



INSTITUTO LATINO-AMERICANO DE  
CIÊNCIAS DA VIDA E DA NATUREZA  
(ILACVN)  
PROGRAMA DE PÓS-GRADUAÇÃO EM  
FÍSICA APLICADA

UNIVERSIDADE FEDERAL DA INTEGRAÇÃO  
LATINO-AMERICANA - UNILA

PROJETO APRESENTADO NA DISCIPLINA: TÓPICOS  
ESPECIAIS EM FÍSICA EXPERIMENTAL

---

## Projeto I: Previsão Preditiva da Severidade do Câncer

---

**Discente:**

Fernando José Zardinello Batistti

**Orientador:**

Prof. Dr. Joylan Nunes Maciel

Foz do Iguaçu  
2025



# Resumo

Este trabalho apresenta o desenvolvimento de um modelo de *Machine Learning* para prever a pontuação de severidade do câncer com base em dados multifatoriais de pacientes. Utilizando o algoritmo *RandomForestRegressor*, um *pipeline* de pré-processamento foi construído, incluindo a normalização de dados numéricos e a codificação de variáveis categóricas. O modelo foi treinado, avaliado e otimizado através da técnica de *RandomizedSearchCV*. O modelo final otimizado alcançou um coeficiente de determinação ( $R^2$ ) de 0,979, demonstrando alta capacidade preditiva e validando a eficácia da abordagem para apoiar a tomada de decisão clínica.

**Palavras-chave:** Aprendizado de Máquina; Regressão; Random Forest; Previsão de Câncer; Ciência de Dados.

# Abstract

This work presents the development of a *Machine Learning* model to predict cancer severity scores based on multifactorial patient data. Using the *RandomForestRegressor* algorithm, a preprocessing *pipeline* was built, including the normalization of numerical data and the encoding of categorical variables. The model was trained, evaluated, and optimized using the *RandomizedSearchCV* technique. The final optimized model achieved a coefficient of determination ( $R^2$ ) of 0,979, demonstrating high predictive capability and validating the effectiveness of the approach in supporting clinical decision-making.

**Keywords:** Machine Learning; Regression; Random Forest; Cancer Prediction; Data Science.

# Sumário

<b>Resumo</b>	<b>i</b>
<b>Abstract</b>	<b>ii</b>
<b>1 Introdução</b>	<b>1</b>
<b>2 Fundamentação Teórica</b>	<b>2</b>
2.1 Aprendizado de Máquina Supervisionado e Regressão . . . . .	2
2.2 Random Forest Regressor . . . . .	2
2.3 Métricas de Avaliação para Regressão . . . . .	2
2.4 Otimização de Hiperparâmetros . . . . .	3
<b>3 Desenvolvimento</b>	<b>4</b>
3.1 Coleta e Análise Exploratória dos Dados . . . . .	4
3.2 Pré-processamento e Construção do <i>Pipeline</i> . . . . .	7
3.3 Treinamento e Otimização do Modelo . . . . .	7
<b>4 Resultados e Discussões</b>	<b>8</b>
4.1 Desempenho dos Modelos . . . . .	8
4.2 Análise Comparativa . . . . .	8
<b>5 Conclusão e Trabalhos Futuros</b>	<b>10</b>
<b>Referências Bibliográficas</b>	<b>11</b>

# Lista de Figuras

3.1	Distribuição da variável alvo (Target Severity Score). . . . .	4
3.2	Relação entre o Estágio do Câncer e a Pontuação de Severidade. . . . .	5
3.3	Matriz de correlação entre as variáveis numéricas do dataset. . . . .	6
3.4	Curva de aprendizado do modelo base. . . . .	7
4.1	Gráfico de dispersão dos valores reais vs. previstos para o modelo otimizado. . . . .	9

