

# UMA ANÁLISE SOBRE FILMES USANDO R E RSTUDIO

Felipe Mateus Evangelista<sup>1</sup>, Filipe Pereira da Silva<sup>2</sup>, William Salvador Antonelli<sup>3</sup>

## 1 INTRODUÇÃO E OBJETIVOS

O principal objetivo deste resumo foi fazer uma análise estatística sobre dados de filmes. A partir de uma base de dados, buscamos investigar informações sobre algumas variáveis dos filmes, classificar e descrever estas variáveis. Também buscamos criar tabelas sobre estes dados, criar gráficos para visualizar informações e entender como esses dados se comportam, criar métricas com os dados e discutir os resultados alcançados.

## 2 MATERIAIS E MÉTODOS

A base de dados escolhida foi uma base com dados relacionados a filmes. O conjunto de dados foi removido do site da Kaggle, que possui uma base do TMDb (The Movie Database) com 5 mil filmes, mas utilizamos somente uma amostra de 350 filmes. Nesta base existiam algumas variáveis que foram desconsideradas, sendo assim filtramos ela na sintaxe R para analisarmos somente os dados desejados. Sendo excluídas as variáveis, homepage, id, keywords, overview, spoken\_languages, status, tagline, original\_title, runtime. Outrossim, outra medida necessária foi ignorar alguns filmes cujo os dados não estavam satisfatórios, sendo por ter dados nulos, zerados ou incongruentes.

As variáveis definidas para análise foram:

- *Custo(budget)*: Informa o valor gasto na construção do filme. Sendo classificada como uma variável quantitativa contínua.  
Foi escolhida porque permite comparar o orçamento com a receita gerada e calcular a rentabilidade do filme. Essa informação é importante para avaliar se o filme foi financeiramente bem-sucedido ou não.
- *Gênero(genre)*: Informa qual é o gênero do filme. Sendo classificada como uma variável qualitativa nominal
- *Língua original(original language)*: Informa a língua original na qual o filme foi produzido. Sendo classificada uma variável qualitativa nominal.
- *Título original (original title)*: Informa o título do filme. Sendo classificada como uma variável qualitativa nominal
- *(Popularidade)popurlaty*: Informa a quantidade de pessoas que assistiram o filme no cinema. Sendo classificada como uma variável quantitativa discreta.
- *Produtora(production\_company)*: Informa a produtora do filme. Sendo classificada como uma variável qualitativa nominal.
- *País(production\_country)*: Informa o país que produziu o filme. Sendo classificada como uma variável qualitativa nominal.

- *Data de lançamento(release\_date)*: Informa a data de lançamento do filme. Sendo classificada como uma variável qualitativa ordinal
- *Receita(revenue)*: Informa a receita do filme. Sendo classificada como uma variável quantitativa contínua.

Foi escolhida porque é importante para avaliar o retorno sobre o investimento e é importante para entender o impacto financeiro da produção.

- *Média dos votos(vote\_average)*: Informa a média da nota do filme. Sendo classificada como uma variável quantidade contínua. É quantitativa contínua.
- *Quantidade de votos(vote\_count)*: Informa a quantidade de votos que o filme recebeu. Sendo classificada com uma variável quantitativa discreta.

### 3 RESULTADOS E DISCUSSÕES

#### 3.1 Tabelas

Primeiro vamos analisar o comportamento de algumas variáveis qualitativas nominais por meio de tabelas de frequência, são elas, Língua original (original language), Gênero (genre), Produtora (production\_company), País (production\_country).

Figura 1 - Tabela de frequência para Língua original.

```

R 4.3.3 ~ /Desktop/UFSC/second_semester/probability_statistics/first_homework/
> #Tabela com os dados usando da língua usada nos filmes
> table(movies_dataset_filtred$original_language)

de en es fr hi id ja te zh
1 338 2 2 1 1 3 1 1
>

```

Conforme podemos ver na imagem acima, é notável uma predominância da língua inglesa, sendo ela majoritariamente compondo a maioria das linguagens dos filmes que são criados de acordo com a nossa amostra.

Figura 2 - Tabela de frequência para Gênero

```

R 4.3.3 ~ /Desktop/UFSC/second_semester/probability_statistics/first_homework/
> #Tabela com os dados usando do gênero usada nos filmes
> table(movies_dataset_filtred$genre)

Action Adventure Animation Comedy Crime Documentary Drama Family Fantasy Horror Music Mystery Romance Science Fiction Thriller
72 38 3 63 12 6 78 4 14 25 1 3 6 4 13
War Western
3 5
>

```

Diferente da Figura 1, onde os dados de língua ficavam concentrados em uma única língua, na tabela sobre gênero podemos ver uma melhor distribuição dos dados, ainda assim é possível notar que os gêneros Drama, Comédia e Ação tem uma maior quantidade de filmes representados.

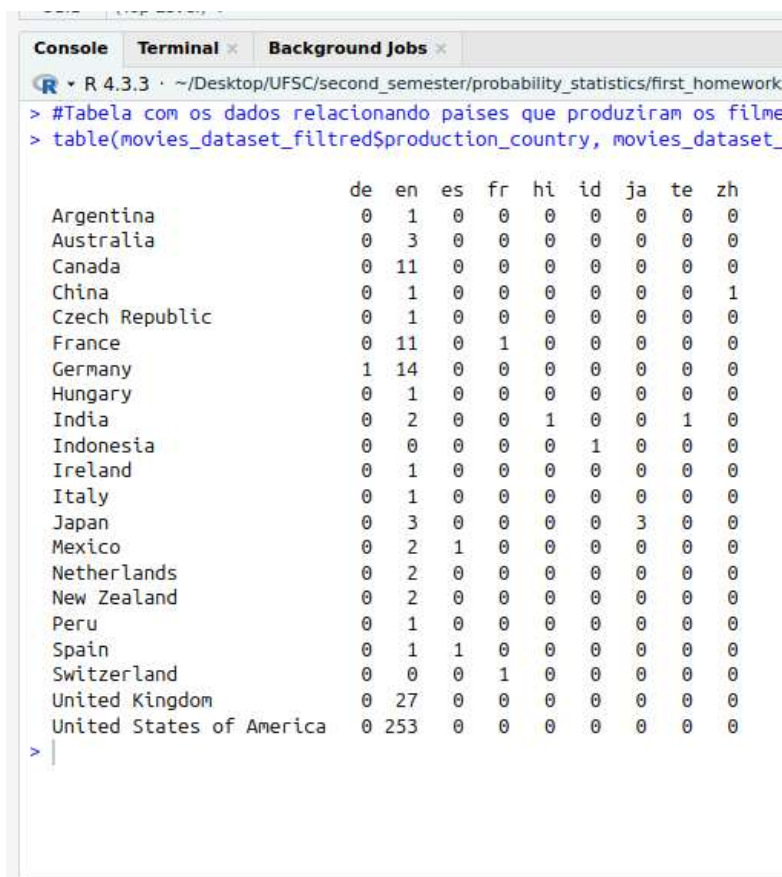
Figura 3 - Tabela de frequência para Produtora



