

# UNIVERSIDADE DE SÃO PAULO

Instituto de Ciências Matemáticas e de Computação

Título

Ramon de Castro Ramos

Monografia - MBA em Ciência de Dados (CeMEAI)



**Ramon de Castro Ramos**

## **Título**

Monografia apresentada ao Centro de Ciências Matemáticas Aplicadas à Indústria do Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo - ICMC/USP, como parte dos requisitos para obtenção do título de Especialista em Ciências de Dados.

Área de concentração: Ciências de Dados

Orientador: Prof. Dr. Adriano Kamimura Suzuki

## **Versão original**

**São Carlos**  
**2025**

É possível elaborar a ficha catalográfica em LaTeX ou incluir a fornecida pela Biblioteca. Para tanto observe a programação contida nos arquivos USPSC-modelo.tex e fichacatalografica.tex e/ou gere o arquivo fichacatalografica.pdf.

A biblioteca da sua Unidade lhe fornecerá um arquivo PDF com a ficha catalográfica definitiva, que deverá ser salvo como fichacatalografica.pdf no diretório do seu projeto.

## **Ramon de Castro Ramos**

### **Title**

Monograph presented to the Centro de Ciências Matemáticas Aplicadas à Indústria do Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo - ICMC/USP, as part of the requirements for obtaining the title of Specialist in Data Science.

Concentration area: Data Science

Advisor: Prof. Dr. Adriano Kamimura Suzuki

### **Original version**

**São Carlos  
2025**

Folha de aprovação em conformidade  
com o padrão definido  
pela Unidade.

No presente modelo consta como  
folhadeaprovacao.pdf

*Dedico este trabalho aos meus pais,  
por todo o amor, apoio, incentivos e sacrifícios  
que me impulsionaram a trilhar o caminho que trilhei.*



## **AGRADECIMENTOS**

Primeira frase do agradecimento ....

Segunda frase ....

Outras frases ....

Última frase ....



*“Be yourself, everyone else is already taken.”*

*Oscar Wilde*



## **RESUMO**

RAMOS, R. C. **Título.** 2025. 108 p. Monografia (MBA em Ciências de Dados) - Centro de Ciências Matemáticas Aplicadas à Indústria, Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo, São Carlos, 2025.

xxx

**Palavras-chave:** x. x. x. x. x. x.



## **ABSTRACT**

RAMOS, R. C. **Title.** 2025. 108 p. Monograph (MBA in Data Sciences) - Centro de Ciências Matemáticas Aplicadas à Indústria, Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo, São Carlos, 2025.

xxx

**Keywords:** x. x. x. x. x. x.



## LISTA DE FIGURAS

Figura 1 – Exemplo de uma Regressão Linear simples com dados fictícios . . . . .	36
Figura 2 – Exemplo de uma Árvore de Decisão com o <i>dataset Iris</i> . . . . .	37
Figura 3 – Exemplo de uma <i>Random Forest</i> com o <i>dataset Iris</i> . . . . .	38
Figura 4 – Modelo CRISP-DM . . . . .	39
Figura 5 – Histograma das notas - Ciências Humanas . . . . .	52
Figura 6 – Histograma das notas - Ciências da Natureza . . . . .	53
Figura 7 – Histograma das notas - Linguagem e Código . . . . .	53
Figura 8 – Histograma das notas - Matemática . . . . .	54
Figura 9 – Histograma das notas - Redação . . . . .	54
Figura 10 – Boxplot das notas por edição - Ciências Humanas . . . . .	56
Figura 11 – Boxplot das notas - Ciências da Natureza . . . . .	57
Figura 12 – Boxplot das notas - Linguagem e Código . . . . .	57
Figura 13 – Boxplot das notas - Matemática . . . . .	58
Figura 14 – Boxplot das notas - Redação . . . . .	58
Figura 15 – Dez maiores <i>Permutation Importance</i> - Humanas . . . . .	63
Figura 16 – Dez maiores <i>Permutation Importance</i> - Natureza . . . . .	64
Figura 17 – Dez maiores <i>Permutation Importance</i> - Linguagem . . . . .	64
Figura 18 – Dez maiores <i>Permutation Importance</i> - Matemática . . . . .	65
Figura 19 – Dez maiores <i>Permutation Importance</i> - Redação . . . . .	65
Figura 20 – Erro do <i>Grid Search - XGBoost</i> . . . . .	69
Figura 21 – Erro do <i>Grid Search - LightGBM</i> . . . . .	70
Figura 22 – Erro do <i>Grid Search - Random Forest</i> . . . . .	70
Figura 23 – Erro MAPE - <i>XGBoost</i> . . . . .	72
Figura 24 – Erro MAPE - <i>LightGBM</i> . . . . .	73
Figura 25 – Erro MAPE - <i>Ensemble (XGBoost + LightGBM)</i> . . . . .	74
Figura 26 – Erro MAPE - <i>Ensemble (XGBoost + LightGBM + Random Forest)</i> . .	74
Figura 27 – Rank de Importância - Humanas . . . . .	78
Figura 28 – Rank de Importância - Natureza . . . . .	78
Figura 29 – Rank de Importância - Linguagem e Código . . . . .	79
Figura 30 – Rank de Importância - Matemática . . . . .	79
Figura 31 – Rank de Importância - Redação . . . . .	80
Figura 32 – Curva de Sensibilidade - Faixa Etária . . . . .	81
Figura 33 – Curva de Sensibilidade - Sexo . . . . .	82
Figura 34 – Curva de Sensibilidade - Cor/Raça . . . . .	82
Figura 35 – Curva de Sensibilidade - Escolaridade do Pai . . . . .	83
Figura 36 – Curva de Sensibilidade - Escolaridade da Mãe . . . . .	83

Figura 37 – Curva de Sensibilidade - Ocupação do Pai . . . . .	84
Figura 38 – Curva de Sensibilidade - Ocupação da Mãe . . . . .	84
Figura 39 – Curva de Sensibilidade - Renda Familiar . . . . .	85

## LISTA DE TABELAS

Tabela 1 – Variáveis socioeconômicas e suas referências . . . . .	45
Tabela 2 – Quantidade de observações e variáveis por edição do ENEM . . . . .	48
Tabela 3 – Percentual de valores nulos por variável . . . . .	49
Tabela 4 – Observações e variáveis por conjunto de dados . . . . .	51
Tabela 5 – Estatísticas descritivas por conjunto de dados . . . . .	52
Tabela 6 – Assimetria, Curtose e Notas zeradas . . . . .	55
Tabela 7 – Teste ANOVA das médias das notas por edição . . . . .	55
Tabela 8 – Quantidade e percentual de outliers nas notas . . . . .	59
Tabela 9 – Cinco maiores concentrações - Humanas . . . . .	59
Tabela 10 – Cinco maiores concentrações - Natureza . . . . .	60
Tabela 11 – Cinco maiores concentrações - Linguagem . . . . .	60
Tabela 12 – Cinco maiores concentrações - Matemática . . . . .	60
Tabela 13 – Cinco maiores concentrações - Redação . . . . .	60
Tabela 14 – Cinco maiores correlações Phik - Humanas . . . . .	61
Tabela 15 – Cinco maiores correlações Phik - Natureza . . . . .	62
Tabela 16 – Cinco maiores correlações Phik - Linguagem . . . . .	62
Tabela 17 – Cinco maiores correlações Phik - Matemática . . . . .	62
Tabela 18 – Cinco maiores correlações Phik - Redação . . . . .	62
Tabela 19 – Concentração cruzada - Redação . . . . .	66
Tabela 20 – <i>Grid Search - XGBoost</i> . . . . .	67
Tabela 21 – <i>Grid Search - LightGBM</i> . . . . .	68
Tabela 22 – <i>Grid Search - Random Forest</i> . . . . .	68
Tabela 23 – Hiperparâmetros Ajustados - <i>XGBoost</i> . . . . .	71
Tabela 24 – Hiperparâmetros Ajustados - <i>LightGBM</i> . . . . .	71
Tabela 25 – Hiperparâmetros Ajustados - <i>Random Forest</i> . . . . .	71
Tabela 26 – Erro MAPE - <i>Random Forest</i> . . . . .	73
Tabela 27 – Cinco melhores modelos - Humanas . . . . .	76
Tabela 28 – Cinco melhores modelos - Natureza . . . . .	76
Tabela 29 – Cinco melhores modelos - Linguagem . . . . .	76
Tabela 30 – Cinco melhores modelos - Matemática . . . . .	76
Tabela 31 – Cinco melhores modelos - Redação . . . . .	77



## **LISTA DE QUADROS**



## LISTA DE ABREVIATURAS E SIGLAS

AdaBoost	<i>Adaptive Boosting</i>
COVID-19	<i>Coronavirus Disease 2019</i> - Doença do Coronavírus 2019
CPU	<i>Central Processing Unit</i> - Unidade Central de Processamento
CRISP-DM	<i>Cross-Industry Standard Process for Data Mining</i> - Processo de Mineração de Dados Padrão entre Indústrias
CSV	<i>Comma-Separated Values</i> - Valores Separados por Vírgula
ENEM	Exame Nacional do Ensino Médio
Fies	Fundo de Financiamento Estudantil
GPU	<i>Graphics Processing Unit</i> - Unidade de Processamento Gráfico
INEP	Instituto Nacional de Estudos e Pesquisas Educacionais Anísio Teixeira
IQR	<i>Interquartile Range</i> - Intervalo Interquartil
ML	<i>Machine Learning</i> - Aprendizado de Máquina
ProUni	Programa Universidade Para Todos
RAM	<i>Random Access Memory</i> - Memória de Acesso Aleatório
RF	<i>Random Forest</i> - Floresta Aleatória
SISU	Sistema de Seleção Unificada
UFABC	Universidade Federal do ABC
XGBoost	<i>Extreme Gradient Boosting</i>



## **LISTA DE SÍMBOLOS**

$\sigma$  Letra Grega sigma minúscula; desvio padrão.

$\mu$  Letra Grega mu minúscula; média.

$\leq$  Menor ou igual

$Md$  Mediana



## SUMÁRIO

<b>1</b>	<b>INTRODUÇÃO . . . . .</b>	<b>29</b>
<b>2</b>	<b>FUNDAMENTAÇÃO TEÓRICA . . . . .</b>	<b>31</b>
<b>2.1</b>	<b>O ENEM no Cenário Educacional Brasileiro . . . . .</b>	<b>31</b>
<b>2.2</b>	<b>Teorias sobre Desigualdades Educacionais: O Capital Cultural de Bourdieu . . . . .</b>	<b>32</b>
<b>2.3</b>	<b>Fatores Socioeconômicos e Desempenho no ENEM . . . . .</b>	<b>32</b>
<b>2.4</b>	<b>Características escolares e o “Efeito Escola” . . . . .</b>	<b>33</b>
<b>2.5</b>	<b>Disparidades Regionais e a Participação no ENEM . . . . .</b>	<b>34</b>
<b>2.6</b>	<b>Aplicações de Ciência de Dados na Análise do ENEM e resultados obtidos . . . . .</b>	<b>34</b>
<b>2.7</b>	<b>Métodos de <i>Machine Learning</i> . . . . .</b>	<b>35</b>
2.7.1	Regressão Linear . . . . .	35
2.7.2	Árvore de Decisão . . . . .	36
2.7.3	<i>Random Forest</i> . . . . .	37
2.7.4	<i>Boosting</i> . . . . .	38
<b>3</b>	<b>METODOLOGIA . . . . .</b>	<b>39</b>
<b>3.1</b>	<b>Entendimento de Negócio . . . . .</b>	<b>39</b>
<b>3.2</b>	<b>Entendimento dos dados . . . . .</b>	<b>40</b>
<b>3.3</b>	<b>Preparação dos dados . . . . .</b>	<b>40</b>
<b>3.4</b>	<b>Modelagem . . . . .</b>	<b>41</b>
3.4.1	Análise Exploratória dos Dados . . . . .	41
<b>3.5</b>	<b>Treinamento dos Modelos . . . . .</b>	<b>42</b>
<b>3.6</b>	<b>Avaliação dos Modelos . . . . .</b>	<b>43</b>
<b>3.7</b>	<b>Influência das Variáveis Preditoras . . . . .</b>	<b>44</b>
<b>4</b>	<b>RESULTADOS . . . . .</b>	<b>45</b>
<b>4.1</b>	<b>Entendimento de Negócio . . . . .</b>	<b>45</b>
<b>4.2</b>	<b>Entendimento dos dados . . . . .</b>	<b>46</b>
4.2.1	Escolha e Coleta dos Dados . . . . .	46
4.2.2	Compreensão Inicial dos Dados . . . . .	46
4.2.2.1	Edição de 2024 do ENEM e LGPD . . . . .	46
4.2.3	Análise dos Dicionários de Dados . . . . .	47
4.2.4	Definição da Variável Resposta . . . . .	47
<b>4.3</b>	<b>Preparação dos dados . . . . .</b>	<b>47</b>

4.3.1	Preparação do Ambiente Tecnológico e Analítico . . . . .	47
4.3.2	Leitura dos Dados . . . . .	48
4.3.3	Integração dos Dados . . . . .	48
4.3.4	Tratamento de Valores Nulos . . . . .	49
4.3.5	Separação dos Conjuntos de Dados por Variável Resposta . . . . .	51
<b>4.4</b>	<b>Modelagem . . . . .</b>	<b>51</b>
4.4.1	Análise Exploratória dos Dados - Variáveis Resposta . . . . .	51
4.4.1.1	Distribuições . . . . .	51
4.4.1.2	Teste de Hipótese . . . . .	55
4.4.1.3	Análise de Outliers . . . . .	56
4.4.2	Análise Exploratória - Variáveis Preditoras . . . . .	59
4.4.2.1	Concentração . . . . .	59
4.4.2.2	Correlação Phik . . . . .	61
4.4.2.3	<i>Permutation Importance</i> . . . . .	63
4.4.2.4	Seleção de Variáveis . . . . .	65
<b>4.5</b>	<b>Treinamento dos Modelos . . . . .</b>	<b>67</b>
4.5.1	Ajuste dos Hiperparâmetros . . . . .	67
4.5.2	Treinamento final dos modelos . . . . .	72
4.5.3	Construção dos modelos de <i>ensemble</i> . . . . .	73
4.5.4	Avaliação dos modelos . . . . .	75
<b>4.6</b>	<b>Influência das Variáveis Preditoras . . . . .</b>	<b>77</b>
4.6.1	Importância . . . . .	77
4.6.2	Sensibilidade das Variáveis Respostas . . . . .	80
<b>5</b>	<b>CONCLUSÃO . . . . .</b>	<b>87</b>
	<b>REFERÊNCIAS . . . . .</b>	<b>89</b>
	<b>APÊNDICES</b>	<b>93</b>
	<b>APÊNDICE A – DICIONÁRIO DE DADOS DOS MICRODADOS DO ENEM . . . . .</b>	<b>95</b>
	<b>APÊNDICE B – DICIONÁRIO DE DADOS DO CENSO ESCOLAR</b>	<b>97</b>
	<b>APÊNDICE C – CONFIGURAÇÃO DO AMBIENTE VIRTUAL . . . . .</b>	<b>99</b>

## 1 INTRODUÇÃO

O Exame Nacional do Ensino Médio (ENEM) consolidou-se, na última década, como a principal avaliação educacional do Ensino Médio no Brasil, transcendendo seu papel inicial de termômetro da qualidade da educação básica para se tornar a porta de entrada para o ensino superior em instituições públicas e privadas, através de programas como o Sistema de Seleção Unificada (SISU), o Programa Universidade Para Todos (ProUni) e o Fundo de Financiamento Estudantil (Fies). Sua relevância reside na capacidade de fornecer um panorama detalhado do desempenho dos estudantes, bem como de aspectos socioeconômicos e contextuais que permeiam o ambiente escolar e familiar dos participantes.

Apesar dos esforços contínuos para aprimorar a qualidade da educação no Brasil, persistem desafios significativos, evidenciados pelas variações no desempenho dos estudantes em avaliações de larga escala como o ENEM. A literatura acadêmica aponta para a influência de múltiplos fatores nesse desempenho, que vão desde as condições socioeconômicas das famílias até as características estruturais e pedagógicas das escolas, além das peculiaridades regionais (1). A análise estatística de microdados do ENEM entre 2021 e 2023, por exemplo, revela desigualdades estruturais marcantes entre estudantes de escolas públicas e privadas (2). A persistência dessas disparidades indica que as desigualdades educacionais no Brasil não são meramente aleatórias, mas profundamente associadas às desigualdades sociais (3).

A análise aprofundada dos microdados do ENEM, portanto, constitui uma oportunidade ímpar para desvendar a complexa interação entre os fatores socioeconômicos, as características do ambiente escolar e as peculiaridades regionais que moldam o desempenho dos estudantes. Isso permite ir além da simples constatação das disparidades, oferecendo um panorama mais claro de como um instrumento concebido para democratizar o acesso ao ensino superior pode, na prática, atuar como um espelho das desigualdades sociais estruturais e, em certos contextos, até mesmo contribuir para a sua perpetuação, um fenômeno consistentemente observado em análises de dados históricos (2). A compreensão desses mecanismos é vital para a formulação de políticas públicas que não apenas mitiguem as lacunas, mas que atuem nas causas-raiz das iniquidades educacionais.

Nesse contexto, este Trabalho de Conclusão de Curso propõe investigar e quantificar a influência dos principais fatores socioeconômicos no desempenho dos estudantes no ENEM. A pergunta central que guia esta pesquisa é: “Quais são os principais fatores socioeconômicos que influenciam o desempenho dos estudantes no ENEM e qual a magnitude da influência de cada um desses conjuntos de fatores nas notas dos participantes?”. O objetivo geral é utilizar os microdados do exame para fornecer *insights* robustos sobre a qualidade da educação básica no Brasil, contribuindo para a identificação de áreas que necessitam de

maior atenção e investimento. A quantificação da influência dos fatores, por meio de modelos preditivos e análise de importância de variáveis (2), é um diferencial crucial. Não se trata apenas de identificar a existência de correlações, mas de medir o grau de impacto, o que é fundamental para a formulação de políticas públicas eficazes e direcionadas.

Para tanto, buscam-se os seguintes objetivos específicos: i) Coletar, pré-processar e realizar uma análise exploratória dos microdados do ENEM (4) selecionando as variáveis relevantes; ii) Identificar padrões, tendências e correlações entre as variáveis selecionadas e o desempenho dos estudantes; iii) Aplicar técnicas de Ciência de Dados para construir modelos preditivos e determinar a importância relativa de cada grupo de fatores; e iv) Discutir os resultados obtidos, correlacionando-os com a literatura existente e extraiendo dados práticos.

A relevância desta pesquisa reside na sua capacidade de oferecer uma análise quantitativa detalhada das correlações entre múltiplos fatores e o desempenho educacional, utilizando uma vasta base de dados. Os dados gerados podem servir como subsídio para educadores, formuladores de políticas públicas e pesquisadores, auxiliando na compreensão das raízes das desigualdades educacionais e na elaboração de estratégias direcionadas para a melhoria do ensino médio no país. A pesquisa não se limita a um exercício acadêmico; ela tem um potencial transformador social ao fornecer dados concretos para subsidiar políticas públicas mais justas e fortalecer a rede pública de ensino (2).

Os próximos capítulos irão apresentar a metodologia adotada neste trabalho, os resultados obtidos e a discussão desses resultados, culminando nas conclusões e recomendações para futuros trabalhos.

## 2 FUNDAMENTAÇÃO TEÓRICA

Este capítulo estabelece o contexto teórico e empírico para o estudo, fundamentando a análise no conhecimento acadêmico existente.

### 2.1 O ENEM no Cenário Educacional Brasileiro

O Exame Nacional do Ensino Médio (ENEM) teve sua primeira edição em 1998, contando com a participação de aproximadamente 115 mil participantes. Na época, suas notas só eram utilizadas por 2 instituições de ensino superior, número que salta para 93 instituições no ano seguinte. A importância do ENEM cresce com o passar dos anos, alcançando a marca de mais de 1 milhão de participantes na sua quarta edição e tornando-se uma das principais formas de acesso ao ensino superior, com a criação do Programa Universidade Para Todos (ProUni) em 2005 (5).

Em 2009, com a criação do Sistema de Seleção Unificada (SISU), o ENEM foi reformulado e assume o formato que tem hoje: 180 questões objetivas divididas em 4 áreas do conhecimento e uma redação. No ano seguinte, os resultados do ENEM passaram a ser adotados pelo Fundo de Financiamento Estudantil (Fies) e em 2013, quase todas as instituições federais adotam o ENEM como critério de seleção. Duas universidades portuguesas, a Universidade de Coimbra e Universidade de Algrave, passar a usar o ENEM como critério de seleção em 2014, número que chega a 35 instituições portuguesas em 2018 (5).

É evidente que o ENEM deixa de ser apenas uma ferramenta de avaliação e transforma-se em um instrumento multifacetado que desempenha um papel central na trajetória educacional dos jovens brasileiros. Além de aferir o desempenho dos estudantes ao final do ensino médio, o ENEM serve como a principal porta de acesso ao ensino superior, sendo a base para o SISU, o ProUni e o Fies (6). Essa centralidade significa que qualquer fator que influencie o desempenho no exame tem um impacto direto e significativo nas oportunidades de acesso ao ensino superior e, consequentemente, na mobilidade social dos indivíduos.

Os microdados do ENEM, disponibilizados anualmente pelo Instituto Nacional de Estudos e Pesquisas Educacionais Anísio Teixeira (INEP), representam uma fonte de informação rica e valiosa para pesquisas educacionais (2). Esses dados detalhados permitem uma compreensão aprofundada dos padrões de desempenho, das características socioeconômicas dos participantes e dos contextos escolares, possibilitando análises complexas sobre as desigualdades educacionais no país.

## 2.2 Teorias sobre Desigualdades Educacionais: O Capital Cultural de Bourdieu

Para compreender a reprodução das desigualdades sociais no sistema educacional, a teoria do capital cultural de Pierre Bourdieu oferece um arcabouço teórico fundamental. Este argumenta que o sucesso escolar não depende apenas do mérito individual ou da capacidade cognitiva, mas também da posse de diferentes formas de capital: o econômico (posses que o indivíduo tem), o social (relacionamentos que podem ser benéficos aos indivíduo), o simbólico (prestígio/honra) e o cultural (conhecimentos reconhecidos por diplomas e títulos) (7).

O capital cultural ainda se divide em três estados: (i) o capital cultural incorporado, composto por elementos pessoais como gostos (musicais, artísticos etc.), domínio de línguas; (ii) o capital cultural objetivado, composto por posses de livros e obras de arte ou acesso a museus, cinema etc.; (ii) o capital cultural instucionalizado, caracterizado por diplomas e títulos de conhecimento (7).

A acumulação de capital cultural é o que influenciará o desempenho escolar do indivíduo e futuramente seu posicionamento no mercado de trabalho. Se os dados do ENEM confirmarem a forte influência de variáveis socioeconômicas e de escolaridade parental, isso reforçará a tese da reprodução escolar das desigualdades, sugerindo que o sistema educacional, em vez de ser um equalizador, pode perpetuar as hierarquias sociais. Isso se manifesta, por exemplo, na forma como a escolaridade da mãe e a renda familiar são fatores relevantes para o desempenho e a dispersão das notas do ENEM (1).

Oliveira e Cruz (2014) argumentam que a escola ao reconhecer os alunos mais inteligentes ou aplicados, na verdade estão selecionando os alunos com o capital cultural mais diverso e amplo, o que propaga a desigualdade social ao criar os “mitos de aluno inteligente-brilhante / aluno fracassado-invisível”, fazendo com que “o próprio oprimido passa a acreditar que não é capaz de ter sucesso por características pessoais e não do sistema.”

## 2.3 Fatores Socioeconômicos e Desempenho no ENEM

A literatura é vasta ao associar variáveis socioeconômicas ao desempenho em avaliações de larga escala e o ENEM não é exceção. As persistentes e quantificáveis desigualdades de desempenho ligadas a fatores socioeconômicos (1) indicam que o acesso a “experiências educacionais muito mais ricas” (8) fora do ambiente escolar formal é um preditor poderoso do sucesso no ENEM. Isso sugere que a escola, por si só, pode não ser capaz de compensar totalmente essas desvantagens de origem e que o campo educacional não é nivelado desde o início.

Estudos sobre o ENEM consistentemente apontam o impacto de diversos fatores:

- **Renda Familiar:** Uma correlação positiva e significativa é observada entre a renda familiar e as notas do ENEM (1). Análises indicam que a diferença na nota de redação pode ser de até 40% entre os grupos de menor e maior renda (8).
- **Raça / Cor:** O desempenho de alunos brancos consistentemente supera o de outros grupos raciais, mesmo quando outras variáveis são controladas (1). Em média, o desempenho de alunos brancos superou o dos demais em menos de 10 pontos nas quatro provas em 2018, controlando outras variáveis (9).
- **Escolaridade dos Pais / Nível Instrucional da Mãe:** Este é um fator relevante para o desempenho e a dispersão das notas dos estudantes (1). Mães com escolaridade a partir do ensino médio e famílias de renda alta têm um impacto positivo no desempenho (10).
- **Sexo:** Diferenças de desempenho por sexo são notadas, especialmente na prova de Matemática, com vantagem para os homens (até 36 pontos a mais) (10).
- **Idade / Atraso Escolar:** O atraso escolar associa-se negativamente ao desempenho. Alunos com pelo menos um ano de atraso escolar tiveram, em média, de 16,7 a 29,0 pontos a menos nas provas (9).

