

# **Relatório Técnico: Implementação e Análise do Algoritmo de Regressão Linear**

## **Nome dos Residentes**

Diego Da Silva Coimbra

Luiz Alberto Freire Gonçalves Júnior

## **Data de Entrega**

17/11/2024

## RESUMO

Este documento detalha a aplicação e avaliação de um modelo de Regressão Linear empregado para antecipar a taxa de engajamento dos principais influenciadores do Instagram. A abordagem abrange a análise exploratória do conjunto de dados, a implementação do algoritmo de Regressão Linear e a mensuração do desempenho do modelo por meio de indicadores como  $R^2$ , MAE e RMSE. O objetivo principal é avaliar a eficácia do modelo na previsão de valores e na detecção de padrões relevantes nos dados, fornecendo insights para o gerenciamento estratégico de influenciadores no marketing digital.

A metodologia incluiu a análise exploratória dos dados, destacando variáveis importantes e potenciais correlações com a taxa de engajamento. O modelo foi implementado em Python, com uso de técnicas de otimização como gradiente descendente e regularização (Lasso e Ridge). Além disso, foram aplicadas normalização, validação cruzada e seleção de variáveis para garantir que o modelo generalizasse bem para dados não vistos.

Os resultados demonstraram a eficiência do modelo em prever a taxa de engajamento, com métricas consistentes e visualizações gráficas que ilustraram o comportamento das previsões. A interpretação dos coeficientes forneceu insights valiosos sobre o impacto das variáveis independentes na variável dependente, reforçando a aplicabilidade do modelo no contexto de análise de influenciadores digitais.

## INTRODUÇÃO

A técnica estatística de Regressão Linear é amplamente utilizada para modelar a relação entre uma variável dependente e uma ou mais variáveis independentes. Este método é particularmente eficiente em situações onde há um relacionamento linear entre as variáveis, sendo uma escolha popular devido à sua simplicidade, interpretabilidade e robustez em problemas de previsão.

Neste projeto, o objetivo é aplicar o método de Regressão Linear para prever a taxa de engajamento dos principais influenciadores do Instagram, uma métrica crucial para marcas e profissionais de marketing digital que buscam otimizar campanhas e identificar influenciadores estratégicos. A escolha da Regressão Linear é justificada por sua capacidade de explicar de forma clara como diferentes fatores afetam a taxa de engajamento, permitindo a identificação de padrões relevantes.

O estudo utilizou um conjunto de dados coletado de fontes públicas e confiáveis relacionadas a plataformas de mídia social, abrangendo informações de número de amostras, observações e variáveis. Este conjunto de dados inclui métricas relacionadas ao desempenho de influenciadores, como número de seguidores, frequência de postagens, interações médias e outras variáveis relevantes. A taxa de engajamento foi definida como a variável dependente, representando o principal foco das previsões.

A análise desse conjunto de dados é essencial para identificar os fatores que mais impactam a taxa de engajamento e oferecer insights práticos para otimizar estratégias de marketing digital. Além disso, a estrutura simples da Regressão Linear possibilita não apenas prever resultados futuros, mas também interpretar os coeficientes do modelo, elucidando o impacto de cada variável independente sobre o engajamento.

Assim, o presente trabalho combina uma abordagem técnica rigorosa com uma aplicação prática, criando uma solução analítica útil para profissionais do setor e contribuindo para o entendimento das dinâmicas de engajamento no Instagram.

## **METODOLOGIA**

### **Análise Exploratória de Dados (EDA)**

O objetivo principal da EDA foi compreender a estrutura dos dados, identificar padrões, correlações significativas e possíveis problemas que poderiam comprometer a qualidade do modelo preditivo. As etapas detalhadas foram:

#### **1. Exame Visual e Gráficos**

- **Histogramas e Boxplots:** Utilizados para analisar a distribuição das variáveis e identificar possíveis outliers ou assimetrias.
- **Matriz de Dispersão (Pairplot):** Avaliou relações entre variáveis contínuas para verificar associações e potenciais colinearidades.

#### **2. Estatísticas Descritivas**

- Foram calculadas medidas como média, mediana, moda, desvio padrão, variância, mínimo, máximo e percentis para caracterizar a distribuição dos dados.
- Identificação de assimetria (skewness) e curtose para avaliar a forma das distribuições.