# Lista de Tabelas

4.1	Comparação de desempenho entre o modelo base e o otimizado. . . . .	8
-----	---	---

# Lista de Códigos

1	Melhores hiperparâmetros encontrados pela busca aleatória. . . . .	8
---	--	---

# 1. Introdução

O câncer representa um dos desafios mais proeminentes para a saúde pública global, mantendo-se como uma das principais causas de morbidade e mortalidade em todo o mundo, com uma estimativa de 20 milhões de novos casos e 9,7 milhões de óbitos apenas em 2022 [6]. O avanço no entendimento da biologia do câncer tem revelado sua natureza heterogênea, onde a progressão e o prognóstico da doença são influenciados por uma complexa interação de fatores demográficos, genéticos, ambientais e de estilo de vida. Essa multifatorialidade torna a avaliação da severidade de cada caso uma tarefa extremamente desafiadora para a prática clínica, limitando a capacidade de personalização dos tratamentos e de alocação eficiente dos recursos de saúde.

Diante dessa complexidade, as abordagens tradicionais de análise podem ser insuficientes para capturar os padrões não lineares e as interdependências sutis entre as variáveis preditoras. É neste cenário que a Inteligência Artificial, especificamente o Aprendizado de Máquina (*Machine Learning*), surge como uma ferramenta poderosa. Modelos de *Machine Learning* são capazes de aprender a partir de grandes volumes de dados, identificando padrões complexos e gerando previsões que podem apoiar de forma significativa a tomada de decisão clínica.

Este trabalho se propõe a explorar essa oportunidade, utilizando um algoritmo clássico de *Machine Learning* para desenvolver um modelo preditivo robusto. A questão central que norteia esta pesquisa é: é possível, utilizando um conjunto de dados multifatoriais de pacientes, construir um modelo de *Machine Learning* com alta acurácia para prever a pontuação de severidade do câncer?

O objetivo principal deste projeto é, portanto, desenvolver, treinar e avaliar um modelo de regressão para prever a pontuação de severidade do câncer (**Target\_Severity\_Score**) com base em características de pacientes. Um modelo bem-sucedido tem o potencial de oferecer um impacto prático substancial, auxiliando na estratificação de risco dos pacientes, na otimização de planos terapêuticos e na gestão estratégica de recursos hospitalares, contribuindo assim para o avanço da medicina de precisão.

Para atingir tal objetivo, este relatório está estruturado da seguinte forma: a Fundamentação Teórica detalhará os conceitos de *Machine Learning* e os algoritmos utilizados; a seção de Desenvolvimento descreverá o *pipeline* completo, desde a coleta de dados até a implementação do modelo; os Resultados apresentarão as métricas de performance; e, por fim, a Conclusão discutirá os achados, desafios e direções futuras do trabalho.



## 2. Fundamentação Teórica

Esta seção apresenta os conceitos fundamentais de Aprendizado de Máquina que foram aplicados no desenvolvimento deste projeto, incluindo a definição da tarefa de regressão, o funcionamento do algoritmo *Random Forest*, as métricas de avaliação e a metodologia de otimização de hiperparâmetros.

### 2.1 Aprendizado de Máquina Supervisionado e Regressão

O *Machine Learning* é um subcampo da Inteligência Artificial que se dedica ao desenvolvimento e estudo de algoritmos capazes de aprender a partir de dados [3]. Dentre suas vertentes, o Aprendizado Supervisionado é utilizado quando se dispõe de um conjunto de dados rotulado, ou seja, um conjunto de exemplos de entrada (*features*) acompanhados de suas respectivas saídas corretas (alvo ou *target*).

O objetivo de um modelo supervisionado é aprender uma função de mapeamento  $f$  que possa, a partir de um novo vetor de *features*  $X$ , prever sua saída  $\hat{y}$  da forma mais acurada possível, tal que  $\hat{y} \approx f(X)$ . O presente projeto se enquadra como uma tarefa de regressão, pois o objetivo é prever a `Target_Severity_Score`, uma variável contínua que quantifica a severidade do câncer.

### 2.2 Random Forest Regressor

O algoritmo escolhido para este trabalho foi o *Random Forest* (Floresta Aleatória), um dos modelos de *ensemble* mais robustos e populares, conforme descrito por Breiman [1]. Suas unidades fundamentais são as Árvores de Decisão, que são combinadas para mitigar o problema de sobreajuste (*overfitting*) através de técnicas como *Bagging* e *Random Subspace Method*. A previsão final de um `RandomForestRegressor` é a média das previsões de todas as árvores individuais que compõem a floresta.

