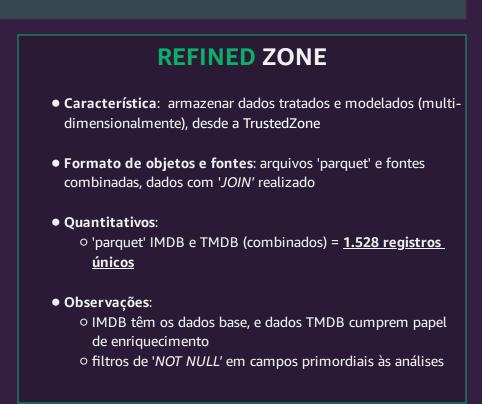
#### **Pipeline dos Dados**

## ZONE • Característica: armazenar dados em estado bruto • Formato de objetos e fontes\*: arquivos .csv e .json, e dados armazenados separados por fonte • Quantitativos: o 'csv' IMDB = 244.543 registros únicos ○ 'json' TMDB = 44.641 registros únicos Observações: o 'csv' fornecido com todos os gêneros de filmes disponíveis ○ 'json' extraído (via API) já com filtro para manter apenas filmes do gênero Crime e/ou Guerra

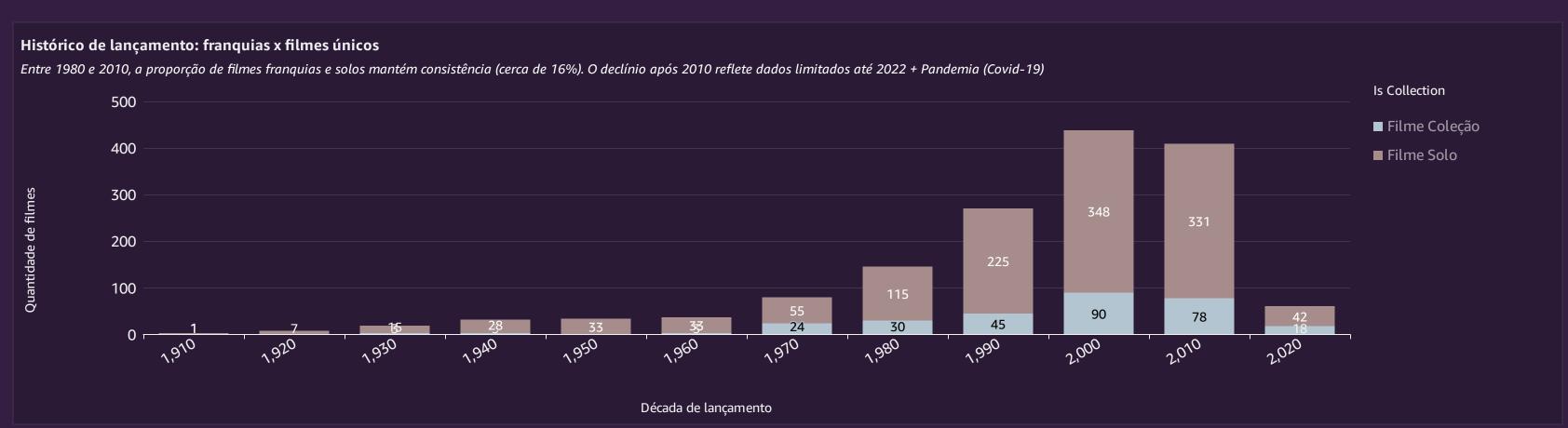
## **TRUSTED ZONE** • Característica: armazenar tratados, desde a RawZone • Formato de objetos e fontes: arquivos 'parquet' e dados armazenados ainda separados por fonte • Quantitativos: o 'parquet' IMDB = 29.128 registros únicos o 'parquet' TMDB = 44.641 registros Observações: o dados do IMDB foram filtrados para gêneros - manter apenas registros que contém Crime e/ou Guerra o dados TMDB tratados sem mais filtros

