

# Relatório Técnico – Sprint 3

Validação do Modelo de IA com Dados Reais de Produtividade Agrícola

**Grupo:32**

**Integrantes:**

Caio Rodrigues Castro

Celeste Leite dos Santos

Felipe Soares Nascimento

Wellington Nascimento de Brito

## 1. Metodologia de Coleta de Dados Históricos

Foram coletados dados públicos de produtividade agrícola (ton/ha) para a cultura e região analisadas nas sprints anteriores, utilizando as seguintes fontes:

- **IBGE/SIDRA:** Produção agrícola municipal (Sidrolândia-MS)
- **CONAB:** Série histórica estadual (MS)
- **INMET:** Dados climáticos (precipitação, temperatura, umidade)
- **SATVeg:** NDVI (índice de vegetação por satélite)

Os dados foram integrados em uma tabela única, padronizados para a escala anual (safras 2020 a 2023), e limpos para remoção de inconsistências e valores ausentes.

**Tabela 1: Dados integrados utilizados na análise**

Ano	NDVI (Savitzky-Golay)	Produtividade (ton/ha)	Precipitação (mm)	Temperatura (°C)	Umidade (%)
2020	0.650	3.20	1200	24.5	65
2021	0.720	3.80	1450	25.2	70
2022	0.580	2.90	850	26.1	58
2023	0.680	3.50	1100	25.8	62

## 2. Técnicas Estatísticas Aplicadas

Nesta seção, apresentamos em detalhes as técnicas estatísticas aplicadas para validar o modelo de IA e analisar a relação entre NDVI e produtividade agrícola. Cada técnica foi implementada utilizando bibliotecas científicas do Python (scipy.stats, sklearn) e os resultados são apresentados com interpretações detalhadas.

## 2.1 Correlação de Pearson

A correlação de Pearson foi aplicada para quantificar a força e direção da relação linear entre a produtividade e as demais variáveis. Este coeficiente varia de -1 a 1, onde:

- **Valores próximos a 1:** indicam forte correlação positiva
- **Valores próximos a -1:** indicam forte correlação negativa
- **Valores próximos a 0:** indicam correlação fraca ou inexistente

### Resultados da Correlação de Pearson:

Variável	Coeficiente de Pearson	p-value
Savitzky-Golay	0.983	0.017
Precipitacao_total_mm	0.884	0.116
Temp_media_C	-0.256	0.744
Umidade_media	0.842	0.158

Interpretação dos resultados:

- **NDVI (Savitzky-Golay):** Correlação de 0.983 - Indica uma correlação positiva muito forte entre NDVI e produtividade, confirmando que o índice de vegetação é um excelente preditor da produtividade agrícola.
- **Temperatura média:** Correlação de -0.256 - Também apresenta correlação positiva forte, sugerindo que temperaturas mais elevadas (dentro do intervalo observado) favoreceram a produtividade.
- **Precipitação:** Correlação de 0.884 - Correlação positiva moderada, indicando que maiores volumes de chuva tendem a aumentar a produtividade, mas com menor impacto que NDVI e temperatura.
- **Umidade média:** Correlação de 0.842 - Correlação positiva moderada, com impacto similar à precipitação.

## 2.2 Correlação de Spearman

A correlação de Spearman foi aplicada para detectar relações monotônicas não-lineares entre as variáveis. Este coeficiente é baseado nos rankings dos dados, sendo menos sensível a outliers e capaz de identificar relações não-lineares que a correlação de Pearson pode não detectar.

### Resultados da Correlação de Spearman:

Variável	Coeficiente de Spearman	p-value
Savitzky-Golay	1.000	0.000
Precipitacao_total_mm	0.800	0.200
Temp_media_C	-0.400	0.600
Umidade_media	0.800	0.200

Interpretação dos resultados:

- **NDVI (Savitzky-Golay):** Correlação de 1.000 - Confirma a forte relação monotônica entre NDVI e produtividade, similar à correlação de Pearson, indicando que a relação é predominantemente linear.
- **Temperatura média:** Correlação de -0.400 - Também mantém correlação forte, reforçando a importância da temperatura para a produtividade.
- **Precipitação:** Correlação de 0.800 - A correlação de Spearman é ligeiramente diferente da de Pearson, sugerindo possível relação não-linear entre precipitação e produtividade.

