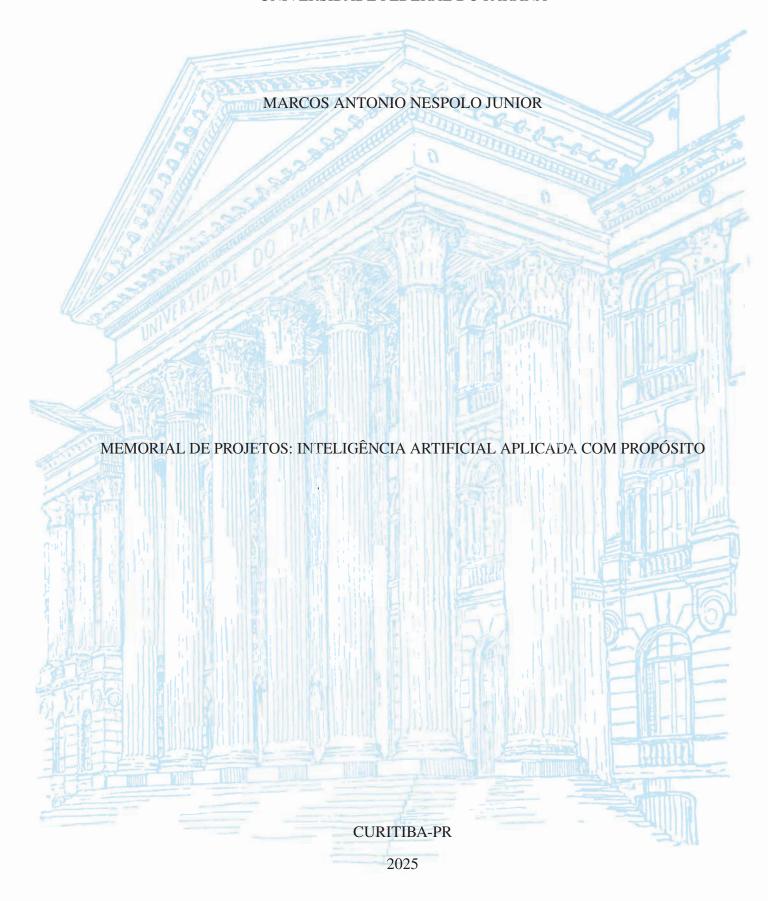
## UNIVERSIDADE FEDERAL DO PARANÁ



# MARCOS ANTONIO NESPOLO JUNIOR MEMORIAL DE PROJETOS: INTELIGÊNCIA ARTIFICIAL APLICADA COM PROPÓSITO Memorial de Projetos apresentado ao curso de Especi-

alização em Inteligência Artificial Aplicada, Setor de Educação Profissional e Tecnológica, Universidade Federal do Paraná, como requisito parcial à obtenção do título de Especialista em Inteligência Artificial Aplicada.

Área de concentração: Inteligência Artificial Aplicada.

Orientador: Prof<sup>a</sup>. Dr<sup>a</sup>. Rafaela Mantovani Fontana.

**CURITIBA-PR** 



MINISTÉRIO DA EDUCAÇÃO
SETOR DE EDUCAÇÃO PROFISSIONAL E TECNOLÓGICA
UNIVERSIDADE FEDERAL DO PARANÁ
PRÓ-REITORIA DE PÓS-GRADUAÇÃO
CURSO DE PÓS-GRADUAÇÃO INTELIGÊNCIA ARTIFICIAL
APLICADA - 40001016399E1

## TERMO DE APROVAÇÃO

Os membros da Banca Examinadora designada pelo Colegiado do Programa de Pós-Graduação Inteligência Artificial Aplicada da Universidade Federal do Paraná foram convocados para realizar a arguição da Monografia de Especialização de MARCOS ANTONIO NESPOLO JUNIOR, intitulada: MEMORIAL DE PROJETOS: INTELIGÊNCIA ARTIFICIAL APLICADA COM PROPÓSITO, que após terem inquirido o aluno e realizada a avaliação do trabalho, são de parecer pela sua <u>aprovação</u> no rito de defesa.

A outorga do título de especialista está sujeita à homologação pelo colegiado, ao atendimento de todas as indicações e correções solicitadas pela banca e ao pleno atendimento das demandas regimentais do Programa de Pós-Graduação.

Curitiba, 04 de Outubro de 2025.

RAFAELA MAN OVANI FONTANA Presidente da Banca Examinadora

Avaliador Interno (UNIVERSIDADE FEDERAL DO PARANÁ)

#### **AGRADECIMENTOS**

Agradeço primeiramente à minha família, em especial aos meus avós, Almerinda e Antonio, que, com amor e sacrifício, dedicaram suas vidas para suavizar o peso da minha e abriram caminhos que sozinho eu não conseguiria trilhar.

Aos meus amigos, que foram apoio essencial nesta caminhada, agradeço pela companhia, pelo incentivo e pelos momentos de descontração que me ajudaram a manter o equilíbrio enquanto conciliava trabalho e duas pós-graduações ao mesmo tempo.

Por fim, agradeço à vida, ao destino — ou a Deus — por ter colocado no meu caminho os desafios que me lapidaram, as oportunidades que me abriram horizontes e as conquistas que deram sentido à jornada, permitindo-me concluir mais esta etapa da minha trajetória.

#### **RESUMO**

O conceito de IA aplicada com propósito consolidou-se como uma resposta ao uso acrítico de modelos de alto desempenho, priorizando decisões sobre quando, como e por quê aplicar IA de forma útil, segura e alinhada a valores. A partir das disciplinas do curso de pós-graduação, foram exploradas diferentes dimensões da IA aplicada, incluindo fundamentos de IA, aprendizado de máquina e deep learning; estatística e arquitetura de dados; big data; programação em Python e R; visualização de dados e storytelling; visão computacional e frameworks; tópicos em IA (lógica fuzzy, algoritmos genéticos e séries temporais); e gestão de projetos de IA. Essas abordagens evidenciaram como a IA aplicada se manifesta tanto no nível técnico, por meio de formulação do problema, preparação e qualidade dos dados, modelagem e validação, quanto no organizacional, com ênfase em métricas alinhadas ao contexto, comunicação com as partes interessadas e acompanhamento de resultados. Além dos aspectos metodológicos, discutiram-se implicações éticas e sociais, como privacidade, viés algorítmico, explicabilidade e transparência. Em síntese, o memorial destaca que IA aplicada com propósito não é apenas construir modelos, mas integrar pessoas, dados, processos e tecnologia para entregar valor verificável, com responsabilidade e impacto real.

**Palavras-chave**: inteligência artificial aplicada; aprendizado de máquina; aprendizado profundo; arquitetura de dados; ética em IA.

#### **ABSTRACT**

Purpose-driven Applied AI emerges as a response to the uncritical use of high-performing models, prioritizing decisions about when, how, and why to apply AI in useful, safe, and value-aligned ways. Drawing on graduate-level coursework, the report covers multiple dimensions of applied AI: foundations in AI, machine learning, and deep learning; statistics and data architecture; big data; programming in Python and R; data visualization and storytelling; computer vision and frameworks; advanced topics (fuzzy logic, genetic algorithms, time series); and AI project management. These tracks show how applied AI operates both technically—through problem framing, data preparation and quality, modeling, and validation—and organizationally, emphasizing context-aligned metrics, stakeholder communication, and results follow-up. Beyond methodology, the study examines ethical and social implications, including privacy, algorithmic bias, explainability, and transparency. In sum, purpose-driven applied AI is not merely about building models, but about integrating people, data, processes, and technology to deliver verifiable value with responsibility and real impact.

**Keywords:** applied artificial intelligence; machine learning; deep learning; data architecture; AI ethics.

# LISTA DE FIGURAS

1.1	Exercício 2: Busca Heurística - Melhor Caminho para Bucharest 16
1.2	Exercício 4: Rede Neural Artificial
1.3	Resolução Exercício 2: Caminho de Lugoj até Bucharest
2.1	Resolução Exercício 2: Distribuição da Quantidade de Carros por Marca 25
2.2	Resolução Exercício 2: Distribuição da Quantidade de Carros por Tipo de Engrenagem
2.3	Resolução Exercício 2: Evolução da Média de Preço dos Carros ao Longo de 2022 27
2.4	Resolução Exercício 2: Média de Preços dos Carros por Marca e Tipo de Engrenagem Agrupados por Tipo
2.5	Resolução Exercício 2: Média de Preços dos Carros por Marca e Tipo de Engrenagem Agrupados por Marca
2.6	Resolução Exercício 2: Média de Preços dos Carros por Marca e Tipo de Combustível Agrupados por Marca
2.7	Resolução Exercício 2: Média de Preços dos Carros por Marca e Tipo de Combustível Agrupados por Tipo
2.8	Resolução Exercício 3: Boxplot da Distribuição dos Preços Médios
2.9	Resolução Exercício 3: Mapa de Correlação entre as Variáveis Numéricas 33
4.1	Exercício 1: Histograma e Box-plot
4.2	Exercício 3: Distribuição Normal e Histograma
4.3	Exercício 3: Box-plots
7.1	Admissão - Gráfico de resíduos (melhor modelo)
7.2	Biomassa - Gráfico de resíduos (melhor modelo)
8.1	Classificação de Imagens com CNN: Gráfico de Perda
8.2	Classificação de Imagens com CNN: Gráfico de Acurácia
8.3	Classificação de Imagens com CNN: Matriz de Confusão
8.4	Classificação de Imagens com CNN: Resultado da Predição 81
8.5	Detector de SPAM com RNN: Gráfico de Perda
8.6	Detector de SPAM com RNN: Gráfico de Acurácia
8.7	Gerando dados com GAN: Geração de Imagem Ruidosa
8.8	Gerando dados com GAN: Checkpoint da Geração de Dígitos
8.9	Tranformer: Codificação Posicional
9.1	Visão geral das Ferramentas Utilizadas pelo AirBnb
13.1	Exercício 1: Gráfico de Perda
13.2	Exercício 1: Gráfico de Acurácia
13.3	Exercício 1: Matriz de Confusão

13.4	Exercício 1: Apresentando Classificação Incorreta
13.5	Exercício 2: Gráfico de Perda
13.6	Exercício 2: Gráfico de RMSE
13.7	Exercício 2: Gráfico de R2
13.8	Exercício 3: Gráfico de Perda
13.9	Exercício 3: Imagem Baixada da URL Fornecida
13.10	Exercício 3: Resultado Obtido no Processo de Deep Dream
13.11	Exercício 3: Resultado Obtido no Processo de Deep Dream à Oitava 146
14.1	Análise de Tendências de Busca por Passagem de avião
15.1	Exercício 1: Rota Inicial Aleatória, Distância = 5159,51
15.2	Exercício 1: Rota Após Evolução, Distância = 1154,43
15.3	Exercício 2: Projeção PCA 3D
15.4	Exercício 2: Projeção PCA 2D
15.5	Exercício 2: Dendograma - Distâncias Originais
15.6	Exercício 2: Dendograma - Distâncias PCA 2D

## LISTA DE TABELAS

1.1	Valores de $h_{DLR}$ — distâncias em linha reta para Bucareste	17
2.1	Exercício 2 de LPA: Metados da Base de Dados Carros Brasil	21
4.1	Tabela de frequência das idades (maridos)	48
4.2	Tabela de frequência das idades (esposas)	48
5.1	Estrutura do data.frame (N = 2574, 16 variáveis)	53
5.2	MAE por modelo (10 reamostragens)	54
5.3	RMSE por modelo (10 reamostragens)	54
5.4	$R^2$ por modelo (10 reamostragens)	54
5.5	Comparativo de modelos	55
6.1	Matriz de confusão - Treino (acurácia: 84,06%)	60
6.2	Relatório de classificação — Treino	61
6.3	Matriz de confusão - Amostra (acurácia: 70%)	61
6.4	Relatório de classificação - Amostra	61
7.1	Veículos - Resultados por técnica com parâmetros, acurácia e matrizes de confusão.	62
7.2	Diabetes - Resultados por técnica com parâmetros, acurácia e matriz de confusão.	63
7.3	Admissão - Resultados por técnica com parâmetros e métricas ( $R^2$ , $S_{yx}$ , coeficiente de Pearson, RMSE e MAE)	64
7.4	Amostra com GRE, TOEFL, rating universitário e métricas preditivas	64
7.5	Biomassa - Resultados por técnica com parâmetros e métricas ( $R^2$ , $S_{yx}$ , coeficiente de Pearson, RMSE e MAE)	65
7.6	Amostra com variáveis <i>dap</i> , <i>h</i> , <i>Me</i> e predição da RNA	65
7.7	K-means (médias por cluster) - Parte 1	66
7.8	K-means (médias por cluster) - Parte 2	66
7.9	Regras de associação — top 10 linhas	66

# SUMÁRIO

1	INTELIGÊNCIA ARTIFICIAL APLICADA
1.1	DEFINIÇÃO
1.2	INTELIGÊNCIA ARTIFICIAL APLICADA COM PROPÓSITO
	REFERÊNCIAS
	APÊNDICE 1 – INTRODUÇÃO À INTELIGÊNCIA ARTIFICIAL 16
	APÊNDICE 2 – LINGUAGEM DE PROGRAMAÇÃO APLICADA 21
	APÊNDICE 3 – LINGUAGEM R
	APÊNDICE 4 – ESTATÍSTICA APLICADA I
	APÊNDICE 5 – ESTATÍSTICA APLICADA II
	APÊNDICE 6 – ARQUITETURA DE DADOS
	APÊNDICE 7 – APRENDIZADO DE MÁQUINA 62
	APÊNDICE 8 – DEEP LEARNING
	<b>APÊNDICE 9 – BIG DATA</b>
	APÊNDICE 10 - VISÃO COMPUTACIONAL 104
	APÊNDICE 11 – ASPECTOS FILOSÓFICOS E ÉTICOS DA IA 126
	APÊNDICE 12 – GESTÃO DE PROJETOS DE IA
	APÊNDICE 13 – FRAMEWORKS DE INTELIGÊNCIA ARTIFICIAL 133
	APÊNDICE 14 – VISUALIZAÇÃO DE DADOS E STORYTELLING 147
	APÊNDICE 15 – TÓPICOS EM INTELIGÊNCIA ARTIFICIAL 149

#### 1 INTELIGÊNCIA ARTIFICIAL APLICADA

A IA Aplicada representa o esforço de transformar métodos computacionais em soluções de valor para problemas concretos. Ao longo do curso, essa perspectiva foi trabalhada tanto em discussões conceituais quanto em atividades práticas, buscando articular teoria, experimentação e impacto real.

### 1.1 DEFINIÇÃO

A Inteligência Artificial Aplicada refere-se à utilização sistemática de métodos de IA para resolver problemas reais em contextos específicos. Seu foco está na conexão entre modelos computacionais, dados e processos organizacionais, com o objetivo de gerar impacto tangível em produtos, serviços e decisões. Nesse sentido, a ênfase recai sobre a transformação de capacidades algorítmicas em soluções operacionais de valor mensurável, seja no setor privado ou em políticas públicas (Russell e Norvig, 2022; Provost e Fawcett, 2013; Kuhn e Johnson, 2013).

Do ponto de vista técnico, essa abordagem envolve diversas estratégias de aprendizado de máquina, incluindo aprendizado supervisionado, não supervisionado e por reforço (Mitchell, 1997). Também são incorporadas arquiteturas modernas de aprendizado profundo, especialmente em aplicações como visão computacional, processamento de linguagem natural e análise de séries temporais (Bishop, 2006; Goodfellow, Bengio e Courville, 2016). Nessas situações, a modelagem não é um fim em si, mas um meio para atingir objetivos como previsão, classificação, recomendação, detecção de anomalias, otimização e suporte à decisão.

O ciclo de vida da IA aplicada é iterativo e orientado à operação. Ele compreende as etapas de descoberta e preparação dos dados, modelagem e validação, implantação e monitoramento contínuo. Esse fluxo exige práticas de engenharia robustas, capazes de lidar com desafios como acoplamentos entre componentes, mudanças nos dados (*data drift*), e dependências de infraestrutura. Esses fatores compõem o que se denomina como *dívida técnica* em sistemas de aprendizado de máquina (Sculley *et al.*, 2015). Além disso, a atuação colaborativa entre equipes de produto, dados e engenharia é essencial para garantir a qualidade e a continuidade das soluções desenvolvidas (Amershi *et al.*, 2019).

A avaliação dos sistemas de IA deve considerar tanto critérios estatísticos (como acurácia, erro, robustez e capacidade de generalização) quanto métricas de impacto organizacional, como

valor entregue, custo de operação, riscos envolvidos e conformidade regulatória. A comparação entre modelos deve ser conduzida com rigor metodológico e validação compatível com o problema em questão (Kuhn e Johnson, 2013; Bishop, 2006). Ao mesmo tempo, é importante reconhecer os limites dos paradigmas adotados e os vieses que podem emergir durante o processo de modelagem (Breiman, 2001).

Finalmente, a aplicação responsável da IA exige atenção a aspectos éticos e de governança. Entre os principais temas estão a privacidade e segurança dos dados, o viés algorítmico, a explicabilidade dos modelos e a documentação adequada de todo o ciclo de vida. Boas práticas incluem o uso de *model cards* e *datasheets* para promover a transparência, bem como técnicas de interpretação local que possibilitam a análise e compreensão das decisões automatizadas (Mitchell *et al.*, 2019; Gebru *et al.*, 2021; Ribeiro, Singh e Guestrin, 2016; Barocas e Selbst, 2016). Em suma, a Inteligência Artificial Aplicada busca articular pessoas, dados, processos e tecnologia para entregar soluções eficazes, seguras e sustentáveis.

#### 1.2 INTELIGÊNCIA ARTIFICIAL APLICADA COM PROPÓSITO

Ao longo do curso de pós-graduação, enfatizou-se a importância de compreender **quando**, **como** e **por quê** aplicar inteligência artificial. Antes de qualquer linha de código, trabalhou-se a formulação do problema e a hipótese de valor: identificar se a aplicação de IA é pertinente, quais ganhos são esperados, quais custos e restrições estão envolvidos e quais partes interessadas serão impactadas. A partir dessa análise inicial, é possível estabelecer um ciclo consistente de desenvolvimento: preparo dos dados, construção e validação do modelo, entrega e monitoramento dos resultados, sempre com critérios explícitos de sucesso e risco.

No que se refere ao **como** aplicar, os conteúdos das disciplinas de Introdução à IA, Aprendizado de Máquina e *Deep Learning* consolidaram tarefas e técnicas fundamentais, como regressão, classificação, agrupamento, associação e aprendizado por reforço, além de arquiteturas contemporâneas voltadas a visão computacional, linguagem natural e séries temporais (Russell e Norvig, 2022; Mitchell, 1997; Goodfellow, Bengio e Courville, 2016; Bishop, 2006). A ênfase recaiu sobre o alinhamento entre objetivo, métrica e contexto: definir *baselines*, escolher indicadores pertinentes ao problema e interpretar os resultados de forma a subsidiar decisões, e não apenas buscar métricas de desempenho isoladas (Kuhn e Johnson, 2013; Provost e Fawcett, 2013).

A dimensão dos dados estruturou a qualidade da aplicação: técnicas de limpeza, codificação, imputação, construção e seleção de atributos, redução de correlações e mitigação de ruído foram integradas à análise estatística descritiva, testes paramétricos e não paramétricos, análise de viés/variância e comparação de modelos (Kuhn e Johnson, 2013; Bishop, 2006). A disciplina de *Big Data* ampliou a perspectiva para cenários de alto volume, variedade e velocidade, abordando estratégias de armazenamento, processamento e integração com *data lakes*, essenciais quando escala e latência impõem restrições técnicas relevantes.

Para assegurar que o **por quê** permaneça claro e comunicável, disciplinas como Visualização de Dados e *Storytelling* foram centrais na transformação de evidências em narrativas acionáveis, voltadas tanto ao público técnico quanto ao executivo. O uso de Python e R permitiu prototipagem ágil e experimentos reprodutíveis. A disciplina de Visão Computacional abordou desde aquisição e pré-processamento de imagens até *transfer learning* e *fine-tuning* com redes convolucionais. Já os *Frameworks* de IA organizaram o uso de bibliotecas consolidadas para tarefas como regressão, classificação, séries temporais, CNNs, RNNs e modelos de linguagem. Tópicos em IA ampliou o repertório com lógica *fuzzy*, algoritmos genéticos e métodos especializados para séries temporais. Por sua vez, Gestão de Projetos de IA conectou os aspectos técnicos à organização, abordando *workflows*, escopo, riscos, critérios de aceite, documentação e a dinâmica entre equipes, encurtando a distância entre protótipo e solução real.

O componente ético e social reforçou que aplicar IA de forma responsável é indissociável das perguntas sobre **quando**, **como** e **por quê**. Foram discutidos temas como privacidade, segurança de dados, viés algorítmico, explicabilidade, autonomia e responsabilização. Instrumentos como *model cards* e *datasheets*, aliados a técnicas de interpretação local, viabilizam maior transparência e inspeção de decisões (Barocas e Selbst, 2016; Mitchell *et al.*, 2019; Gebru *et al.*, 2021; Ribeiro, Singh e Guestrin, 2016). Em aplicações com grande impacto, como visão computacional, linguagem natural e análise em larga escala, é fundamental equilibrar os benefícios de eficiência e suporte à decisão com salvaguardas éticas e comunicação clara de limitações e incertezas a clientes, usuários e stakeholders.

Em síntese, os projetos realizados ao longo do curso formaram uma base sólida para aplicar IA com propósito. Um profissional capacitado deve ser capaz de responder, com criticidade e responsabilidade, às três questões fundamentais: **quando** aplicar IA (viabilidade e valor), **como** aplicá-la (dados, métodos e validação) e **por quê** aplicá-la (impacto e responsabilidade). Mais

do que construir modelos de alto desempenho, trata-se de decidir e agir de forma útil, segura e alinhada aos valores da organização e da sociedade.

#### REFERÊNCIAS

AMERSHI, S. et al. Software engineering for machine learning: A case study. In: *Proceedings of the 41st International Conference on Software Engineering: Software Engineering in Practice (ICSE-SEIP)*. [S.l.: s.n.], 2019. p. 291–300.

BAROCAS, S.; SELBST, A. D. Big data's disparate impact. *California Law Review*, v. 104, n. 3, p. 671–732, 2016.

BISHOP, C. M. Pattern Recognition and Machine Learning. [S.l.]: Springer, 2006.

BREIMAN, L. Statistical modeling: The two cultures. *Statistical Science*, v. 16, n. 3, p. 199–231, 2001.

GEBRU, T. et al. Datasheets for datasets. Communications of the ACM, v. 64, n. 12, p. 86–92, 2021.

GOODFELLOW, I.; BENGIO, Y.; COURVILLE, A. Deep Learning. [S.l.]: MIT Press, 2016.

KUHN, M.; JOHNSON, K. Applied Predictive Modeling. [S.l.]: Springer, 2013.

MITCHELL, M. et al. Model cards for model reporting. In: *Proceedings of the Conference on Fairness, Accountability, and Transparency (FAT\*)*. [S.l.: s.n.], 2019. p. 220–229.

MITCHELL, T. M. Machine Learning. [S.l.]: McGraw-Hill, 1997.

PROVOST, F.; FAWCETT, T. Data Science for Business. [S.l.]: O'Reilly Media, 2013.

RIBEIRO, M. T.; SINGH, S.; GUESTRIN, C. "why should i trust you?": Explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*. [S.l.: s.n.], 2016. p. 1135–1144.

RUSSELL, S.; NORVIG, P. Artificial Intelligence: A Modern Approach, 4th US ed. [S.l.]: Pearson, 2022.

SCULLEY, D. *et al.* Hidden technical debt in machine learning systems. In: *Advances in Neural Information Processing Systems (NeurIPS)*. [S.l.: s.n.], 2015. p. 2503–2511.

# APÊNDICE 1 - INTRODUÇÃO À INTELIGÊNCIA ARTIFICIAL

#### A - ENUNCIADO 1 ChatGPT

- a) (6,25 pontos) Pergunte ao ChatGPT o que é Inteligência Artificial e cole aqui o resultado.
- b) **(6,25 pontos)** Dada essa resposta do ChatGPT, classifique usando as 4 abordagens vistas em sala. Explique o porquê.
- c) (6,25 pontos) Pesquise sobre o funcionamento do ChatGPT (sem perguntar ao próprio ChatGPT) e escreva um texto contendo no máximo 5 parágrafos. Cite as referências.
- d) **(6,25 pontos)** Entendendo o que é o ChatGPT, classifique o próprio ChatGPT usando as 4 abordagens vistas em sala. Explique o porquê.

#### 2 Busca Heurística

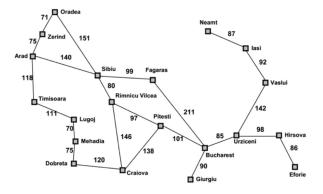
Realize uma busca utilizando o algoritmo  $A^*$  para encontrar o melhor caminho para chegar a **Bucharest** partindo de **Lugoj**. Construa a árvore de busca criada pela execução do algoritmo apresentando os valores de f(n), g(n) e h(n) para cada nó. Utilize a heurística de distância em linha reta, que pode ser observada na tabela abaixo.

Essa tarefa pode ser feita em uma **ferramenta de desenho**, ou até mesmo no **papel**, desde que seja digitalizada (foto) e convertida para PDF.

a) (25 pontos) Apresente a árvore final, contendo os valores, da mesma forma que foi apresentado na disciplina e nas práticas. Use o formato de árvore, não será permitido um formato em blocos, planilha, ou qualquer outra representação.

# NÃO É NECESSÁRIO IMPLEMENTAR O ALGORITMO.

Figura 1.1 – EXERCÍCIO 2: BUSCA HEURÍSTICA - MELHOR CAMINHO PARA BUCHAREST



Arad 366 Mehadia 241 234 Bucareste 0 Neamt Craiova 160 Oradea 380 Drobeta Pitesti 242 100 Eforie 161 Rimnicu Vilcea 193 253 Fagaras 176 Sibiu Giurgiu 77 Timisoara 329

Urziceni

Vaslui

Zerind

80

199

374

151

226

244

Tabela 1.1: Valores de  $h_{DLR}$  — distâncias em linha reta para Bucareste.

#### 3 Lógica

Verificar se o argumento lógico é válido:

Se as uvas caem, então a raposa as come

Hirsova

Iasi

Lugoj

Se a raposa as come, então estão maduras

As uvas estão verdes ou caem

Logo

A raposa come as uvas se e somente se as uvas caem

Deve ser apresentada uma prova, no mesmo formato mostrado nos conteúdos de aula e nas práticas.

#### Dicas:

1. Transformar as afirmações para lógica:

p: as uvas caem

q: a raposa come as uvas

r: as uvas estão maduras

2. Transformar as três primeiras sentenças para formar a base de conhecimento:

 $R1: p \rightarrow q$ 

 $R2: q \rightarrow r$ 

 $R3: \neg r \lor p$ 

3. Aplicar equivalências e regras de inferência para se obter o resultado esperado. Isto é, com essas três primeiras sentenças devemos derivar  $q \leftrightarrow p$ . Cuidado com a ordem em que as fórmulas são geradas.

**Equivalência Implicação:**  $(\alpha \to \beta)$  equivale a  $(\neg \alpha \lor \beta)$ 

**Silogismo Hipotético:**  $\alpha \rightarrow \beta, \beta \rightarrow \gamma \vdash \alpha \rightarrow \gamma$ 

**Conjunção:**  $\alpha, \beta \vdash \gamma \land \beta$ 

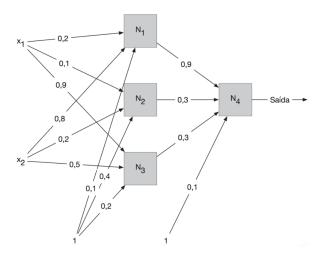
**Equivalencia Bicondicional:**  $(\alpha \leftrightarrow \beta)$  equivale a  $(\alpha \to \beta) \land (\beta \to \alpha)$ 

a) (25 pontos) Deve-se mostrar todos os passos e regras aplicadas, no mesmo formato apresentado nas aulas e nas práticas. As equivalências e regras necessárias estão descritas acima e no material.

#### 4 Redes Neurais Artificiais

Seja a RNA da figura abaixo:

Figura 1.2 – EXERCÍCIO 4: REDE NEURAL ARTIFICIAL



Os neurônios  $N_1$ ,  $N_2$  e  $N_3$  possuem função de ativação linear. Já N4 possui função de ativação tangente hiperbólica (pesquise a fórmula e aplique). Dada a entrada  $x_1 = -3$ ,  $x_2 = 1$ , dê:

- a) (6,25 pontos) Valor de saída do neurônio  $N_1$
- b) (6,25 pontos) Valor de saída do neurônio  $N_2$
- c) (6,25 pontos) Valor de saída do neurônio  $N_3$
- d) (6,25 pontos) Valor de saída da rede como um todo

# **B - RESOLUÇÃO**

#### 1 ChatGPT

- a) Inteligência Artificial (IA) refere-se à capacidade de uma máquina ou programa de computador realizar tarefas que normalmente exigiria inteligência humana. Essas tarefas incluem aprendizado, raciocínio, percepção, reconhecimento de padrões, compreensão da linguagem natural e tomada de decisões. O objetivo da Inteligência Artificial é criar sistemas que possam executar funções complexas de forma autônoma, imitando ou replicando a inteligência humana. Existem dois tipos principais de IA: IA fraca (ou estreita) e IA forte (ou geral).
  - 1. IA fraca (estreita): Refere-se a sistemas especializados em realizar tarefas específicas sem possuir a capacidade de generalizar para diferentes domínios. Exemplos incluem assistentes virtuais, reconhecimento de voz, carros autônomos e sistemas de recomendação.

2. IA forte (geral): Representa um nível mais avançado de inteligência, comparável à inteligência humana em termos de habilidades cognitivas e adaptação a diferentes contextos. Até o momento, a IA forte permanece mais uma aspiração do que uma realidade prática.

As técnicas comuns em IA incluem aprendizado de máquina (*machine learning*), redes neurais, processamento de linguagem natural, visão computacional e algoritmos de otimização. O campo da Inteligência Artificial continua a evoluir rapidamente, com aplicações em diversas áreas, como saúde, finanças, transporte, manufatura e entretenimento.

- b) É categorizável na abordagem "agir como seres humanos", pois no texto o chat GPT ressaltou a imitação de habilidades inerentemente humanas ao tentar definir o conceito de Inteligência Artificial.
- c) Como o ChatGPT funciona: ChatGPT Sigla para *Generative Pré-Trained Transformer* (Transformador pré-treinado generativo, em tradução livre) é um sistema ou um modelo de linguagem baseado em inteligência artificial mais especificamente com as técnicas de aprendizado de máquinas, redes neurais, aprendizado por reforço com *feedback* humano (sigla em inglês RHLF) e processamento de linguagem natural usados com o foco em diálogos virtuais. A seguir vamos falar sobre cada uma das partes do algoritmo GPT. O ChatGPT usa como fonte de dados para manter os seus diálogos virtuais e o contexto das mesmas um conjunto de informações disponíveis e acessíveis principalmente através da internet.

O generativo pré-treinado, o GP de GPT, é uma técnica onde ChatGPT recebe uma grande quantidade de regras de linguagens e dados não rotulados. O GPT é então deixado sem supervisão para aprender sozinho sobre as regras e os relacionamentos que governam as linguagens.

A letra T de GPT, *Transfomer architecture*, refere-se à arquitetura do modelo, a qual é baseada na arquitetura *Transformer*. Agora, todos os dados pré-treinados são usados para criar uma rede neural de aprendizagem profunda, permitindo ao ChatGPT aprender padrões e relacionamentos de textos, dando ao ChatGPT a habilidade de criar respostas como humanos e de predizer qual texto vem em seguida em uma dada sentença. O *Transformer* é capaz de ler cada palavra de uma sentença e comparar cada palavra com as outras palavras da mesma sentença. Isso permite que ele direcione a sua atenção à palavra mais importante da sentença. Esse processo é conhecido como *self-attention*. Com isso, é importante ressaltar que o ChatGPT não entende o que é dito e o que ele responde.

No início, a rede neural do ChatGPT não estava inteiramente segura para ser lançada ao público, pois ela foi treinada e sumariamente através de dados da internet aberta sem nenhuma orientação. Então para que o ChatGPT pudesse dar respostas coerentes, sensatas e seguras ao público, foi otimizado com uma técnica chamada aprendizado por reforço com *feedback* humano (sigla em inglês RHLF). Basicamente são demonstrados dados que guiam a rede neural em como responder adequadamente às solicitações humanas. O RHLF, mesmo não sendo uma aprendizado supervisionado puro, permite ao GPT ser efetivamente *fine-tuned*.

Outra técnica igualmente utilizada é o *natural language processing* (NLP), que é o processo de fazer um inteligência artificial a entender as regras e a sintaxe de uma linguagem.

Referências

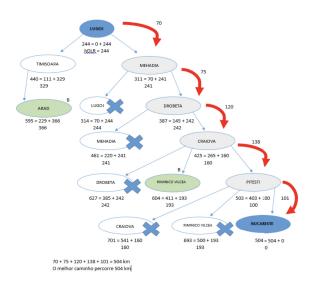
https://zapier.com/blog/how-does-chatgpt-work/

https://investnews.com.br/guias/chatgpt/

d) Pensar racionalmente, pois fornece respostas com base em padrões aprendidos a partir de uma vasta gama de dados. Ele não possui emoções, experiências pessoais ou compreensão profunda das nuances humanas. E ele não entende o que é dito e o que ele responde.

#### 2 Busca Heurística

Figura 1.3 – RESOLUÇÃO EXERCÍCIO 2: CAMINHO DE LUGOJ ATÉ BUCHAREST



#### 3 Lógica

 $R4: r \rightarrow p$  - Equivalência da implicação em R3

 $R5: q \rightarrow p$  - Silogismo hipotético em R2 e R4

 $R6: p \rightarrow q \land q \rightarrow p$  - Conjunção em R1 e R5

 $R7: q \leftrightarrow p$  - Equivalência bicondicional em R6

#### 4 Redes Neurais Artificiais

a) 
$$N_1 = 0.3$$
 função linear onde  $fa(u) = u - 3 * 0.2 + 1 * 0.8 + 1 * 0.1 = 0.3$ 

b) 
$$N_2 = 0, 3|-3*0, 1+1*0, 2+1*0, 4=0, 3$$

c) 
$$N_3 = -2|-3*0,9+1*0,5+1*0,2=-2$$

d) 
$$Sada = -0.139$$

$$0.3 * 0.9 + 0.3 * 0.3 - 2 * 0.3 + 1 * 0.1 = -0.14$$

$$TangH = (EXP(-0, 14) - EXP(-(-0, 14)))/(EXP(-0, 14) + EXP(-(-0, 14))) = -0, 139$$

# APÊNDICE 2 - LINGUAGEM DE PROGRAMAÇÃO APLICADA

#### A - ENUNCIADO

Nome da base de dados do exercício: precos carros brasil.csv

#### Informações sobre a base de dados:

Dados dos preços médios dos carros brasileiros, das mais diversas marcas, no ano de 2021, de acordo com dados extraídos da tabela FIPE (Fundação Instituto de Pesquisas Econômicas). A base original foi extraída do site Kaggle (Acesse aqui a base original). A mesma foi adaptada para ser utilizada no presente exercício.

Observação: As variáveis fuel, gear e engine\_size foram extraídas dos valores da coluna model, pois na base de dados original não há coluna dedicada a esses valores. Como alguns valores do modelo não contêm as informações do tamanho do motor, este conjunto de dados não contém todos os dados originais da tabela FIPE.

#### **Metadados:**

Tabela 2.1: Exercício 2 de LPA: Metados da Base de Dados Carros Brasil

Nome do campo	Descrição
year_of_reference	O preço médio corresponde a um mês de ano de referência.
month_of_reference	O preço médio corresponde a um mês de referência, ou seja, a
	FIPE atualiza sua tabela mensalmente.
fipe_code	Código único da FIPE.
authentication	Código de autenticação único para consulta no site da FIPE.
brand	Marca do carro.
model	Modelo do carro.
fuel	Tipo de combustível do carro.
gear	Tipo de engrenagem do carro.
engine_size	Tamanho do motor em centímetros cúbicos.
year_model	Ano do modelo do carro. Pode não corresponder ao ano de
	fabricação.
avg_price	Preço médio do carro, em reais.

**Atenção:** ao fazer o download da base de dados, selecione o formato .csv. É o formato que será considerado correto na resolução do exercício.

#### 1 Análise Exploratória dos Dados

A partir da base de dados precos\_carros\_brasil.csv, execute as seguintes tarefas:

- a) Carregue a base de dados media\_precos\_carros\_brasil.csv;
- b) Verifique se há valores faltantes nos dados. Caso haja, escolha uma tratativa para resolver o problema de valores faltantes;
- c) Verifique se há dados duplicados nos dados;

- d) Crie duas categorias, para separar colunas numéricas e categóricas. Imprima o resumo de informações das variáveis numéricas e categóricas (estatística descritiva dos dados);
- e) Imprima a contagem de valores por modelo (model) e marca do carro (brand);
- f) Dê uma breve explicação (máximo de quatro linhas) sobre os principais resultados encontrados na Análise Exploratória dos dados.

#### 2 Visualização dos Dados

A partir da base de dados precos\_carros\_brasil.csv, execute as seguintes tarefas:

- a) Gere um gráfico da distribuição da quantidade de carros por marca;
- b) Gere um gráfico da distribuição da quantidade de carros por tipo de engrenagem do carro;
- c) Gere um gráfico da evolução da média de preço dos carros ao longo dos meses de 2022 (variável de tempo no eixo X);
- d) Gere um gráfico da distribuição da média de preço dos carros por marca e tipo de engrenagem;
- e) Dê uma breve explicação (máximo de quatro linhas) sobre os resultados gerados no item d:
- f) Gere um gráfico da distribuição da média de preço dos carros por marca e tipo de combustível;
- g) Dê uma breve explicação (máximo de quatro linhas) sobre os resultados gerados no item f.

#### 3 Aplicação de Modelos de Machine Learning para Prever o Preço Médio dos Carros

A partir da base de dados precos\_carros\_brasil.csv, execute as seguintes tarefas:

- a) Escolha as variáveis **numéricas** (modelos de Regressão) para serem as variáveis independentes do modelo. A variável *target* é avg\_price.
  - Observação: caso julgue necessário, faça a transformação de variáveis categóricas em variáveis numéricas para inputar no modelo. Indique foram transformadas e **como** foram transformadas;
- b) Crie partições contendo 75% dos dados para treino e 25% para teste;
- c) Treine modelos RandomForestRegressor e XGBRegressor para predição dos preços dos carros. Observação: caso julgue necessário, mude os parâmetros dos modelos e rode novos modelos. Indique quais parâmetros foram inputados e indique o treinamento de cada modelo;
- d) Grave os valores preditos em variáveis criadas;

- e) Realize a análise de importância das variáveis para estimar a variável *target*, para cada modelo treinado;
- f) Dê uma breve explicação (máximo de quatro linhas) sobre os resultados encontrados na análise de importância de variáveis;
- g) Escolha o melhor modelo com base nas métricas de avaliação MSE, MAE e R<sup>2</sup>;
- h) Dê uma breve explicação (máximo de quatro linhas) sobre qual modelo gerou o melhor resultado e a métrica de avaliação utilizada.