### 2.3 Métricas de Avaliação para Regressão

Para avaliar quantitativamente o desempenho do modelo de regressão, foram utilizadas as seguintes métricas, onde  $y_i$  é o valor real,  $\hat{y}_i$  é o valor previsto e  $\bar{y}$  é a média dos valores reais:

- *Mean Absolute Error* (MAE): Mede a média dos erros absolutos.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

- *Root Mean Squared Error* (RMSE): É a raiz quadrada do erro quadrático médio, penalizando erros maiores de forma mais significativa.

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

- *R-squared* ( $R^2$ ): O Coeficiente de Determinação mede a proporção da variância no alvo que é previsível a partir das *features*.

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

## 2.4 Otimização de Hiperparâmetros

Os hiperparâmetros são configurações externas de um modelo definidas antes do treinamento. Neste projeto, foi utilizada a técnica *RandomizedSearchCV*, disponível na biblioteca *Scikit-learn* [5], para buscar de forma eficiente a combinação ótima de hiperparâmetros.

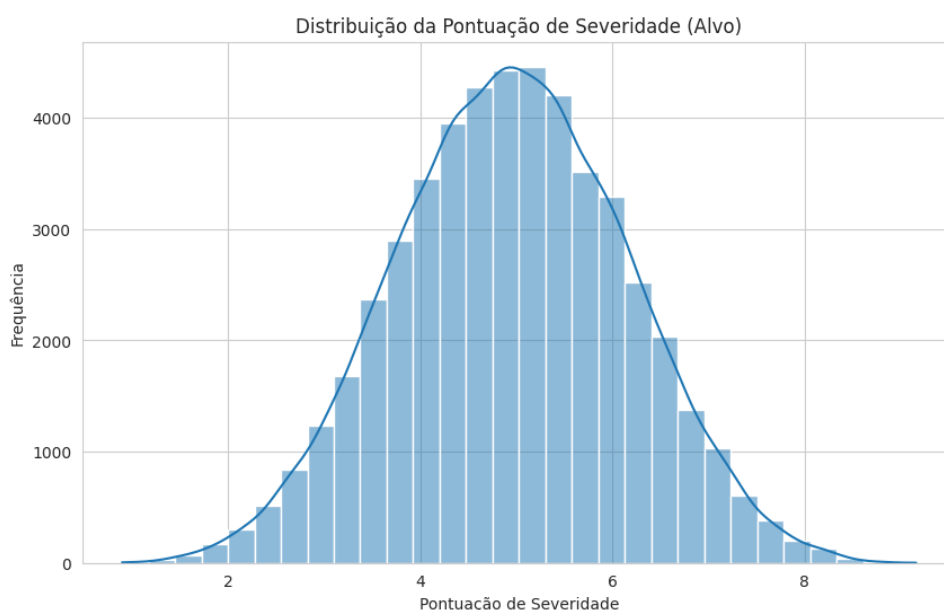
## 3. Desenvolvimento

Esta seção detalha o processo prático de construção do modelo, desde a configuração do ambiente até o treinamento e a otimização do algoritmo, implementados em Python.