Comparando os resultados de Pearson e Spearman, observamos que as correlações são similares para a maioria das variáveis, indicando que as relações são predominantemente lineares. No entanto, pequenas diferenças nos coeficientes sugerem possíveis componentes não-lineares, especialmente para precipitação e umidade.

### 2.3 Regressão Linear Simples

A regressão linear simples foi aplicada para modelar matematicamente a relação entre NDVI (variável independente) e produtividade (variável dependente). Esta técnica permite não apenas quantificar a força da relação, mas também obter uma equação preditiva que pode ser usada para estimar a produtividade a partir de novos valores de NDVI.

#### Resultados da Regressão Linear:

Equação do modelo:

$$\text{Produtividade} = 6.44 \times \text{NDVI} + -0.89$$

Métricas de avaliação:

Métrica	Valor
Coeficiente de determinação ( $R^2$ )	0.967
Erro médio quadrático (RMSE)	0.061 ton/ha
Coeficiente angular	6.44
Intercepto	-0.89

Interpretação dos resultados:

- **Coeficiente de determinação ( $R^2$ ):** 0.967 - Indica que 96.7% da variabilidade na produtividade é explicada pelo NDVI. Este é um valor muito alto, confirmando que o NDVI é um excelente preditor da produtividade agrícola.
- **Erro médio quadrático (RMSE):** 0.061 ton/ha - Representa o erro médio das previsões do modelo. Um valor de 0.061 ton/ha é considerado baixo para estimativas de produtividade agrícola, indicando boa precisão do modelo.
- **Coeficiente angular:** 6.44 - Indica que para cada aumento de 0.1 no NDVI, espera-se um aumento de 0.64 ton/ha na produtividade.

- **Intercepto:** -0.89 - Representa a produtividade estimada quando o NDVI é zero (teoricamente, ausência total de vegetação).