## 2.4 Características escolares e o “Efeito Escola”

As características das escolas também exercem influência no desempenho dos estudantes e o conceito de “efeito escola” busca mensurar a contribuição da instituição de ensino para o desempenho do aluno, além dos fatores individuais e familiares (10). Achados relevantes incluem:

- **Dependência Administrativa (Pública vs. Privada):** Alunos de escolas privadas consistentemente superam os de escolas públicas (10). Em Matemática, a diferença pode ser de aproximadamente 83,9 pontos entre alunos de escolas privadas e estaduais (9). Um estudo da UFABC, por exemplo, mostrou que em Matemática, apenas 2,9% dos estudantes da rede pública atingiram 720 pontos, contra 20% da rede privada (2).
- **Atributos Escolares:** Fatores como complexidade de gestão, média de horas-aula, número de alunos por turma, qualidade dos professores (esforço e adequação docente) e o nível socioeconômico médio da escola são importantes (10). O nível socioeconômico médio da escola e a regularidade docente destacam-se como os mais significativos, aumentando a nota em 22,7 pontos para cada nível socioeconômico e em 14,6 para cada nível de regularidade docente em escolas privadas (10).

Embora o “Efeito Escola” seja um fator, a literatura sugere que uma grande parte da explicação das notas do ENEM reside em fatores externos ao controle escolar (10). Isso significa que, embora a qualidade da escola seja importante, as disparidades socioeconômicas dos alunos e o ambiente familiar podem ter um peso ainda maior. Isso desafia a ideia de que a escola, por si só, pode reverter completamente as desigualdades de origem, apontando para a necessidade de políticas holísticas que abordem tanto os fatores intra-escolares quanto os extra-escolares.

## 2.5 Disparidades Regionais e a Participação no ENEM

O desempenho no ENEM também exibe variações significativas entre diferentes regiões e unidades da federação (1). As disparidades regionais não são apenas geográficas, mas refletem a heterogeneidade socioeconômica e a capacidade de resposta dos sistemas educacionais locais a crises, como a pandemia de COVID-19 (11).

O período pós-pandemia, em particular, evidenciou um agravamento das desigualdades regionais na participação e no desempenho, com quedas não homogêneas nas taxas de inscrição (12). A maior queda proporcional na taxa de inscrição ocorreu na região Sudeste, que de um pico de 63% em 2016, chegou a apenas 26% em 2021, tornando-se a região com o menor indicador naquele ano (11).

## 2.6 Aplicações de Ciência de Dados na Análise do ENEM e resultados obtidos

A aplicação de técnicas de Ciência de Dados e *Machine Learning* na análise dos microdados do ENEM tem se mostrado uma abordagem poderosa para aprofundar a compreensão dos fatores que influenciam o desempenho (1). Estudos têm utilizado regressão linear, árvores de decisão, *Random Forest*, *Boosting* entre outras técnicas para predição de notas e identificação de fatores relevantes (1, 3, 9, 10, 13–15).

Em seu trabalho, Melo *et al.* (1) utilizaram o método de regressão linear múltipla para modelar a média da prova objetiva, média da redação e as respectivas variâncias. Seus resultados indicam fortemente que o nível de escolaridade e profissionalização da mãe, a raça do estudante e a renda média da família são relevantes para o desempenho na prova objetiva. Ao adicionar uma componente espacial, os modelos apresentaram uma melhora, indicando que fatores regionais também influenciam o desempenho do estudante.

Moraes *et al.* (10) também aplicaram o método de regressão linear múltipla para analisar o efeito escola no desempenho em matemática, considerando variáveis como a quantidade média de alunos por turma, a média de horas-aula por dia e mais algumas variáveis que caracterizam a escola. Em sua análise exploratória, os autores identificaram as diferenças e similares entre as escolas públicas e privadas, a exemplo do nível socioeconômico médio dos alunos da escola, onde “87% das escolas privadas estão nos níveis 5 e 6, enquanto

90% das escolas públicas possui nível socioeconômico entre os níveis 3 ou 4. Assim, as escolas públicas lidam [...] com alunos com níveis socioeconômico menores.”

O nível socioeconômico médio dos alunos da escola chega “a aumentar a nota em 22,7 pontos para cada nível socioeconômico [...] nas escolas privadas e 12,3 pontos [...] nas escolas públicas.” Essa variável foi construída pelos autores e separada em 6 grupos, onde o grupo 6 reúne as escolas com os alunos de maior nível socioeconômico e o grupo 1 reúne as escolas com os alunos de menor nível socioeconômico.

Os Trabalhos de Conclusão de Curso de Amanda Ferraz (14) e Mayra Romero (13), para este mesmo MBA, aplicaram técnicas mais robustas. Ferraz utilizou *Random Forest* e *Boosting* para prever a aprovação de participantes do ENEM no SISU para o curso de Medicina, obtendo resultados satisfatórios com Coeficiente de Correlação de Matthews superior a 0,9. Já Romero desenvolveu e comparou modelos de classificação, incluindo o *Random Forest*, para identificar características socioeconômicas que indicam maior chance de o candidato atingir uma pontuação média acima de 500 pontos no ENEM. Ela concluiu que o *Random Forest* teve o melhor desempenho e que a renda familiar e o número de computadores são informações que impactam a previsibilidade do modelo.

## 2.7 Métodos de *Machine Learning*

Essa seção pretende apresentar, de forma não exaustiva, alguns dos métodos de *Machine Learning* utilizados em trabalhos anteriores relacionados ao tema deste trabalho. Para isso, foram usadas as referências (16–19) como base para a descrição dos métodos.

### 2.7.1 Regressão Linear

A Regressão Linear é um dos pilares do *Machine Learning*, sendo um método fundamental para a modelagem preditiva. Trata-se de um método paramétrico de aprendizado supervisionado que busca definir um modelo para uma relação linear entre a variável resposta e uma ou mais variáveis preditoras, tendo como objetivo central encontrar a melhor reta (ou hiperplano), em termos de erro na previsão, que descreva essa relação.

A implementação mais básica é expressa pela equação

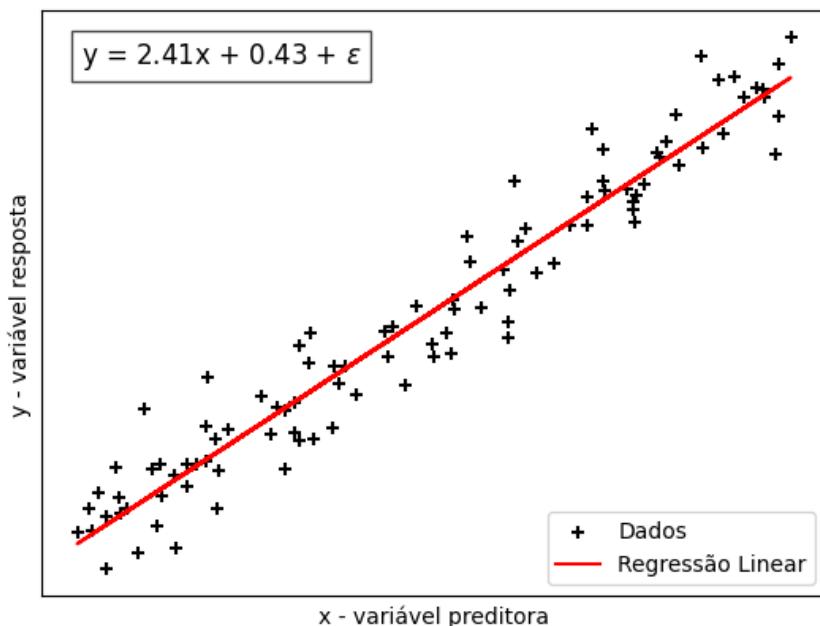
$$Y = \beta_0 + \beta_1 \times X + \epsilon \quad (2.1)$$

onde  $Y$  denota a variável resposta,  $X$  a variável preditora,  $\beta_0$  o intercepto (o valor de  $Y$  quando  $X = 0$ ),  $\beta_1$  o coeficiente angular (indicando o impacto de  $X$  sobre  $Y$ ) e  $\epsilon$  o termo de erro. Em uma regressão múltipla, diversas variáveis independentes são consideradas, cada uma com o seu  $\beta_i$  correspondente.

Por trás da regressão linear, há algumas premissas adotadas, como a linearidade da relação entre  $X$  e  $Y$ , a independência dos erros, a homocedasticidade e a normalidade

dos resíduos. Essas premissas podem ser interpretadas como desvantagens do modelo de regressão linear, por restringir ou até mesmo a inviabilizar a sua aplicação. Já a fácil interpretação, simplicidade e eficiência computacional são algumas das vantagens desse método, que também é muito utilizado como *benchmark* de métodos mais complexos.

Figura 1 – Exemplo de uma Regressão Linear simples com dados fictícios



Fonte: elaborado pelo autor.

### 2.7.2 Árvore de Decisão

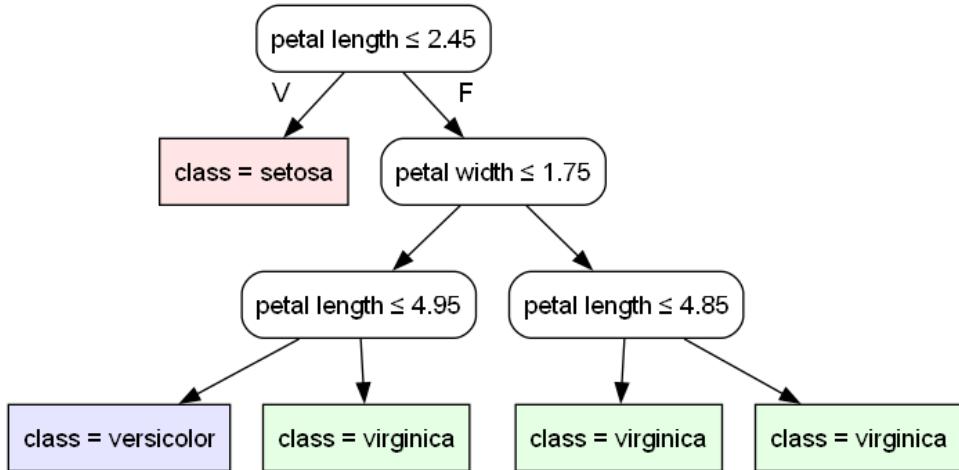
A Árvore de Decisão é um método paramétrico de aprendizado supervisionado que utiliza uma abordagem intuitiva de separação dos dados em grupos semelhantes, através de regras hierárquicas simples e de forma recursiva. Pode ser utilizado para resolver problemas de regressão, com a média da variável resposta em cada folha, ou de classificação, com a classe mais frequente em cada folha.

O processo de divisão segue uma lógica de “se-então”: se o dado de entrada tem o valor de uma variável preditora menor ou igual a um limite, então este segue pelo caminho a esquerda; se não, então este segue pelo caminho a direita. É dessa lógica que surge a analogia com árvore, já que as regras usadas para definir o modelo, podem ser representadas em um gráfico de árvore binária. A seleção das melhores divisões é baseada, para os problemas de classificação, em alguma medida de impureza, como a Entropia ou o Índice de Gini. Já para os problemas de regressão, as divisões são baseadas na redução de alguma medida de erro, como o erro quadrático médio (*Mean Squared Error* - MSE).

Assim como a Regressão Linear, a Árvore de Decisão é um modelo de fácil interpretação, já que as regras de decisão são explícitas e podem ser visualizadas graficamente. É capaz de lidar com variáveis categóricas e contínuas, o que a torna versátil, não requer

normalização dos dados e é robusta a outliers. No entanto, ela é propensa ao *overfitting*, se não aplicadas técnicas de poda, e são instáveis, já que pequenas variações nos dados podem levar a grandes mudanças na estrutura da árvore.

Figura 2 – Exemplo de uma Árvore de Decisão com o dataset *Iris*



Fonte: elaborado pelo autor.

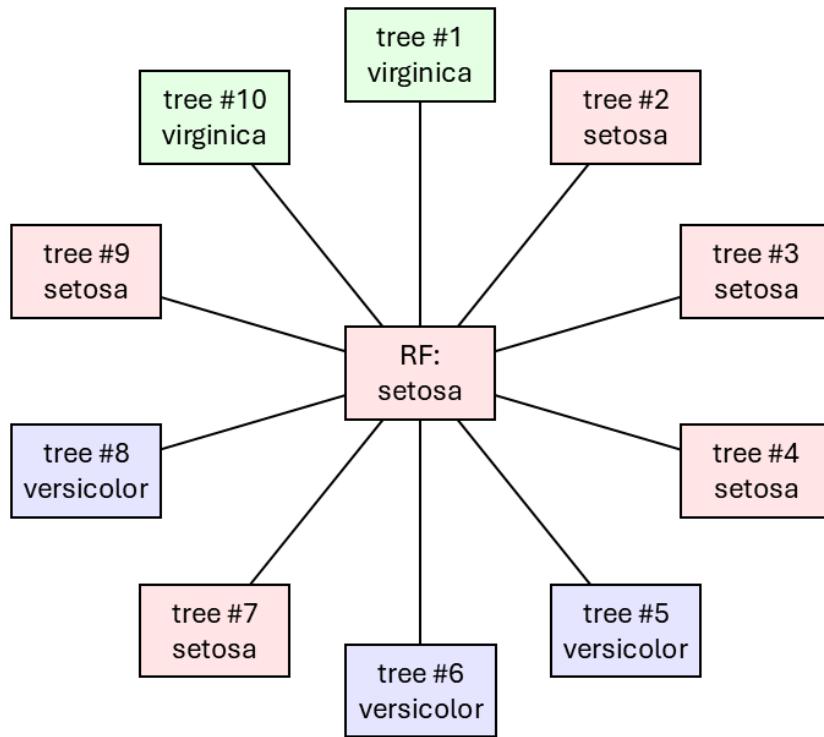
### 2.7.3 Random Forest

O *Random Forest* é um método derivado da Árvore de Decisão, sendo um dos algoritmos mais populares e eficazes em *Machine Learning*. Ele adota uma abordagem de *ensemble*, ou seja, combina múltiplos modelos para melhorar a precisão e a robustez das previsões. A ideia central é criar uma “floresta” de Árvores de Decisão, onde a decisão final é feita pela média/mediana das previsões para um problema de regressão ou pela classe mais frequente entre todas as árvores no caso de um problema de classificação.

O seu processo de construção envolve duas etapas principais: (i) a amostragem aleatória dos dados, onde cada árvore é treinada em um subconjunto diferente dos dados originais, e (ii) a seleção aleatória de variáveis em cada divisão, o que reduz a correlação entre as árvores e melhora a generalização do modelo. Essa aleatoriedade é crucial para evitar o *overfitting* e aumentar a diversidade entre as árvores.

O *Random Forest* é conhecido por sua alta precisão, capacidade de lidar com grandes conjuntos de dados e variáveis de diferentes tipos, resistência a *outliers* e facilidade de interpretação através da análise da importância das variáveis. No entanto, ele pode ser computacionalmente intensivo e menos interpretável do que uma única árvore de decisão, já que a combinação de múltiplas árvores torna mais difícil entender as regras subjacentes.

Figura 3 – Exemplo de uma *Random Forest* com o dataset *Iris*



Fonte: elaborado pelo autor.

#### 2.7.4 Boosting

O *Boosting* é uma técnica de *ensemble*, combinando múltiplos modelos fracos para criar um modelo forte. A ideia central é treinar sequencialmente uma série de modelos, onde cada novo modelo foca em corrigir os erros cometidos pelos modelos anteriores. Alguns algoritmos populares de *Boosting* incluem o *AdaBoost*, *Gradient Boosting* e *XGBoost*.

O *AdaBoost* (*Adaptive Boosting*) foi um dos primeiros algoritmos de *Boosting* e funciona aumentando o peso dos dados de treinamento que foram classificados incorretamente pelos modelos anteriores. Ao final, as previsões de todos os modelos são combinadas, ponderadas pela precisão de cada modelo.

O *Gradient Boosting* usa uma abordagem de otimização, onde cada novo modelo é treinado especificamente nos resíduos do modelo anterior, buscando minimizá-los. Os novos aprendizes são adicionados de forma iterativa e geralmente são árvores de decisão de pequeno porte.

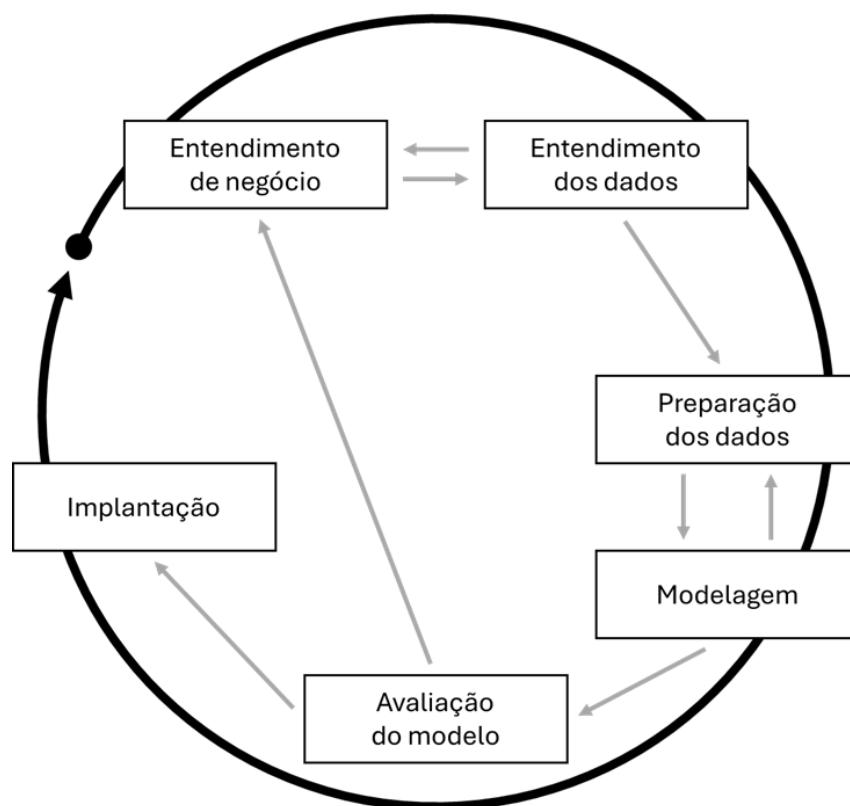
O *XGBoost* (*Extreme Gradient Boosting*) é uma implementação otimizada do *Gradient Boosting*, que oferece melhorias significativas em termos de velocidade e desempenho, implementando técnicas de regularização (L1 e L2), tratamento de valores ausentes, paralelização e outras otimizações.

### 3 METODOLOGIA

Este capítulo detalha a metodologia de trabalho utilizada, apresentando o delineamento da pesquisa, da coleta e processamento dos dados e as técnicas analíticas empregadas para responder às perguntas de pesquisa.

A estrutura metodológica adotada será baseada no modelo CRISP-DM (*Cross-Industry Standard Process for Data Mining*) (20), contendo as etapas de (i) entendimento de negócio, (ii) entendimento dos dados, (iii) preparação dos dados, (iv) modelagem, (v) avaliação e (vi) implantação.

Figura 4 – Modelo CRISP-DM



Fonte: modificado de Chapman *et al.* (20).

#### 3.1 Entendimento de Negócio

A etapa de entendimento de negócio envolve a definição clara dos objetivos do projeto, a compreensão do contexto em que a pesquisa está inserida, a identificação das partes interessadas, a formulação das perguntas de pesquisa que guiarão a análise dos dados e os resultados que espera-se alcançar.

Como o foco principal deste trabalho será a formulação de hipóteses relacionadas aos fatores que influenciam o desempenho dos estudantes no ENEM, através da análise

dos dados de performance dos participantes e suas características socioeconômicas, foi necessário formular perguntas de pesquisa específicas que possam ser respondidas através da análise dos dados disponíveis.

### **3.2 Entendimento dos dados**

Com as perguntas de pesquisa definidas, a próxima etapa foi encontrar dados que fossem adequados para responder a essas perguntas e estabelecer uma forma consistente de coleta e armazenamento desses dados para em seguida realizar uma compreensão de sua estrutura e significado.

Após essa etapa, foi possível identificar quais os arquivos são relevantes para a análise e quais variáveis dentro desses arquivos serão utilizadas como variáveis preditoras e como variáveis resposta.

Há depender dos tipos de dados a serem utilizados, é necessário submeter o projeto a um comitê de ética em pesquisa para aprovação, garantindo que todos os aspectos éticos relacionados ao uso dos dados sejam devidamente considerados.

### **3.3 Preparação dos dados**

Com os arquivos relevantes selecionados, passa-se para a etapa de preparação dos dados, que envolve a leitura, limpeza, transformação e integração dos dados para torná-los adequados para a modelagem. Essa etapa é crucial, pois a qualidade dos dados impacta diretamente na eficácia dos modelos preditivos que serão construídos posteriormente.

Para a execução dessa e das etapas posteriores, é necessário preparar um ambiente tecnológico e analítico adequado que permita a manipulação eficiente dos dados e a construção dos modelos preditivos.

Após a leitura dos dados, estes foram integrados em único conjunto de dados para facilitar a análise e a modelagem. Foram realizados os ajustes necessários no esquema dos dados para garantir a consistência e a integridade das informações.

Em seguida, as variáveis foram renomeadas para nomes mais intuitivos e de fácil compreensão e foi analisado se necessário transformar os valores das variáveis originais em valores mais comprehensíveis, por exemplo, transformar os códigos numéricos de variáveis categóricas em rótulos textuais.

Posteriormente, foi realizada uma análise para identificar e tratar valores nulos, removendo ou imputando valores conforme apropriado. Após o tratamento dos valores nulos, os dados foram separados em diferentes conjuntos de dados, cada um correspondente a uma variável resposta específica, garantindo que cada conjunto contenha apenas as observações relevantes para a análise daquela variável e sem valores nulos.

### 3.4 Modelagem

Com os dados preparados, a etapa de modelagem envolve a análise exploratória dos dados para entender suas características, avaliação de correlações entre as variáveis preditoras, seleção e treinamento dos modelos preditivos e otimização dos hiperparâmetros. Além disso, é importante medir a magnitude da influência das variáveis preditoras nas variáveis resposta, utilizando técnicas de interpretação de modelos e análise de sensibilidade.

#### 3.4.1 Análise Exploratória dos Dados

A análise exploratória foi iniciada com o entendimento da variáveis respostas, buscando entender a distribuição das notas do ENEM. Em seguida, foi feito um teste de hipótese para se avaliar se as médias das notas variam significativamente entre as edições selecionadas para esse trabalho.

Esse análise é necessária uma vez que as edições do ENEM podem apresentar variações no nível de dificuldade das provas, impactando as notas dos estudantes. Para isso, será utilizada a Estatística F de ANOVA (*Analysis of Variance*) (21) para comparar as médias das notas entre as diferentes edições do ENEM e o cálculo do tamanho do efeito para quantificar a magnitude das diferenças encontradas, utilizando a métrica SMD (*Standardized Mean Difference*) utilizando os intervalos definidos por Jacob Cohen (22).

A próxima etapa da análise exploratória foi a identificação de outliers nas variáveis respostas, utilizando o critério de 1.5 vezes o intervalo interquartil (IQR) (21). A identificação dos outliers é importante para entender a distribuição das notas e avaliar se esses valores extremos podem influenciar os resultados dos modelos preditivos.

Finalizadas as análises das variáveis respostas, foram realizadas três análises sobre as variáveis preditoras: (i) análise de concentração de categorias, (ii) análise de correlação entre as variáveis preditoras e (iii) cálculo do *Permutation Importance* para cada variável preditora.

Utilizando essas três informações, foi feita uma análise qualitativa para identificar possíveis variáveis preditoras a serem removidas do conjunto de dados, seja por apresentarem baixa correlação com as variáveis respostas, por apresentarem alta correlação com outras variáveis preditoras (multicolinearidade) ou por apresentarem baixa importância na predição das notas do ENEM.

Para a correlação, foi utilizada a biblioteca `phik` (23) (24), que permite calcular a correlação entre variáveis categóricas e numéricas, além de fornecer métricas para avaliar a força da correlação.

Para o cálculo do *Permutation Importance*, foi utilizado o modelo de *Random Forest Regressor* do `cuml` (25) e o método de *permutation importance* da biblioteca `sklearn` (26).

Foi necessário já separar os conjuntos de dados em treino e teste, afim de evitar vazamento de dados e garantir que a avaliação da importância das variáveis seja feita de forma justa e realista.

Para essa separação, foi utilizado o método `train_test_split` da biblioteca `sklearn`, com uma proporção de 80% dos dados para treino e 20% para teste. A partir desse momento, não usamos mais o conjunto de dados completo, mas sim apenas o conjunto de treino para as análises das variáveis preditoras.