### 3.1 Coleta e Análise Exploratória dos Dados

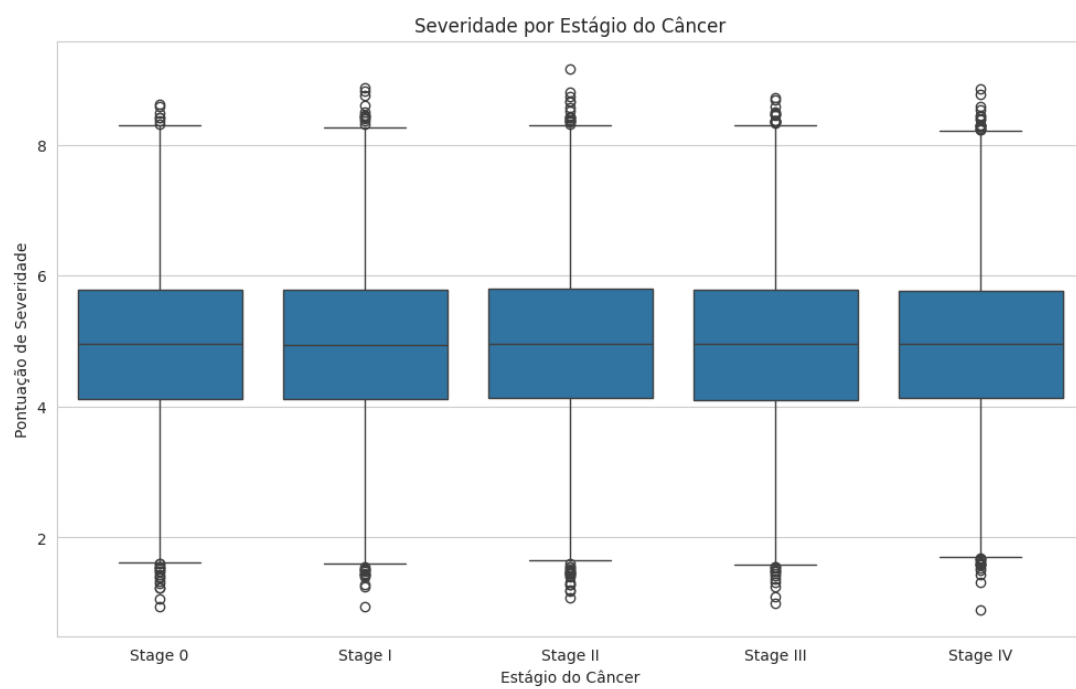
O ponto de partida foi a utilização de um conjunto de dados público [4]. É fundamental ressaltar que se trata de um *dataset* sintético, gerado para simular tendências globais de câncer com base em relatórios de saúde pública, não contendo informações de pacientes reais. Após a coleta a partir do arquivo `global_cancer_patients_2015_2024.csv` e a remoção do identificador `Patient_ID`, a Análise Exploratória de Dados (EDA) revelou uma forte correlação positiva entre o estágio do câncer e a severidade, conforme ilustrado na Figura 3.2. Um desafio técnico identificado foi a presença de numerais romanos na coluna `Cancer_Stage`, o que exigiu uma função de conversão para permitir a correta ordenação visual. A Figura 3.3 apresenta a matriz de correlação, que apontou outras relações lineares de interesse.

Figura 3.1: Distribuição da variável alvo (Target Severity Score).