Produtora	Frequência
2929 Productions	1
Ames Ra Films	1
Bac Films	1
Boran Entertainment Inc.	1
Celluloid Dreams	1
Conundrum Entertainment	1
Dark Horse Entertainment	1
Dimension Films	1
Duplass Brothers Productions	1
Films	1
Fox Searchlight Pictures	1
Graviter Productions	1
Itaca Films	1
Kennedy/Marshall Company, The	1
Legendary Pictures	1
Mike Zoss Productions	1
Nick Wechsler Productions	1
Plan B Entertainment	1
Revolution Studios	1
Scott Rudin Productions	1
Sony Pictures Classics	1
Sunswep Entertainment	1
Touchstone Pictures	1
Universal Pictures	1
Warner Bros.	9
40 Acres & A Mule Filmworks	1
American Zoetrope	1
BBC	1
Brooksfilms	1
Central Cinema Company Film	1
Cookout Productions	1
Davis Entertainment	1
Dinovi Pictures	1
Electric Entertainment	1
Fine Line Features	1
France 3 Cinéma	1
Grindstone Entertainment Group	1
Jerry Bruckheimer Films	1
Killer Films	1
Lions Gate Films	1
Miramax Films	1
NULL	1
Playtone	1
Rhombus Media	1
Screen Gems	1
Spyglass Entertainment	1
Tapestry Films	1
Tree Line Films	1
Universal Studios	1
Warner Bros. Pictures	1
A Loopy Production LLC	1
Anchor Bay Films	1
BBC Films	1
Canal+	1
Charles K. Feldman Group	1
Corrino Media Corporation	1
Davis-Films	1
DreamWorks	1
Ennott/Furla Films	1
Flicks Motion Pictures	1
France 3 Cinéma	1
Imagine Entertainment	1
Jerry Weintraub Productions	1
Lakeshore Entertainment	1
Lionsgate	1
National Geographic Entertainment	1
Original Film	1
Punch 21 Productions	1
River Road Entertainment	1
Screen Gems, Inc.	1
Studio Babelsberg	1
The Weinstein Company	1
Trilogy Entertainment Group	1
USA Films	1
Warner Brothers/Seven Arts	1
Alcon Entertainment	1
Applan Way	1
Beijing New Picture Film Co. Ltd.	1
Cannon Films	1
Cineplex Odeon Films	1
Crossbow Productions	1
DC Comics	1
DreamWorks Animation	1
Eon Productions	1
Focus Features	1
Franchise Pictures	1
Incentive Filmed Entertainment	1
Journeyman Pictures	1
Larger Than Life Productions	1
Live Entertainment	1
New Amsterdam Entertainment	1
Paramount Pictures	1
Punch Productions	1
Robert Simonds Productions	1
Shaw Brothers	1
Studio Ghibli	1
ThinkFilm	1
Tristar Pictures	1
Twentieth Century Fox Film Corporation	1
Vertigo Films	1
Worldview Entertainment	1
Alliance Atlantis Communications	1
Arka Media Works	1
Blumhouse Productions	1
Castle Rock Entertainment	1
Columbia Pictures	1
Columbia Pictures Corporation	1
Crystal Sky Worldwide	1
Destination Films	1
DreamWorks SKG	1
Estudios Picasso	1
Fox 2000 Pictures	1
Gaumont	1
Ingenious Film Partners	1
Keith Barish Productions	1
Laurel Group	1
Marvel Studios	1
New Line Cinema	1
Paramount Vantage	1
Red Envelope Entertainment	1
Ruby in Paradise	1
Sidney Kimmel Entertainment	1
StudioCanal	1
This Is That Productions	1
Village Roadshow Pictures	1
X-Film Creative Pool	1
Anblin Entertainment	1
Artisan Entertainment	1
Bold Films	1
CBS Films	1
Cube Vision	1
Di Bonaventura Pictures	1
Dune Entertainment	1
Filmnet Films	1
Fox Entertainment Group	1
Gracie Films	1
Ingenious Media	1
Kennedy Miller Productions	1
Lawrence Bender Productions	1
Metro-Goldwyn-Mayer (MGM)	1
New World Pictures	1
Participant Media	1
Regency Enterprises	1
Saturn Films	1
Silver Pictures	1
Summit Entertainment	1
Tig Productions	1
United Artists	1
Walt Disney Pictures	12

Na tabela para a variável produtora podemos notar uma heterogeneidade dos dados, ficando eles bem espalhados e tendo um número bem alto de produtoras, o que mostra que no mundo dos filmes temos uma grande variedade de empresas que os produzem.

Figura 4 - Tabela de frequência para línguas e países



País	de	en	es	fr	hi	id	ja	te	zh
Argentina	0	1	0	0	0	0	0	0	0
Australia	0	3	0	0	0	0	0	0	0
Canada	0	11	0	0	0	0	0	0	0
China	0	1	0	0	0	0	0	0	1
Czech Republic	0	1	0	0	0	0	0	0	0
France	0	11	0	1	0	0	0	0	0
Germany	1	14	0	0	0	0	0	0	0
Hungary	0	1	0	0	0	0	0	0	0
India	0	2	0	0	1	0	0	1	0
Indonesia	0	0	0	0	0	1	0	0	0
Ireland	0	1	0	0	0	0	0	0	0
Italy	0	1	0	0	0	0	0	0	0
Japan	0	3	0	0	0	0	3	0	0
Mexico	0	2	1	0	0	0	0	0	0
Netherlands	0	2	0	0	0	0	0	0	0
New Zealand	0	2	0	0	0	0	0	0	0
Peru	0	1	0	0	0	0	0	0	0
Spain	0	1	1	0	0	0	0	0	0
Switzerland	0	0	0	1	0	0	0	0	0
United Kingdom	0	27	0	0	0	0	0	0	0
United States of America	0	253	0	0	0	0	0	0	0

Figura 5 - Tabela de frequência para Países

```

R - R 4.3.3 - ~/Desktop/KFSCsecond_semester/probability_statistics/first_homework/ >
#Tabela com os dados usando dos países que produziram os filmes
> table(movies_dataset_filtred$production_country)

```

Argentina	1	Australia	3	Canada	11	China	2	Czech Republic	1	France	12	Germany	15	Hungary	1	India	4
Indonesia	1	Ireland	1	Italy	1	Japan	6	Mexico	3	Netherlands	2	New Zealand	2	Peru	1	Spain	2
Switzerland	1	United Kingdom	27	United States of America	253												

```

> |

```

A tabela de países segue a mesma dinâmica da tabela de línguas dos filmes, podemos ver que a maioria dos dados ficam concentrados em um país, os Estados Unidos, tendo os outros países um número muito menor de frequência.

Podemos também analisar estas variáveis conjuntas para poder ver correlações entre os dados. Por exemplo, ao analisarmos a figura 4, notamos que a quantidade de filmes em língua inglesa se dá pelo fato da maioria dos filmes serem produzidos pelos Estados Unidos. Cabe também destacar um fato curioso, a França e a Alemanha, apesar de terem línguas próprias, francês e alemão, respectivamente, possuem mais filmes feitos em língua inglesa do que em suas línguas nativas.

Por último, podemos analisar a variável gênero do filme por país, como nas figuras abaixo. Na figura 6 estamos analisando os gêneros dos filmes produzidos pela França e na Paramount Pictures, em ambas podemos perceber que os gêneros são bem distribuídos, sem haver gêneros que se sobressaíam expressivamente perante aos outros.