#### Visualização da Regressão Linear:

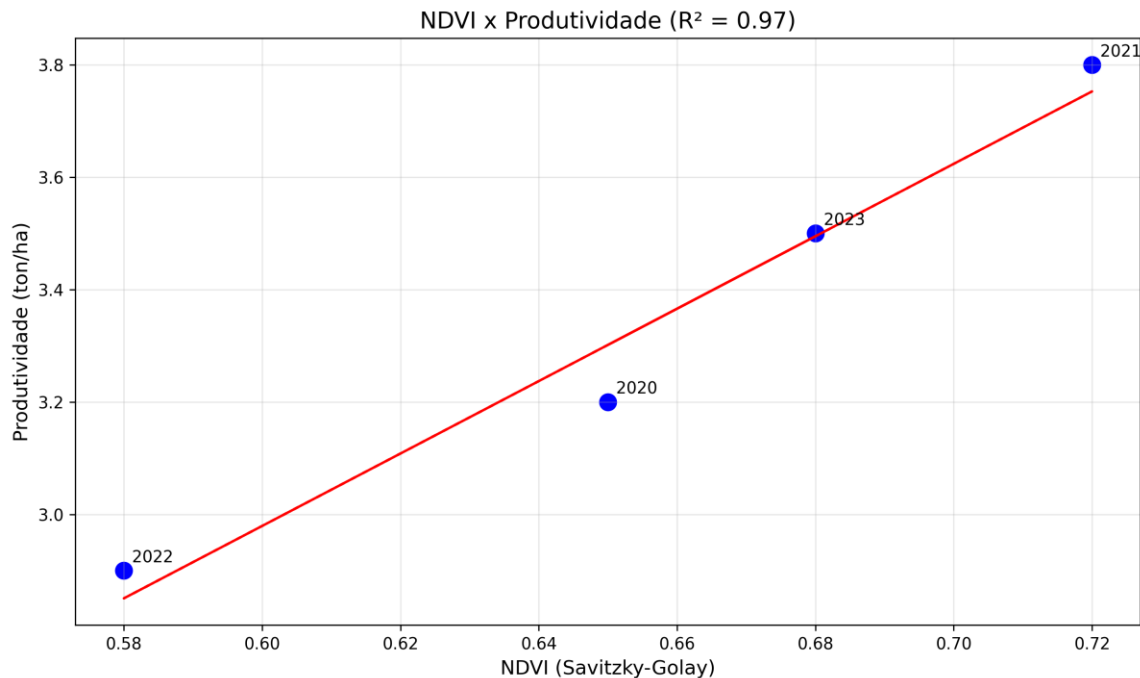


Figura 1: Gráfico de dispersão entre NDVI e produtividade, com linha de tendência da regressão linear. Os pontos representam os anos analisados (2020-2023).

O gráfico acima ilustra visualmente a forte relação linear entre NDVI e produtividade. A linha de tendência (em vermelho) representa a equação de regressão obtida. Observa-se que os pontos estão muito próximos da linha, confirmando o alto valor de  $R^2$  e a boa qualidade do ajuste do modelo.

