

## PLANO DE ENSINO

<b>IDENTIFICAÇÃO</b>		
<b>HABILITAÇÃO:</b>	Profissionalizante em Machine Learning (Aprendizagem de Máquina)	<b>C. H. TOTAL:</b> 400 horas
<b>PROFESSOR (A):</b>	Gerson do Nascimento Silva (PhD Student, UnB)	

<b>PLANEJAMENTO</b>
<b>BASES TECNOLÓGICAS</b>
<p><b>Ementa:</b> Desenvolvimento de código-fonte (<i>script</i>) utilizando linguagem <i>python</i> para implementação, mostrando como:</p> <ol style="list-style-type: none"> <li>1. <b>Entendimento do Ecossistema de Aprendizagem de Máquina:</b> <ol style="list-style-type: none"> <li>a) Preparação do Ecossistema de Aprendizagem de Máquina (ML);</li> <li>b) Começando com o Python;</li> <li>c) Começando com a biblioteca <i>p/</i> manipulação de dados - Pandas;</li> <li>d) Começando com a biblioteca <i>p/</i> arranjos multidimensionais - NumPy;</li> <li>e) Começando com a biblioteca <i>p/</i> criação de gráficos e visualizações de dados em geral - Matplotlib;</li> <li>f) Estatística descritiva - Entendendo os dados;</li> <li>g) Preparação de dados;</li> <li>h) Seleção de Características (RFE - Recursive Feature Elimination);</li> <li>i) Reamostragem;</li> <li>j) Métricas de Desempenho (performance);</li> <li>k) Algoritmos de Classificação;</li> <li>l) Algoritmos de Regressão;</li> <li>m) Comparando Algoritmos de Aprendizagem de Máquina;</li> <li>n) Fluxos de trabalho (workflows) com "pipelines" (processo contínuo);</li> <li>o) Uso de "Ensembles" (conjunto de algoritmos);</li> <li>p) Ajuste de "Hiper Parâmetros";</li> <li>q) Como "Salvar" e/ou "Carregar";</li> </ol> </li> <li>2. <b>Uso da biblioteca "Seaborn", PCA, ANOVA, Curva ROC, Hiper-parâmetros e Modelos de Regressão</b> <ol style="list-style-type: none"> <li>a) Reindexar <i>dataframes</i>;</li> <li>b) Substituir múltiplos valores usando biblioteca computacional <i>Pandas</i>;</li> <li>c) Plotar <i>dataframes</i> com biblioteca computacional <i>Seaborn</i>;</li> <li>d) Limpeza, estruturação e enriquecimento de dados brutos – "Data Wrangling";</li> <li>e) Plotar gráfico de barras usando <i>dataframes Pandas</i>;</li> <li>f) Utilizar séries temporais em <i>Pandas</i>;</li> <li>g) Gerar séries temporais com <i>Pandas</i> e <i>Seaborn</i>;</li> <li>h) Gerar "grouped bar";</li> <li>i) Determinar e plotar coeficientes de correlação;</li> <li>j) Reduzir dimensionalidade de matriz esparsa;</li> <li>k) Usar PCA (<i>Principal Component Analysis</i>) para reduzir dimensionalidade;</li> <li>l) Extrair características usando PCA;</li> <li>m) Usar "ANOVA <i>F-values</i>" para selecionar características;</li> <li>n) Usar "<i>Chi squared</i>" para selecionar características (<i>features</i>);</li> <li>o) Remover características altamente correlacionadas;</li> <li>p) Verificar acurácia de um modelo com validação cruzada;</li> <li>q) Verificar "pontuação AUC" de um modelo;</li> <li>r) Plotar curva de aprendizado de máquina;</li> <li>s) Plotar "curva ROC";</li> <li>t) Usar o algoritmo "<i>Random Forest</i>";</li> <li>u) Sintonizar hiper-parâmetros com algoritmo "<i>GridSearchCV</i>" e "<i>Random Search</i>";</li> <li>v) Otimizar hiper-parâmetros de um modelo de regressão logística;</li> <li>w) Otimizar hiper-parâmetros de um modelo de árvore de decisão;</li> <li>x) Criar e otimizar "baseline" de um modelo de regressão linear;</li> <li>y) Criar e otimizar "baseline" de um modelo de regressão ridge;</li> <li>z) Criar e otimizar "baseline" de um modelo de regressão lasso e "ElasticNet".</li> </ol> </li> <li>3. <b>Classificação, Clusterização e Regressão</b> <ol style="list-style-type: none"> <li>a) Criar e otimizar modelo para regressão e classificação;</li> <li>b) Utilizar "<i>nearest neighbours</i>" para regressão e classificação;</li> <li>c) Fazer agrupamento aglomerativo (<i>Agglomerative Clustering</i>);</li> <li>d) Fazer clusterização com o "<i>Kmeans</i>";</li> </ol> </li> </ol>