### 3.5 Treinamento dos Modelos

Nesta seção serão detalhadas as técnicas de modelagem empregadas e os critérios utilizados para a seleção dos modelos. Como as nossas variáveis respostas são numéricas e contínuas, estamos em um problema de regressão de aprendizado supervisionado. Assim, é preciso selecionar modelos de regressão, que são os adequados para prever variáveis contínuas.

Usaremos três modelos de regressão: (i) *Random Forest Regressor*, (ii) *XGBoost Regressor* e (iii) *LightGBM Regressor*. Esses modelos são amplamente utilizados devido à sua eficácia e interpretabilidade, além de serem capazes de lidar com grandes volumes de dados e variáveis preditoras.

Para otimizar o desempenho dos modelos selecionados, será realizada uma busca em grade (*Grid Search*) para identificar os melhores hiperparâmetros. A busca não será realizada utilizando validação cruzada, devido ao alto volume de dados, mas sim um subconjunto do conjunto de treino equivalente a 10% dos dados originais, garantindo que os resultados sejam robustos e generalizáveis.

A implementação da busca em grade será feita manualmente, utilizando loops para iterar sobre os diferentes valores dos hiperparâmetros e avaliando o desempenho dos modelos utilizando o conjunto de validação. Isso se dá para evitar o alto custo computacional associado à utilização de bibliotecas como `GridSearchCV` da `sklearn`, que realizam a busca em grade utilizando validação cruzada, o que pode ser inviável para grandes volumes de dados e modelos complexos.

Definidos os hiperparâmetros ótimos, os modelos serão treinados utilizando o conjunto de dados preparado na Seção 3.3, com uma quantidade maior de estimadores fracos (através do hiperparâmetro `n_estimators`) para garantir um melhor desempenho dos modelos. O treinamento será realizado utilizando o conjunto de treino completo, garantindo que os modelos sejam treinados com a maior quantidade possível de dados para melhorar sua capacidade de generalização.

Finalizando a etapa de treinamento, serão criados modelos de ensemble utilizando a técnica de *bagging*, onde os modelos de regressão individuais serão combinados para

criar um modelo mais robusto e preciso. A combinação será feita utilizando a média das previsões dos modelos individuais, garantindo que o modelo de ensemble aproveite as forças de cada modelo individual para melhorar a precisão das previsões.

### 3.6 Avaliação dos Modelos

Feito o treinamento final dos modelos, iremos avaliá-los utilizando o conjunto de teste e duas métricas de desempenho apropriadas para problemas de regressão: a Raiz do Erro Quadrático Médio (*Root Mean Squared Error* - RMSE) e Erro Percentual Absoluto Médio (*Mean Absolute Percentage Error* - MAPE). Essas métricas fornecem uma visão clara da precisão das previsões dos modelos em relação aos valores reais das notas do ENEM.

Para a nota da Redação, como estas apresentam uma quantidade limitada de valores possíveis (de 0 a 1000, com incrementos de 20 pontos), será realizado um tratamento adicional para arredondar as previsões dos modelos para os valores possíveis antes do cálculo das métricas, garantindo que as previsões sejam coerentes com a escala de notas do ENEM, com o código abaixo:

```
def arredonda_redacao(
    nota: float
) -> int:

    """
    Arredonda a nota de redação para o múltiplo de 20 mais próximo
    """

    passo = 20

    return np.round(nota / passo) * passo
```

Após o cálculo das métricas de desempenho, duas análise serão realizadas: (i) uma análise entre o erro de treinamento e o erro de teste para avaliar se os modelos estão sofrendo de *overfitting* e (ii) uma análise entre os modelos para identificar qual modelo apresenta o melhor desempenho na previsão das notas do ENEM.

Na análise de *overfitting*, será adotado um critério sobre a razão do erro RMSE de teste em relação ao erro RMSE de treinamento. Uma razão maior que 15% será considerada um indicativo de *overfitting*, sugerindo que o modelo está se ajustando demais aos dados de treinamento e não generalizando bem para os dados de teste.

Na análise comparativa entre os modelos, para se escolher o melhor modelo para cada variável resposta, será considerado o modelo que apresentar o menor erro MAPE no conjunto de teste, garantindo que o modelo selecionado seja aquele que tem a melhor capacidade de prever as notas do ENEM com precisão.

### **3.7 Influência das Variáveis Preditoras**

Nesta seção, buscamos responder às perguntas de pesquisa que foram formuladas na seção 4.1, utilizando técnicas de interpretação de modelos e análise de sensibilidade para medir a magnitude da influência das variáveis preditoras nas variáveis resposta.

Primeiramente, será realizada uma análise de importância das variáveis preditoras do modelo final selecionado para cada variável resposta. Essa análise permitirá identificar quais variáveis preditoras têm a maior influência na previsão das notas do ENEM, fornecendo insights sobre os fatores socioeconômicos que mais impactam o desempenho dos estudantes.

Em seguida, será realizada uma análise de sensibilidade para avaliar como as variações nas variáveis preditoras afetam as previsões dos modelos. Será construída uma base sintética de dados, onde algumas variáveis preditoras selecionadas serão as variáveis de interesse e as demais variáveis preditoras serão preenchidas com o valor mais frequente do conjunto de treino. A ideia é buscar entender como as variações nas variáveis de interesse impactam as previsões dos modelos, mantendo as demais variáveis constantes.

## 4 RESULTADOS

Este capítulo apresenta os resultados obtidos a partir da aplicação da metodologia descrita no Capítulo 3 - Metodologia. Os resultados serão apresentados na mesma ordem das etapas descritas na metodologia.

### 4.1 Entendimento de Negócio

Conforme mencionado no Capítulo 2 - Fundamentação Teórica, o ENEM é um exame de grande relevância no contexto educacional brasileiro e compreender os fatores que impactam o desempenho dos estudantes é crucial para a formulação de políticas educacionais eficazes.

Trabalhos anteriores citam algumas variáveis socioeconômicas como discriminadores de performance no ENEM. A Tabela 1 apresenta essas variáveis identificadas na literatura, juntamente com suas respectivas referências.

Tabela 1 – Variáveis socioeconômicas e suas referências

Variável socioeconômica	Referência
Renda familiar	Melo <i>et al.</i> (1) Vasconcellos (8)
Raça / Cor	Melo <i>et al.</i> (1) Moraes <i>et al.</i> (10)
Sexo	Moraes <i>et al.</i> (10)
Idade / Atraso Escolar	Jaloto e Primi (9)
Administração: Pública vs. Privada	Moraes <i>et al.</i> (10) Jaloto e Primi (9) Ortega <i>et al.</i> (2)
Atributos Escolares	Moraes <i>et al.</i> (10)

Fonte: elaborado pelo autor.

Assim, ao avaliarmos os trabalhos anteriores disponíveis, concluímos que há uma variedade de fatores socioeconômicos que podem influenciar o desempenho dos estudantes no ENEM. Com base nisso, foram formuladas as seguintes perguntas de pesquisa:

- **Pergunta 1:** Quais são os principais fatores socioeconômicos que influenciam o desempenho dos estudantes no ENEM?
- **Pergunta 2:** Qual é a magnitude da influência de cada um desses conjuntos de fatores nas notas dos participantes?

## 4.2 Entendimento dos dados

### 4.2.1 Escolha e Coleta dos Dados

Como descrito no Capítulo 3 - Metodologia, foi necessário identificar dados que fossem relevantes para responder as perguntas de pesquisa formuladas. Foi realizada uma busca por bases de dados públicas que contivessem informações detalhadas sobre os participantes do ENEM, incluindo suas características socioeconômicas e desempenho no exame.

Os microdados do ENEM, disponibilizados anualmente pelo Instituto Nacional de Estudos e Pesquisas Educacionais Anísio Teixeira (INEP), foram escolhidos como a principal fonte de dados para este trabalho e podem ser acessados através do portal do INEP (4).

No mesmo portal, também estão disponíveis os dados do Censo Escolar, que fornecem informações adicionais sobre as escolas de todo o território nacional (15). Esses foram escolhidos como fonte complementar por fornecerem um contexto mais amplo sobre o ambiente educacional.

Foram então selecionadas as edições de 2020 a 2024 (as últimas cinco edições disponíveis) de ambos os conjuntos de dados e os arquivos disponibilizados foram baixados através de download simples e armazenados localmente para posterior leitura e manipulação.

Como os dados escolhidos são públicos e anonimizados por quem os distribui, entendeu-se que não há limitações éticas para o uso desses dados neste trabalho e não foi necessário submeter o projeto a um comitê de ética em pesquisa.

### 4.2.2 Compreensão Inicial dos Dados

Os arquivos de microdados do ENEM e do Censo Escolar são disponibilizados em formato compactado (.zip), separados pelo ano de aplicação do exame/censo.

Dentre os arquivos existentes nos arquivos compactados dos microdados do ENEM, foram selecionados os arquivos CSV (*Comma-Separated Values*) que contêm as informações dos participantes e suas notas e os dicionários de dados de cada edição em formato XLSX (Formato nativo do *Microsoft Excel*), que foi utilizado para interpretar os valores categóricos e identificar variáveis importantes.

Para os arquivos compactados do Censo Escolar, foram selecionados os arquivos CSV que contêm as informações das escolas e os dicionários de dados em formato XLSX.

#### 4.2.2.1 Edição de 2024 do ENEM e LGPD

Na edição de 2024 dos microdados do ENEM, foi feita uma alteração no formato de disponibilização dos dados dos participantes e das notas, que passaram a ser disponibilizados

em arquivos separados.

Isso se deu "Devido à vigência da Lei Geral de Proteção de Dados (LGPD), incorporada ao ordenamento jurídico brasileiro por meio da Lei nº 13.709, de 14 de agosto de 2018"(27), conforme descrito no arquivo auxiliar "Leia-Me"(27) disponibilizado junto com os microdados do ENEM 2024 .

Assim, o formato dos arquivos de microdados do ENEM 2024 difere das edições anteriores, por mais que as informações contidas permanecem as mesmas. Houve a separação dos dados dos participantes e das notas em arquivos distintos e sem uma chave primária que permita a junção dos dois conjuntos de dados. Dessa forma, os dados da edição de 2024 do ENEM não puderam ser utilizados para este trabalho

#### 4.2.3 Análise dos Dicionários de Dados

Foram analisados os dicionários de dados dos microdados do ENEM e do Censo Escolar para identificar as variáveis disponíveis em cada conjunto de dados. Os dicionários completos estão disponíveis no Apêndice A e B. A partir dessa análise, foi possível identificar as variáveis que seriam relevantes para responder às perguntas de pesquisa formuladas na Seção 4.1.

Não foi possível localizar uma variável que permitisse a identificação única das escolas dos participantes do ENEM nos microdados do ENEM, o que impossibilitou correlacionar diretamente os dados dos participantes do ENEM com os dados das escolas do Censo Escolar para agregar informações das escolas aos dados dos participantes. Dessa forma, optou-se por utilizar apenas os dados dos microdados do ENEM para a realização deste trabalho.

#### 4.2.4 Definição da Variável Resposta

Como esse trabalho pretende avaliar o desempenho dos estudantes no ENEM e os fatores que influenciam esse desempenho, a variável resposta deve refletir esse objetivo. Assim, foram utilizadas como variáveis resposta as notas obtidas pelos estudantes nas quatro provas objetivas e na redação do ENEM.

Ou seja, usamos cinco variáveis resposta distintas para análise: (i) Nota da prova de Ciências da Natureza; (ii) Nota da prova de Ciências Humanas; (iii) Nota da prova de Linguagem e Códigos; (iv) Nota da prova de Matemática; e (v) Nota da Redação.

### 4.3 Preparação dos dados

#### 4.3.1 Preparação do Ambiente Tecnológico e Analítico

Para a execução desse trabalho, foi utilizado um ambiente baseado em Python versão 3.11 através do gerenciador de ambientes virtuais Miniconda3 (28). O computador

utilizado possui uma CPU AMD Ryzen 7 9800X3D, 32 GB de memória RAM e uma GPU NVIDIA GeForce RTX 4070 Ti Super, com 16 GB de memória dedicada com sistemas operacionais Ubuntu 24.04 LTS e Windows 11 Pro.

O ambiente foi especificamente configurado com o ecossistema NVIDIA CUDAX (29) para possibilitar a execução utilizando a GPU do equipamento, visando acelerar o processamento dos dados e a modelagem. Esta suíte de bibliotecas de software permite executar pipelines de Ciência de Dados e análises inteiramente na GPU, minimizando a transferência de dados entre a CPU e a GPU.

Foram utilizados seus principais componentes: `cudf` (30), uma biblioteca para manipulação de `DataFrames` na GPU, análoga ao `pandas` (31), e `cuml` (25), que fornece implementações de algoritmos de *Machine Learning* acelerados por GPU, análoga ao `scikit-learn` (26). Todo o ambiente foi construído sobre a plataforma CUDA 13.0, com as bibliotecas e dependências gerenciadas diretamente pelo Conda.

O arquivo YML de configuração do ambiente virtual utilizado está disponível no Apêndice C.

#### 4.3.2 Leitura dos Dados

Os arquivos CSV dos microdados do ENEM foram lidos utilizando o método `read_csv` da biblioteca `pandas` especificando o separador como ponto e vírgula (`sep = ' ; '`).

Serão utilizados os dados das edições de 2020 a 2023 e possuem as quantidade de observações e variáveis descritas na Tabela 2. As tabelas foram carregadas já desconsiderando colunas que não agregam ao modelo, como o número de inscrição do participante, por exemplo.

Tabela 2 – Quantidade de observações e variáveis por edição do ENEM

Edição	Observações	Variáveis
2020	5.783.109	52
2021	3.389.832	52
2022	3.476.105	52
2023	3.933.955	52

Fonte: microdados do INEP; elaborado pelo autor.

#### 4.3.3 Integração dos Dados

Analizando o dicionário de dados de cada edição, foi possível observar que todas as edições possuem o mesmo esquema, ou seja, as mesmas variáveis com os mesmos nomes estão presentes em todas as edições selecionadas. Assim, a integração entre edições foi

realizada por meio da concatenação vertical dos quatro conjuntos de dados, utilizando o método `concat` da biblioteca `pandas`.

Em seguida, foi feita uma modificação no nome das variáveis para nomes que fossem mais intuitivos e de compreensão rápida do conteúdo. Essa modificação foi realizada utilizando o método `rename`, a partir de um dicionário que mapeava os nomes originais para os novos nomes desejados.

Com o dicionário de dados analisado, foi diagnosticado que algumas variáveis categóricas estavam codificadas com valores numéricos que não eram intuitivos. Assim, seus valores foram transformados, substituindo os códigos numéricos por descrições textuais mais comprehensíveis através do método `map` da biblioteca `pandas`, utilizando dicionários de mapeamento construídos especificamente para cada variável categórica que necessitava de transformação.

#### 4.3.4 Tratamento de Valores Nulos

Inicialmente, foi feito o cálculo do percentual de valores nulos por variável. A Tabela 3 apresenta estes valores do conjunto de dados integrados, antes dos tratamentos.

Tabela 3 – Percentual de valores nulos por variável

Variável	Percentual de nulos
sigla_uf_escola	78,1%
cod_municipio_escola	78,1%
tp_adm_escola	78,1%
funcionamento_escola	78,1%
tp_local_escola	78,1%
tp_ensino	69,8%
tp_escola	68,2%
ano_conclusao_ensino_medio	51,6%
nota_ciencias_natureza	40,4%
nota_matematica	40,4%
nota_ciencias_humanas	37,0%
nota_redacao	37,0%
nota_linguagem_codigos	37,0%
03_ocupacao_pai	12,5%
01_escolaridade_pai	9,8%
04_ocupacao_mae	9,0%
estado_civil	4,2%
02_escolaridade_mae	3,5%

*Continua na próxima página...*

Variável	Percentual de nulos
cor_raca	1,8%
10_qtde_carro	0,6%
05_qtde_moradores	0,6%
06_renda_familiar	0,6%
07_dias_trabalhador_domestico	0,6%
08_qtde_banheiro	0,6%
09_qtde_quarto	0,6%
18_flag_aspirador_po	0,6%
11_qtde_motocicleta	0,6%
12_qtde_geladeira	0,6%
13_qtde_freezer	0,6%
14_qtde_maq_lavar_roupa	0,6%
15_qtde_maq_secar_roupa	0,6%
16_qtde_micro_ondas	0,6%
17_qtde_maq_lavar_louca	0,6%
22_qtde_celular	0,6%
19_qtde_tv	0,6%
20_flag_aparelho_dvd	0,6%
21_flag_tv_assinatura	0,6%
24_qtde_computadores	0,6%
23_flag_telefone_fixo	0,6%
25_flag_internet	0,6%
nacionalidade	0,05%

Fonte: elaborado pelo autor.

O percentual de valores nulos varia significativamente, com algumas variáveis apresentando mais de 70% de valores nulos, enquanto outras possuem menos de 1%. Variáveis com uma alta proporção de valores nulos podem comprometer a análise se for realizado alguma imputação de valores. Sendo assim, foi decidido remover as variáveis que apresentavam mais de 50% de valores nulos, resultando na remoção de nove variáveis do conjunto de dados.

Para as variáveis das notas, por serem as variáveis resposta deste trabalho, foi realizada uma análise mais detalhada. Primeiro, verificou-se a existência de valores zerados nessas variáveis e se possuem significado diferentes de valores nulos. Para isso, foi feita uma análise com a presença nas provas e o status da redação. Foi identificado que a nota zerada significa que o participante esteve presente na prova, mas obteve nota zero, enquanto o valor nulo indica que o participante ou não realizou a prova, ou foi eliminado, ou teve

sua redação anulada. Dessa forma, optou-se por manter as observações com notas zeradas no conjunto de dados, removendo apenas as observações com notas nulas. A decisão de incorporar ou não as notas zero na análise será discutida na Seção 4.4.

Para as demais variáveis com valores nulos, foi realizada uma análise consolidada, ou seja, foram retiradas as observações que possuíam valor nulo em qualquer uma das variáveis restantes, o que resultou na retirada de 4.341.559 observações do conjunto de dados.

Assim, restaram 12.241.442 observações e 45 variáveis no conjunto de dados após os tratamentos.

#### 4.3.5 Separação dos Conjuntos de Dados por Variável Resposta

Após a separação dos conjuntos de dados por variável resposta, conforme descrito na Seção 3.3, foram criados cinco conjuntos de dados distintos. A Tabela 4 apresenta a quantidade de observações e variáveis em cada conjunto de dados.

Tabela 4 – Observações e variáveis por conjunto de dados

Conjunto de Dados	Observações	Variáveis
Ciências Humanas	7.895.093	35
Ciências da Natureza	7.500.050	35
Linguagem e Código	7.895.093	35
Matemática	7.500.050	35
Redação	7.895.093	35

Fonte: elaborado pelo autor.

A estrutura de dicionários foi utilizada para manter o controle dos conjuntos de dados, suas respectivas variáveis resposta e variáveis preditoras ao longo do trabalho.

### 4.4 Modelagem

#### 4.4.1 Análise Exploratória dos Dados - Variáveis Resposta

##### 4.4.1.1 Distribuições

O primeiro passo da análise exploratória foi entender o domínio das variáveis respostas e foi constatado que para as notas das provas objetivas (Ciências Humanas, Ciências da Natureza, Linguagem e Códigos e Matemática) haviam mais de cinco mil notas distintas, com variações pequenas entre elas (décimos de pontos). Já para a nota da redação, o número de notas distintas era significativamente menor (apenas 50 notas) com variação de pontos de 20 em 20 pontos.

A tabela 5 apresenta as estatísticas descritivas das notas de cada prova, obtido através do método `describe`.

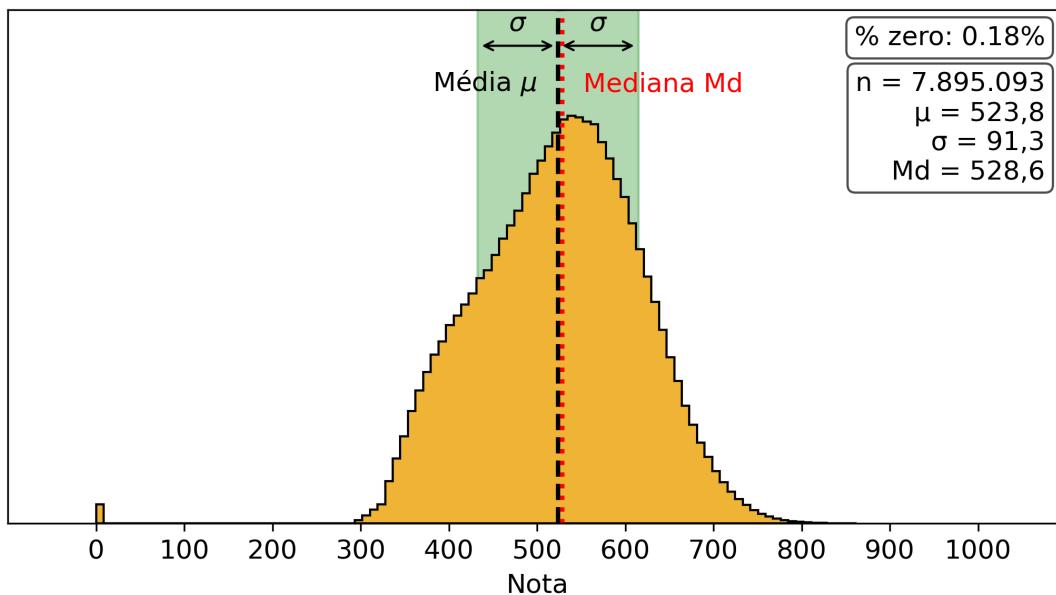
Tabela 5 – Estatísticas descritivas por conjunto de dados

Estatística	Humanas	Natureza	Linguagem	Matemática	Redação
Contagem	7.895.093	7.500.050	7.895.093	7.500.050	7.895.093
Média	523,8	496,9	519,0	538,7	616,6
Desvio Padrão	91,3	81,4	91,3	121,5	204,7
Mínimo	0,0	0,0	0,0	0,0	0,0
25º Percentil	460,1	437,3	460,1	441,6	520
50º Percentil	528,6	490,3	528,6	526	620
75º Percentil	588,2	551,9	588,2	624,9	760
Máximo	862,6	875,3	826,1	985,7	1000

Fonte: elaborado pelo autor.

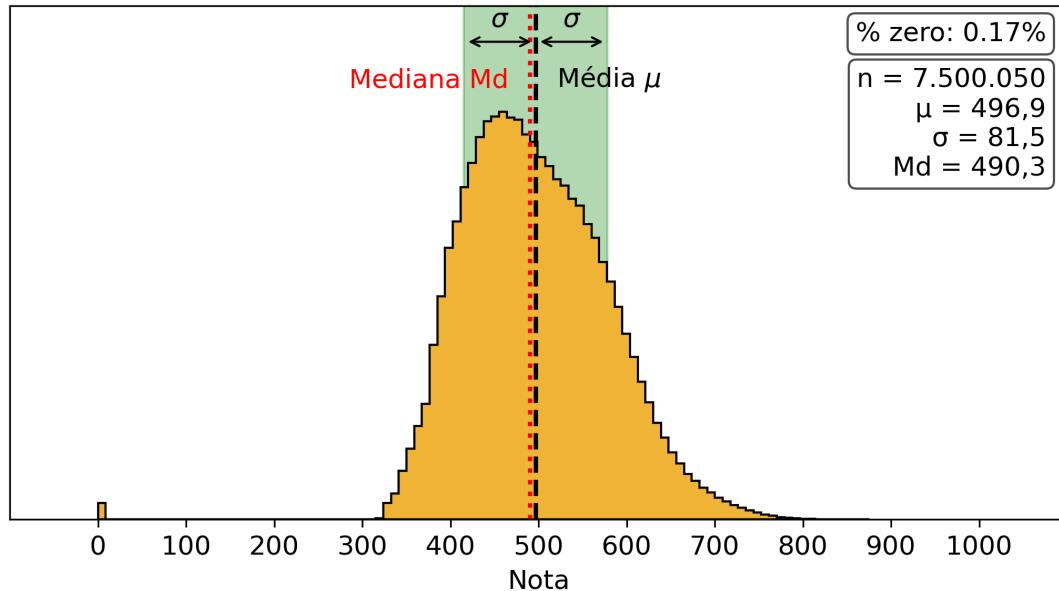
Em seguida, foram construídos histogramas de cada nota para entender a distribuição das notas. As Figuras 5, 6, 7, 8 e 9 apresentam os histogramas das notas de cada prova.