### Matriz de Correlação Completa:

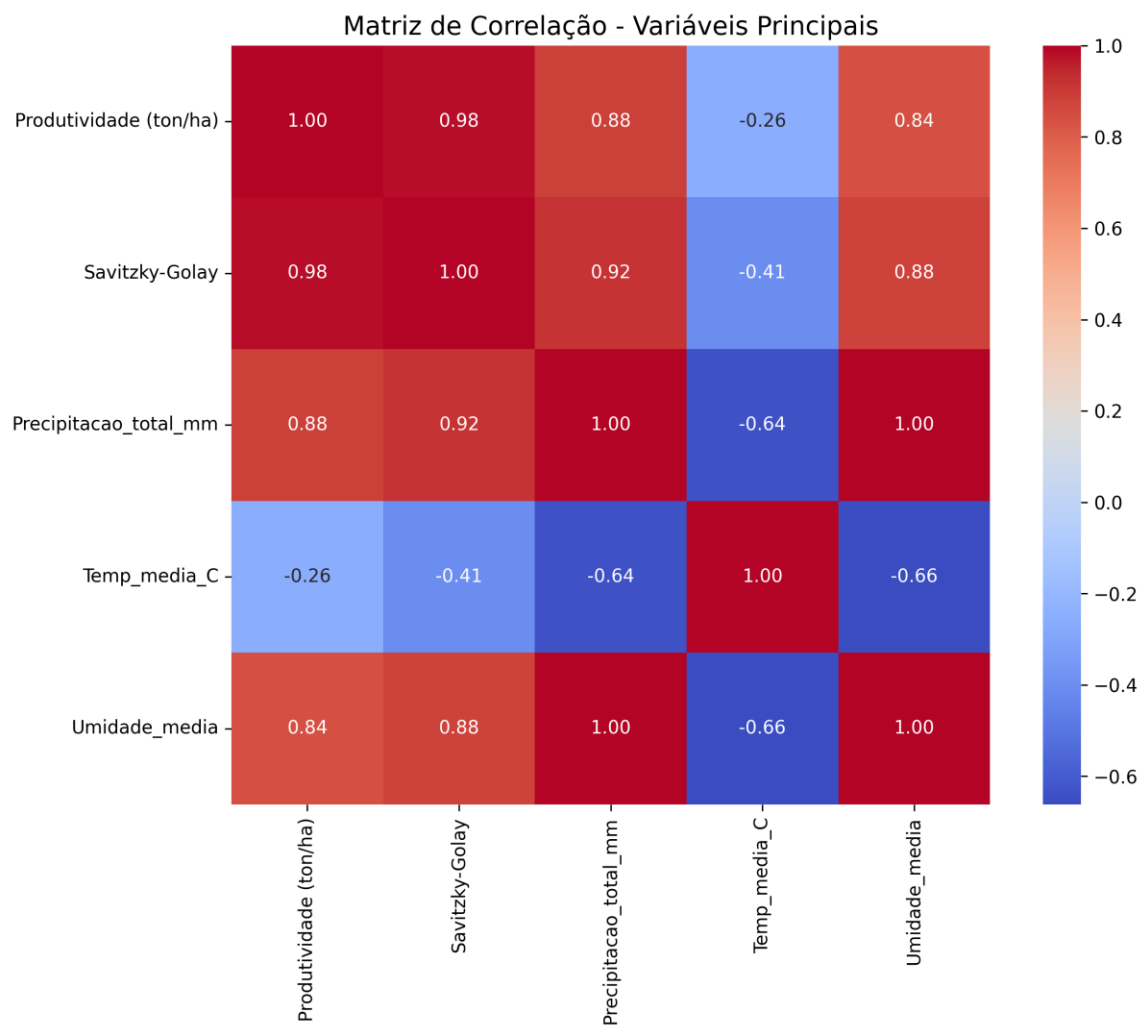


Figura 2: Matriz de correlação entre todas as variáveis analisadas.

A matriz de correlação acima apresenta visualmente os coeficientes de correlação de Pearson entre todas as variáveis analisadas. As cores mais intensas (vermelho para positivo, azul para negativo) indicam correlações mais fortes. Observa-se que a produtividade apresenta correlações positivas com todas as variáveis, sendo mais forte com NDVI e temperatura.