Figura 6 - Tabela de frequência dos gênero dos filmes produzidos pela França

```
Console Terminal x Background Jobs x
R 4.3.3 · ~/Desktop/UFSC/second_semester/probability_statistics/first_home
> france_movies = movies_dataset_filtred[movies_dataset_filtred$production_country == "France", ]
> table(france_movies$production_country, france_movies$genre)
```

	Action	Adventure	Drama	Thriller	War	
France	4		2	3	1	2

```
> |
```

Figura 7 - Tabela de frequência dos gênero dos filmes produzidos pela Paramount Pictures

**Console** **Terminal** **Background Jobs**

R 4.3.3 · ~/Desktop/UFSC/second\_semester/probability\_statistics/first\_homework/

```
> #Tabela com os dados gêneros dos filmes produzidos pela Paramount Pictures
> paroumunt_movies = movies_dataset_filtred[movies_dataset_filtred$production_co
> table(paroumunt_movies$production_company, paroumunt_movies$genre)
```

	Action	Adventure	Comedy	Crime	Drama	Horror	Western
Paramount Pictures	5	6	6	1	9	2	1

```
> |
```

### 3.2 Gráficos

Com a base de dados também buscamos criar gráficos para conseguir visualizar o comportamento dos dados e como eles se relacionam individualmente ou agrupados, os primeiro gráficos a serem analisados serão os de densidade.

Figura 8 - Gráfico de densidade para Orçamento (budget)

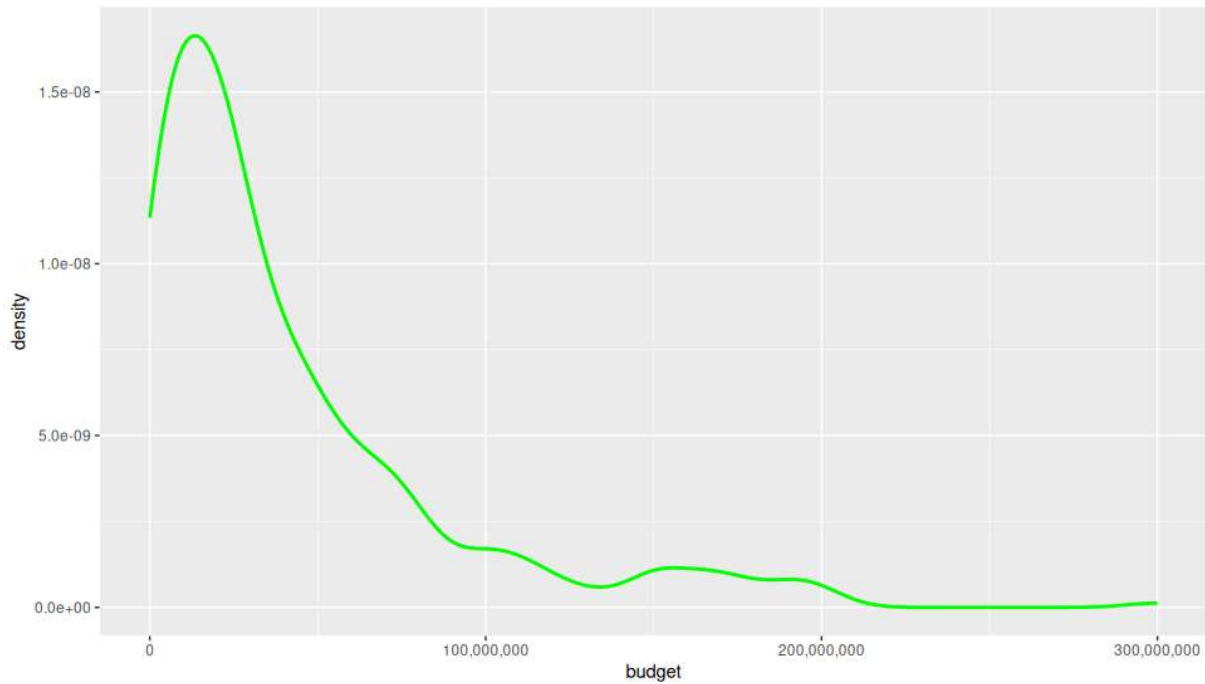


Figura 9 - Gráfico de densidade para Receita (revenue)