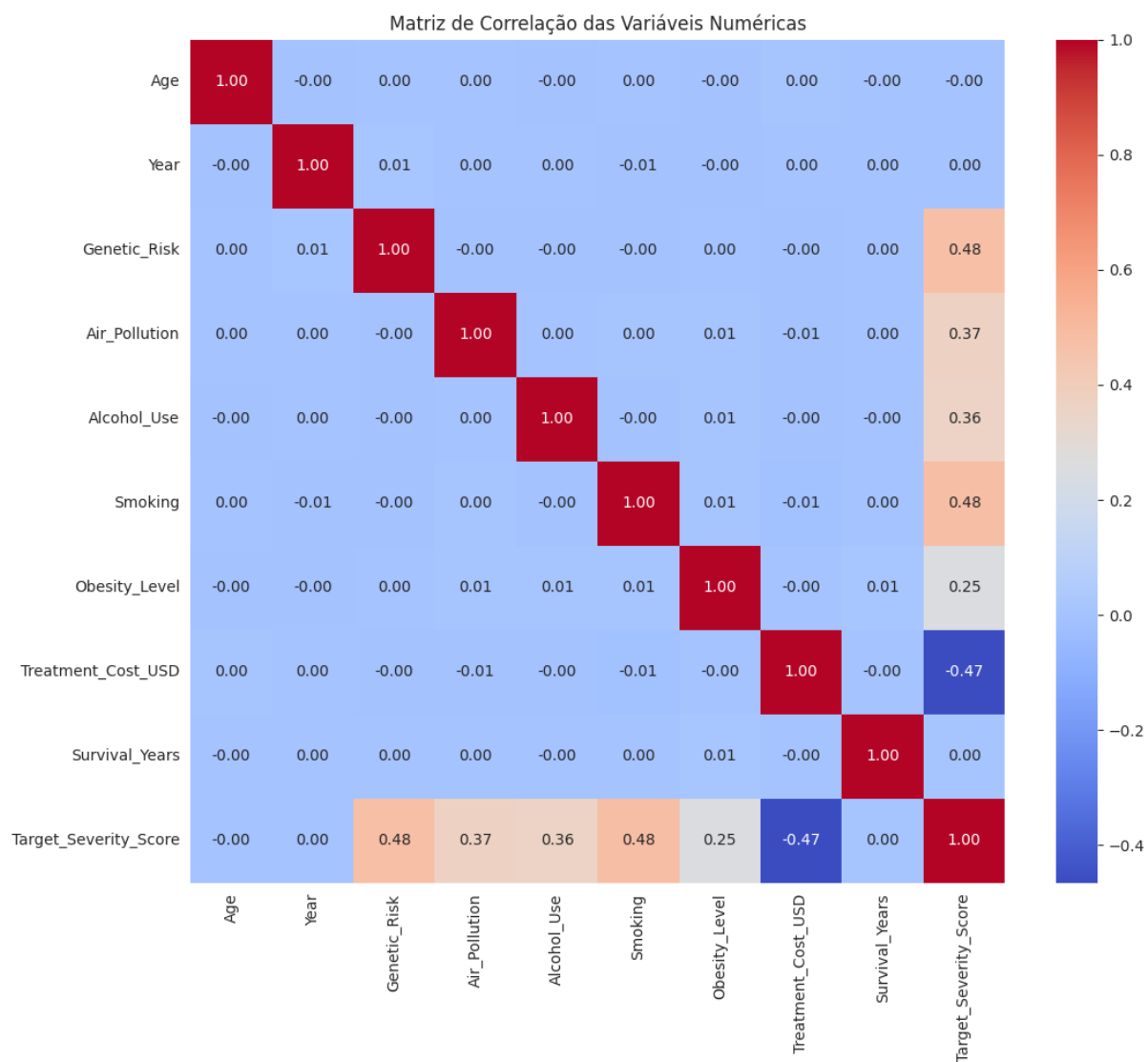
*Fonte: Autor (2025).*

Figura 3.2: Relação entre o Estágio do Câncer e a Pontuação de Severidade.



Fonte: Autor (2025).

Figura 3.3: Matriz de correlação entre as variáveis numéricas do dataset.



Fonte: Autor (2025).

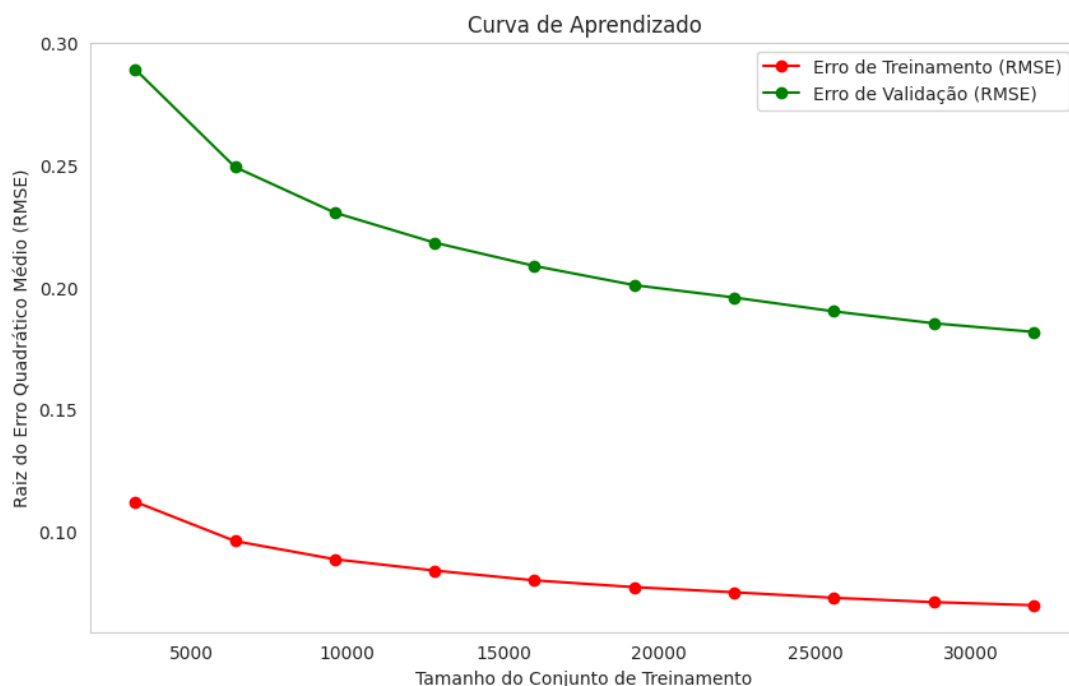
## 3.2 Pré-processamento e Construção do *Pipeline*