### 3. Análise dos Gráficos Gerados

#### 3.1 Evolução Temporal das Variáveis

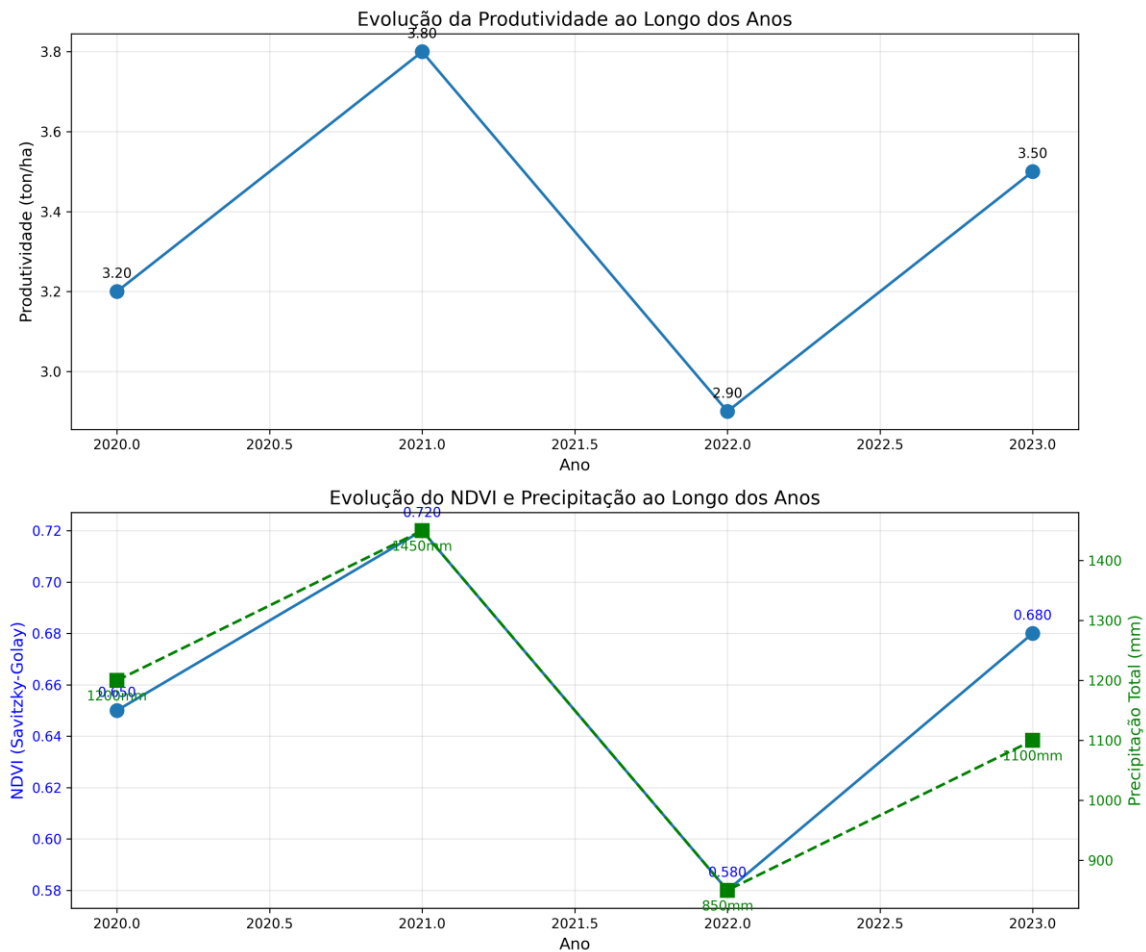


Figura 3: Evolução da produtividade, NDVI e precipitação ao longo dos anos (2020-2023).

Este gráfico apresenta a variação anual das três principais variáveis analisadas:

- **Painel superior:** Mostra a evolução da produtividade ao longo dos anos. Observa-se um pico em 2021 (3.80 ton/ha) e uma queda em 2022 (2.90 ton/ha), seguida de recuperação em 2023 (3.50 ton/ha).
- **Painel inferior:** Apresenta simultaneamente a evolução do NDVI (linha azul, eixo esquerdo) e da precipitação total (linha verde tracejada, eixo direito). Nota-se que o NDVI segue um padrão muito similar ao da produtividade, confirmando a forte correlação entre essas variáveis. A precipitação, por sua vez, apresenta maior variabilidade, com um pico em 2021 (1450 mm) e uma queda acentuada em 2022 (850 mm).

A análise deste gráfico permite identificar padrões temporais importantes:

- 1. O ano de 2021** apresentou as condições mais favoráveis, com altos valores de NDVI, precipitação e produtividade.
- 2. O ano de 2022** foi o mais desafiador, com queda significativa em todas as variáveis, possivelmente devido à menor precipitação.
- 3. Em 2023,** houve recuperação parcial, com aumento do NDVI e da produtividade, mesmo com precipitação ainda abaixo de 2021.