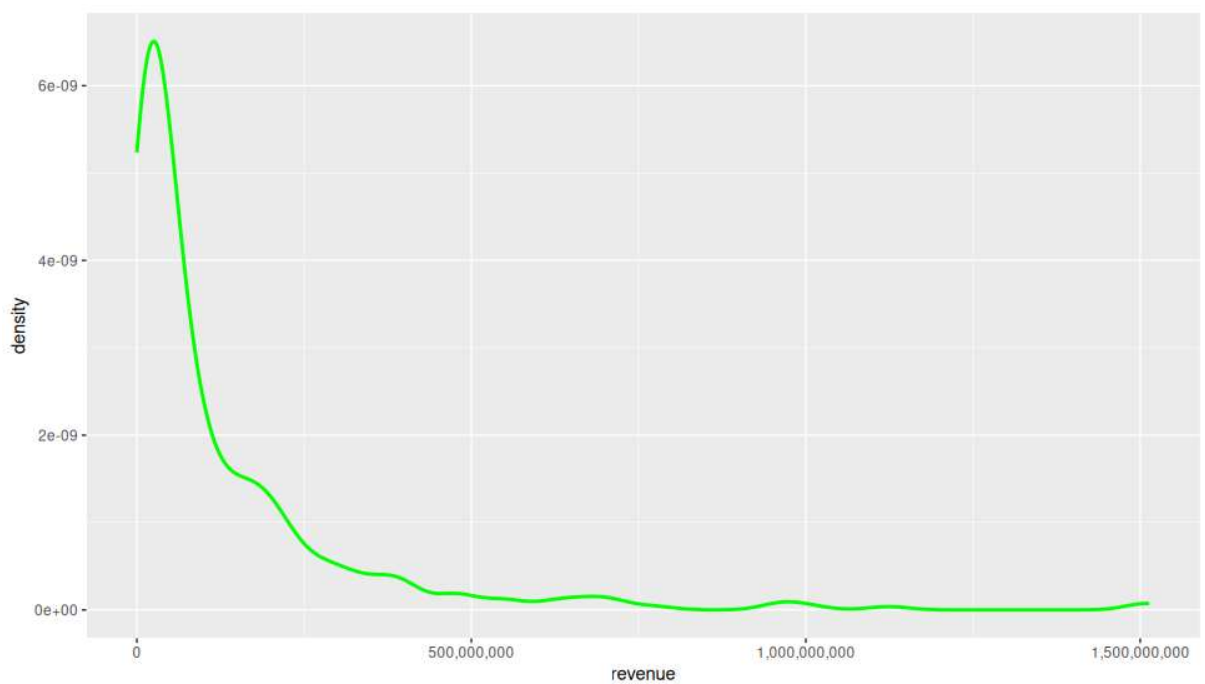
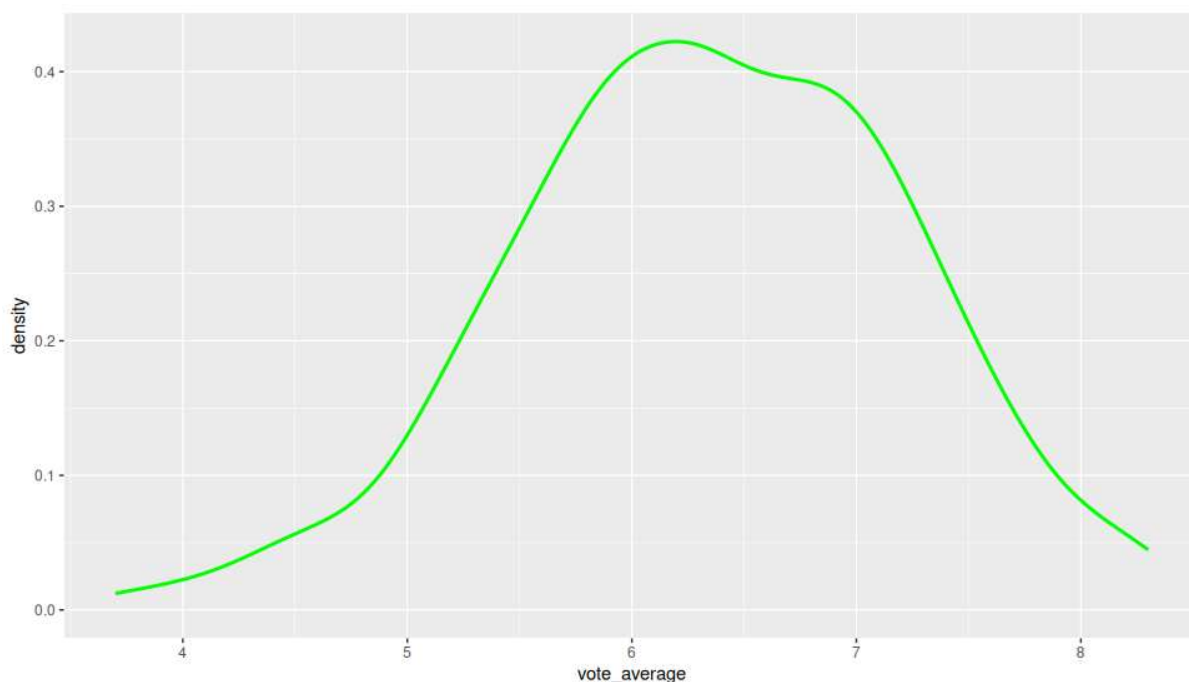
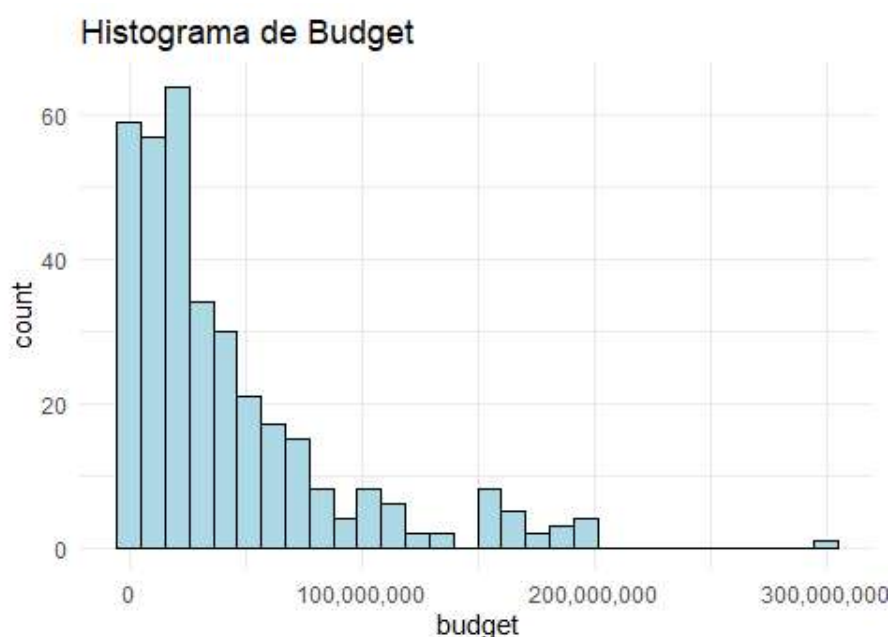


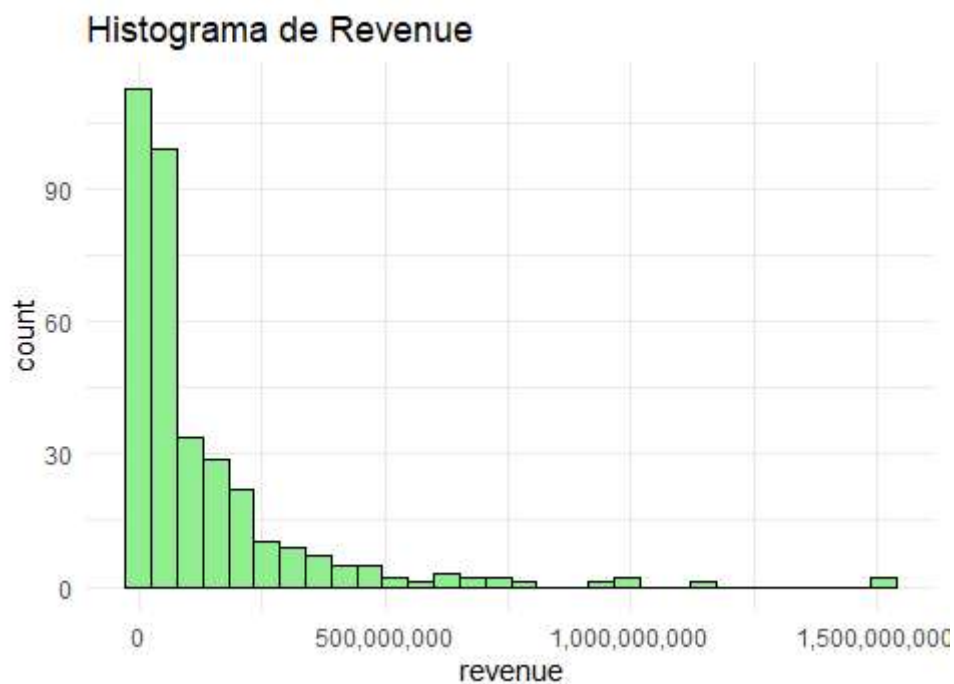
Figura 10 - Gráfico de densidade para Notas dos filmes (vote\_average)