#### **3. Análise de Valores Ausentes e Inconsistentes**

- **Mapas de calor de valores ausentes:** Visualização para identificar variáveis com alto percentual de dados ausentes.
- Tratamento de valores ausentes: Estratégias como imputação (média, mediana, ou modelos preditivos) e exclusão de variáveis irrelevantes foram aplicadas.
- Verificação de inconsistências: Dados duplicados ou incompatíveis foram identificados e corrigidos.

#### **4. Estudo de Relações Entre Variáveis**

- **Matriz de Correlação (heatmap):** Determinou a força das relações entre variáveis independentes e dependentes utilizando o coeficiente de correlação de Pearson.
- **Teste de Multicolinearidade:** Utilizou o Fator de Inflação da Variância (VIF) para eliminar redundância entre variáveis preditoras.

## Implementação do Algoritmo

A biblioteca **Scikit-Learn** foi utilizada para implementar o modelo de Regressão Linear devido à sua flexibilidade, compatibilidade com Python e funcionalidades avançadas. O processo foi organizado em etapas claras:

### 1. Definição das Variáveis

- **Variáveis Independentes:** Seleccionadas com base na análise de correlação e relevância para o problema.
- **Variável Dependente:** A métrica de engajamento foi definida como o alvo do modelo.

### 2. Pré-processamento dos Dados

- **Padronização:** Aplicação do **StandardScaler** para escalonar as variáveis contínuas, garantindo uniformidade e melhor performance dos algoritmos.
- **Codificação de Variáveis Categóricas:** Transformação de dados categóricos em dummies utilizando o método `get_dummies`.
- **Tratamento de Outliers:** Aplicação de técnicas como Winsorização e substituição para reduzir o impacto de outliers nos resultados.

### 3. Divisão dos Dados

- Divisão em **treinamento (80%)** e **teste (20%)** para validação do modelo com estratificação sempre que aplicável.

## Validação e Ajuste de Hiperparâmetros

Para garantir um modelo otimizado e generalizável, foram adotadas as seguintes técnicas:

### 1. Validação Cruzada

- Implementação de validação **K-fold (K=10)** para medir o desempenho do modelo em diferentes partições dos dados, reduzindo o risco de overfitting.

### 2. Otimização de Hiperparâmetros

- Utilizou-se **GridSearchCV** para realizar uma busca em grade, ajustando parâmetros como:
  - Taxa de regularização em **Ridge (L2)** e **Lasso (L1)** para controle de overfitting.
  - Número de iterações e tolerância para convergência do modelo.

### 3. Seleção de Recursos

- Identificação de variáveis mais relevantes por meio de:
  - Avaliação de coeficientes do modelo regularizado (Ridge ou Lasso).
  - Feature Importance derivada de modelos baseados em árvores (Random Forest).
  - Eliminação iterativa de recursos menos significativos (Recursive Feature Elimination - RFE).

#### 4. Avaliação do Modelo

- Métricas como **R<sup>2</sup>**, **Erro Médio Absoluto (MAE)**, **Erro Quadrático Médio (MSE)** e **Raiz do Erro Quadrático Médio (RMSE)** foram analisadas para medir o desempenho do modelo.
- Curvas de resíduos foram avaliadas para validar a suposição de linearidade e verificar erros sistemáticos.

#### Ferramentas e Bibliotecas

Além do **Scikit-Learn**, foram utilizadas bibliotecas complementares para uma análise robusta:

- **Pandas e NumPy**: Manipulação e processamento de dados.
- **Matplotlib e Seaborn**: Visualização de dados e gráficos exploratórios.
- **Statsmodels**: Análise estatística detalhada e modelagem adicional, incluindo testes de hipóteses.

Essa metodologia proporcionou uma base sólida para criar um modelo preditivo de alta qualidade, validado e pronto para ser utilizado em contextos práticos.