# B - RESOLUÇÃO

#### 1 Análise Exploratória dos Dados

a. Carregue a base de dados media\_precos\_carros\_brasil.csv

```
# Importando bibliotecas necessárias
  import pandas as pd
2
  import matplotlib.pyplot as plt
  import seaborn as sns
4
5
  import warnings
  from sklearn.preprocessing import LabelEncoder
   from sklearn.model_selection import train_test_split
7
  from sklearn.ensemble import RandomForestRegressor
   from xgboost import XGBRegressor
9
  from sklearn.metrics import mean_squared_error, mean_absolute_error,
10
      r2_score
11
  warnings.filterwarnings('ignore')
12
13
14
   # Carregando dados do arquivo precos_carros_brasil.csv
15
   dados = pd.read_csv('precos_carros_brasil.csv')
16
17
   # Listando o nome das colunas
18
   dados.columns
19
   # Imprimindo somente as cinco primeiras linhas
20
  dados.head()
21
22
23
   # Imprime o tipo de dado de cada coluna
  dados.dtypes
25
26 # Número de linhas e colunas
27
  dados.shape
```

b. Verifique se há valores faltantes nos dados. Caso haja, escolha uma tratativa para resolver o problema de valores faltantes

```
# Verificando se há valores faltantes nos dados
dados.isna().any()
dados.isna().sum()

# Verificando se há linhas inteiramente vazias e quantas existem
```

```
dados.isnull().all(axis=1).sum()

# Removendo todas as linhas que têm todos os campos vazios
dados.dropna(axis=0, how='all', inplace=True)

# Verificando novamente se há valores faltantes
dados.isna().sum()
```

#### c. Verifique se há dados duplicados nos dados

```
# Verificando se há dados duplicados
dados.duplicated().sum()

# Removendo dados duplicados
dados.drop_duplicates(inplace=True)
```

# d. Crie duas categorias, para separar colunas numéricas e categóricas. Imprima o resumo de informações das variáveis numéricas e categóricas (estatística descritiva dos dados)

```
# Convertendo colunas numéricas para tipos adequados
2
  dados['engine_size'] = dados['engine_size'].str.replace(',', '.').
      astype(float)
  dados['year_of_reference'] = pd.to_numeric(dados['year_of_reference']
      |, errors='coerce').astype('Int64')
  dados['year_model'] = pd.to_numeric(dados['year_model'], errors='
      coerce').astype('Int64')
5
  dados['avg_price_brl'] = pd.to_numeric(dados['avg_price_brl'], errors
      ='coerce')
   # Criando duas categorias: numéricas e categóricas
7
  numericas_cols = [col for col in dados.columns if dados[col].dtype !=
8
       'object']
  categoricas_cols = [col for col in dados.columns if dados[col].dtype
10
      == 'object']
11
   # Estatística descritiva das variáveis numéricas
  dados[numericas_cols].describe()
13
14
   # Estatística descritiva das variáveis categóricas
15
  dados[categoricas_cols].describe()
```

## e. Imprima a contagem de valores por modelo (model) e marca do carro (brand)

```
# Contagem de valores por modelo e marca
dados.groupby(['model', 'brand']).size().reset_index(name='contagem')
```

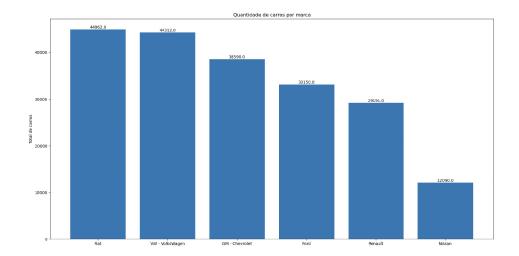
# f) Dê uma breve explicação (máximo de quatro linhas) sobre os principais resultados encontrados na Análise Exploratória dos dados.

No geral, o preço médio de um carro no Brasil custa aproximadamente R\$ 52.756,91. Um dos modelos mais vendidos é Palio Week. Adv/Adv TRYON 1.8 mpi Flex, com câmbio manual e motor 1.6. O menor preço médio de carros foi de R\$6.647,00, e o maior preço médio é R\$ 979.358.00.

#### 2 Visualização dos Dados

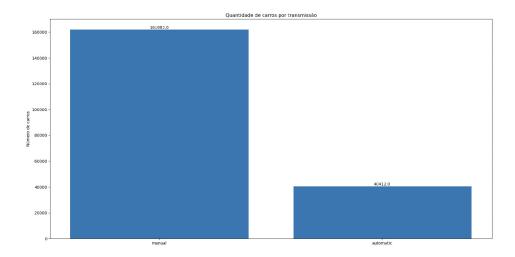
#### a. Gere um gráfico da distribuição da quantidade de carros por marca

Figura 2.1 – RESOLUÇÃO EXERCÍCIO 2: DISTRIBUIÇÃO DA QUANTIDADE DE CARROS POR MARCA



#### b. Gere um gráfico da distribuição da quantidade de carros por tipo de engrenagem do carro

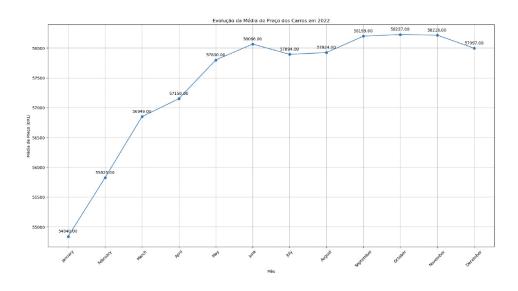




c. Gere um gráfico da evolução da média de preço dos carros ao longo dos meses de 2022 (variável de tempo no eixo X)

```
# Evolução da média de preço dos carros ao longo de 2022
  dados['year_of_reference'] = dados['year_of_reference'].round().
2
      astype(int)
  cars_in_2022 = dados[dados['year_of_reference'] == 2022]
3
  month_order = ['January','February','March','April','May','June','
4
      July',
                  'August','September','October','November','December']
5
  cars_in_2022['month_of_reference'] = pd.Categorical(cars_in_2022['
6
      month_of_reference'],
7
                                                       categories=
      month_order, ordered=True)
   cars_avg_price_in_2022 = cars_in_2022.groupby(['month_of_reference'])
8
      ['avg_price_brl'].mean().round(0).reset_index(name='average_price'
9
10
  plt.figure(figsize=(20,10))
11
  plt.plot(cars_avg_price_in_2022['month_of_reference'],
            cars_avg_price_in_2022['average_price'], marker='o',
12
      linestyle='-')
13
   for x, y in zip(cars_avg_price_in_2022['month_of_reference'],
      cars_avg_price_in_2022['average_price']):
      plt.annotate(f"{y:.2f}", (x,y), textcoords="offset points",
14
      xytext=(0,10), ha='center')
  plt.title('Evolução da Média de Preço dos Carros em 2022')
15
16 plt.xlabel('Mês')
17 | plt.ylabel('Média de Preço (BRL)')
  plt.xticks(cars_avg_price_in_2022.index, month_order, rotation=45)
  plt.grid(True)
```



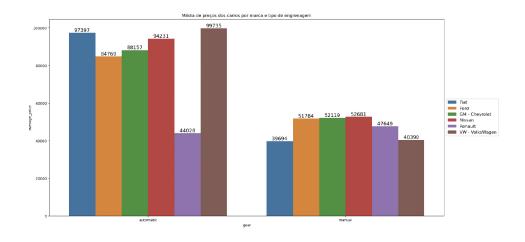


#### d. Gere um gráfico da distribuição da média de preço dos carros por marca e tipo de engrenagem

```
# Agrupar média de preços médios dos carros por marcar e tipo de
      engrenagem
  car_gear_brand_avg_price = dados.groupby(['brand', 'gear'])['
2
      avg_price_brl'].mean().round(0)
   car_gear_brand_avg_price.head(12)
3
4
5
   # Utilizando a função reset_index para criar uma ordem e facilitar a
      criação do gráfico
   car_gear_brand_avg_price = car_gear_brand_avg_price.reset_index(name=
      'average_price')
7
   car_gear_brand_avg_price.head(12)
8
   # Gerar o gráfico de distruibuição média dos preços dos carros por
      marca e tipo de engrenagem
10 plt.figure(figsize=(20,10))
11 | barplot = sns.barplot(
12 | x=' gear',
  y='average_price',
14
   hue='brand',
  data=car_gear_brand_avg_price,
  hue_order=car_gear_brand_avg_price['brand'].unique(),
16
17
   )
18
  barplot.legend(loc='center left', bbox_to_anchor=(1, 0.5), ncol=1,
      fontsize=12)
20 barplot.bar_label(barplot.containers[0], fontsize=14);
21 | barplot.bar_label(barplot.containers[1], fontsize=14);
22 barplot.bar_label(barplot.containers[2], fontsize=14);
  barplot.bar_label(barplot.containers[3], fontsize=14);
24 barplot.bar_label(barplot.containers[4], fontsize=14);
25 | barplot.bar_label(barplot.containers[5], fontsize=14);
```

```
plt.title('Média de preços dos carros por marca e tipo de engrenagem')
```

Figura 2.4 – RESOLUÇÃO EXERCÍCIO 2: MÉDIA DE PREÇOS DOS CARROS POR MARCA E TIPO DE ENGRENAGEM AGRUPADOS POR TIPO



```
# Gerar o segundo tipo de gráfico de distruibuição média dos preços
      dos carros por marca e tipo de engrenagem
2 plt.figure(figsize=(20,10))
3
  |barplot = sns.barplot(
4
  x='brand',
  y='average_price',
5
  hue='gear',
6
7
  data=car_gear_brand_avg_price,
8
  hue_order=car_gear_brand_avg_price['gear'].unique(),)
  barplot.legend(loc='center left', bbox_to_anchor=(1, 0.5), ncol=1,
      fontsize=12)
10 barplot.bar_label(barplot.containers[0], fontsize=14);
  |barplot.bar_label(barplot.containers[1], fontsize=14);
12
  plt.title('Média de Preços dos Carros por Marca e Tipo de engrenagem'
      )
```

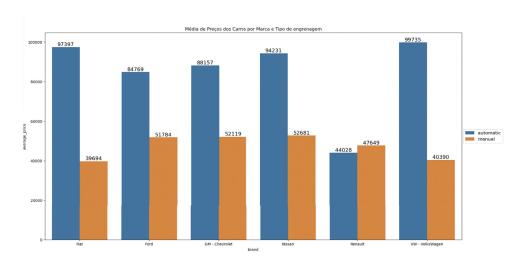


Figura 2.5 – RESOLUÇÃO EXERCÍCIO 2: MÉDIA DE PREÇOS DOS CARROS POR MARCA E TIPO DE ENGRENAGEM AGRUPADOS POR MARCA

#### e. Dê uma breve explicação (máximo de quatro linhas) sobre os resultados gerados no item d:

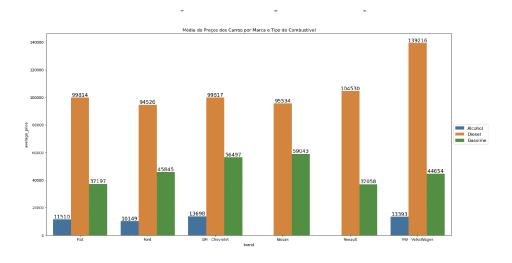
O preço médio do carro com câmbio automático é relativamente mais alto quando comparado aos de câmbio manual para todas as marcas, com exceção da Renault. Volkswagen e Fiat lideram as marcas que possuem maior preço médio para carros de câmbio automático, ao passo que lideram também as marcas com menor preço médio dos carros com câmbio manual.

## f. Gere um gráfico da distribuição da média de preço dos carros por marca e tipo de combustível

```
# Agrupar média de preços médios dos carros por marcar e tipo de
      combustível
  cars_avg_price_brand_fuel = dados.groupby(['brand', 'fuel'])['
      avg_price_brl'].mean().round(0)
3
  cars_avg_price_brand_fuel.head(16)
4
5
   # Utilizando a função reset_index para criar uma ordem e facilitar a
      criação do gráfico
6
   cars_avg_price_brand_fuel = cars_avg_price_brand_fuel.reset_index(
      name='average_price')
7
  cars_avg_price_brand_fuel.head(16)
8
  # Gerar gráfico de média de preços médios de carros por marcar e tipo
9
       de combustível
  plt.figure(figsize=(20, 10))
10
  barplot = sns.barplot(
11
  x='brand',
12
13
  y='average_price',
14
  hue='fuel',
  data=cars_avg_price_brand_fuel,
15
  hue_order=cars_avg_price_brand_fuel['fuel'].unique()
16
17
  barplot.legend(loc='center left', bbox_to_anchor=(1, 0.5), ncol=1,
18
      fontsize=12)
19
  |barplot.bar_label(barplot.containers[0], fontsize=14)
  barplot.bar_label(barplot.containers[1], fontsize=14)
  |barplot.bar_label(barplot.containers[2], fontsize=14)
```

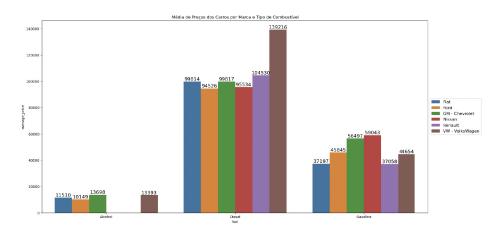
```
22 plt.title('Média de Preços dos Carros por Marca e Tipo de Combustível')
```

Figura 2.6 – RESOLUÇÃO EXERCÍCIO 2: MÉDIA DE PREÇOS DOS CARROS POR MARCA E TIPO DE COMBUSTÍVEL AGRUPADOS POR MARCA



```
# Gerar gráfico de média de preços médios de carros por marcar e tipo
       de combustível - Segundo tipo
2
  plt.figure(figsize=(20, 10))
3
  barplot = sns.barplot(
  x='fuel',
4
5
   y='average_price',
  hue='brand',
  data=cars_avg_price_brand_fuel,
7
  hue_order=cars_avg_price_brand_fuel['brand'].unique())
8
9
10
  barplot.legend(loc='center left', bbox_to_anchor=(1, 0.5), ncol=1,
      fontsize=12)
11 | barplot.bar_label(barplot.containers[0], fontsize=14)
12 | barplot.bar_label(barplot.containers[1], fontsize=14)
13
  barplot.bar_label(barplot.containers[2], fontsize=14)
  barplot.bar_label(barplot.containers[3], fontsize=14)
14
15 | barplot.bar_label(barplot.containers[4], fontsize=14)
  barplot.bar_label(barplot.containers[5], fontsize=14)
17
  plt.title ('Média de Preços dos Carros por Marca e Tipo de Combustível
      ′)
```





- g. Dê uma breve explicação (máximo de quatro linhas) sobre os resultados gerados no item f:
  - 1. O preço médio do carro à diesel é muito mais elevado para todas as marcas, sendo a Volkswagen a marca com maior preço médio para esta categoria.
  - 2. O preço médio do carro à álcool lidera com os menores valores.
  - 3. Os carros à gasolina com menor preço médio são das marcas Fiat e Renault.

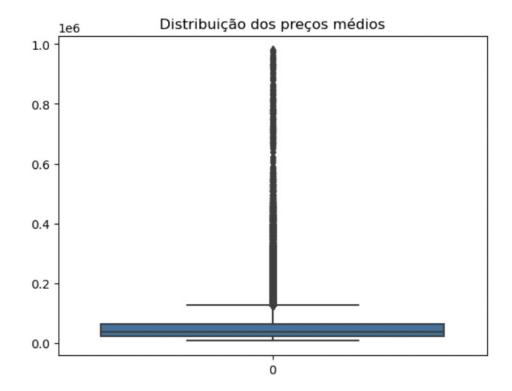
#### 3 Aplicação de Modelos de Machine Learning para Prever o Preço Médio dos Carros

a) Escolha as variáveis numéricas (modelos de Regressão) para serem as variáveis independentes do modelo. A variável target é avg\_price.

```
# Gráfico de distribuição boxplot dos preços médios dos carros
1
   sns.boxplot(dados['avg_price_brl']).set_title('Distribuição dos preç
2
      os médios')
3
   # Transformar mês de referência em uma variável numérica
4
   dados['month_of_reference'] = LabelEncoder().fit_transform(dados['
5
      month of reference' 1)
6
7
   # Transformar marca de carro em uma variável numérica
   dados['brand'] = LabelEncoder().fit transform(dados['brand'])
8
9
10
   # Transformar modelo de carro em uma variável numérica
11
   dados['model'] = LabelEncoder().fit_transform(dados['model'])
12
13
   # Transformar tipo de combustível em uma variável numérica
   dados['fuel'] = LabelEncoder().fit_transform(dados['fuel'])
14
15
   # Transformar tipo de transmissão em uma variável numérica
16
17
   dados['gear'] = LabelEncoder().fit_transform(dados['gear'])
18
19
   # Remover as variáveis de entrada que não são consideradas
      importantes para a predição
20 dados_interesting = dados.drop([
```

```
21 'fipe_code',
   'authentication',
22
23
   ], axis=1)
24
   # Mapa de correlação das variáveis numéricas com variável Target
25
26
   sns.heatmap(dados_interesting.corr("spearman"), annot = True)
   plt.title("Mapa de Correlação das Variáveis Numéricas\n", fontsize =
27
      18)
28
  plt.show()
29
30 | # Variável X contém apenas variáveis numéricas de interesse para a an
      álise excluindo a variável target
31 | X = dados_interesting.drop(['avg_price_brl'], axis=1)
32
  X.head()
33
34 | # Variável Y contém apenas a variável target - avg_price_brl
  Y = dados_interesting['avg_price_brl']
35
36 Y.head()
```

Figura 2.8 – RESOLUÇÃO EXERCÍCIO 3: BOXPLOT DA DISTRIBUIÇÃO DOS PREÇOS MÉDIOS



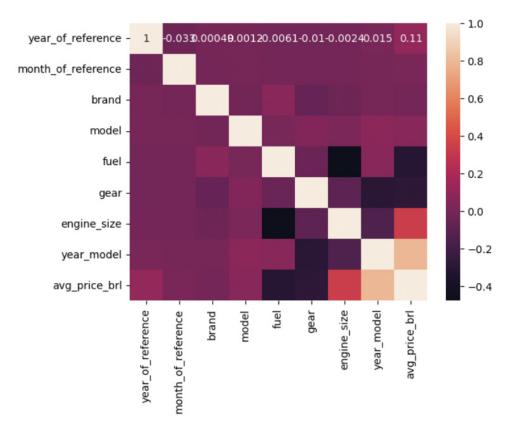


Figura 2.9 – RESOLUÇÃO EXERCÍCIO 3: MAPA DE CORRELAÇÃO ENTRE AS VARIÁVEIS NUMÉRICAS

#### b. Crie partições contendo 75% dos dados para treino e 25% para teste

```
1 # Divisão: 25% dos dados são de teste e 75% de treinamento
2 X_train, x_test, Y_train, y_test = train_test_split(X, Y, test_size = 0.25, random_state = 42)
```

# c. Treine modelos RandomForest (biblioteca RandomForestRegressor) e XGBoost (biblioteca XGBRegressor) para predição dos preços dos carros

```
#Algoritmo Random Forest, sem especificar nenhum parâmetro (número de
       árvores, número de ramificações, etc)
  model_rf = RandomForestRegressor()
2
3
4
  # Random Forest com os hiperparametros:
  model_rf_paramenters = RandomForestRegressor(
  max_depth=29,
6
  min_samples_leaf=32,
7
  min_samples_split=28,
8
9
  n_estimators=208,
10
  random_state=43
11
12
  # Ajuste do modelo, de acordo com as variáveis de treinamento
  model_rf_paramenters.fit(X_train, Y_train)
13
14
15
   # Ajuste do modelo, de acordo com as variáveis de treinamento
16
  model_rf.fit(X_train, Y_train)
17
18
```

```
19 | # Algoritmo XGBoost sem parametros
20
  model_xg = XGBRegressor()
21
22 | # Ajuste do modelo, de acordo com as variáveis de treinamento
   model xq.fit(X train, Y train)
23
24
   #Algoritmo XGBoost com parametros
25
26
   model_parameters_xg = XGBRegressor(
27
       learning rate=1,
28
       n estimators=100,
29
       max_depth=25,
30
       min_child_weight=1,
31
       qamma=0,
32
       subsample=0.8,
       colsample_bytree=0.8,
33
34
       reg_alpha=0,
35
       reg_lambda=1,
36
       random_state=42,
37
       n_{jobs}=-1
38
39
   # Ajuste do modelo, de acordo com as variáveis de treinamento
40
  model_parameters_xg.fit(X_train, Y_train)
```

#### d. Grave os valores preditos em variáveis criadas

```
# Predição dos valores de preço médio dos carros com base nos dados
      de teste Random Forest sem parametros
  predicted_values_rf = model_rf.predict(x_test)
3
  # Predição dos valores de preço médio dos carros com base nos dados
4
      de teste Random Forest com parametros
5
  predicted_values_parameters_rf = model_rf_paramenters.predict(x_test)
6
7
  # Predição dos valores de preços médios dos carros com base nos dados
       de teste XGBoost sem parametros
8
  predicted_values_xg = model_xg.predict(x_test)
10
  # Predição dos valores de preços médios dos carros com base nos dados
       de teste XGBoost com parametros
11 | predicted_values_parameters_xg = model_parameters_xg.predict(x_test)
```

# e. Realize a análise de importância das variáveis para estimar a variável target, para cada modelo treinado

```
brand 0.016695
         fuel 0.032713
         gear 0.035363
         model 0.058640
         year model 0.388397
         engine_size 0.450472
  # Análise da importancia das variáveis para estimar a variável target
      - Random Forest com parametros
2 model_rf_paramenters.feature_importances_
  feature importances parameters rf = pd.DataFrame(
 model_rf_paramenters.feature_importances_,
  index = X_train.columns,
  columns = ['importances']
  ).sort_values('importances', ascending=True)
7
  feature_importances_parameters_rf
         month of reference 0.000690
         year_of_reference 0.010383
         brand 0.016098
         gear 0.021959
         fuel 0.034132
         model 0.037028
         year_model 0.410473
         engine_size 0.469236
 # Análise da importancia das variáveis para estimar a variável target
      - XGBoost com parametros
 model_parameters_xq.feature_importances_
3
  feature_importances_parameters_xg = pd.DataFrame(
 model_rf_paramenters.feature_importances_,
5
  index = X_train.columns,
  columns = ['importances']
  ).sort_values('importances', ascending=True)
  feature_importances_parameters_xg
         month_of_reference 0.000690
         year_of_reference 0.010383
         brand 0.016098
         gear 0.021959
         fuel 0.034132
         model 0.037028
         year_model 0.410473
         engine_size 0.469236
  # Analisando a importância das variáveis para estimar a variável
     target - XGBoost sem parametros
 model_xq.feature_importances_
 |feature_importances_xg = pd.DataFrame(
4 model_xg.feature_importances_,
```

month\_of\_reference 0.005168 year\_of\_reference 0.012553

```
5 index=X_train.columns,
6 columns=['importances']
7 ).sort_values('importances', ascending=True)
feature_importances_xg
```

month\_of\_reference 0.004271 year\_of\_reference 0.017475 model 0.023331 brand 0.056276 gear 0.122657 fuel 0.154437 year\_model 0.190992 engine\_size 0.430561

f. Dê uma breve explicação (máximo de quatro linhas) sobre os resultados encontrados na análise de importância de variáveis:

As variáveis que tiveram maior relevância na medida da performance dos modelos treinados foram o tamanho do motor (engine\_size) e ano do modelo (year\_model) respectivamente. O mês de referência foi o que obteve menor relevância em todos o modelos treinados.

g. Escolha o melhor modelo com base nas métricas de avaliação MSE, MAE e R2:

```
1 # MSE - calcula o erro quadrático médio das predições do modelo
  # Random Forest sem parametro
3
  |mse_rf = mean_squared_error(y_test, predicted_values_rf) #12181571.08
4
5
  # O MAE calcula a média da diferença absoluta entre o valor predito e
       o valor real
  # Random Forest sem parametro
7
  mea_rf = mean_absolute_error(y_test, predicted_values_rf) # 1765.88
8
9
  # O R2 é uma métrica que varia entre O e 1 e é uma razão que indica o
       quão bom o nosso modelo
  # Random Forest sem parametros
10
  r2_score(y_test, predicted_values_rf) # 0.995
11
```

```
# MSE - XGBoost sem parametros
mse_xg = mean_squared_error(y_test, predicted_values_xg) #
30952404.229

# MAE - XGBoost sem parametros
mae_xg = mean_absolute_error(y_test, predicted_values_xg) # 3245.31
```

```
6
7 # R2 score - XGBoost sem parametros
8 r2_score(y_test, predicted_values_xg) # 0.988
```

O melhor modelo com base nessas medidas de acurácia foi o modelo RandomForest sem parâmetros.

h. Dê uma breve explicação (máximo de quatro linhas) sobre qual modelo gerou o melhor resultado e a métrica de avaliação utilizada:

O modelo que apresentou os melhores resultado foi o RandomForest sem parametros para medir a sua acurácia foi utilizado as médidas de acurácia MSE (*mean squared error*) com o valor 12165425.33231638, MAE (*mean absolute error*) com o valor 12165425.33 e R2\_score com 0.995 de precisão.

## APÊNDICE 3 – LINGUAGEM R

### A - ENUNCIADO

### 1 Pesquisa com Dados de Satélite (Satellite)

O banco de dados consiste nos valores multiespectrais de pixels em vizinhanças 3x3 em uma imagem de satélite, e na classificação associada ao pixel central em cada vizinhança. O objetivo é prever esta classificação, dados os valores multiespectrais.

Um quadro de imagens do Satélite Landsat com MSS (Multispectral Scanner System) consiste em quatro imagens digitais da mesma cena em diferentes bandas espectrais. Duas delas estão na região visível (correspondendo aproximadamente às regiões verde e vermelha do espectro visível) e duas no infravermelho (próximo). Cada pixel é uma palavra binária de 8 bits, com 0 correspondendo a preto e 255 a branco. A resolução espacial de um pixel é de cerca de 80m x 80m. Cada imagem contém 2340 x 3380 desses pixels. O banco de dados é uma subárea (minúscula) de uma cena, consistindo de 82 x 100 pixels. Cada linha de dados corresponde a uma vizinhança quadrada de pixels 3x3 completamente contida dentro da subárea 82x100. Cada linha contém os valores de pixel nas quatro bandas espectrais (convertidas em ASCII) de cada um dos 9 pixels na vizinhança de 3x3 e um número indicando o rótulo de classificação do pixel central.

As classes são: solo vermelho, colheita de algodão, solo cinza, solo cinza úmido, restolho de vegetação, solo cinza muito úmido.

Os dados estão em ordem aleatória e certas linhas de dados foram removidas, portanto você não pode reconstruir a imagem original desse conjunto de dados. Em cada linha de dados, os quatro valores espectrais para o pixel superior esquerdo são dados primeiro, seguidos pelos quatro valores espectrais para o pixel superior central e, em seguida, para o pixel superior direito, e assim por diante, com os pixels lidos em sequência, da esquerda para a direita e de cima para baixo. Assim, os quatro valores espectrais para o pixel central são dados pelos atributos 17, 18, 19 e 20. Se você quiser, pode usar apenas esses quatro atributos, ignorando os outros. Isso evita o problema que surge quando uma vizinhança 3x3 atravessa um limite.

O banco de dados se encontra no pacote mlbench e é completo (não possui dados faltantes).

Tarefas:

- 1. Carregue a base de dados Satellite.
- 2. Crie partições contendo 80% para treino e 20% para teste.
- 3. Treine modelos RandomForest, SVM e RNA para predição destes dados.
- 4. Escolha o melhor modelo com base em suas matrizes de confusão.
- 5. Indique qual modelo dá o melhor o resultado e a métrica utilizada.

### 2 Estimativa de Volumes de Árvores

Modelos de aprendizado de máquina são bastante usados na área da engenharia florestal (mensuração florestal) para, por exemplo, estimar o volume de madeira de árvores sem ser necessário abatê-las.

O processo é feito pela coleta de dados (dados observados) através do abate de algumas árvores, onde sua altura, diâmetro na altura do peito (dap), etc, são medidos de forma exata.

Com estes dados, treina-se um modelo de AM que pode estimar o volume de outras árvores da população.

Os modelos, chamados alométricos, são usados na área há muitos anos e são baseados em regressão (linear ou não) para encontrar uma equação que descreve os dados. Por exemplo, o modelo de Spurr é dado por:

$$Volume = b0 + b1 * dap2 * Ht$$
(3.1)

Onde dap é o diâmetro na altura do peito (1,3 metros), Ht é a altura total. Tem-se vários modelos alométricos, cada um com uma determinada característica, parâmetros, etc. Um modelo de regressão envolve aplicar os dados observados e encontrar b0 e b1 no modelo apresentado, gerando assim uma equação que pode ser usada para prever o volume de outras árvores.

Dado o arquivo Volumes.csv, que contém os dados de observação, escolha um modelo de aprendizado de máquina com a melhor estimativa, a partir da estatística de correlação.

Tarefas

- 1. Carregar o arquivo Volumes.csv (http://www.razer.net.br/datasets/Volumes.csv)
- 2. Eliminar a coluna NR, que só apresenta um número sequencial
- 3. Criar partição de dados: treinamento 80%, teste 20%
- 4. Usando o pacote "caret", treinar os modelos: Random Forest (rf), SVM (svmRadial), Redes Neurais (neuralnet) e o modelo alométrico de SPURR. O modelo alométrico é dado por:

$$Volume = b0 + b1 * dap2 * Ht$$
 (3.2)

```
1    alom <- nls(VOL ~ b0 + b1*DAP*DAP*HT, dados, start=list(b0 = 0.5, b1=0.5))</pre>
2
```

- 5. Efetue as predições nos dados de teste
- 6. Crie suas próprias funções (UDF) e calcule as seguintes métricas entre a predição e os dados observados
  - Coeficiente de determinação :  $R^2$

$$R^{2} = 1 - \frac{\sum_{i=1}^{n} (y_{i} - \hat{y}_{i})^{2}}{\sum_{i=1}^{n} (y_{i} - \bar{y})^{2}}$$
(3.3)

onde yi é o valor observado, yi é o valor predito e y é a média dos valores yi observados. Quanto mais perto de 1 melhor é o modelo;

- Erro padrão da estimativa:  $S_{yx}$ 

$$S_{yx} = \sqrt{\frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{n-2}}$$
 (3.4)

esta métrica indica erro, portanto quanto mais perto de 0 melhor é o modelo;

- Syx%:

$$S_{yx}\% = \frac{S_{yx}}{\bar{y}} \times 100 \tag{3.5}$$

esta métrica indica porcentagem de erro, portanto quanto mais perto de 0 melhor é o modelo;

7. Escolha o melhor modelo.

# B - RESOLUÇÃO

## 1 Pesquisa com Dados de Satélite (Satellite)

```
#### 1.1 Carreque a base de dados Satellite
1
2
   #instalar o pacotes necessários
3
   #install.packages("mlbench", repos = "http://cran.us.r-project.org")
5
   #install.packages("e1017", repos = "http://cran.us.r-project.org")
  #install.packages("randomForest", repos = "http://cran.us.r-project.
6
      orq")
   #install.packages("kernlab", repos = "http://cran.us.r-project.org")
7
   #install.packages("caret", repos = "http://cran.us.r-project.org")
8
10 # usar os pacotes necessários
  library (mlbench)
11
12 | library(caret)
13
14 # carregar os dados Satellite
   data(Satellite)
15
16
17
  # exibir estrutura dos dados Satellite
   str(Satellite)
18
19
   # apresentar alguma medidas estatísticas do dados Satellite
20
   summary(Satellite)
21
2.2.
23
   # exibir alguns dados do Satellite
   head(Satellite, n = 6)
24
25
   #### 1.2 Crie partições contendo 80% para treino e 20% para teste
26
27
28 | # Para reproductibilidade
29
   set.seed(7)
30
31
  | # particionar em 80% para treino e 20% para teste
   indices <- createDataPartition(Satellite[["classes"]], p=0.8, list=F)</pre>
32
   treino <- Satellite[indices, ]</pre>
33
  teste <- Satellite[-indices, ]</pre>
34
35
   #### 1.3 Treine modelos RandomForest, SVM e RNA para predição destes
36
      dados.
37
38 | # treinar modelos RandomForest, SVM e RNA
```

```
39 | rf <- train(classes ~ ., data=treino, method="rf")
  svm <- train(classes ~ ., data=treino, method="svmRadial")</pre>
41
   rna <- train(classes ~ ., data=treino, method="nnet", trace=F)</pre>
42
43
   #### 1.4 Escolha o melhor modelo com base em suas matrizes de confusã
44
45
   # predições
46 | predict.rf <- predict(rf, teste)
   predict.svm <- predict(svm, teste)</pre>
47
48
   predict.rna <- predict(rna, teste)</pre>
49
50 | # matrizes de confusões de cada uma das predições
51
52
   # matriz de confusão para o modelo RF
53
   conf_matrix.rf <- confusionMatrix(predict.rf, teste[["classes"]])</pre>
54
55 | print (conf_matrix.rf)
56
   cat('\n')
57
58 # matriz de confusão para o modelo SVM
59
   conf_matrix.svm <- confusionMatrix(predict.svm, teste[["classes"]])</pre>
60
   print(conf_matrix.svm)
61
  cat('\n')
62
63
64
   # matriz de confusão para o modelo RNA
   conf_matrix.rna <- confusionMatrix(predict.rna, teste[["classes"]])</pre>
66
67
   print(conf_matrix.rna)
   cat('\n')
68
69
70
   #### 1.5 Indique qual modelo dá o melhor o resultado e a métrica
      utilizada
71
72 | #0 melhor modelo foi random forest com acurácia de 0.9182 e kappa de
      0.8987. A métrica utilizada foram a acurácia e kappa.
```

### 2 Estimativa de Volume de Árvores

```
#### 2.1 Carregar o arquivo Volumes.csv (<http://www.razer.net.br/
      datasets/Volumes.csv>)
2
  dados <- read.csv("http://www.razer.net.br/datasets/Volumes.csv", sep
3
      =";", dec=",")
4
  head(dados)
5
6
  #### 2.2 Eliminar a coluna NR, que só apresenta um número sequencial
7
   dados[["NR"]] <- NULL</pre>
8
9
10 | #### 2.3 Criar partição de dados: treinamento 80%, teste 20%
```

```
11
   regression.indices <- caret::createDataPartition(dados[["VOL"]], p</pre>
12
      =0.8, list=F)
   regression.treino <- dados[regression.indices, ]</pre>
13
   regression.teste <- dados[-regression.indices, ]</pre>
14
15
  | #### 2.4 Usando o pacote "caret", treinar os modelos: Random Forest (
16
      rf), SVM (svmRadial), Redes Neurais (neuralnet) e o modelo alomé
      trico de SPURR
17
   # Para reproductibilidade
18
19
   set.seed(7)
20
2.1
  regression.rf <- caret::train(VOL ~ ., data=regression.treino, method
22
   regression.svm <- caret::train(VOL ~ ., data=regression.treino,
      method="svmRadial")
   regression.rna <- caret::train(VOL ~ ., data=regression.treino,
23
      method="nnet", trControl=trainControl(method = "LOOCV"), trace=F)
24
25
   ##### treino do modelo Spurr
26
27
   regression.spurr <- nls( VOL ~ b0 + b1*DAP*DAP*HT, data=regression.
      treino, start=list(b0=0.5, b1=0.5))
28
29
   ##### visualizar os resultados de Spurr
30
31
   summary(regression.spurr)
32
33
   #### 2.5 Efetue as predições nos dados de teste
34
35
   # predições
36
37
   # Para reproductibilidade
38
   set.seed(7)
39
40
   predict.regression.rf <- predict(regression.rf, regression.teste)</pre>
   predict.regression.svm <- predict(regression.svm, regression.teste)</pre>
41
42
   predict.regression.rna <- predict(regression.rna, regression.teste)</pre>
43
   predict.regression.suprr <- predict(regression.spurr, regression.</pre>
      teste)
44
45
   #### 2.6 Crie suas próprias funções (UDF) e calcule as seguintes mé
      tricas entre a predição e os dados observados
46
         Erro padrão de estimativa: Syx
47
48
49
   Syx <- function(reals, predicteds, n) {</pre>
   return (sqrt(sum((reals - predicteds)^2)/(n - 2)))
50
51
   }
52
```

```
53
   # - Erro padrão de estimativa em porcentagem: Syx%
54
   SyxPercent <- function(reals, predicteds, n) {</pre>
55
   return ((Syx(reals, predicteds, n)/mean(reals)) *100)
56
57
58
59
         Coeficientededeterminação: R2
60
61
   R2 <- function (reals, predicteds) {
     return (1 - sum((reals - predicteds)^2)/sum((reals - mean(reals))
       ^2))
63
64
65
   #### 2.7 Escolha o melhor modelo.
66
67
   ##### métrica de estimativas para o modelo RandomForest - Regressão
68
69
        coeficiente de determinação
70
71
   R2 (regression.teste[["VOL"]], predict.regression.rf)
72
73
         Erro padrão estimativa
74
75
   n <- nrow(regression.teste)</pre>
   Syx(regression.teste[["VOL"]], predict.regression.rf, n)
76
77
78
        Erro padrão estimativa em porcentagem
79
80
   n <- nrow(regression.teste)</pre>
   SyxPercent(regression.teste[["VOL"]], predict.regression.rf, n)
81
82
83
   ##### métrica de estimativas para o modelo SVM - Regressão
84
    # - coeficiente de determinação
85
86
   R2 (regression.teste[["VOL"]], predict.regression.svm)
87
88
89
   # - Erro padrão estimativa
90
91
   n <- nrow(regression.teste)</pre>
   Syx(regression.teste[["VOL"]], predict.regression.svm, n)
92
93
94
   # - Erro padrão estimativa em porcentagem
95
96
   n <- nrow(regression.teste)</pre>
   SyxPercent (regression.teste[["VOL"]], predict.regression.svm, n)
97
98
99
   ##### métricas de estimativas para o modelo nnet - Regressão
100
101
          coeficiente de determinação
102
```

```
103 R2 (regression.teste[["VOL"]], predict.regression.rna)
104
105
         Erro padrão estimativa
106
107
    n <- nrow(regression.teste)</pre>
    Syx(regression.teste[["VOL"]], predict.regression.rna, n)
108
109
110
         Erro padrão estimativa em porcentagem
111
112
    n <- nrow(regression.teste)</pre>
    SyxPercent(regression.teste[["VOL"]], predict.regression.rna, n)
113
114
115
    ##### métricas de estimativas para o modelo Spurr
116
         coeficiente de determinação
117
118
   R2 (regression.teste[["VOL"]], predict.regression.suprr)
119
120
121
   # - Erro padrão estimativa
122
123
   n <- nrow(regression.teste)</pre>
124
    Syx(regression.teste[["VOL"]], predict.regression.suprr, n)
125
126
         Erro padrão estimativa em porcentagem
127
   n <- nrow(regression.teste)</pre>
128
129
    SyxPercent (regression.teste[["VOL"]], predict.regression.suprr, n)
130
131
    ##### 2.7 escolha o melhor modelo
132
    ##### Resumo dos resultados
133
134
135
    ##### RF:
136
137
    # 1. coeficiente de determinação: 0.8223603.
    # 2. Erro padrão estimativa: 0.1376052.
138
139
    # 3. Erro padrão estimativa em porcentagem: 10.42195
140
    ##### SVM:
141
142
143
    # 1. coeficiente de determinação: 0.6254546
144
    # 2. Erro padrão estimativa: 0.19981
145
    # 3. Erro padrão estimativa em porcentagem: 15.13322
146
147
    ##### nnet:
148
149
    # 1. coeficiente de determinação: -1.069672
    # 2. Erro padrão estimativa: 0.4696948
150
151
    # 3. Error padrão estimativa em porcentagem: 35.57377
152
153 | ###### spurr:
```

154	
155	# 1. coeficiente de determinação: 0.7734018
156	# 2. Erro padrão estimativa: 0.1554151
157	# 3. Error padrão estimativa em porcentagem: 11.77084
158	
159	# Com base nas métricas, o modelo que se saiu melhor foi o
	RandomForest, com R2 igual 0.8223603, Erro padrão estimativa de
	0.1376052 e Erro padrão de estima- tiva em porcentagem de 10.42195

## APÊNDICE 4 – ESTATÍSTICA APLICADA I

### A - ENUNCIADO

### 1. Gráficos e tabelas

- (15 pontos) Elaborar os gráficos box-plot e histograma das variáveis "age" (idade da esposa) e "husage" (idade do marido) e comparar os resultados
- (15 pontos) Elaborar a tabela de frequencias das variáveis "age" (idade da esposa) e "husage" (idade do marido) e comparar os resultados

## 2. Medidas de posição e dispersão

- (15 pontos) Calcular a média, mediana e moda das variáveis "age" (idade da esposa) e "husage" (idade do marido) e comparar os resultados
- (15 pontos) Calcular a variância, desvio padrão e coeficiente de variação das variáveis "age" (idade da esposa) e "husage" (idade do marido) e comparar os resultados

## 3. Testes paramétricos ou não paramétricos

(40 pontos) Testar se as médias (se você escolher o teste paramétrico) ou as medianas (se você escolher o teste não paramétrico) das variáveis "age" (idade da esposa) e "husage" (idade do marido) são iguais, construir os intervalos de confiança e comparar os resultados.

