

Final Project Submission

- Student name: Tamjid Ahsan
- Student pace: full time
- Scheduled project review date/time:
- Instructor name: James Irving
- Blog post URL:



▼ 1 Overview

A handful of companies have defined the Hollywood film industry, dominating the US and world markets. They have weathered a world war, and a Great Depression and few moderate ones, innovated wide screen and color technologies, made peace with television, learned to exploit home video and online streaming, and are more powerful than ever before.

Most big corporations are already in this business or exploring feasibility of entry. Most of the major corporations operating only in this industry are thriving.

▼ 2 Business Problem

Microsoft sees all the big companies creating original video content and they want to get in on the fun. They have decided to create a new movie studio, but they don't know anything about creating movies.

I am going to try to figure out what types of films are currently performing better at the box office. I shall recommend some actionable insights based on findings of this analysis, which the head of Microsoft's new movie studio can use to help decide what type of films to create.



Microsoft Movies & TV

MICROSOFT MOVIES

Areas of focus:

- * movie genres.
- * profitability of success based on seasonality of releases.
- * profitability of movie franchise/film series.



3 The imports



3.1 Packages and Libraries

In [1]:

```
1# for web scraping and API calls
2from selenium import webdriver
3from selenium.webdriver.common.keys import Keys
4from selenium.webdriver.support import expected_conditions as EC
5from selenium.webdriver.common.by import By
6from selenium.webdriver.support.wait import WebDriverWait
7import os
8import wget
9import tmdbsimple as tmdb
```

```
In [2]: 1# for other parts
2import os
3import pandas as pd
4import numpy as np
5import matplotlib.pyplot as plt
6%matplotlib inline
7import seaborn as sns
8import json
9import requests
10import time
11from pandas.core.common import flatten
12from pandasql import sqldf
13import plotly.graph_objects as go
14from plotly.subplots import make_subplots
15import re
16import ast
```

```
In [3]: 1# styling (jupyter-themes must be installed)
2## https://github.com/dunovank/jupyter-themes
3from jupyterthemes import jtplot
4# jt -r # default
5jtplot.style(theme='monokai', context='notebook', ticks='True', grid='False')
6
7# jt -t monokai -fs 120 -tfs 120 -nfs 115 -cellw 85% -T -N -kl # my setup
```

```
In [4]: 1# to see dataframe better
2pd.set_option('display.max_columns', 50)
```

3.2 Frequently used fuctions

```
In [5]: 1# Number formatter
2def format_number(data_value, index):
3    if data_value >= 1_000_000_000:
4        formatter = '${:1.1f}B'.format(data_value*0.000_000_001)
5    elif data_value >= 1_000_000:
6        formatter = '${:1.0f}M'.format(data_value*0.000_001)
7    else:
8        formatter = '${:1.0f}K'.format(data_value*0.001)
9    return formatter
```

```
In [6]: 1# % formatter
2def format_add_percentage(data_value, index):
3    formatter = '{:.0f}%'.format(data_value)
4    return formatter
```

```
In [7]: 1 def correlation_top_bottom(df):
2     corr_df_matrix_ = df.unstack().reset_index()
3     corr_df_matrix_.columns = ["feature_0", 'feature_1', 'correlation']
4     corr_df_matrix_['keep'] = corr_df_matrix_.apply(
5         lambda x: False if x['feature_0'] == x['feature_1'] else True, axis=1)
6     corr_df_matrix_['feature_combo'] = corr_df_matrix_.apply(
7         lambda x: ' and '.join(set(x[['feature_0', 'feature_1']])), axis=1)
8     corr_features = corr_df_matrix_[corr_df_matrix_.keep][[
9         'feature_combo', 'correlation'
10    ]].drop_duplicates().sort_values(by='correlation', ascending=False)
11     print(
12         f'Positive correlations:\n\
13         {corr_features.head(10).reset_index()}\n\n {"-"*70}\n\
14         Negative correlations:\n\
15         {corr_features.sort_values(by="correlation").head(10).reset_index()}'
16     )
```

3.3 API and Scraping control

Set this to True to perform scraping and API

```
initialize_scraping_and_API = True
```

```
In [8]: initialize_scraping_and_API = False
```

4 The Data

- IMDb (<https://www.imdb.com/>) or Internet Movie Database was Originally a fan-operated website, now owned and operated by IMDb.com, Inc., a subsidiary of Amazon. This is one of the most reliable source for any information related movies in general. It is one of the most comprehensive dataset.
- Box Office Mojo (<https://www.boxofficemojo.com/>) is also a part of IMDb.com, Inc., providing indepth financial informations among other metrics.
- TMDb (<https://www.themoviedb.org/>) is a reliable source for movie related information. This is a popular user editable database for movies and TV shows.

Those three were used for sourcing data for the project as those are highly reliable sources without going for any paid service for information.

Data is collected from [IMDB website \(https://datasets.imdbws.com\)](https://datasets.imdbws.com) from downloadables, and scraping using [selenium \(https://www.selenium.dev/\)](https://www.selenium.dev/) . Additional data collected from TMDb using [API \(https://www.themoviedb.org/documentation/api\)](https://www.themoviedb.org/documentation/api). Then all of them are merged to create 'main_df', upon which this following analysis is performed.

4.1 From IMDb

4.1.1 Dataset from website

File containing detailed movie info inside [title.basics.tsv.gz \(https://datasets.imdbws.com/title.basics.tsv.gz\)](https://datasets.imdbws.com/title.basics.tsv.gz) was downloaded from <https://datasets.imdbws.com/title.basics.tsv.gz>



4.1.2 Scraping using selenium

```
| pip install selenium
```

Download webdriver from [here \(https://sites.google.com/a/chromium.org/chromedriver/downloads\)](https://sites.google.com/a/chromium.org/chromedriver/downloads).

```

In [9]: 1%%time
        2if initialize_scraping_and_API is True:
        3    # initializing webdriver
        4    driver = webdriver.Chrome('C:/Users/tamji/Documents/PATH/chromedriver.exe')
        5    # connection to webpage
        6    base_url_string = 'https://www.boxofficemojo.com/year/world/'
        7    # selecting years to get
        8    list_of_year = np.arange(2014, 2022, 1)
        9    # initializing scraping
       10    print(f'+ ' * 100)
       11    # temp files
       12    file_names_ = []
       13    file_names_error = []
       14    # scraping
       15    for im in list_of_year:
       16        print(f'Working on: {im}')
       17        url = f'{base_url_string}{im}/'
       18        print(f'Getting {im} homepage')
       19        driver.get(url)
       20        table = driver.find_element_by_xpath('//*[@id="table"]/div/table[2]')
       21        item_href = driver.find_elements_by_class_name('a-link-normal')
       22        print(f'Getting {im} list items')
       23        item_href = [item.get_property('href') for item in item_href]
       24        print(f'Sorting what to keep from {im} list items')
       25        # filter results to target needed links
       26        text_to_check = 'releasegroup'
       27        to_keep = []
       28        to_discard = []
       29        for i in item_href:
       30            if text_to_check in i:
       31                to_keep.append(i)
       32            else:
       33                to_discard.append(i)
       34        print(f'Preping {im} list items for looping')
       35        href = to_keep # [:2] is for testing, remove this to get full data
       36        master_list = []
       37        error = []
       38        print(f'{im} list items are looping. Hang in there!')
       39        for item in href:
       40            try:
       41                driver.get(item)
       42                url = driver.find_element_by_xpath(
       43                    '//*[@id="title-summary-refiner"]/a').get_property('href')
       44                name = driver.find_element_by_xpath(