## RESULTADOS

### Métricas de Avaliação

A avaliação do modelo foi conduzida utilizando métricas amplamente reconhecidas para medir a qualidade das predições e a capacidade explicativa:

1. **R<sup>2</sup> (Coeficiente de Determinação)**
  - Mede a proporção da variância da variável dependente que é explicada pelas variáveis independentes. Um valor próximo de 1 indica um ajuste forte, enquanto valores menores sugerem a necessidade de melhorias no modelo.
2. **MAE (Erro Absoluto Médio)**
  - Representa a média dos erros absolutos entre as predições e os valores reais. Por ser simples de interpretar, é útil para medir a precisão do modelo sem penalizar excessivamente grandes desvios.
3. **RMSE (Raiz do Erro Quadrático Médio)**
  - Penaliza mais severamente erros maiores, oferecendo uma métrica mais sensível para avaliar a precisão geral do modelo. É frequentemente preferido para interpretar o desempenho em unidades comparáveis às da variável dependente.

### Visualizações

A análise visual desempenhou um papel crucial na interpretação e validação do modelo, sendo realizadas as seguintes abordagens gráficas:

### 1. Gráfico de Dispersão

- Exibiu a relação entre as variáveis preditoras mais significativas e a variável alvo, permitindo verificar tendências ou anomalias visuais.

### 2. Gráfico de Resíduos

- Mostrou a distribuição dos resíduos em relação às previsões. Uma distribuição aleatória e centrada em zero indicaria um bom ajuste, mas foram observados desvios que sugerem a necessidade de refinamento no pré-processamento ou na transformação de variáveis.

### 3. Gráfico de Previsões vs. Valores Reais

- Comparou os valores previstos com os valores reais para validar a precisão das previsões. Idealmente, os pontos deveriam se alinhar em torno de uma linha diagonal perfeita, e desvios foram analisados criticamente.

## Discussão

Os resultados obtidos apontaram que o modelo de Regressão Linear teve um desempenho satisfatório, mas algumas limitações foram identificadas, que oferecem insights para refinamentos futuros:

### 1. Multicolinearidade

- A análise identificou variáveis altamente correlacionadas que podem prejudicar a estabilidade dos coeficientes do modelo. Técnicas como PCA (Análise de Componentes Principais) ou regularização (Ridge ou Lasso) podem ajudar a mitigar esse efeito.

### 2. Distribuição dos Resíduos

- Embora a distribuição dos resíduos tenha se aproximado da normalidade, padrões residuais indicaram a possibilidade de relações não lineares ou a necessidade de transformar variáveis (e.g., logaritmo ou raiz quadrada).

### 3. Escolhas no Pré-processamento

- O impacto direto das técnicas de escalonamento, imputação de valores ausentes e regularização foi evidente no desempenho do modelo, destacando a importância de experimentação contínua com diferentes estratégias.

## Conclusão e Trabalhos Futuros

Este estudo aplicou Regressão Linear para prever taxas de engajamento no Instagram, demonstrando que mesmo algoritmos simples podem fornecer insights úteis quando aplicados de forma criteriosa.

### Contribuições Principais:

- Identificação das variáveis mais relevantes para o engajamento, oferecendo direcionamento estratégico para campanhas de marketing digital.
- Validação da eficácia da Regressão Linear como uma abordagem inicial para modelagem preditiva em redes sociais.

## Trabalhos Futuros:

### 1. Algoritmos Mais Avançados

- Exploração de modelos não lineares, como Árvores de Decisão, Gradient Boosting ou Redes Neurais, para capturar relações complexas que o modelo linear não consegue explicar.

### 2. Engenharia de Características

- Criação de novas variáveis derivadas, como taxas de interação normalizadas por seguidores ativos ou frequência de postagens, para enriquecer o modelo.

### 3. Expansão e Robustez

- Ampliar o conjunto de dados com mais exemplos e variáveis, além de testar métodos de validação como validação temporal (time-series split), para assegurar que o modelo se generalize bem para novos dados.

### 4. Integração de Dados Contextuais

- Incorporar variáveis externas, como tendências de hashtags ou dados demográficos dos seguidores, para agregar contexto às previsões de engajamento.

## Referências

- Documentação do **Scikit-Learn**: <https://scikit-learn.org>
- **Kaggle**: Fonte do conjunto de dados e benchmark de práticas em modelagem preditiva.