Obs: Você deve fazer os testes necessários (e mostra-los no documento pdf) para saber se você deve usar o unpaired test (paramétrico) ou o teste U de Mann-Whitney (não paramétrico), justifique sua resposta sobre a escolha. Lembre-se de que os intervalos de confiança já são mostrados nos resultados dos testes citados no item 1 acima.

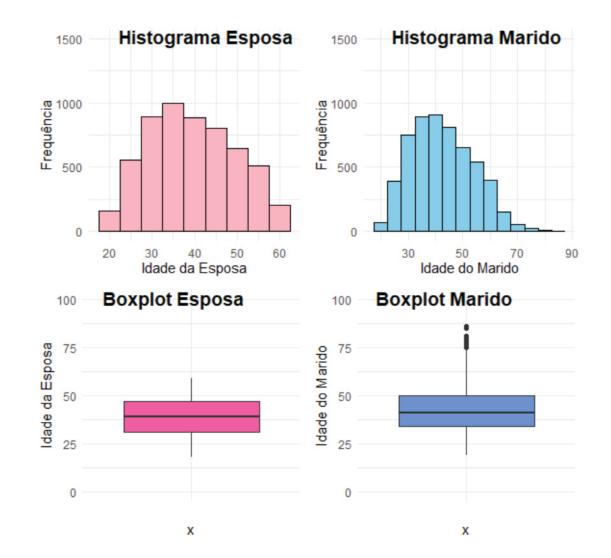
# **B - RESOLUCÃO**

#### 1. Gráficos e tabelas

```
# Criando o histograma da idade do marido
1
  histograma_marido <- ggplot(salarios, aes(x = husage)) +</pre>
2
3
     geom_histogram(binwidth = 5, bins = 9, fill = "#87CEEB", color =
      black") +
     labs(x = "Idade do Marido", y = "Frequência") +
4
5
     theme_minimal() +
6
     coord\_cartesian(ylim = c(0, 1500))
7
   # Criando o boxplot da idade do marido
8
9
  boxplot_marido <- ggplot(salarios, aes(x = "", y = husage)) +
     geom_boxplot(fill = "#6495ED") +
10
     labs(y = "Idade do Marido") +
11
    theme_minimal() +
12
     coord\_cartesian(ylim = c(0, 100))
13
14
  # Criando o histograma da idade da esposa
15
16 histograma_esposa <- ggplot(salarios, aes(x = age)) +
```

```
17
     geom_histogram(binwidth = 5, bins = 9, fill = "#FFB6C1", color =
      black") +
18
     labs(x = "Idade da Esposa", y = "Frequência") +
19
     theme_minimal() +
20
     coord\_cartesian(ylim = c(0, 1500))
21
22
   # Criando o boxplot da idade da esposa
23
   boxplot_esposa <- ggplot(salarios, aes(x = "", y = age)) +</pre>
24
     geom_boxplot(fill = "#F25EA2") +
     labs(y = "Idade da Esposa") +
25
26
     theme_minimal() +
     coord_cartesian(ylim = c(0, 100))
27
28
29
   # Combinando os gráficos em uma única visualização
   plot_comparativo <- plot_grid(histograma_esposa, histograma_marido,</pre>
      boxplot_esposa, boxplot_marido, labels = c("Histograma Esposa",
      Histograma Marido", "Boxplot Esposa", "Boxplot Marido"), ncol = 2)
```

Figura 4.1 – EXERCÍCIO 1: HISTOGRAMA E BOX-PLOT



A análise mostra que os maridos são, em média, 3 anos mais velhos que as esposas (42,46 anos contra 39,43 anos). A variação de idade é maior entre os maridos, com idades entre 19 e 86 anos, enquanto as esposas têm idades entre 18 e 59 anos.

```
1  # Tabela de frequência das idades dos maridos
2  tabela.freq.M <- fdt(salarios$husage)</pre>
```

Tabela 4.1: Tabela de frequência das idades (maridos)

Class limits	f	rf	rf(%)	cf	cf(%)
[18.81, 23.671)	102	0.02	1.81	102	1.81
[23.671, 28.531)	466	0.08	8.27	568	10.08
[28.531, 33.392)	809	0.14	14.36	1377	24.44
[33.392, 38.253)	895	0.16	15.89	2272	40.33
[38.253, 43.114)	917	0.16	16.28	3189	56.60
[43.114, 47.974)	629	0.11	11.16	3818	67.77
[47.974, 52.835)	649	0.12	11.52	4467	79.29
[52.835, 57.696)	541	0.10	9.60	5008	88.89
[57.696, 62.556)	394	0.07	6.99	5402	95.88
[62.556, 67.417)	152	0.03	2.70	5554	98.58
[67.417, 72.278)	51	0.01	0.91	5605	99.49
[72.278, 77.139)	21	0.00	0.37	5626	99.86
[77.139, 81.999)	6	0.00	0.11	5632	99.96
[81.999, 86.86)	2	0.00	0.04	5634	100.00

```
# Tabela de frequência das idades das esposas
tabela.freq.E <- fdt(salarios$age)</pre>
```

Tabela 4.2: Tabela de frequência das idades (esposas)

Class limits	f	rf	rf(%)	cf	cf(%)
[17.82, 20.804)	61	0.01	1.08	61	1.08
[20.804, 23.787)	161	0.03	2.86	222	3.94
[23.787, 26.771)	312	0.06	5.54	534	9.48
[26.771, 29.754)	505	0.09	8.96	1039	18.44
[29.754, 32.738)	562	0.10	9.98	1601	28.42
[32.738, 35.721)	571	0.10	10.13	2172	38.55
[35.721, 38.705)	624	0.11	11.08	2796	49.63
[38.705, 41.689)	510	0.09	9.05	3306	58.68
[41.689, 44.672)	542	0.10	9.62	3848	68.30
[44.672, 47.656)	432	0.08	7.67	4280	75.97
[47.656, 50.639)	389	0.07	6.90	4669	82.87
[50.639, 53.623)	358	0.06	6.35	5027	89.23
[53.623, 56.606)	304	0.05	5.40	5331	94.62
[56.606, 59.59)	303	0.05	5.38	5634	100.00

Aproximadamente metade das esposas tem entre 29 e 41 anos, enquanto quase metade dos maridos está entre 28 e 43 anos.

# 2. Medidas de posição e dispersão

```
# Média da idade dos maridos = 42.45
   media.marido <- mean(salarios$husage)</pre>
2
3
4
   # Média da idade das esposas = 39.43
  media.esposa <- mean(salarios$age)</pre>
7
   # Mediana da idade dos maridos = 41
  mediana.marido <- median(salarios$husage)</pre>
8
10
   # Mediana da idade das esposas = 39
   mediana.esposa <- median(salarios$age)</pre>
11
12
13
  # Moda da idade dos maridos = a idade 44 aparece 201 vezes
  moda.marido <- subset(table(salarios$husage), table(salarios$husage)</pre>
      == max(table(salarios$husage)))
   moda.marido.calc <- names(moda.marido)</pre>
16 | moda.marido.calc <- as.numeric(moda.marido.calc)</pre>
17
   # Moda da idade das esposas = a idade 37 aparece 217 vezes
18
19 | moda.esposa <- subset(table(salarios$age), table(salarios$age) == max
       (table(salarios$age)))
20 moda.esposa.calc <- names(moda.esposa)
   moda.esposa.calc <- as.numeric(moda.esposa.calc)</pre>
```

Os maridos são 7,68% mais velhos, em média, que as esposas. 5,13%, considerando a mediana da amostra, e, a idade mais comum entre eles (moda) é 18,92% maior que a delas.

```
# Variância da idade dos maridos = 126,07
1
  var(salarios$husage)
2
3
  # Variância da idade das esposas = 99,75
4
5
  var(salarios$age)
6
7
  # Desvio padrão da idade dos maridos = 11,23
  sd(salarios$husage)
8
  # Desvio padrão da idade das esposas = 9,99
10
  sd(salarios$age)
11
12
   # Coeficiente de variação da idade dos maridos = 26,45%_
  (sd(salarios$husage)/mean(salarios$husage))*100
14
15
  # Coeficiente de variação da idade das esposas = 25,33%_
16
   (sd(salarios$age)/mean(salarios$age))*100
```

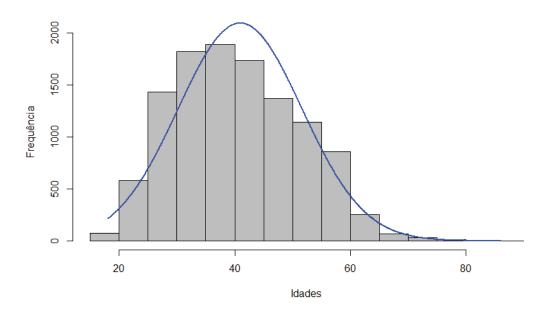
Em resumo, a variação de idade é um pouco maior entre os maridos, mas no geral, as idades tanto das esposas quanto dos maridos não variam muito dentro da amostra.

## 3. Testes paramétricos ou não paramétricos

```
# Agrupando os dados
1
2
   idades.df <- data.frame(</pre>
3
     group = rep(c("Esposas", "Maridos"), each = 5634),
     idades = c(salarios$age, salarios$husage)
4
5
6
7
   # Media e distancia interquartílica
8
   group_by(idades.df, group) %>%
9
     summarise(
10
       count = n(),
       mean = mean(idades, na.rm = TRUE),
11
12
       sd = sd(idades, na.rm = TRUE)
13
    )
14
15
   # Média e mediana do grupo
   mean(idades.df$idades)
17
   median(idades.df$idades)
18
19
   # Desvio padrao
20 sd(idades.df$idades)
21
22 # Histograma com curva normal
23 | plotNormalHistogram(idades.df$idades, prob = FALSE, main = "Distribui
      ção Normal sobreposta ao Histograma", length = 1000, xlab="Idades"
      , ylab = "Frequência" )
```

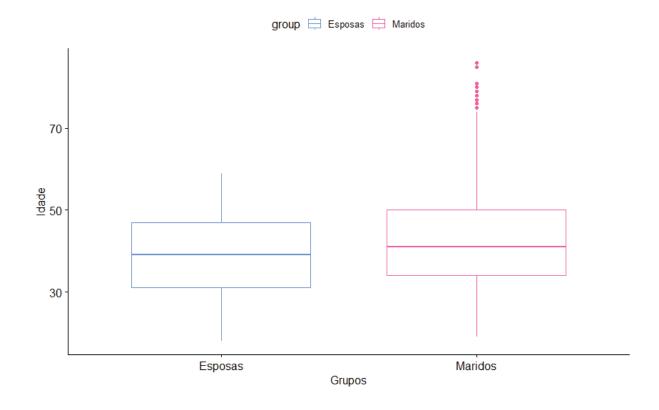
Figura 4.2 – EXERCÍCIO 3: DISTRIBUIÇÃO NORMAL E HISTOGRAMA

### Distribuição Normal sobreposta ao Histograma



```
lillie.test(idades.df$idades)
  D = 0.060083
3
4
   p-value < 2.2e-16
5
   # Sumario estatistico
6
7
   group_by(idades.df, group) %>%
8
     summarise(
9
       count = n(),
10
       median = median(idades, na.rm = TRUE),
11
       IQR = IQR(idades, na.rm = TRUE)
12
   # Box-plots
13
   ggboxplot(idades.df, x = "group", y = "idades", color = "group",
      palette=c("#6495ED", "#F25EA2"), ylab = "Idade", xlab = "Grupos")
```

Figura 4.3 – EXERCÍCIO 3: BOX-PLOTS



```
# Teste se a idade mediana dos maridos é igual a mediana das esposas wilcox.test(idades ~ group, data = idades.df, exact = FALSE, conf.int =TRUE)
```

Estatisticamente, as idades medianas dos maridos e das esposas são diferentes. O resultado do teste estatístico é altamente significativo (p-valor = 0,0000000000000000022), o que significa que podemos rejeitar a hipótese de medianas iguais. O intervalo de confiança da diferença entre as medianas também reforça essa conclusão, indicando que a mediana dos maridos é entre 2 e 3 anos menor que a mediana das esposas.

# APÊNDICE 5 - ESTATÍSTICA APLICADA II

### A - ENUNCIADO

### Regressões Ridge, Lasso e ElasticNet

(100 pontos) Fazer as regressões Ridge, Lasso e ElasticNet com a variável dependente "lwage" (salário-hora da esposa em logaritmo neperiano) e todas as demais variáveis da base de dados são variáveis explicativas (todas essas variáveis tentam explicar o salário-hora da esposa). No pdf você deve colocar a rotina utilizada, mostrar em uma tabela as estatísticas dos modelos (RMSE e  $\mathbb{R}^2$ ) e concluir qual o melhor modelo entre os três, e mostrar o resultado da predição com intervalos de confiança para os seguintes valores:

- husage = 40 (anos idade do marido)
- husunion = 0 (marido não possui união estável)
- husearns = 600 (US\$ renda do marido por semana)
- huseduc = 13 (anos de estudo do marido)
- husblck = 1 (o marido é preto)
- hushisp = 0 (o marido não é hispânico)
- hushrs = 40 (horas semanais de trabalho do marido)
- kidge6 = 1 (possui filhos maiores de 6 anos)
- age = 38 (anos idade da esposa)
- black = 0 (a esposa não é preta)
- educ = 13 (anos de estudo da esposa)
- hispanic = 1 (a esposa é hispânica)
- union = 0 (esposa não possui união estável)
- exper = 18 (anos de experiência de trabalho da esposa)
- kidlt6 = 1 (possui filhos menores de 6 anos)

Obs: lembre-se de que a variável dependente "lwage" já está em logarítmo, portanto voçê não precisa aplicar o logaritmo nela para fazer as regressões, mas é necessário aplicar o antilog para obter o resultado da predição.

# B - RESOLUÇÃO

### 1. Regressões Ridge, Lasso e ElasticNet

```
# Carregar os dados
setwd("C:/Users/marco/UFPR/IAA/IAA005-Estatística 2")
load("trabalhosalarios.RData")

# Verificar a estrutura dos dados
str(trabalhosalarios)
```

Tabela 5.1: Estrutura do data.frame (N = 2574, 16 variáveis)

variável	tipo	exemplos (início da coluna)
husage	num	56, 31, 33, 34, 42, 45, 33, 31, 31, 44,
husunion	num	$0, 0, 0, 0, 0, 0, 0, 0, 0, \dots$
husearns	num	1500, 800, 950, 1000, 730,
huseduc	num	14, 17, 13, 14, 14, 16, 16, 18, 12, 12,
husblck	num	$0, 0, 0, 0, 0, 0, 0, 0, 0, \dots$
hushisp	num	$0, 0, 0, 0, 0, 0, 0, 0, 0, \dots$
hushrs	num	40, 40, 60, 50, 40, 38, 40, 55, 40, 40,
kidge6	num	$1, 0, 0, 1, 1, 1, 0, 0, 0, 1, \dots$
age	num	49, 29, 30, 31, 41, 45, 32, 27, 30, 42,
black	num	$0, 0, 0, 0, 0, 0, 0, 0, 0, \dots$
educ	num	12, 14, 12, 12, 12, 18, 12, 14, 15, 12,
hispanic	num	$0, 0, 0, 0, 0, 0, 0, 0, 0, \dots$
union	num	$0, 0, 0, 0, 0, 0, 0, 0, 0, \dots$
exper	num	31, 9, 12, 13, 23, 21, 14, 7, 9, 24,
kidlt6	num	$0, 0, 1, 0, 0, 0, 0, 1, 1, 0, \dots$
lwage	num	1.90, 2.48, 2.43, 1.63, 2.30,

```
1 | # Separar a variável dependente (lwage) e as variáveis independentes
  y <- trabalhosalarios$lwage
3 | X <- trabalhosalarios[, setdiff(names(trabalhosalarios), "lwage")]</pre>
   # Instalar e carregar os pacotes necessários
5
6 | if(!require(glmnet)) install.packages("glmnet", dependencies=TRUE)
7 | if(!require(caret)) install.packages("caret", dependencies=TRUE)
8 | library(glmnet)
  library(caret)
9
10
  # Definir os parâmetros para os modelos
   train_control <- trainControl(method = "cv", number = 10)</pre>
12
13
14 # Modelo Ridge
15
  set.seed(123)
16 model_ridge <- train(</pre>
     lwage ~ ., data = trabalhosalarios, method = "glmnet",
17
   trControl = train_control,
18
     tuneGrid = expand.grid(alpha = 0, lambda = seq(0.0001, 1, length =
19
      100))
20
   )
21
22 # Modelo Lasso
23 | set.seed(123)
```

```
24 model_lasso <- train(
     lwage ~ ., data = trabalhosalarios, method = "glmnet",
25
26
     trControl = train_control,
     tuneGrid = expand.grid(alpha = 1, lambda = seq(0.0001, 1, length =
27
      100))
28
29
30
  # Modelo ElasticNet
31
  set.seed(123)
   model_elastic <- train(</pre>
32
     lwage ~ ., data = trabalhosalarios, method = "glmnet",
33
     trControl = train_control,
34
     tuneLength = 10
35
36
   )
37
38 # Comparar os modelos
   results <- resamples(list(Ridge = model_ridge, Lasso = model_lasso,
39
      ElasticNet = model_elastic))
40
   summary(results)
```

Tabela 5.2: MAE por modelo (10 reamostragens)

Model	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
Ridge	0.308 521 7	0.317 837 4	0.3337157	0.3307714	0.3404411	0.3604150	0
Lasso	0.309 219 9	0.3174310	0.333 607 5	0.3306203	0.3398527	0.3598834	0
ElasticNet	0.308 933 6	0.318 400 0	0.333 550 7	0.330 532 9	0.338 915 8	0.360 162 7	0

Tabela 5.3: RMSE por modelo (10 reamostragens)

Model	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
Ridge	0.3907595	0.409 507 5	0.434 672 9	0.441 576 2	0.448 293 3	0.549 436 2	0
Lasso	0.3916782	0.409 944 1	0.434 323 5	0.441 793 9	0.449 545 3	0.549 327 3	0
ElasticNet	0.391 631 7	0.4107487	0.433 751 9	0.441 500 8	0.4466197	0.5503250	0

Tabela 5.4: R<sup>2</sup> por modelo (10 reamostragens)

Model	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
Ridge	0.208 528 8	0.260 878 7	0.2702706	0.283 137 8	0.3204610	0.344 000 5	0
Lasso	0.2106588	0.261 484 9	0.2697880	0.2825007	0.3203778	0.3432997	0
ElasticNet	0.209 416 1	0.264 155 0	0.273 555 9	0.283 372 5	0.321 804 2	0.338 680 8	0

```
# Coletar as estatísticas dos modelos
rmse_ridge <- min(model_ridge$results$RMSE)
r2_ridge <- max(model_ridge$results$Rsquared)
mae_ridge <- min(model_ridge$results$MAE)

rmse_lasso <- min(model_lasso$results$RMSE)
r2_lasso <- max(model_lasso$results$Rsquared, na.rm = TRUE)</pre>
```

```
8
  mae_lasso <- min(model_lasso$results$MAE)</pre>
10 rmse_elastic <- min(model_elastic$results$RMSE)</pre>
  r2_elastic <- max(model_elastic$results$Rsquared)
11
   mae_elastic <- min(model_elastic$results$MAE)</pre>
13
14 | # Criar a tabela das estatísticas dos modelos
   stats <- data.frame(</pre>
15
   Model = c("Ridge", "Lasso", "ElasticNet"),
16
17
     RMSE = c(rmse_ridge, rmse_lasso, rmse_elastic),
   R2 = c(r2\_ridge, r2\_lasso, r2\_elastic),
18
19
     MAE = c(mae_ridge, mae_lasso, mae_elastic)
20
21
  print(stats)
```

Tabela 5.5: Comparativo de modelos

Model	RMSE	$\mathbb{R}^2$	MAE
Ridge	0.441 576 2	0.283 226 9	0.3307714
Lasso	0.441 793 9	0.2828095	0.3306203
ElasticNet	0.4415008	0.2833890	0.3305329

```
# Valores fornecidos para a predição
  new_data <- data.frame(</pre>
2
3
     husage = 40, husunion = 0, husearns = 600, huseduc = 13, husblck =
      1, hushisp = 0, hushrs = 40, kidge6 = 1, age = 38, black = 0, educ
       = 13, hispanic = 1, union = 0, exper = 18, kidlt6 = 1
4
5
6
  # Previsão com o melhor modelo (ElasticNet)
  pred <- predict(model_elastic, new_data)</pre>
7
8
9 | # Aplicar o antilogaritmo para obter o salário-hora
10
  pred_salary <- exp(pred)</pre>
11
12 | # Exibir o resultado
13 | print (pred_salary)
```

Resultado: 7.998988

# APÊNDICE 6 - ARQUITETURA DE DADOS

### A - ENUNCIADO

## 1 Construção de Características: Identificador automático de idioma

O problema consiste em criar um modelo de reconhecimento de padrões que dado um texto de entrada, o programa consegue classificar o texto e indicar a língua em que o texto foi escrito.

Parta do exemplo (notebook produzido no Colab) que foi disponibilidade e crie as funções para calcular as diferentes características para o problema da identificação da língua do texto de entrada.

Nessa atividade é para "construir características".

Meta: a acurácia deverá ser maior ou igual a 70%.

Essa tarefa pode ser feita no Colab (Google) ou no Jupiter, em que deverá exportar o notebook e imprimir o notebook para o formato PDF. Envie no UFPR Virtual os dois arquivos.

### 2 Melhore uma base de dados ruim

Escolha uma base de dados pública para problemas de classificação, disponível ou com origem na UCI Machine Learning.

Use o mínimo de intervenção para rodar a SVM e obtenha a matriz de confusão dessa base.

O trabalho começa aqui, escolha as diferentes tarefas discutidas ao longo da disciplina, para melhorar essa base de dados, até que consiga efetivamente melhorar o resultado.

Considerando a acurácia para bases de dados balanceadas ou quase balanceadas, se o percentual da acurácia original estiver em até 85%, a meta será obter 5%. Para bases com mais de 90% de acurácia, a meta será obter a melhora em pelo menos 2 pontos percentuais (92% ou mais).

Nessa atividade deverá ser entregue o script aplicado (o notebook e o PDF correspondente).

# B - RESOLUÇÃO

### 1 Construção de Características: Identificador automático de idioma

```
ingles = [
1
2
       "Hello, how are you?",
       "I love to read books.",
3
       "The weather is nice today.",
4
5
       "Where is the nearest restaurant?",
       "What time is it?",
6
       "I enjoy playing soccer.",
7
       "Can you help me with this?",
8
9
       "I'm going to the movies tonight.",
10
       "This is a beautiful place.",
       "I like listening to music.",
11
       "Do you speak English?",
12
13
       "What is your favorite color?",
       "I'm learning to play the guitar.",
14
       "Have a great day!",
15
```

```
16
       "I need to buy some groceries.",
17
       "Let's go for a walk.",
       "How was your weekend?",
18
19
       "I'm excited for the concert.",
       "Could you pass me the salt, please?",
20
       "I have a meeting at 2 PM.",
21
22
       "I'm planning a vacation.",
23
       "She sings beautifully.",
24
       "The cat is sleeping.",
25
       "I want to learn French.",
26
       "I enjoy going to the beach.",
27
       "Where can I find a taxi?",
28
       "I'm sorry for the inconvenience.",
29
       "I'm studying for my exams.",
30
       "I like to cook dinner at home.",
31
       "Do you have any recommendations for restaurants?",
32
33
   espanhol = [
34
       "Hola, cómo estás?",
35
       "Me encanta leer libros.",
36
       "El clima está agradable hoy.",
       "Dónde está el restaurante más cercano?",
37
38
       "Qué hora es?",
39
       "Voy al parque todos los días.",
40
       "Puedes ayudarme con esto?",
41
       "Me gustaría ir de vacaciones.",
42
       "Este es mi libro favorito.",
43
       "Me gusta bailar salsa.",
44
       "Hablas espaol?",
45
       "Cuál es tu comida favorita?",
46
       "Estoy aprendiendo a tocar el piano.",
47
       "Que tengas un buen día!",
48
       "Necesito comprar algunas frutas.",
49
       "Vamos a dar un paseo.",
50
       "Cómo estuvo tu fin de semana?",
51
       "Estoy emocionado por el concierto.",
52
       "Me pasas la sal, por favor?",
53
       "Tengo una reunión a las 2 PM.",
54
       "Estoy planeando unas vacaciones.",
55
       "Ella canta hermosamente.",
56
       "El perro está jugando.",
57
       "Quiero aprender italiano.",
58
       "Disfruto ir a la playa.",
59
       "Dónde puedo encontrar un taxi?",
60
       "Lamento las molestias.",
       "Estoy estudiando para mis exámenes.",
61
62
       "Me gusta cocinar la cena en casa.",
63
       "Tienes alguna recomendación de restaurantes?",
64
65
   portugues = [
       "Estou indo para o trabalho agora.",
```

```
67
        "Adoro passar tempo com minha família.",
        "Preciso comprar leite e pão.",
68
69
        "Vamos ao cinema no sábado.",
70
        "Gosto de praticar esportes ao ar livre.",
71
        "O trânsito está terrível hoje.",
72
        "A comida estava deliciosa!",
73
        "Você já visitou o Rio de Janeiro?",
74
        "Tenho uma reunião importante amanhã.",
75
        "A festa começa às 20h.",
76
        "Estou cansado depois de um longo dia de trabalho.",
77
        "Vamos fazer um churrasco no final de semana.",
78
        "O livro que estou lendo é muito interessante.",
79
        "Estou aprendendo a cozinhar pratos novos.",
80
        "Preciso fazer exercícios físicos regularmente.",
81
        "Vou viajar para o exterior nas férias.",
82
        "Você gosta de dançar?",
83
        "Hoje é meu aniversário!",
84
        "Gosto de ouvir música clássica.",
85
        "Estou estudando para o vestibular.",
86
        "Meu time de futebol favorito ganhou o jogo.",
87
        "Quero aprender a tocar violão.",
88
        "Vamos fazer uma viagem de carro.",
89
        "O parque fica cheio aos finais de semana.",
90
        "O filme que assisti ontem foi ótimo.",
91
        "Preciso resolver esse problema o mais rápido possível.",
92
        "Adoro explorar novos lugares.",
93
        "Vou visitar meus avós no domingo.",
        "Estou ansioso para as férias de verão.",
94
95
        "Gosto de fazer caminhadas na natureza.",
96
        "O restaurante tem uma vista incrível.",
97
        "Vamos sair para jantar no sábado."
98
99
100
   pre_padroes = []
101
   for frase in ingles:
102
   pre_padroes.append( [frase, 'inglês'])
103
104
   for frase in espanhol:
105
      pre_padroes.append( [frase, 'espanhol'])
106
107
   for frase in portugues:
108
    pre_padroes.append( [frase, 'português'])
109
110
   random.shuffle(pre_padroes)
111
112
   dados = pd.DataFrame(pre_padroes)
113
114
   def tamanhoMedioFrases(texto):
115
      palavras = re.split("\s", texto)
116
      tamanhos = [len(s) for s in palavras if len(s)>0]
117
```

```
118
    for t in tamanhos:
119
        soma=soma+t
120
      return soma / len(tamanhos)
121
122
    # Calcula a frequência relativa de caracteres acentuados no texto
123
    def frequenciaAcentos(texto):
        texto_lower = texto.lower()
124
125
        acentuados = 'ãoáéióúàâêôç'
        total_letras = sum(c.isalpha() for c in texto_lower)
126
127
        if total letras == 0:
128
            return 0
129
        count = sum(c in acentuados for c in texto_lower)
130
        return count / total_letras
131
    # Calcula a frequência relativa de vogais no texto
132
133
    def frequenciaVogais(texto):
134
        texto_lower = texto.lower()
135
        vogais = 'aeiou'
136
        total_letras = sum(c.isalpha() for c in texto_lower)
137
        if total_letras == 0:
138
            return 0
        count = sum(c in vogais for c in texto_lower)
139
140
        return count / total_letras
141
142
    def extraiCaracteristicas(frase):
143
      # frase é um vetor [ 'texto', 'lingua' ]
144
      texto = frase[0]
145
      pattern_regex = re.compile('[^\w+]', re.UNICODE)
146
      texto = re.sub(pattern_regex,' ',texto)
147
      caracteristical=tamanhoMedioFrases(texto)
148
      caracteristica2=frequenciaAcentos(texto)
149
      caracteristica3=frequenciaVogais(texto)
150
151
      # Contagem de palavras típicas de cada idioma
152
      palavras = texto.lower().split()
153
      caracteristica4 = sum(w in palavras_portugues for w in palavras)
154
      caracteristica5 = sum(w in palavras_ingles
                                                    for w in palavras)
155
      caracteristica6 = sum(w in palavras_espanhol for w in palavras)
156
157
      # acrescente as suas funcoes no vetor padrao
158
159
      padrao = [caracteristica1, caracteristica2, caracteristica3,
       caracteristica4, caracteristica5, caracteristica6, frase[1] ]
160
      return padrao
161
162
    # Conjuntos de palavras comuns em cada idioma
   palavras_portugues = {'que','de','em','um','uma','por','não','vou','j
163
       á','estou','é','no','nos','sua','dos','das','os','as','pra'}
   palavras_ingles = {'the','and','to','of','a','in','that','i','is','
164
       it','for', 'you','he','she','they','do','we','not','are','have','
       just','so'}
```

```
165 | palavras_espanhol = {'los', 'las', 'en', 'no', 'una', 'un', 'lo', 'eso', '
      muy','sin','pero','él','ella'}
166
167
   def geraPadroes(frases):
    padroes = []
168
169
     for frase in frases:
        padrao = extraiCaracteristicas(frase)
170
171
        padroes.append(padrao)
172
    return padroes
173
174
   # converte o formato [frase classe] em
175
    # [caracteristica_1, caracteristica_2,... caracteristica n, classe]
176
   padroes = geraPadroes(pre_padroes)
177
   dados = pd.DataFrame(padroes)
178
179
180 ##Treinando o modelo
181
182 | ####Separando o conjunto de treinamento do conjunto de testes
183
   vet = np.array(padroes)
   classes = vet[:,-1]
184
   padroes sem classe = vet[:,0:-1]
185
186
   X_train, X_test, y_train, y_test = train_test_split(
       padroes_sem_classe, classes, test_size=0.25, stratify=classes)
187
   ####Treinando o modelo com SVM
188
189
   treinador = svm.SVC() #algoritmo escolhido
190 | modelo = treinador.fit(X_train, y_train)
191
192 # score com os dados de treinamento
193
   |acuracia = modelo.score(X_train, y_train)
194
   y_pred = modelo.predict(X_train)
195
   cm = confusion_matrix(y_train, y_pred)
196
197
   # com dados de teste que não foram usados no treinamento
198 y_pred2 = modelo.predict(X_test)
199
   cm = confusion_matrix(y_test, y_pred2)
```

Tabela 6.1: Matriz de confusão - Treino (acurácia: 84,06%)

	predito					
verdadeiro	espanhol	inglês	português			
espanhol	19	0	3			
inglês	2	21	0			
português	5	1	18			

Tabela 6.2: Relatório de classificação — Treino

classe	precisão	revocação	f1-score	suporte
espanhol	0.73	0.86	0.79	22
inglês	0.95	0.91	0.93	23
português	0.86	0.75	0.80	24
accuracy		0.84		69
macro avg	0.85	0.84	0.84	69
weighted avg	0.85	0.84	0.84	69

Tabela 6.3: Matriz de confusão - Amostra (acurácia: 70%)

	predito					
verdadeiro	espanhol	inglês	português			
espanhol	5	1	2			
inglês	1	6	0			
português	3	0	5			

Tabela 6.4: Relatório de classificação - Amostra

classe	precisão	revocação	f1-score	suporte
espanhol	0.56	0.62	0.59	8
inglês	0.86	0.86	0.86	7
português	0.71	0.62	0.67	8
accuracy		0.70		23
macro avg	0.71	0.70	0.70	23
weighted avg	0.70	0.70	0.70	23

# APÊNDICE 7 - APRENDIZADO DE MÁQUINA

### A - ENUNCIADO

Para cada uma das tarefas abaixo (Classificação, Regressão etc.) e cada base de dados (Veículo, Diabetes etc.), fazer os experimentos com todas as técnicas solicitadas (KNN, RNA etc.) e preencher os quadros com as estatísticas solicitadas, bem como os resultados pedidos em cada experimento.

# B - RESOLUÇÃO

# Classificação - Veículos

Tabela 7.1: Veículos - Resultados por técnica com parâmetros, acurácia e matrizes de confusão.

Técnica	Parâmetro	Acurácia	Matriz de Confusão							
			Reference							
			Prediction	bus	opel	saab	van			
SVM – CV	C=100, Sigma=0.01	0.8443	bus	40	0	1	0			
3 V IVI – C V	C=100, Sigilia=0.01	0.0443	opel	0	31	10	0			
			saab	0	10	31	0			
			van	3	1	1	39			
			Reference							
			Prediction	bus	opel	saab	van			
RNA – CV	size=31, decay=0.7	0.8204	bus	41	0	0	0			
	size=31, decay=0.7	0.0204	opel	0	27	13	0			
			saab	0	13	30	0			
			van	2	0	0	39			
			Reference							
			Prediction	bus	opel	saab	van			
SVM – Hold-out	C=1, Sigma=0.07	0.7485	bus	40	0	0	1			
S v Ivi – Hold-out	C=1, Sigilia=0.07	0.7403	opel	0	20	13	0			
			saab	0	20	26	0			
			van	3	0	0	39			
			Reference							
			Prediction	bus	opel	saab	van			
RF – Hold-out		0.7425	bus	40	0	1	0			
Ki – Holu-out	mtry=10	0.7423	opel	0	19	13	0			
			saab	0	20	26	0			
			van	3	3	3	39			

Continua na próxima página

# Continuação da Tabela 7.1

Técnica	Parâmetro	Acurácia	Matriz de Confusão				
			Reference				
			Prediction	bus	opel	saab	van
RF – CV  RNA – Hold-out	mtry_5	0.7305	bus	40	0	1	0
	mtry=5	0.7303	opel	0	18	14	0
			saab	0	21	25	0
RNA – Hold-out			van	3	3	0	39
			Reference				
			Prediction	bus	opel	saab	van
DNA Hold out	size=5, decay=0.1	0.6707	bus	35	1	0	0
KNA – Hold-out	S12e=3, decay=0.1	0.0707	opel	0	0	0	0
			saab	2	9	43	5
			van	6	2	0	34
			Reference				
			Prediction	bus	opel	saab	van
VNIN	k=1	0.6353	bus	39	1	3	2
IZININ	K=1	0.0353	opel	1	18	19	2
			saab	7	23	15	0
			van	0	4	0	36

# Classificação - Diabetes

Tabela 7.2: Diabetes - Resultados por técnica com parâmetros, acurácia e matriz de confusão.

