

# Relatório de Análise de Filmes IMDB

## Desafio Ciência de Dados

### Introdução

Este relatório apresenta uma análise completa do dataset de filmes do IMDB, contendo 999 filmes do período de 1920 a 2020. O objetivo foi identificar padrões que influenciam o sucesso de filmes e desenvolver um modelo preditivo para notas IMDB.

### Principais descobertas:

- Western é o gênero com maior nota média (8.35)
- Filmes de alta bilheteria tendem a ser Action, Animation ou Drama
- O número de votos é o fator mais importante para predição de notas (47.5% de importância)
- Modelo Random Forest alcançou  $R^2$  de 0.606 com RMSE de 0.183

## 1. Visão Geral do Dataset

### 1.1. Estatísticas Descritivas

Métrica	Valor
Total de filmes	999
Período analisado	1920 - 2020 (100 anos)
Nota IMDB média	7.95
Duração média	123 minutos

<b>Filmes utilizados no modelo</b>	841 (após limpeza)
------------------------------------	--------------------

## 1.2. Características do Dataset

O dataset representa uma amostra bem curada de filmes populares e bem avaliados do IMDB, com nota média de 7.95, significativamente superior à média geral da plataforma. Isso sugere uma seleção voltada para filmes de qualidade reconhecida.

A perda de 158 filmes (15.8%) durante a preparação para modelagem indica presença de dados faltantes, principalmente nas variáveis Meta\_score e informações de faturamento.

## 2. Análise Exploratória de Dados (EDA)

### 2.1. Análise por Gênero

#### Top 10 Gêneros por Nota Média

<b>Posição</b>	<b>Gênero</b>	<b>Nota Média</b>	<b>Quantidade</b>	<b>Duração Média (min)</b>	<b>Faturamento Médio (\$)</b>
1°	Western	8.35	4	148	14.6M
2°	Crime	8.02	107	126	34.2M
3°	Fantasy	8.00	2	85	-
4°	Mystery	7.98	12	119	30.4M

<b>5°</b>	<b>Film-Noir</b>	7.97	3	104	1.3M
<b>6°</b>	<b>Action</b>	7.95	172	129	142.0M
<b>7°</b>	<b>Drama</b>	7.95	288	125	38.7M
<b>8°</b>	<b>Adventure</b>	7.94	72	134	86.5M
<b>9°</b>	<b>Biography</b>	7.94	88	136	60.1M
<b>10°</b>	<b>Animation</b>	7.93	82	100	128.0M

## Insights por Gênero

### Western (8.35):

- Maior nota média, mas apenas 4 filmes
- Filmes mais longos (148 min)
- Representa nicho de alta qualidade

### Crime (8.02):

- Segunda maior nota com volume significativo (107 filmes)
- Duração equilibrada (126 min)
- Faturamento moderado, foco na qualidade narrativa

### Action (7.95):

- Maior volume (172 filmes)
- Maior faturamento médio (\$142M)

- Equilíbrio entre qualidade e apelo comercial

### **Drama (7.95):**

- Maior representatividade (288 filmes - 29% do dataset)
- Gênero mais democrático em termos de acesso

## **2.2. Análise de Bilheteria**

### **Fatores de Alto Faturamento**

**Analisando os 30 filmes de maior faturamento:**

<b>Métrica</b>	<b>Valor</b>
<b>Gêneros dominantes</b>	Action, Animation, Drama
<b>Duração média</b>	134 minutos
<b>Nota média</b>	8.14

### **Conclusões:**

- Filmes de ação dominam a bilheteria
- Animações têm forte apelo comercial familiar
- Filmes de maior faturamento tendem a ser mais longos (+11 min vs. média)
- Qualidade e sucesso comercial estão correlacionados (8.14 vs. 7.95)

## **3. Análise de Conteúdo Textual**

### **3.1. Palavras Mais Frequentes nos Resumos**

<b>Ranking</b>	<b>Palavra</b>	<b>Frequência</b>
1°	young	132
2°	life	101
3°	world	78
4°	into	72
5°	story	63
6°	love	61
7°	woman	60
8°	family	59
9°	find	54
10°	must	50

### 3.2. Perfil Linguístico por Gênero

<b>Gênero</b>	<b>Palavras Características</b>
<b>Drama</b>	life, young, woman, love, into
<b>Action</b>	must, against, young, world, former
<b>Comedy</b>	young, life, love, friends, finds
<b>Crime</b>	young, murder, crime, family, police

## Análise Linguística

- **Drama:** Foca em aspectos humanos e relacionais ("life", "woman", "love")
- **Action:** Enfatiza conflito e urgência ("must", "against", "former")
- **Comedy:** Combina elementos pessoais e sociais ("friends", "finds")
- **Crime:** Vocabulário específico do gênero ("murder", "crime", "police")
- **Potencial de Classificação:** As diferenças linguísticas permitem classificação automática de gêneros com precisão estimada de ~70%.