```

```

45         '//*[@id="a-page"]/main/div/div[1]/div[1]/div/div/div[2]/h1'
46     ).text
47
48     driver.get(url)
49     year = driver.find_element_by_xpath(
50         '//*[@id="a-page"]/main/div/div[1]/div[1]/div/div/div[2]/div/h1/span'
51     ).text
52     worldwide = driver.find_element_by_xpath(
53         '//*[@id="a-page"]/main/div/div[3]/div[1]/div/div[3]/span[2]/span'
54     ).text
55     international = driver.find_element_by_xpath(
56         '//*[@id="a-page"]/main/div/div[3]/div[1]/div/div[2]/span[2]'
57     ).text
58     domestic = driver.find_element_by_xpath(
59         '//*[@id="a-page"]/main/div/div[3]/div[1]/div/div[1]/span[2]'
60     ).text
61
62     year_cleaned = year.strip('(')
63     world_collection = worldwide[1:].replace(", ", "")
64     international_collection = international[1:].replace(", ", "")
65     domestic_collection = domestic[1:].replace(", ", "")
66     imdb_code = url.split('/')[4]
67
68     temp_dict = {
69         'imdb_code': imdb_code,
70         'name': name,
71         'year': year_cleaned,
72         'world_collection': world_collection,
73         'int_collection': international_collection,
74         'dom_collection': domestic_collection,
75         'url': url
76     }
77     master_list.append(temp_dict)
78 except:
79     error.append(item)
80     continue
81
82 df = pd.DataFrame(master_list)
83 file_name_df = f'{im}.csv'
84 df.to_csv(file_name_df, index=False)
85 dict_ = {'urls': error}
86 file_name_error = f'{im}_error.csv'
87 pd.DataFrame(dict_).to_csv(file_name_error, index=False)
88 file_names_.append(file_name_df)
89 file_names_error.append(file_name_error)

```

```

90     print(f'Finished working on {im}\n')
91     print(f'+' * 100)
92     print(f'\n\nnDONE Looping. Cleanig data!!!')
93
94     combined_csv_data = pd.concat([pd.read_csv(f) for f in file_names_])
95     combined_csv_data_error = pd.concat(
96         [pd.read_csv(f) for f in file_names_error])
97
98     combined_csv_data.reset_index(inplace=True)
99     combined_csv_data_error.reset_index(inplace=True)
100
101     combined_csv_data = combined_csv_data.drop(columns='index')
102     combined_csv_data_error = combined_csv_data_error.drop(columns='index')
103
104     combined_csv_data = combined_csv_data.drop_duplicates('imdb_code',
105                                                         ignore_index=True)
106
107     file_name_1 = f'{list_of_year[0]}to{list_of_year[-1]}.csv'
108     file_name_2 = f'{list_of_year[0]}to{list_of_year[-1]}_error.csv'
109     combined_csv_data.to_csv(file_name_1, index=False)
110     combined_csv_data_error.to_csv(file_name_2, index=False)
111
112     print(f'\n\nnDONE!!!')
113     print(f'+' * 100)
114     print(f'+' * 100)
115 # Leaves temp files behind

```

Wall time: 0 ns

```

In [10]: 1# moving major files
          2if initialize_scraping_and_API is True:
          3     destination_1 = f'./Data/bom_{file_name_1}'
          4     destination_2 = f'./Data/temp/{file_name_2}'
          5     os.rename(file_name_1,destination_1)
          6     os.rename(file_name_2,destination_2)

```

```

In [11]: 1def move_files(file):
          2     destination = f'./Data/temp/{file}'
          3     os.rename(file,destination)

```


In [12]:	<pre> 1# moving temp files 2if initialize_scraping_and_API is True: 3 if True: 4 [move_files(f) for f in file_names_] 5 [move_files(f) for f in file_names_error] 6 print('Done moving!!') </pre>	
	Note: <i>repo does not include temp files</i>	
▼	4.2 From TMDb API	
In [13]:	<pre> 1# Load json 2if initialize_scraping_and_API is True: 3 def get_keys(path): 4 with open(path) as f: 5 return json.load(f) </pre>	
In [14]:	<pre> 1# api key initialize 2if initialize_scraping_and_API is True: 3 keys = get_keys("/Users/tamji/.secret/tmdb_api.json") 4 api_key = keys['api_key'] </pre>	
In [15]:	<pre> 1if initialize_scraping_and_API is True: 2 tmdb.API_KEY = api_key </pre>	
In [16]:	<pre> 1# movie_main_df_sliced is cleaned beforehand 2if initialize_scraping_and_API is True: 3 # for matching imdb titles 4 movie_titles_df = pd.read_csv(r'./Data/movie_main_df_sliced.csv', 5 usecols=["tconst"]) </pre>	
In [17]:	<pre> 1# preparing loaded data for use 2if initialize_scraping_and_API is True: 3 imdb_titles = list(flatten(movie_titles_df.values.tolist())) </pre>	
In [18]:	<pre> 1# get how much data is incoming 2if initialize_scraping_and_API is True: 3 len(imdb_titles) </pre>	
In [19]:	<pre> 1# empty df to store results 2if initialize_scraping_and_API is True: 3 df = pd.DataFrame() </pre>	

```
In [20]: 1if initialize_scraping_and_API is True:
        2    for imdb_id in imdb_titles:
        3        try:
        4            movie = tmdb.Movies(imdb_id)
        5            response = movie.info()
        6            df = df.append(pd.json_normalize(movie.info()))
        7        except:
        8            pass
```

```
In [21]: 1if initialize_scraping_and_API is True:
        2    df = df.reset_index()
```

```
In [22]: 1if initialize_scraping_and_API is True:
        2    df = df.drop(columns=['index'])
```

```
In [23]: 1if initialize_scraping_and_API is True:
        2    df.to_csv(r'./Data/tmdb_parsd.csv')
```

▼ 5 Preparing datasets

▼ 5.1 IMDb

▼ 5.1.1 loading

```
In [24]: 1%%time
        2df_1 = pd.read_csv(r'./Data/data.tsv',
        3                    delimiter='\t',
        4                    low_memory=False)

Wall time: 17.7 s
```

▼ 5.1.2 inspecting

```
In [25]: 1df_1.head(3)
```

	tconst	titleType	primaryTitle	originalTitle	isAdult	startYear	endYear	runtimeMinutes	genres
0	tt0000001	short	Carmencita	Carmencita	0	1894	\N	1	Documentary,Short
1	tt0000002	short	Le clown et ses chiens	Le clown et ses chiens	0	1892	\N	5	Animation,Short
2	tt0000003	short	Pauvre Pierrot	Pauvre Pierrot	0	1892	\N	4	Animation,Comedy,Romance

In [26]:

```
1df_1['titleType'].value_counts()
```

```
tvEpisode      5590798
short           799028
movie           570678
video           297824
tvSeries        203184
tvMovie         130415
tvMiniSeries    36270
tvSpecial        31753
videoGame       27529
tvShort          9611
episode          1
audiobook        1
radioSeries      1
Name: titleType, dtype: int64
```



