(11):7-17.</p> <p>Rajkomar A et al. Ensuring Fairness in Machine Learning to Advance Health</p>

		Equity. Ann Intern Med 2018;169(12):866-872. Chen IY et al. Ethical Machine Learning in Healthcare. Annual Review of Biomedical Data Science 2021 4:1, 123- 144
26/10	Aula 20 - Monitoria.	
31/10	Aula 21 – Apresentação dos trabalhos.	