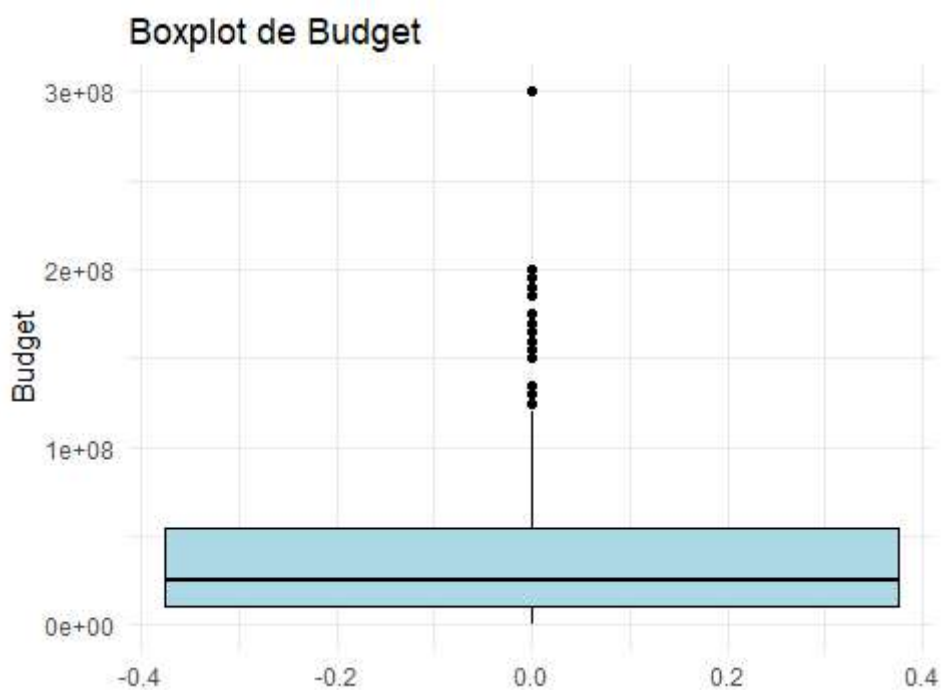
Com os gráficos das figuras 8, 9 e 10 podemos perceber alguns comportamentos, como, nos gráficos de orçamento e receita temos uma maior pico de densidade para filmes que tem um baixo orçamento e filmes com uma baixa receita, também percebemos a presença de outliers, o que faz com que estes gráficos fiquem com as “caudas” mais esticadas e que podem influenciar métricas como cálculo da média. Já no gráfico da figura 10, temos um comportamento diferente, nesse gráfico podemos ver que as notas são bem distribuídas entre a média central, fazendo com que o pico fique próximo no centro e também mostrando que os dados estão bem distribuídos pelo gráfico, sem um local central de concentração.



O gráfico acima mostra que a maior parte dos filmes têm orçamentos baixos, poucos filmes têm orçamentos muito altos, o que gera a cauda longa à direita.

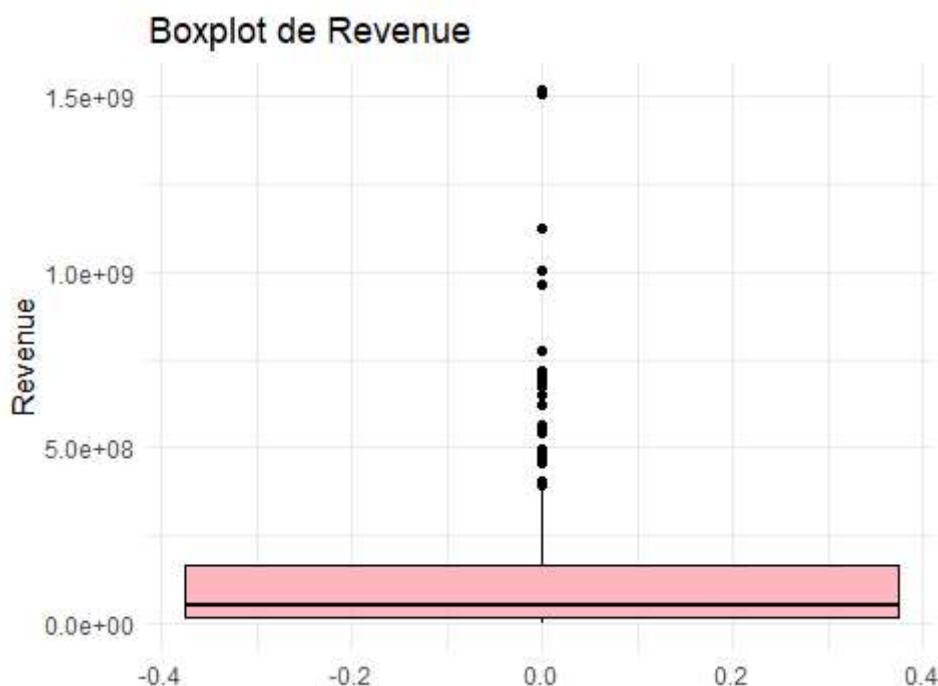


O gráfico acima mostra que a grande maioria dos filmes gera receitas menores. Existem alguns poucos filmes com faturamentos altíssimos, a distribuição é assimétrica à direita.





O gráfico acima mostra que a maioria dos valores de budget está concentrada entre valores relativamente baixos. Existem muitos outliers, o que indica que há filmes com orçamentos muito maiores que a média. A mediana (linha preta dentro da caixa) é bem próxima de 0 comparada ao máximo, o que mostra que a maioria dos filmes tem orçamento baixo em relação a poucos filmes com orçamento muito alto. O budget é altamente assimétrico para a direita. Há uma grande desigualdade no orçamento dos filmes: poucos filmes gastam muito, a maioria gasta pouco.



O gráfico acima mostra que assim como o orçamento, o revenue também apresenta muitos outliers. A maioria dos filmes arrecada valores baixos, mas alguns poucos arrecadam muito. A mediana também é baixa comparada ao valor máximo. O faturamento dos filmes também é muito desigual, com a maior parte dos filmes faturando pouco e poucos filmes faturando muito. A distribuição é assimétrica à direita, indicando que existem filmes que se destacam muito em receita.

### 3.3 Medidas de Resumo

```
> summary(movies_dataset_filtred$budget)
```

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
4	10.875.000	25.000.000	41.818.896	54.750.000	300.000.000	

```
[1] "Variância: 2.160.957.918.892.042"
```

```
[1] "Desvio Padrão: 46.486.104,5785947"
```

```
[1] "Coeficiente de Variação: 111,16052600433 %"
```



O desvio padrão é muito alto em relação à média, o que indica que os valores de budget são muito dispersos. O coeficiente de variação acima de 100% mostra que a variabilidade é maior que a média. Como a média é maior que a mediana, a distribuição é assimétrica à direita.

```
> summary(movies_dataset_filtred$revenue)
```

```
Min. 1st Qu.  Median    Mean  3rd Qu.    Max.
 7 17.520.000 50.770.000 128.200.000 164.300.000 1.514.000.000
```