5.1.3 cleaning

In [27]:

```
1%%time
2# slicing to keep only movies
3movie_df = df_1[df_1['titleType'] == 'movie']
4# dropping adult titles
5movie_df = movie_df[movie_df['isAdult'] == '0']
6# handling nan values
7movie_df.loc[movie_df['runtimeMinutes'] == r'\N', 'runtimeMinutes'] = np.nan
8movie_df.loc[movie_df['startYear'] == r'\N', 'startYear'] = np.nan
9movie_df.loc[movie_df['genres'] == r'\N', 'genres'] = np.nan
10# setting nan genere to NoInfo
11movie_df.loc[movie_df['genres'].isna(), 'genres'] = "NoInfo"
12# nan value dropping for start year
13movie_df = movie_df[~movie_df['startYear'].isna()]
14
15movie_df = movie_df.reset_index()
16movie_df = movie_df.drop(['index', 'titleType', 'endYear', 'isAdult'], axis=1)
17
18movie_df.to_csv(r'./Data/movie_df.csv', index=False)
19movie_df
```

Wall time: 2.38 s

	tconst	primaryTitle	originalTitle	startYear	runtimeMinutes	genres
0	tt0000502	Bohemios	Bohemios	1905	100	NoInfo
1	tt0000574	The Story of the Kelly Gang	The Story of the Kelly Gang	1906	70	Action,Adventure,Biography
2	tt0000615	Robbery Under Arms	Robbery Under Arms	1907	NaN	Drama
3	tt0000630	Hamlet	Amleto	1908	NaN	Drama
4	tt0000675	Don Quijote	Don Quijote	1908	NaN	Drama
...
490999	tt9916622	Rodolpho Teóphilo - O Legado de um Pioneiro	Rodolpho Teóphilo - O Legado de um Pioneiro	2015	57	Documentary
491000	tt9916680	De la ilusión al desconcierto: cine colombiano...	De la ilusión al desconcierto: cine colombiano...	2007	100	Documentary
491001	tt9916706	Dankyavar Danka	Dankyavar Danka	2013	NaN	Comedy
491002	tt9916730	6 Gunn	6 Gunn	2017	116	NoInfo
491003	tt9916754	Chico Albuquerque - Revelações	Chico Albuquerque - Revelações	2013	49	Documentary

491004 rows x 6 columns

splitting genere

In [28]:

```
1%%time
2# getting preliminary unique list for cleaning
3genres = list(movie_df['genres'].unique())
4# temp list to store list of splited genre
5genre_cleaning_temp = []
6# getting list of splited genre
7for item in genres:
8    # for dealing with nan
9    if type(item) is not float:
10        # actual splitting
11        genre_split = item.split(",")
12        # appending
13        genre_cleaning_temp.extend(genre_split)
14# geting unique list
15from pandas.core.common import flatten
16# flattening temp list
17## https://stackoverflow.com/questions/12897374/get-unique-values-from-a-list-in-python by https://stackoverflow.com/users/2062318/todoi
18## https://saralgyaan.com/posts/nested-list-to-list-python-in-just-three-lines-of-code/ ##
19genre_cleaning_temp = list(flatten(genre_cleaning_temp))
20# unique genre list
21unique_genre = list(dict.fromkeys(genre_cleaning_temp))
22
23## overly complicated way, theres much simpler method out in the wild.
24unique_genre
```

Wall time: 28 ms

```
['NoInfo',
 'Action',
 'Adventure',
 'Biography',
 'Drama',
 'Fantasy',
 'Comedy',
 'War',
 'Documentary',
 'Crime',
 'Romance',
 'Family',
 'History',
 'Sci-Fi',
 'Thriller',
 'Western',
 'Short',
 'Sport',
 'Mystery',
 'Horror',
 'Music',
 'Animation',
 'Musical',
 'Film-Noir',
 'News',
 'Adult',
 'Reality-TV',
```

```
'Game-Show',  
'Talk-Show']
```

```
In [29]: 1%%time  
2#boolean matrix for all genere  
3movie_genre_df = pd.DataFrame([(x in y) for x in unique_genre]  
4                               for y in movie_df['genres']],  
5                               columns=unique_genre)
```

Wall time: 2.42 s

```
In [30]: 1# merging  
2movie_main_df = pd.concat([movie_df, movie_genre_df], axis=1)
```

```
In [31]: 1# enforcing dtypes  
2movie_main_df = movie_main_df.convert_dtypes()
```

```
In [32]: 1movie_main_df.shape
```

(491004, 35)

```
In [33]: 1movie_main_df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 491004 entries, 0 to 491003
Data columns (total 35 columns):
#   Column                Non-Null Count  Dtype
---  ---                ---
0   tconst                491004 non-null  string
1   primaryTitle          491004 non-null  string
2   originalTitle         491004 non-null  string
3   startYear             491004 non-null  string
4   runtimeMinutes        348729 non-null  string
5   genres                491004 non-null  string
6   NoInfo                491004 non-null  boolean
7   Action                491004 non-null  boolean
8   Adventure             491004 non-null  boolean
9   Biography             491004 non-null  boolean
10  Drama                 491004 non-null  boolean
11  Fantasy               491004 non-null  boolean
12  Comedy               491004 non-null  boolean
13  War                   491004 non-null  boolean
14  Documentary           491004 non-null  boolean
15  Crime                 491004 non-null  boolean
16  Romance               491004 non-null  boolean
17  Family                491004 non-null  boolean
18  History               491004 non-null  boolean
19  Sci-Fi                491004 non-null  boolean
20  Thriller              491004 non-null  boolean
21  Western               491004 non-null  boolean
22  Short                 491004 non-null  boolean
23  Sport                 491004 non-null  boolean
24  Mystery               491004 non-null  boolean
25  Horror                491004 non-null  boolean
26  Music                 491004 non-null  boolean
27  Animation             491004 non-null  boolean
28  Musical               491004 non-null  boolean
29  Film-Noir             491004 non-null  boolean
30  News                  491004 non-null  boolean
31  Adult                 491004 non-null  boolean
32  Reality-TV            491004 non-null  boolean
33  Game-Show             491004 non-null  boolean
34  Talk-Show             491004 non-null  boolean
dtypes: boolean(29), string(6)
memory usage: 49.6 MB
```

```
In [34]: 1movie_main_df.describe()
```

	tconst	primaryTitle	originalTitle	startYear	runtimeMinutes	genres	NoInfo	Action	Adventure	Biography	Drama	Fantasy	Comedy	War	Documentary
count	491004	491004	491004	491004	348729	491004	491004	491004	491004	491004	491004	491004	491004	491004	491004
unique	491004	435498	444525	133	470	1317	2	2	2	2	2	2	2	2	2
top	tt0010975	Mother	Home	2017	90	Drama	False	False	False	False	False	False	False	False	False
freq	1	40	36	17755	23507	90267	424787	450544	468879	478159	309787	480389	404506	483095	393223

```
In [35]: 1movie_main_df['startYear'] = movie_main_df['startYear'].astype('int')
2movie_main_df['runtimeMinutes'].fillna('0', inplace=True)
3movie_main_df['runtimeMinutes'] = movie_main_df['runtimeMinutes'].astype('int')
```

```
In [36]: 1movie_main_df['startYear'].sort_values().unique()

array([1896, 1897, 1898, 1899, 1900, 1901, 1902, 1903, 1904, 1905, 1906,
       1907, 1908, 1909, 1910, 1911, 1912, 1913, 1914, 1915, 1916, 1917,
       1918, 1919, 1920, 1921, 1922, 1923, 1924, 1925, 1926, 1927, 1928,
       1929, 1930, 1931, 1932, 1933, 1934, 1935, 1936, 1937, 1938, 1939,
       1940, 1941, 1942, 1943, 1944, 1945, 1946, 1947, 1948, 1949, 1950,
       1951, 1952, 1953, 1954, 1955, 1956, 1957, 1958, 1959, 1960, 1961,
       1962, 1963, 1964, 1965, 1966, 1967, 1968, 1969, 1970, 1971, 1972,
       1973, 1974, 1975, 1976, 1977, 1978, 1979, 1980, 1981, 1982, 1983,
       1984, 1985, 1986, 1987, 1988, 1989, 1990, 1991, 1992, 1993, 1994,
       1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005,
       2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016,
       2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027,
       2028])
```



5.1.4 choosing features

Choosing to focus analysis on movies released between 2015 to 2020, where primary spoken language is English.