- e) Fazer clusterização baseado em afinidade;
- f) Utilizar "DBSCAN Clustering";
- g) Utilizar a abordagem do deslocamento médio (MinShift);
- h) Utilizar a árvore de classificação e regressão;
- i) Utilizar "AdaBoost";
- j) Utilizar "RandomForest";
- k) Utilizar "GradientBoosting";
- l) Utilizar classificador e regressor multicamadas – "MLP - Multi Layer Perceptron";
- m) Utilizar classificador e regressor de reforço gradual "XgBoost";
- n) Utilizar classificador e regressor "CatBoost";
- o) Utilizar classificador e regressor "LightGBM";
- p) Utilizar classificador e regressor "SVM";
- q) Classificar com modelos lineares – "Multiclass Classification";
- r) Classificar com modelos lineares – "Naive Bayes";
- s) Classificar com modelos lineares – "Nearest Neighbors";
- t) Classificar com modelos lineares – "LDA e QDA";
- u) Classificar com modelos lineares – "Tree Model";
- v) Classificar com modelos lineares – "Ensemble Bagging Model";
- w) Classificar com modelos lineares – "Ensemble Boosting Model";
- x) Utilizar métrica de classificação e regressão;
- y) Comparar algoritmos de classificação;
- z) Implementar "Ensemble Model";
- aa) Salvar modelos treinados;
- bb) Avaliar modelos com curvas de aprendizagem;
- cc) Paralelizar execução e validação cruzada no "XGBoost";
- dd) Otimizar número de árvores no "XGBoost".

**4. Visualização de dados: Plotagem de gráficos**

- a) Autocorrelação (ACF) e Autocorrelação Parcial (PACF);
- b) Pizza com destacamento;
- c) Plotagem de textos;
- d) Divergências de escala;
- e) Densidade;
- f) Series temporais múltiplas com escalas;
- g) Boxplot;
- h) Correlograma – "Correlogram";
- i) Curvas de densidade – "Cross Correlation";
- j) Decomposição de serie temporal;
- k) Dispersão com linha de regressão linear de melhor ajuste;
- l) Área;
- m) Barras;
- n) Bolha – "Bubble";
- o) Cascata- "Waterfall";
- p) Lotes em par – "Pairwise";
- q) Histograma empilhado;
- r) Variáveis contínuas (Histograma);
- s) Histogramas marginais – "Marginal Boxplot";
- t) Quadro de marcadores – "Lollipop";
- u) Pirâmide populacional;
- v) Diagramas de dispersão – "Scatter";
- w) Anotações de picos e vales em series temporais;
- x) Cilindros – "Cylinder";
- y) Feixes sólidos – "Joy";
- z) Dado sazonal.

TEORIA	PRÁTICA
<p><b>-video aula:</b>  <b>(playlist do Dr.Rafael Izbicki – UFSCar)</b></p> <ul style="list-style-type: none"> <li>• Aprendizado de Máquina: uma abordagem estatística;</li> </ul> <p><b>(playlist do Msc. Eduardo Ferreira – UNICAMP)</b></p> <ul style="list-style-type: none"> <li>• Machine Learning;</li> </ul> <p><b>(playlist M.Sc Afrânio Melo – COPPE-UFRJ)</b></p> <ul style="list-style-type: none"> <li>• Data Science e Machine Learning;</li> </ul> <p><b>(playlist do PhD Francisco A. Rodrigues – USP)</b></p> <ul style="list-style-type: none"> <li>• "Estatísticas";</li> <li>• "Probabilidades";</li> <li>• "Redes Complexas";</li> </ul>	<p><b>- Entendimento do Ecossistema de Aprendizagem de Máquina:</b>  passo importante no processo de trabalhar dados. A frase "garbage in, garbage out" é particularmente aplicável aos projetos da trilha de IA (machine learning , data mining e data science). Os métodos de coleta de dados geralmente são pouco controlados, resultando em valores de intervalo out-of-range (por exemplo, renda: -100), combinações de dados impossíveis (por exemplo, sexo: masculino, grávidas: sim), missing values, etc. Tal ecossistema é um ambiente de análise de dados que automatiza a construção de modelos analíticos. É um ramo da inteligência artificial baseado na ideia de que sistemas podem aprender com dados, identificar padrões e tomar decisões com o mínimo de intervenção humana.</p>