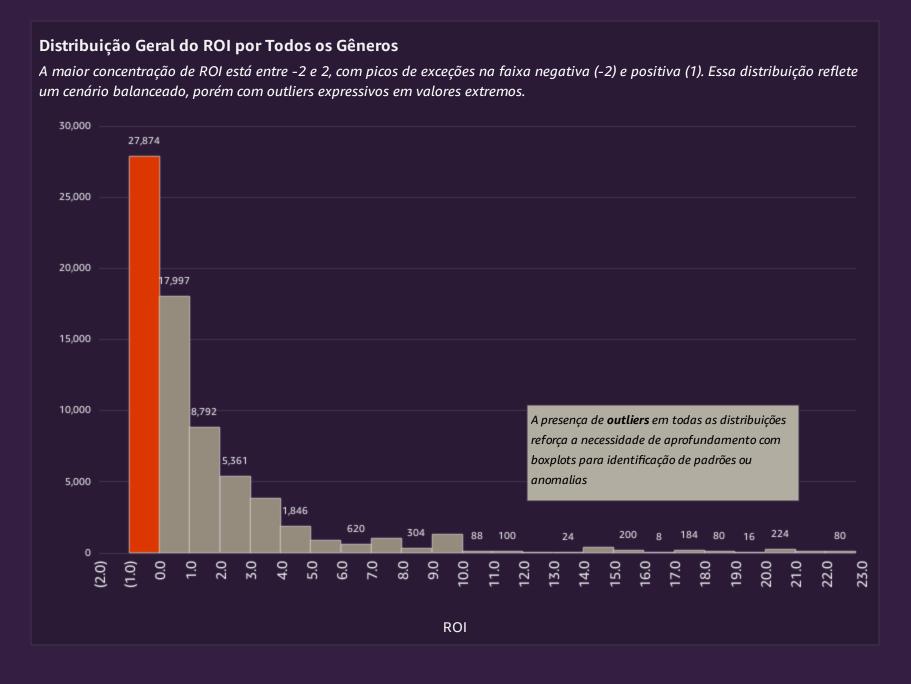


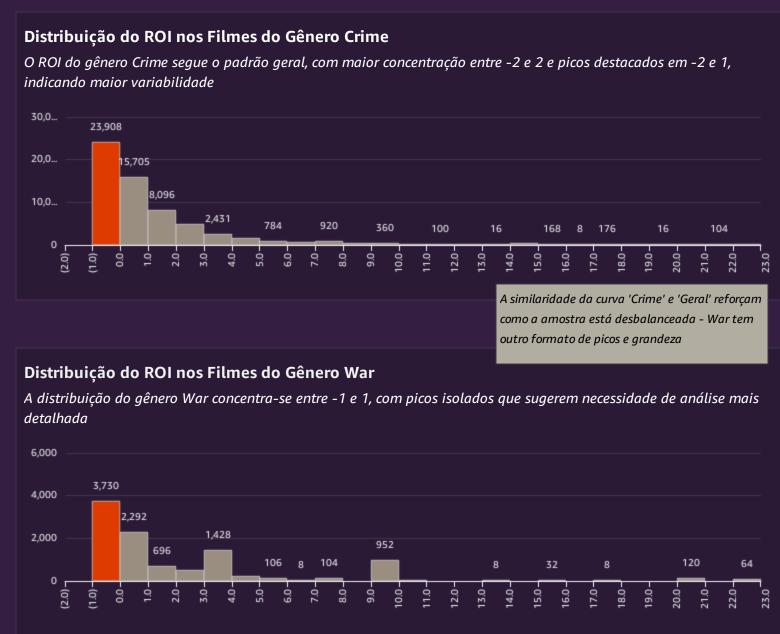
\*Fontes dos dados: 1) IMDB - Internet Movie Database (https://www.imdb.com) | 2) TMDB - The Movie Database (https://www.themoviedb.org/)

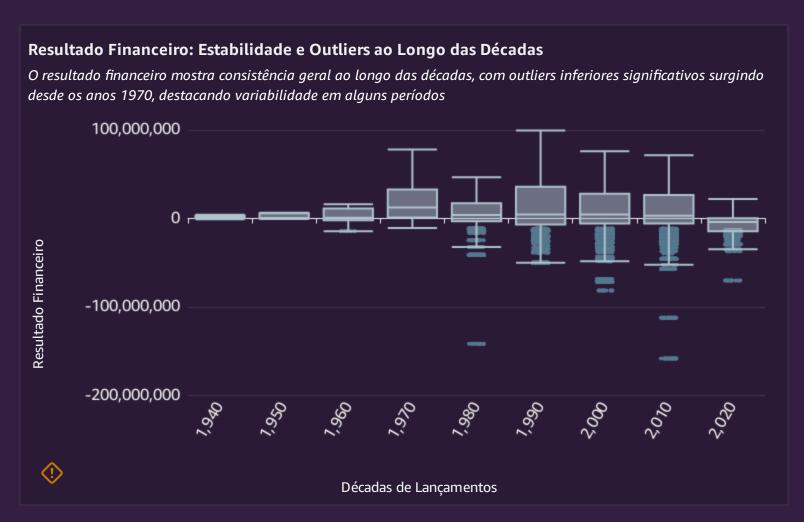
### Verificação de consistência da amostra (estatística descritiva básica)