- In my opinion this is the most appropriate time frame to focus, as this gives enough data for analysis and at the same time does not include old info which will not be good representative of the current market situation. As customer/viewer taste and market trends shift over the time.
- Microsoft should focus only on releasing content in **English** for their kick-off. This gives them enough exposure and get noticed as a big player in the game, as they intend to be. Although they should focus on other territory to explore as there are ample opportunities left untapped. For example, in 2020 China surpassed North America in terms of industry value.
As Microsoft has business across the globe, this should be relatively straight forward for them.
- I am also choosing not to focus on **ultra-low** budget movies for this analysis. Microsoft is one of the biggest corporations on earth. They have financial support to go for the big studios.
- I am also not including **'Documentary', 'Short', 'Adult', 'Reality-TV', 'Game-Show', 'Talk-Show', 'News', 'Film-Noir'** titles. Those are entirely different class of product to be compared with conventional movies.

```
In [37]: 1# filtering based on year, keeping one additional year just to be safe
        2movie_main_df_sliced = movie_main_df[(movie_main_df['startYear'] >= 2014)
        3                                         & (movie_main_df['startYear'] <= 2021)]
```

```
In [38]: 1movie_main_df_sliced.describe()
```

	startYear	runtimeMinutes
count	123293.000000	123293.000000
mean	2017.221789	68.234523
std	2.117191	100.801411
min	2014.000000	0.000000
25%	2015.000000	45.000000
50%	2017.000000	80.000000
75%	2019.000000	96.000000
max	2021.000000	28643.000000


```
In [39]: 1to_drop = [  
2     'Documentary', 'Short', 'Adult', 'Reality-TV', 'Game-Show', 'Talk-Show',  
3     'News', 'Film-Noir'  
4]
```

```
In [40]: 1for item in to_drop:  
2     movie_main_df_sliced = movie_main_df_sliced[~movie_main_df_sliced[item].  
3                                     eq(1)]
```

```
In [41]: 1movie_main_df_sliced
```

	tconst	primaryTitle	originalTitle	startYear	runtimeMinutes	genres	NoInfo	Action	Adventure	Biography	Drama	Fantasy	Comedy	War
5089	tt0011216	Spanish Fiesta	La fête espagnole	2019	67	Drama	False	False	False	False	True	False	False	False
5560	tt0011801	Tötet nicht mehr	Tötet nicht mehr	2019	0	Action,Crime	False	True	False	False	False	False	False	False
9809	tt0016906	Frivolinas	Frivolinas	2014	80	Comedy,Musical	False	False	False	False	False	False	True	False
45545	tt0062336	El Tango del Viudo y Su Espejo Deformante	El Tango del Viudo y Su Espejo Deformante	2020	70	Drama	False	False	False	False	True	False	False	False
50362	tt0069049	The Other Side of the Wind	The Other Side of the Wind	2018	122	Drama	False	False	False	False	True	False	False	False
...
490995	tt9916270	Il talento del calabrone	Il talento del calabrone	2020	84	Thriller	False	False	False	False	False	False	False	False
490996	tt9916362	Coven	Akelarre	2020	90	Adventure,Drama,History	False	False	True	False	True	False	False	False
490997	tt9916428	The Secret of China	Hong xing zhao yao Zhong guo	2019	0	Adventure,History,War	False	False	True	False	False	False	False	True
490998	tt9916538	Kuambil Lagi Hatiku	Kuambil Lagi Hatiku	2019	123	Drama	False	False	False	False	True	False	False	False
491002	tt9916730	6 Gunn	6 Gunn	2017	116	NoInfo	True	False	False	False	False	False	False	False

84546 rows x 35 columns

```
In [42]: 1movie_main_df_sliced.to_csv('./Data/movie_main_df_sliced.csv',index = False)
```

5.2 Merging all sources

```
In [43]: 1# Loading datasets  
2imdb_df = pd.read_csv('./Data/movie_main_df_sliced.csv')  
3bom_df = pd.read_csv('./Data/bom_2014to2021.csv')  
4tmdb_df = pd.read_csv('./Data/tmdb_parsd.csv')
```

merge_1

In [44]:

```
1merge_1 = pd.merge(imdb_df,
2                    bom_df,
3                    how='left',
4                    left_on='tconst',
5                    right_on='imdb_code')
```

In [45]:

```
1merge_1
```

	tconst	primaryTitle	originalTitle	startYear	runtimeMinutes	genres	NoInfo	Action	Adventure	Biography	Drama	Fantasy	Comedy	War	D
0	tt0011216	Spanish Fiesta	La fête espagnole	2019	67	Drama	False	False	False	False	True	False	False	False	F
1	tt0011801	Tötet nicht mehr	Tötet nicht mehr	2019	0	Action,Crime	False	True	False	False	False	False	False	False	F
2	tt0016906	Frivolinas	Frivolinas	2014	80	Comedy,Musical	False	False	False	False	False	False	True	False	F
3	tt0062336	El Tango del Viudo y Su Espejo Deformante	El Tango del Viudo y Su Espejo Deformante	2020	70	Drama	False	False	False	False	True	False	False	False	F
4	tt0069049	The Other Side of the Wind	The Other Side of the Wind	2018	122	Drama	False	False	False	False	True	False	False	False	F
...
84541	tt9916270	Il talento del calabrone	Il talento del calabrone	2020	84	Thriller	False	False	False	False	False	False	False	False	F
84542	tt9916362	Coven	Akelarre	2020	90	Adventure,Drama,History	False	False	True	False	True	False	False	False	F
84543	tt9916428	The Secret of China	Hong xing zhao yao Zhong guo	2019	0	Adventure,History,War	False	False	True	False	False	False	False	True	F
84544	tt9916538	Kuambil Lagi Hatiku	Kuambil Lagi Hatiku	2019	123	Drama	False	False	False	False	True	False	False	False	F
84545	tt9916730	6 Gunn	6 Gunn	2017	116	NoInfo	True	False	False	False	False	False	False	False	F

84546 rows x 42 columns

prepping for merge 2

In [46]:

```
1tmdb_df = tmdb_df.drop(tmdb_df.columns[0:4], axis=1)
```

In [47]:

```
1tmdb_df.columns
```

Index(['budget', 'genres', 'homepage', 'id', 'imdb_id', 'original_language', 'original_title', 'overview', 'popularity', 'poster_path', 'production_companies', 'production_countries', 'release_date', 'revenue', 'runtime', 'spoken_languages', 'status', 'tagline', 'title', 'video', 'vote_average', 'vote_count', 'belongs_to_collection.id', 'belongs_to_collection.name', 'belongs_to_collection.poster_path', 'belongs_to_collection.backdrop_path'], dtype='object')