Figura 5 – Histograma das notas - Ciências Humanas



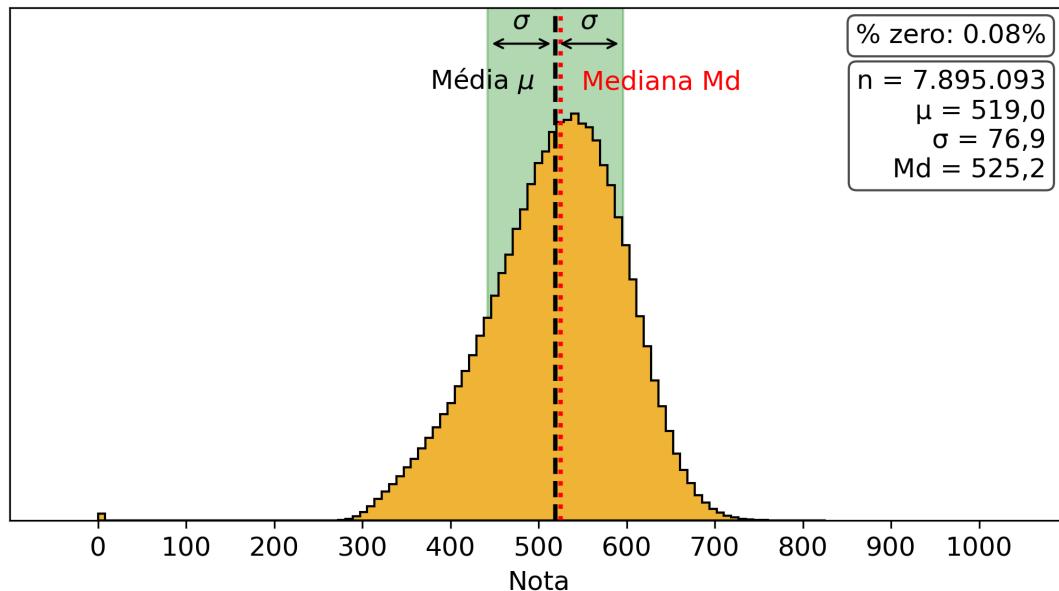
Fonte: elaborado pelo autor.

Figura 6 – Histograma das notas - Ciências da Natureza



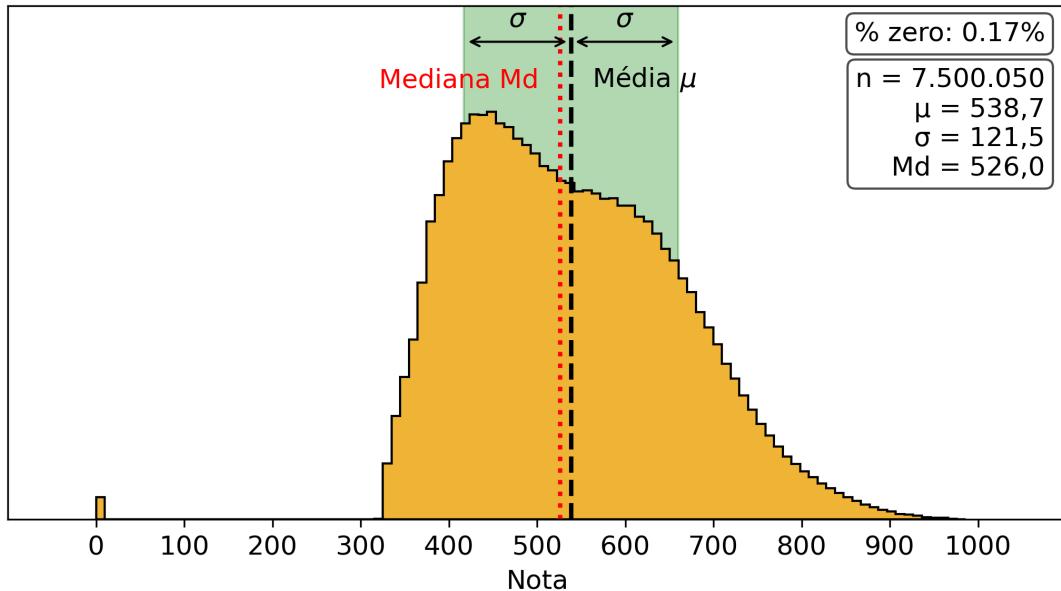
Fonte: elaborado pelo autor.

Figura 7 – Histograma das notas - Linguagem e Código



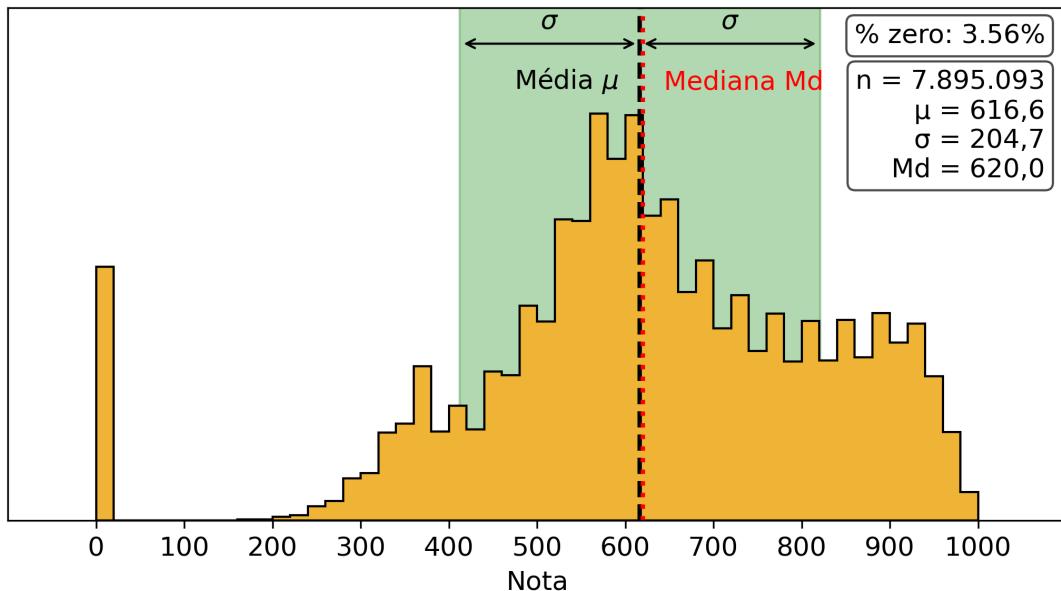
Fonte: elaborado pelo autor.

Figura 8 – Histograma das notas - Matemática



Fonte: elaborado pelo autor.

Figura 9 – Histograma das notas - Redação



Fonte: elaborado pelo autor.

A tabela 6 apresenta os valores de assimetria e curtose das notas de cada prova, obtidos através dos métodos `skew` e `kurtosis` da biblioteca `pandas`, e o percentual de notas zero em cada conjunto de dados.

Tabela 6 – Assimetria, Curtose e Notas zeradas

Variável	Assimetria	Curtose	Notas zeradas
Humanas	-0,3408	1,1269	0,18%
Natureza	0,0321	1,8372	0,17%
Linguagem	-0,5113	1,1780	0,08%
Matemática	0,3138	0,0850	0,17%
Redação	-0,7457	1,0488	3,56%

Fonte: elaborado pelo autor.

Analizando os valores de assimetria e curtose, é possível observar que as distribuições das notas possuem diferentes características.

A assimetria da nota de redação é a mais negativa, o que indica que os alunos tiveram, em geral, o melhor desempenho, assim como nas provas de Linguagem e Códigos e Ciências Humanas, que também apresentam assimetria negativa, porém com valores menores. Na prova de Ciências da Natureza, a assimetria é praticamente nula, indicando uma distribuição mais simétrica das notas, enquanto a nota de Matemática apresenta a assimetria mais positiva, indicando um desempenho relativamente pior dos alunos nessa prova.

Analizando os valores da curtose, a prova de Matemática foi a única a apresentar uma curtose próxima de zero, indicando uma distribuição mais próxima da normalidade. As outras provas apresentaram valores maiores que 1 indicando distribuições com caudas mais pesadas e picos mais acentuados.

#### 4.4.1.2 Teste de Hipótese

Foi realizado o teste de hipótese ANOVA com nível de significância de 0,1% para comparar as médias das notas por edição do ENEM em cada conjunto de dados, onde a hipótese nula  $H_0$  é de que as médias são iguais entre as edições. A tabela 7 apresenta os valores de F, p-valor e a métrica SMD (*Standardized Mean Difference*) obtidos para cada conjunto de dados.

Tabela 7 – Teste ANOVA das médias das notas por edição

Variável	Valor F	Rejeita-se $H_0$ ?	p-valor	SMD	Tamanho do efeito
Humanas	13.161	Sim	0,0000	0,185	Insignificante
Natureza	3.431	Sim	0,0000	0,087	Insignificante
Linguagem	26.506	Sim	0,0000	0,269	Pequeno
Matemática	13.041	Sim	0,0000	0,197	Insignificante
Redação	27.498	Sim	0,0000	0,242	Pequeno

---

Fonte: elaborado pelo autor.

Foi decidido manter todas as edições do ENEM no conjunto de dados para a modelagem preditiva e sem a necessidade de segmentação por edição, uma vez que o tamanho do efeito é insignificante ou pequeno.

#### 4.4.1.3 Análise de Outliers

Para a análise dos outliers, foram utilizados os boxplots das notas de cada prova, apresentados nas Figuras 10, 11, 12, 13 e 14.

Figura 10 – Boxplot das notas por edição - Ciências Humanas

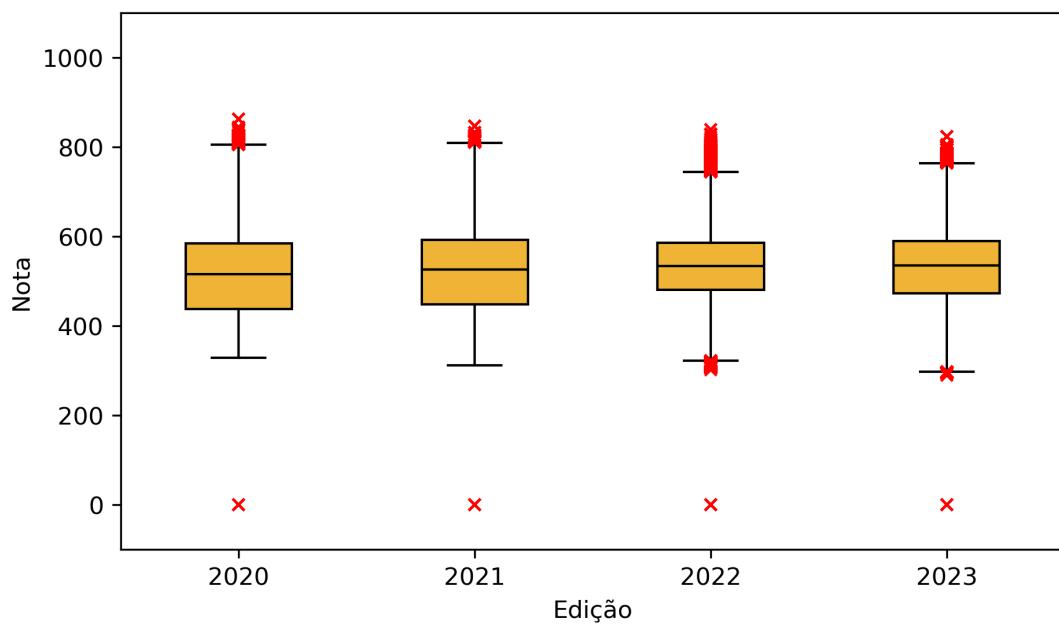


Figura 11 – Boxplot das notas - Ciências da Natureza

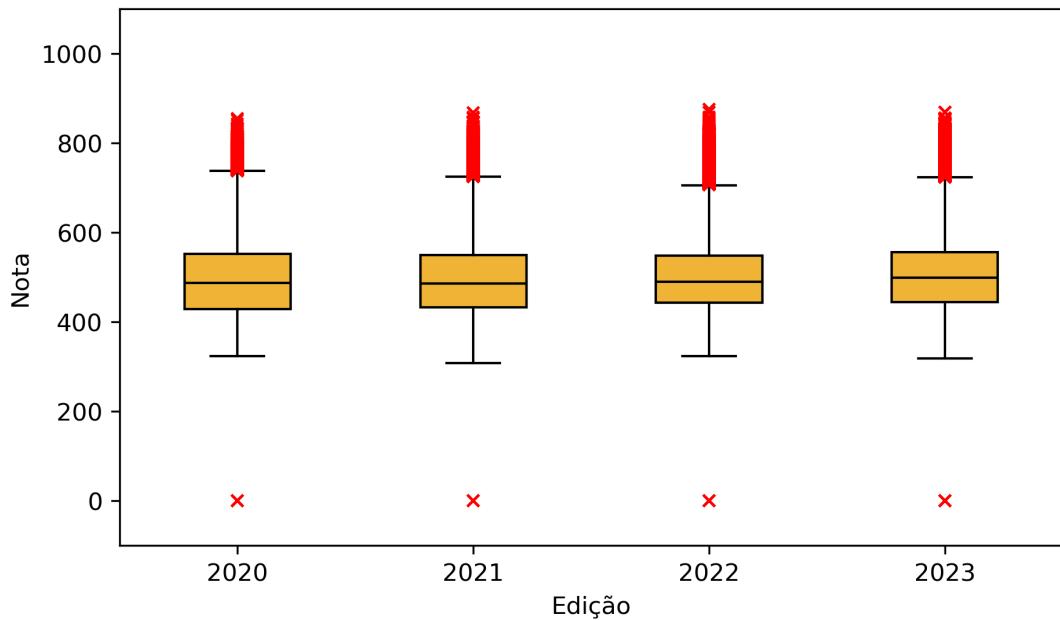


Figura 12 – Boxplot das notas - Linguagem e Código

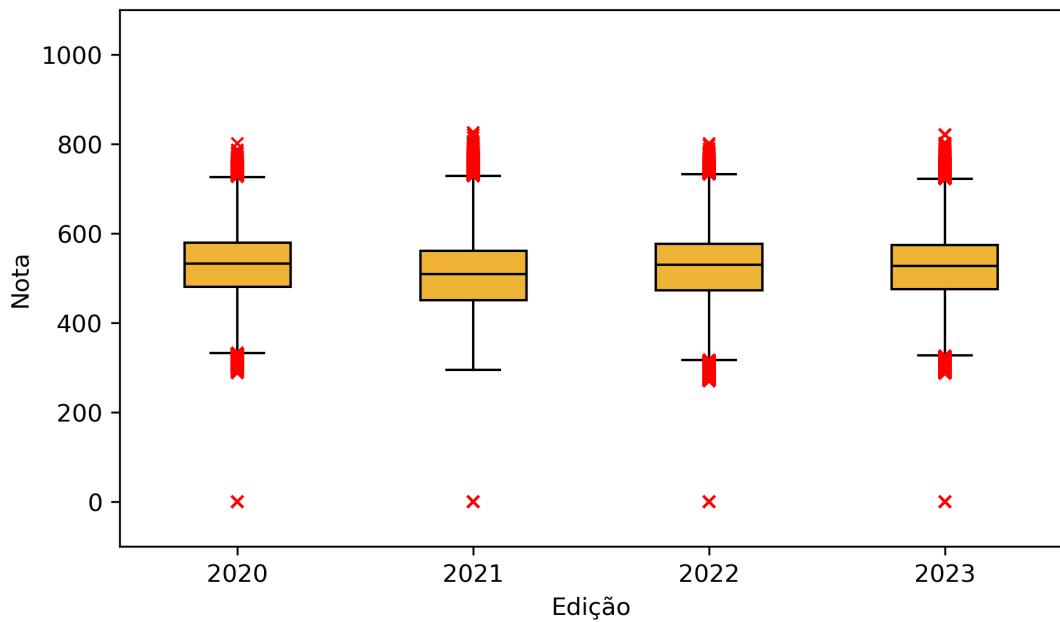


Figura 13 – Boxplot das notas - Matemática

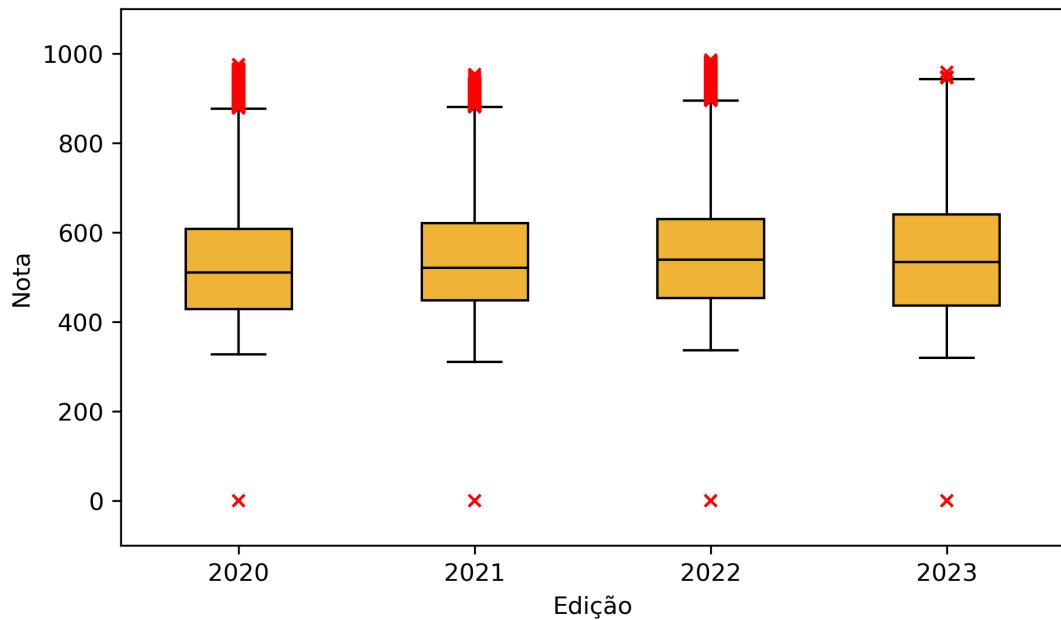
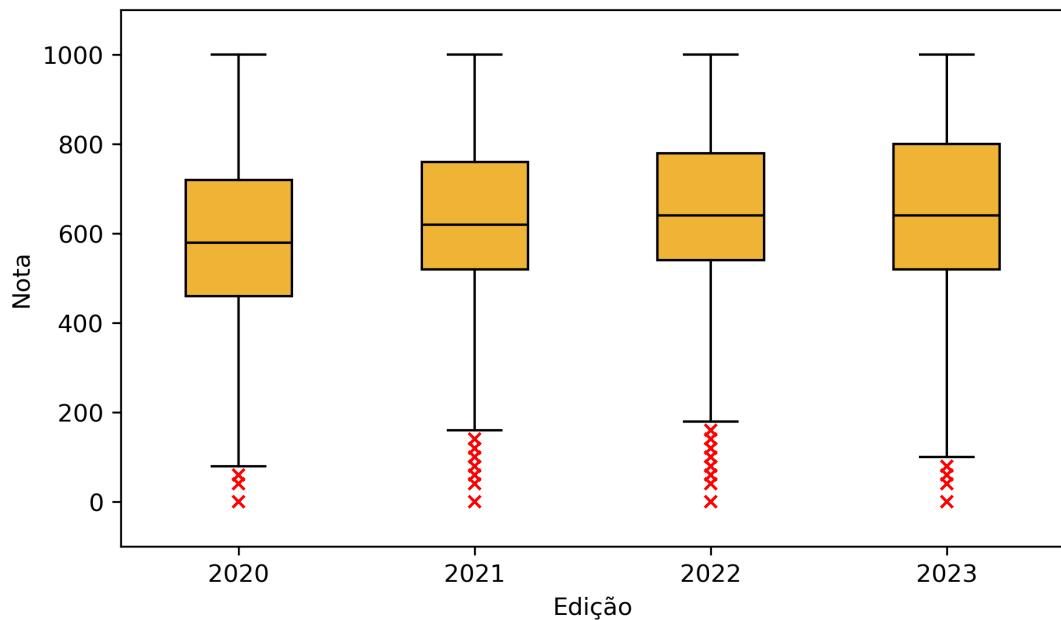


Figura 14 – Boxplot das notas - Redação



Fonte: elaborado pelo autor.

Foi utilizado o critério do intervalor interquartil (*Interquartile Range - IQR*) para identificar os outliers nas notas de cada prova. Foram considerados outliers os valores que estavam abaixo de  $Q1 - 1,5 \times IQR$  ou acima de  $Q3 + 1,5 \times IQR$ , onde Q1 é o primeiro quartil, Q3 é o terceiro quartil e  $IQR = Q3 - Q1$ . A Tabela 8 apresenta o limite inferior, o limite superior, a quantidade e o percentual de outliers identificados em cada conjunto de dados.

Tabela 8 – Quantidade e percentual de outliers nas notas

Variável	Limite Inferior	Limite Superior	Outliers
Humanas	267,95	780,35	19.596 (0,25%)
Natureza	265,40	723,80	49.017 (0,65%)
Linguagem	314,85	728,85	35.998 (0,46%)
Matemática	166,65	899,85	27.944 (0,37%)
Redação	160,00	1000,00	282.438 (3,58%)

Fonte: elaborado pelo autor.

Ao analisarmos os intervalos dos outliers da nota da Redação, vemos que a nota máxima (1.000 pontos) não foi considerada outlier. Porém, tendo o contexto do ENEM em mente e analisando a distribuição das notas, é razoável retirarmos as notas máximas do conjunto de dados, uma vez que são notas extremamente raras. Isso resulta em mais 116 observações removidas do conjunto de dados da nota de Redação.

#### 4.4.2 Análise Exploratória - Variáveis Preditoras

A análise exploratória das variáveis preditoras seguiu três etapas: (i) concentração, (ii) correlação e (iii) *Permutation Importance*, conforme descrito na Seção 3.4.1.

##### 4.4.2.1 Concentração

Ao calcularmos a proporção de observações para cada categoria das variáveis categóricas, foi possível identificar três variáveis que apresentavam uma concentração acima de 93%. Devido a essa alta concentração, foi decidido remover essas variáveis do conjunto de dados. As tabelas 9 a 13 apresentam as cinco variáveis categóricas com maior concentração e suas respectivas proporções da categoria de maior Concentração para cada conjunto de dados.

Tabela 9 – Cinco maiores concentrações - Humanas

Variável	Maior Concentração
12_qtde_geladeira	92.8%
07_dias_trabalhador_domestico	90.2%
25_flag_internet	89.9%
15_qtde_maq_secar_roupa	86.5%
23_flag_telefone_fixo	84.1%

Fonte: elaborado pelo autor.

Tabela 10 – Cinco maiores concentrações - Natureza

Variável	Maior Concentração
12_qtde_geladeira	92.9%
07_dias_trabalhador_domestico	90.3%
25_flag_internet	89.9%
15_qtde_maq_secar_roupa	86.6%
23_flag_telefone_fixo	84.0%

Fonte: elaborado pelo autor.

Tabela 11 – Cinco maiores concentrações - Linguagem

Variável	Maior Concentração
12_qtde_geladeira	92.9%
07_dias_trabalhador_domestico	90.2%
25_flag_internet	89.9%
15_qtde_maq_secar_roupa	86.5%
23_flag_telefone_fixo	84.1%

Fonte: elaborado pelo autor.

Tabela 12 – Cinco maiores concentrações - Matemática

Variável	Maior Concentração
12_qtde_geladeira	92.8%
07_dias_trabalhador_domestico	90.2%
25_flag_internet	90.0%
15_qtde_maq_secar_roupa	86.5%
23_flag_telefone_fixo	84.0%

Fonte: elaborado pelo autor.

Tabela 13 – Cinco maiores concentrações - Redação

Variável	Maior Concentração
12_qtde_geladeira	92.8%
07_dias_trabalhador_domestico	90.1%
25_flag_internet	90.1%
15_qtde_maq_secar_roupa	86.4%
23_flag_telefone_fixo	83.9%

Fonte: elaborado pelo autor.

#### 4.4.2.2 Correlação Phik

A próxima etapa da análise exploratória das variáveis preditoras foi a análise de correlação utilizando a métrica Phik. As tabelas 14 a 18 apresentam as cinco variáveis com maior correlação Phik com a variável resposta em cada conjunto de dados.

Tabela 14 – Cinco maiores correlações Phik - Humanas

Variável	Correlação Phik
24_qtde_computadores	44,5%
03_ocupacao_pai	39,7%
04_ocupacao_mae	37,6%
08_qtde_banheiro	35,7%
18_flag_aspirador_po	35,4%

Fonte: elaborado pelo autor.

Tabela 15 – Cinco maiores correlações Phik - Natureza

Variável	Correlação Phik
24_qtde_computadores	44,6%
03_ocupacao_pai	40,2%
04_ocupacao_mae	38,0%
08_qtde_banheiro	37,4%
18_flag_aspirador_po	36,4%

Fonte: elaborado pelo autor.

Tabela 16 – Cinco maiores correlações Phik - Linguagem

Variável	Correlação Phik
24_qtde_computadores	44,1%
03_ocupacao_pai	41,7%
04_ocupacao_mae	39,9%
08_qtde_banheiro	36,8%
18_flag_aspirador_po	36,1%

Fonte: elaborado pelo autor.

Tabela 17 – Cinco maiores correlações Phik - Matemática

Variável	Correlação Phik
24_qtde_computadores	47,9%
03_ocupacao_pai	44,7%
04_ocupacao_mae	42,5%
08_qtde_banheiro	41,9%
18_flag_aspirador_po	40,5%

Fonte: elaborado pelo autor.