Técnica	Parâmetro	Acurácia	Matriz de	Confu	são	
			Reference			
CVM Hold out	C_0.5 Sigma_0.14	0.7647	Prediction	neg	pos	
SVM – Hold-out	C=0.5, Sigma=0.14	5, Sigma=0.14 0.7647 neg	92	28		
			pos	8	25	
			Reference			
RF – Hold-out	mtry=2	0.7647	Prediction	neg	pos	
		0.7047	neg	86	22	
			pos	14	31	
			Reference			
RF – CV	mtry=9	0.7647	Prediction	neg	pos	
KI'-CV	111ti y=9	0.7047	neg	83	19	
			pos	17	34	
			Reference			
SVM – CV	C-2 Sigmo-0.015	0.7582	Prediction	neg	pos	
3 V IVI — C V	C=2, Sigma=0.015	0.7362	neg	91	28	
			pos	9	25	

Continua na próxima página

# Continuação da Tabela 7.2

Técnica	Parâmetro	Acurácia	Matriz de Confusão			
			Reference			
KNN	k=7	0.7273	Prediction	neg pos 85 31 11 27 neg pos 81 26 19 27	pos	
KININ	K=/	0.7273	neg	85	31	
			pos	11	27	
			Reference			
RNA – CV	size=21, decay=0.4	0.7059	Prediction	neg	pos	
KIVA – C V	Size=21, decay=0.4	0.7039	neg	81	26	
			pos	19	27	
			Reference			
RNA – Hold-out	size=5, decay=0.1	0.6732	Prediction	neg	pos	
KINA – HOIU-OUL	\$126=3, decay=0.1	0.0732	neg	78	28	
			pos	22	25	

# Regressão - Admissão

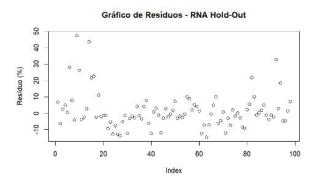
Tabela 7.3: Admissão - Resultados por técnica com parâmetros e métricas ( $R^2$ ,  $S_{yx}$ , coeficiente de Pearson, RMSE e MAE).

Técnica	Parâmetro	$R^2$	$S_{yx}$	Pearson	RMSE	MAE
RNA – Hold-out	size=3, decay=1e-4	0.8234	0.0606	0.9083	0.0583	0.0434
RF – Hold-out	mtry=2	0.8231	0.0606	0.9089	0.0584	0.0431
RF – CV	mtry=2	0.8196	0.0612	0.9072	0.0590	0.0436
SVM – CV	C=50, Sigma=0.01	0.8160	0.0618	0.9063	0.0595	0.0430
SVM – Hold-out	C=1, Sigma=0.18	0.8014	0.0642	0.9007	0.0619	0.0438
KNN	k=9	0.7384	0.0737	0.8600	0.0710	0.0547
RNA – CV	size=41, decay=0.1	0.7344	0.0743	0.8703	0.0716	0.0534

Tabela 7.4: Amostra com GRE, TOEFL, rating universitário e métricas preditivas.

GRE.Score	TOEFL.Score	<b>University.Rating</b>	SOP	LOR	CGPA	Research	predict.rna
316.99	110.72	4.66	1.19	2.20	9.70	0.71	0.89
303.68	99.41	2.67	4.11	2.01	8.08	0.41	0.59
295.93	97.60	4.72	3.48	2.84	9.80	0.80	0.82

Figura 7.1 – ADMISSÃO - GRÁFICO DE RESÍDUOS (MELHOR MODELO)



# Regressão - Biomassa

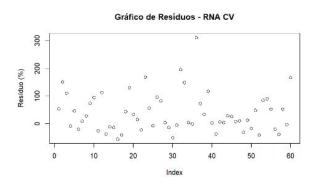
Tabela 7.5: Biomassa - Resultados por técnica com parâmetros e métricas ( $R^2$ ,  $S_{yx}$ , coeficiente de Pearson, RMSE e MAE).

Técnica	Parâmetro	$R^2$	$S_{yx}$	Pearson	RMSE	MAE
RNA – CV	size=11, decay=0.4	0.9877	141.04	0.9939	137.47	64.407
RF – CV	mtry=5	0.9846	163.73	0.9923	153.88	61.807
RF – Hold-out	mtry=2	0.9843	165.45	0.9921	155.50	63.321
SVM – CV	C=100, Sigma=0.01	0.9448	310.55	0.9881	291.87	94.531
KNN	k=3	0.8576	498.83	0.9578	468.83	124.52
RNA – Hold-out	size=5, decay=0.1	0.5433	861.62	0.8205	839.80	186.86
SVM – Hold-out	C=1, Sigma=0.60	0.3435	1071.3	0.6175	1006.89	210.16

Tabela 7.6: Amostra com variáveis dap, h, Me e predição da RNA.

dap	h	Me	predict.rna
132.31429	23.76053	0.4378101	16744.171
55.74442	48.14359	0.6876391	2724.243
45.74574	14.74192	0.3821269	14206.252

Figura 7.2 – BIOMASSA - GRÁFICO DE RESÍDUOS (MELHOR MODELO)



K-means clustering with 10 clusters of size 75, 67, 54, 131, 24, 137, 78, 109, 70, 101

Tabela 7.7: K-means (médias por cluster) - Parte 1

Cl	Comp	Circ	DCirc	RadRa	PrAxisRa	MaxLRa	ScatRa	Elong	PrAxisRect	MaxLRect
1	89.96	39.17	75.71	168.43	65.49	9.40	149.69	44.17	19.00	135.45
2	95.99	45.96	90.64	196.25	64.34	8.31	182.42	35.90	21.49	148.67
3	85.67	36.74	57.31	122.44	56.19	5.44	120.81	55.72	17.22	128.41
4	85.61	43.20	70.41	138.50	58.93	7.20	148.49	45.30	18.94	143.68
5	107.25	55.54	103.42	190.25	55.79	5.79	248.71	26.88	27.08	168.08
6	104.12	53.66	102.67	201.88	62.42	10.47	217.18	30.70	24.42	169.28
7	102.14	49.71	100.51	204.21	63.67	9.37	200.67	32.81	22.95	156.59
8	88.88	38.89	66.83	138.25	57.81	7.27	134.41	49.92	18.03	135.41
9	93.47	41.33	80.15	188.79	66.96	9.14	165.79	39.61	20.14	139.43
10	91.05	45.54	80.15	157.95	63.08	10.02	156.26	43.02	19.56	152.19

Tabela 7.8: K-means (médias por cluster) - Parte 2

Cl	ScVarMaxis	ScVarmaxis	RaGyr	SkewMaxis	Skewmaxis	Kurtmaxis	KurtMaxis	HolIRa
1	173.88	335.92	142.24	69.19	4.72	14.37	193.47	200.41
2	203.27	502.81	179.22	68.01	6.24	12.97	193.78	200.34
3	140.31	215.50	134.94	74.96	7.00	12.19	185.94	190.83
4	170.03	324.98	173.22	77.76	6.00	9.14	183.64	188.79
5	271.29	910.54	247.17	83.63	6.50	14.21	182.67	183.63
6	227.10	697.25	212.86	71.92	7.58	15.88	187.91	197.39
7	218.59	603.91	199.73	69.38	6.46	14.76	191.15	199.00
8	156.91	267.55	142.82	71.35	6.09	11.05	189.72	194.93
9	192.69	414.83	155.57	69.83	5.40	15.00	194.31	200.69
10	176.25	360.31	178.36	73.12	7.11	9.27	187.40	195.52

# Associação - Musculação

Tabela 7.9: Regras de associação — top 10 linhas.

Lhs	Rhs	Support	Confidence	Coverage	Lift	Count
{V2=Gêmeos}	{V3=Bicicleta}	0.026	1.0	0.026	3.2	1
{V2=Gêmeos}	{V4=Esteira}	0.026	1.0	0.026	3.2	1
{V2=Gêmeos}	{V1=Extensor}	0.026	1.0	0.026	1.0	1
{V3=Gêmeos}	{V2=AgachamentoSmith}	0.051	1.0	0.051	3.9	2
{V3=Gêmeos}	{V4=Esteira}	0.051	1.0	0.051	3.2	2
{V4=Crucifixo}	{V3=Afundo}	0.051	1.0	0.051	6.5	2
{V4=Crucifixo}	{V2=Gêmeos}	0.051	1.0	0.051	5.5	2
{V4=Crucifixo}	{V1=LegPress}	0.051	1.0	0.051	2.6	2
{V2=Flexor}	{V3=Bicicleta}	0.051	1.0	0.051	3.2	2
{V2=Flexor}	{V4=Esteira}	0.051	1.0	0.051	3.2	2

```
1 | install.packages("e1071")
  2 | install.packages("caret")
  3
       install.packages("Metrics")
  4
  5 | library("caret")
  6
       library("Metrics")
  7
  8
       % CLASSIFICACAO - VEICULOS
  9
10 | setwd("/Users/blb/Documents/ESPECIALIZAÇÃO EM IA/APRENDIZADO DE
               MAQUINA/CLASSIFICACAO/")
11
      |dados <- read.csv("6 - Veiculos - Dados.csv")
12
       dados[[a]] <- NULL
13
14 % KNN
15
16 set.seed(202410)
       ran <- sample(1:nrow(dados), 0.8 * nrow(dados))</pre>
18 | treino <- dados[ran,]</pre>
       teste <- dados[-ran,]</pre>
19
20 set.seed(202410)
21 | tuneGrid <- expand.grid(k = c(1, 3, 5, 7, 9))
22 knn <- train(tipo ~ ., data = treino, method = "knn", tuneGrid=
               tuneGrid)
23 | predict.knn <- predict(knn, teste)
24 | confusionMatrix(predict.knn, as.factor(teste[["tipo"]]))
25
26 % RNA HOLD OUT
27
28 | set.seed(202410)
29 | indices <- createDataPartition(dados[["tipo"]], p=0.80, list=FALSE)
30 | treino <- dados[indices,]
31 | teste <- dados[-indices,]
32 | set.seed(202410)
       rna <- train(tipo ~ ., data=treino, method="nnet",trace=FALSE)</pre>
33
34 | predicoes.rna <- predict(rna, teste)
35
       confusionMatrix(predicoes.rna, as.factor(teste[["tipo"]]))
36
37
      % RNA - Cross-Validation
38
39 | ctrl <- trainControl(method = "cv", number = 10)
40 | grid \leftarrow expand.grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid \leftarrow expand.grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid \leftarrow expand.grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid \leftarrow expand.grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid \leftarrow expand.grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40 | grid(size = seq(from = 1, to = 45, by = 10), decay = 40
               seq(from = 0.1, to = 0.9,
41 |by = 0.3)
       set.seed(202410)
42
43
       rna <- train(
44
                  form = tipo\sim.,
45
                  data = treino,
                 method = "nnet",
46
47
                  tuneGrid = grid,
                  trControl = ctrl,
48
```

```
49
       maxit = 2000,trace=FALSE)
50 | predicoes.rna <- predict(rna, teste)
51
  confusionMatrix(predicoes.rna, as.factor(teste[["tipo"]]))
53
  % SVM - Hold-out
54
55 | set.seed(202410)
56 | indices <- createDataPartition(dados[["tipo"]], p=0.80, list=FALSE)
57 | treino <- dados[indices,]
   teste <- dados[-indices,]</pre>
59 | set.seed(202410)
60 | svm <- train(tipo~., data=treino, method="svmRadial")
61 | predicoes.svm <- predict(svm, teste)
62 | confusionMatrix(predicoes.svm, as.factor(teste[["tipo"]]))
63
64 % SVM - Cross-Validation
65
66 set.seed(202410)
67 | indices <- createDataPartition(dados[["tipo"]], p=0.80, list=FALSE)
68 | treino <- dados[indices,]
69 teste <- dados[-indices,]
70 | ctrl <- trainControl(method = "cv", number = 10)
71 |tuneGrid = expand.grid(C=c(1, 2, 10, 50, 100), sigma=c(.01, .015,
      0.2))
72 | set.seed(202410)
73 | svm <- train(tipo~., data=treino, method="svmRadial", trControl=ctrl,
       tuneGrid=tuneGrid)
74 | predicoes.svm <- predict(svm, teste)
   confusionMatrix(predicoes.svm, as.factor(teste[["tipo"]]))
76
77
  % RANDOM FOREST - Hold-out
78
79 | dados_aux <- dados
   dados_aux[["tipo"]] <- NULL</pre>
81 | newdata <-lapply (dados_aux, function (coluna) {
82 | return(runif(3, min = min(coluna), max = max(coluna)))
83
   })
84 | predicoes.svm <- predict(svm, data.frame(newdata))
   resultado <- cbind(data.frame(newdata), predicoes.svm)</pre>
86 | View (resultado)
87
88 set.seed(202410)
89 | indices <- createDataPartition(dados[["tipo"]], p=0.80, list=FALSE)
90
   treino <- dados[indices,]</pre>
91 | teste <- dados[-indices,]
92 | set.seed(202410)
93 | rf <- train(tipo~., data=treino, method="rf")
94
95 | predicoes.rf <- predict(rf, teste)
96 | confusionMatrix(predicoes.rf, as.factor(teste[["tipo"]]))
97
```

```
98 % RANDOM FOREST - Cross-Validation
99
100 | set.seed(202410)
101 | indices <- createDataPartition(dados[["tipo"]], p=0.80, list=FALSE)
102
   treino <- dados[indices,]</pre>
103 | teste <- dados[-indices,]
104 ctrl <- trainControl(method = "cv", number = 10)
105 | tuneGrid = expand.grid(mtry=c(2, 5, 7, 9)) |
106 | set.seed(202410)
   rf <- train(tipo~., data=treino, method="rf", trControl=ctrl,
107
       tuneGrid=tuneGrid)
108 | predicoes.rf <- predict(rf, teste)</pre>
109
   confusionMatrix(predicoes.rf, as.factor(teste[["tipo"]]))
```

```
% CLASSIFICACAO - DIABETES
1
2
3
  setwd("/Users/blb/Documents/ESPECIALIZAÇÃO EM IA/APRENDIZADO DE
      MAQUINA/CLASSIFICACAO/")
4
   dados <- read.csv("10 - Diabetes - Dados.csv")</pre>
5
   dados[["num"]] <- NULL</pre>
6
7
   % KNN
8
9
  set.seed(202410)
10 | ran <- sample(1:nrow(dados), 0.8 * nrow(dados))</pre>
11 | treino <- dados[ran,]</pre>
12 | teste <- dados[-ran,]</pre>
   set.seed(202410)
14 | tuneGrid <- expand.grid(k = c(1, 3, 5, 7, 9))
15 knn <- train(diabetes ~ ., data = treino, method = "knn", tuneGrid=
      tuneGrid)
16 | predict.knn <- predict(knn, teste)</pre>
17
   confusionMatrix(predict.knn, as.factor(teste[["diabetes"]]))
18
19 | % RNA - Hold-out
20
21 | set.seed(202410)
22 | indices <- createDataPartition(dados[["diabetes"]], p=0.80, list=
      FALSE)
23 | treino <- dados[indices,]
24 | teste <- dados[-indices,]
25
  set.seed(202410)
  rna <- train(diabetes ~ ., data=treino, method="nnet",trace=FALSE)</pre>
  predicoes.rna <- predict(rna, teste)</pre>
   confusionMatrix(predicoes.rna, as.factor(teste[["diabetes"]]))
29
30 % RNA - Cross-Validation
31
32 ctrl <- trainControl(method = "cv", number = 10)
33
       grid \leftarrow expand.grid(size = seq(from = 1, to = 45, by = 10), decay
      = seq(from = 0.1, to = 0.9,
```

```
34
       by = 0.3)
35
       set.seed(202410)
36
       rna <- train(
37
       form = diabetes~.,
38
       data = treino,
39
       method = "nnet",
       tuneGrid = grid,
40
41
       trControl = ctrl,
42
       maxit = 2000,trace=FALSE)
43
   predicoes.rna <- predict(rna, teste)</pre>
   confusionMatrix(predicoes.rna, as.factor(teste[["diabetes"]]))
44
45
  % SVM - Hold-out
46
47
48 set.seed(202410)
49 | indices <- createDataPartition(dados[["diabetes"]], p=0.80, list=
      FALSE)
50 | treino <- dados[indices,]
51 | teste <- dados[-indices,]
52 set.seed(202410)
53 | svm <- train(diabetes~., data=treino, method="svmRadial")
   predicoes.svm <- predict(svm, teste)</pre>
  confusionMatrix(predicoes.svm, as.factor(teste[["diabetes"]]))
56
57
  dados_aux <- dados
58 | dados_aux[["diabetes"]] <- NULL
  newdata <-lapply(dados_aux, function(coluna) {</pre>
60 | return(runif(3, min = min(coluna), max = max(coluna)))
61
   })
62
  predicoes.svm <- predict(svm, data.frame(newdata))</pre>
63
  resultado <- cbind(data.frame(newdata), predicoes.svm)
64
  % SVM - Cross-Validation
65
66
67
  set.seed(202410)
68 indices <- createDataPartition(dados[["diabetes"]], p=0.80, list=
      FALSE)
69 | treino <- dados[indices,]
  teste <- dados[-indices,]</pre>
71
  |ctrl <- trainControl(method = "cv", number = 10)</pre>
0.2))
73 | set.seed(202410)
  svm <- train(diabetes~., data=treino, method="svmRadial", trControl=</pre>
      ctrl, tuneGrid=tuneGrid)
75 | predicoes.svm <- predict(svm, teste)
76 | confusionMatrix(predicoes.svm, as.factor(teste[["diabetes"]]))
77
  % RANDOM FOREST - Hold-out
78
79
80 | set.seed(202410)
```

```
81 | indices <- createDataPartition(dados[["diabetes"]], p=0.80, list=
      FALSE)
82 | treino <- dados[indices,]</pre>
83 | teste <- dados[-indices,]
   set.seed(202410)
85 | rf <- train(diabetes~., data=treino, method="rf")
86 | predicoes.rf <- predict(rf, teste)
   confusionMatrix(predicoes.rf, as.factor(teste[["diabetes"]]))
87
88
89
   % RANDOM FOREST - Cross-Validation
90
91 | set.seed(202410)
92 | indices <- createDataPartition(dados[["diabetes"]], p=0.80, list=
      FALSE)
93 | treino <- dados[indices,]
94 | teste <- dados[-indices,]
   ctrl <- trainControl(method = "cv", number = 10)
96 | tuneGrid = expand.grid(mtry=c(2, 5, 7, 9))
97 | set.seed(202410)
98
   rf <- train(diabetes~., data=treino, method="rf", trControl=ctrl,
      tuneGrid=tuneGrid)
99 | predicoes.rf <- predict(rf, teste)
   confusionMatrix(predicoes.rf, as.factor(teste[["diabetes"]]))
100
```

```
% REGRESSAO - ADMISSAO
1
2
  setwd("/Users/blb/Documents/ESPECIALIZAÇÃO EM IA/APRENDIZADO DE MÁ
      QUINA/REGRESSÃO/")
4
  dados <- read.csv("9 - Admissao - Dados.csv", header=T)
5
   dados[["num"]] <- NULL</pre>
6
7
  r2 <- function(predito, observado) {
   return(1 - (sum((predito-observado)^2) / sum((observado-mean(
8
      observado))^2)))
9
10
  syx <- function(predito, observado, num_variaveis) {</pre>
11
   return(sqrt((sum((predito-observado)^2) / (length(predito) -
12
      num_variaveis))))
13
14
15
   regression_metrics <- function(predito, observado){</pre>
16
       r2_aux <- r2(predito, observado)</pre>
17
       syx_aux <- syx(predito, observado, 7)</pre>
       pearson_aux <- cor(predito, observado, method='pearson')</pre>
18
19
       rsme_aux <- rmse(predito, observado)</pre>
20
       mae_aux <- mae(predito, observado)</pre>
21
       return(c(r2_aux, syx_aux, pearson_aux, rsme_aux, mae_aux))
22
23
24 | set.seed(202410)
```

```
25 | ind <- createDataPartition(dados[["ChanceOfAdmit"]], p=0.80, list =
      FALSE)
26 | treino <- dados[ind,]
27 | teste <- dados[-ind,]
   tuneGrid \leftarrow expand.grid(k = c(1,3,5,7,9))
29 | set.seed(202410)
30 knn <- train(ChanceOfAdmit ~ ., data = treino, method = "knn",
      tuneGrid=tuneGrid)
31 | predict.knn <- predict(knn, teste)
   regression_metrics(predict.knn,teste[["ChanceOfAdmit"]])
33
34 % RNA - Hold-out
35
36 | set.seed(202410)
37
  |indices <- createDataPartition(dados[["ChanceOfAdmit"]], p=0.80, list
      =FALSE)
38 | treino <- dados[indices,]
39 | teste <- dados[-indices,]
40 | set.seed(202410)
41 | rna <- train(ChanceOfAdmit ~ ., data=treino, method="nnet", trace=
      FALSE)
42 | predict.rna <- predict(rna, teste)
43 | regression_metrics(predict.rna, teste[["ChanceOfAdmit"]])
   residuos <- function(predito, observado){</pre>
  return(((predito - observado)/observado) * 100)
45
46 }
   plot(residuos(predict.rna, teste[["ChanceOfAdmit"]]), ylab='Resíduo
       (%)', main='Gráfico de Resíduo
   s - RNA Hold-Out')
48
49
50 dados_aux <- dados
   dados_aux[["ChanceOfAdmit"]] <- NULL</pre>
51
52 newdata <-lapply(dados_aux, function(coluna) {
53
       return(runif(3, min = min(coluna), max = max(coluna)))
54
   })
55
   predict.rna <- predict(rna, data.frame(newdata))</pre>
57
   resultado <- cbind(data.frame(newdata), predict.rna)</pre>
58
59 | RNA - Cross-Validation
60
61
  ctrl <- trainControl(method = "cv", number = 10)</pre>
  grid <- expand.grid(size = seq(from = 1, to = 45, by = 10), decay = 10
62
      seq(from = 0.1, to = 0.9,
63 | by = 0.3)
64
   set.seed(202410)
65
   rna <- train(</pre>
66
       form = ChanceOfAdmit~.,
67
       data = treino,
68
       method = "nnet",
69
       tuneGrid = grid,
```

```
70
        trControl = ctrl,
71
        maxit = 2000,trace=FALSE)
72
73
   predict.rna <- predict(rna, teste)</pre>
74
    regression_metrics(predict.rna, teste[["ChanceOfAdmit"]])
75
76 % SVM - Hold-out
77
78 | set.seed(202410)
79
   indices <- createDataPartition(dados[["ChanceOfAdmit"]], p=0.80, list</pre>
       =FALSE)
80 treino <- dados[indices,]
81 | teste <- dados[-indices,]
82 set.seed(202410)
83
   svm <- train(ChanceOfAdmit~., data=treino, method="svmRadial")</pre>
84 | predict.svm <- predict(svm, teste)
   regression_metrics(predict.svm,teste[["ChanceOfAdmit"]])
86
   |% SVM - Cross-Validation
87
88
89 | set.seed(202410)
90 | indices <- createDataPartition(dados[["ChanceOfAdmit"]], p=0.80, list
       =FALSE)
91 | treino <- dados[indices,]
92 | teste <- dados[-indices,]
93 ctrl <- trainControl(method = "cv", number = 10)
   tuneGrid = expand.grid(C=c(1, 2, 10, 50, 100), sigma=c(.01, .015, .015)
       0.2))
95 | set.seed(202410)
96 | svm <- train(ChanceOfAdmit~., data=treino, method="svmRadial",
       trControl=ctrl, tuneGrid=tuneG
97 | rid)
98 predict.svm <- predict(svm, teste)
    regression_metrics(predict.svm,teste[["ChanceOfAdmit"]])
100
   % RANDOM FOREST - Hold-out
101
102
103 | set.seed(202410)
104
   indices <- createDataPartition(dados[["ChanceOfAdmit"]], p=0.80, list</pre>
       =FALSE)
105 | treino <- dados[indices,]</pre>
106 | teste <- dados[-indices,]
107
   set.seed(202410)
   rf <- train(ChanceOfAdmit~., data=treino, method="rf")</pre>
108
109 | predict.rf <- predict(rf, teste)</pre>
110
   regression_metrics(predict.rf, teste[["ChanceOfAdmit"]])
111
112 % RANDOM FOREST - Cross-Validation
113
114 | set.seed(202410)
```

```
|% REGRESSAO - BIOMASSA
 1
 2
 3
   setwd("/Users/blb/Documents/ESPECIALIZAÇÃO EM IA/APRENDIZADO DE MÁ
      QUINA/REGRESSÃO/")
   dados <- read.csv("5 - Biomassa - Dados.csv", header=T)
 4
 5
   dados[["num"]] <- NULL</pre>
 6
 7
   % KNN
 8
 9
   r2 <- function(predito, observado) {
   return(1 - (sum((predito-observado)^2) / sum((observado-mean()))
10
      observado))^2)))
11
12
   syx <- function(predito, observado, num_variaveis) {</pre>
13
       return(sqrt((sum((predito-observado)^2) / (length(predito) -
      num_variaveis))))
14
15
   regression_metrics <- function(predito, observado){</pre>
16
17
       r2_aux <- r2(predito, observado)</pre>
18
       syx_aux <- syx(predito, observado, 3)</pre>
19
       pearson_aux <- cor(predito, observado, method='pearson')</pre>
       rsme_aux <- rmse(predito, observado)</pre>
20
21
       mae_aux <- mae(predito, observado)</pre>
22
       return(c(r2_aux, syx_aux, pearson_aux, rsme_aux, mae_aux))
23
24
25
   residuos <- function(predito, observado) {</pre>
       return(((predito - observado)/observado) * 100)
26
27
28
29
   set.seed(202410)
30 | ind <- createDataPartition(dados[["biomassa"]], p=0.80, list = FALSE)
31 | treino <- dados[ind,]
32
   teste <- dados[-ind,]</pre>
   |tuneGrid \leftarrow expand.grid(k = c(1,3,5,7,9))
33
34 set.seed(202410)
   knn <- train(biomassa ~ ., data = treino, method = "knn", tuneGrid=
35
      tuneGrid)
```

```
36 | predict.knn <- predict(knn, teste)
37
      regression_metrics(predict.knn,teste[["biomassa"]])
38
39
      % RNA - Hold-out
40
41
42 set.seed(202410)
43 | indices <- createDataPartition(dados[["biomassa"]], p=0.80, list=
              FALSE)
44 | treino <- dados[indices,]
45 | teste <- dados[-indices,]
46 set.seed(202410)
47 | rna <- train(biomassa ~ ., data=treino, method="nnet", trace=FALSE,
              linout=T)
48 | predict.rna <- predict(rna, teste)
49 regression_metrics(predict.rna, teste[["biomassa"]])
50
51 | RNA - Cross-Validation
52
53 | ctrl <- trainControl(method = "cv", number = 10)
|qrid| < -expand.qrid(size = seq(from = 1, to = 45, by = 10), decay = |qrid| < -expand.qrid(size = seq(from = 1, to = 45, by = 10), decay = |qrid| < -expand.qrid(size = seq(from = 1, to = 45, by = 10), decay = |qrid| < |qrid| 
              seq(from = 0.1, to = 0.9,
55 | by = 0.3)
56 set.seed(202410)
57 | rna <- train(
58 | form = biomassa~.,
      data = treino,
60 method = "nnet",
61
      tuneGrid = grid,
62 trControl = ctrl,
63 | maxit = 2000, trace=FALSE, linout=T)
64 | predict.rna <- predict(rna, teste)
65 regression_metrics(predict.rna,teste[["biomassa"]])
66
67 | dap = runif(3, min = min(dados[["dap"]]), max = max(dados[["dap"]]))
68 h = runif(3, min = min(dados[["h"]]), max = max(dados[["h"]]))
69 | Me = runif(3, min = min(dados[["Me"]]), max = max(dados[["Me"]]))
70 | newdata <- data.frame(dap=dap, h=h, Me=Me)
      predict.rna <- predict(rna, newdata)</pre>
72 | resultado <- cbind(newdata, predict.rna)
73
74 | % SVM - Hold-out
75
      set.seed(202410)
76
77 | indices <- createDataPartition(dados[["biomassa"]], p=0.80, list=
              FALSE)
78 | treino <- dados[indices,]
79 | teste <- dados[-indices,]
80 set.seed(202410)
81 | svm <- train(biomassa~., data=treino, method="svmRadial")
82 | predict.svm <- predict(svm, teste)
```

```
83
    regression_metrics(predict.svm, teste[["biomassa"]])
84
   % SVM - Cross-Validation
85
86
87
   set.seed(202410)
   indices <- createDataPartition(dados[["biomassa"]], p=0.80, list=</pre>
       FALSE)
89 | treino <- dados[indices,]
90 | teste <- dados[-indices,]
91
   ctrl <- trainControl(method = "cv", number = 10)
92 | tuneGrid = expand.grid(C=c(1, 2, 10, 50, 100), sigma=c(.01, .015,
       0.2))
93 | set.seed(202410)
94 | svm <- train(biomassa~., data=treino, method="svmRadial", trControl=
       ctrl, tuneGrid=tuneGrid)
95
   predict.svm <- predict(svm, teste)</pre>
96
    regression_metrics(predict.svm, teste[["biomassa"]])
97
98
   % RANDOM FOREST - Hold-out
99
100
   set.seed(202410)
101
   indices <- createDataPartition(dados[["biomassa"]], p=0.80, list=</pre>
       FALSE)
102 | treino <- dados[indices,]</pre>
103 | teste <- dados[-indices,]</pre>
104 | set.seed(202410)
105
   rf <- train(biomassa~., data=treino, method="rf")</pre>
106 | predict.rf <- predict(rf, teste)
   regression_metrics(predict.rf,teste[["biomassa"]])
107
108
109
   % RANDOM FOREST - Cross-Validation
110
   set.seed(202410)
111
    indices <- createDataPartition(dados[["biomassa"]], p=0.80, list=
112
       FALSE)
113 | treino <- dados[indices,]</pre>
114 teste <- dados[-indices,]</pre>
   ctrl <- trainControl(method = "cv", number = 10)
115
   tuneGrid = expand.grid(mtry=c(2, 5, 7, 9))
116
117
   set.seed(202410)
118 rf <- train(biomassa~., data=treino, method="rf", trControl=ctrl,
       tuneGrid=tuneGrid, linout=
119 T)
   predict.rf <- predict(rf, teste)</pre>
120
121 | regression_metrics(predict.rf, teste[["biomassa"]])
```

```
5 | dados[["a"]] <- NULL
6
   dados[["tipo"]] <- NULL</pre>
7
8 set.seed(202410)
   dados_normalizados <- scale(dados)</pre>
9
10 km.res = kmeans(dados, 10)
1 |% REGRAS DE ASSOCIACAO - MUSCULACAO
2
3 | install.packages('arules', dep=T)
4 library(arules)
5 | library (datasets)
6
7
   setwd("/Users/blb/Documents/ESPECIALIZAÇÃO EM IA/APRENDIZADO DE MÁ
      QUINA/ASSOCIAÇÃO/")
   dados <- read.csv("2 - Musculacao - Dados(in).csv", header=0, sep=';'</pre>
      )
9
10 set.seed(202410)
11 | rules <- apriori(dados, parameter = list(supp = 0.001, conf = 0.7,
      minlen=2))
12 | summary(rules)
13 options (digits=2)
```

14 | inspect(sort(rules[1:20], by="confidence"))

### APÊNDICE 8 – DEEP LEARNING

#### A - ENUNCIADO

# 1 Classificação de Imagens (CNN)

Implementar o exemplo de classificação de objetos usando a base de dados CIFAR10 e a arquitetura CNN vista no curso.

#### 2 Detector de SPAM (RNN)

Implementar o detector de spam visto em sala, usando a base de dados SMS Spam e arquitetura de RNN vista no curso.

#### 3 Gerador de Dígitos Fake (GAN)

Implementar o gerador de dígitos fake usando a base de dados MNIST e arquitetura GAN vista no curso.

### 4 Tradutor de Textos (Transformer)

i = Input(shape=x\\_train[0].shape)

 $5 \mid x = \text{Conv2D}(32, (3,3), \text{strides=2, activation="relu"})(i)$ 

3 4

Implementar o tradutor de texto do português para o inglês, usando a base de dados e a arquitetura Transformer vista no curso.

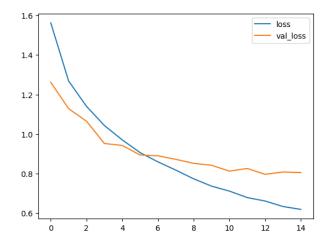
# B - RESOLUÇÃO

### 1 Classificação de Imagens

```
1 import tensorflow as tf
2 import numpy as np
3 import matplotlib.pyplot as plt
  from tensorflow.keras.layers import Input, Conv2D, Dense, Flatten,
      Dropout
  from tensorflow.keras.models import Model
  from mlxtend.plotting import plot_confusion_matrix
7
   from sklearn.metrics import confusion_matrix
8
9
  cifar10 = tf.keras.datasets.cifar10
10
   (x_train, y_train), (x_test, y_test) = cifar10.load_data()
11
12
   x_{train}, x_{test} = x_{train}/255.0, x_{test}/255.0
13 | y_train, y_test = y_train.flatten(), y_test.flatten()
14
15
  print("x_train.shape", x_train.shape)
16 | print("y_train.shape", y_train.shape)
17
  print("x_test.shape", x_test.shape)
18 | print("y_test.shape", y_test.shape)
   x_train.shape (50000, 32, 32, 3)
   y_train.shape (50000,)
   x_test.shape (10000, 32, 32, 3)
   y_test.shape (10000,)
1 \mid K = len(set(y \setminus train))
```

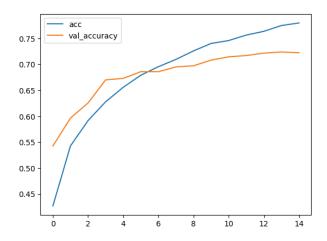
```
x = Conv2D(64, (3,3), strides=2, activation="relu")(x)
7
  x = Conv2D(128, (3,3), strides=2, activation="relu")(x)
8
9
  x = Flatten()(x)
10
   x = Dropout(0.5)(x)
  |x| = Dense(1024, activation="relu")(x)
12 \mid x = Dropout(0.2)(x)
13 x = Dense(K, activation="softmax")(x)
14 \mod = Model(i, x)
15
16 model.summary()
17
  model.compile(optimizer="adam", loss="sparse_categorical_crossentropy
18
      ", metrics=["accuracy"])
19
20
  r = model.fit(x_train, y_train, validation_data=(x_test, y_test),
      epochs=15)
21
22 plt.plot(r.history["loss"], label="loss")
23
  |plt.plot(r.history["val_loss"], label="val_loss")
24 plt.legend()
  plt.show()
```

Figura 8.1 – CLASSIFICAÇÃO DE IMAGENS COM CNN: GRÁFICO DE PERDA



```
plt.plot(r.history["accuracy"], label="acc")
plt.plot(r.history["val_accuracy"], label="val_accuracy")
plt.legend()
plt.show()
```





```
1  y_pred = model.predict(x_test).argmax(axis=1)
2  cm = confusion_matrix(y_test, y_pred)
plot_confusion_matrix(conf_mat=cm, figsize=(7,7), show_normed=True)
```

Figura 8.3 – CLASSIFICAÇÃO DE IMAGENS COM CNN: MATRIZ DE CONFUSÃO

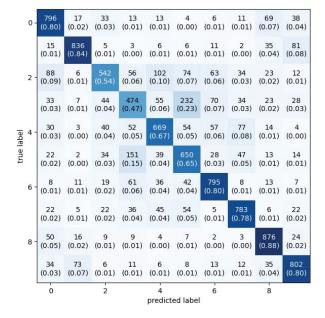
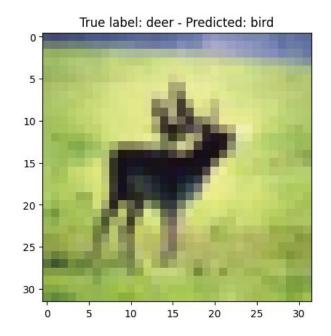


Figura 8.4 – CLASSIFICAÇÃO DE IMAGENS COM CNN: RESULTADO DA PREDIÇÃO



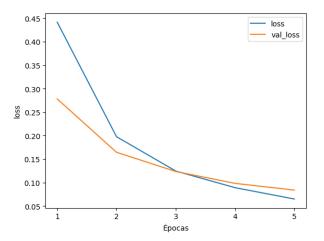
#### 2 Detector de SPAM

```
1 import tensorflow as tf
2 import numpy as np
3 import matplotlib.pyplot as plt
  import pandas as pd
  from sklearn.model_selection import train_test_split
  from tensorflow.keras.layers import Input, Embedding, LSTM, Dense
  from tensorflow.keras.layers import GlobalMaxPooling1D
7
  from tensorflow.keras.models import Model
   from tensorflow.keras.preprocessing.sequence import pad_sequences
9
10 from tensorflow.keras.preprocessing.text import Tokenizer
   from google.colab import files
  import io
12
13
14
  uploaded = files.upload()
15
  df = pd.read_csv(io.BytesIO(uploaded['spam.csv']), encoding="ISO
      -8859-1")
17 df.head()
18
19 df = df.drop(["Unnamed: 2", "Unnamed: 3", "Unnamed: 4"], axis=1)
  df.columns = ["labels", "data"]
  df["b_labels"] = df["labels"].map({"ham": 0, "spam": 1})
22
   y = df["b_labels"].values
23
  x_train, x_test, y_train, y_test = train_test_split(df["data"], y,
24
      test_size=0.33)
25
26 \mid \text{num\_words} = 2000
27 | tokenizer = Tokenizer(num_words=num_words)
```

```
28 tokenizer.fit_on_texts(x_train)
29 | sequences_train = tokenizer.texts_to_sequences(x_train)
30 | sequences_test = tokenizer.texts_to_sequences(x_test)
31 | word2index = tokenizer.word_index
   V = len(word2index)
33 | print("%s tokens" % V)
34
35 | data_train = pad_sequences(sequences_train)
36 T = data_train.shape[1]
   data_test = pad_sequences(sequences_test, maxlen=T)
38 | print("data_train.shape = ", data_train.shape)
39 print("data_test.shape = ", data_test.shape)
40
41 D = 20
   M = 5
42
43
44 | i = Input (shape=(T,))
45 \mid x = \text{Embedding}(V + 1, D)(i)
46 \mid x = LSTM(M)(x)
47
   x = Dense(1, activation="sigmoid")(x)
48
49 \mod = Model(i,x)
50
51 model.compile(
52
       loss="binary_crossentropy",
53
       optimizer="adam",
54
       metrics=["accuracy"]
55
56
57
  epochs = 5
58
59
   r = model.fit(
60
       data_train, y_train,
       validation_data=(data_test, y_test),
61
62
       epochs=epochs
63
64
65 | plt.plot(r.history["loss"], label="loss")
   plt.plot(r.history["val_loss"], label="val_loss")
67 | plt.xlabel("Épocas")
68 plt.ylabel("loss")
69 plt.xticks(np.arange(0, epochs, 1), labels=range(1, epochs + 1))
70 plt.legend()
71 plt.show()
```

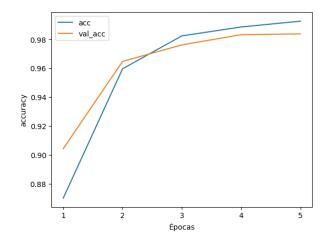
7242 tokens data\_train.shape = (3733, 175) data\_test.shape = (1839, 175)





```
plt.plot(r.history["accuracy"], label="acc")
plt.plot(r.history["val_accuracy"], label="val_acc")
plt.xlabel("Épocas")
plt.ylabel("accuracy")
plt.xticks(np.arange(0, epochs, 1), labels=range(1, epochs + 1))
plt.legend()
plt.show()
```