<ul style="list-style-type: none"> <li>• “Processos Estocásticos”;</li> </ul> <p><b>-video aula: (playlist do Dr. Alexandre L. M. Levada – USP)</b></p> <ul style="list-style-type: none"> <li>• “Introdução a Teoria dos Grafos”;</li> <li>• “Reconhecimento de Padrões”;</li> </ul> <p><b>-video aula: (playlist da Dra. Cibele Russo – USP)</b></p> <ul style="list-style-type: none"> <li>• “Visualização e Exploração de Dados”;</li> <li>• “Análise Multivariada e Aprendizado Não-Supervisionado”.</li> </ul>	<p><b>- Análise de dados:</b> é um processo de inspeção, limpeza, transformação e modelagem de dados com o objetivo de descobrir informações úteis, informar conclusões e apoiar a tomada de decisões. A análise de dados tem múltiplas facetas e abordagens, abrangendo diversas técnicas sob uma variedade de nomes, e é usada em diferentes domínios dos negócios, ciências e ciências sociais. Desempenha um papel tornando a tomada de decisões mais científicas e ajudando no processo de operar com mais eficácia.</p> <p><b>- Classificação:</b> baseia-se em prever a categoria de uma observação dada. Procura-se estimar um “classificador” que gere como saída a classificação qualitativa de um dado não observado com base em dados de entrada (que abrangem observações com classificações já definidas). <u>Exemplo:</u> um classificador que utilize dados não observados de um paciente e classifique-o como doente ou não-doente.</p> <p><b>- Regressão:</b> de forma similar a classificação, utiliza dados de entrada (preditores) já observados para prever uma resposta. A grande diferença é que, neste caso, procura-se estimar um valor numérico e não uma classificação de uma observação. <u>Exemplo:</u> estimar um modelo que utilize a idade e os anos de escolaridade de um indivíduo não-observado anteriormente para tentar prever seu salário. Utiliza-se como base desse modelo: idades, anos de escolaridades e salários de diversos indivíduos já observados anteriormente.</p> <p><b>- Agrupamento:</b> também conhecido como “<i>Clustering</i>”, tem como objetivo agrupar observações em grupos conhecidos como “<i>clusters</i>”. Essas observações apresentam similaridades dentro de seu <i>cluster</i> e diferenças em relação aos demais <i>clusters</i> formados. Diferente da Classificação, não é realizada a rotulação dos <i>clusters</i>, fazendo com que não exista uma clusterização errada ou certa. A clusterização utilizada resulta em diferentes tipos de <i>clusters</i>, e a escolha dessas técnicas deve ser previamente analisada pelo pesquisador. <u>Exemplo:</u> agrupar fotos de animais similares em <i>clusters</i>, sem ter o conhecimento prévio de qual animal está sendo apresentado.</p> <p><b>- Visualização:</b> Tornar a apresentação dos dados atraente e de fácil entendimento; Identificar tendências; Perceber situações atípicas em um conjunto de dados; Contar uma história encontrada nos dados; Reforçar um argumento ou opinião; Destacar um ponto importante em um conjunto de dados.</p>
--	--

<b>ANÁLISE DA REALIDADE</b>
Pré-Requisito
<ol style="list-style-type: none"> <li>1. Lógica de Programação;</li> <li>2. Programação Orientada a Objetos;</li> <li>3. Análise de Sistemas Orientada a Objeto;</li> <li>4. Fundamentos e Modelagem de banco de Dados.</li> </ol>
Necessidade da turma:
- Programação intermediário em <i>python</i> ou outra linguagem de programação que traga subsídios para entendimento dos conceitos computacionais.