Tabela 18 – Cinco maiores correlações Phik - Redação

Variável	Correlação Phik
03_ocupacao_pai	35,9%
24_qtde_computadores	35,6%
04_ocupacao_mae	34,9%
08_qtde_banheiro	33,2%
10_qtde_carro	29,3%

Fonte: elaborado pelo autor.

Ao analisarmos as matrizes de correlação Phik completas para cada conjunto de dados, foi possível identificar um par de variáveis que apresentavam uma correlação perfeita ( $\text{Phik} = 1.0$ ): flag de treineiro e status da conclusão do ensino médio.

Devido a essa correlação perfeita, analisamos a distribuição cruzadas das categorias dessas duas variáveis, onde foi possível observar que 100% das observações da categoria "Treineiro" da variável "flag\_treineiro" estavam associadas à categoria "Termina o ensino médio após o ano da prova" da variável "conclusao\_ensino\_medio". Diante isso, dado que a variável "conclusao\_ensino\_medio" apresenta mais categorias e, portanto, mais informações, foi decidido manter essa variável no conjunto de dados e remover a variável "flag\_treineiro".

#### 4.4.2.3 *Permutation Importance*

Conforme descrito na Seção 3.4.1, a última etapa da análise exploratória das variáveis preditoras foi a análise de importância utilizando a métrica *Permutation Importance*. Realizadas as separações dos conjuntos de dados em treino e teste, foi treinado um modelo de *Random Forest Regressor* em cada conjunto de dados para em seguida calcularmos o *Permutation Importance* de cada variável preditora.

As Figuras 15 a 19 apresentam os gráficos de importância das dez variáveis mais importantes para cada conjunto de dados.

Figura 15 – Dez maiores *Permutation Importance* - Humanas

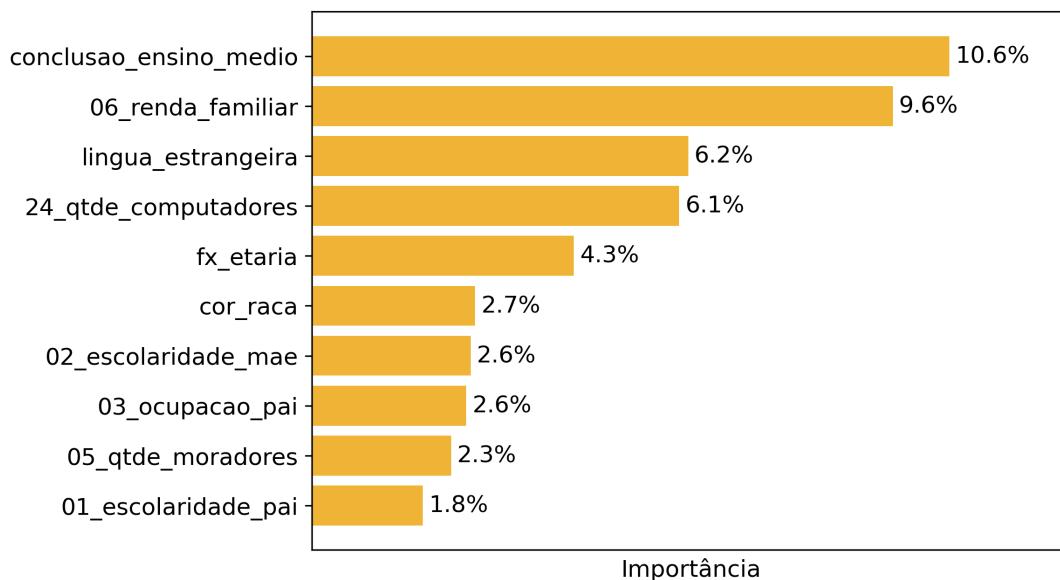


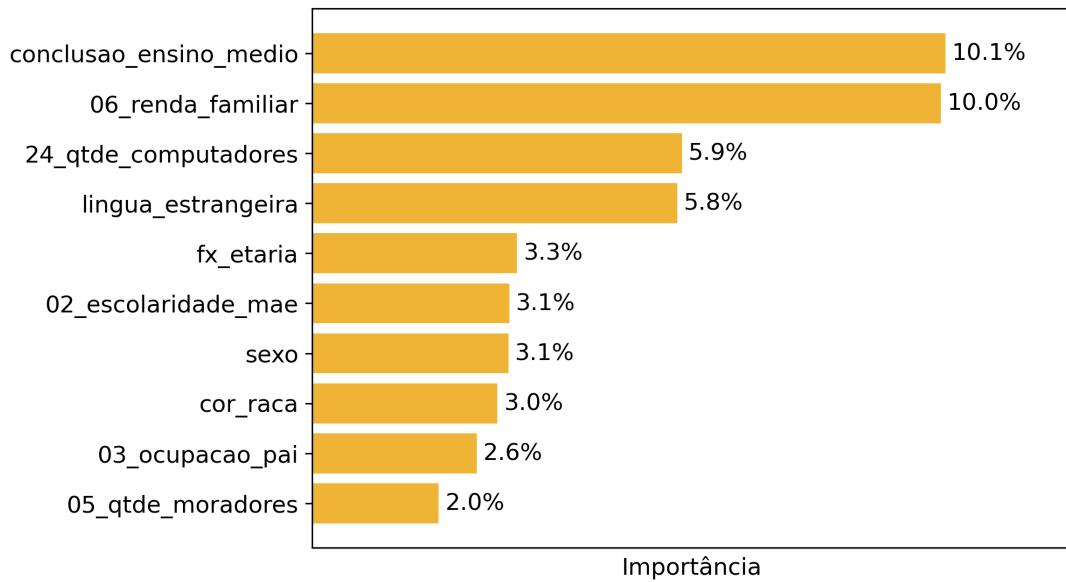
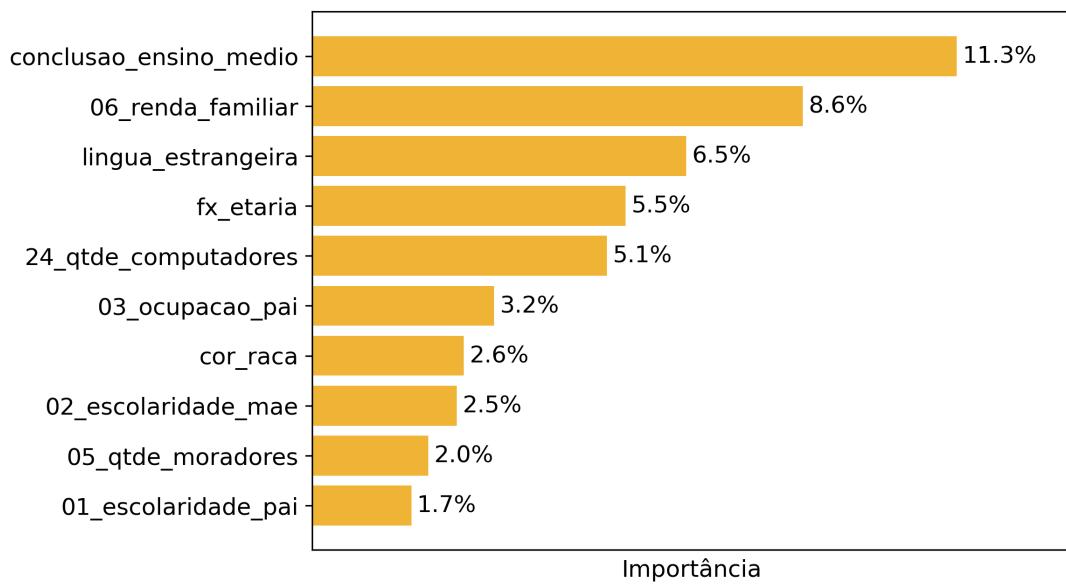
Figura 16 – Dez maiores *Permutation Importance* - NaturezaFigura 17 – Dez maiores *Permutation Importance* - Linguagem

Figura 18 – Dez maiores *Permutation Importance* - Matemática

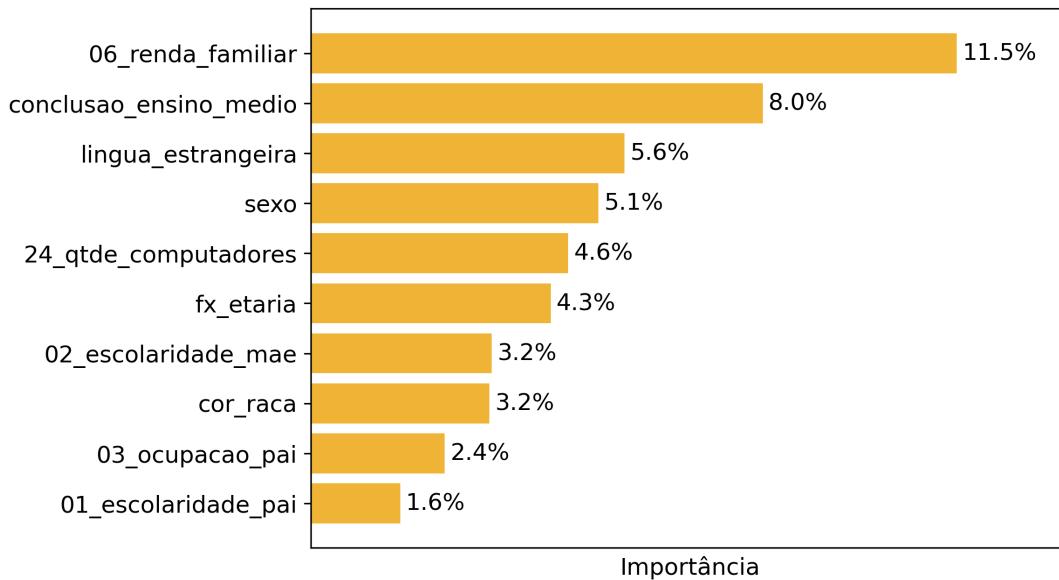
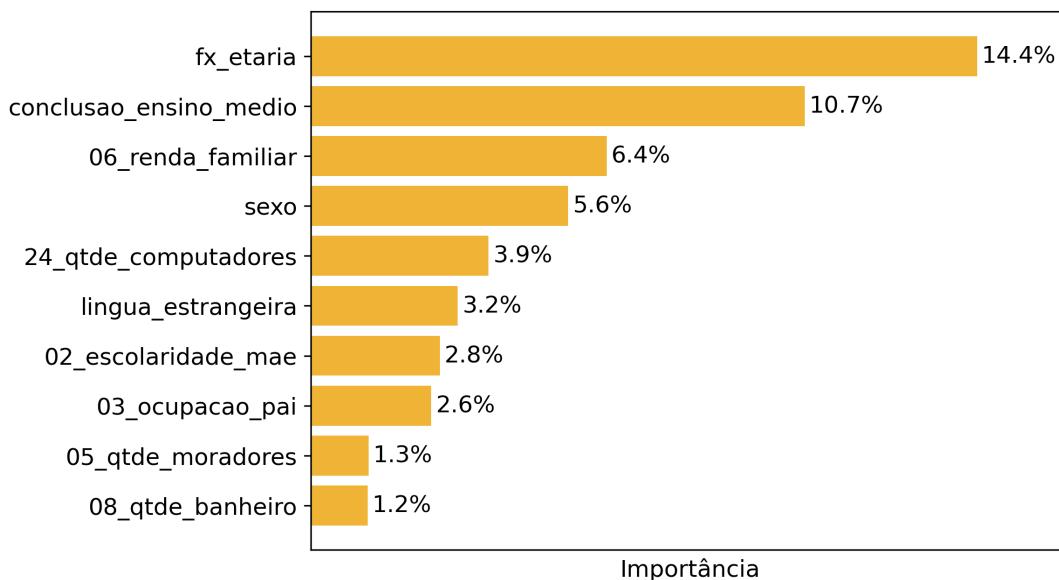


Figura 19 – Dez maiores *Permutation Importance* - Redação



Fonte: elaborado pelo autor.

#### 4.4.2.4 Seleção de Variáveis

Até o momento já aplicamos dois critérios para a seleção das variáveis preditoras: (i) concentração, removendo variáveis com concentração de categoria maior que 93% e (ii) correlação perfeita com outra variável preditora, resultando em quatro variáveis removidas do conjunto de dados: (i) nacionalidade, 17\_qtde\_maq\_lavar\_louca e estado\_civil e (ii) flag\_treineiro.

Com as informações de correlação Phik e *Permutation Importance*, aplicaremos mais dois critérios para a seleção das variáveis preditoras: (i) correlação baixa com a

variável resposta ( $\text{Phik} < 0,05\%$ ) e (ii) um critério duplo de correlação alta ( $\text{Phik} > 85\%$ ) com outras variáveis preditoras e menor *Permutation Importance* entre as variáveis correlacionadas.

No critério de baixa correlação com a variável resposta, nenhuma variável preditora foi removida, uma vez que todas apresentaram Phik maior que 0,05%, sendo o menor valor, entre todos os conjuntos de dados, de 3.64% (variável de quantidade de motociclistas com a nota da Redação).

No critério duplo, foi apresentado apenas uma dupla de variáveis com correlação alta ( $\text{Phik} > 85\%$ ) em cada conjunto de dados: status da conclusão do ensino médio e a faixa etária do participante. Apenas no conjunto de dados da Redação, a variável de faixa etária apresentou maior *Permutation Importance* e então seria mantida no conjunto de dados, enquanto a variável de status da conclusão do ensino médio seria removida.

Porém, considerando a concentração cruzada das duas variáveis (apresentada na Tabela 19 para o conjunto de dados da Redação), foi decidido manter ambas as variáveis no conjunto de dados por escolha do autor.

Tabela 19 – Concentração cruzada - Redação

Faixa etária	Código 1	Código 2	Código 3	Código 4
menor de 17 anos	-	0.4%	10.1%	0.2%
17 anos	0.3%	15.9%	6.6%	-
18 anos	6.3%	16.0%	0.6%	-
19 anos	8.8%	2.4%	0.1%	-
20 anos	6.2%	0.6%	-	-
21 anos	4.3%	0.2%	-	-
22 anos	3.1%	0.1%	-	-
23 anos	2.4%	-	-	-
24 anos	1.9%	-	-	-
25 anos	1.5%	-	-	-
26 a 30 anos	4.5%	0.1%	-	-
31 a 35 anos	2.4%	-	-	-
36 a 40 anos	1.7%	-	-	-
41 a 45 anos	1.2%	-	-	-
46 a 50 anos	0.7%	-	-	-
51 a 55 anos	0.4%	-	-	-
56 a 60 anos	0.2%	-	-	-
61 a 65 anos	0.1%	-	-	-
66 a 70 anos	-	-	-	-

Continua na próxima página...

Faixa etária	Código 1	Código 2	Código 3	Código 4
maior de 70 anos	-	-	-	-

Fonte: elaborado pelo autor.

## 4.5 Treinamento dos Modelos

### 4.5.1 Ajuste dos Hiperparâmetros

A primeira etapa do treinamento dos modelos foi o ajuste dos hiperparâmetros utilizando a técnica de *Grid Search*. Inicialmente, foi utilizado o método `GridSearchCV` da biblioteca `scikit-learn` (26) para realizar o ajuste dos hiperparâmetros dos modelos. Porém, a execução do código foi interrompida subitamente algumas vezes, possivelmente devido ao alto consumo de memória da GPU.

Assim, para contornar esse problema, o *Grid Search* foi implementado manualmente. Foi estabelecido um dicionário com os hiperparâmetros e seus respectivos valores a serem testados para cada modelo e gerada uma lista com todas as combinações possíveis desses hiperparâmetros. Em seguida, através de um *loop*, cada combinação de hiperparâmetros foi utilizada para instanciar cada modelo, treinar o modelo com os dados de treino e avaliar o desempenho do modelo com os dados de validação.

Para gerar o conjunto de dados de validação, foi feita uma separação no conjunto de dados de treino de forma que o conjunto de validação tenha 10% dos dados originais, considerando o conjunto de teste que já foi separado. O desempenho do modelo foi avaliado utilizando a métrica *Root Mean Squared Error* - RMSE (raiz quadrada do erro quadrático médio).

Para cada algoritmo, foi estabelecido um conjunto de hiperparâmetros e seus respectivos valores a serem testados, assim como hiperparâmetros fixos em valores pré-definidos. As tabelas 20 a 22 apresentam os hiperparâmetros e seus respectivos valores a serem testados para cada modelo, bem como os hiperparâmetros fixados em valores pré-definidos.

Tabela 20 – *Grid Search - XGBoost*

Hiperparâmetro	Valores
<code>learning_rate</code>	[0.05, 0.10, 0.20]
<code>max_depth</code>	[6, 8, 10]
<code>min_child_weight</code>	[1, 5, 10]
<code>colsample_bytree</code>	[0.70, 0.85, 1.0]
<code>subsample</code>	[0.70, 0.85, 1.0]
<code>n_estimators</code>	100

Continua na próxima página...

Hiperparâmetro	Valores
objective	"reg:squarederror"
tree_method	"hist"
device	"cuda"
eval_metric	"rmse"

Fonte: elaborado pelo autor.

Tabela 21 – *Grid Search - LightGBM*

Hiperparâmetro	Valores
num_leaves	[31, 63, 127]
learning_rate	[0.05, 0.10, 0.20]
min_child_samples	[20, 50, 100]
colsample_bytree	[0.70, 0.85, 1.0]
subsample	[0.70, 0.85, 1.0]
n_estimators	100
objective	"regression"
metric	"rmse"
device	"cpu"
n_jobs	-1

Fonte: elaborado pelo autor.

Tabela 22 – *Grid Search - Random Forest*

Hiperparâmetro	Valores
max_depth	[10, 15, 20]
max_features	[0.7, 0.9, 1.0]
max_samples	[0.8, 0.9, 1.0]
split_criterion	"mse"
bootstrap	True
n_bins	256
min_samples_leaf	15
n_streams	4
n_estimators	100

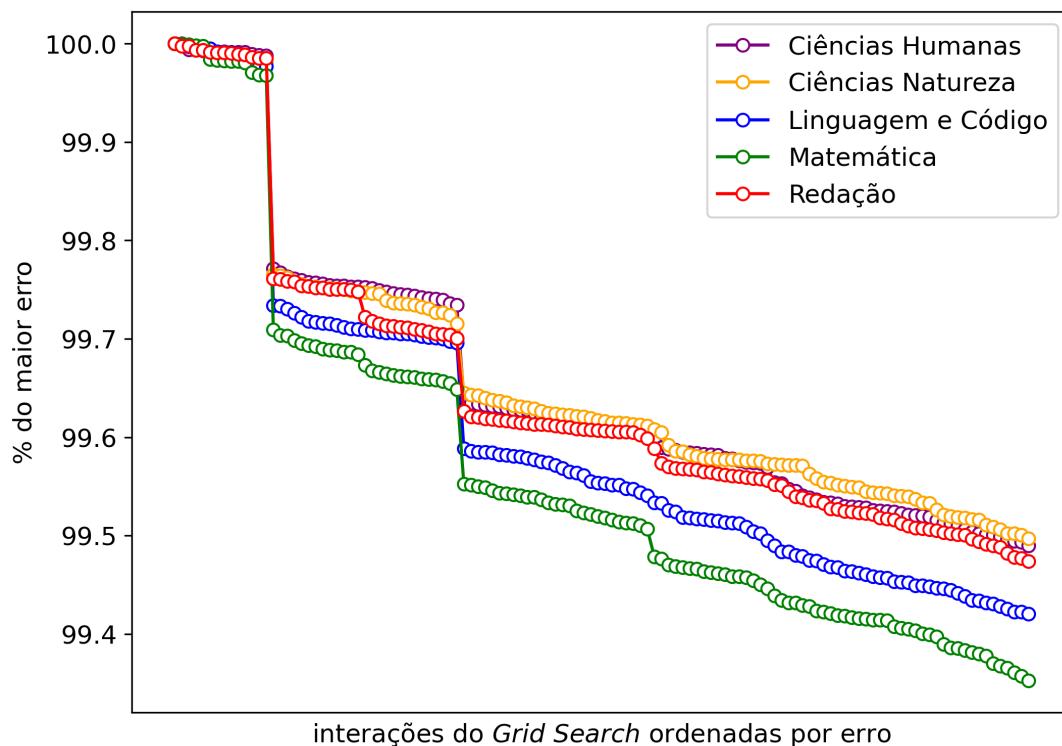
Fonte: elaborado pelo autor.

Foram então treinados 243 combinações de hiperparâmetros para os algoritmos de *XGBoost* e *LightGBM* (5 hiperparâmetros com 3 valores cada) e 27 combinações de hiperparâmetros para o algoritmo de *Random Forest* (3 hiperparâmetros com 3 valores cada)

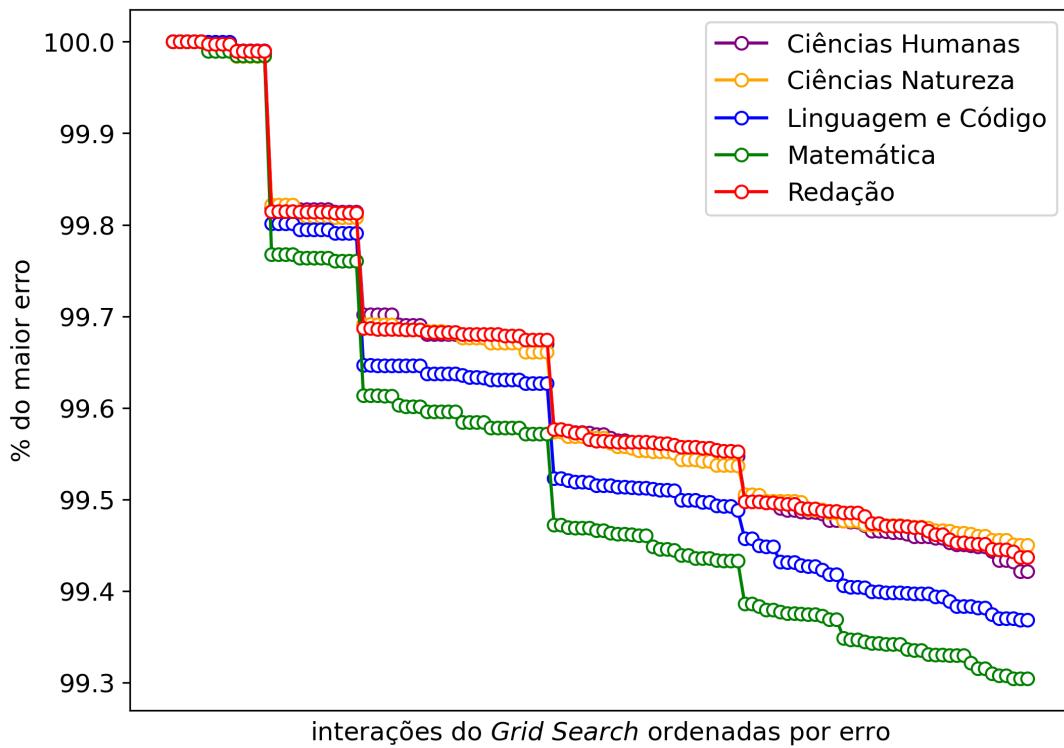
totalizando 270 modelos treinados para cada conjunto de dados (Humanas, Natureza, Linguagem, Matemática e Redação) e 1.350 modelos treinados no total, o que levou aproximadamente 5 horas para ser executado.

Os gráficos 20 a 22 apresentam o erro RMSE para as combinações de hiperparâmetros testadas para cada modelo e conjunto de dados, com as interações do *Grid Search* ordenadas em ordem decrescente de erro. Para melhor visualização, os valores do erro foram relativizados ao maior erro RMSE encontrado para cada modelo e conjunto de dados, ou seja, o maior erro RMSE encontrado para cada modelo e conjunto de dados foi considerado como 100% e os demais erros foram relativizados em relação a esse valor.