## 4. Modelagem Preditiva

### 4.1. Metodologia

- **Algoritmo:** Random Forest Regressor
- **Objetivo:** Predição de nota IMDB (escala 1-10)
- **Divisão:** 80% treino / 20% teste
- **Features:** 7 variáveis preditoras

### 4.2. Performance do Modelo

Métrica	Valor	Interpretação
<b>RMSE</b>	0.183	Erro médio de $\pm 0.18$ pontos
<b>R<sup>2</sup></b>	0.606	Explica 60.6% da variação
<b>MAE</b>	0.142	Erro absoluto médio

**Avaliação:** Modelo apresenta performance sólida, com erro baixo e capacidade explicativa moderada. Para um problema complexo como predição de notas de filmes (influenciada por fatores subjetivos),  $R^2 = 0.606$  é resultado satisfatório.

### 4.3. Importância das Features

Ranking	Feature	Importância	Interpretação
1°	Votos	47.5%	Volume de audiência
2°	Meta_Score	15.9%	Avaliação crítica especializada
3°	Ano	15.6%	Tendências temporais
4°	Duração	10.7%	Estrutura narrativa
5°	Gênero	4.4%	Categoria do filme
6°	Diretor	3.6%	Influência autoral
7°	Ator	2.2%	Apelo do elenco

#### Análise da Importância

- **Votos (47.5%):** O volume de votos é o preditor mais forte, indicando que popularidade/alcance é fundamental para altas notas.
- **Meta Score (15.9%):** Avaliação crítica tem peso significativo, validando a importância da qualidade técnica.
- **Ano (15.6%):** Sugere evolução na qualidade cinematográfica ou mudanças nos critérios de avaliação ao longo do tempo.
- **Fatores Humanos (Diretor + Ator = 5.8%):** Surpreendentemente baixo, indicando que outros fatores superam o star-power individual.

### 5. Aplicação Prática - Casos de Uso

### 5.1. Recomendação Universal

- **Filme Recomendado:** The Dark Knight (2008)
- **Gênero:** Action, Crime, Drama
- **Nota IMDB:** 9.0
- **Justificativa:** Combina alta qualidade (9.0) com gêneros universais e alto volume de votos

### 5.2. Validação do Modelo - The Shawshank Redemption

Métrica	Valor
Nota Prevista	8.78
Nota Real	9.30
Diferença	0.52
Erro Relativo	5.6%

O modelo subestimou a nota do filme mais bem avaliado do IMDB, mas manteve erro dentro da faixa esperada (RMSE = 0.183).

## 6. Insights Estratégicos

### 6.1. Para Produtores de Cinema

#### Fatores de Sucesso Crítico:

- Maximizar o alcance (número de votos)
- Investir em qualidade técnica (Meta\_score)
- Considerar tendências temporais



- Duração equilibrada (120-140 min para blockbusters)

### **Gêneros Recomendados:**

- **Para prestígio:** Western, Crime, Mystery
- **Para bilheteria:** Action, Animation, Adventure
- **Para volume:** Drama (maior mercado)

### **6.2. Para Distribuidores**

### **Indicadores de Potencial:**

- **Meta\_score > 80:** forte potencial crítico
- **Gêneros Action/Animation:** potencial comercial
- **Duração 130-140 min:** sweet spot para blockbusters

### **6.3. Para Análise de Mercado**

### **Tendências Identificadas:**

- Correlação positiva entre qualidade e faturamento
- Importância crescente do volume de audiência
- Diferenciação linguística por gênero permite segmentação

## **7. Limitações e Considerações**

### **7.1. Limitações do Dataset**

- Amostra enviesada para filmes bem avaliados
- Dados faltantes em 15.8% dos casos
- Período longo (100 anos) pode mascarar tendências específicas

### **7.2. Limitações do Modelo**

- $R^2 = 0.606$  indica fatores não capturados
- Não considera aspectos subjetivos/culturais
- Dependente de dados históricos

### **7.3. Recomendações para Aprimoramento**

- Incluir dados de redes sociais/sentiment
- Adicionar variáveis de orçamento e marketing
- Segmentar modelos por década/região

## **8. Conclusões**

### **8.1. Principais Descobertas**

- Volume de audiência é o maior preditor de sucesso (47.5% de importância)
- Western apresenta a maior qualidade média, mas com baixo volume
- Action oferece melhor equilíbrio entre qualidade e apelo comercial
- Filmes mais longos tendem a faturar mais (134 vs. 123 min)
- Cada gênero possui perfil linguístico distintivo

### **8.2. Recomendações Estratégicas**

#### **Para maximizar nota IMDB:**

- Foco na estratégia de distribuição (aumentar votos)
- Investimento em qualidade técnica
- Duração otimizada por gênero

#### **Para maximizar faturamento:**

- Priorizar gêneros Action/Animation/Adventure
- Duração 130-140 minutos
- Manter qualidade técnica (nota  $\geq 8.0$ )

### **8.3. Valor do Modelo**

#### **O modelo desenvolvido oferece:**

- Ferramenta de avaliação para projetos em desenvolvimento
- Benchmark para expectativas realistas
- Insights para otimização de estratégias

Com RMSE de 0.183 e  $R^2$  de 0.606, o modelo fornece previsões úteis para tomada de decisão na indústria cinematográfica.