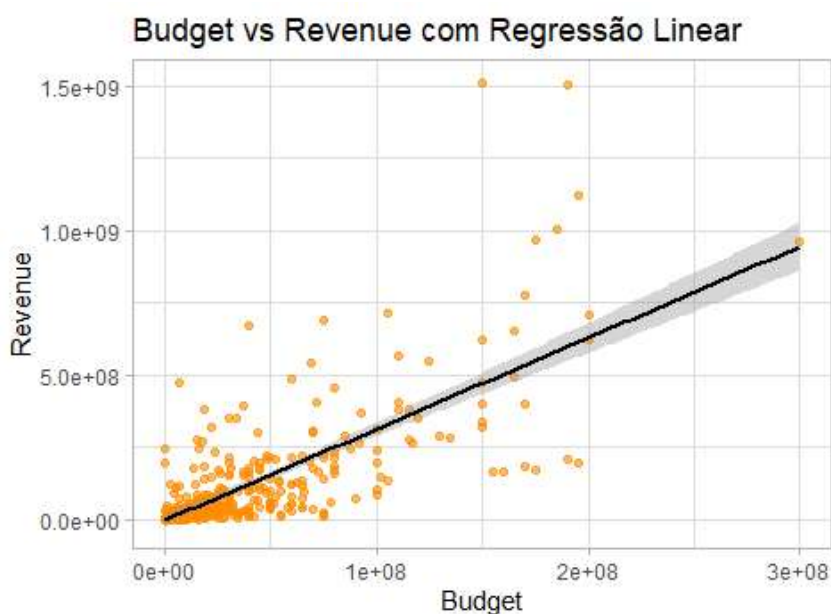
```
[1] "Variância: 41.004.528.060.731.096"
```

```
[1] "Desvio Padrão: 202.495.748,253466"
```

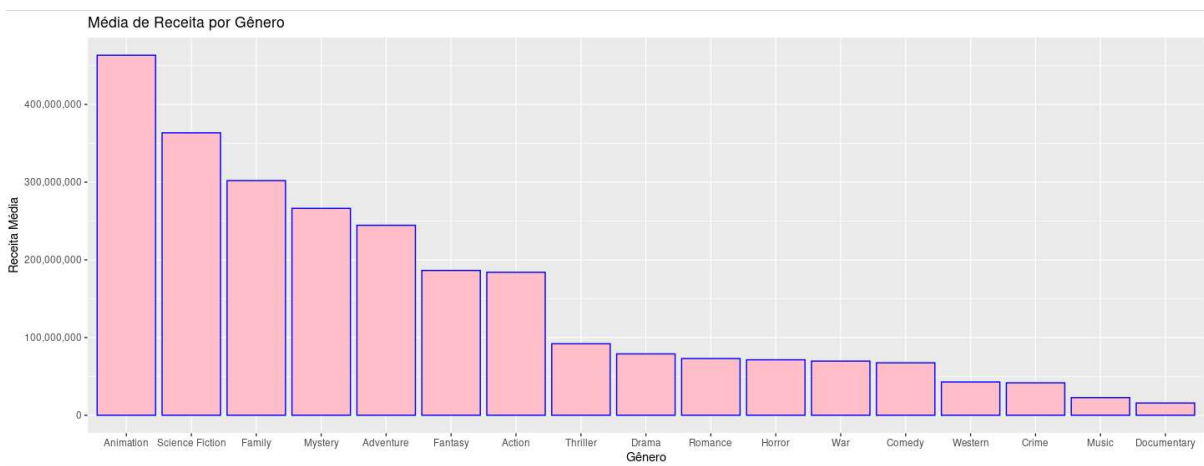
```
[1] "Coeficiente de Variação: 157,950738665499 %"
```

O desvio padrão é ainda mais alto que no caso do budget, reforçando que o faturamento dos filmes é altamente variável. O coeficiente de variação acima de 150% mostra uma variabilidade extrema, muito maior que a média. A diferença grande entre média e mediana também indica uma assimetria forte à direita, com poucos filmes faturando bilhões e muitos faturando valores baixos.

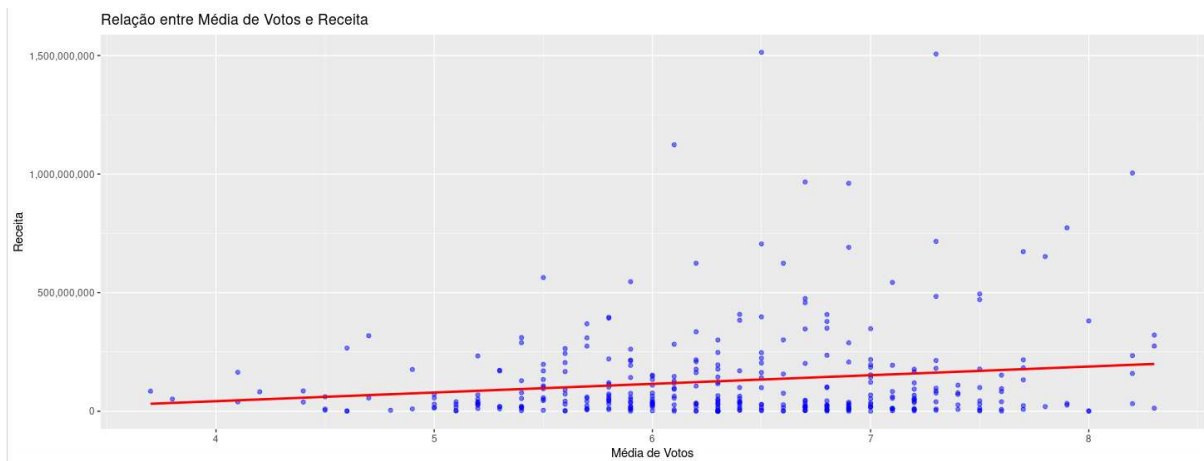
### 3.4 Gráficos Relacionando Variáveis



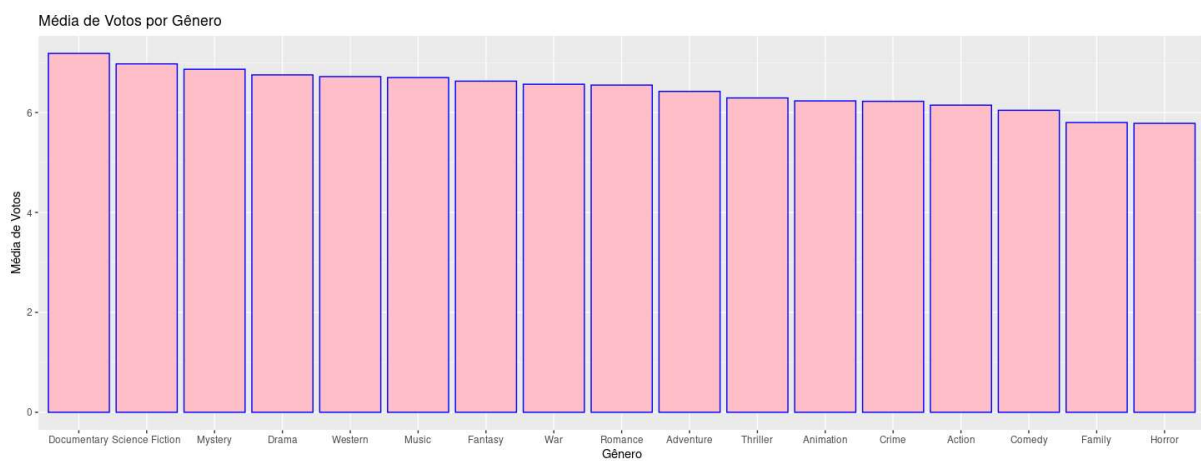
O gráfico mostra que existe uma correlação positiva entre o budget e o revenue: em geral, quanto maior o orçamento, maior a receita. Apesar disso, a dispersão é alta e nem todo filme caro fatura muito, e nem todo filme barato fatura pouco. A linha de regressão é inclinada positivamente, confirmando a tendência. Investir mais tende a aumentar a receita, mas não é garantia de sucesso.