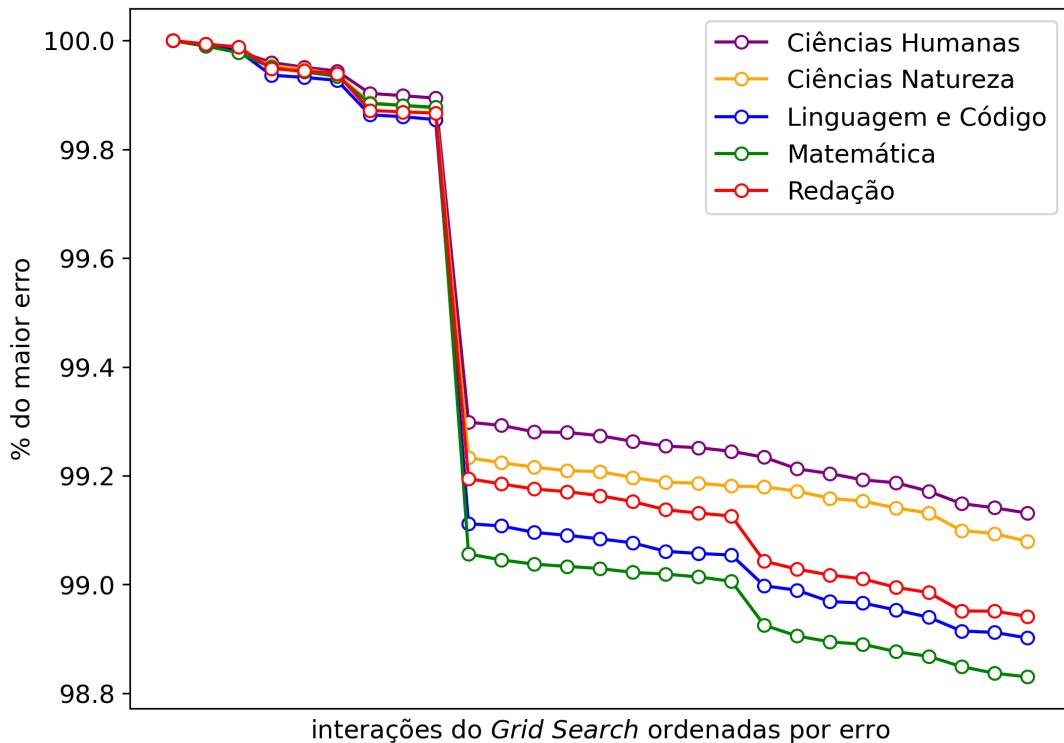
Figura 20 – Erro do *Grid Search* - *XGBoost*



Obs.: apenas metade dos dados foi plotada para melhor visualização.

Figura 21 – Erro do *Grid Search - LightGBM*

Obs.: apenas metade dos dados foi plotada para melhor visualização.

Figura 22 – Erro do *Grid Search - Random Forest*

Fonte: elaborado pelo autor.

A redução no erro RMSE para as melhores combinações de hiperparâmetros acabou não sendo tão significativa para os modelos de *XGBoost* e *LightGBM*, apresentando uma redução menor que 1% em relação ao maior erro RMSE encontrado na *Grid Search*. Já para o modelo de *Random Forest*, a redução no erro RMSE foi mais significativa, apresentando uma redução de aproximadamente 2,2% em relação ao maior erro RMSE encontrado na *Grid Search*.

Até este momento, não ajustamos o hiperparâmetro `n_estimators` (número de estimadores), uma vez que a execução do código para as combinações com esse parâmetro maior que 100 foi interrompida subitamente algumas vezes, também possivelmente devido ao alto consumo de memória da GPU. Dessa forma, o ajuste desse hiperparâmetro será realizado na próxima etapa, juntamente com o treinamento final dos modelos.

As tabelas 23 a 25 apresentam os melhores hiperparâmetros encontrados para cada modelo em cada conjunto de dados, bem como o erro RMSE correspondente a cada combinação de hiperparâmetros.

Tabela 23 – Hiperparâmetros Ajustados - *XGBoost*

Hiperparâmetro	Humanas	Natureza	Linguagem	Matemática	Redação
<code>learning_rate</code>	0.1	0.1	0.1	0.1	0.1
<code>max_depth</code>	10.0	10.0	10.0	10.0	10.0
<code>min_child_weight</code>	10.0	10.0	1.0	5.0	5.0
<code>colsample_bytree</code>	0.7	0.7	0.7	0.7	0.7
<code>subsample</code>	1.0	1.0	1.0	1.0	1.0
RMSE	75,4	65,1	63,1	96,0	148,1

Tabela 24 – Hiperparâmetros Ajustados - *LightGBM*

Hiperparâmetro	Humanas	Natureza	Linguagem	Matemática	Redação
<code>num_leaves</code>	127.0	127.0	127.0	127.0	127.0
<code>learning_rate</code>	0.2	0.2	0.2	0.2	0.2
<code>min_child_samples</code>	100.0	50.0	50.0	100.0	100.0
<code>colsample_bytree</code>	0.7	0.7	0.85	0.7	0.7
<code>subsample</code>	1.0	0.7	1.0	1.0	1.0
RMSE	75,4	65,1	63,1	96,0	148,2

Tabela 25 – Hiperparâmetros Ajustados - *Random Forest*

Hiperparâmetro	Humanas	Natureza	Linguagem	Matemática	Redação
<code>max_depth</code>	20.0	20.0	20.0	20.0	20.0

*Continua na próxima página...*

Hiperparâmetro	Humanas	Natureza	Linguagem	Matemática	Redação
max_features	0.7	0.7	0.7	0.7	0.7
max_samples	0.8	0.8	0.8	0.8	0.8
RMSE	75,6	65,4	63,3	96,4	148,6

Fonte: elaborado pelo autor.

#### 4.5.2 Treinamento final dos modelos

Conforme descrito na Seção 3.5, o ajuste do hiperparâmetro `n_estimators` (número de estimadores) será realizado na próxima etapa, juntamente com o treinamento final dos modelos. Para os modelos de *XGBoost* e *LightGBM*, serão testados os valores de 100, 500, 1.000 e 3.000 estimadores. Para o modelo de *Random Forest*, só foi possível treinar o modelo com 100 e 500 estimadores, uma vez que a execução do código para as combinações com 1.000 e 3.000 estimadores foi interrompida subitamente algumas vezes, também possivelmente devido ao alto consumo de memória da GPU.

As Figuras 23 e 24 apresentam o erro MAPE do conjunto de teste, relativo ao maior erro MAPE encontrado para cada modelo e conjunto de dados, para os modelos de *XGBoost* e *LightGBM* com os melhores hiperparâmetros encontrados na etapa de ajuste dos hiperparâmetros, considerando o número de estimadores como 100, 500, 1.000 e 3.000. A tabela 26 apresenta o erro MAPE do conjunto de teste para o modelo de *Random Forest* com os melhores hiperparâmetros encontrados na etapa de ajuste dos hiperparâmetros, considerando o número de estimadores como 100 e 500.

Figura 23 – Erro MAPE - *XGBoost*

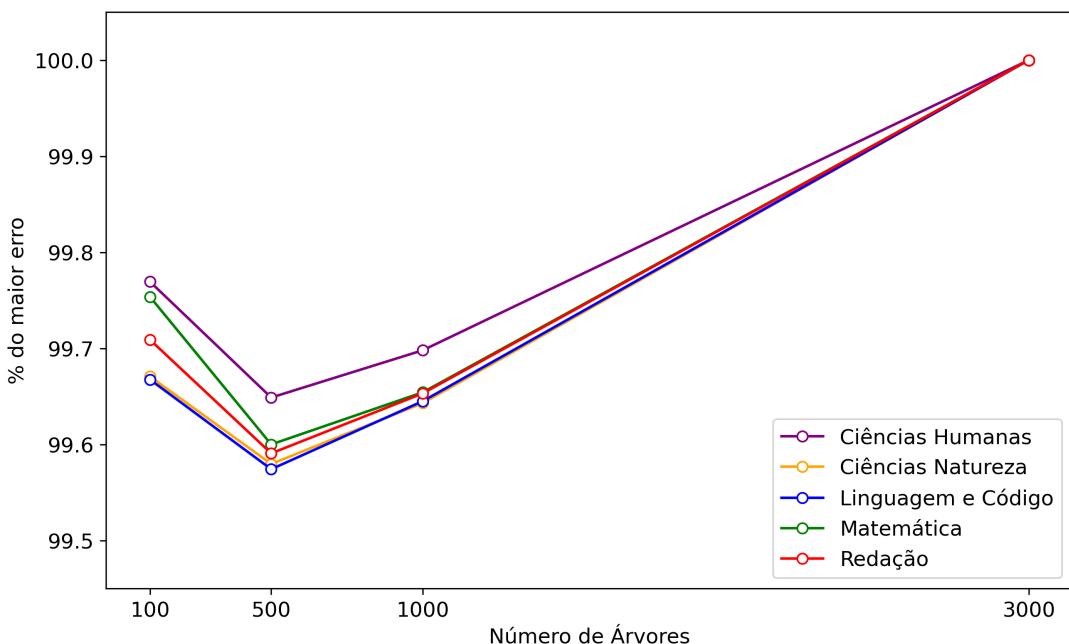
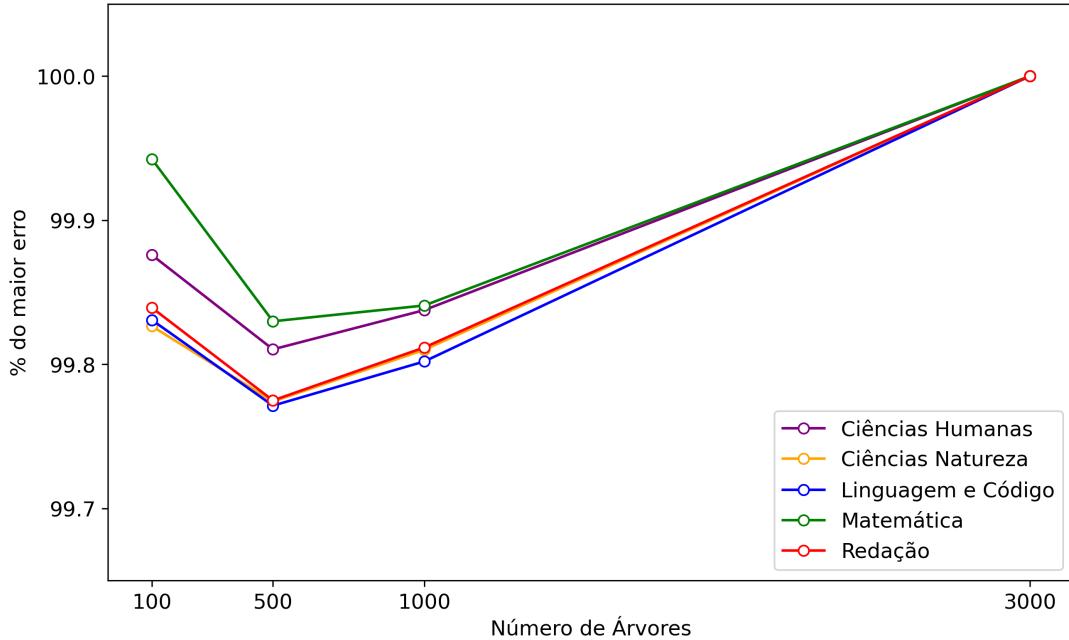


Figura 24 – Erro MAPE - *LightGBM*Tabela 26 – Erro MAPE - *Random Forest*

Área	100 estimadores	500 estimadores
Ciências Humanas	12.369%	12.364%
Ciências Natureza	10.924%	10.920%
Linguagem e Código	10.264%	10.260%
Matemática	15.160%	15.154%
Redação	21.327%	21.320%

Fonte: elaborado pelo autor.

É possível notar que assim como aconteceu no ajuste dos hiperparâmetros, a redução no erro MAPE para diferentes números de estimadores acabou não sendo tão significativa para os três modelos. Para o modelo de *XGBoost*, a maior redução no erro MAPE foi de 0.4258% na prova de Linguagem e Código. Para o modelo de *LightGBM*, a maior redução no erro MAPE foi de 0.2286% na prova de Linguagem e Código. Para o modelo de *Random Forest*, a maior redução no erro MAPE foi de 0.0068% na prova de Redação.

#### 4.5.3 Construção dos modelos de *ensemble*

A partir dos três modelos treinados, foram construídos dois modelos de *ensemble* usando a técnica de *bagging*: (i) um modelo de *ensemble* utilizando os modelos de *XGBoost* e *LightGBM* e (ii) um modelo de *ensemble* utilizando os modelos de *XGBoost*, *LightGBM* e *Random Forest*. Para ambos os modelos de *ensemble*, foi utilizada a média aritmética das previsões dos modelos individuais para obter a previsão final do modelo de *ensemble*.

As Figuras 25 e 26 apresentam o erro MAPE do conjunto de teste, relativo ao maior erro MAPE encontrado para cada modelo e conjunto de dados, para os modelos de *ensemble* construídos a partir dos modelos de *XGBoost*, *LightGBM* e *Random Forest* com os melhores hiperparâmetros encontrados na etapa de ajuste dos hiperparâmetros, considerando o número de estimadores como 100, 500, 1.000 e 3.000 para os modelos de *XGBoost* e *LightGBM* e 100 e 500 para o modelo de *Random Forest*.

Figura 25 – Erro MAPE - *Ensemble* (*XGBoost* + *LightGBM*)

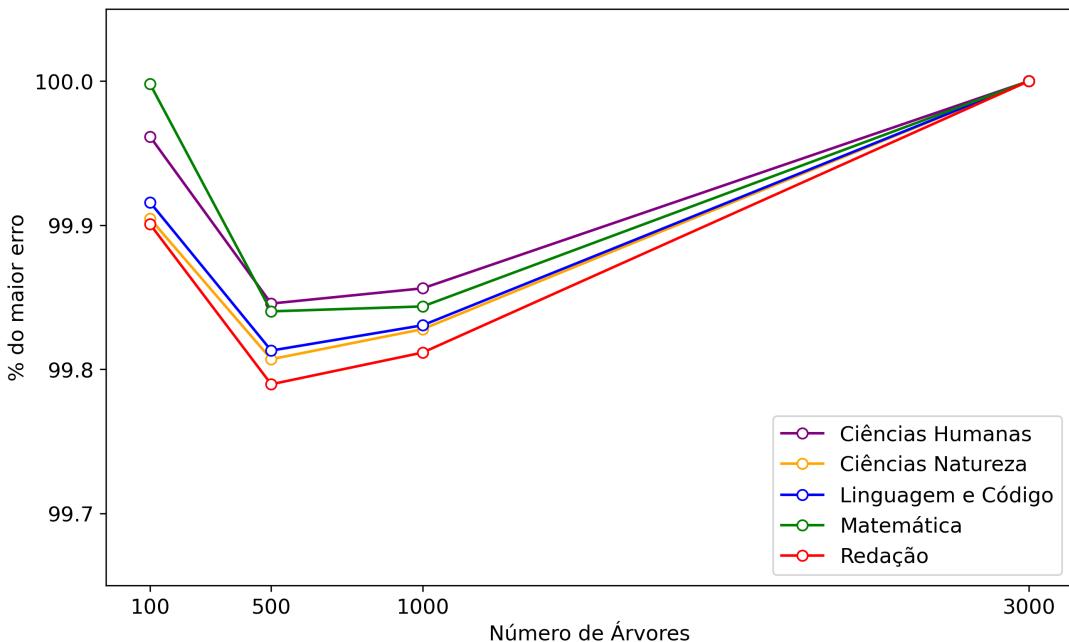
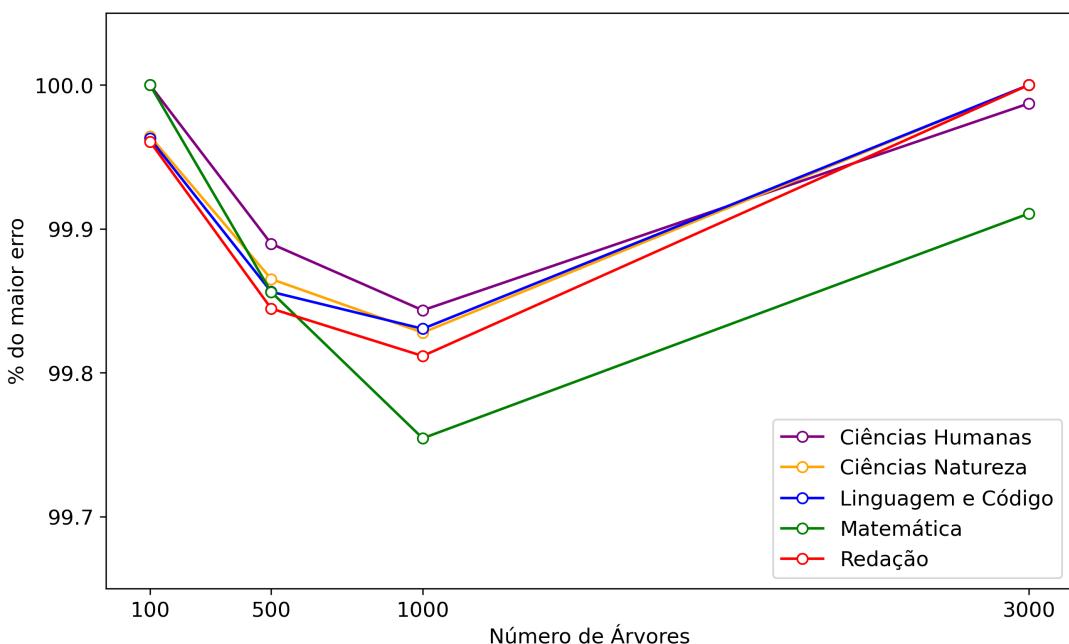


Figura 26 – Erro MAPE - *Ensemble* (*XGBoost* + *LightGBM* + *Random Forest*)



Fonte: elaborado pelo autor.

Novamente, a redução no erro MAPE para os modelos de *ensemble* acabou não sendo tão significativa para os diferentes números de estimadores. Para o modelo de *ensemble* utilizando os modelos de *XGBoost* e *LightGBM*, a maior redução no erro MAPE foi de 0.2103% na prova de Redação. Para o modelo de *ensemble* utilizando os modelos de *XGBoost*, *LightGBM* e *Random Forest*, a maior redução no erro MAPE foi de 0.1883% na prova de Redação.

#### 4.5.4 Avaliação dos modelos

Nessa etapa, a primeira análise realizada foi sobre um possível *overfitting* dos modelos. Conforme o critério adotado e explicado na Seção 3.6, nenhum dos modelos apresentou *overfitting* segundo esse critério, ou seja, nenhum dos modelos apresentou um erro do conjunto de teste 15% maior que o erro do conjunto de treino.

Em seguida, foi feita a comparação dos modelos entre si para cada conjunto de dados, a fim de se encontrar o melhor modelo para cada conjunto de dados. Para isso, foi adotado o critério do menor erro MAPE do conjunto de teste, ou seja, o modelo com o menor erro MAPE do conjunto de teste foi considerado o melhor modelo para aquele conjunto de dados. As Tabelas 27 a 31 apresentam os cinco modelos que apresentaram os menores erros MAPE do conjunto de teste para cada conjunto de dados, ordenados do menor para o maior erro MAPE.

Tabela 27 – Cinco melhores modelos - Humanas

<b>Modelo</b>	<b>Qtd. estimadores</b>	<b>MAPE teste</b>
XGB + LGBM	500	12.302290%
XGB + LGBM	1000	12.303594%
XGBoost	500	12.306396%
XGB + LGBM + RF	500	12.309288%
XGBoost	1000	12.312470%

Tabela 28 – Cinco melhores modelos - Natureza

<b>Modelo</b>	<b>Qtd. estimadores</b>	<b>MAPE teste</b>
XGB + LGBM	500	10.862886%
XGB + LGBM	1000	10.865162%
XGBoost	500	10.867356%
XGB + LGBM + RF	500	10.869213%
LightGBM	500	10.872770%

Tabela 29 – Cinco melhores modelos - Linguagem

<b>Modelo</b>	<b>Qtd. estimadores</b>	<b>MAPE teste</b>
XGB + LGBM	500	10.205704%
XGB + LGBM	1000	10.207500%
XGBoost	500	10.209998%
XGB + LGBM + RF	500	10.210127%
XGB + LGBM	100	10.216225%

Tabela 30 – Cinco melhores modelos - Matemática

<b>Modelo</b>	<b>Qtd. estimadores</b>	<b>MAPE teste</b>
XGB + LGBM	500	15.045981%
XGB + LGBM	1000	15.046495%
XGBoost	500	15.052611%
LightGBM	500	15.059630%
XGBoost	1000	15.060808%

Fonte: elaborado pelo autor.

Tabela 31 – Cinco melhores modelos - Redação

Modelo	Qtd. estimadores	MAPE teste
XGB + LGBM	500	21.208775%
XGB + LGBM	1000	21.213456%
XGBoost	500	21.214043%
LightGBM	500	21.219782%
XGB + LGBM + RF	500	21.220464%

Fonte: elaborado pelo autor.

Foram então escolhidos os seguintes modelos como os melhores para cada conjunto de dados, considerando o modelo com o menor erro MAPE do conjunto de teste:

- Humanas: modelo de *ensemble* utilizando os modelos de *XGBoost* e *LightGBM* com 500 estimadores;
- Natureza: modelo de *ensemble* utilizando os modelos de *XGBoost* e *LightGBM* com 500 estimadores;
- Linguagem e Código: modelo de *ensemble* utilizando os modelos de *XGBoost* e *LightGBM* com 500 estimadores;
- Matemática: modelo de *ensemble* utilizando os modelos de *XGBoost* e *LightGBM* com 500 estimadores;
- Redação: modelo de *ensemble* utilizando os modelos de *XGBoost* e *LightGBM* com 500 estimadores.

Essa escolha vai em linha com o apresentado até o momento, onde os modelos individuais apresentaram o menor erro MAPE com 500 estimadores e os modelos de *ensemble* apresentaram o menor erro MAPE em relação aos modelos individuais.

## 4.6 Influência das Variáveis Preditoras

### 4.6.1 Importância

Com o modelo final em mãos, foi possível analisar a importância das variáveis preditoras para cada modelo. Para isso, foi utilizado o método `feature_importances_` dos algoritmos para se extrair a importância de cada variável preditora para cada modelo. Como o modelo final é um modelo de *ensemble* utilizando os modelos de *XGBoost* e *LightGBM*, a importância de cada variável preditora para o modelo de *ensemble* foi calculada usando um método de ranking, onde a importância de cada variável preditora para o modelo de *ensemble* foi calculada como a média do ranking da importância de

cada variável preditora para os modelos de *XGBoost* e *LightGBM*. As Figuras 27 a 31 apresentam a importância das variáveis preditoras para cada modelo final, ordenadas da maior para a menor importância.