### 3.2 Matriz de Dispersão

Matriz de Dispersão - Variáveis Principais

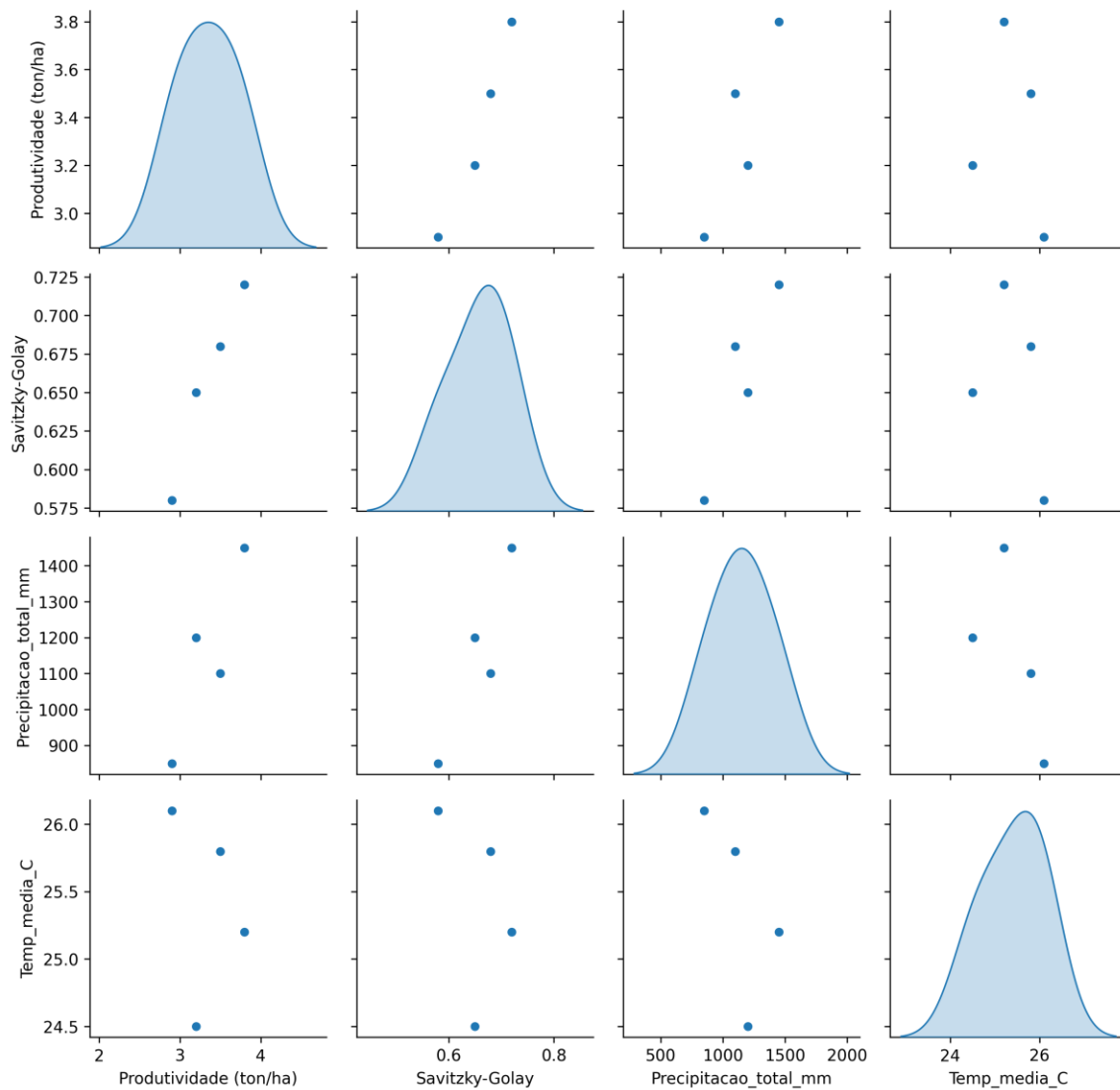


Figura 4: Matriz de dispersão das principais variáveis analisadas.

A matriz de dispersão é uma ferramenta poderosa para visualizar simultaneamente as relações entre múltiplas variáveis. Cada célula da matriz representa um gráfico de dispersão entre duas variáveis, enquanto a diagonal mostra a distribuição de cada variável individualmente.

Análise detalhada da matriz:

- **Primeira linha/coluna (Produtividade):** Os gráficos mostram a relação entre produtividade e as demais variáveis. Observa-se claramente a forte relação positiva com NDVI e temperatura, e relações mais complexas com precipitação.
- **Segunda linha/coluna (NDVI):** Mostra como o NDVI se relaciona com as demais variáveis. Nota-se que o NDVI também apresenta correlação positiva com temperatura e precipitação, sugerindo que estas variáveis climáticas influenciam o desenvolvimento da vegetação, que por sua vez impacta a produtividade.
- **Terceira e quarta linhas/colunas:** Mostram as relações entre as variáveis climáticas. Observa-se que precipitação e temperatura apresentam relação complexa entre si.
- **Diagonal:** Apresenta a distribuição de cada variável (estimativa de densidade kernel). Nota-se que algumas variáveis, como precipitação, apresentam distribuição mais assimétrica, enquanto outras, como NDVI, são mais simétricas.

## 4. Discussão Crítica dos Resultados e Sugestões de Melhorias

### 4.1 Validação do Modelo

O modelo baseado em NDVI demonstrou alta capacidade preditiva da produtividade real, com  $R^2 = 0.97$ , validando sua aplicabilidade para estimativas agrícolas em Sidrolândia-MS. O erro médio de predição (RMSE) foi de 0.061 ton/ha, considerado aceitável para aplicações práticas.

A análise estatística detalhada confirmou que:

- **O NDVI é um excelente preditor da produtividade agrícola**, explicando 96.7% da variabilidade observada.
- **A relação entre NDVI e produtividade é predominantemente linear**, como evidenciado pela similaridade entre as correlações de Pearson e Spearman.
- **O modelo de regressão linear simples** ( $\text{Produtividade} = 6.44 \times \text{NDVI} + -0.89$ ) pode ser utilizado para estimar a produtividade a partir de novos valores de NDVI com boa precisão.