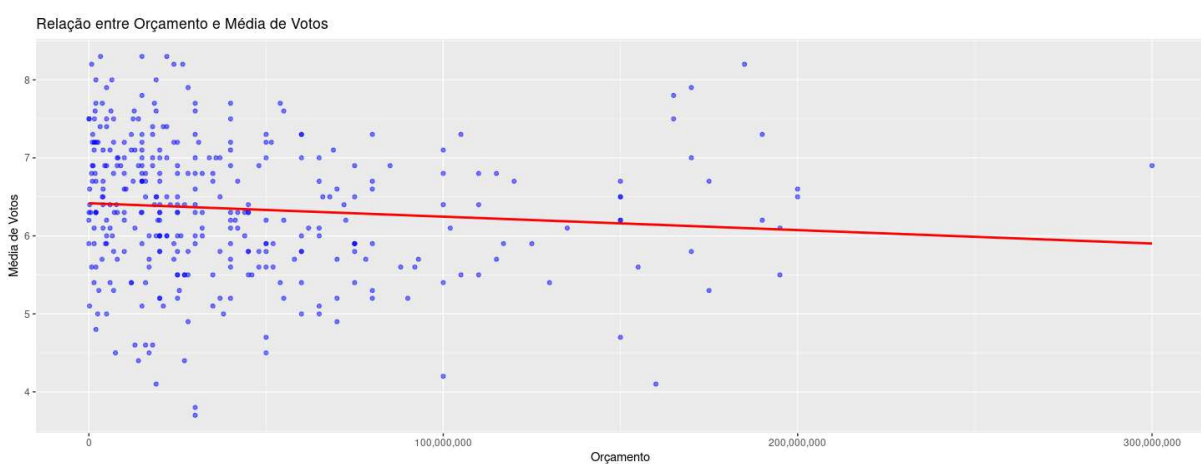
Ao analisar o gráfico que compara a receita média dos filmes por gênero fica evidente que filmes de nichos mais específicos possuem o faturamento médio muito menor quando comparados a filmes cujo público alvo é mais abrangente. Os filmes de animação são os que possuem o maior faturamento enquanto documentários seguem por último com uma diferença exorbitante. Essa disparidade pode ser atribuída a diversos fatores. Filmes de animação costumam atrair públicos de diferentes faixas etárias e contam com investimentos massivos em marketing com uma ampla distribuição internacional. Já os documentários geralmente têm públicos muito mais restritos e são lançados com campanhas de divulgação mais modestas.



O gráfico demonstra uma concentração grande de filmes com receita baixa independentemente da média de votos, entretanto, há uma leve tendência positiva na linha de regressão onde filmes com média de votos mais alta tendem a gerar receitas maiores. Apesar da presença de muitos outliers com receitas muito altas, isso demonstra que melhores avaliações podem estar associadas a maiores receitas, ainda que a correlação também pareça fraca.



O gráfico que compara a média de votos por gênero apresentou resultados bem equilibrados e com pouca variação. O gênero Documentário foi o que conquistou a maior média de votos, enquanto o gênero Horror ficou em último lugar. Isso se deve ao fato dos documentários possuírem um teor mais educacional e não necessariamente precisarem agradar o público quando comparado aos outros gêneros.



O gráfico acima mostra que a maioria dos pontos se concentra na parte esquerda do gráfico onde há orçamentos mais baixos. Apesar de haver filmes com orçamentos altos eles não apresentam uma média de votos consideravelmente maior. A linha de tendência tem uma leve inclinação negativa, indicando uma fraca correlação entre orçamento e a média de votos, entretanto, essa relação é muito fraca e provavelmente não é estatisticamente significativa.

## 4 CONSIDERAÇÕES FINAIS

A análise estatística realizada permitiu obter diversas informações relevantes sobre os filmes presentes na base de dados. Foi possível identificar padrões de produção e distribuição, como a predominância da língua inglesa e dos Estados Unidos como país produtor, além da diversidade de gêneros e produtoras envolvidas.

Através das medidas de tendência central e dispersão, observamos que tanto o orçamento quanto a receita dos filmes apresentam distribuições altamente assimétricas, com grande quantidade de outliers. Isso evidencia uma forte desigualdade no investimento e no retorno financeiro entre os filmes analisados.

Os gráficos relacionando as variáveis quantitativas demonstraram que há uma correlação positiva entre orçamento e receita, ou seja, filmes com maior investimento tendem a gerar maiores receitas, embora existam exceções. Já a relação entre média de votos e receita também mostrou uma tendência positiva, mas fraca, indicando que filmes bem avaliados podem, em média, ter melhores desempenhos financeiros. Por outro lado, não foi possível observar uma correlação significativa entre orçamento e média de votos, sugerindo que um orçamento elevado não garante necessariamente uma boa recepção por parte do público.

Por fim, o uso da linguagem R e da plataforma RStudio mostrou-se eficaz para o tratamento, análise e visualização dos dados, permitindo uma exploração detalhada e acessível de informações complexas. Esta análise evidencia como a estatística pode ser uma ferramenta poderosa para interpretar e compreender melhor os dados no contexto do entretenimento e da indústria cinematográfica.

## REFERÊNCIAS BIBLIOGRÁFICAS

CHUAN, August Sun. **TMDB 5000 Movie Dataset**. 2017. Disponível em:

<[https://www.kaggle.com/datasets/tmdb/tmdb-movie-metadata?select=tmdb\\_5000\\_movies.csv](https://www.kaggle.com/datasets/tmdb/tmdb-movie-metadata?select=tmdb_5000_movies.csv)>.

Acesso em: 27 abr. 2025

R Foundation. **Documentation**. 2025. Disponível em: <<https://www.r-project.org/other-docs.html>>.

Acesso em: 27 abr. 2025.