Figura 8.6 – DETECTOR DE SPAM COM RNN: GRÁFICO DE ACURÁCIA



```
texto = "Is your car dirty? Discover our new product. Free for all.
    Click the link."

seq_texto = tokenizer.texts_to_sequences([texto])
data_texto = pad_sequences(seq_texto, maxlen=T)

pred = model.predict(data_texto)
print(pred)
print(pred)
print("SPAM" if pred >= 0.5 else "OK")
```

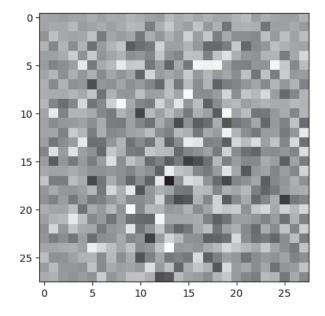
# [[0.6281646]] SPAM

3 Gerador de Dígitos Fake

```
%pip install imageio
2
  %pip install git+https://github.com/tensorflow/docs
3
4 import glob
  import imageio
   import numpy as np
6
7
  import os
  import PIL
8
9
  import time
10 import tensorflow as tf
   import matplotlib.pyplot as plt
11
12
   from tensorflow.keras import layers
13
14 from IPython import display
15
16
  (train_images, train_labels) , (_,_) = tf.keras.datasets.mnist.
      load_data()
  train_images = train_images.reshape(train_images.shape[0], 28, 28, 1)
17
      .astype('float32')
18
   train_images = (train_images - 127.5) / 127.5
19
20 | BUFFER_SIZE = 60000
21
  BATCH_SIZE = 256
22
23
   train_dataset = tf.data.Dataset.from_tensor_slices(train_images).
      shuffle(BUFFER_SIZE).batch(BATCH_SIZE)
24
25
   def make_generator_model():
   model = tf.keras.Sequential()
26
     model.add(layers.Dense(7*7*256, use_bias=False, input_shape=(100,))
27
      )
28
29
     model.add(layers.BatchNormalization())
30
     model.add(layers.LeakyReLU())
31
32
     model.add(layers.Reshape((7, 7, 256)))
33
     assert model.output_shape == (None, 7, 7, 256)
34
35
     model.add(layers.Conv2DTranspose(128, (5, 5), strides=1, padding='
      same', use_bias=False))
36
     assert model.output_shape == (None, 7, 7, 128)
37
38
     model.add(layers.BatchNormalization())
39
     model.add(layers.LeakyReLU())
40
41
     model.add(layers.Conv2DTranspose(64, (5, 5), strides=(2, 2),
      padding='same', use bias=False))
42
     assert model.output_shape == (None, 14, 14, 64)
43
44
     model.add(layers.BatchNormalization())
45
     model.add(layers.LeakyReLU())
```

```
46
47
     model.add(layers.Conv2DTranspose(1, (5, 5), strides=(2, 2), padding
      ='same', use_bias=False, activation='tanh'))
48
     assert model.output_shape == (None, 28, 28, 1)
49
50
     return model
51
52
   generator = make_generator_model()
   noise = tf.random.normal([1, 100])
53
   generated_image = generator(noise, training=False)
54
   plt.imshow(generated_image[0, :, :, 0], cmap='gray')
55
```

Figura 8.7 – GERANDO DADOS COM GAN: GERAÇÃO DE IMAGEM RUIDOSA

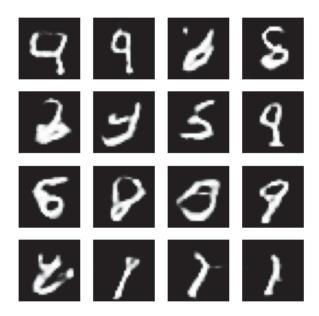


```
def make_discriminator_model():
1
2
     model = tf.keras.Sequential()
3
     model.add(layers.Conv2D(64, (5, 5), strides=(2, 2), padding='same',
       input_shape=[28, 28, 1]))
4
     model.add(layers.LeakyReLU())
5
     model.add(layers.Dropout(0.3))
6
7
     model.add(layers.Conv2D(128, (5, 5), strides=(2, 2), padding='same'
      ) )
8
     model.add(layers.LeakyReLU())
9
     model.add(layers.Dropout(0.3))
10
11
     model.add(layers.Flatten())
     model.add(layers.Dense(1))
12
13
14
     return model
15
   discriminator = make_discriminator_model()
16
  decision = discriminator(generated_image)
```

```
18
19
   cross_entropy = tf.keras.losses.BinaryCrossentropy(from_logits=True)
20
21
  def discriminator_loss(real_output, fake_output):
     real_loss = cross_entropy(tf.ones_like(real_output), real_output)
22
23
     fake_loss = cross_entropy(tf.zeros_like(fake_output), fake_output)
     total_loss = real_loss + fake_loss
24
25
     return total_loss
26
27
   def generator_loss(fake_output):
   return cross_entropy(tf.ones_like(fake_output), fake_output)
28
29
30
  generator_optimizer = tf.keras.optimizers.Adam(1e-4)
31
   discriminator_optimizer = tf.keras.optimizers.Adam(1e-4)
32
33
  checkpoint_dir = './training_checkpoints'
   checkpoint_prefix = os.path.join(checkpoint_dir, "ckpt")
34
   checkpoint = tf.train.Checkpoint(generator_optimizer=
35
      generator_optimizer,
36
                                     discriminator_optimizer=
      discriminator_optimizer,
37
                                     generator=generator,
                                     discriminator=discriminator)
38
39
40
  EPOCHS = 100
41
  noise dim = 100
42
   num_examples_to_generate = 16
43
44
   seed = tf.random.normal([num_examples_to_generate, noise_dim])
45
46
   @tf.function
   def train_step(images):
47
48
     noise = tf.random.normal([BATCH_SIZE, noise_dim])
49
50
     with tf.GradientTape() as gen_tape, tf.GradientTape() as disc_tape:
51
       generated_images = generator(noise, training=True)
52
53
       real_output = discriminator(images, training=True)
       fake_output = discriminator(generated_images, training=True)
54
55
56
       gen_loss = generator_loss(fake_output)
57
       disc_loss = discriminator_loss(real_output, fake_output)
58
59
       gradients_of_generator = gen_tape.gradient(gen_loss, generator.
      trainable_variables)
       gradients_of_discriminator = disc_tape.gradient(disc_loss,
60
      discriminator.trainable_variables)
61
62
       generator_optimizer.apply_gradients(zip(gradients_of_generator,
      generator.trainable_variables))
```

```
63
       discriminator_optimizer.apply_gradients(zip(
      gradients_of_discriminator, discriminator.trainable_variables))
64
65
   def train(dataset, epochs):
66
     for epoch in range (epochs):
67
       start = time.time()
68
69
       for image_batch in dataset:
70
         train_step(image_batch)
71
72
       display.clear_output (wait=True)
73
       generate_and_save_images(generator, epoch + 1, seed)
74
75
       if (epoch + 1) % 15 == 0:
76
         checkpoint.save(file_prefix = checkpoint_prefix)
77
78
       print('Time for epoch {} is {} sec' .format(epoch + 1, time.time
      () - start))
79
80
     display.clear_output (wait=True)
81
     generate_and_save_images(generator, epochs, seed)
82
83
   def generate_and_save_images(model, epoch, test_input):
84
     predictions = model(test_input, training=False)
85
     fig = plt.figure(figsize=(4, 4))
86
     for i in range(predictions.shape[0]):
87
       plt.subplot(4, 4, i + 1)
88
89
       plt.imshow(predictions[i, :, :, 0] \star 127.5 + 127.5, cmap='qray')
90
       plt.axis('off')
91
92
     plt.savefig('image_at_epoch_{:04d}.png'.format(epoch))
93
     plt.show()
94
95
  train(train_dataset, EPOCHS)
   checkpoint.restore(tf.train.latest_checkpoint(checkpoint_dir))
96
```

Figura 8.8 – GERANDO DADOS COM GAN: CHECKPOINT DA GERAÇÃO DE DÍGITOS



```
def display_image(epoch_no):
1
2
   return PIL.Image.open('image_at_epoch_{:04d}.png'.format(epoch_no))
3
4
  display_image(EPOCHS)
5
6
  anim_file = 'dcgan.gif'
7
8
  with imageio.get_writer(anim_file, mode='I') as writer:
9
     filenames = glob.glob('image*.png')
   filenames = sorted(filenames)
10
     for filename in filenames:
11
12
       image = imageio.imread(filename)
13
       writer.append_data(image)
14
    image = imageio.imread(filename)
15
     writer.append_data(image)
16
17
  import tensorflow_docs.vis.embed as embed
  embed.embed_file(anim_file)
```

### 4 Transformer

```
!pip uninstall tensorflow
!pip install tensorflow==2.15.0

!pip install tensorflow_datasets
!pip install -U tensorflow-text==2.15.0
import collections
import logging
import os
import pathlib
import re
import string
import sys
```

```
12 import time
13
  import numpy as np
   import matplotlib.pyplot as plt
  import tensorflow_datasets as tfds
   import tensorflow_text as text
17
   import tensorflow as tf
  logging.getLogger('tensorflow').setLevel(logging.ERROR)
18
19
  examples, metadata = tfds.load('ted_hrlr_translate/pt_to_en',
21
   with_info=True, as_supervised=True)
22
23
  train_examples, val_examples = examples['train'], examples['
      validation']
24
25
   for pt_examples, en_examples in train_examples.batch(3).take(1):
26
   for pt in pt_examples.numpy():
27
       print(pt.decode('utf-8'))
28
29
     print()
30
31
     for en in en_examples.numpy():
32
       print(en.decode('utf-8'))
```

e quando melhoramos a procura , tiramos a única vantagem da impressão , que é a serendipidade .

mas e se estes fatores fossem ativos ?

mas eles não tinham a curiosidade de me testar .

and when you improve searchability, you actually take away the one

advantage of print, which is serendipity.

but what if it were active?

but they did n't test for curiosity.

```
# Tokenização e Destokenização do texto
2
  model_name = "ted_hrlr_translate_pt_en_converter"
3
  tf.keras.utils.get_file(f"{model_name}.zip",
4
5
  f"https://storage.googleapis.com/download.tensorflow.org/models/{
      model_name}.zip", cache_dir='.', cache_subdir='', extract=True)
   # Tem 2 tokenizers: um pt outro em en
6
7
  # tokenizers.en tokeniza e detokeniza
  tokenizers = tf.saved_model.load(model_name)
8
9
  # PIPELINE DE ENTRADA
10
   # Codificar/tokenizar lotes de texto puro
11
12
  def tokenize_pairs(pt, en):
    pt = tokenizers.pt.tokenize(pt)
13
14
    # Converte ragged (irregular, tam variável) para dense
15
     # Faz padding com zeros.
   pt = pt.to_tensor()
16
17
     en = tokenizers.en.tokenize(en)
18
19
     # ragged -> dense
```

```
20
   en = en.to_tensor()
21
22
    return pt, en
23
24
   # Pipeline simpes: processa, embaralha, agrupa os dados, prefetch
25
   # Datasets de entrada terminam com prefetch
26 BUFFER_SIZE = 20000
27
   BATCH_SIZE = 64
28
29
   def make batches(ds):
30
   return (
31
       ds
32
       .cache()
33
       .shuffle(BUFFER_SIZE)
34
       .batch(BATCH_SIZE)
35
       .map(tokenize_pairs, num_parallel_calls=tf.data.AUTOTUNE)
36
       .prefetch(tf.data.AUTOTUNE))
37
38
   train_batches = make_batches(train_examples)
39
   val_batches = make_batches(val_examples)
40
   # CODIFICAÇÃO POSICIONAL
41
42
   def get_angles(pos, i, d_model):
43
     angle_rates = 1 / \text{np.power}(10000, (2 * (i//2)) / \text{np.float32}(d_model)
      ) )
44
     return pos * angle_rates
45
   def positional_encoding(position, d_model):
46
47
     angle_rads = get_angles(np.arange(position)[:, np.newaxis], np.
      arange(d_model)[np.newaxis, :], d_model)
48
49
     # sin em índices pares no array; 2i
50
     angle_rads[:, 0::2] = np.sin(angle_rads[:, 0::2])
51
52
    # cos em índices ímpares no array; 2i+1
53
     angle_rads[:, 1::2] = np.cos(angle_rads[:, 1::2])
54
55
     # newaxis, aumenta a dimensão [] -> [ [] ]
56
     pos_encoding = angle_rads[np.newaxis, ...]
57
58
     return tf.cast(pos_encoding, dtype=tf.float32)
59
   # CODIFICAÇÃO POSICIONAL
60
   n, d = 2048, 512
61
  pos_encoding = positional_encoding(n, d)
63
   print (pos_encoding.shape)
  pos_encoding = pos_encoding[0]
64
65
66 # Arrumar as dimensões
67 pos_encoding = tf.reshape(pos_encoding, (n, d//2, 2))
68 pos_encoding = tf.transpose(pos_encoding, (2, 1, 0))
```

```
pos_encoding = tf.reshape(pos_encoding, (d, n))

plt.pcolormesh(pos_encoding, cmap='RdBu')

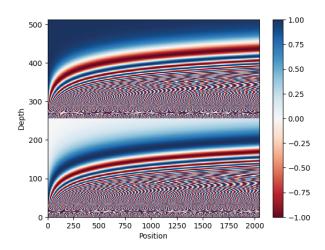
plt.ylabel('Depth')

plt.xlabel('Position')

plt.colorbar()

plt.show()
```

Figura 8.9 – TRANFORMER: CODIFICAÇÃO POSICIONAL



```
# Cria uma máscara de 0 e 1, 0 para quando há valor e 1 quando não há
1
2
   def create_padding_mask(seq):
3
     seq = tf.cast(tf.math.equal(seq, 0), tf.float32)
4
5
     # add extra dimensions to add the padding
6
   # to the attention logits.
7
     return seq[:, tf.newaxis, tf.newaxis, :] # (batch_size, 1, 1, seq_1
8
9
   # Máscara futura, usada no decoder
10
   def create_look_ahead_mask(size):
11
     # zera o triângulo inferior
     mask = 1 - tf.linalg.band_part(tf.ones((size, size)), -1, 0)
12
13
     return mask # (seq_len, seq_len)
14
15
   # Função de Atenção
16
   def scaled_dot_product_attention(q, k, v, mask):
17
     # O K^T
18
     matmul_qk = tf.matmul(q, k, transpose_b=True) # (..., seq_len_q,
      seq_len_k)
19
     # converte matmul_qk para float32
20
21
     dk = tf.cast(tf.shape(k)[-1], tf.float32)
22
23
     # divide por sqrt(d_k)
     scaled_attention_logits = matmul_qk / tf.math.sqrt(dk)
24
25
```

```
26
   # Soma a máscara, e os valores faltantes serão um número próximo a
      -inf
27
     if mask is not None:
28
       scaled_attention_logits += (mask * -1e9)
29
     # softmax normaliza os dados, soman 1. // (..., seq_len_q,
30
      seq_len_k)
31
     attention_weights = tf.nn.softmax(scaled_attention_logits, axis=-1)
32
33
     output = tf.matmul(attention_weights, v) # (..., seq_len_q, depth_v
      )
34
35
     return output, attention_weights
```

```
# Atenção Multi-cabeças
1
   class MultiHeadAttention(tf.keras.layers.Layer):
2
3
     def __init__(self, d_model, num_heads):
       super(MultiHeadAttention, self).__init__()
4
5
       self.num_heads = num_heads
6
       self.d_model = d_model
7
8
       assert d_model % self.num_heads == 0
9
       self.depth = d_model // self.num_heads
10
11
12
       self.wq = tf.keras.layers.Dense(d_model)
13
       self.wk = tf.keras.layers.Dense(d_model)
       self.wv = tf.keras.layers.Dense(d_model)
14
15
       self.dense = tf.keras.layers.Dense(d_model)
16
17
     def split_heads(self, x, batch_size):
18
19
       """Separa a última dimensão em (num_heads, depth).
20
       Transpõe o resultado para o shape (batch_size, num_heads, seq_len
      , depth)
21
       11 11 11
22
       x = tf.reshape(x, (batch_size, -1, self.num_heads, self.depth))
23
24
       return tf.transpose(x, perm=[0, 2, 1, 3])
25
26
     def call(self, v, k, q, mask):
27
       batch\_size = tf.shape(q)[0]
28
29
       q = self.wq(q) # (batch_size, seq_len, d_model)
30
       k = self.wk(k) # (batch_size, seq_len, d_model)
31
       v = self.wv(v) # (batch_size, seq_len, d_model)
32
33
       q = self.split_heads(q, batch_size) # (batch_size, num_heads,
      seq_len_q, depth)
34
       k = self.split_heads(k, batch_size) # (batch_size, num_heads,
      seq_len_k, depth)
```

```
seq_len_v, depth)
36
37
       # Calcula a atenção para cada cabeça (de forma matricial)
       # scaled_attention.shape == (batch_size, num_heads, seq_len_q,
38
      depth)
       # attention_weights.shape == (batch_size, num_heads, seq_len_q,
39
      seq_len_k)
40
       scaled_attention, attention_weights =
      scaled_dot_product_attention(q, k, v, mask)
41
42
       # Troca a dimensão 2 com 1, para acertar o num_heads
43
       # (batch_size, seq_len_q, num_heads, depth)
44
       scaled_attention = tf.transpose(scaled_attention, perm=[0, 2, 1,
      31)
45
       # Concatena os valores em: (batch_size, seq_len_q, d_model)
46
       concat_attention = tf.reshape(scaled_attention, (batch_size, -1,
47
      self.d_model))
48
49
       output = self.dense(concat_attention) # (batch_size, seq_len_q,
      d model)
50
51
       return output, attention_weights
52
   def point_wise_feed_forward_network(d_model, dff):
53
54
     return tf.keras.Sequential([
       tf.keras.layers.Dense(dff, activation='relu'), # (batch_size,
55
      seq_len, dff)
56
       tf.keras.layers.Dense(d_model) # (batch_size, seq_len, d_model)
57
     ])
   class EncoderLayer(tf.keras.layers.Layer):
1
2
     def __init__(self, d_model, num_heads, dff, rate=0.1):
3
       super(EncoderLayer, self).__init__()
4
5
       self.mha = MultiHeadAttention(d_model, num_heads)
6
       self.ffn = point_wise_feed_forward_network(d_model, dff)
7
8
       self.layernorm1 = tf.keras.layers.LayerNormalization(epsilon=1e
9
       self.layernorm2 = tf.keras.layers.LayerNormalization(epsilon=1e
      -6)
10
       self.dropout1 = tf.keras.layers.Dropout(rate)
11
12
       self.dropout2 = tf.keras.layers.Dropout(rate)
13
14
     def call(self, x, training, mask):
15
       attn_output, _ = self.mha(x, x, x, mask) # (batch_size,
      input_seq_len, d_model)
       attn_output = self.dropout1(attn_output, training=training)
16
```

v = self.split\_heads(v, batch\_size) # (batch\_size, num\_heads,

35

```
class DecoderLayer(tf.keras.layers.Layer):
1
2
     def __init__(self, d_model, num_heads, dff, rate=0.1):
3
       super(DecoderLayer, self).__init__()
4
5
       self.mha1 = MultiHeadAttention(d model, num heads)
       self.mha2 = MultiHeadAttention(d_model, num_heads)
6
7
8
       self.ffn = point_wise_feed_forward_network(d_model, dff)
9
10
       self.layernorm1 = tf.keras.layers.LayerNormalization(epsilon=1e
11
       self.layernorm2 = tf.keras.layers.LayerNormalization(epsilon=1e
       self.layernorm3 = tf.keras.layers.LayerNormalization(epsilon=1e
12
13
       self.dropout1 = tf.keras.layers.Dropout(rate)
14
15
       self.dropout2 = tf.keras.layers.Dropout(rate)
       self.dropout3 = tf.keras.layers.Dropout(rate)
16
17
18
     def call(self, x, enc_output, training, look_ahead_mask,
      padding_mask):
19
       # enc_output.shape == (batch_size, input_seq_len, d_model)
       # (batch_size, target_seq_len, d_model)
20
21
       attn1, attn_weights_block1 = self.mha1(x, x, x, look_ahead_mask)
       attn1 = self.dropout1(attn1, training=training)
22
23
       out1 = self.layernorm1(attn1 + x)
24
25
       # (batch_size, target_seq_len, d_model)
       attn2, attn_weights_block2 = self.mha2(enc_output, enc_output,
26
      out1, padding_mask)
       attn2 = self.dropout2(attn2, training=training)
27
28
       out2 = self.layernorm2(attn2 + out1) # (batch_size,
      target seg len, d model)
29
       ffn_output = self.ffn(out2) # (batch_size, target_seq_len,
30
      d_model)
31
       ffn_output = self.dropout3(ffn_output, training=training)
32
       out3 = self.layernorm3(ffn_output + out2) # (batch_size,
      target_seq_len, d_model)
```

```
33 | return out3, attn_weights_block1, attn_weights_block2
```

```
class Encoder(tf.keras.layers.Layer):
1
     def __init__(self, num_layers, d_model, num_heads, dff,
2
      input_vocab_size, maximum_position_encoding, rate=0.1):
3
       super(Encoder, self).__init__()
4
5
       self.d_model = d_model
6
       self.num_layers = num_layers
7
       self.embedding = tf.keras.layers.Embedding(input_vocab_size,
      d model)
8
       self.pos_encoding = positional_encoding(maximum_position_encoding
      , self.d_model)
9
       self.enc_layers = [EncoderLayer(d_model, num_heads, dff, rate)
      for in range (num layers) ]
       self.dropout = tf.keras.layers.Dropout(rate)
10
11
12
    def call(self, x, training, mask):
13
       seq_len = tf.shape(x)[1]
14
15
       # adding embedding and position encoding.
       x = self.embedding(x) # (batch_size, input_seq_len, d_model)
16
17
       x *= tf.math.sqrt(tf.cast(self.d_model, tf.float32))
       x += self.pos_encoding[:, :seq_len, :]
18
19
       x = self.dropout(x, training=training)
20
       for i in range(self.num_layers):
21
22
         x = self.enc_layers[i](x, training, mask)
23
24
       return x # (batch_size, input_seq_len, d_model)
```

```
1
   class Decoder(tf.keras.layers.Layer):
     def __init__(self, num_layers, d_model, num_heads, dff,
2
      target_vocab_size, maximum_position_encoding, rate=0.1):
3
       super(Decoder, self).__init__()
4
5
       self.d_model = d_model
6
       self.num_layers = num_layers
7
8
       self.embedding = tf.keras.layers.Embedding(target_vocab_size,
      d_model)
9
       self.pos_encoding = positional_encoding(maximum_position_encoding
      , d_model)
10
       self.dec_layers = [DecoderLayer(d_model, num_heads, dff, rate)
11
      for _ in range(num_layers)]
12
       self.dropout = tf.keras.layers.Dropout(rate)
13
     def call(self, x, enc_output, training, look_ahead_mask,
14
      padding_mask):
```

```
15
       seq_len = tf.shape(x)[1]
16
       attention_weights = {}
17
18
       x = self.embedding(x) # (batch_size, target_seq_len, d_model)
19
       x *= tf.math.sqrt(tf.cast(self.d model, tf.float32))
20
       x += self.pos_encoding[:, :seq_len, :]
21
22
       x = self.dropout(x, training=training)
23
24
       for i in range(self.num_layers):
25
         x, block1, block2 = self.dec_layers[i](x, enc_output, training,
       look_ahead_mask, padding_mask)
26
27
         attention_weights[f'decoder_layer{i+1}_block1'] = block1
         attention_weights[f'decoder_layer{i+1}_block2'] = block2
28
29
       # x.shape == (batch_size, target_seq_len, d_model)
30
       return x, attention_weights
31
1
   class Transformer(tf.keras.Model):
2
    def __init__(self, num_layers, d_model, num_heads, dff,
      input_vocab_size, target_vocab_size, pe_input, pe_target, rate
      =0.1):
3
       super().__init__()
4
       self.encoder = Encoder(num_layers, d_model, num_heads, dff,
      input_vocab_size,
5
       pe_input, rate)
       self.decoder = Decoder(num_layers, d_model, num_heads, dff,
6
      target_vocab_size,
7
       pe_target, rate)
8
       self.final_layer = tf.keras.layers.Dense(target_vocab_size)
9
10
     def call(self, inputs, training):
11
       # Keras models prefer if you pass all your inputs in the first
      argument
12
       inp, tar = inputs
13
       enc_padding_mask, look_ahead_mask, dec_padding_mask = self.
14
      create_masks(inp, tar)
15
16
       # (batch_size, inp_seq_len, d_model)
17
       enc_output = self.encoder(inp, training, enc_padding_mask)
18
19
       # dec_output.shape == (batch_size, tar_seq_len, d_model)
20
       dec_output, attention_weights = self.decoder(tar, enc_output,
      training, look_ahead_mask, dec_padding_mask)
21
22
       # (batch_size, tar_seq_len, target_vocab_size)
23
       final_output = self.final_layer(dec_output)
24
25
       return final_output, attention_weights
```

```
26
27
     def create_masks(self, inp, tar):
28
       # Encoder padding mask
29
       enc_padding_mask = create_padding_mask(inp)
       # Used in the 2nd attention block in the decoder.
30
31
       # This padding mask is used to mask the encoder outputs.
       dec_padding_mask = create_padding_mask(inp)
32
33
       # Used in the 1st attention block in the decoder.
34
       # It is used to pad and mask future tokens in the input received
35
       # the decoder.
36
       look_ahead_mask = create_look_ahead_mask(tf.shape(tar)[1])
37
       dec_target_padding_mask = create_padding_mask(tar)
38
       look_ahead_mask = tf.maximum(dec_target_padding_mask,
      look_ahead_mask)
39
40
       return enc_padding_mask, look_ahead_mask, dec_padding_mask
41
42
   # Hiperparâmetros
43
   num_layers = 4
  d \mod el = 128
   dff = 512
45
46 \mid \text{num heads} = 8
47
   dropout_rate = 0.1
1
2
   class CustomSchedule(tf.keras.optimizers.schedules.
      LearningRateSchedule):
     def __init__(self, d_model, warmup_steps=4000):
3
       super(CustomSchedule, self).__init__()
4
5
       self.d_model = d_model
       self.d_model = tf.cast(self.d_model, tf.float32)
6
7
       self.warmup_steps = warmup_steps
8
9
     def __call__(self, step):
10
       step = tf.cast(step, tf.float32) # Adicionado para evitar ERRO
11
       arg1 = tf.math.rsqrt(step)
       arg2 = step * (self.warmup_steps ** -1.5)
12
13
       return tf.math.rsqrt(self.d_model) * tf.math.minimum(arg1, arg2)
1
2
   learning_rate = CustomSchedule(d_model)
3
   optimizer = tf.keras.optimizers.Adam(learning_rate, beta_1=0.9,
      beta_2=0.98, epsilon=1e-9)
4
5
   loss_object = tf.keras.losses.SparseCategoricalCrossentropy(
      from_logits=True, reduction='none')
6
   def loss_function(real, pred):
7
8
     mask = tf.math.logical_not(tf.math.equal(real, 0))
     loss_ = loss_object(real, pred)
```

```
10
     mask = tf.cast(mask, dtype=loss_.dtype)
11
     loss\_ \star= mask
12
     return tf.reduce_sum(loss_)/tf.reduce_sum(mask)
13
14
   def accuracy_function(real, pred):
15
     accuracies = tf.equal(real, tf.argmax(pred, axis=2))
     mask = tf.math.logical_not(tf.math.equal(real, 0))
16
17
     accuracies = tf.math.logical_and(mask, accuracies)
18
     accuracies = tf.cast(accuracies, dtype=tf.float32)
19
     mask = tf.cast(mask, dtype=tf.float32)
     return tf.reduce_sum(accuracies)/tf.reduce_sum(mask)
20
21
22
   train_loss = tf.keras.metrics.Mean(name='train_loss')
23
   train_accuracy = tf.keras.metrics.Mean(name='train_accuracy')
24
25
  transformer = Transformer(
26
   num_layers=num_layers,
27
     d_model=d_model,
28
   num_heads=num_heads,
29
     dff=dff,
   input_vocab_size=tokenizers.pt.get_vocab_size().numpy(),
30
31
     target_vocab_size=tokenizers.en.get_vocab_size().numpy(),
   pe_input=1000,
32
33
     pe_target=1000,
   rate=dropout_rate)
34
35
36
   # Checkpoint
37
  checkpoint_path = "./checkpoints/train"
   ckpt = tf.train.Checkpoint(transformer=transformer, optimizer=
38
      optimizer)
39
   ckpt_manager = tf.train.CheckpointManager(ckpt, checkpoint_path,
      max_to_keep=5)
40
   # if a checkpoint exists, restore the latest checkpoint.
41
   if ckpt_manager.latest_checkpoint:
42
43
     ckpt.restore(ckpt_manager.latest_checkpoint)
44
   print('Latest checkpoint restored!!')
45
  EPOCHS = 20
46
47
   train_step_signature = [
     tf.TensorSpec(shape=(None, None), dtype=tf.int64),
48
49
     tf.TensorSpec(shape=(None, None), dtype=tf.int64),
50
   @tf.function(input_signature=train_step_signature)
51
52
  def train_step(inp, tar):
53
     tar_inp = tar[:, :-1]
   tar_real = tar[:, 1:]
54
55
     with tf.GradientTape() as tape:
56
       predictions, _ = transformer([inp, tar_inp], training = True)
57
       loss = loss_function(tar_real, predictions)
58
```

```
59
     gradients = tape.gradient(loss, transformer.trainable_variables)
     optimizer.apply_gradients(zip(gradients, transformer.
60
      trainable_variables))
61
62
     train loss(loss)
     train_accuracy(accuracy_function(tar_real, predictions))
63
64
65
   for epoch in range (EPOCHS):
66
     start = time.time()
67
     train loss.reset state()
68
     train_accuracy.reset_state()
     # inp -> portuguese, tar -> english
69
70
     for (batch, (inp, tar)) in enumerate(train_batches):
71
       train_step(inp, tar)
72
       if batch % 50 == 0:
73
         print(f'Epoch {epoch + 1} Batch {batch} Loss {train_loss.result
      ():.4f} Accuracy {train_accuracy.result():.4f}')
74
     if (epoch + 1) % 5 == 0:
75
       ckpt_save_path = ckpt_manager.save()
76
       print(f'Saving checkpoint for epoch {epoch+1} at {ckpt_save_path}
      ′ )
77
     print(f'Epoch {epoch + 1} Loss {train_loss.result():.4f} Accuracy {
78
      train_accuracy.result():.4f}')
79
     print(f'Time taken for 1 epoch: {time.time() - start:.2f} secs\n')
   class Translator(tf.Module):
     def __init__(self, tokenizers, transformer):
2
3
       self.tokenizers = tokenizers
4
       self.transformer = transformer
5
     def __call__(self, sentence, max_length=20):
6
       # input sentence is portuguese, hence adding the start and end
      token
8
       assert isinstance (sentence, tf.Tensor)
9
       if len(sentence.shape) == 0:
```

```
10
         sentence = sentence[tf.newaxis]
11
12
       sentence = self.tokenizers.pt.tokenize(sentence).to_tensor()
13
       encoder_input = sentence
14
15
       # as the target is english, the first token to the transformer
      should be the
16
       # english start token.
       start end = self.tokenizers.en.tokenize([''])[0]
17
18
       start = start_end[0][tf.newaxis]
19
       end = start_end[1][tf.newaxis]
20
21
       output_array = tf.TensorArray(dtype=tf.int64, size=0,
      dynamic_size=True)
22
       output_array = output_array.write(0, start)
```

```
23
       for i in tf.range(max_length):
24
25
         output = tf.transpose(output_array.stack())
26
         predictions, _ = self.transformer([encoder_input, output],
      training=False)
27
         predictions = predictions[:, -1:, :] # (batch_size, 1,
      vocab_size)
28
         predicted_id = tf.argmax(predictions, axis=-1)
29
         output_array = output_array.write(i+1, predicted_id[0])
         if predicted_id == end:
30
31
           break
32
33
       output = tf.transpose(output_array.stack())
34
       # output.shape (1, tokens)
35
       text = tokenizers.en.detokenize(output)[0]
36
       tokens = tokenizers.en.lookup(output)[0]
       _, attention_weights = self.transformer([encoder_input, output
37
      [:,:-1]], training=False)
38
39
       return text, tokens, attention_weights
40
   translator = Translator(tokenizers, transformer)
41
   sentence = "Eu li sobre triceratops na enciclopédia."
42
   translated_text, translated_tokens, attention_weights = translator(tf
      .constant(sentence))
44
45
  print(f'{"Prediction":15s}: {translated_text}')
```

Prediction: b'i read about trieps in the encycloody.'

#### APÊNDICE 9 – BIG DATA

### A - ENUNCIADO

Enviar um arquivo PDF contendo uma descrição breve (2 páginas) sobre a implementação de uma aplicação ou estudo de caso envolvendo Big Data e suas ferramentas (NoSQL e NewSQL). Caracterize os dados e Vs envolvidos, além da modelagem necessária dependendo dos modelos de dados empregados.

# B - RESOLUÇÃO

### Estudo de Caso: A Aplicação de Big Data no Airbnb

## Introdução

O Airbnb é uma plataforma online que conecta pessoas que desejam alugar suas casas com viajantes que buscam acomodação. Desde a sua fundação em 2008, a empresa cresceu exponencialmente, operando em mais de 190 países. Esse crescimento é sustentado pela capacidade do Airbnb de coletar, processar e analisar grandes volumes de dados, permitindo uma experiência personalizada para usuários e anfitriões. Neste estudo de caso, analisaremos como os dados do Airbnb podem ser classificados como Big Data, explorando os "Vs"do Big Data presentes em suas operações, e detalharemos as ferramentas utilizadas pela empresa no gerenciamento desses dados.

### Classificação dos Dados do Airbnb como Big Data

Os dados gerados e utilizados pelo Airbnb apresentam características que os classificam como Big Data. Os cinco "Vs"do Big Data, Volume, Velocidade, Variedade, Veracidade e Valor, estão presentes de forma evidente nos dados da empresa.

- 1. Volume: O Airbnb lida com uma quantidade massiva de dados. São milhões de usuários, reservas, avaliações, mensagens e interações diariamente. Esse grande volume requer sistemas capazes de armazenar e processar dados em escala petabyte.
- 2. Velocidade: Os dados são gerados e precisam ser processados em tempo real ou quase real. Reservas, cancelamentos e atualizações de disponibilidade ocorrem constantemente, exigindo processamento rápido para manter a plataforma atualizada.
- 3. Variedade: Os dados do Airbnb são diversos, incluindo dados estruturados (informações de reservas, perfis de usuários), semiestruturados (logs de atividades) e não estruturados (comentários, imagens). Essa variedade demanda sistemas flexíveis para gerenciamento.
- 4. Veracidade: A qualidade e a confiabilidade dos dados são cruciais. O Airbnb precisa garantir que as informações sejam precisas para manter a confiança entre hóspedes e anfitriões.
- 5. Valor: A análise desses dados gera insights valiosos que permitem melhorar a experiência do usuário, otimizar preços e expandir os negócios. O valor extraído dos dados é fundamental para a estratégia da empresa.

#### Ferramentas Utilizadas

O Airbnb utiliza uma arquitetura sofisticada de Big Data, incorporando diversas ferramentas e tecnologias para atender às suas necessidades. Abaixo, detalhamos cada ferramenta e seu papel no processo de dados.

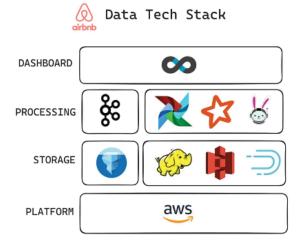
1. Amazon Web Services (AWS): S3 (Simple Storage Service): Utilizado para armazenamento de objetos, o S3 permite que o Airbnb armazene grandes volumes de dados de forma escalável e segura. E EC2 (Elastic Compute Cloud): Fornece capacidade computacional elástica, permitindo que a empresa escale seus recursos de processamento conforme a demanda.

- 2. Hadoop: O Hadoop é um framework de software que facilita o processamento distribuído de grandes conjuntos de dados. No Airbnb, é usado para armazenar e processar dados em clusters de servidores, permitindo a análise de grandes volumes de informação.
- HBase: Um banco de dados NoSQL que funciona em cima do Hadoop, o HBase é usado para armazenamento em tempo real de dados grandes e dispersos. Ele permite acessos rápidos a dados específicos dentro de conjuntos massivos.
- 4. Hive: O Hive permite consultas tipo SQL em dados armazenados no Hadoop. Isso facilita para os analistas de dados executarem consultas complexas sem a necessidade de programação em baixo nível.
- Presto: Presto é um motor de consultas distribuídas que permite ao Airbnb executar consultas SQL interativas e de baixa latência em grandes volumes de dados, espalhados por várias fontes de dados.
- 6. Apache Spark: O Spark é utilizado para processamento em lote e em tempo real. Com sua capacidade de processamento na memória, o Spark acelera o processamento de tarefas como análise de dados, aprendizado de máquina e processamento de gráficos.
- 7. Apache Airflow: Desenvolvido pelo próprio Airbnb, o Airflow é uma plataforma de orquestração de fluxos de trabalho que permite a automação de pipelines de dados. Ele gerencia a execução de tarefas, dependências e agendamento de processos.
- 8. Apache Superset: Outra criação do Airbnb, o Superset é uma plataforma de visualização de dados. Permite a criação de dashboards interativos e gráficos, auxiliando na interpretação e análise dos dados.
- 9. Druid: Druid é um sistema de banco de dados analítico projetado para consultas rápidas em grandes conjuntos de dados. No Airbnb, é usado para análise em tempo real e interativa, especialmente em dados de séries temporais.

### Aplicação das Ferramentas

O processo de dados no Airbnb pode ser dividido em várias etapas: coleta, armazenamento, processamento, análise e visualização. A Figura abaixo mostra uma visão geral das ferramentas utilizadas em cada etapa, sendo estas especificadas nas seções seguintes.

Figura 9.1 – VISÃO GERAL DAS FERRAMENTAS UTILIZADAS PELO AIRBNB



- 1. Coleta: Dados são coletados de diversas fontes, incluindo interações de usuários, informações de reservas, avaliações, logs de atividades e dados de terceiros (como eventos locais e clima).
- 2. Armazenamento: S3 e Hadoop são utilizados para armazenar os dados coletados. O S3 armazena objetos de dados, enquanto o Hadoop permite armazenamento distribuído e escalável.
- 3. Processamento: Apache Spark processa dados em lote e em tempo real, permitindo análises complexas e treinamento de modelos de machine learning. Hive e Presto são usados para consultas e manipulação de dados, facilitando o acesso a informações específicas. HBase oferece acesso rápido a dados em tempo real, essencial para aplicações que requerem baixa latência.
- 4. Análise: Druid permite análises rápidas e interativas, especialmente útil para monitoramento em tempo real e análise de séries temporais. Modelos de machine learning são implementados para diversas finalidades, como previsão de demanda, definição de preços dinâmicos e recomendações personalizadas.
- Orquestração e Automação: Apache Airflow gerencia e automatiza os fluxos de trabalho de dados, garantindo que os processos sejam executados corretamente e no tempo certo.
- 6. Visualização: Apache Superset é utilizado para criar dashboards e relatórios que auxiliam na tomada de decisões baseadas em dados.