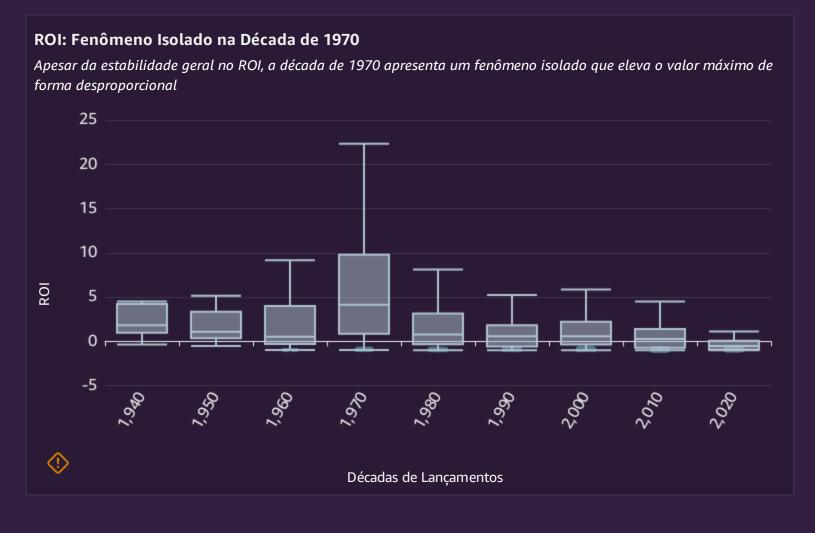




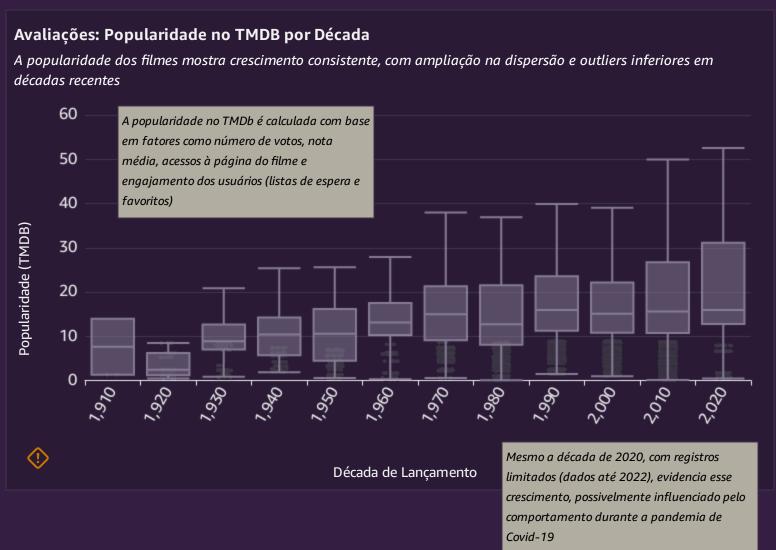












# Conclusões sobre a Consistência Estatística + Tratamentos Complementares Aplicados

# Conclusões

- Os dados da RefinedZone, como estão, podem gerar interpretações enviesadas e conclusões imprecisas. Há uma necessidade evidente de tratamentos adicionais para assegurar maior robustez analítica.
- Massas de outliers foram identificadas, e sua análise separada, por meio de clustering ou outros métodos, poderia fornecer insights úteis. Entretanto, essas métricas devem ser excluídas das análises centrais para evitar distorções.
- Para cumprir nosso objetivo de gerar análises históricas consistentes que forneçam padrões e insights estratégicos, optaremos por descartar outliers, reduzindo os riscos associados às métricas financeiras e de avaliação do público.

inferiores que não interessam como objetivo da análise).

- **Tratamentos Complementares** • Filtro Temporal: focaremos em produções lançadas entre 1980 e 2020, evitando picos ou vales (extremos) de dados, e que não reflitam o contexto cultural e histórico atual.
  - Aplicação de Filtros Estatísticos: adotaremos apenas os dados entre os quartis Q1 e Q3 dos boxplots das métricas de Média de Notas, ROI e Resultado Financeiro, assegurando uma base mais estável para análise, desprezando outliers (especialmente os