<b>PROJEÇÃO DE FINALIDADES</b>
--------------------------------

<b>Objetivos Gerais (Competências)</b>
Compreender os conceitos de IA e seus recursos e capacidades para implementar código-fonte reproduzível.
<b>Objetivos Específicos (Habilidades)</b>
<ul style="list-style-type: none"> <li>-Capacidade para implementar código-fonte (uso de <i>python</i> ou outras tecnologias com mesma finalidade);</li> <li>-Capacidade para implementar API'S para o pré-processamento de dados;</li> <li>-Capacidade para implementar API'S para análise de dados;</li> <li>-Integração do uso de API'S em manipulação de banco de dados;</li> <li>-Entender como aplicar <i>frameworks python</i> focados em – Classificação;</li> <li>-Entender como aplicar <i>frameworks python</i> focados em – Clusterização;</li> <li>-Entender como aplicar <i>frameworks python</i> focados em – Regressão;</li> <li>-Criar, manipular e ajustar modelos de aprendizagem de máquina;</li> <li>-Criar métricas de desempenho e fluxos de trabalho em aprendizagem de máquina;</li> <li>-Utilizar conjunto de algoritmos (<i>Ensembles</i>) em aprendizagem de máquina</li> <li>-Ser capaz de implementar código-fonte reproduzível utilizando conceitos anteriores.</li> </ul>

<b>FORMAS DE MEDIAÇÃO</b>
<b>Procedimentos Metodológicos</b>
<ul style="list-style-type: none"> <li>-Elaborar modelos de implementação de código-fonte que valide as teorias, com uso de exemplos;</li> <li>-Propor trabalhos práticos sobre o assunto;</li> <li>-Corrigir os trabalhos práticos;</li> <li>-Elaborar avaliação do conhecimento sobre o assunto ministrado;</li> <li>-Uso de computadores para implementar exemplos;</li> </ul>

<b>AVALIAÇÃO*</b>	
<b>Instrumento(s)</b>	<b>Data</b>
<ul style="list-style-type: none"> <li>• Avaliação única por meio do instrumento <i>Quiz</i>;</li> <li>• Criação de um projeto final “<i>End to End</i>” implementado em código-fonte reproduzível <i>python</i>.</li> </ul>	Ao fim do módulo estudado.

\* No Exame Quiz não é permitido qualquer espécie de acréscimo de nota (trabalhos, participações e outros) estranhos à avaliação.

\* Quiz é um jogo mental no qual os jogadores tentam responder corretamente a questões que lhes são colocadas. A palavra também é utilizada como sinônimo de avaliação de aquisição de conhecimentos ou capacidades em ambientes de aprendizagem.

<b>BIBLIOGRAFIA</b>					
<b>BIBLIOGRAFIA BÁSICA</b> (Títulos, periódicos, etc)					
Título/Periódico	Autor	Edição	Editora	Ano	Livro Texto
Álgebra Linear com Python	Rafael F. V. C. Santos	1ª edição	Amazon (e-book)	2018	-
Trilhas Python	Eduardo Pereira	1ª edição	Casa do Código	2018	-
Python para Desenvolvedores	Luiz Eduardo Borges	3ª edição	Novatec	2014	Sim
<b>BIBLIOGRAFIA COMPLEMENTAR</b> (Títulos, periódicos, etc)					
Título/Periódico	Autor	Edição	Editora	Ano	
Curso Intensivo de Python	Eric Matthes	1ª edição	Novatec	2017	
Aprendizado de Máquina para Leigos	John Paul Mueller, Luca Massaron	1ª edição	Alta Books	2019	
Introduction to Applied Linear Algebra	Stephen Boyd	1ª edição	Cambridge University Press	2018	
<b>Outros:</b>	<p><b>Essa é uma inspiração para Machine Learning (sugestão p/ aquisição):</b></p> <ul style="list-style-type: none"> <li>• 1. Machine Learning: Introdução à classificação. Autor(es) - Guilherme Silveira, Bennett Bullock, 2017. Editora: Casa do Código.</li> </ul>				