#### Conclusão

O Airbnb é um exemplo notável de como Big Data pode ser utilizado para impulsionar negócios digitais. Através da coleta e análise de grandes volumes de dados, a empresa consegue oferecer experiências personalizadas, otimizar operações e inovar continuamente. Os cinco "Vs"do Big Data estão profundamente integrados nas operações do Airbnb, e o uso eficiente de ferramentas como AWS, Hadoop, Spark e outras permite que a empresa transforme dados brutos em insights valiosos. Como estudante de Big Data, observar o caso do Airbnb destaca a importância de uma arquitetura de dados robusta e da integração de diversas tecnologias para lidar com os desafios e oportunidades que o Big Data apresenta.

#### Referências

https://www.linkedin.com/pulse/what-technology-stack-airbnb-john-doe?utm\_source=share utm\_medium=member\_android&utm\_campaign=share\_via

https://www.junaideffendi.com/p/airbnb-data-tech-stack

https://vutr.substack.com/p/groupby-40-data-infrastructure-at

https://medium.com/airbnb-engineering/data-infrastructure-at-airbnb-8adfb34f169c

### APÊNDICE 10 - VISÃO COMPUTACIONAL

#### A - ENUNCIADO

### Extração de Características

Os bancos de imagens fornecidos são conjuntos de imagens de 250x250 pixels de imunohistoquímica (biópsia) de câncer de mama. No total são 4 classes (0, 1+, 2+ e 3+) que estão divididas em diretórios. O objetivo é classificar as imagens nas categorias correspondentes. Uma base de imagens será utilizada para o treinamento e outra para o teste do treino.

As imagens fornecidas são recortes de uma imagem maior do tipo WSI (Whole Slide Imaging) disponibilizada pela Universidade de Warwick (link). A nomenclatura das imagens segue o padrão XX\_HER\_YYYY.png, onde XX é o número do paciente e YYYY é o número da imagem recortada. Separe a base de treino em 80% para treino e 20% para validação. **Separe por pacientes (XX), não utilize a separação randômica! Pois, imagens do mesmo paciente não podem estar na base de treino e de validação, pois isso pode gerar um viés.** No caso da CNN VGG16 remova a última camada de classificação e armazene os valores da penúltima camada como um vetor de características. Após o treinamento, os modelos treinados devem ser validados na base de teste.

Tarefas:

- a) Carregue a base de dados de Treino.
- b) Crie partições contendo 80% para treino e 20% para validação (atenção aos pacientes).
- c) Extraia características utilizando LBP e a CNN VGG16 (gerando um csv para cada extrator).
- d) Treine modelos Random Forest, SVM e RNA para predição dos dados extraídos.
- e) Carregue a base de Teste e execute a tarefa 3 nesta base.
- f) Aplique os modelos treinados nos dados de treino
- g) Calcule as métricas de Sensibilidade, Especificidade e F1-Score com base em suas matrizes de confusão.
- h) Indique qual modelo dá o melhor o resultado e a métrica utilizada

#### **Redes Neurais**

Utilize as duas bases do exercício anterior para treinar as Redes Neurais Convolucionais VGG16 e a Resnet50. Utilize os pesos pré-treinados (Transfer Learning), refaça as camadas Fully Connected para o problema de 4 classes. Compare os treinos de 15 épocas com e sem Data Augmentation. Tanto a VGG16 quanto a Resnet50 têm como camada de entrada uma imagem 224x224x3, ou seja, uma imagem de 224x224 pixels coloridos (3 canais de cores). Portanto, será necessário fazer uma transformação de 250x250x3 para 224x224x3. Ao fazer o Data Augmentation cuidado para não alterar demais as cores das imagens e atrapalhar na classificação.

Tarefas:

- a) Utilize a base de dados de Treino já separadas em treino e validação do exercício anterior
- b) Treine modelos VGG16 e Resnet50 adaptadas com e sem Data Augmentation
- c) Aplique os modelos treinados nas imagens da base de Teste
- d) Calcule as métricas de Sensibilidade, Especificidade e F1-Score com base em suas matrizes de confusão.

e) Indique qual modelo dá o melhor o resultado e a métrica utilizada

# B - RESOLUÇÃO

```
2
   import os
  import cv2
3
4
  import shutil
5
  import numpy as np
  import pandas as pd
   import matplotlib.pyplot as plt
7
8
9
   from tqdm import tqdm
10 from pathlib import Path
11
  from sklearn.svm import SVC
12
13
  from sklearn.metrics import recall_score
  from sklearn.preprocessing import LabelEncoder
14
  from sklearn.ensemble import RandomForestClassifier
15
   from sklearn.model_selection import train_test_split
17
   from sklearn.metrics import classification_report, confusion_matrix
18
19
   from tensorflow.keras.layers import Dense
20
  from tensorflow.keras.models import Model
21
   from tensorflow.keras.models import Sequential
   from tensorflow.keras.applications import VGG16
23
   from tensorflow.keras.preprocessing import image
24
   from tensorflow.keras.applications.vgg16 import preprocess_input
25
26 | !unzip Train_Warwick.zip
27
28
  # Diretório das amostras de treino
29
   images_dir = Path('/content/Train_4cls_amostra')
30
   # Diretórios de destino para treino, validação e teste
31
32
  train_dir = images_dir.parent / 'train'
  val_dir = images_dir.parent / 'val'
33
34
35
  # Classes disponíveis
   classes = ['0', '1', '2', '3']
37
38
  patient_dict = {}
39
40
   # Percorrer cada classe e organizar pacientes
   for cls in classes:
41
42
       class_dir = images_dir / cls
43
       for img_name in os.listdir(class_dir):
           if img_name.endswith('.png'):
44
45
               # Extrair o ID do paciente do nome do arquivo (
      XX_HER_YYYY.png)
               patient_id = img_name.split('_')[0]
46
```

```
47
               img_path = class_dir / img_name
48
               if cls not in patient_dict:
49
                   patient_dict[cls] = {}
50
               if patient_id not in patient_dict[cls]:
51
                    patient dict[cls][patient id] = []
52
               patient_dict[cls][patient_id].append((img_path, cls))
53
54
   # Função para dividir pacientes em 80% treino e 20% validação
55
   def split_patients(patients):
56
       patient_ids = list(patients.keys())
57
       train_size = int(0.8 * len(patient_ids))
       train_patient_ids = patient_ids[:train_size]
58
59
       val_patient_ids = patient_ids[train_size:]
60
       return train_patient_ids, val_patient_ids
61
   # Função para mover as imagens para os diretórios de treino e validaç
   def move_images(patient_ids, patient_dict, target_dir):
63
64
       for pid in patient_ids:
65
           for img_path, cls in patient_dict[pid]:
               # Diretório da classe no conjunto alvo
66
67
               class dir = target dir / cls
               class_dir.mkdir(parents=True, exist_ok=True)
68
69
                # Caminho de destino da imagem
70
               dest_path = class_dir / img_path.name
71
                # Mover a imagem
72
               shutil.move(str(img_path), str(dest_path))
73
74
   # Para cada classe, dividir pacientes e mover as imagens
75
   for cls in classes:
76
       print(f"Processando classe {cls}...")
77
78
       # Obter os pacientes dessa classe
79
       patients_in_class = patient_dict[cls]
80
81
       # Dividir em treino e validação
82
       train_patient_ids, val_patient_ids = split_patients(
      patients_in_class)
83
      print('Número total de pacientes:', len(train_patient_ids) + len(
84
      val_patient_ids))
85
       print('Número de pacientes para treino:', len(train_patient_ids))
86
       print('Número de pacientes para validação:', len(val_patient_ids)
87
       # Mover as imagens de treino
88
89
       move_images(train_patient_ids, patients_in_class, train_dir)
90
91
       # Mover as imagens de validação
92
       move_images(val_patient_ids, patients_in_class, val_dir)
93
```

```
94 print ("Divisão de treino e validação concluída.")
95
96
   # Processando classe 0...
97
   # Número total de pacientes: 5
    # Número de pacientes para treino: 4
98
99
   # Número de pacientes para validação: 1
100 # Processando classe 1...
101
   # Número total de pacientes: 5
102 # Número de pacientes para treino: 4
103
    # Número de pacientes para validação: 1
104
   # Processando classe 2...
   # Número total de pacientes: 5
105
106 # Número de pacientes para treino: 4
107
   # Número de pacientes para validação: 1
   # Processando classe 3...
108
109
   # Número total de pacientes: 5
110 # Número de pacientes para treino: 4
   # Número de pacientes para validação: 1
111
112
   # Divisão de treino e validação concluída.
113
114
   def get_pixel(img, center, x, y):
115
116
            new_value = 0
117
118
            try:
119
                     # If local neighbourhood pixel
120
                     # value is greater than or equal
121
                     # to center pixel values then
                     # set it to 1
122
123
                     if img[x][y] >= center:
124
                             new_value = 1
125
126
            except:
127
                     # Exception is required when
128
                     # neighbourhood value of a center
129
                     # pixel value is null i.e. values
130
                     # present at boundaries.
131
                    pass
132
133
            return new_value
134
135
    # Function for calculating LBP
136
137
    def lbp_calculated_pixel(img, x, y, radius):
138
139
            center = imq[x][y]
140
141
            val_ar = []
142
143
            # top_left
144
            val_ar.append(get_pixel(img, center, x-radius, y-radius))
```

```
145
146
            # top
147
            val_ar.append(get_pixel(img, center, x-radius, y))
148
149
            # top right
150
            val_ar.append(get_pixel(img, center, x-radius, y + radius))
151
152
            # right
153
            val_ar.append(get_pixel(img, center, x, y + radius))
154
155
            # bottom_right
            val_ar.append(get_pixel(img, center, x + radius, y + radius))
156
157
158
            # bottom
159
            val_ar.append(get_pixel(img, center, x + radius, y))
160
161
            # bottom_left
162
            val_ar.append(get_pixel(img, center, x + radius, y-radius))
163
164
165
            val_ar.append(get_pixel(img, center, x, y-radius))
166
167
            # Now, we need to convert binary
168
            # values to decimal
169
            power_val = [1, 2, 4, 8, 16, 32, 64, 128]
170
171
            val = 0
172
173
            for i in range(len(val_ar)):
174
                     val += val_ar[i] * power_val[i]
175
176
            return val
177
178
    # Função para calcular o LBP
179
    def extract_lbp_image(image):
180
        img_gray = cv2.cvtColor(image,
181
                     cv2.COLOR_BGR2GRAY)
182
        img_gray = cv2.bilateralFilter(img_gray, 7, 75, 75)
183
184
        height, width, _ = image.shape
185
186
        img_lbp = np.zeros((height, width),
187
                np.uint8)
188
189
        for i in range(0, height):
190
          for j in range(0, width):
191
            img_lbp[i, j] = lbp_calculated_pixel(img_gray, i, j, 1)
192
193
        # Calcular o histograma
194
        h = cv2.calcHist([img_lbp], [0], None, [256], [0, 256])
195
        hist, bins = np.histogram(h.flatten(), 256, [0, 256])
```

```
196
        h_normalized = 100*((hist - hist.min()) /(hist.max() - hist.min())
       ) )
197
198
        return img_lbp, h_normalized
199
    def extract_lbp_from_directory(directory, classes):
200
201
        features = []
202
        labels = []
203
        print(f"Extraindo características LBP das imagens no diretório: {
204
       directory } ")
205
206
        # Iterar sobre cada classe
207
        for cls in classes:
            class_dir = os.path.join(directory, cls)
208
209
            image_names = os.listdir(class_dir)
210
            for img_name in tqdm(image_names, desc=f"Classe {cls}"):
                img_path = os.path.join(class_dir, img_name)
211
212
                # Ler a imagem
                image = cv2.imread(img_path, 1)
213
214
                 # Verificar se a imagem tem 4 canais
                if len(image.shape) == 3 and image.shape[2] == 4:
215
216
                     # Remover o canal alfa
217
                     image = image[:, :, :3]
218
                 # Extrair características LBP
219
                img_lbp, h_normalized = extract_lbp_image(image)
220
                features.append(h_normalized)
221
                 # Armazenar o label
222
                labels.append(cls)
223
224
        return features, labels
225
226
    # Extrair características do conjunto de treino
227
    train_features, train_labels = extract_lbp_from_directory(train_dir,
       classes)
228
229
    # Extrair características do conjunto de validação
    val_features, val_labels = extract_lbp_from_directory(val_dir,
230
       classes)
231
232
   # Extraindo características LBP das imagens no diretório: /content/
       val
233
234
    # Converter as listas em arrays numpy
235
   train_features = np.array(train_features)
    train_labels = np.array(train_labels)
236
237
    val_features = np.array(val_features)
238
   val_labels = np.array(val_labels)
239
240
   # Criar DataFrames
241 | train_df = pd.DataFrame(train_features)
```

```
242
    train_df['label'] = train_labels
243
244
   val_df = pd.DataFrame(val_features)
245
   val_df['label'] = val_labels
246
247
   # Salvar em CSV
   train_df.to_csv('train_lbp_features.csv', index=False)
248
249
   val_df.to_csv('val_lbp_features.csv', index=False)
250
251
   print("Características LBP salvas nos arquivos 'train_lbp_features.
       csv' e 'val_lbp_features.csv'")
252
253
   # Características LBP salvas nos arquivos 'train_lbp_features.csv' e
       'val_lbp_features.csv'
254
255
    # Carregar o modelo VGG16 pré-treinado, sem a camada de classificação
        final
256
    # pooling='avg' para obter um vetor de características de saída
257
   base_model = VGG16(weights='imagenet', include_top=False, pooling='
       avq')
258
259
    # Função para carregar uma imagem, redimensioná-la e pré-processá-la
    def process_image(img_path):
260
261
        img = image.load_img(img_path, target_size=(224, 224))
262
        img_array = image.img_to_array(img)
263
        img_array = np.expand_dims(img_array, axis=0)
264
        img_array = preprocess_input(img_array) # Pré-processamento
       VGG16
265
        return img_array
266
267
    # Função para extrair as características usando o modelo VGG16
268
    def extract_vgg16_features(img_path, model):
269
        img_array = process_image(img_path)
270
        features = model.predict(img_array)
271
        return features.flatten() # Converter para vetor de caracterí
       sticas
272
273
      # Extrair características de um diretório de imagens
    def extract_features_from_directory(directory, model):
274
275
        features = []
276
        labels = []
277
        for class_label in os.listdir(directory):
278
279
            class_dir = os.path.join(directory, class_label)
280
281
            if os.path.isdir(class_dir):
282
                for img_name in tqdm(os.listdir(class_dir), desc=f"Classe
        {class_label}"):
                    img_path = os.path.join(class_dir, img_name)
283
284
                     # Extrair características
285
```

```
286
                    img_features = extract_vgg16_features(img_path, model
       )
287
288
                    # Armazenar as características e o rótulo (classe)
289
                    features.append(img features)
290
                    labels.append(class_label)
291
292
        return np.array(features), np.array(labels)
293
294
    # Extração das características das imagens de treino usando o modelo
       VGG16
   print("Extraindo características do conjunto de treino usando o
295
       modelo VGG16")
296
   train_features, train_labels = extract_features_from_directory(
       train_dir, base_model)
297
298
   # Extração das características das imagens de validação usando o
      modelo VGG16
299
   print("Extraindo características do conjunto de validação usando o
       modelo VGG16")
   val_features, val_labels = extract_features_from_directory(val_dir,
300
       base model)
301
302
   # Converter para DataFrame e salvar como CSV
303
304
   train_df = pd.DataFrame(train_features)
305
   train_df['label'] = train_labels
   train_df.to_csv('train_vgg16_features.csv', index=False)
306
307
308
   val_df = pd.DataFrame(val_features)
   val_df['label'] = val_labels
309
310 | val_df.to_csv('val_vgg16_features.csv', index=False)
311
   print("Características VGG16 salvas nos arquivos '
312
       train_vgg16_features.csv' e 'val_vgg16_features.csv'")
313
   # Carregar características LBP e CNN VGG16
314
315
   train_lbp = pd.read_csv('train_lbp_features.csv')
316
   train_vgg16 = pd.read_csv('train_vgg16_features.csv')
317
318
   # Separar características e rótulos (labels)
319
   X_train_lbp = train_lbp.drop(columns=['label'])
320
   y_train_lbp = train_lbp['label']
321
322 | X_train_vgg16 = train_vgg16.drop(columns=['label'])
   y_train_vgg16 = train_vgg16['label']
323
324
325
   # Codificar as labels para valores numéricos
326
   label_encoder = LabelEncoder()
327
   |y_train_lbp = label_encoder.fit_transform(y_train_lbp)
328 | y_train_vgg16 = label_encoder.fit_transform(y_train_vgg16)
```

```
329
330
   # Treinando o modelo SVM com características LBP
331
    svm_lbp = SVC(kernel='linear', random_state=42)
332
   svm_lbp.fit(X_train_lbp, y_train_lbp)
333
334
    # Treinando o modelo SVM com características VGG16
335
    svm_vgg16 = SVC(kernel='linear', random_state=42)
336
    svm_vgg16.fit(X_train_vgg16, y_train_vgg16)
337
338
    # Função para criar o modelo RNA
339
    def create_rna_model(input_dim):
340
        model = Sequential()
        model.add(Dense(128, activation='relu', input_dim=input_dim))
341
342
        model.add(Dense(64, activation='relu'))
343
        model.add(Dense(32, activation='relu'))
344
       model.add(Dense(len(label_encoder.classes_), activation='softmax'
           # Saída com número de classes
345
       model.compile(optimizer='adam', loss='
       sparse_categorical_crossentropy', metrics=['accuracy'])
346
        return model
347
348
    # Treinando a RNA com características LBP
349
    rna_lbp = create_rna_model(X_train_lbp.shape[1])
350
    rna_lbp.fit(X_train_lbp, y_train_lbp, epochs=10, batch_size=32,
       validation_split=0.2)
351
352
   # Treinando a RNA com características VGG16
353
   rna_vgg16 = create_rna_model(X_train_vgg16.shape[1])
354
    rna_vgg16.fit(X_train_vgg16, y_train_vgg16, epochs=10, batch_size=32,
        validation_split=0.2)
355
356
   # Avaliando Random Forest (LBP)
357
   y_pred_rf_lbp = rf_lbp.predict(X_train_lbp)
   print("Relatório de Classificação (Random Forest - LBP):")
358
359
   print(classification_report(y_train_lbp, y_pred_rf_lbp))
360
361
   # Avaliando SVM (LBP)
   y_pred_svm_lbp = svm_lbp.predict(X_train_lbp)
362
363
   print("Relatório de Classificação (SVM - LBP):")
364
   print(classification_report(y_train_lbp, y_pred_svm_lbp))
365
366
   # Avaliando RNA (LBP)
367
   y_pred_rna_lbp = rna_lbp.predict(X_train_lbp)
   y_pred_rna_lbp = np.argmax(y_pred_rna_lbp, axis=1)
368
369
   print("Relatório de Classificação (RNA - LBP):")
   print(classification_report(y_train_lbp, y_pred_rna_lbp))
370
371
372 | # Avaliando Random Forest (VGG16)
373
   y_pred_rf_vgg16 = rf_vgg16.predict(X_train_vgg16)
374
   print("Relatório de Classificação (Random Forest - VGG16):")
375 | print (classification_report (y_train_vgg16, y_pred_rf_vgg16))
```

```
376
377
    # Avaliando SVM (VGG16)
378
    y_pred_svm_vgg16 = svm_vgg16.predict(X_train_vgg16)
379
    print("Relatório de Classificação (SVM - VGG16):")
    print(classification_report(y_train_vgg16, y_pred_svm_vgg16))
380
381
    # Avaliando RNA (VGG16)
382
383
    y_pred_rna_vgg16 = rna_vgg16.predict(X_train_vgg16)
    y_pred_rna_vgg16 = np.argmax(y_pred_rna_vgg16, axis=1)
384
385
    print("Relatório de Classificação (RNA - VGG16):")
386
    print(classification_report(y_train_vgg16, y_pred_rna_vgg16))
387
388
    #Relatório de Classificação (Random Forest - LBP):
389
                    precision
                                   recall
                                            f1-score
                                                        support
390
391
                          1.00
                                     1.00
                                                1.00
                                                            116
                                                1.00
                                                            117
392
                 2
                                                            120
393
                          1.00
                                     1.00
                                                1.00
                                                            120
394
                                                1.00
395
396
                                                1.00
                                                            473
         accuracy
                                                1.00
                                                            473
397
        macro avq
                          1.00
                                     1.00
                                     1.00
                                                1.00
                                                            473
398
                          1.00
    #weighted avg
399
400
    #Relatório de Classificação (SVM - LBP):
401
                    precision
                                   recall
                                            f1-score
                                                        support
402
403
                 0
                          1.00
                                     1.00
                                                1.00
                                                            116
                 1
404
                                                1.00
                 2
                                                            120
405
                          1.00
                                     1.00
                                                1.00
                 3
406
                                                1.00
                                                            120
407
                                                            473
408
                                                1.00
         accuracy
409
                          1.00
                                     1.00
                                                1.00
                                                            473
        macro avg
    #weighted avg
410
                                                1.00
                                                            473
411
412
    #Relatório de Classificação (RNA -
413
                    precision
                                   recall
                                            f1-score
                                                        support
414
415
                 0
                          0.76
                                     0.73
                                                0.75
                                                            116
416
                 1
                          0.71
                                     0.89
                                                0.79
                                                            117
                 2
                                                            120
417
                          0.43
                                     0.74
                                                0.54
                          1.00
                                                0.12
                                                             120
418
419
420
         accuracy
                                                0.60
                                                            473
                          0.73
421
                                     0.61
                                                0.55
                                                            473
        macro avg
422
                                      0.60
                                                             473
    # weighted avg
423
424
    #Relatório de Classificação (Random Forest - VGG16):
425
                    precision
                                   recall
                                           f1-score
                                                        support
426
```

```
1.00
427
                                     1.00
                 0
                                                1.00
                                                            116
428
                                     1.00
                                                1.00
                                                            117
429
                 2
                          1.00
                                     1.00
                                                1.00
                                                            120
                 3
430
                                                1.00
431
432
                                                1.00
                                                            473
        accuracy
433
        macro avg
                          1.00
                                     1.00
                                                1.00
                                                            473
434
    #weighted avg
                          1.00
                                     1.00
                                                1.00
                                                            473
435
    # Relatório de Classificação (SVM - VGG16):
436
437
                   precision
                                  recall
                                          f1-score
                                                       support
438
439
                 0
                          1.00
                                     1.00
                                                1.00
                                                            116
440
                          1.00
                                     1.00
                                                1.00
                                                            117
441
                 2
                          1.00
                                     1.00
                                                1.00
                                                            120
442
                 3
                                                1.00
443
                                                1.00
                                                            473
444
         accuracy
445
        macro avg
                          1.00
                                     1.00
                                                1.00
                                                            473
    #weighted avg
446
                                                1.00
                                                            473
447
    #Relatório de Classificação (RNA - VGG16):
448
449
                    precision
                                   recall
                                            f1-score
                                                        support
450
451
                          0.98
                                     1.00
                                                0.99
                                                            116
452
                 1
                          0.92
                                     1.00
                                                0.96
                                                            117
                 2
                                     0.85
453
                          1.00
                                                0.92
                                                            120
                          0.95
                                     1.00
                                                0.98
                                                            120
454
455
                                                            473
456
        accuracy
457
        macro avg
                          0.96
                                     0.96
                                                0.96
                                                            473
458
    #weighted avg
                                                            473
459
460
    # !unzip Test_Warwick.zip
461
    # Diretório das amostras de teste
462
    test_dir = Path('/content/Test_4cl_amostra')
463
464
    # Extrair características do conjunto de teste
465
    test_features, test_labels = extract_lbp_from_directory(test_dir,
466
       classes)
467
    # Converter as listas em arrays numpy
468
469
    test_features = np.array(test_features)
470
    test_labels = np.array(test_labels)
471
472
    # Criar DataFrame
473
    test_df = pd.DataFrame(test_features)
474
    test_df['label'] = test_labels
475
476 | # Salvar em CSV
```

```
477 | test_df.to_csv('test_lbp_features.csv', index=False)
478
479
   print("Características LBP salvas no arquivo 'test_lbp_features.csv'"
       )
480
481
    # Extração das características das imagens de teste
   print("Extraindo características do conjunto de teste...")
482
   test_features, test_labels = extract_features_from_directory(test_dir
483
       , base_model)
484
485
    # Converter para DataFrame e salvar como CSV
   test_df = pd.DataFrame(test_features)
486
487
   test_df['label'] = test_labels
488
   test_df.to_csv('test_vgg16_features.csv', index=False)
489
490
   print("Características VGG16 salvas nos arquivos 'test_vgg16_features
       .csv'")
491
492
    # Carregar os dados de teste
493
    test_lbp = pd.read_csv('test_lbp_features.csv')
494
   test_vgg16 = pd.read_csv('test_vgg16_features.csv')
495
496
   # Separar as características e os rótulos (labels)
497
   X_test_lbp = test_lbp.drop(columns=['label'])
498
   y_test_lbp = test_lbp['label']
499
500 | X_test_vgg16 = test_vgg16.drop(columns=['label'])
501
   y_test_vgg16 = test_vgg16['label']
502
503
   # Codificar os rótulos de teste (mesma codificação usada para os
       dados de treino)
504
   label_encoder = LabelEncoder()
505
   y_test_lbp = label_encoder.fit_transform(y_test_lbp)
506
507
   label_encoder = LabelEncoder()
   y_test_vgg16 = label_encoder.fit_transform(y_test_vgg16)
508
509
510
   # Fazer previsões com o modelo Random Forest
511
   y_pred_rf_lbp = rf_lbp.predict(X_test_lbp)
512
513
   y_pred_rf_vgg16 = rf_vgg16.predict(X_test_vgg16)
514
515
   # Fazer previsões com o modelo SVM
516
   y_pred_svm_lbp = svm_lbp.predict(X_test_lbp)
517
   y_pred_svm_vgg16 = svm_vgg16.predict(X_test_vgg16)
518
519
520
   # Fazer previsões com a RNA
521
   y_pred_rna_lbp = rna_lbp.predict(X_test_lbp)
522
   y_pred_rna_lbp = np.argmax(y_pred_rna_lbp, axis=1)
523
```

```
524
    y_pred_rna_vgg16 = rna_vgg16.predict(X_test_vgg16)
525
    y_pred_rna_vgg16 = np.argmax(y_pred_rna_vgg16, axis=1)
526
527
    # Calcular Especificidade a partir da Matriz de Confusão
528
    def calculate_specificity(conf_matrix):
529
        TN = conf_matrix[0, 0]
        FP = conf_matrix[0, 1]
530
531
        specificity = TN / (TN + FP)
532
        return specificity
533
534
    print("Random Forest\n")
535
536
    # Avaliar o desempenho do Random Forest (LBP)
537
    print("LBP")
538
    print(classification_report(y_test_lbp, y_pred_rf_lbp))
539
    conf_matrix_rf_lbp = confusion_matrix(y_test_lbp, y_pred_rf_lbp)
    specificity_rf_lbp = calculate_specificity(conf_matrix_rf_lbp)
540
541
    print(f"Especificidade: {specificity_rf_lbp}")
542
543
    # Avaliar o desempenho do Random Forest (VGG16)
544
   print("\n\nVGG16")
545
    print(classification_report(y_test_vgg16, y_pred_rf_vgg16))
    conf_matrix_rf_vgg16 = confusion_matrix(y_test_vgg16, y_pred_rf_vgg16
546
547
    specificity_rf_vgg16 = calculate_specificity(conf_matrix_rf_vgg16)
548
    print(f"Especificidade: {specificity_rf_vgg16}")
549
550
551
    #Random Forest
552
553
    #LBP
554
                    precision
                                 recall f1-score
                                                    support
555
556
                                    0.50
                                              0.43
                                                          101
557
                 1
                         0.29
                                    0.27
                                               0.28
                                                           90
558
                 2
                         0.43
                                    0.29
559
                 3
                         0.77
                                    0.82
                                               0.80
                                                           90
560
561
                                               0.47
                                                          371
         accuracy
562
       macro avg
                         0.47
                                    0.47
                                              0.46
                                                          371
563
    #weighted avg
                         0.47
                                    0.47
                                               0.46
                                                          371
564
    #Especificidade: 0.6493506493506493
565
566
567
    #VGG16
568
569
                                  recall
                                          f1-score
                    precision
                                                      support
570
571
                         0.80
                                    0.99
                                              0.88
                                                          101
                 1
572
                                              0.26
573
                         0.59
                                    0.83
                                               0.69
```

```
0.98
574
                          0.95
                                                0.96
575
576
         accuracy
577
                          0.73
                                     0.74
                                                0.70
                                                            371
        macro avg
                                                0.70
                                                            371
                          0.73
                                     0.75
578
    #weighted avg
579
    #Especificidade: 0.9900990099009901
580
581
582
    print("SVM\n")
583
584
    # Avaliar o desempenho do SVM (LBP)
585
    print("LBP")
586
    print(classification_report(y_test_lbp, y_pred_svm_lbp))
587
    conf_matrix_svm_lbp = confusion_matrix(y_test_lbp, y_pred_svm_lbp)
    specificity_svm_lbp = calculate_specificity(conf_matrix_svm_lbp)
588
589
    print(f"Especificidade: {specificity_svm_lbp}")
590
591
    # Avaliar o desempenho do SVM (LBP)
592
    print("\n\nVGG16")
593
    print(classification_report(y_test_vgg16, y_pred_svm_vgg16))
594
    conf_matrix_svm_vgg16 = confusion_matrix(y_test_vgg16,
       y_pred_svm_vgg16)
595
    specificity_svm_vgg16 = calculate_specificity(conf_matrix_svm_vgg16)
596
    print(f"Especificidade: {specificity_svm_vgg16}")
597
598
    #SVM
599
    #LBP
600
601
                    precision
                                   recall
                                            f1-score
                                                        support
602
603
                 0
                          0.38
                                     0.50
                                                0.43
                                                            101
604
                                                0.34
605
                 2
                          0.50
                                     0.43
                                                0.46
                                                             90
                          0.74
606
                                     0.63
                                                0.68
607
                                                            371
608
                                                0.47
         accuracy
                                     0.47
                                                0.48
                                                            371
609
        macro avg
                          0.49
                          0.49
                                     0.47
                                                0.48
                                                            371
610
    #weighted avg
611
    #Especificidade: 0.6538461538461539
612
613
614
    #VGG16
615
                    precision
                                   recall
                                            f1-score
                                                        support
616
617
                          0.85
                                     0.92
                                                0.88
                                                            101
                          0.77
                                                0.43
618
619
                 2
                          0.64
                                     0.96
                                                0.77
                                                             90
620
                                     1.00
                                                0.99
621
622
                                                            371
                                     0.79
                                                0.77
                                                            371
623
        macro avg
                          0.81
```

```
624
    #weighted avg 0.81 0.80
625
    #Especificidade: 0.9393939393939394
626
627
628
    print("RNA\n")
629
    # Avaliar o desempenho da RNA (LBP)
630
631
    print("LBP")
632
   print(classification_report(y_test_lbp, y_pred_rna_lbp, zero_division
633
    conf_matrix_rna_lbp = confusion_matrix(y_test_lbp, y_pred_rna_lbp)
    specificity_rna_lbp = calculate_specificity(conf_matrix_rna_lbp)
634
    print(f"Especificidade: {specificity_rna_lbp}")
635
636
637
    # Avaliar o desempenho da RNA (LBP)
638
    print("\n\nVGG16")
639
    print(classification_report(y_test_vgg16, y_pred_rna_vgg16,
       zero_division=1))
    conf_matrix_rna_vgg16 = confusion_matrix(y_test_vgg16,
640
       y_pred_rna_vgg16)
641
    specificity_rna_vgg16 = calculate_specificity(conf_matrix_rna_vgg16)
642
    print(f"Especificidade: {specificity_rna_vgg16}")
643
    #RNA
644
645
646
    #LBP
647
                    precision
                                  recall
                                           f1-score
                                                       support
648
649
                         0.41
                                    0.45
                                               0.43
                                                           101
                 1
650
                                    0.42
651
                 2
                         0.28
                                    0.46
                                               0.35
                                                            90
652
                                               0.02
653
                                                           371
654
                                               0.34
655
                         0.51
                                    0.33
                                               0.29
                                                           371
        macro avg
    #weighted avg
                                    0.34
656
657
658
    #Especificidade: 0.5357142857142857
659
660
    #VGG16
661
                                  recall
                    precision
                                           f1-score
                                                       support
662
                         0.92
                                    0.95
                                               0.94
                                                           101
663
664
                         0.81
                                    0.58
                                               0.68
665
                 2
                         0.69
                                    0.73
                                               0.71
                                                            90
                         0.84
                                    1.00
                                               0.91
666
667
                                                           371
                                               0.82
668
         accuracy
                         0.82
                                    0.82
                                                           371
669
        macro avg
                                               0.81
670
    #weighted avg
671
```

```
import tensorflow as tf
  from tensorflow.keras.preprocessing.image import ImageDataGenerator
  from tensorflow.keras.applications import VGG16, ResNet50
3
  from tensorflow.keras.layers import Dense, Flatten,
      GlobalAveragePooling2D
  from tensorflow.keras.models import Model
5
   from tensorflow.keras.optimizers import Adam
6
7
   import numpy as np
   import matplotlib.pyplot as plt
9
   from sklearn.metrics import confusion_matrix, classification_report
10 import seaborn as sns
11
   import os
12
  img_height, img_width = 224, 224
13
  batch_size = 32
15 | num_classes = 4  # Classes: 0, 1, 2, 3
16
17
  train_datagen = ImageDataGenerator(
       preprocessing_function=tf.keras.applications.vgg16.
18
      preprocess_input # Pré-processamento específico
19
20
   val_datagen = ImageDataGenerator(
21
22
       preprocessing_function=tf.keras.applications.vgg16.
      preprocess_input
23
   )
24
25
  train_generator = train_datagen.flow_from_directory(
26
       train_dir,
27
       target_size=(img_height, img_width),
28
       batch_size=batch_size,
29
       class_mode='categorical'
30
   )
31
32
  validation_generator = val_datagen.flow_from_directory(
33
       val dir,
34
       target_size=(img_height, img_width),
35
       batch_size=batch_size,
36
       class_mode='categorical'
37
38
39 | def create_vgg16_model():
```

```
40
       # Carregar o modelo base VGG16 sem as camadas totalmente
      conectadas (top)
       base_model_vgg = VGG16(weights='imagenet', include_top=False,
41
      input_shape=(img_height, img_width, 3))
42
43
       # Congelar as camadas convolucionais
       for layer in base_model_vgg.layers:
44
45
           layer.trainable = False
46
47
       # Adicionar novas camadas
       x = base\_model\_vgg.output
48
49
       x = Flatten()(x)
50
       x = Dense(1024, activation='relu')(x)
51
       predictions = Dense(num_classes, activation='softmax')(x)
52
53
       # Criar o modelo final
       model_vgg = Model(inputs=base_model_vgg.input, outputs=
54
      predictions)
55
56
       return model_vgg
57
58
   model_vgg = create_vgg16_model()
59
60
   model_vgg.compile(optimizer=Adam(learning_rate=1e-4),
                      loss='categorical crossentropy',
61
62
                      metrics=['accuracy'])
63
   epochs = 5
64
65
   history_vgg = model_vgg.fit(
66
67
       train_generator,
68
       steps_per_epoch=train_generator.samples // batch_size,
69
       epochs=epochs,
70
       validation_data=validation_generator,
71
       validation_steps=validation_generator.samples // batch_size
72
73
74
   def create resnet50 model():
75
       base_model_resnet = ResNet50(weights='imagenet', include_top=
      False, input_shape=(img_height, img_width, 3))
       for layer in base_model_resnet.layers:
76
77
           layer.trainable = False
78
       x = base_model_resnet.output
79
       x = GlobalAveragePooling2D()(x)
80
       x = Dense(1024, activation='relu')(x)
       predictions = Dense(num_classes, activation='softmax')(x)
81
82
       model_resnet = Model(inputs=base_model_resnet.input, outputs=
      predictions)
83
       return model_resnet
84
  |model_resnet = create_resnet50_model()
85
```

```
86
87
   model_resnet.compile(optimizer=Adam(learning_rate=1e-4),
88
                          loss='categorical_crossentropy',
89
                          metrics=['accuracy'])
90
91
    history_resnet = model_resnet.fit(
92
        train_generator,
93
        steps_per_epoch=train_generator.samples // batch_size,
94
        epochs=epochs,
95
        validation data=validation generator,
96
        validation_steps=validation_generator.samples // batch_size
97
98
99
    train_datagen_aug = ImageDataGenerator(
        preprocessing_function=tf.keras.applications.vgg16.
100
       preprocess_input,
101
        rotation_range=20,
102
        width_shift_range=0.1,
103
        height_shift_range=0.1,
104
        horizontal_flip=True
105
106
107
    train_generator_aug = train_datagen_aug.flow_from_directory(
108
        train_dir,
109
        target_size=(img_height, img_width),
110
        batch_size=batch_size,
111
        class_mode='categorical'
112
113
114
    # Treinamento com Data Augmentation
   model_vgg_aug = create_vgg16_model()
115
116
   model_vgg_aug.compile(optimizer=Adam(learning_rate=1e-4),
117
                       loss='categorical_crossentropy',
118
                       metrics=['accuracy'])
119
    history_vgg_aug = model_vgg_aug.fit(
120
121
        train_generator_aug,
122
        steps_per_epoch=train_generator_aug.samples // batch_size,
123
        epochs=epochs,
124
        validation_data=validation_generator,
125
        validation_steps=validation_generator.samples // batch_size
126
127
128
    # Treinamento com Data Augmentation
129
    model_resnet_aug = create_vgg16_model()
    model_resnet_aug.compile(optimizer=Adam(learning_rate=1e-4),
130
131
                          loss='categorical_crossentropy',
132
                          metrics=['accuracy'])
133
134
   history_resnet_aug = model_resnet_aug.fit (
135
        train_generator_aug,
```

```
136
        steps_per_epoch=train_generator_aug.samples // batch_size,
137
        epochs=epochs,
138
        validation_data=validation_generator,
139
        validation_steps=validation_generator.samples // batch_size
140
141
142
   test_datagen = ImageDataGenerator(
143
        preprocessing_function=tf.keras.applications.vgg16.
       preprocess_input
144
145
    test_generator = test_datagen.flow_from_directory(
146
147
        test_dir,
148
        target_size=(img_height, img_width),
149
        batch_size=1,
                       # Para avaliar imagem por imagem
150
        class_mode='categorical',
        shuffle=False # Importante para manter a ordem das imagens
151
152
153
154
   # VGG sem Data Augmentation
155
   test_generator.reset()
156
   preds_vgg = model_vgg.predict(test_generator, steps=test_generator.
       samples)
157
158
   # VGG com Data Augmentation
159
   test_generator.reset()
160
   preds_vgg_aug = model_vgg_aug.predict(test_generator, steps=
       test_generator.samples)
161
162
   # Converter as predições para classes
163
   y_pred_vgg = np.argmax(preds_vgg, axis=1)
164
   y_pred_vgg_aug = np.argmax(preds_vgg_aug, axis=1)
   y_pred_resnet = np.argmax(preds_resnet, axis=1)
165
   y_pred_resnet_aug = np.argmax(preds_resnet_aug, axis=1)
166
167
168 # Obter as classes verdadeiras
169
   y_true = test_generator.classes
170
   class_labels = list(test_generator.class_indices.keys())
171
172
   from sklearn.metrics import confusion_matrix
173
174
    # VGG16 sem Data Augmentation
175
   cm_vgg = confusion_matrix(y_true, y_pred_vgg)
176
177
    # VGG16 com Data Augmentation
178
   cm_vgg_aug = confusion_matrix(y_true, y_pred_vgg_aug)
179
   # ResNet50 sem Data Augmentation
180
   cm_resnet = confusion_matrix(y_true, y_pred_resnet)
181
182
183 | # ResNet50 com Data Augmentation
```

```
184
    cm_resnet_aug = confusion_matrix(y_true, y_pred_resnet_aug)
185
186
    from sklearn.metrics import classification_report, confusion_matrix
187
188
    def calculate_specificity(cm):
189
        # cm é a matriz de confusão
        FP = cm.sum(axis=0) - np.diag(cm)
190
191
        TN = cm.sum() - (FP + cm.sum(axis=1) - np.diag(cm) + np.diag(cm))
192
        specificity = TN / (TN + FP)
193
        return specificity
194
195
    # Avaliar o desempenho do modelo VGG16 sem Data Augmentation
196
   print("VGG16 sem Data Augmentation\n")
197
   print(classification_report(y_true, y_pred_vgg, target_names=
       class_labels))
198
   conf_matrix_vgg = confusion_matrix(y_true, y_pred_vgg)
199
    specificity_vqq = calculate_specificity(conf_matrix_vqq)
   print(f"Especificidade: {specificity_vgg}")
200
201
    # Avaliar o desempenho do modelo VGG16 com Data Augmentation
202
203
   print("\nVGG16 com Data Augmentation\n")
204
   print(classification_report(y_true, y_pred_vgg_aug, target_names=
       class_labels))
205
   conf_matrix_vgg_aug = confusion_matrix(y_true, y_pred_vgg_aug)
206
   specificity_vgg_aug = calculate_specificity(conf_matrix_vgg_aug)
207
   print(f"Especificidade: {specificity_vgg_aug}")
208
209
    # Avaliar o desempenho do modelo ResNet50 sem Data Augmentation
210
    print("\nResNet50 sem Data Augmentation\n")
211
   print(classification_report(y_true, y_pred_resnet, target_names=
       class_labels))
212
   conf_matrix_resnet = confusion_matrix(y_true, y_pred_resnet)
213
   specificity_resnet = calculate_specificity(conf_matrix_resnet)
214
    print(f"Especificidade: {specificity_resnet}")
215
216
   # Avaliar o desempenho do modelo ResNet50 com Data Augmentation
217
   print("\nResNet50 com Data Augmentation\n")
218
   print(classification_report(y_true, y_pred_resnet_aug, target_names=
       class labels))
219
   conf_matrix_resnet_aug = confusion_matrix(y_true, y_pred_resnet_aug)
    specificity_resnet_aug = calculate_specificity(conf_matrix_resnet_aug
220
221
   print (f"Especificidade: {specificity_resnet_aug}")
222
223
    #VGG16 sem Data Augmentation
224
225
                                 recall
                                         f1-score
                   precision
                                                    support
226
227
                0
                        0.77
                                   0.90
                                             0.83
                                                         101
                1
228
                                   0.18
                                             0.26
229
                2
                        0.63
                                   0.91
                                             0.75
```