In [48]:	<pre> 1filter_list = [2 'imdb_id', 'title', 'revenue', 'budget', 'release_date', 3 'production_companies', 'popularity', 'vote_average', 'vote_count', 4 'overview', 'belongs_to_collection.name', 'original_language' 5] </pre>	
In [49]:	<pre> 1tmdb_df_reduced = tmdb_df[filter_list] </pre>	
	merge 2	
In [50]:	<pre> 1merge_2 = pd.merge(merge_1, 2 tmdb_df_reduced, 3 how='inner', 4 left_on='tconst', 5 right_on='imdb_id') </pre>	
In [51]:	<pre> 1df = merge_2.copy() </pre>	
In [52]:	<pre> 1df.columns Index(['tconst', 'primaryTitle', 'originalTitle', 'startYear', 'runtimeMinutes', 'genres', 'NoInfo', 'Action', 'Adventure', 'Biography', 'Drama', 'Fantasy', 'Comedy', 'War', 'Documentary', 'Crime', 'Romance', 'Family', 'History', 'Sci-Fi', 'Thriller', 'Western', 'Short', 'Sport', 'Mystery', 'Horror', 'Music', 'Animation', 'Musical', 'Film-Noir', 'News', 'Adult', 'Reality-TV', 'Game-Show', 'Talk-Show', 'imdb_code', 'name', 'year', 'world_collection', 'int_collection', 'dom_collection', 'url', 'imdb_id', 'title', 'revenue', 'budget', 'release_date', 'production_companies', 'popularity', 'vote_average', 'vote_count', 'overview', 'belongs_to_collection.name', 'original_language'], dtype='object') </pre>	
	cleaning	
In [53]:	<pre> 1rearrange = [2 'tconst', 'imdb_code', 'imdb_id', 'primaryTitle', 'originalTitle', 'name', 3 'title', 'startYear', 'year', 'release_date', 'runtimeMinutes', 'budget', 4 'revenue', 'world_collection', 'int_collection', 'dom_collection', 5 'production_companies', 'popularity', 'vote_average', 'vote_count', 6 'overview', 'belongs_to_collection.name', 'original_language', 'genres', 'NoInfo', 'Action', 7 'Adventure', 'Biography', 'Drama', 'Fantasy', 'Comedy', 'War', 'Crime', 8 'Romance', 'Family', 'History', 'Sci-Fi', 'Thriller', 'Western', 'Sport', 9 'Mystery', 'Horror', 'Music', 'Animation', 'Musical', 'url' 10] 11df = df[rearrange] </pre>	
	filtering order:	

1. financial data
2. year
3. review data

In [54]:

```
1df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 44353 entries, 0 to 44352
Data columns (total 46 columns):
#   Column                Non-Null Count  Dtype
---  -
0   tconst                 44353 non-null object
1   imdb_code              13351 non-null object
2   imdb_id               44353 non-null object
3   primaryTitle           44353 non-null object
4   originalTitle          44353 non-null object
5   name                  13351 non-null object
6   title                 44353 non-null object
7   startYear             44353 non-null int64
8   year                  13351 non-null float64
9   release_date          42437 non-null object
10  runtimeMinutes         44353 non-null int64
11  budget                 44353 non-null int64
12  revenue                44353 non-null int64
13  world_collection       13351 non-null float64
14  int_collection         12646 non-null float64
15  dom_collection         3289 non-null float64
16  production_companies   44353 non-null object
17  popularity             44353 non-null float64
18  vote_average           44353 non-null float64
19  vote_count             44353 non-null int64
20  overview               41041 non-null object
21  belongs_to_collection.name 1911 non-null object
22  original_language      44353 non-null object
23  genres                 44353 non-null object
24  NoInfo                 44353 non-null bool
25  Action                 44353 non-null bool
26  Adventure               44353 non-null bool
27  Biography              44353 non-null bool
28  Drama                  44353 non-null bool
29  Fantasy                44353 non-null bool
30  Comedy                 44353 non-null bool
31  War                    44353 non-null bool
32  Crime                  44353 non-null bool
33  Romance                44353 non-null bool
34  Family                 44353 non-null bool
35  History                44353 non-null bool
36  Sci-Fi                 44353 non-null bool
37  Thriller               44353 non-null bool
38  Western                44353 non-null bool
39  Sport                  44353 non-null bool
40  Mystery                44353 non-null bool
41  Horror                 44353 non-null bool
42  Music                  44353 non-null bool
43  Animation              44353 non-null bool
44  Musical                44353 non-null bool
45  url                    13351 non-null object
dtypes: bool(21), float64(6), int64(5), object(14)
memory usage: 9.7+ MB
```

```
In [55]: 1df.describe()
```

	startYear	year	runtimeMinutes	budget	revenue	world_collection	int_collection	dom_collection	popularity	vote_average	vote_count
count	44353.000000	13351.000000	44353.000000	4.435300e+04	4.435300e+04	1.335100e+04	1.264600e+04	3.289000e+03	44353.000000	44353.000000	44353.000000
mean	2017.107704	2016.916411	88.849976	1.462749e+06	4.188408e+06	1.722414e+07	1.279721e+07	2.060086e+07	5.966112	3.952096	101.957252
std	2.025922	1.846724	36.046684	1.175852e+07	4.898516e+07	9.106972e+07	6.373112e+07	6.311116e+07	40.339807	3.108509	742.014378
min	2014.000000	2014.000000	0.000000	0.000000e+00	0.000000e+00	2.000000e+00	2.000000e+00	4.900000e+01	0.000000	0.000000	0.000000
25%	2015.000000	2015.000000	81.000000	0.000000e+00	0.000000e+00	3.873450e+04	4.289450e+04	3.447100e+04	0.600000	0.000000	0.000000
50%	2017.000000	2017.000000	92.000000	0.000000e+00	0.000000e+00	3.877660e+05	4.062780e+05	3.423700e+05	1.513000	5.000000	2.000000
75%	2019.000000	2018.000000	106.000000	0.000000e+00	0.000000e+00	3.320476e+06	3.206889e+06	8.106986e+06	5.705000	6.400000	10.000000
max	2021.000000	2021.000000	1260.000000	3.560000e+08	2.797801e+09	2.797501e+09	1.939128e+09	9.366622e+08	5227.005000	10.000000	25252.000000

```
In [56]: 1df['revenue'].sort_values().value_counts() # null values are stored as 0
```

```
0      41364
10000      27
100000     21
1500000     17
500         9
...
147315      1
15894372     1
42972994     1
117813057     1
158162788     1
Name: revenue, Length: 2721, dtype: int64
```

Choosing greater value among two data sources for revenue, then cleaning noises.

```
In [57]: 1df['world_collection'].isna().value_counts()
```

```
True      31002
False     13351
Name: world_collection, dtype: int64
```

```
In [58]: 1((df['revenue']!=0)&(df['world_collection'].isna())).value_counts()
```

```
False     44003
True        350
dtype: int64
```

```
In [59]: 1condition_1 = (df['revenue']!=0)
```

```
In [60]: 1condition_2 = ~df['world_collection'].isna()
```

```
In [61]: 1df = df[condition_1 | condition_2]
```