Para garantir a robustez e a reprodutibilidade, foi construído um *Pipeline* com um *ColumnTransformer*, aplicando *StandardScaler* às *features* numéricas e *OneHotEncoder* às categóricas. Esta abordagem previne o vazamento de dados e encapsula todo o fluxo de transformação.

## 3.3 Treinamento e Otimização do Modelo

O conjunto de dados foi dividido em 80% para treinamento e 20% para teste. O *RandomForestRegressor* foi treinado primeiramente com seus hiperparâmetros padrão (`n_estimators: 100`; `max_depth: None`; `min_samples_split: 2`; `min_samples_leaf: 1`). A Figura 3.4 ilustra a curva de aprendizado do modelo base. Subsequentemente, a técnica *RandomizedSearchCV* foi utilizada para otimizar os hiperparâmetros, resultando em um modelo final ajustado.

Figura 3.4: Curva de aprendizado do modelo base.



Fonte: Autor (2025).

## 4. Resultados e Discussões

Neste capítulo, são apresentados os resultados quantitativos da performance dos modelos.

### 4.1 Desempenho dos Modelos

O modelo inicial (base) foi avaliado no conjunto de teste para estabelecer uma linha de base de performance. Após o processo de otimização, o modelo com os melhores hiperparâmetros (Listagem 1) foi igualmente avaliado.

```
{
  'regressor__n_estimators': 300,
  'regressor__min_samples_split': 2,
  'regressor__min_samples_leaf': 1,
  'regressor__max_depth': None
}
```

Listing 1: Melhores hiperparâmetros encontrados pela busca aleatória.

### 4.2 Análise Comparativa

A Tabela 4.1 apresenta uma comparação direta das métricas de performance entre o modelo base e o modelo otimizado. A Figura 4.1 ilustra a precisão do modelo otimizado.

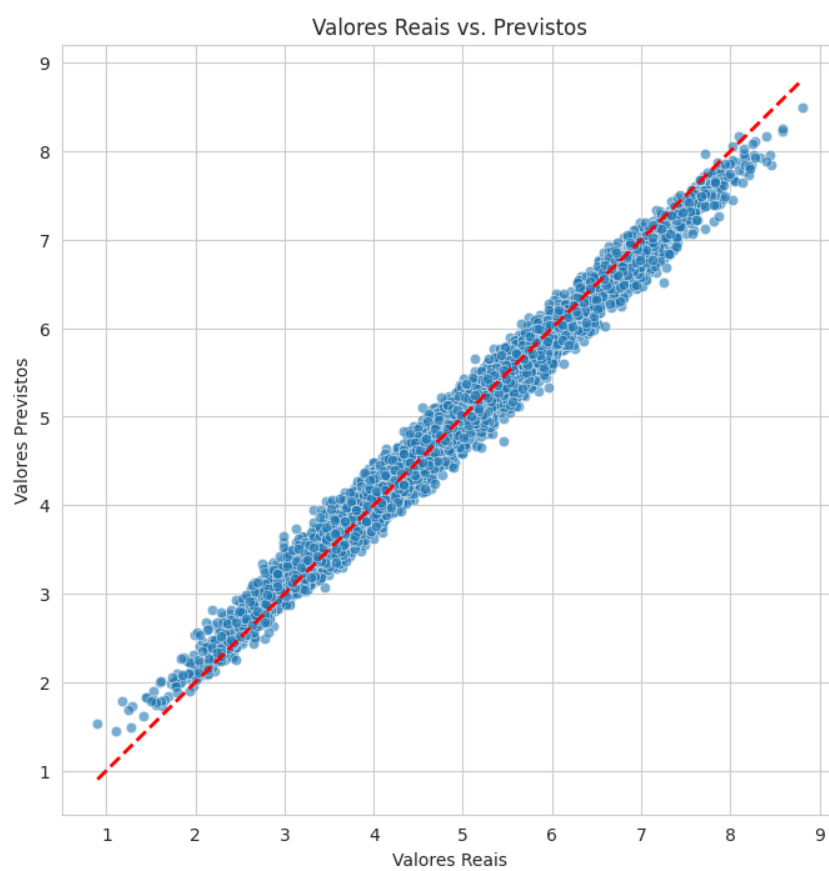
Tabela 4.1: Comparação de desempenho entre o modelo base e o otimizado.