220	П	2	0 07	0 00	0.00	0.0	
230	#	3	0.97	0.99	0.98	90	
231	ш				0 75	0.01	
232	# accur				0.75	371	
233	# macro	_	0.72	0.74	0.70	371	
234	#weighted	avg	0.72	0.75	0.71	371	
235							
236	#Especific	idade:	[0.9	0.94661	922 0.82918	3149 0.989323	84]
237							
238	#VGG16 com	Data A	ugmentatic	n			
239							
240	#	pr	ecision	recall	f1-score	support	
241							
242	#	0	0.81	0.99	0.89	101	
243	#	1	0.81	0.29	0.43	90	
244	#	2	0.65	0.84	0.73	90	
245	#	3	0.92	1.00	0.96	90	
246							
247	# accur	acy			0.79	371	
248	# macro	avg	0.80	0.78	0.75	371	
249	#weighted	avg	0.80	0.79	0.76	371	
250							
251	#Especific	idade:	[0.9111111	1 0.97864	769 0.85409	9253 0.9715302	25]
252							
253	#ResNet50 sem Data Augmentation						
254							
255	#	pr	ecision	recall	f1-score	support	
256							
257	#	0	0.99	0.97	0.98	101	
258	#	1	0.93	0.56	0.69	90	
259	#	2	0.69	0.97	0.81	90	
260	#	3	0.98	1.00	0.99	90	
261							
262	# accur	acy			0.88	371	
263	# macro	avg	0.90	0.87	0.87	371	
264	#weighted	avg	0.90	0.88	0.87	371	
265							
266	#Especific	idade:	[0.9962963	0.98576	5512 0.86120	0.992882	56]
267							
268	#ResNet50	com Dat	a Augmenta	ition			
269							
270	#	pr	ecision	recall	f1-score	support	
271							
272	#	0	0.75	0.99	0.85	101	
273	#	1	0.58	0.12	0.20	90	
274	#	2	0.61	0.81	0.70	90	
275	#	3	0.91	1.00	0.95	90	
276							
277	# accur				0.74	371	
278	# macro	_	0.71	0.73	0.68	371	
279	#weighted	avg	0.71	0.74	0.68	371	
280							

### APÊNDICE 11 - ASPECTOS FILOSÓFICOS E ÉTICOS DA IA

### A - ENUNCIADO

Título do Trabalho: "Estudo de Caso: Implicações Éticas do Uso do ChatGPT"

**Trabalho em Grupo:** O trabalho deverá ser realizado em grupo de alunos de no máximo seis (06) integrantes.

**Objetivo do Trabalho:** Investigar as implicações éticas do uso do ChatGPT em diferentes contextos e propor soluções responsáveis para lidar com esses dilemas.

## Parâmetros para elaboração do Trabalho:

- 1. Relevância Ética: O trabalho deve abordar questões éticas significativas relacionadas ao uso da inteligência artificial, especialmente no contexto do ChatGPT. Os alunos devem identificar dilemas éticos relevantes e explorar como esses dilemas afetam diferentes partes interessadas, como usuários, desenvolvedores e a sociedade em geral.
- **2. Análise Crítica:** Os alunos devem realizar uma análise crítica das implicações éticas do uso do ChatGPT em estudos de caso específicos. Eles devem examinar como o algoritmo pode influenciar a disseminação de informações, a privacidade dos usuários e a tomada de decisões éticas. Além disso, devem considerar possíveis vieses algorítmicos, discriminação e questões de responsabilidade.
- **3. Soluções Responsáveis:** Além de identificar os desafios éticos, os alunos devem propor soluções responsáveis e éticas para lidar com esses dilemas. Isso pode incluir sugestões para políticas, regulamentações ou práticas de design que promovam o uso responsável da inteligência artificial. Eles devem considerar como essas soluções podem equilibrar os interesses de diferentes partes interessadas e promover valores éticos fundamentais, como transparência, justiça e privacidade.
- **4.** Colaboração e Discussão: O trabalho deve envolver discussões em grupo e colaboração entre os alunos. Eles devem compartilhar ideias, debater diferentes pontos de vista e chegar a conclusões informadas através do diálogo e da reflexão mútua. O estudo de caso do ChatGPT pode servir como um ponto de partida para essas discussões, incentivando os alunos a aplicar conceitos éticos e legais aprendidos ao analisar um caso concreto.
- **5. Limite de Palavras:** O trabalho terá um limite de 6 a 10 páginas teria aproximadamente entre 1500 e 3000 palavras.
- **6. Estruturação Adequada:** O trabalho siga uma estrutura adequada, incluindo introdução, desenvolvimento e conclusão. Cada seção deve ocupar uma parte proporcional do total de páginas, com a introdução e a conclusão ocupando menos espaço do que o desenvolvimento.
- **7. Controle de Informações:** Evitar incluir informações desnecessárias que possam aumentar o comprimento do trabalho sem contribuir significativamente para o conteúdo. Concentre-se em informações relevantes, argumentos sólidos e evidências importantes para apoiar sua análise.
- **8. Síntese e Clareza:** O trabalho deverá ser conciso e claro em sua escrita. Evite repetições desnecessárias e redundâncias. Sintetize suas ideias e argumentos de forma eficaz para transmitir suas mensagens de maneira sucinta.
- **9. Formatação Adequada:** O trabalho deverá ser apresentado nas normas da ABNT de acordo com as diretrizes fornecidas, incluindo margens, espaçamento, tamanho da fonte e estilo de citação. Deve-se seguir o seguinte template de arquivo: hfps://bibliotecas.ufpr.br/wp- content/uploads/2022/03/template-artigo-de-periodico.docx **B RESOLUÇÃO**

**Estudo de Caso:** Implicações Éticas do Uso do ChatGPT **RESUMO** 

Este estudo de caso investiga as implicações éticas do uso do ChatGPT em diferentes contextos e propõe soluções responsáveis para lidar com esses dilemas. A análise aborda questões éticas significativas, incluindo a disseminação de informações, privacidade dos usuários, vieses algorítmicos, discriminação e

responsabilidade. Soluções responsáveis são propostas para equilibrar os interesses das partes interessadas e promover valores fundamentais como transparência, justiça e privacidade.

Palavras-chave: Inteligência Artificial, Ética, ChatGPT, Privacidade, Vieses.

**ABSTRACT** This case study investigates the ethical implications of using ChatGPT in different contexts and proposes responsible solutions to address these dilemmas. The analysis addresses significant ethical issues, including information dissemination, user privacy, algorithmic biases, discrimination, and responsibility. Responsible solutions are proposed to balance stakeholder interests and promote fundamental values such as transparency, fairness, and privacy.

**Keywords:** Artificial Intelligence, Ethics, ChatGPT, Privacy, Biases.

# 1. INTRODUÇÃO

De acordo com Aruda (2024), em novembro de 2022, houve uma certa euforia e também assombro quando a empresa OpenAI tornou público e acessível a plataforma ChatGPT. Segundo a empresa OpenAI (2022), o ChatGPT é um modelo de inteligência artificial (IA) que interage de maneira conversacional, ou seja, ele foi programado para que haja um diálogo entre o usuário e a máquina, permitindo que a plataforma responda a perguntas, admita erros, conteste premissas incorretas e rejeite solicitações inadequadas.

Com a definição dada pelos próprios criadores da ferramenta, é possível identificar inúmeras lacunas relacionadas a questões éticas provenientes do processo de geração de respostas por meio de IA. Nesse sentido, alguns dilemas abordam problemáticas éticas relacionadas a: (1) privacidade e segurança dos dados; (2) desinformação e uso indevido; (3) viés discriminatório; (4) impactos nas formas de trabalho; (5) responsabilidade e transparência; (6) interações humanas; e, por último, (7) robustez e precisão das informações geradas automaticamente. A discussão acerca de cada item acima gera um cenário multifacetado, com perspectivas provenientes de diversas áreas. No âmbito da jurisprudência, o contexto torna-se ainda mais complexo, exigindo uma compreensão mais aprofundada para que as questões éticas sejam cuidadosamente examinadas.

Neste trabalho abordaremos três questões éticas, sendo elas: (1) privacidade e segurança dos dados; (2) desinformação e uso indevido e (7) robustez e precisão das informações geradas de forma automática. Quando se fala em privacidade dos dados no uso do ChatGPT, é importante ressaltar que a funcionalidade da ferramenta requer um grande volume de dados para o treinamento do modelo da IA por trás dela. Além disso, é necessário que as requisições sejam feitas pelo usuário. Desta forma, fica evidente que os dados são elementos protagonistas e participam de todas as camadas desta aplicação.

No entanto, muitos debates têm sido levantados acerca da responsabilização pela posse dos dados, sejam eles sensíveis ou não, que a plataforma possui. Uma vez que os dados são fornecidos como entrada, eles se tornam parte do sistema. Neste sentido, há abertura de brechas relacionadas ao uso e compartilhamento de dados sensíveis de forma indevida, causando danos quase que irreversíveis aos afetados. Segundo Wu et al. (2024) embora existam mecanismos de proteção de privacidade no ChatGPT, como o bloqueio de acesso a dados pessoais de indivíduos, não é garantido que não ocorra nenhum vazamento.

Aliado a isso, o modelo de IA (LLM - Large Language Model, tradução para o português: Modelo de Linguagem de Grande Escala) usado para a geração de informações pelo ChatGPT não é simples de entender, o que torna o processo pouco transparente e suscetível a gerar informações pouco precisas sobre determinado assunto e até informações falsas. Essa falta de transparência pode manipular a opinião pública, causando desinformação e influenciando decisões de forma negativa. Para elucidar estas questões o trabalho discutirá sobre cada um dos itens de forma mais aprofundada trazendo alguns exemplos reais dos impactos gerados pelos dilemas éticos.

#### 2. DESENVOLVIMENTO

A revisão de literatura foca nos estudos prévios relacionados ao uso ético da IA e os desafios associados à implementação destes sistemas. Neste sentido, esta seção tem como objetivo trazer as definições necessárias presentes na literatura e os debates acerca dos impactos gerados por cada um dos dilemas éticos latentes no assunto.

#### 2.1 Privacidade e segurança dos dados

Para Westin (1967): A privacidade na era digital é um direito fundamental que deve ser protegido com rigor. A coleta massiva de dados pessoais sem o consentimento adequado dos usuários representa uma ameaça significativa à autonomia e à dignidade individual.

A privacidade dos usuários é uma questão em evidência neste momento em que a IA está sendo inserida na sociedade e as pessoas começam a conhecer e utilizar. A partir do momento em que interagimos com o ChatGPT, ele armazena todas as informações em sua base de dados para treinamento e isso implica que esses dados se tornarão públicos, a partir do momento em que se faz uma pergunta e a ferramenta entender que aqueles dados inseridos por outra pessoa são os melhores para uma resposta a qualquer usuário. Então, surge a grande preocupação quanto a privacidade dos dados, que dados serão utilizados para a análise, que dados serão apresentados nas respostas, existe algum contato humano com esta base de dados? Não basta somente informar ao usuário de que ele deve evitar inserir dados sensíveis, e sim, o ChatGPT deve utilizar critérios éticos que aperfeiçoem a ferramenta para que ela analise esses dados e gere uma nova conclusão e resposta, considerando os critérios da ética e segurança e dando a eles o máximo de importância.

No estudo apresentado por Khowaja et. al (2024) são discutidas as preocupações acerca da privacidade dos dados usados para o treinamento de LLMs, em especial do ChatGPT. Os autores mencionam que os dados utilizados para o treinamento do ChatGPT são sistematicamente extraídos de mensagens, websites, artigos, livros e informações pessoais sem a devida autorização. Além disso, de acordo com a política de privacidade do ChatGPT (OpenAI, 2024b), são coletadas informações como dados de interação do usuário com o site, configurações e tipo de navegador, e endereço IP, juntamente com o tipo de conteúdo com o qual os usuários interagem no ChatGPT. Eles também coletam informações sobre atividades de navegação em diferentes sites e ao longo de um determinado período de tempo. A política de privacidade também afirma que "Além disso, de tempos em tempos, podemos analisar o comportamento geral e as características dos usuários de nossos serviços e compartilhar informações agregadas, como estatísticas gerais de usuários, com terceiros, publicar essas informações agregadas ou torná-las geralmente disponíveis. Podemos coletar informações agregadas através dos Serviços, através de cookies e por outros meios descritos nesta política de privacidade"e ainda que "em certas circunstâncias, podemos fornecer suas informações pessoais a terceiros sem aviso prévio, a menos que exigido por lei".

### 2.2 Desinformação e uso indevido

Segundo Wardle e Derakhshan (2017): A desinformação, especialmente em plataformas digitais, pode ter consequências graves para a sociedade. Ela pode manipular opiniões, influenciar decisões eleitorais e até mesmo causar pânico em situações de crise.

Atualmente, a internet é a principal fonte de informação. Desta forma, o desafio não está apenas em obter um conteúdo relevante, mas também em filtrar as informações incorretas da grande quantidade disponível. Antes do lançamento do ChatGPT, as pessoas utilizavam outros métodos para confirmar a veracidade das informações, avaliando o nível de conhecimento do criador do conteúdo pelo rigor da linguagem, correção do formato e a fonte do texto, sendo esses indicadores de confiabilidade. No entanto, ao utilizar o ChatGPT, é possível observar que esses indicadores são excelentes nesses aspectos, dando a falsa sensação de confiabilidade. Os usuários podem cair na armadilha de confiar cegamente no conteúdo gerado pelo ChatGPT depois de executar testes simples. Essa confiança cega pode levar a julgamentos errados devido a hábitos arraigados. Em áreas críticas, como documentos de pesquisa médica, informações básicas de experimentos em larga escala e conteúdo de políticas, fazer referência a informações errôneas do ChatGPT e tirar conclusões incorretas pode ter consequências imprevisíveis. Esses riscos não decorrem de fraude intencional, mas sim dos erros do ChatGPT, apesar de sua aparência inicialmente confiável (Wu et al., 2024).

Outro aspecto muito importante relacionado ao uso do ChatGPT é a geração de notícias falsas, que têm implicações éticas de bastante relevância social e política e já são um dos principais problemas que vêm sendo enfrentados no Brasil. Estes Modelos de Linguagem Grande (LLMs) são muito populares e podem ser facilmente acessados para a geração de informações falsas ou enganosas que podem facilmente serem disseminadas nas redes sociais e plataforma digitais, podendo ter com consequências graves, como

por exemplo a manipulação de eleições e a incitação à violências (Khowaja et. al, 2024). Até que ponto as informações geradas pelo ChatGPT são verdadeiras? Esse é um grande desafio, pelo motivo de que em muitos momentos as informações necessitam ser validadas, muitas vezes depende de decisões que ainda precisam ser formuladas e votadas pela sociedade, categorias e governos. O fato de que as informações geradas pelo ChatGPT são originadas de uma base de dados e esta base foi construída pelo próprio ChatGPT, leva a conclusão de que esta base de dados deve ser constantemente validada, se possível por um órgão neutro externo que controle a legitimidade das informações e as distribua de forma a guardar a ética e os direitos de todas as partes envolvidas. Atualmente existe a versão gratuita e a versão paga do ChatGPT, e a principal diferença entre elas é a atualização da base de dados e funcionalidades. Já foi constatado que em determinados questionamentos ao ChatGPT, a ferramenta informou erroneamente e isso é muito sério, pois, a maior fatia de utilização é pela opção gratuita, que é a com a base de dados desatualizada e que fornece grandes vulnerabilidades. Uma regulamentação que defina a forma que esses dados são validados é de fundamental importância para garantir a segurança e confiabilidade das informações.

## 2.3 Robustez e precisão das informações geradas

Para O'Neil (2016): A automação na tomada de decisões promete eficiência e objetividade, mas também apresenta riscos consideráveis. Algoritmos podem perpetuar vieses existentes nos dados de treinamento, levando a decisões injustas ou discriminatórias.

Outro desafio é a tomada de decisões automatizadas, que dá poder para a ferramenta decidir o que fazer com base nas regras que foram criadas por seus programadores. Um grande exemplo são as IAs de recrutamento e seleção, que fazem a triagem com base em critérios definidos e podem muitas vezes serem injustas. Com relação ao ChatGPT, eu posso desde solicitar a criação de um código em determinada linguagem que diante de um perfil pré-selecionado busque a opção ideal em sua base, posso pedir para que ele gere uma imagem como por exemplo pedir uma imagem de uma pessoa bonita, e ele irá criar uma imagem de uma pessoa branca, loira, de olhos claros, ou seja, seguindo critérios já definidos pelos programadores ou seguindo uma tendência das imagens buscadas na internet. Deixar que a IA faça a escolha sozinha oferece riscos, pois ela foi configurada para buscar a melhor solução em sua base de dados e ela pode estar repleta de vieses.

# 3. CONCLUSÃO

Através de discussões em grupo e colaboração entre os membros da equipe, chegamos a algumas propostas para lidar com as implicações éticas do uso do ChatGPT. Nossa análise destacou a importância de criar regulamentações claras para garantir que o uso da IA respeite a privacidade dos usuários e evite a disseminação de informações falsas. Essas regulamentações devem definir padrões rigorosos para a coleta, armazenamento e uso dos dados, assegurando a proteção das informações pessoais.

Além disso, identificamos a necessidade de promover a transparência nos processos de decisão do ChatGPT. Isso pode ser alcançado através da implementação de mecanismos de auditoria e monitoramento contínuo, permitindo verificar e corrigir vieses algorítmicos, reduzindo o risco de discriminação.

Durante nossas discussões, também enfatizamos a importância de práticas de design ético na criação e atualização dos modelos de IA. Envolver equipes multidisciplinares e incluir feedback de usuários e especialistas externos pode contribuir para a criação de um sistema mais robusto e justo.

Finalmente, destacamos a necessidade de educar os usuários sobre os potenciais riscos e limitações do ChatGPT. Programas de conscientização e treinamentos específicos podem capacitar os usuários a utilizarem a IA de forma crítica e responsável, reconhecendo a importância da verificação das informações fornecidas pela ferramenta.

Em resumo, através da colaboração entre os membros da equipe e discussões aprofundadas, conseguimos desenvolver soluções práticas para lidar com as implicações éticas do uso do ChatGPT. As soluções propostas incluem regulamentações claras, transparência nos processos de decisão, práticas de design ético e educação do usuário, promovendo um uso mais seguro e justo dessa tecnologia.

### 4. REFERÊNCIAS

Aruda, E. P. Inteligência Artificial Generativa No Contexto Da Transformação do Trabalho Docente. 2024. Disponível em: <a href="https://doi.org/10.1590/0102-469848078">https://doi.org/10.1590/0102-469848078</a>>.

Deng, J.; Lin, Y. The benefits and challenges of ChatGPT: An overview. Frontiers in Computing and Intelligent Systems, v. 2, n. 2, p. 81-83, 2022.

Khowaja, S. A.; Khuwaja, P.; Dev, K.; Wang, W.; Nkenyereye, L. Chatgpt needs spade (sustainability, privacy, digital divide, and ethics) evaluation: A review. Cognitive Computation, p. 1-23, 2024.

O'Neil, C. Armas de Destruição Matemática: Como o Big Data Aumenta a Desigualdade e Ameaça a Democracia. Nova York: Crown Publishing Group, 2016.

OpenAI [internet]. Introducing ChatGPT. Disponível em: <a href="https://openai.com/index/chatgpt/">https://openai.com/index/chatgpt/</a>. Acessado em 01 de Julho de 2024a.

OpenAi [internet]. Privacy policy. Disponível em: <a href="https://openai.com/policies/privacy-policy/">https://openai.com/policies/privacy-policy/</a>. Acesso em 8 de Julho de 2024b.

Wardle, C.; Derakhshan, H. Desordem da Informação: Rumo a uma Estrutura Interdisciplinar para Pesquisa e Formulação de Políticas. Conselho da Europa, 2017.

Westin, A. F. Privacidade e Liberdade. Nova York: Atheneum, 1967. Wu, X.; Ran, D.; Jianbing, N. Unveiling security, privacy, and ethical concerns of ChatGPT. Journal of Information and Intelligence, v. 2, n. 2, p. 102-115, 2024.

### APÊNDICE 12 – GESTÃO DE PROJETOS DE IA

#### A - ENUNCIADO

### 1 Objetivo

Individualmente, ler e resumir – seguindo o template fornecido – um dos artigos abaixo: AHMAD, L.; ABDELRAZEK, M.; ARORA, C.; BANO, M; GRUNDY, J. Requirements practices and gaps when engineering human-centered Artificial Intelligence systems. Applied Soft Computing. 143. 2023. DOI https://doi.org/10.1016/j.asoc.2023.110421

NAZIR, R.; BUCAIONI, A.; PELLICCIONE, P.; Architecting ML-enabled systems: Challenges, best practices, and design decisions. The Journal of Systems & Software. 207. 2024. DOI https://doi.org/10.1016/j.jss.2023.111860

SERBAN, A.; BLOM, K.; HOOS, H.; VISSER, J. Software engineering practices for machine learning – Adoption, effects, and team assessment. The Journal of Systems & Software. 209. 2024. DOI https://doi.org/10.1016/j.jss.2023.111907

STEIDL, M.; FELDERER, M.; RAMLER, R. The pipeline for continuous development of artificial intelligence models – Current state of research and practice. The Journal of Systems & Software. 199. 2023. DOI https://doi.org/10.1016/j.jss.2023.111615

XIN, D.; WU, E. Y.; LEE, D. J.; SALEHI, N.; PARAMESWARAN, A. Whither AutoML? Understanding the Role of Automation in Machine Learning Workflows. In CHI Conference on Human Factors in Computing Systems (CHI'21), Maio 8-13, 2021, Yokohama, Japão. DOI https://doi.org/10.1145/3411764.3445306

## 2 Orientações adicionais

Escolha o artigo que for mais interessante para você. Utilize tradutores e o Chat GPT para entender o conteúdo dos artigos – caso precise, mas escreva o resumo em língua portuguesa e nas suas palavras. Não esqueça de preencher, no trabalho, os campos relativos ao seu nome e ao artigo escolhido. No template, você deverá responder às seguintes questões:

- Qual o objetivo do estudo descrito pelo artigo?
- Qual o problema/oportunidade/situação que levou a necessidade de realização deste estudo?
- Qual a metodologia que os autores usaram para obter e analisar as informações do estudo?
- Quais os principais resultados obtidos pelo estudo?

Responda cada questão utilizando o espaço fornecido no template, sem alteração do tamanho da fonte (Times New Roman, 10), nem alteração do espaçamento entre linhas (1.0). Não altere as questões do template. Utilize o editor de textos de sua preferência para preencher as respostas, mas entregue o trabalho em PDF.

## B - RESOLUÇÃO

Artigo: Requirements practices and gaps when engineering human-centered Artificial Intelligence systems

### Qual o objetivo do estudo descrito pelo artigo?

O objetivo do estudo é entender como os profissionais estão desenvolvendo sistemas de Inteligência Artificial que são centrados no ser humano. O artigo busca identificar quais práticas de Engenharia de Requisitos estão sendo usadas na indústria para esses sistemas de IA e quais aspectos centrados no ser humano estão sendo considerados. Além disso, pretende descobrir as lacunas entre

as diretrizes existentes, fornecidas por empresas como Google, Microsoft e Apple, e as práticas reais adotadas.

### Qual o problema/oportunidade/situação que levou à necessidade de realização desse estudo?

O estudo foi realizado porque existe uma necessidade de compreender melhor como os profissionais estão desenvolvendo sistemas de Inteligência Artificial que sejam centrados nas pessoas, especialmente na fase de Engenharia de Requisitos. Embora grandes empresas como Google, Microsoft e Apple tenham criado diretrizes para auxiliar nesse processo, ainda há pouca informação sobre como essas práticas estão sendo aplicadas na indústria. Além disso, há lacunas entre o que a pesquisa acadêmica propõe e o que é realmente feito nas empresas.

### Qual a metodologia que os autores usaram para obter e analisar as informações do estudo?

Os autores utilizaram uma metodologia em duas etapas. Primeiro, mapearam as diretrizes industriais existentes para o desenvolvimento de IA centrada no ser humano e as compararam com estudos acadêmicos sobre Engenharia de Requisitos para IA. Em seguida, realizaram uma pesquisa com profissionais para entender as práticas atuais na indústria em relação à Engenharia de Requisitos para IA e aos aspectos centrados no ser humano considerados.

#### Quais os principais resultados obtidos pelo estudo?

Os principais resultados obtidos pelo estudo foram que a maioria dos profissionais pesquisados concordou que todos os aspectos centrados no ser humano devem ser considerados durante a fase de Engenharia de Requisitos ao desenvolver sistemas de Inteligência Artificial. Além disso, identificou que muitas das ferramentas e métodos atualmente utilizados, como UML e Microsoft Office, não são adequados para gerenciar requisitos de software baseado em IA. Foram encontradas divergências entre a pesquisa acadêmica e as práticas da indústria na Engenharia de Requisitos para IA. Com base nesses resultados, os autores propuseram cinco recomendações para aprimorar as práticas e ferramentas na Engenharia de Requisitos de sistemas de IA centrados no ser humano.

### APÊNDICE 13 - FRAMEWORKS DE INTELIGÊNCIA ARTIFICIAL

#### A - ENUNCIADO

#### 1 Classificação (RNA)

Implementar o exemplo de Classificação usando a base de dados Fashion MNIST e a arquitetura RNA vista na aula **FRA - Aula 10 - 2.4 Resolução de exercício de RNA - Classificação**. Além disso, fazer uma breve explicação dos seguintes resultados:

- Gráficos de perda e de acurácia;
- Imagem gerada na seção "Mostrar algumas classificações erradas", apresentada na aula prática.

Informações:

- Base de dados: Fashion MNIST Dataset
- **Descrição:** Um dataset de imagens de roupas, onde o objetivo é classificar o tipo de vestuário. É semelhante ao famoso dataset MNIST, mas com peças de vestuário em vez de dígitos.
- Tamanho: 70.000 amostras, 784 features (28x28 pixels).
- Importação do dataset: Copiar código abaixo.

```
data = tf.keras.datasets.fashion_mnist
(x_train, y_train), (x_test, y_test) = fashion_mnist.load_data()
```

## 2 Regressão (RNA)

Implementar o exemplo de Classificação usando a base de dados Wine Dataset e a arquitetura RNA vista na aula FRA - Aula 12 - 2.5 Resolução de exercício de RNA - Regressão. Além disso, fazer uma breve explicação dos seguintes resultados:

- Gráficos de avaliação do modelo (loss);
- Métricas de avaliação do modelo (pelo menos uma entre MAE, MSE, R<sup>2</sup>).

Informações:

- Base de dados: Wine Quality
- **Descrição:** O objetivo deste dataset prever a qualidade dos vinhos com base em suas características químicas. A variável target (y) neste exemplo será o score de qualidade do vinho, que varia de 0 (pior qualidade) a 10 (melhor qualidade)
- Tamanho: 1599 amostras, 12 features.
- Importação: Copiar código abaixo.

```
1
2
   url = "https://archive.ics.uci.edu/ml/machine-learning-databases/wine
      -quality/winequality-red.csv"
   data = pd.read_csv(url, delimiter=';')
3
4
5
   Dica 1. Para facilitar o trabalho, renomeie o nome das colunas para
      português, dessa forma:
6
7
   data.columns = [
8
       'acidez_fixa',
                                   # fixed acidity
       'acidez_volatil',
9
                                   # volatile acidity
       'acido_citrico',
10
                                    citric acid
       'acucar_residual',
                                     residual sugar
11
       'cloretos',
12
                                   # chlorides
13
       'dioxido_de_enxofre_livre', # free sulfur dioxide
14
       'dioxido_de_enxofre_total', # total sulfur dioxide
       'densidade',
15
                                   # density
16
       'pH',
                                   # pH
17
       'sulfatos',
                                   # sulphates
       'alcool',
                                   # alcohol
18
19
       'score_qualidade_vinho'
                                                # quality
20
21
22
   Dica 2. Separe os dados (x e y) de tal forma que a última coluna (í
      ndice -1), chamada score_qualidade_vinho, seja a variável target (
      y)
```

#### 3 Sistemas de Recomendação

Implementar o exemplo de Sistemas de Recomendação usando a base de dados Base\_livos.csv e a arquitetura vista na aula **FRA - Aula 22 - 4.3 Resolução do Exercício de Sistemas de Recomendação**. Além disso, fazer uma breve explicação dos seguintes resultados:

- Gráficos de avaliação do modelo (loss);
- Exemplo de recomendação de livro para determinado Usuário.

Informações:

- Base de dados: Base\_livros.csv
- Descrição: Esse conjunto de dados contém informações sobre avaliações de livros (Notas), nomes de livros (Titulo), ISBN e identificação do usuário (ID\_usuario)
- Importação: Base de dados disponível no Moodle (UFPR Virtual), chamada Base\_livros (formato .csv).

#### 4 Deepdream

Implementar o exemplo de implementação mínima de Deepdream usando uma imagem de um felino - retirada do site Wikipedia - e a arquitetura Deepdream vista na aula **FRA - Aula 23 - Prática Deepdream**. Além disso, fazer uma breve explicação dos seguintes resultados:

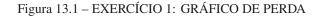
• Imagem onírica obtida por Main Loop;

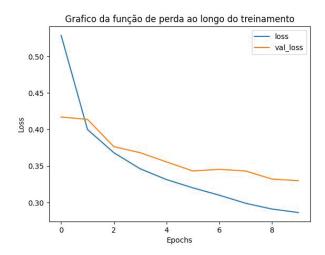
- Imagem onírica obtida ao levar o modelo até uma oitava;
- Diferenças entre imagens oníricas obtidas com Main Loop e levando o modelo até a oitava.
   Informações:
- Base de dados: https://commons.wikimedia.org/wiki/File:Felis\_catus-cat\_on\_snow.jpg
- Importação da imagem: Copiar código abaixo.