Figura 27 – Rank de Importância - Humanas

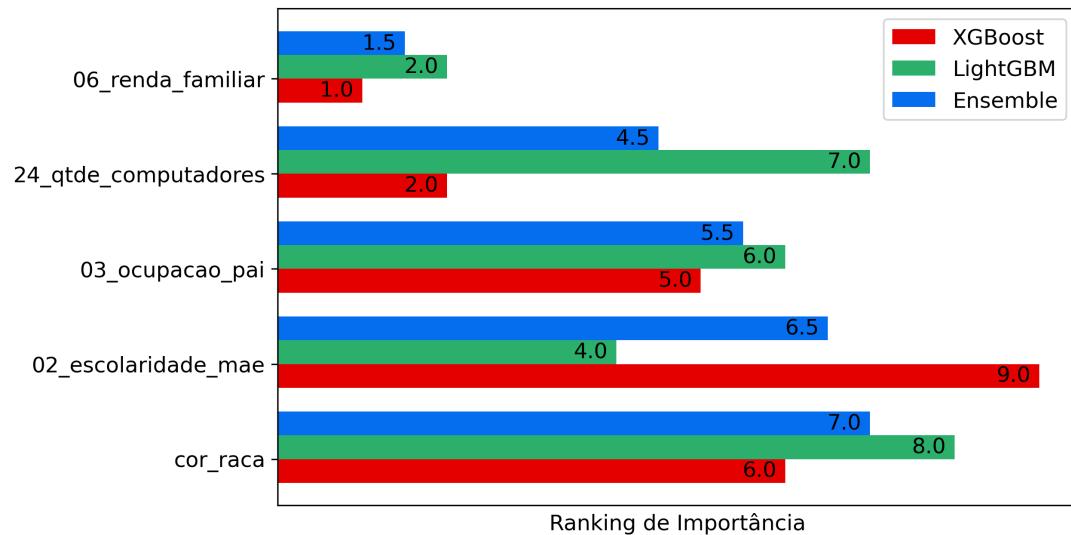


Figura 28 – Rank de Importância - Natureza

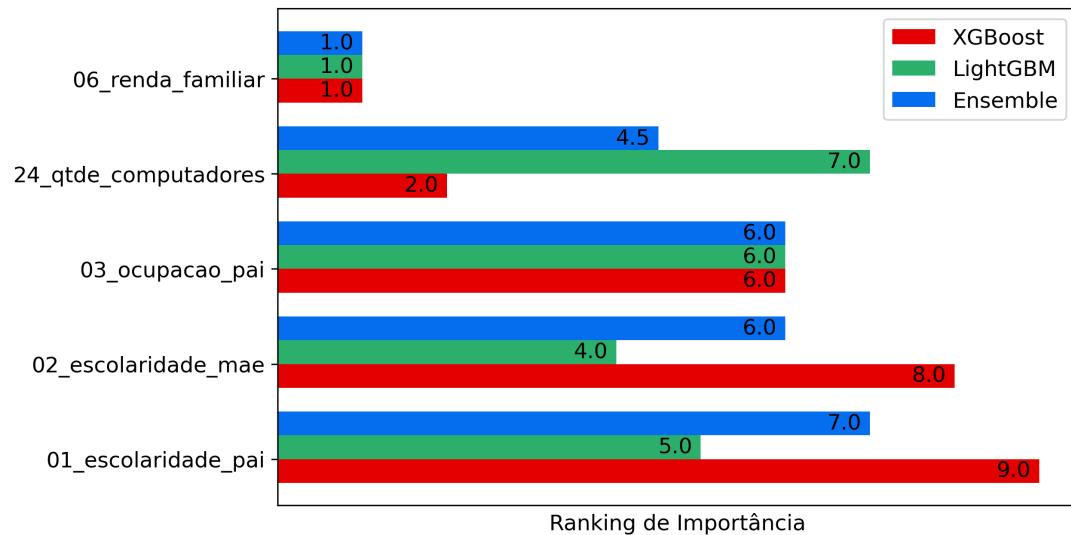


Figura 29 – Rank de Importância - Linguagem e Código

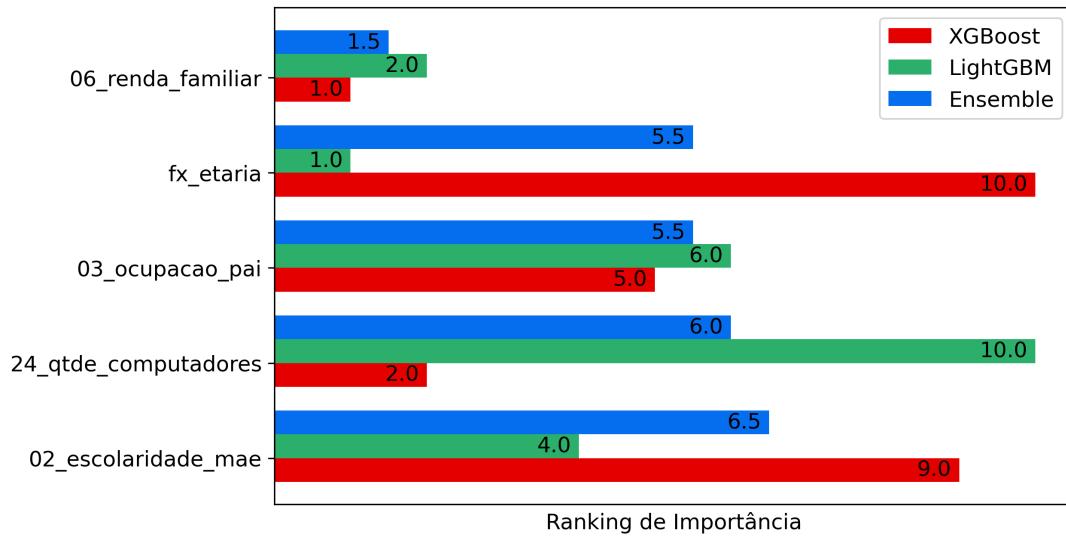


Figura 30 – Rank de Importância - Matemática

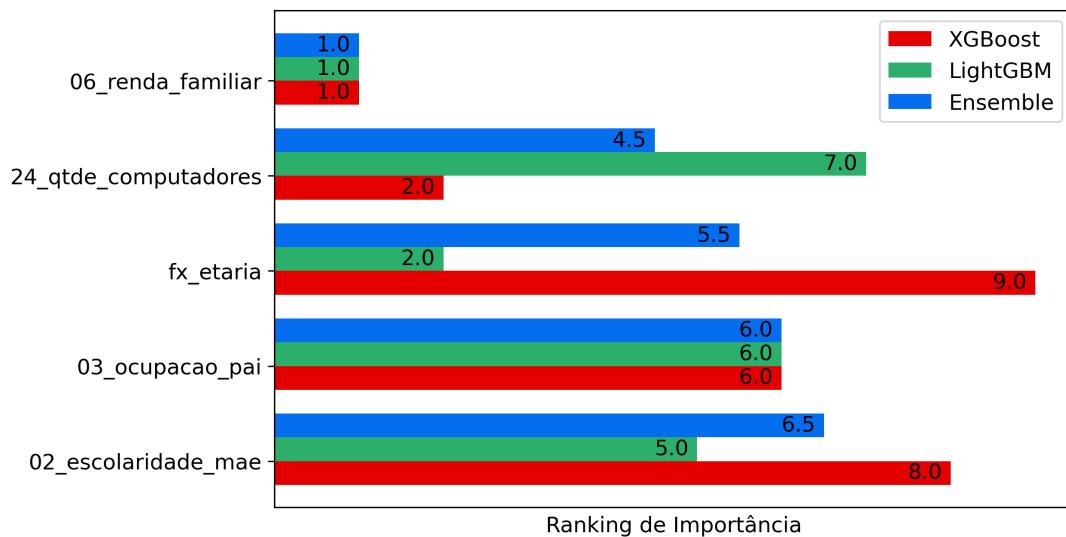
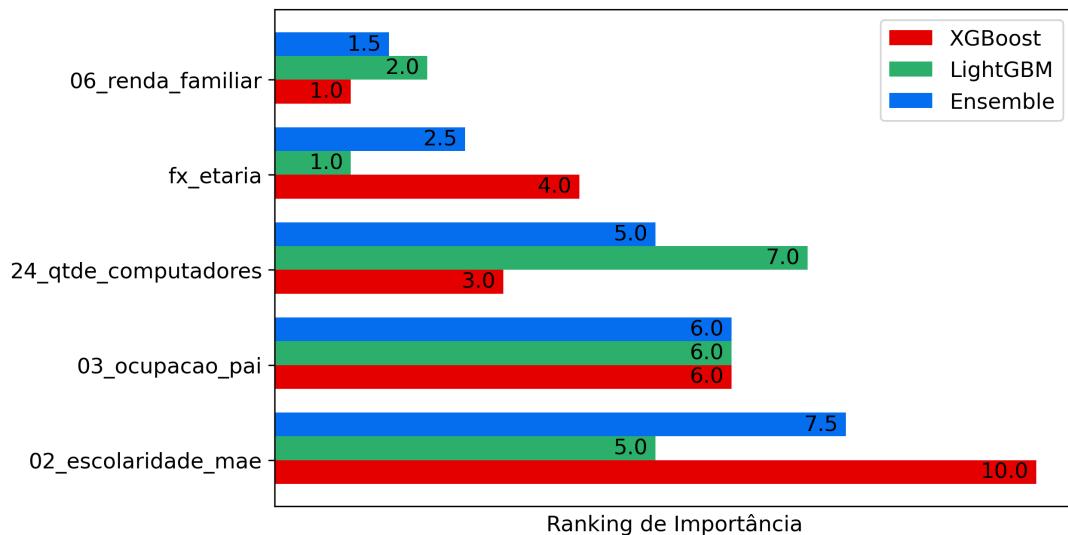


Figura 31 – Rank de Importância - Redação



Fonte: elaborado pelo autor.

As variáveis preditoras de renda familiar, quantidade de computadores em casa, ocupação do pai e escolaridade da mãe apareceram entre as 5 variáveis mais importantes para os cinco modelos finais. A variável da faixa etária do estudante apareceu para os modelos finais de Linguagem e Código, Matemática e Redação, mas não apareceu para os demais. A quinta variável para o modelo final de Ciências da Natureza foi a variável de escolaridade do pai e para o modelo final de Ciências Humanas foi a variável de cor/raça do estudante.

#### 4.6.2 Sensibilidade das Variáveis Respostas

Para se analisar a sensibilidade das variáveis respostas em relação às alterações nas variáveis preditoras, foi criada uma base sintética de dados, onde oito variáveis preditoras foram preenchidas com todos os seus possíveis valores e as demais 21 variáveis preditoras foram preenchidas com os seus valores mais frequentes. As variáveis preditoras selecionadas foram:

- **fx\_etaria**: faixa etária do estudante;
- **sexo**: sexo do estudante;
- **cor\_raca**: cor/raça do estudante;
- **01\_escolaridade\_pai**: escolaridade do pai do estudante;
- **02\_escolaridade\_mae**: escolaridade da mãe do estudante;
- **03\_ocupacao\_pai**: ocupação do pai do estudante;

- 04\_ocupacao\_mae: ocupação da mãe do estudante;
- 06\_renda\_familiar: renda familiar do estudante.

Foram selecionadas essas oito variáveis preditoras para a análise de sensibilidade por serem de interesse para a análise e por apresentarem uma quantidade razoável de valores distintos, o que possibilita uma análise mais detalhada da sensibilidade das variáveis respostas em relação às alterações nessas variáveis preditoras.

As figuras 32 a 39 apresentam as curvas de sensibilidade de cada Variável resposta em relação às alterações nas variáveis preditoras, onde cada curva representa a média das previsões do modelo final para cada valor da variável preditora, mantendo as demais variáveis preditoras fixas em seus valores mais frequentes.

As curvas de sensibilidade foram normalizadas ao primeiro ponto da curva, ou seja, o valor da curva para o primeiro ponto da variável preditora foi definido como 0 e os demais pontos foram calculados como o aumento percentual da curva para cada ponto em relação ao primeiro ponto.

Figura 32 – Curva de Sensibilidade - Faixa Etária

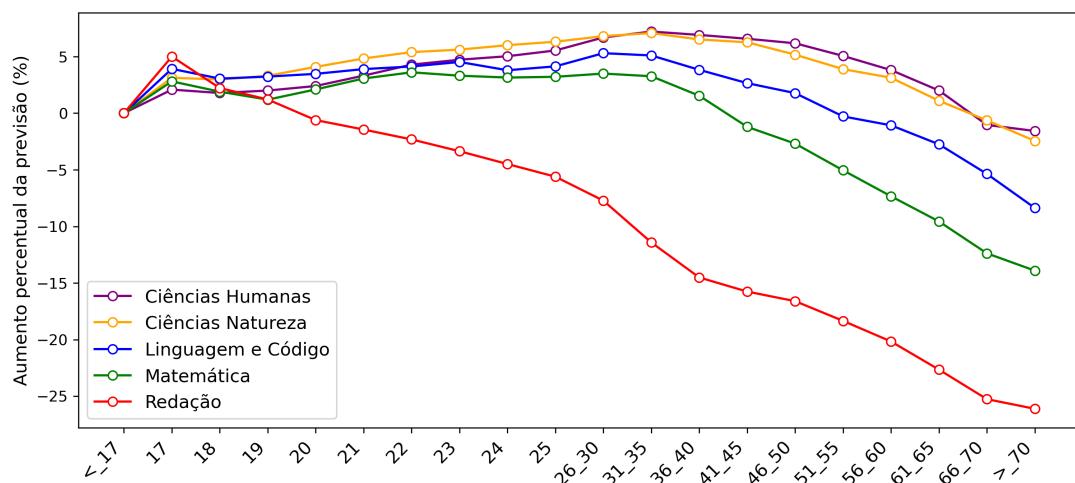


Figura 33 – Curva de Sensibilidade - Sexo

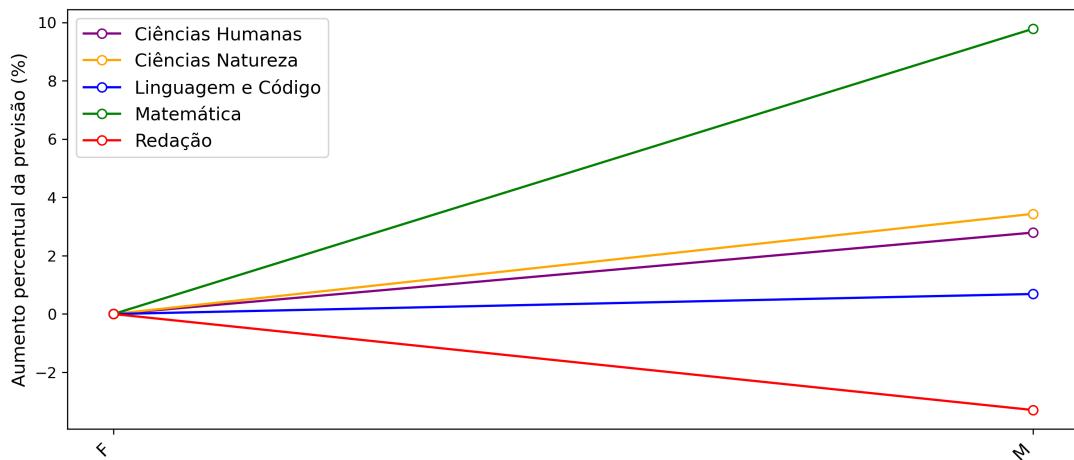


Figura 34 – Curva de Sensibilidade - Cor/Raça

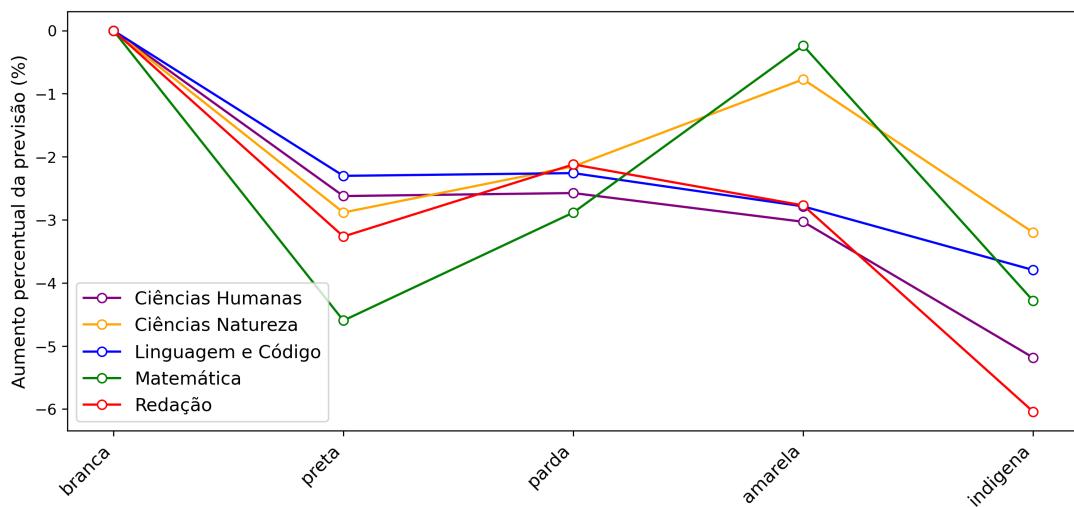


Figura 35 – Curva de Sensibilidade - Escolaridade do Pai

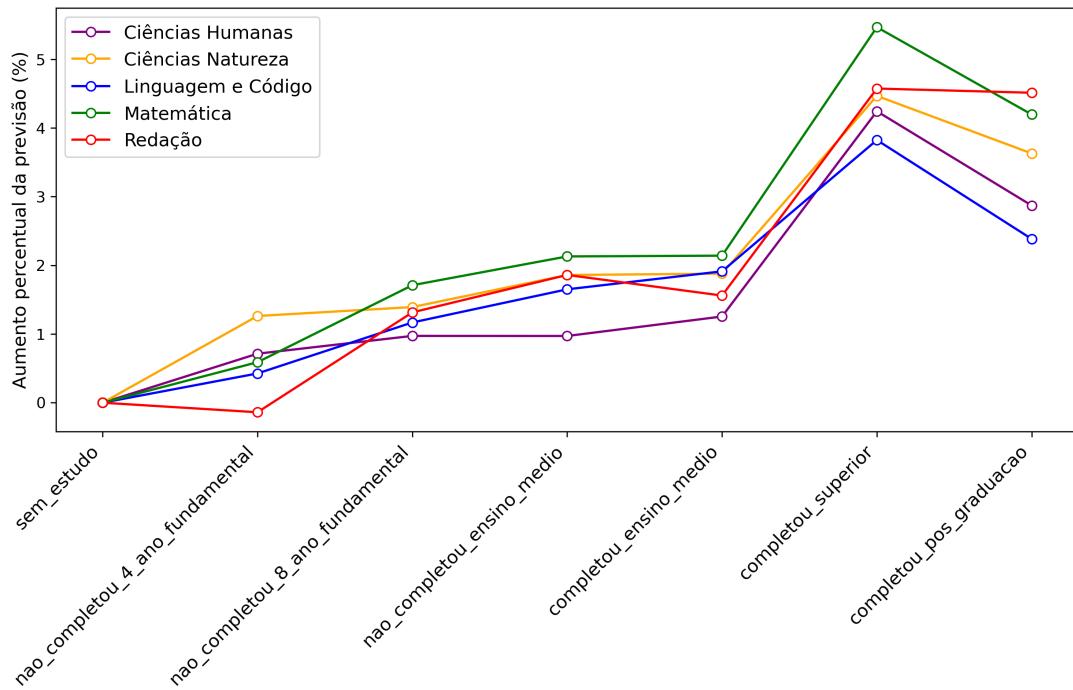


Figura 36 – Curva de Sensibilidade - Escolaridade da Mãe

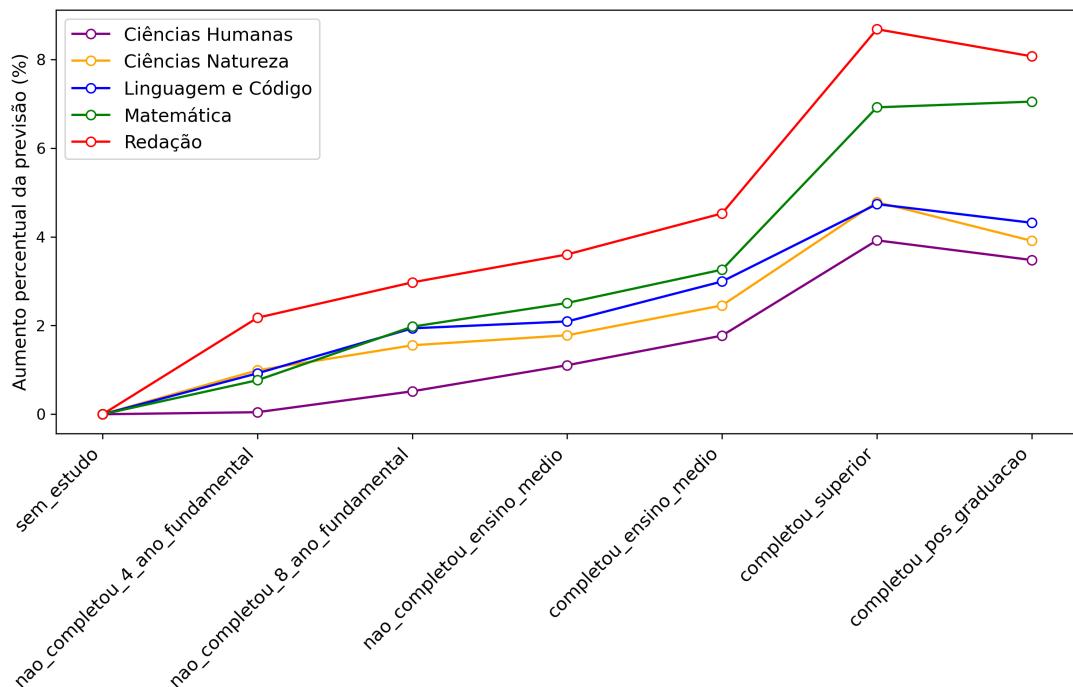


Figura 37 – Curva de Sensibilidade - Ocupação do Pai

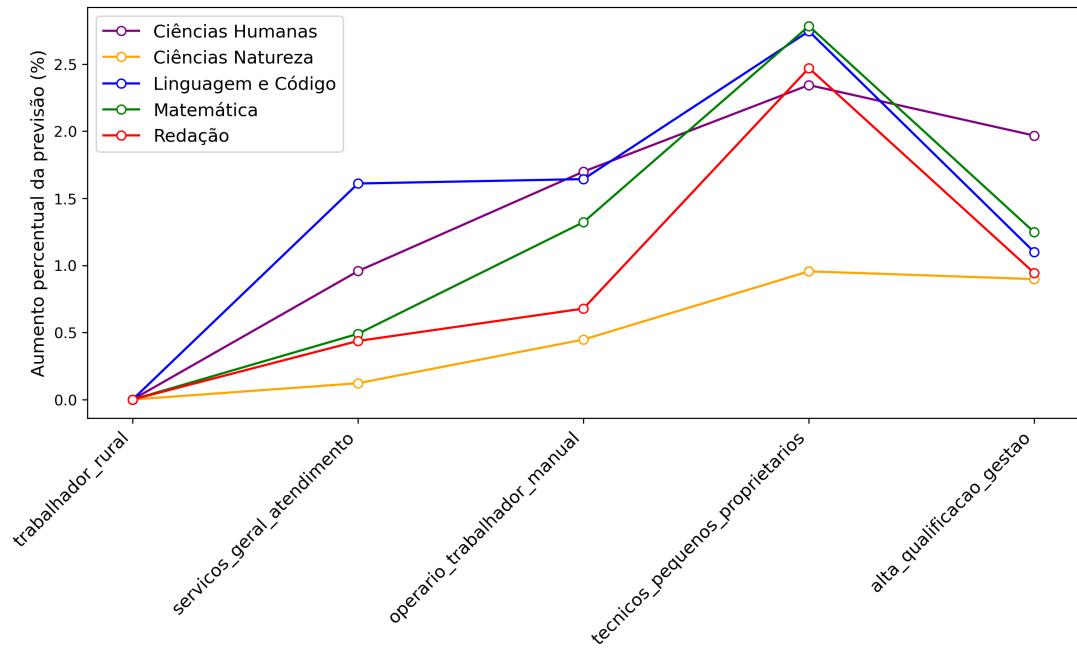


Figura 38 – Curva de Sensibilidade - Ocupação da Mãe

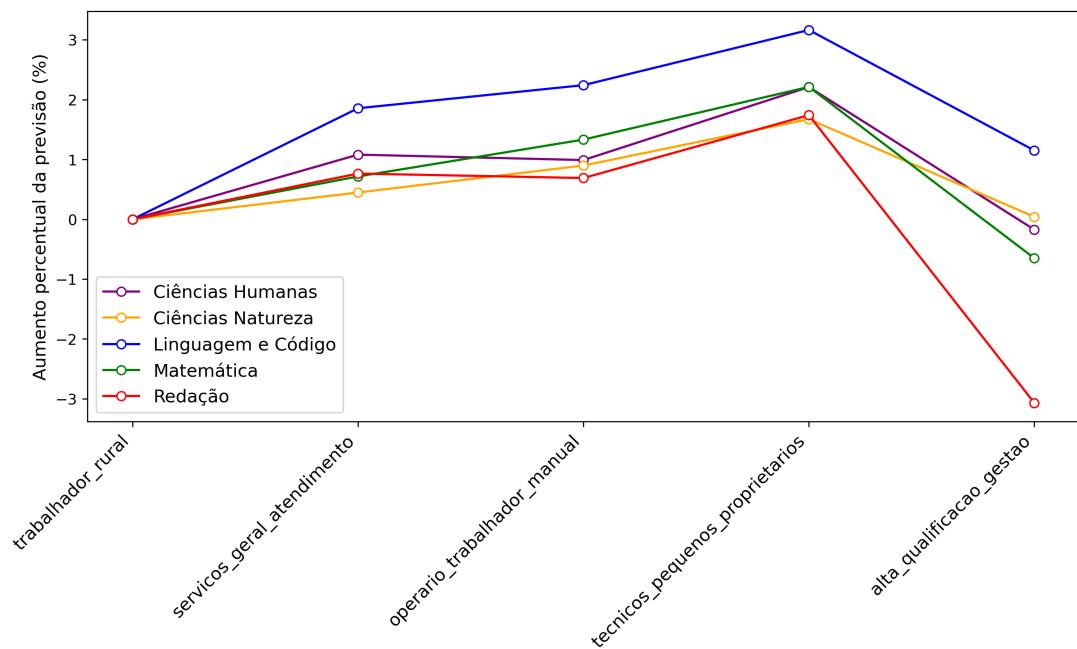
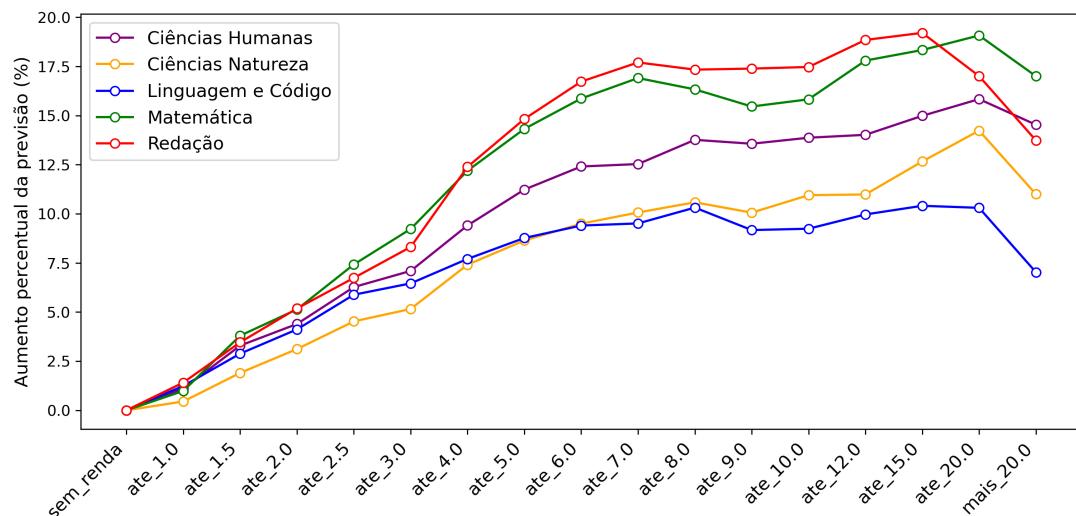


Figura 39 – Curva de Sensibilidade - Renda Familiar



Fonte: elaborado pelo autor.