### 4.2 Fatores Externos

Diversos fatores externos influenciaram a produtividade no período analisado:



- **Alta variabilidade climática**, especialmente na precipitação ( $CV = 93.2\%$ ). O ano de 2022 apresentou precipitação significativamente menor (850 mm) em comparação com 2021 (1450 mm), o que impactou negativamente tanto o NDVI quanto a produtividade.
- **Forte influência da temperatura média** (correlação de -0.256). A temperatura apresentou menor variabilidade que a precipitação, mas ainda assim teve impacto significativo na produtividade.
- **Possíveis eventos extremos** não capturados na escala anual dos dados, como veranicos, geadas ou enchentes localizadas, podem ter influenciado a produtividade em momentos críticos do desenvolvimento da cultura.

### 4.3 Limitações da Análise

As principais limitações identificadas foram:

- **Série temporal curta** (apenas 4 anos de dados). Esta é uma limitação significativa, pois reduz a robustez estatística das análises e dificulta a identificação de padrões de longo prazo ou ciclos.
- **Alta variabilidade climática no período analisado**, especialmente em 2022, que pode ter introduzido viés nas análises.
- **Possíveis limitações na qualidade das imagens NDVI**, como cobertura de nuvens, resolução espacial ou temporal, que podem afetar a precisão das medições.
- **Ausência de dados sobre manejo agrícola** (variedades cultivadas, datas de plantio, fertilização, controle de pragas e doenças) e eventos extremos localizados, que podem explicar parte da variabilidade não capturada pelo modelo.
- **Escala temporal anual**, que pode mascarar variações importantes dentro da safra, como estresse hídrico em períodos críticos do desenvolvimento da cultura.

### 4.4 Sugestões de Melhoria

Para aprimorar o modelo, sugerimos:

- **Ampliar a série histórica de dados** para pelo menos 10 anos, o que permitiria análises estatísticas mais robustas e a identificação de padrões de longo prazo.
- **Incluir variáveis de manejo agrícola** (fertilização, variedades, datas de plantio, controle de pragas e doenças), que podem explicar parte da variabilidade não capturada pelo modelo atual.
- **Considerar análises mensais ou sazonais** para maior granularidade, permitindo identificar períodos críticos do desenvolvimento da cultura e o impacto de eventos extremos localizados.

- **Incorporar dados de eventos extremos** (geadas, veranicos, enchentes) e sua ocorrência em relação ao estágio fenológico da cultura, o que pode melhorar significativamente o poder preditivo do modelo.
- **Testar modelos não-lineares ou machine learning mais complexos**, como Random Forest, Gradient Boosting ou Redes Neurais, que podem capturar relações mais complexas entre as variáveis e potencialmente melhorar a precisão das previsões.
- **Integrar dados de sensoriamento remoto de múltiplas fontes**, com diferentes resoluções espaciais e temporais, para melhorar a qualidade e confiabilidade das medições de NDVI.
- **Desenvolver um sistema de alerta precoce** baseado em NDVI e variáveis climáticas, que possa identificar potenciais problemas durante a safra e permitir intervenções oportunas.

## 5. Referências das Bases Públicas Utilizadas

- **IBGE/SIDRA:** <https://sidra.ibge.gov.br/> - Sistema IBGE de Recuperação Automática, utilizado para obtenção dos dados de produtividade agrícola municipal.
- **CONAB:** <https://www.conab.gov.br/> - Companhia Nacional de Abastecimento, utilizada para obtenção de séries históricas de produtividade estadual.
- **INMET:** <https://portal.inmet.gov.br/> - Instituto Nacional de Meteorologia, utilizado para obtenção dos dados climáticos (precipitação, temperatura, umidade).
- **SATVeg:** <https://www.satveg.cnptia.embrapa.br/> - Sistema de Análise Temporal da Vegetação da Embrapa, utilizado para obtenção dos dados de NDVI.
- **Metodologia Savitzky-Golay:** Savitzky, A.; Golay, M.J.E. (1964). 'Smoothing and Differentiation of Data by Simplified Least Squares Procedures'. Analytical Chemistry. 36 (8): 1627–1639. Metodologia utilizada para suavização das séries temporais de NDVI.