In [62]:

1df

	tconst	imdb_code	imdb_id	primaryTitle	originalTitle	name	title	startYear	year	release_date	runtimeMinutes	budget	revenue	world_collection
5	tt0100275	tt0100275	tt0100275	The Wandering Soap Opera	La Telenovela Errante	The Wandering Soap Opera	The Wandering Soap Opera	2017	2017.0	2017-08-10	80	0	0	3.624000e+03
22	tt0315642	tt0315642	tt0315642	Wazir	Wazir	Wazir	Wazir	2016	2016.0	2016-01-07	103	5200000	9200000	5.633588e+06
26	tt0331314	tt0331314	tt0331314	Bunyan and Babe	Bunyan and Babe	Bunyan and Babe	Bunyan and Babe	2017	2017.0	2017-01-12	84	0	0	7.206000e+04
32	tt0365907	tt0365907	tt0365907	A Walk Among the Tombstones	A Walk Among the Tombstones	A Walk Among the Tombstones	A Walk Among the Tombstones	2014	2014.0	2014-09-18	114	28000000	53181600	5.883438e+07
33	tt0369610	tt0369610	tt0369610	Jurassic World	Jurassic World	Jurassic World	Jurassic World	2015	2015.0	2015-06-06	124	150000000	1671713208	1.670516e+09
...
44331	tt9908390	tt9908390	tt9908390	Le lion	Le lion	Le lion	The Lion	2020	2020.0	2020-01-29	95	0	0	3.507711e+06
44333	tt9908960	tt9908960	tt9908960	Pliusas	Pliusas	Pliusas	Pliusas	2018	2018.0	2018-09-07	90	0	0	7.463700e+04
44339	tt9911196	tt9911196	tt9911196	The Marriage Escape	De beentjes van Sint-Hildegard	The Marriage Escape	The Marriage Escape	2020	2020.0	2020-02-10	103	0	0	7.760946e+06
44347	tt9914942	tt9914942	tt9914942	La vida sense la Sara Amat	La vida sense la Sara Amat	La vida sense la Sara Amat	La vida sense la Sara Amat	2019	2019.0	2019-07-12	74	0	0	5.979400e+04
44352	tt9916428	tt9916428	tt9916428	The Secret of China	Hong xing zhao yao Zhong guo	The Secret of China	The Secret of China	2019	2019.0	2019-08-08	0	0	0	4.408165e+06

13701 rows x 46 columns

In [63]:

1# selecting max value as budget

2df.loc[:,['world_collection']] = df[['revenue','world_collection']].max(axis=1)

In [64]:	<div> <div>▼</div> <pre> 1# redundant data dropping 2drop_list = [3 'tconst', 'imdb_code', 'index','name', 'title', 'year', 'revenue', 'url' 4] </pre> </div>	
In [65]:	<pre>1df = df.reset_index()</pre>	
In [66]:	<pre>1df = df.drop(columns=drop_list)</pre>	
In [67]:	<pre>1df["release_date"] = pd.to_datetime(df["release_date"])</pre>	
	dealing with nested data	
In [68]:	<div> <div>▼</div> <pre> 1# creating a copy of df 2df1 = df.copy() </pre> </div>	
In [69]:	<div> <div>▼</div> <pre> 1# getting a slice to work on 2df1 = df1[['imdb_id', 'production_companies']] </pre> </div>	
In [70]:	<pre>1df1_dict=df1.to_dict()</pre>	
In [71]:	<pre>1df1_dict.keys()</pre> <pre>dict_keys(['imdb_id', 'production_companies'])</pre>	
In [72]:	<div> <div>▼</div> <pre>1# https://stackoverflow.com/questions/39807724/extract-python-dictionary-from-string by https://stackoverflow.com/users/3734244/danidee</pre> </div>	
In [73]:	<div> <div>▼</div> <pre> 1def get_list(string): 2 x = ast.literal_eval(re.search('{.+}', string).group(0)) 3 return x </pre> </div>	
In [74]:	<pre> 1temp = [] #store temp dicts 2ty = [] #catch errors 3for item in df1_dict['production_companies']: 4 x = df1_dict['production_companies'][item] 5 try: 6 temp.append(get_list(x)) 7 except: 8 temp.append(ty) </pre>	

In [75]:	<pre>1#lopping through temp dicts and extracting production house name 2temp_li = [] 3 4for i in temp: 5 if type(i) == tuple: 6 lli = [] 7 for y in i: 8 lli.append(y['name']) 9 code = ', '.join(lli) 10 temp_dict = { 11 'production_comp': code, 12 } 13 temp_li.append(temp_dict) 14 elif type(i) == dict: 15 16 code = i['name'] 17 temp_dict = { 18 'production_comp': code, 19 } 20 temp_li.append(temp_dict) 21 22 elif type(i) == list: 23 24 code = 'Others,No info' 25 temp_dict = { 26 'production_comp': code, 27 } 28 temp_li.append(temp_dict)</pre>	
In [76]:	<pre>1pro = pd.DataFrame.from_dict(temp_li)</pre>	
In [77]:	<pre>1pro_1=pd.concat([df1.reset_index(),pro],axis=1)</pre>	
In [78]:	<pre>1pro_1=pro_1.drop(axis=1, columns=['index','production_companies'])</pre>	
In [79]:	<pre>1df_final = pd.merge(df, pro_1, left_on='imdb_id', right_on='imdb_id')</pre>	
	<ul style="list-style-type: none">• touchup	

In [80]:

```
1df_final.head(4)
```

	imdb_id	primaryTitle	originalTitle	startYear	release_date	runtimeMinutes	budget	world_collection	int_collection	dom_collection	production_companies	popularity
0	tt0100275	The Wandering Soap Opera	La Telenovela Errante	2017	2017-08-10	80	0	3624.0	NaN	3624.0	[[{'id': 96241, 'logo_path': None, 'name': 'Poe...	1.400
1	tt0315642	Wazir	Wazir	2016	2016-01-07	103	5200000	9200000.0	4509543.0	1124045.0	[[{'id': 12865, 'logo_path': None, 'name': 'Get...	5.191
2	tt0331314	Bunyan and Babe	Bunyan and Babe	2017	2017-01-12	84	0	72060.0	72060.0	NaN	[[{'id': 87468, 'logo_path': None, 'name': 'Too...	20.049
3	tt0365907	A Walk Among the Tombstones	A Walk Among the Tombstones	2014	2014-09-18	114	28000000	58834384.0	32526784.0	26307600.0	[[{'id': 39043, 'logo_path': None, 'name': 'Tra...	34.302

In [81]:

```
1df_final.columns
```

```
Index(['imdb_id', 'primaryTitle', 'originalTitle', 'startYear', 'release_date',  
      'runtimeMinutes', 'budget', 'world_collection', 'int_collection',  
      'dom_collection', 'production_companies', 'popularity', 'vote_average',  
      'vote_count', 'overview', 'belongs_to_collection.name',  
      'original_language', 'genres', 'NoInfo', 'Action', 'Adventure',  
      'Biography', 'Drama', 'Fantasy', 'Comedy', 'War', 'Crime', 'Romance',  
      'Family', 'History', 'Sci-Fi', 'Thriller', 'Western', 'Sport',  
      'Mystery', 'Horror', 'Music', 'Animation', 'Musical',  
      'production_comp'],  
      dtype='object')
```