## 5 CONCLUSÃO

XX



## REFERÊNCIAS

- 1 MELO, R. O. *et al.* Impacto das variáveis socioeconômicas no desempenho do ENEM: uma análise espacial e sociológica. **Revista de Administração Pública**, v. 55, n. 6, p. 1271–1294, nov./dez. 2021.
- 2 ORTEGA, A. *et al.* Análise comparativa: Escola pública x escola privada no ENEM. In: **Primeiro Hackthon de Dados pela Universidade Federal do ABC**. São Paulo: [S.l.: s.n.], 2025. Relatório.
- 3 NASCIMENTO, M. M. *et al.* Análise estatística e pluriescalar das desigualdades educacionais: aspirações científicas e desempenho de estudantes no ENEM. **Sociologias**, v. 27, 2025. Disponível em: <https://doi.org/10.1590/1807-0337/e130399>.
- 4 INEP. **Microdados ENEM**. Disponível em: <https://www.gov.br/inep/pt-br/acesso-a-informacao/dados-abertos/microdados/enem>.
- 5 INEP. **Histórico do ENEM**. Disponível em: <https://www.gov.br/inep/pt-br/areas-de-atuacao/avaliacao-e-exames-educacionais/enem/historico>.
- 6 INEP. **ENEM**. Disponível em: <https://www.gov.br/inep/pt-br/areas-de-atuacao/avaliacao-e-exames-educacionais/enem>.
- 7 OLIVEIRA, L. K. S.; CRUZ, R. C. Capital cultural e educação: uma análise da obra de bordieu. In: **XIII Encontro Cearense de Historiadores da Educação - ECHE, III Encontro Nacional do Núcleo de História e Memória da Educação - ENHIME, III Simpósio Nacional de Estudos Culturais e Geoeducacionais - SINECGEO**. [S.l.: s.n.], 2014. p. 1247–1255. ISBN: 978-85-8126-065-5. Documento de evento, sem data e local de publicação explícitos.
- 8 VASCONCELLOS, F. **Resultados do ENEM refletem desigualdades comuns no país**. 2013. Disponível em: <https://oglobo.globo.com/brasil/educacao/resultados-do-enem-refletem-desigualdades-comuns-no-pais-10445682>.
- 9 JALOTO, A.; PRIMI, R. Fatores socioeconômicos associados ao desempenho no ENEM. **Em Aberto**, v. 34, n. 112, p. 125–141, dec 2021. Disponível em: <https://www.researchgate.net/publication/357656960>.
- 10 MORAES, C. P. d. *et al.* Efeito escola a partir de indicadores educacionais: análise entre escolas públicas e privadas no ENEM. **REVISTA META: AVALIAÇÃO**, v. 14, n. 42, p. 67–93, mar 2022.
- 11 BARTHOLO, T. *et al.* **Oportunidades educacionais de estudantes concluintes do Ensino Médio: Relatório 1-Inscrição e Participação no ENEM entre 2013 e 2021**. Rio de Janeiro, 2023.
- 12 HIROMI, F. ENEM mais desigual requer atenção dos gestores. **Aprendizagem em Foco**, n. 92, oct 2023. Disponível em: <https://www.institutounibanco.org.br/boletim/enem-mais-desigual-requer-atencao-dos-gestores/>.

- 13 ROMERO, M. C. **Aplicando técnicas de Machine Learning para avaliar resultados do ENEM**. 2021. 72 p. Dissertação (Trabalho de Conclusão de Curso (MBA em Ciências de Dados)) — Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo, São Carlos, 2021.
- 14 FERRAZ, A. P. **Prevendo a aprovação de um participante do ENEM no SISU para o curso de Medicina**. 2020. 70 p. Dissertação (Trabalho de Conclusão de Curso (MBA em Ciências de Dados)) — Instituto de Ciências Matemáticas e de Computação, Universidade de São Paulo, São Carlos, 2020.
- 15 INEP. **Microdados Censo Escolar**. Local: Brasília, DF. [s.d.]. Disponível em: <https://www.gov.br/inep/pt-br/acesso-a-informacao/dados-abertos/microdados/censo-escolar>.
- 16 JAMES, G. *et al.* **An Introduction to Statistical Learning: with Applications in Python**. Boca Raton: CRC Press, 2023.
- 17 GRUS, J. **Data Science from Scratch: First Principles with Python**. 2. ed. Sebastopol, CA: O'Reilly Media, Inc., 2019.
- 18 LINDHOLM, A. *et al.* **Machine Learning: A First Course for Engineers and Scientists**. Cambridge, UK; New York, NY: Cambridge University Press, 2022.
- 19 BREIMAN, L. Random forests. **Machine Learning**, v. 45, n. 1, p. 5–32, oct 2001.
- 20 CHAPMAN, P. *et al.* **CRISP-DM 1.0: Step-by-step Data Mining Guide**. [S.l.], 2000. Disponível em: <https://mineracaodedados.wordpress.com/wp-content/uploads/2012/12/crisp-dm-1-0.pdf>.
- 21 BUSSAB, W. d. O.; MORETTIN, P. A. **Estatística Básica**. 9. ed. São Paulo: Saraiva, 2017.
- 22 COHEN, J. **Statistical Power Analysis for the Behavioral Sciences**. 2. ed. Hillsdale: Lawrence Erlbaum Associates, 1988.
- 23 KPMG Advisory N.V. **Phi\_K Correlation Constant**. 2024. Disponível em: <https://phik.readthedocs.io/en/latest/>.
- 24 BAAK, M. *et al.* A new correlation coefficient between categorical, ordinal and interval variables with Pearson characteristics. **Computational Statistics and Data Analysis**, v. 152, p. 107043, 2020. ISSN 0167-9473.
- 25 NVIDIA. **Welcome to cuML's documentation!** 2023. Disponível em: <https://docs.rapids.ai/api/cuml/stable/>.
- 26 PEDREGOSA, F. *et al.* Scikit-learn: Machine learning in Python. **Journal of Machine Learning Research**, v. 12, p. 2825–2830, 2011.
- 27 INEP. **Microdados do Enem 2024: Leia-Me**. Brasília, 2025. Disponível em: <https://www.gov.br/inep/pt-br/acesso-a-informacao/dados-abertos/microdados/enem>.
- 28 Anaconda. **Getting started with Miniconda**. n.d. Disponível em: <https://www.anaconda.com/docs/getting-started/miniconda/main>.

- 29 NVIDIA. **CUDA-X Data Science Libraries**. n.d. Disponível em: <https://developer.nvidia.com/topics/ai/data-science/cuda-x-data-science-libraries>.
- 30 NVIDIA. **Welcome to the cuDF documentation!** n.d. Disponível em: <https://docs.rapids.ai/api/cudf/stable/>.
- 31 The pandas development team. **pandas documentation**. 2025. Disponível em: <https://pandas.pydata.org/docs/>.



## **APÊNDICES**



**APÊNDICE A – DICIONÁRIO DE DADOS DOS MICRODADOS DO ENEM**



**APÊNDICE B – DICIONÁRIO DE DADOS DO CENSO ESCOLAR**



## APÊNDICE C – CONFIGURAÇÃO DO AMBIENTE VIRTUAL

```

name: kernel_2
channels:
- conda-forge
- rapidsai
- defaults
dependencies:
- __libgcc_mutex=0.1
- __openmp_mutex=4.5
- __python_abi3_support=1.0
- aiohappyeyeballs=2.6.1
- aiohttp=3.13.3
- aiosignal=1.4.0
- alsa-lib=1.2.15.3
- anyio=4.12.1
- aom=3.9.1
- argon2-cffi=25.1.0
- argon2-cffi-bindings=25.1.0
- arrow=1.4.0
- asttokens=3.0.1
- async-lru=2.1.0
- attr=2.5.2
- attrs=25.4.0
- aws-c-auth=0.9.3
- aws-c-cal=0.9.13
- aws-c-common=0.12.6
- aws-c-compression=0.3.1
- aws-c-event-stream=0.5.7
- aws-c-http=0.10.7
- aws-c-io=0.23.3
- aws-c-mqtt=0.13.3
- aws-c-s3=0.11.3
- aws-c-sdkutils=0.2.4
- aws-checksums=0.2.7
- aws-crt-cpp=0.35.4
- aws-sdk-cpp=1.11.606
- azure-core-cpp=1.16.1
- azure-identity-cpp=1.13.2
- azure-storage-blobs-cpp=12.16.0
- azure-storage-common-cpp=12.12.0
- azure-storage-files-datalake-cpp=12.14.0
- babel=2.17.0
- backports.zstd=1.3.0
- beautifulsoup4=4.14.3
- binutils_impl_linux-64=2.45
- binutils_linux-64=2.45
- bleach=6.3.0
- bleach-with-css=6.3.0
- blosc=1.21.6
- bokeh=3.8.2
- branca=0.8.2
- brotli=1.2.0
- brotli-bin=1.2.0
- brotli-python=1.2.0
- brunslis=0.1
- bzip2=1.0.8

```

```
- c-ares=1.34.6
- c-blosc2=2.23.0
- ca-certificates=2026.1.4
- cached-property=1.5.2
- cached_property=1.5.2
- cachetools=7.0.0
- cairo=1.18.4
- certifi=2026.1.4
- cffi=2.0.0
- charls=2.4.2
- charset-normalizer=3.4.4
- click=8.3.1
- cloudpickle=3.1.2
- cmake=4.2.3
- colorcet=3.1.0
- comm=0.2.3
- contourpy=1.3.3
- cpython=3.11.14
- cucim=25.12.00
- cuda-bindings=12.9.5
- cuda-cccl_linux-64=12.5.39
- cuda-core=0.3.2
- cuda-crt-dev_linux-64=12.5.82
- cuda-crt-tools=12.5.82
- cuda-cudart=12.5.82
- cuda-cudart-dev=12.5.82
- cuda-cudart-dev_linux-64=12.5.82
- cuda-cudart-static=12.5.82
- cuda-cudart-static_linux-64=12.5.82
- cuda-cudart_linux-64=12.5.82
- cuda-nvcc-dev_linux-64=12.5.82
- cuda-nvcc-impl=12.5.82
- cuda-nvcc-tools=12.5.82
- cuda-nvrtc=12.5.82
- cuda-nvvm-dev_linux-64=12.5.82
- cuda-nvvm-impl=12.5.82
- cuda-nvvm-tools=12.5.82
- cuda-pathfinder=1.3.3
- cuda-profiler-api=12.5.39
- cuda-python=12.9.5
- cuda-version=12.5
- cudf=25.12.00
- cudf-polars=25.12.00
- cudf_kafka=25.12.00
- cugraph=25.12.02
- cuml=25.12.00
- cupy=13.6.0
- cupy-core=13.6.0
- custreamz=25.12.00
- cuvs=25.12.00
- cuxfilter=25.12.00
- cycler=0.12.1
- cyrus-sasl=2.1.28
- cytoolz=1.1.0
- dask=2025.9.2
- dask-core=2025.9.2
- dask-cuda=25.12.00
- dask-cudf=25.12.00
- datashader=0.18.2
- dav1d=1.2.1
```

---

```
- dbus=1.16.2
- debugpy=1.8.20
- decorator=5.2.1
- defusedxml=0.7.1
- distributed=2025.9.2
- distributed-ucxx=0.47.00
- dlpack=0.8
- double-conversion=3.4.0
- exceptiongroup=1.3.1
- executing=2.2.1
- fastrllock=0.8.3
- folium=0.20.0
- font-ttf-dejavu-sans-mono=2.37
- font-ttf-inconsolata=3.000
- font-ttf-source-code-pro=2.038
- font-ttf-ubuntu=0.83
- fontconfig=2.15.0
- fonts-conda-ecosystem=1
- fonts-conda-forge=1
- fonttools=4.61.1
- fqdn=1.5.1
- freetype=2.14.1
- freexl=2.0.0
- frozenlist=1.7.0
- fsspec=2026.1.0
- gcc_impl_linux-64=13.4.0
- gcc_linux-64=13.4.0
- geopandas=1.1.2
- geopandas-base=1.1.2
- geos=3.14.1
- gflags=2.2.2
- giflib=5.2.2
- glog=0.7.1
- graphite2=1.3.14
- gxx_impl_linux-64=13.4.0
- gxx_linux-64=13.4.0
- h11=0.16.0
- h2=4.3.0
- harfbuzz=12.3.2
- holoviews=1.22.1
- hpack=4.1.0
- httpcore=1.0.9
- httpx=0.28.1
- hyperframe=6.1.0
- icu=78.2
- idna=3.11
- imagecodecs=2026.1.14
- imageio=2.37.0
- importlib-metadata=8.7.0
- ipykernel=7.1.0
- ipython=9.10.0
- ipython_pygments_lexers=1.1.1
- isoduration=20.11.0
- jedi=0.19.2
- jinja2=3.1.6
- joblib=1.5.3
- json-c=0.18
- json5=0.13.0
- jsonpointer=3.0.0
- jsonschema=4.26.0
```

```
- jsonschema-specifications=2025.9.1
- jsonschema-with-format-nongpl=4.26.0
- jupyter-lsp=2.3.0
- jupyter-server-proxy=4.4.0
- jupyter_client=8.8.0
- jupyter_core=5.9.1
- jupyter_events=0.12.0
- jupyter_server=2.17.0
- jupyter_server_terminals=0.5.4
- jupyterlab=4.5.3
- jupyterlab_pygments=0.3.0
- jupyterlab_server=2.28.0
- jxrlib=1.1
- kernel-headers_linux-64=6.12.0
- keyutils=1.6.3
- kiwisolver=1.4.9
- krb5=1.21.3
- lark=1.3.1
- lazy-loader=0.4
- lazy_loader=0.4
- lcms2=2.18
- ld_impl_linux-64=2.45
- lerc=4.0.0
- libabseil=20250512.1
- libaec=1.1.5
- libarchive=3.8.5
- libarrow=21.0.0
- libarrow-acero=21.0.0
- libarrow-compute=21.0.0
- libarrow-dataset=21.0.0
- libarrow-substrait=21.0.0
- libavif16=1.3.0
- libblas=3.11.0
- libbrotlicommon=1.2.0
- libbrotlidec=1.2.0
- libbrotlienc=1.2.0
- libcap=2.77
- libcblas=3.11.0
- libclang-cpp21.1=21.1.8
- libclang13=21.1.8
- libcrc32c=1.1.2
- libcublas=12.5.3.2
- libcublas-dev=12.5.3.2
- libcucim=25.12.00
- libcudf=25.12.00
- libcudf_kafka=25.12.00
- libcufft=11.2.3.61
- libcufile=1.10.1.7
- libcufile-dev=1.10.1.7
- libcugraph=25.12.02
- libcugraph_etl=25.12.02
- libcuml=25.12.00
- libcumlprims=25.12.00
- libcurl=8.18.0
- libcusolver=11.6.3.83
- libcusolver-dev=11.6.3.83
- libcusparse=12.5.1.3
```

---

- `libcusparse-dev=12.5.1.3`
- `libcuvs=25.12.00`
- `libcuvs-headers=25.12.00`
- `libdeflate=1.25`
- `libdrm=2.4.125`
- `libedit=3.1.20250104`
- `libegl=1.7.0`
- `libev=4.33`
- `libevent=2.1.12`
- `libexpat=2.7.3`
- `libffi=3.5.2`
- `libfreetype=2.14.1`
- `libfreetype6=2.14.1`
- `libgcc=15.2.0`
- `libgcc-devel_linux-64=13.4.0`
- `libgcc-ng=15.2.0`
- `libgdal-core=3.12.1`
- `libgfortran=15.2.0`
- `libgfortran5=15.2.0`
- `libgl=1.7.0`
- `libglib=2.86.3`
- `libglvnd=1.7.0`
- `libglx=1.7.0`
- `libgomp=15.2.0`
- `libgoogle-cloud=2.39.0`
- `libgoogle-cloud-storage=2.39.0`
- `libgrpc=1.73.1`
- `libhwy=1.3.0`
- `libiconv=1.18`
- `libjpeg-turbo=3.1.2`
- `libjxl=0.11.1`
- `libkml=1.3.0`
- `libkvikio=25.12.00`
- `liblapack=3.11.0`
- `libllvm21=21.1.8`
- `liblzma=5.8.2`
- `libnghttp2=1.67.0`
- `libnl=3.11.0`
- `libnsl=2.0.1`
- `libntlm=1.8`
- `libnuma=2.0.18`
- `libnvcomp=5.0.0.6`
- `libnvcomp-dev=5.0.0.6`
- `libnvimimgcodec=0.6.0`
- `libnvimimgcodec0=0.6.0`
- `libnvjitlink=12.9.86`
- `libnvjpeg=12.3.2.81`
- `libnvjpeg2k0=0.9.0.43`
- `libnvtiff=0.5.1.75`
- `libnvtiff0=0.5.1.75`
- `libopenblas=0.3.30`
- `libopengl=1.7.0`
- `libopentelemetry-cpp=1.21.0`
- `libopentelemetry-cpp-headers=1.21.0`
- `libparquet=21.0.0`
- `libpciaccess=0.18`
- `libpng=1.6.54`
- `libpq=18.1`
- `libprotobuf=6.31.1`
- `libraft=25.12.00`

```
- libraft-headers=25.12.00
- libraft-headers-only=25.12.00
- librdkafka=2.8.0
- libre2-11=2025.11.05
- librmm=25.12.00
- librttopo=1.1.0
- libsanitizer=13.4.0
- libsodium=1.0.20
- libspatialite=5.1.0
- libsqlite=3.51.2
- libssh2=1.11.1
- libstdcxx=15.2.0
- libstdcxx-devel_linux-64=13.4.0
- libstdcxx-ng=15.2.0
- libsystemd0=258.3
- libthrift=0.22.0
- libtiff=4.7.1
- libucxx=0.47.00
- libudev1=258.3
- libutf8proc=2.11.3
- libuuid=2.41.3
- libuv=1.51.0
- libvulkan-loader=1.4.328.1
- libwebp-base=1.6.0
- libxcb=1.17.0
- libxcrypt=4.4.36
- libxgboost=3.1.2
- libxkbcommon=1.13.1
- libxml2=2.15.1
- libxml2-16=2.15.1
- libxml2-devel=2.15.1
- libxslt=1.1.43
- libzlib=1.3.1
- libzopfli=1.0.3
- linkify-it-py=2.0.3
- llvmlite=0.44.0
- locket=1.0.0
- lz4=4.4.5
- lz4-c=1.10.0
- lzo=2.10
- make=4.4.1
- mapclassify=2.10.0
- markdown=3.10.1
- markdown-it-py=4.0.0
- markupsafe=3.0.3
- matplotlib=3.10.8
- matplotlib-base=3.10.8
- matplotlib-inline=0.2.1
- mdit-py-plugins=0.5.0
- mdurl=0.1.2
- minizip=4.0.10
- mistune=3.2.0
- msgpack-python=1.1.2
- multidict=6.7.0
- multipledispatch=0.6.0
- munkres=1.1.4
- muparser=2.3.5
- narwhals=2.16.0
- nbclient=0.10.4
- nbconvert-core=7.17.0
```

---

```
- nbformat=5.10.4
- nccl=2.29.2.1
- ncurses=6.5
- nest-asyncio=1.6.0
- networkx=3.6.1
- nlohmann_json=3.12.0
- nodejs=25.2.1
- notebook-shim=0.2.4
- numba=0.61.2
- numba-cuda=0.19.2
- numpy=2.2.6
- nvidia-ml-py=13.590.48
- nvtx=0.2.14
- nx-cugraph=25.12.00
- openjpeg=2.5.4
- openjph=0.26.0
- openldap=2.6.10
- openssl=3.6.1
- orc=2.2.2
- overrides=7.7.0
- packaging=26.0
- pandas=2.3.3
- pandocfilters=1.5.0
- panel=1.8.7
- param=2.3.1
- parso=0.8.5
- partd=1.4.2
- patsy=1.0.2
- pcre2=10.47
- pexpect=4.9.0
- phik=0.12.5
- pillow=12.1.0
- pip=26.0
- pixman=0.46.4
- platformdirs=4.5.1
- polars=1.34.0
- polars-runtime-32=1.34.0
- proj=9.7.1
- prometheus-cpp=1.3.0
- prometheus_client=0.24.1
- prompt-toolkit=3.0.52
- propcache=0.3.1
- psutil=7.2.2
- pthread-stubs=0.4
- ptyprocess=0.7.0
- pure_eval=0.2.3
- py-xgboost=3.1.2
- pyarrow=21.0.0
- pyarrow-core=21.0.0
- pycparser=2.22
- pyct=0.6.0
- pygments=2.19.2
- pylibcudf=25.12.00
- pylibcugraph=25.12.02
- pylibraft=25.12.00
- pyogrio=0.12.1
- pyparsing=3.3.2
- pyproj=3.7.2
- pyside6=6.10.1
- pysocks=1.7.1
```

```
- python=3.11.14
- python-confluent-kafka=2.8.0
- python-dateutil=2.9.0.post0
- python-fastjsonschema=2.21.2
- python-gil=3.11.14
- python-json-logger=2.0.7
- python-tzdata=2025.3
- python_abi=3.11
- pytz=2025.2
- pyviz_comms=3.0.6
- pywavelets=1.9.0
- pyyaml=6.0.3
- pyzmq=27.1.0
- qhull=2020.2
- qt6-main=6.10.1
- raft-dask=25.12.00
- rapids=25.12.00
- rapids-dask-dependency=25.12.01
- rapids-logger=0.2.3
- rapids-xgboost=25.12.00
- rav1e=0.7.1
- rdma-core=61.0
- re2=2025.11.05
- readline=8.3
- referencing=0.37.0
- requests=2.32.5
- rfc3339-validator=0.1.4
- rfc3986-validator=0.1.1
- rfc3987-syntax=1.1.0
- rhash=1.4.6
- rich=14.3.2
- rmm=25.12.00
- rpds-py=0.30.0
- s2n=1.6.2
- scikit-image=0.24.0
- scikit-learn=1.8.0
- scipy=1.16.3
- seaborn=0.13.2
- seaborn-base=0.13.2
- send2trash=2.1.0
- setuptools=80.10.2
- shapely=2.1.2
- supervisor=1.0.0
- six=1.17.0
- snappy=1.2.2
- sniffio=1.3.1
- sortedcontainers=2.4.0
- soupsieve=2.8.3
- sqlite=3.51.2
- stack_data=0.6.3
- statsmodels=0.14.6
- streamz=0.6.5
- svt-av1=4.0.0
- sysroot_linux-64=2.39
- tblib=3.2.2
- terminado=0.18.1
- threadpoolctl=3.6.0
- tifffile=2026.1.28
- tinyccs2=1.5.1
- tk=8.6.13
```

---

```
- tomli=2.4.0
- toolz=1.1.0
- tornado=6.5.4
- tqdm=4.67.2
- traitlets=5.14.3
- treeelite=4.6.1
- typing-extensions=4.15.0
- typing_extensions=4.15.0
- typing_utils=0.1.0
- tzdata=2025c
- uc-micro-py=1.0.3
- ucx=1.19.1
- ucxx=0.47.00
- unicodedata2=17.0.0
- uri-template=1.3.0
- uriparser=0.9.8
- urllib3=2.6.3
- wayland=1.24.0
- wcwidth=0.5.3
- webcolors=25.10.0
- webencodings=0.5.1
- websocket-client=1.9.0
- wheel=0.46.3
- xarray=2026.1.0
- xcb-util=0.4.1
- xcb-util-cursor=0.1.6
- xcb-util-image=0.4.0
- xcb-util-keysyms=0.4.1
- xcb-util-renderutil=0.3.10
- xcb-util-wm=0.4.2
- xerces-c=3.3.0
- xgboost=3.1.2
- xkeyboard-config=2.46
- xorg-libice=1.1.2
- xorg-libsm=1.2.6
- xorg-libx11=1.8.12
- xorg-libxau=1.0.12
- xorg-libxcomposite=0.4.7
- xorg-libxcursor=1.2.3
- xorg-libxdamage=1.1.6
- xorg-libxdmcp=1.1.5
- xorg-libxext=1.3.7
- xorg-libxfixes=6.0.2
- xorg-libxi=1.8.2
- xorg-libxrandr=1.5.5
- xorg-libxrender=0.9.12
- xorg-libxtst=1.2.5
- xorg-libxxf86vm=1.1.7
- xyzservices=2025.11.0
- yaml=0.2.5
- yarl=1.22.0
- zeromq=4.3.5
- zfp=1.0.1
- zict=3.0.0
- zipp=3.23.0
- zlib=1.3.1
- zlib-ng=2.3.2
- zstd=1.5.7
- pip:
    - duckdb==1.4.4
```

```
- lightgbm==4.6.0
- mpmath==1.3.0
- sympy==1.14.0
prefix: /home/ramon/miniconda3/envs/kernel_2
```