## ANEXOS

### Sintaxe do R utilizada na análise:

```
#Carrega bibliotecas necessárias
library(dplyr)
library(ggplot2)
library(jsonlite)

#Carrega dataset
movies_dataset = read.csv("tmdb_5000_movies.csv")

#Filtrando informações irrelevantes no dataset e tratando dados
columns_to_ignore = c("homepage", "id", "keywords", "overview",
"spoken_languages", "status", "tagline", "original_title")
movies_dataset_filtred <- movies_dataset[, !(names(movies_dataset) %in%
columns_to_ignore)]
movies_dataset_filtred <-
movies_dataset_filtred[movies_dataset_filtred$budget != 0, ]
movies_dataset_filtred <-
movies_dataset_filtred[movies_dataset_filtred$revenue != 0, ]

set.seed(15112000) # data de nascimento como parametro para criar amostra
movies_dataset_filtred =
movies_dataset_filtred[sample(nrow(movies_dataset_filtred), 350), ]

movies_dataset_filtred <- movies_dataset_filtred |>
  mutate(
    genre = as.character(sapply(lapply(genres, fromJSON), function(x)
x$name[1])),
    production_company = as.character(sapply(lapply(production_companies,
fromJSON), function(x) x$name[1])),
    production_country = as.character(sapply(lapply(production_countries,
fromJSON), function(x) x$name[1])),
  )

movies_dataset_filtred <-
  movies_dataset_filtred[, !(names(movies_dataset_filtred) %in% c("genres",
"production_companies", "production_countries"))]

#Tabelas individuais

#Tabela com os dados usando da língua usada nos filmes
table(movies_dataset_filtred$original_language)

#Tabela com os dados usando do Genêro usada nos filmes
```

```

table(movies_dataset_filtred$genre)

#Tabela com os dados usando do Empresas que produziram os filmes
table(movies_dataset_filtred$production_company)

#Tabela com os dados usando dos países que produziram os filmes
table(movies_dataset_filtred$production_country)

#Tabelas agrupando duas variáveis

#Tabela com os dados relacionando países que produziram os filmes e a língua
do do filme
table(movies_dataset_filtred$production_country,
movies_dataset_filtred$original_language)

#Tabela com os dados gêneros dos filmes produzidos pela França
france_movies =
movies_dataset_filtred[movies_dataset_filtred$production_country == "France",
]
table(france_movies$production_country, france_movies$genre)

#Tabela com os dados gêneros dos filmes produzidos pela Paramount Pictures
paroumunt_movies =
movies_dataset_filtred[movies_dataset_filtred$production_company ==
"Paramount Pictures", ]
table(paroumunt_movies$production_company, paroumunt_movies$genre)


# medidas de resumo
summary(movies_dataset_filtred$budget)
summary(movies_dataset_filtred$revenue)

# Variância budget
variancia_budget <- var(movies_dataset_filtred$budget)
print(paste("Variância: ", variancia_budget))

# Desvio Padrão budget
desvio_padrao_budget <- sd(movies_dataset_filtred$budget)
print(paste("Desvio Padrão: ", desvio_padrao_budget))

# Coeficiente de Variação budget
coeficiente_variacao_budget <- desvio_padrao_budget /
mean(movies_dataset_filtred$budget) * 100

```

```

print(paste("Coeficiente de Variação: ", coeficiente_variacao_budget, "%"))

# Variância revenue
variancia_revenue <- var(movies_dataset_filtred$revenue, na.rm = TRUE)
print(paste("Variância: ", variancia_revenue))

# Desvio Padrão revenue
desvio_padrao_revenue <- sd(movies_dataset_filtred$revenue, na.rm = TRUE)
print(paste("Desvio Padrão: ", desvio_padrao_revenue))

# Coeficiente de Variação
coeficiente_variacao_revenue <- desvio_padrao_revenue /
mean(movies_dataset_filtred$revenue, na.rm = TRUE) * 100
print(paste("Coeficiente de Variação: ", coeficiente_variacao_revenue, "%"))

#Gráfico de densidade para budget
ggplot(movies_dataset_filtred, aes(x = budget, y =
                                after_stat(density))) +
  geom_density(color = "green", linewidth = 1) +
  scale_x_continuous(labels = scales::label_number(big.mark = ",",
                                                    decimal.mark = "."))

#Gráfico de densidade de revenue
ggplot(movies_dataset_filtred, aes(x = revenue, y =
                                after_stat(density))) +
  geom_density(color = "green", linewidth = 1) +
  scale_x_continuous(labels = scales::label_number(big.mark = ",",
                                                    decimal.mark = "."))

#Gráfico de densidade de média das notas do filmes
ggplot(movies_dataset_filtred, aes(x = vote_average, y =
                                after_stat(density))) +
  geom_density(color = "green", linewidth = 1) +
  scale_x_continuous(labels = scales::label_number(big.mark = ",",
                                                    decimal.mark = "."))

# Gráfico Gênero vs Receita
media_receita_genero <- movies_dataset_filtred |>
  group_by(genre) |>
  summarise(media_receita = mean(revenue)) |>
  arrange(desc(media_receita))

ggplot(media_receita_genero, aes(x = reorder(genre, -media_receita), y =
media_receita)) +
  geom_col(color = "blue", fill = "pink") +
  labs(title = "Média de Receita por Gênero", x = "Gênero", y = "Receita
Média") +

```



```

    scale_y_continuous(labels = scales::label_number(big.mark = ",",
decimal.mark = "."))

# Gráfico Média de Votos vs Receita
ggplot(movies_dataset_filtred, aes(x = vote_average, y = revenue)) +
  geom_point(color = "blue", alpha = 0.5) +
  geom_smooth(method = "lm", color = "red", se = FALSE) +
  labs(title = "Relação entre Média de Votos e Receita", x = "Média de
Votos", y = "Receita") +
  scale_y_continuous(labels = scales::label_number(big.mark = ",",
decimal.mark = "."))

# Gráfico Gênero vs Média de Votos
media_votos_genero <- movies_dataset_filtred |>
  group_by(genre) |>
  summarise(media_votos = mean(vote_average)) |>
  arrange(desc(media_votos))

ggplot(media_votos_genero, aes(x = reorder(genre, -media_votos), y =
media_votos)) +
  geom_col(color = "blue", fill = "pink") +
  labs(title = "Média de Votos por Gênero", x = "Gênero", y = "Média de
Votos") +
  scale_y_continuous(labels = scales::label_number(big.mark = ",",
decimal.mark = "."))

# Gráfico Orçamento vs Média de Votos
ggplot(movies_dataset_filtred, aes(x = budget, y = vote_average)) +
  geom_point(color = "blue", alpha = 0.5) +
  geom_smooth(method = "lm", color = "red", se = FALSE) +
  labs(title = "Relação entre Orçamento e Média de Votos", x = "Orçamento", y
= "Média de Votos") +
  scale_x_continuous(labels = scales::label_number(big.mark = ",",
decimal.mark = "."))

# Gráfico com linha de tendência entre budget e revenue
ggplot(movies_dataset_filtred, aes(x = budget, y = revenue)) +
  geom_point(alpha = 0.6, color = "darkorange") +
  geom_smooth(method = "lm", se = TRUE, color = "black") +
  labs(
    title = "Budget vs Revenue com Regressão Linear",
    x = "Budget",
    y = "Revenue"
  ) +
  theme_light()

# Histograma para budget
ggplot(movies_dataset_filtred, aes(x = budget)) +

```

```

    geom_histogram(fill = "lightblue", color = "black", bins = 30) + # Cor das
barras e borda
    scale_x_continuous(labels = label_comma()) + # Formata os números no eixo
X
    scale_y_continuous(labels = label_comma()) + # Formata os números no eixo
Y
    theme_minimal() +
    labs(title = "Histograma de Budget")

# Histograma para revenue
ggplot(movies_dataset_filtred, aes(x = revenue)) +
    geom_histogram(fill = "lightgreen", color = "black", bins = 30) + # Cor
das barras e borda
    scale_x_continuous(labels = label_comma()) + # Formata os números no eixo
X
    scale_y_continuous(labels = label_comma()) + # Formata os números no eixo
Y
    theme_minimal() +
    labs(title = "Histograma de Revenue")

# Boxplot de budget
ggplot(movies_dataset_filtred, aes(y = budget)) +
    geom_boxplot(fill = "lightblue", color = "black") +
    labs(
        title = "Boxplot de Budget",
        y = "Budget"
    ) +
    theme_minimal()

# Boxplot de revenue
ggplot(movies_dataset_filtred, aes(y = revenue)) +
    geom_boxplot(fill = "lightpink", color = "black") +
    labs(
        title = "Boxplot de Revenue",
        y = "Revenue"
    ) +
    theme_minimal()

```