Métrica	Modelo Base	Modelo Otimizado
Mean Absolute Error (MAE)	0,136	0,136
Root Mean Squared Error (RMSE)	0,172	0,171
R-squared ( $R^2$ )	0,979	0,979

*Fonte: Autor (2025).*

A análise demonstra que o processo de otimização resultou em uma melhoria marginal, quase insignificante. O coeficiente de determinação ( $R^2$ ) permaneceu em 0,979. Isso indica que o modelo base, com seus hiperparâmetros padrão, já era extremamente eficaz para este conjunto de dados, atingindo um desempenho próximo do ideal e deixando pouco espaço para melhorias incrementais através do ajuste fino.

Figura 4.1: Gráfico de dispersão dos valores reais vs. previstos para o modelo otimizado.



*Fonte: Autor (2025).*



## 5. Conclusão e Trabalhos Futuros

Este trabalho desenvolveu e avaliou com sucesso uma solução de Inteligência Artificial para prever a severidade do câncer. O modelo `RandomForestRegressor` otimizado atingiu um coeficiente de determinação ( $R^2$ ) de 0,979, demonstrando uma capacidade preditiva excepcional e respondendo afirmativamente à pergunta de pesquisa.

O sucesso do modelo reforça o potencial do *Machine Learning* como uma ferramenta de apoio à decisão na oncologia, contribuindo para a medicina de precisão. Um desafio técnico notável foi a manipulação de dados não padronizados (numerais romanos), superado com a criação de uma função de conversão, o que ressalta a importância da etapa de pré-processamento.

A principal limitação do trabalho reside no uso de dados sintéticos, o que impede a generalização direta do modelo para populações de pacientes reais sem uma validação prévia com dados clínicos autênticos. Como trabalhos futuros, sugere-se a exploração de outros algoritmos de *ensemble*, como o *Gradient Boosting* [2], e uma análise aprofundada da importância das *features* para gerar *insights* clínicos. A implantação do modelo como uma *API* seria o passo subsequente para sua validação em um ambiente real.

No fim o projeto entregou uma ferramenta preditiva de alto desempenho e demonstrou o potencial transformador do *Machine Learning* no cuidado quantitativo e personalizado ao paciente com câncer.

# Referências Bibliográficas

- [1] Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001.
- [2] Jerome H. Friedman. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29(5):1189–1232, 2001.
- [3] Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, New York, 2 edition, 2009.
- [4] Zahid Mughal. Global cancer patients (2015-2024). <https://www.kaggle.com/datasets/zahidmughal2343/global-cancer-patients-2015-2024>, 2025. Acesso em: 13 jun. 2025.
- [5] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [6] Hyuna Sung, Jacques Ferlay, Rebecca L. Siegel, Mathieu Laversanne, Isabelle Soerjomataram, Ahmedin Jemal, and Freddie Bray. Global cancer statistics 2022: Globocan estimates of incidence and mortality worldwide for 36 cancers in 185 countries. *CA: A Cancer Journal for Clinicians*, 74(3):209–249, 2024.