## **B-RESOLUÇÃO**

#### 1 Classificação (RNA)

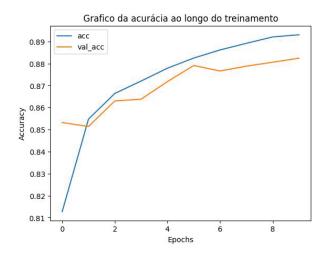
```
import tensorflow as tf
1
2
  import numpy as np
  import matplotlib.pyplot as plt
4
   from mlxtend.plotting import plot_confusion_matrix
5
   from sklearn.metrics import confusion_matrix
6
7
   (X_train, Y_train), (X_test, Y_test) = tf.keras.datasets.
      fashion_mnist.load_data()
8
9
  x_{train}, x_{test} = X_{train}/255.0, X_{test}/255.0
10
  i = tf.keras.layers.Input(shape=(28,28))
11
   x = tf.keras.layers.Flatten()(i)
12
  |x = tf.keras.layers.Dense(128, activation='relu')(x)
13
14 | x = tf.keras.layers.Dropout(0.2)(x)
   x = tf.keras.layers.Dense(10, activation='softmax')(x)
15
16
  model = tf.keras.Model(i, x)
17
  model.compile(optimizer='adam', loss='sparse_categorical_crossentropy
18
      ', metrics=['accuracy'])
19
20
   r = model.fit(x_train, Y_train, validation_data=(x_test, Y_test),
      epochs=10)
21
22 plt.plot(r.history['loss'], label='loss')
23
  plt.plot(r.history['val_loss'], label='val_loss')
24 plt.xlabel('Epochs')
25
  plt.ylabel('Loss')
  plt.title('Grafico da função de perda ao longo do treinamento')
26
27
  plt.legend()
```





```
plt.plot(r.history['accuracy'], label='acc')
plt.plot(r.history['val_accuracy'], label='val_acc')
plt.xlabel('Epochs')
plt.ylabel('Accuracy')
plt.title('Grafico da acurácia ao longo do treinamento')
plt.legend()
```

Figura 13.2 – EXERCÍCIO 1: GRÁFICO DE ACURÁCIA



```
1  y_pred = model.predict(x_test).argmax(axis=1)
2  cm = confusion_matrix(Y_test, y_pred)
3  plot_confusion_matrix(conf_mat=cm, figsize=(7,7), show_normed=True)
```

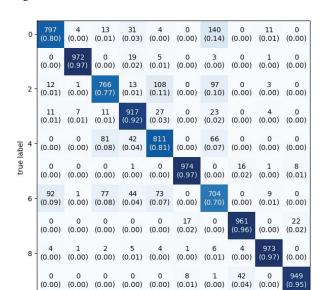
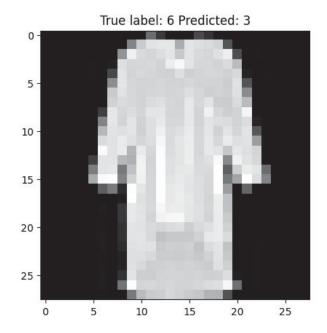


Figura 13.3 – EXERCÍCIO 1: MATRIZ DE CONFUSÃO

```
misclassified = np.where(y_pred != Y_test)[0]
i = np.random.choice(misclassified)
plt.imshow(x_test[i], cmap='gray')
plt.title("True label: %s Predicted: %s" % (Y_test[i], y_pred[i]));
```

predicted label

Figura 13.4 – EXERCÍCIO 1: APRESENTANDO CLASSIFICAÇÃO INCORRETA



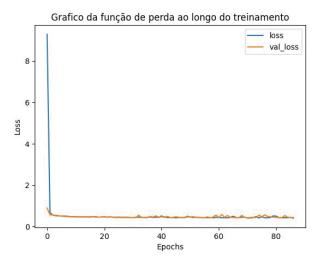
### 2 Regressão (RNA)

```
import tensorflow as tf
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
```

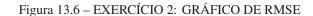
```
5 from tensorflow.python.keras import backend
  from sklearn.model_selection import train_test_split
7
  from sklearn.metrics import r2_score, mean_squared_error
  from math import sqrt
10 | url = "https://archive.ics.uci.edu/ml/machine-learning-databases/wine
      -quality/winequality-red.csv"
11
12 data = pd.read_csv(url, delimiter=';')
13
14 | data.columns = [
15 'acidez_fixa', # fixed acidity
16 'acidez_volatil', # volatile acidity
17
   'acido_citrico', # citric acid
   'acucar_residual', # residual sugar
18
19 'cloretos', # chlorides
20 'dioxido_de_enxofre_livre', # free sulfur dioxide
21 'dioxido_de_enxofre_total', # total sulfur dioxide
22
   'densidade', # density
23
  'pH', # pH
24 'sulfatos', # sulphates
   'alcool', # alcohol
26 'score_qualidade_vinho' # quality
27
28
29 data = data.values
30
  # Dica 2. Separe os dados (x e y) de tal forma que a última coluna (í
31
      ndice -1), chamada score_qualidade_vinho
32
33 X = data[:, 0:-1].astype(float)
|Y| = data[:, -1].astype(float)
35
  X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size
36
      =0.25)
37
38 | i = tf.keras.layers.Input(shape=(11,))
39 | x = tf.keras.layers.Dense(100, activation='relu')(i)
   x = tf.keras.layers.Dense(1)(x)
41
  | model = tf.keras.Model(i, x) |
42
43
  # Criação de funções para as métricas r2 e RMSE serem inseridas no
44
   def rmse(y_true, y_pred):
45
       return backend.sqrt(backend.mean(backend.square(y_true- y_pred),
      axis=-1))
46
47
   def r2(y_true, y_pred):
48
       media = backend.mean(y_true)
49
       num = backend.sum(backend.square(y_true - y_pred))
       den = backend.sum(backend.square(y_true - media))
50
```

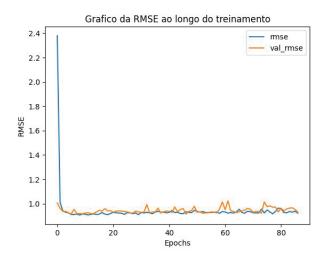
```
51
       return (1.0 - num/den)
52
53
   # Compilação
   optimizer = tf.keras.optimizers.Adam(learning_rate=0.001)
   model.compile(optimizer=optimizer, loss='mse', metrics=[r2, rmse])
56
57
   # Early Stop para Epochs
58
   early_stop = tf.keras.callbacks.EarlyStopping(monitor='val_loss',
      patience=20, restore_best_weights=True)
59
  r = model.fit(X_train, Y_train, validation_data=(X_test, Y_test),
60
      epochs=1500, callbacks=[early_stop])
61
62 plt.plot(r.history['loss'], label='loss')
  plt.plot(r.history['val_loss'], label='val_loss')
63
64 plt.xlabel('Epochs')
65 plt.ylabel('Loss')
66 plt.title('Grafico da função de perda ao longo do treinamento')
67
  plt.legend()
```

Figura 13.5 – EXERCÍCIO 2: GRÁFICO DE PERDA



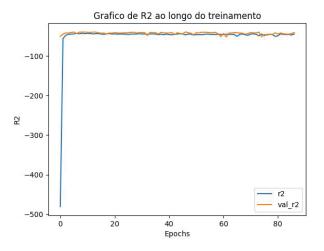
```
plt.plot(r.history['rmse'], label='rmse')
plt.plot(r.history['val_rmse'], label='val_rmse')
plt.xlabel('Epochs')
plt.ylabel('RMSE')
plt.title('Grafico da RMSE ao longo do treinamento')
plt.legend()
```





```
plt.plot(r.history['r2'], label='r2')
plt.plot(r.history['val_r2'], label='val_r2')
plt.xlabel('Epochs')
plt.ylabel('R2')
plt.title('Grafico de R2 ao longo do treinamento')
plt.legend()
```

Figura 13.7 – EXERCÍCIO 2: GRÁFICO DE R2



```
1  y_pred = model.predict(X_test).flatten()
2  mse = mean_squared_error(Y_test, y_pred)
3  rmse = sqrt(mse)
4  r2 = r2_score(Y_test, y_pred)
5  
6  print("mse:", mse)
7  print("rmse:", rmse)
8  print("r2:", r2)
```

mse: 0.4207298387142231 rmse: 0.6486369082269549

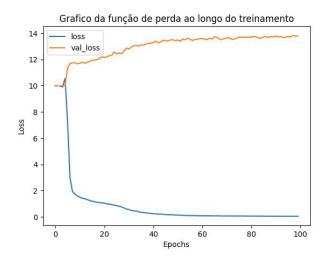
#### r2: 0.38145588854025325

### 3 Sistemas de Recomendação

```
import numpy as np
2 import pandas as pd
  import tensorflow as tf
3
4 import matplotlib.pyplot as plt
   from tensorflow.keras.models import Model
  from tensorflow.keras.layers import Input, Embedding, Dense,
      Embedding, Flatten, Concatenate
7
  from tensorflow.keras.optimizers import SGD, Adam
8
  from sklearn.utils import shuffle
9
10
  df = pd.read_csv('Base_livros(in).csv', sep=',', on_bad_lines='skip')
  df['Notas;;;;;'] = df['Notas;;;;;;'].str.replace(r'[^0-9]', '',
11
      regex=True).astype(float)
12 | df.rename(columns={'Notas;;;;;;': 'Notas'}, inplace=True)
13
  df.dropna(inplace=True)
14
  categorical_columns = ['ID_usuario', 'Titulo', 'Autor', 'Editora']
15
16 for cat_column in categorical_columns:
17
       df[cat_column] = pd.Categorical(df[cat_column])
18
       df['new_' + cat_column] = df[cat_column].cat.codes
19
20
  N = len(set(df['ID_usuario']))
  M = len(set(df['new_Titulo']))
22
23 | # dimensão do embedding
24 | K = 10
25
26 # usuário
   u = Input(shape=(1,))
28 |u_{mb}| = Embedding(N, K)(u) # saída : num_samples, 1, K
29 | u_emb = Flatten()(u_emb) # saída : num_samples, K
30 # filme
31 \mid m = Input (shape=(1,))
   m_emb = Embedding(M, K)(m) # saída : num_samples, 1, K
33 | m_emb = Flatten()(m_emb) # saída : num_samples, K
34 x = Concatenate()([u_emb, m_emb])
35 |x| = Dense(1024, activation="relu")(x)
36 \mid x = Dense(1)(x)
  model = Model(inputs=[u, m], outputs=x)
37
38
39 | model.compile(loss="mse", optimizer=SGD(learning_rate=0.08, momentum
      =0.9)
40
41
  |user_ids, movie_ids, ratings = shuffle(df['new_ID_usuario'], df['
      new_Titulo'], df['Notas'])
42 Ntrain = int(0.8 * len(ratings)) # separar os dados 80% x 20%
43
  |train_user = user_ids[:Ntrain]
44 | train_movie = movie_ids[:Ntrain]
```

```
45 | train_ratings = ratings[:Ntrain]
46
  test_user = user_ids[Ntrain:]
47
   test_movie = movie_ids[Ntrain:]
48
  test_ratings = ratings[Ntrain:]
49
50 # centralizar as notas
51 | avg_rating = train_ratings.mean()
52 train_ratings = train_ratings - avg_rating
53
  test_ratings = test_ratings - avg_rating
54
55 \mid \text{epochs} = 100
56 | r = model.fit(
57 | x=[train_user, train_movie],
58
   y=train_ratings,
59
   epochs=epochs,
60 batch_size=1024,
  verbose=2, # não imprime o progresso
62 | validation_data=([test_user, test_movie], test_ratings)
63
   )
64
65 | plt.plot(r.history['loss'], label='loss')
   plt.plot(r.history['val_loss'], label='val_loss')
67 | plt.xlabel('Epochs')
68
  plt.ylabel('Loss')
69
   plt.title('Grafico da função de perda ao longo do treinamento')
70
   plt.legend()
```

Figura 13.8 – EXERCÍCIO 3: GRÁFICO DE PERDA



```
input_usuario = np.repeat(a=8263, repeats=M)
film = np.array(list(set(movie_ids)))
preds = model.predict([input_usuario, film])

# descentraliza as predições
rat = preds.flatten() + avg_rating
# indice da maior nota
```

```
8 | idx = np.argmax(rat)
9 | print("Recomendação: Filme - ", film[idx], " / ", rat[idx] , "*")
```

Recomendação: Filme - 66366 / 19.101202 \*

## 4 Deepdream

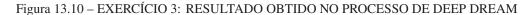
```
1 | import tensorflow as tf
2 import numpy as np
  import matplotlib as mpl
   import IPython.display as display
5
  import PIL.Image
6
7
  url = "https://commons.wikimedia.org/wiki/Special:FilePath/
      Felis_catus-cat_on_snow.jpg"
8
9
   # Download da imagem e gravação em array Numpy
  def download(url, max_dim=None):
10
       name = url.split('/')[-1]
11
12
       image_path = tf.keras.utils.get_file(name, origin=url)
13
       img = PIL.Image.open(image_path)
14
       if max dim:
       img.thumbnail((max_dim, max_dim))
15
16
       return np.array(img)
17
18
   # Normalização da imagem
19
   def deprocess(img):
20
       img = 255*(img + 1.0)/2.0
21
       return tf.cast(img, tf.uint8)
22
23
   # Mostra a imagem
24
  def show(img):
25
       display.display(PIL.Image.fromarray(np.array(img)))
26
27
   # Redução do tamanho da imagem para facilitar o trabalho da RNN
  original_img = download(url, max_dim=500)
28
  show(original_img)
```

Figura 13.9 – EXERCÍCIO 3: IMAGEM BAIXADA DA URL FORNECIDA



```
base_model = tf.keras.applications.InceptionV3(include_top=False,
      weights='imagenet')
2
3
   # Maximizando as ativações das camadas
   names = ['mixed3', 'mixed5']
   layers = [base_model.get_layer(name).output for name in names]
5
7
   # Criação do modelo
   dream_model = tf.keras.Model(inputs=base_model.input, outputs=layers)
8
9
10
   def calc_loss(img, model):
11
       # Passe a imagem pelo modelo para recuperar as ativações.
       # Converte a imagem em um batch de tamanho 1.
12
13
       img_batch = tf.expand_dims(img, axis=0)
14
       layer_activations = model(img_batch)
15
16
       if len(layer_activations) == 1:
17
           layer_activations = [layer_activations]
18
19
       losses = []
20
       for act in layer_activations:
21
           loss = tf.math.reduce_mean(act)
22
           losses.append(loss)
23
24
       return tf.reduce_sum(losses)
25
26
   class DeepDream(tf.Module):
27
       def __init__(self, model):
28
           self.model = model
29
       @tf.function(
           input_signature=(
30
31
                tf.TensorSpec(shape=[None, None, 3], dtype=tf.float32),
32
               tf.TensorSpec(shape=[], dtype=tf.int32),
33
               tf.TensorSpec(shape=[], dtype=tf.float32),)
34
35
       def __call__(self, img, steps, step_size):
           print("Tracing")
36
37
           loss = tf.constant(0.0)
38
39
           for n in tf.range(steps):
40
                with tf.GradientTape() as tape:
                    # Gradientes relativos a img
41
42
                    tape.watch(img)
43
                    loss = calc_loss(img, self.model)
44
45
            # Calculo do gradiente da perda em relação aos pixels da
      imagem de entrada.
46
            gradients = tape.gradient(loss, img)
47
48
           # Normalizacao dos gradintes
49
           gradients /= tf.math.reduce_std(gradients) + 1e-8
```

```
50
           # Na subida gradiente, a "perda" é maximizada.
51
52
           # Você pode atualizar a imagem adicionando diretamente os
      gradientes (porque eles têm o mesmo f
           img = img + gradients*step size
53
54
           img = tf.clip_by_value(img, -1, 1)
55
56
           return loss, img
57
   deepdream = DeepDream(dream_model)
58
59
60
   def run_deep_dream_simple(img, steps=100, step_size=0.01):
       img = tf.keras.applications.inception_v3.preprocess_input(img)
61
62
       img = tf.convert_to_tensor(img)
63
       step_size = tf.convert_to_tensor(step_size)
64
       steps_remaining = steps
65
       step = 0
       while steps_remaining:
66
67
           if steps_remaining>100:
68
               run_steps = tf.constant(100)
69
           else:
               run_steps = tf.constant(steps_remaining)
70
71
72
           steps_remaining -= run_steps
73
           step += run_steps
74
           loss, img = deepdream(img, run_steps, tf.constant(step_size))
75
76
77
           display.clear_output(wait=True)
78
           show(deprocess(img))
79
           print ("Step {}, loss {}".format(step, loss))
80
       result = deprocess(img)
81
82
       display.clear_output(wait=True)
83
       show(result)
84
       return result
85
86
   dream_img = run_deep_dream_simple(img=original_img, steps=100,
87
      step_size=0.01)
```





```
import time
1
2
  start = time.time()
3
4 OCTAVE_SCALE = 1.30
5 | img = tf.constant(np.array(original_img))
6 base_shape = tf.shape(img)[:-1]
7
  float_base_shape = tf.cast(base_shape, tf.float32)
  for n in range (-2, 3):
8
9
       new_shape = tf.cast(float_base_shape*(OCTAVE_SCALE**n), tf.int32)
       img = tf.image.resize(img, new_shape).numpy()
10
11
       img = run_deep_dream_simple(img=img, steps=50, step_size=0.01)
12
13
  display.clear_output(wait=True)
14 | img = tf.image.resize(img, base_shape)
15 | img = tf.image.convert_image_dtype(img/255.0, dtype=tf.uint8)
16 | show(img)
17
  end = time.time()
  end-start
```

Figura 13.11 – EXERCÍCIO 3: RESULTADO OBTIDO NO PROCESSO DE DEEP DREAM À OITAVA



## APÊNDICE 14 - VISUALIZAÇÃO DE DADOS E STORYTELLING

#### A - ENUNCIADO

Escolha um conjunto de dados brutos (ou uma visualização de dados que você acredite que possa ser melhorada) e faça uma visualização desses dados (de acordo com os dados escolhidos e com a ferramenta de sua escolha) Desenvolva uma narrativa/storytelling para essa visualização de dados considerando os conceitos e informações que foram discutidas nesta disciplina. Não esqueça de deixar claro para seu possível público alvo qual o objetivo dessa visualização de dados, o que esses dados significam, quais possíveis ações podem ser feitas com base neles.

# Entregue em um PDF:

- O **conjunto de dados brutos (ou uma visualização de dados** que você acredite que possa ser melhorada);
- Explicação do **contexto e o publico-alvo** da visualização de dados e do storytelling que será desenvolvido;
- A visualização desses dados (de acordo com os dados escolhidos e com a ferramenta de sua escolha) explicando a escolha do tipo de visualização e da ferramenta usada; (50 pontos)

# B - RESOLUÇÃO

#### Contexto e Público-Alvo

Os dados de pesquisa pelo termo "Passagem de avião" mostram um aumento expressivo em dezembro e janeiro. Esses picos coincidem com as férias no Brasil, indicando que muitos consumidores deixam para comprar passagens de última hora, em vez de se planejarem com antecedência. O público alvo serão gestores e analistas de marketing de companhias aéreas e agências de viagens, que podem ajustar campanhas e promoções.

## Storytelling

A visualização traça a evolução das buscas pelo termo "passagem de avião" ao longo do tempo e revela um padrão sazonal marcante. Ao longo do ano, os níveis de interesse permanecem relativamente constantes, mas observam-se picos significativos em dezembro e janeiro. Esses aumentos coincidem com as férias de fim de ano e o início das férias escolares, períodos em que os consumidores tendem a adiar o planejamento da viagem e optam por comprar passagens em cima da hora.

A narrativa constrói a ideia de que essa concentração de buscas reflete um comportamento de compra impulsivo, em que a ausência de um planejamento prévio leva os viajantes a buscarem alternativas somente quando a viagem se aproxima. Esse insight sugere oportunidades para que companhias aéreas e agências de viagens desenvolvam estratégias de marketing mais eficazes, incentivando a compra antecipada por meio de promoções e ofertas exclusivas.

Ao explorar essa visualização, o público-alvo poderá compreender não apenas o padrão sazonal das buscas, mas também as implicações desse comportamento para o setor de turismo, possibilitando a adoção de estratégias que contribuam para uma distribuição mais equilibrada da demanda ao longo do ano.

Figura 14.1 – ANÁLISE DE TENDÊNCIAS DE BUSCA POR PASSAGEM DE AVIÃO

Análise de Tendências de Busca por Indicação de Compras de Última Hora Passagem de avião Conjunto de Dados Contexto e Público-Alvo Os dados de pesquisa pelo termo "Passagem de avião" mostram um aumento expressivo em dezembro e janeiro. Esses picos coincidem com as férias no Brasil, indicando que muitos consumidores deixam para comprar passagens de última hora, em vez de se planejarem com antecedência. O público alvo serão gestores e analistas de marketing de companhias aéreas e agências de viagens, que podem ajustar campanhas e promoções. Interesse (0 a 100) 9 Feb 2025 2 Feb 2025 26 Ian 2025 19 Jan 2025 8 Dec 2024 1 Dec 2024 Storytelling A visualização traça a evolução das buscas pelo termo "passagem de avião" ao longo do tempo e revela um padrão sazonal marcante. Ao longo do ano, os níveis de interesse permanecem relativamente constantes, mas observam-se picos significativos em dezembro e janeiro. Esses aumentos coincidem com as férias de fim de ano e o início das férias escolares, períodos em que os consumidores tendem a adiar o planejamento da viagem e optam por comprar passagens em cima da hora. 1... 17 Nov 2024 1... 10 Nov 2024 1... 3 Nov 2024 A narrativa constrói a ideia de que essa concentração de buscas reflete um comportamento de compra impulsivo, em que a ausência de um planejamento prévio leva os viajantes a buscarem alternativas somente quando a viagem se aproxima. Esse insight sugere oportunidades para que companhias aéreas e agências de viagens desenvolvam estratégias de marketing mais eficazes, incentivando a compra antecipada por meio de promoções e ofertas exclusivas. Ao explorar essa visualização, o público-alvo poderá compreender não apenas o padrão sazonal das buscas, mas também as implicações desse comportamento para o setor de turismo, possibilitando a adoção de estratégias que contribuam para uma distribuição mais equilibrada da demanda ao longo do ano.

## APÊNDICE 15 - TÓPICOS EM INTELIGÊNCIA ARTIFICIAL

### A - ENUNCIADO

### 1) Algoritmo Genético

Problema do Caixeiro Viajante

A Solução poderá ser apresentada em: Python (preferencialmente), ou em R, ou em Matlab, ou em C ou em Java. Considere o seguinte problema de otimização (a escolha do número de 100 cidades foi feita simplesmente para tornar o problema intratável. A solução ótima para este problema não é conhecida). Suponha que um caixeiro deva partir de sua cidade, visitar clientes em outras 99 cidades diferentes, e então retornar à sua cidade. Dadas as coordenadas das 100 cidades, descubra o percurso de menor distância que passe uma única vez por todas as cidades e retorne à cidade de origem.

Para tornar a coisa mais interessante, as coordenadas das cidades deverão ser sorteadas (aleatórias), considere que cada cidade possui um par de coordenadas (x e y) em um espaço limitado de 100 por 100 pixels. O relatório deverá conter no mínimo a primeira melhor solução (obtida aleatoriamente na geração da população inicial) e a melhor solução obtida após um número mínimo de 1000 gerações. Gere as imagens em 2d dos pontos (cidades) e do caminho. Sugestão:

- a) considere o cromossomo formado pelas cidades, onde a cidade de início (escolhida aleatoriamente) deverá estar na posição 0 e 100 e a ordem das cidades visitadas nas posições de 1 a 99 deverão ser definidas pelo algoritmo genético.
- b) A função de avaliação deverá minimizar a distância euclidiana entre as cidades (os pontos).
- c) Utilize no mínimo uma população com 100 indivíduos;
- d) Utilize no mínimo 1% de novos indivíduos obtidos pelo operador de mutação;
- e) Utilize no mínimo de 90% de novos indivíduos obtidos pelo método de cruzamento (crossoverox);
- f) Preserve sempre a melhor solução de uma geração para outra.

Importante: A solução deverá implementar os operadores de "cruzamento" e "mutação".

## 2) Compare a representação de dois modelos vetoriais

Pegue um texto relativamente pequeno, o objetivo será visualizar a representação vetorial, que poderá ser um vetor por palavra ou por sentença. Seja qual for a situação, considere a quantidade de palavras ou sentenças onde tenha no mínimo duas similares e no mínimo 6 textos, que deverão produzir no mínimo 6 vetores. Também limite o número máximo, para que a visualização fique clara e objetiva.

O trabalho consiste em pegar os fragmentos de texto e codificá-las na forma vetorial. Após obter os vetores, imprima-os em figuras (plot) que demonstrem a projeção desses vetores usando a PCA.

O PDF deverá conter o código-fonte e as imagens obtidas.

## B - RESOLUÇÃO

### 1) Algoritmo Genético

```
import math
import random
import matplotlib.pyplot as plt

pefinindo parâmetros do problema
```

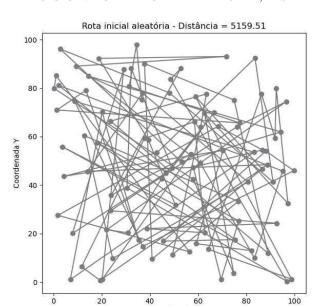
```
6 NUM_CIDADES = 100
7
   TAMANHO_POPULACAO = 100
   NUM_GERACOES = 5000
8
9
   # Função para calcular a distancia total de uma rota (caminho fechado
10
       que passa
11
  def calcular_distancia_total(rota, coordenadas):
12
       11 11 11
       Calcula a distância total percorrida seguindo a sequência de
13
      cidades em 'rota'.
       A distância é a soma das distâncias euclidianas entre cada par
14
      consecutivo.
15
       *Importante: A rota começa e termina na cidade de origem (rota
       [0]).
       11 11 11
16
17
       distancia = 0.0
       # Soma as distâncias entre cada cidade e a próxima na rota
18
19
       for i in range(len(rota) - 1):
20
           cidade_atual = rota[i]
21
           proxima_cidade = rota[i+1]
22
           # Obtém as coordenadas das duas cidades
23
           x1, y1 = coordenadas[cidade atual]
           x2, y2 = coordenadas[proxima_cidade]
24
25
           # Calcula a distancia Euclidiana entre cidade atual e proxima
       cidade
           distancia += math.hypot(x2-x1, y2-y1)
26
27
       return distancia
28
29
   # Função de seleção por torneio para escolher um individuo (rota)
      como pai
30
   def selecionar_pai_por_torneio(populacao, distancias, k=3):
31
32
       Seleciona e retorna um individuo da população usando seleção por
      torneio.
33
       Aleatoriamente escolhe 'k' indivíduos e retorna o que tiver a
      menor distância.
       11 11 11
34
35
       melhor indice = -1
       melhor_distancia = float('inf')
36
37
       # Realiza o torneio entre individuos aleatórios
38
       for i in range(k):
39
           indice = random.randrange(len(populacao)) # escolhe um indice
       aleatorio
40
           # Verifica se este individuo é o melhor do torneio até agora
41
           if distancias[indice] < melhor_distancia:</pre>
               melhor_distancia = distancias[indice]
42
43
               melhor_indice = indice
44
       # Retorna uma cópia da rota vencedora (melhor distância)
45
       return list(populacao[melhor_indice])
46
   def crossover_OX(pail, pai2):
47
```

```
11 11 11
48
49
       Realiza o crossover OX (Order Crossover) entre dois cromossomos
50
       Retorna um cromossomo filho válido, com cidade inicial fixa.
51
52
       tamanho = len(pai1)
       # Inicializa o filho garantindo cidade inicial e final fixas
53
54
       filho = [None] * tamanho
55
       filho[0] = pail[0]
       filho[-1] = pail[-1]
56
57
58
       # Escolhe dois pontos de corte válidos e diferentes
59
       corte1, corte2 = sorted(random.sample(range(1, tamanho - 1), 2))
60
61
       # Copia segmento do pail diretamente para o filho
62
       segmento_pai1 = pai1[corte1:corte2 + 1]
       filho[corte1:corte2+1] = segmento_pai1
63
64
65
       # Guarda as cidades ja inseridas no filho
66
       cidades_no_filho = set(segmento_pail)
67
       # Preenche o restante das posições com cidades de pai2 (na ordem
68
      correta)
69
       pos_filho = (corte2 + 1) % (tamanho - 1)
70
       pos_pai2 = (corte2 + 1) % (tamanho - 1)
71
72
       # Repete até o filho ser totalmente preenchido (sem None)
       while None in filho[1:-1]:
73
74
           cidade = pai2[pos_pai2]
75
           if cidade not in cidades_no_filho:
76
               filho[pos_filho] = cidade
77
                cidades_no_filho.add(cidade)
78
               pos_filho = pos_filho + 1 if pos_filho + 1 < tamanho - 1</pre>
      else 1
79
           pos_pai2 = pos_pai2 + 1 if pos_pai2 + 1 < tamanho - 1 else 1
80
       return filho
81
82
83
   # Função de mutação que troca duas cidades de lugar na rota (mutação
      por troca)
   def mutacao(rota):
84
85
       Realiza uma mutação simples em uma rota, trocando de posição duas
86
       cidades aleatórias.
87
       A cidade de origem (posição 0 e última posição) não é alterada.
       Retorna uma nova rota mutada.
88
89
90
       # Cria uma cópia da rota para mutação (para não alterar a rota
      original diretamente)
91
       rota_mutada = list(rota)
92
```

```
93
        # Seleciona aleatoriamente duas posições diferentes na parte
       intermediaria
94
        i = random.randint(1, len(rota_mutada) - 2)
95
        j = random.randint(1, len(rota_mutada) - 2)
96
97
        while i == j:
            j = random.randint(1, len(rota_mutada) - 2)
98
99
        # Troca as cidades das posições i e j
100
101
        rota_mutada[i], rota_mutada[j] = rota_mutada[j], rota_mutada[i]
102
        return rota_mutada
103
104
    # Gera coordenadas aleatórias para cada cidade dentro de um espaço
       100x100
    coordenadas = [] # Lista para armazenar tuples (x, y) de coordenadas
105
       de cada cidade
106
    for i in range(NUM_CIDADES):
107
        x = random.random() * 100 # coordenada x aleatoria entre 0 e 100
108
        y = random.random() * 100 # coordenada y aleatoria entre 0 e 100
109
        coordenadas.append((x, y))
   print("Coordenadas:", coordenadas)
110
111
112
   # Gera a população inicial de rotas aleatórias
113
   populacao = []
   distancias = []
114
   cidade_origem = 0 # definindo a cidade de indice 0 como cidade de
115
       origem
116
117
   # Inicializa as variaveis para acompanhar a melhor solução encontrada
   melhor_rota = None
118
119
   melhor_distancia = float('inf')
120
121
    for i in range(TAMANHO_POPULACAO):
122
        # Cria uma ordem aleatória das cidades (exceto a cidade de origem
123
        outras_cidades = list(range(1, NUM_CIDADES))
124
        random.shuffle(outras_cidades)
125
        # Monta a rota completa: cidade de origem no início, seguida das
       outras cidades
126
        rota = [cidade_origem] + outras_cidades + [cidade_origem]
        # Adiciona a rota gerada à população
127
128
        populacao.append(rota)
129
        # Calcula a distância dessa rota e guarda
        distancias.append(calcular_distancia_total(rota, coordenadas))
130
131
    # Seleciona uma rota inicial aleatória (por exemplo, a primeira da
132
       população)
   rota_inicial = populacao[0]
133
    distancia_inicial = distancias[0]
134
135
   print("Distancia inicial:", distancia_inicial)
136
```

```
137
    # Identifica o melhor indivíduo da população inicial (com menor
       distancia)
    for i, rota in enumerate (populacao):
138
        d = distancias[i]
139
        if d < melhor distancia:</pre>
140
141
            melhor distancia = d
            melhor_rota = rota
142
143
144
    # Loop principal do Algoritmo Genético
    for geracao in range (NUM GERACOES):
145
146
        nova_populacao = []
147
        nova_distancias = []
148
149
        # Elitismo: preserva a melhor rota da geração atual na próxima
        nova_populacao.append(melhor_rota)
150
151
        nova_distancias.append(melhor_distancia)
152
153
        # Determina quantos novos indivíduos serão gerados por crossover
       e por mutação
        num_por_crossover = int(TAMANHO_POPULACAO * 0.90)
154
155
        num_por_mutacao = int(TAMANHO_POPULACAO * 0.09)
156
        if num_por_mutacao < 1:</pre>
157
            num_por_mutacao = 1
158
159
        total_novos = 1 + num_por_crossover + num_por_mutacao # 1
       elitismo + novos
160
        if total_novos > TAMANHO_POPULACAO:
161
            num_por_crossover -= (total_novos - TAMANHO_POPULACAO)
        elif total_novos < TAMANHO_POPULACAO:</pre>
162
163
            num_por_crossover += (TAMANHO_POPULACAO - total_novos)
164
165
        # Gera novos individuos por crossover
166
        for _ in range(num_por_crossover):
            pai1 = selecionar_pai_por_torneio(populacao, distancias, k=3)
167
            pai2 = selecionar_pai_por_torneio(populacao, distancias, k=3)
168
169
            filho = crossover_OX(pai1, pai2)
170
            nova_populacao.append(filho)
            nova_distancias.append(calcular_distancia_total(filho,
171
       coordenadas))
172
173
        # Gera novos individuos por mutação
174
        for _ in range(num_por_mutacao):
            base = selecionar_pai_por_torneio(populacao, distancias, k=3)
175
176
            individuo_mutado = mutacao(base)
177
            nova_populacao.append(individuo_mutado)
            nova_distancias.append(calcular_distancia_total(
178
       individuo_mutado, coordenadas))
179
        # Se por algum motivo excedemos o tamanho da população, cortamos
180
       o excedente
181
        if len(nova_populacao) > TAMANHO_POPULACAO:
```

```
182
            nova_populacao = nova_populacao[:TAMANHO_POPULACAO]
183
            nova_distancias = nova_distancias[:TAMANHO_POPULACAO]
184
185
        # Atualiza a população para a proxima geração
186
        populacao = nova_populacao
187
        distancias = nova_distancias
188
189
        # Atualiza o melhor individuo
190
        for i, d in enumerate(distancias):
191
            if d < melhor distancia:</pre>
192
                melhor_distancia = d
193
                melhor_rota = populacao[i]
194
195
        # A cada 100 gerações, imprime a evolução
196
        if (geracao + 1) % 100 == 0 or geracao == 0:
197
            print(f"Geração {geracao + 1}: melhor distância encontrada {
       melhor_distancia}")
198
199
   # Imprime as distâncias da rota inicial e da melhor rota final
       encontrada
200
   print(f"\nDistancia total da rota inicial: {distancia_inicial:.2f}")
201
   print(f"Distância total da melhor rota final: {melhor_distancia:.2f}\
       n")
202
203
   # Gera o gráfico da rota inicial aleatória
   plt.figure(figsize=(6,6))
   xs = [coordenadas[cidade][0] for cidade in rota_inicial]
205
206 ys = [coordenadas[cidade][1] for cidade in rota_inicial]
207
   plt.plot(xs, ys, marker='o', color='gray')
208
   plt.scatter(xs, ys, color="blue")
209 | plt.title(f"Rota inicial aleatória - Distancia = {distancia_inicial
       :.2f}")
210 plt.xlabel("Coordenada X")
   plt.ylabel("Coordenada Y")
211
212
   plt.tight_layout()
```



Coordenada X

Figura 15.1 – EXERCÍCIO 1: ROTA INICIAL ALEATÓRIA, DISTÂNCIA = 5159,51

```
# Gera o gráfico da melhor rota final apos a evolução
  plt.figure(figsize=(6,6))
2
  xs_best = [coordenadas[cidade][0] for cidade in melhor_rota]
3
  ys_best = [coordenadas[cidade][1] for cidade in melhor_rota]
  plt.plot(xs_best, ys_best, marker='o', color='green')
5
6 plt.scatter(xs_best, ys_best, color='red')
  plt.title(f"Melhor rota após evolução - Distância = {melhor_distancia
      :.2f}")
8
  plt.xlabel("Coordenada X")
  plt.ylabel("Coordenada Y")
9
10 plt.tight_layout()
11
12
  # Mostra os gráficos gerados
13
  plt.show()
```

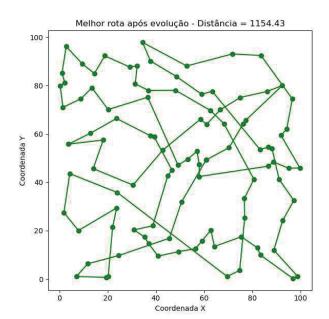


Figura 15.2 – EXERCÍCIO 1: ROTA APÓS EVOLUÇÃO, DISTÂNCIA = 1154,43

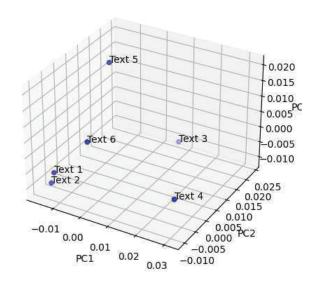
### 2) Compare a representação de dois modelos vetoriais

```
# Comandos para manipulação de pacotes
  !pip uninstall numpy -y
2
  !pip install numpy==1.25.2
3
  !pip uninstall gensim -y
   !pip install --no-binary :all: gensim
5
6
   import numpy as np
7
  import re
8
   from gensim.models import Word2Vec
  from sklearn.decomposition import PCA
10
  from scipy.spatial.distance import pdist, squareform
   from scipy.cluster.hierarchy import dendrogram, linkage
13
   import matplotlib.pyplot as plt
14
15
   # Textos curtos de exemplo
16
   textos = [
17
       "O gato está dormindo no sofá.",
                                              # Texto 1
       "O cachorro está dormindo no tapete.", # Texto 2 (similar to
18
19
       "A tecnologia esta mudando nossas vidas.", # Texto 3
20
       "Estamos vivendo na era da tecnologia.",  # Texto 4 (similar to
      Texto 3)
21
       "O dia está ensolarado e claro.",
                                                   # Texto 5
       "O sol está brilhando no céu azul."
                                                 # Texto 6 (similar to
22
      Texto 5)
23
24
   # Pre-processamento: impor pontuação e criar lista de tokens por
25
  sentencas_tokens = []
26
27 | for txt in textos:
```

```
28
       # Converter para minúsculas e remover pontuação
29
       texto_limpo = re.sub(r'[^\w\s]', '', txt.lower())
30
       tokens = texto_limpo.split()
31
       if tokens:
32
           sentencas_tokens.append(tokens)
33
34
   # Treinar modelo Word2Vec nos tokens
  model = Word2Vec(sentencas_tokens, vector_size=50, window=5,
      min count=1, workers=4)
36
37
  # Obter vetor de cada sentenca fazendo média dos vetores de palavras
  vetores_texto = []
38
39
   for tokens in sentencas_tokens:
40
       vetores_palavras = [model.wv[word] for word in tokens if word in
      model.wv.key_to_index]
41
       if len(vetores_palavras) > 0:
           # Média dos vetores de palavra para representar a sentença
42
43
           vetor_sentenca = np.mean(vetores_palavras, axis=0)
44
       else:
45
           # Caso sentença sem palavras conhecidas
46
           vetor_sentenca = np.zeros(model.vector_size)
47
       vetores_texto.append(vetor_sentenca)
48
   vetores_texto = np.array(vetores_texto)
49
50
   # Calcular vetor medio e centralizar os vetores de texto
51
  vetor_medio = np.mean(vetores_texto, axis=0)
52
  vetores_centralizados = vetores_texto - vetor_medio
53
54
  # Aplicar PCA para reduzir dimensão dos vetores
55 | pca_3d = PCA(n_components=3)
56 pca_2d = PCA(n_components=2)
57
   vetores_3d = pca_3d.fit_transform(vetores_centralizados)
58 | vetores_2d = pca_2d.fit_transform(vetores_centralizados)
59
60 | # Visualização dos vetores em 3D e 2D
  labels = [f"Text {i+1}" for i in range(len(textos))]
61
62
63
  # Plot 3D dos vetores originais (após PCA 3D)
   fig = plt.figure(figsize=(6,5))
65
  ax = fig.add_subplot(111, projection='3d')
66 ax.scatter(vetores_3d[:,0], vetores_3d[:,1], vetores_3d[:,2], c='blue
      ', marker='o')
67
  for i, label in enumerate(labels):
       ax.text(vetores_3d[i,0], vetores_3d[i,1], vetores_3d[i,2], label)
68
  ax.set_title('Vetores originais (projeção PCA 3D)')
70 | ax.set_xlabel('PC1')
71 | ax.set_ylabel('PC2')
72 ax.set_zlabel('PC3')
73
  plt.show()
```

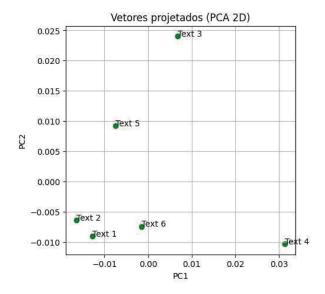
Figura 15.3 – EXERCÍCIO 2: PROJEÇÃO PCA 3D

## Vetores originais (projeção PCA 3D)



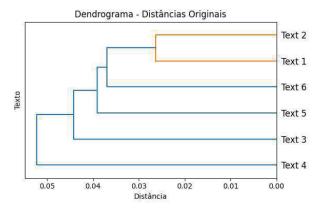
```
# Plot 2D dos vetores projetados (PCA 2D)
  plt.figure(figsize=(5,5))
  plt.scatter(vetores_2d[:,0], vetores_2d[:,1], c='green', marker='o')
3
4
  for i, label in enumerate(labels):
       plt.annotate(label, (vetores_2d[i,0], vetores_2d[i,1]))
5
6
  plt.title('Vetores projetados (PCA 2D)')
7
  plt.xlabel('PC1')
  plt.ylabel('PC2')
9
  plt.grid(True)
  plt.show()
10
```

Figura 15.4 – EXERCÍCIO 2: PROJEÇÃO PCA 2D



```
2 | dist_original = squareform(pdist(vetores_centralizados, metric='
      euclidean'))
  dist_projetado = squareform(pdist(vetores_2d, metric='euclidean'))
3
4
5
  # Gerar dendrogramas de agrupamento hierárquico
  # Dendrograma para distâncias originais
6
  Z_orig = linkage(vetores_centralizados, method='ward', metric='
7
      euclidean')
  plt.figure(figsize=(6,4))
8
  dendrogram(Z_orig, labels=labels, orientation='left')
9
10 | plt.title('Dendrograma - Distâncias Originais')
11 plt.xlabel('Distância')
12 plt.ylabel('Texto')
13 plt.tight_layout()
14 plt.show()
```

Figura 15.5 – EXERCÍCIO 2: DENDOGRAMA - DISTÂNCIAS ORIGINAIS



```
# Dendrograma para distâncias após projeção PCA 2D
Z_proj = linkage(vetores_2d, method='ward', metric='euclidean')
plt.figure(figsize=(6,4))
dendrogram(Z_proj, labels=labels, orientation='left')
plt.title('Dendrograma - Distâncias PCA 2D')
plt.xlabel('Distância')
plt.ylabel('Texto')
plt.tight_layout()
plt.show()
```

Figura 15.6 – EXERCÍCIO 2: DENDOGRAMA - DISTÂNCIAS PCA 2D