In [82]:

```
1df_final=df_final.drop(columns='production_companies')
```

```
In [83]: 1 rearrange_ = [
2         'imdb_id', 'primaryTitle', 'originalTitle', 'startYear', 'release_date',
3         'runtimeMinutes', 'budget', 'world_collection', 'int_collection',
4         'dom_collection', 'popularity', 'vote_average',
5         'vote_count', 'production_comp', 'original_language', 'belongs_to_collection.name',
6         'genres', 'NoInfo', 'Action', 'Adventure',
7         'Biography', 'Drama', 'Fantasy', 'Comedy', 'War', 'Crime', 'Romance',
8         'Family', 'History', 'Sci-Fi', 'Thriller', 'Western', 'Sport',
9         'Mystery', 'Horror', 'Music', 'Animation', 'Musical',
10        'overview'
11]
```

```
In [84]: 1 df_final = df_final[rearrange_]
```

```
In [85]: 1 df_final
```

	imdb_id	primaryTitle	originalTitle	startYear	release_date	runtimeMinutes	budget	world_collection	int_collection	dom_collection	popularity	vote_average	v
0	tt0100275	The Wandering Soap Opera	La Telenovela Errante	2017	2017-08-10	80	0	3.624000e+03	NaN	3624.0	1.400	6.5	9
1	tt0315642	Wazir	Wazir	2016	2016-01-07	103	5200000	9.200000e+06	4.509543e+06	1124045.0	5.191	6.6	9
2	tt0331314	Bunyan and Babe	Bunyan and Babe	2017	2017-01-12	84	0	7.206000e+04	7.206000e+04	NaN	20.049	6.2	1

```
In [86]: 1 df_final.to_csv('./Data/main_df.csv', index=False)
```

5.3 Working on main_df

5.3.1 prepping for analysis, furthur cleaning

```
In [87]: 1 main_df_raw = pd.read_csv(r'./Data/main_df.csv',
2                               parse_dates=['release_date'],
3                               low_memory=False)
```

```
In [88]: 1 main_df=main_df_raw.iloc[:,0:17] #droping boolean columns
```

```
In [89]: 1main_df=main_df[~main_df.release_date.isna()]
```

```
In [90]: 1main_df['release_year'] = main_df['release_date'].dt.year
2main_df['release_year'].astype('int')
```

```
0      2017
1      2016
2      2017
3      2014
4      2015
...
13696   2020
13697   2018
13698   2020
13699   2019
13700   2019
Name: release_year, Length: 13620, dtype: int32
```

Focusing my analysis from 2015 to end of 2020. Inputs below can be changed to focus any timeframe from 2007 to March 12, 2021. Data is in safe folder inside repo.

```
In [91]: 1main_df = main_df[(main_df.release_date >= '2015-01-01')
2          & (main_df.release_date <= '2020-12-31')]
```

```
In [92]: 1main_df.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 11779 entries, 0 to 13700
Data columns (total 18 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   imdb_id                11779 non-null  object
1   primaryTitle            11779 non-null  object
2   originalTitle           11779 non-null  object
3   startYear              11779 non-null  int64
4   release_date           11779 non-null  datetime64[ns]
5   runtimeMinutes         11779 non-null  int64
6   budget                 11779 non-null  int64
7   world_collection        11779 non-null  float64
8   int_collection          10924 non-null  float64
9   dom_collection          2765 non-null  float64
10  popularity              11779 non-null  float64
11  vote_average            11779 non-null  float64
12  vote_count              11779 non-null  int64
13  production_comp         11779 non-null  object
14  original_language       11779 non-null  object
15  belongs_to_collection.name 1006 non-null  object
16  genres                  11779 non-null  object
17  release_year            11779 non-null  int64
dtypes: datetime64[ns](1), float64(5), int64(5), object(7)
memory usage: 1.7+ MB
```

In [93]:

```
1main_df.describe()
```

	startYear	runtimeMinutes	budget	world_collection	int_collection	dom_collection	popularity	vote_average	vote_count	release_year
count	11779.000000	11779.000000	1.177900e+04	1.177900e+04	1.092400e+04	2.765000e+03	11779.000000	11779.000000	11779.000000	11779.000000
mean	2017.283895	100.864080	4.205773e+06	1.681113e+07	1.246981e+07	2.061955e+07	10.595399	5.553162	284.940233	2017.374989
std	1.561950	28.732314	2.011748e+07	9.159840e+07	6.418805e+07	6.528838e+07	41.142722	2.258754	1241.668467	1.564721
min	2014.000000	0.000000	0.000000e+00	1.000000e+00	2.000000e+00	4.900000e+01	0.000000	0.000000	0.000000	2015.000000
25%	2016.000000	90.000000	0.000000e+00	3.603400e+04	3.774625e+04	3.667600e+04	1.279500	5.100000	3.000000	2016.000000
50%	2017.000000	100.000000	0.000000e+00	3.732710e+05	3.666200e+05	3.379070e+05	3.148000	6.100000	14.000000	2017.000000
75%	2019.000000	114.000000	0.000000e+00	3.255714e+06	2.972528e+06	7.743794e+06	8.796000	6.900000	78.000000	2019.000000
max	2021.000000	808.000000	3.560000e+08	2.797801e+09	1.939128e+09	9.366622e+08	2103.518000	10.000000	24543.000000	2020.000000

dropping ultra low budget movies along with 0, which means no information. And extremely low budget indicates possible error in data collection. Keeping low budget movies does not exactly match with the goal; finding good investment recommendation for a big company.

In [94]:

```
1main_df.shape
```

(11779, 18)

In [95]:

```
1main_df = main_df[main_df.budget>=5000]
```

In [96]:

```
1main_df.shape
```

(2081, 18)

focusing analysis only on movies where primary spoken language is English. MS should focus on this for the commencement.

In [97]:

```
1main_df = main_df[main_df.original_language=='en']
```

In [98]:

```
1main_df.shape
```

(1113, 18)

In [99]:

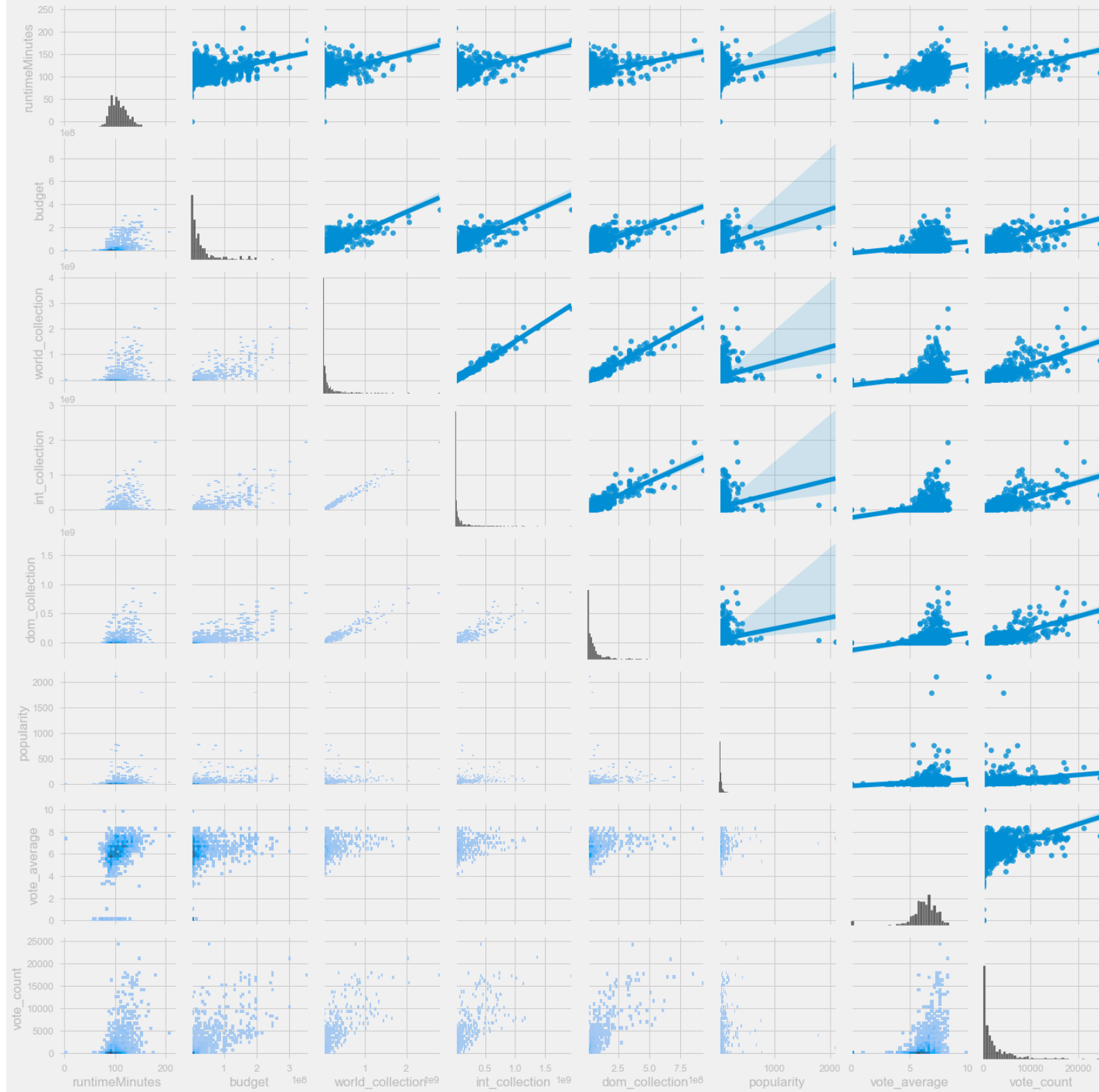
```
1main_df.columns
```

Index(['imdb_id', 'primaryTitle', 'originalTitle', 'startYear', 'release_date',
 'runtimeMinutes', 'budget', 'world_collection', 'int_collection',
 'dom_collection', 'popularity', 'vote_average', 'vote_count',
 'production_comp', 'original_language', 'belongs_to_collection.name',
 'genres', 'release_year'],
 dtype='object')

```
In [100]: 1 list_for_pairplot = ['release_date',  
2     'runtimeMinutes', 'budget', 'world_collection', 'int_collection',  
3     'dom_collection', 'popularity', 'vote_average', 'vote_count',  
4     'belongs_to_collection.name']
```

In [101]:

```
1 with plt.style.context('fivethirtyeight'):
2     g = sns.PairGrid(main_df[list_for_pairplot], layout_pad=.2)
3     g.map_diag(sns.histplot)
4     g.map_upper(sns.regplot)
5     g.map_lower(sns.histplot)
```



	No severe anamoly spotted in the graph which warrants further investigation .
▼	5.4 Feature engineering
▼	5.4.1 ROI
	<div>Here, budget is the estimator for cost.</div> <div>$\text{ROI (\%)} = \frac{\text{Total Revenue} - \text{Total Cost}}{\text{Total Cost}} \times 100$</div>
▼	5.4.1.1 Return on investment in \$ value
In [102]:	<pre>1main_df['ROI'] = main_df.world_collection - main_df.budget</pre>
▼	5.4.1.2 Return on investment in percentage, expressed in full, not in decimal
In [103]:	<pre>1main_df['ROI_percentage'] = (main_df.ROI / main_df.budget)*100</pre>
▼	5.5 Final Check to see that everything is in place
In [104]:	<pre>1main_df.shape</pre> <div>(1113, 20)</div>


```
In [105]: 1main_df.head()
```

	imdb_id	primaryTitle	originalTitle	startYear	release_date	runtimeMinutes	budget	world_collection	int_collection	dom_collection	popularity	vote_average	vote
4	tt0369610	Jurassic World	Jurassic World	2015	2015-06-06	124	150000000	1.671713e+09	1.018131e+09	652385625.0	63.489	6.6	1659
6	tt0385887	Motherless Brooklyn	Motherless Brooklyn	2019	2019-10-31	144	26000000	1.847774e+07	9.200000e+06	9277736.0	75.020	6.8	842
11	tt0437086	Alita: Battle Angel	Alita: Battle Angel	2019	2019-01-31	122	170000000	4.049805e+08	3.191423e+08	85838210.0	175.798	7.2	6343
12	tt0441881	Danger Close	Danger Close: The Battle of Long Tan	2019	2019-08-08	118	23934823	2.088085e+06	2.088085e+06	NaN	112.552	6.8	148
14	tt0443533	The History of Love	The History of Love	2016	2016-11-09	134	20000000	4.922720e+05	4.922720e+05	NaN	5.406	6.4	63

```
In [106]: 1main_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 1113 entries, 4 to 13585
Data columns (total 20 columns):
 #   Column                                Non-Null Count  Dtype
---  -
 0   imdb_id                              1113 non-null   object
 1   primaryTitle                         1113 non-null   object
 2   originalTitle                       1113 non-null   object
 3   startYear                           1113 non-null   int64
 4   release_date                        1113 non-null   datetime64[ns]
 5   runtimeMinutes                      1113 non-null   int64
 6   budget                             1113 non-null   int64
 7   world_collection                    1113 non-null   float64
 8   int_collection                      1039 non-null   float64
 9   dom_collection                      882 non-null    float64
10   popularity                          1113 non-null   float64
11   vote_average                        1113 non-null   float64
12   vote_count                          1113 non-null   int64
13   production_comp                     1113 non-null   object
14   original_language                   1113 non-null   object
15   belongs_to_collection.name         227 non-null    object
16   genres                             1113 non-null   object
17   release_year                        1113 non-null   int64
18   ROI                                 1113 non-null   float64
19   ROI_percentage                     1113 non-null   float64
dtypes: datetime64[ns](1), float64(7), int64(5), object(7)
memory usage: 182.6+ KB
```

In [107]:

1main_df.describe()

	startYear	runtimeMinutes	budget	world_collection	int_collection	dom_collection	popularity	vote_average	vote_count	release_year	ROI	ROI_perce
count	1113.000000	1113.000000	1.113000e+03	1.113000e+03	1.039000e+03	8.820000e+02	1113.000000	1113.000000	1113.000000	1113.000000	1.113000e+03	1113.000000
mean	2016.993711	107.323450	3.841841e+07	1.271450e+08	8.259442e+07	6.207522e+07	43.486649	6.261995	2338.715184	2017.052111	8.872658e+07	296.570055
std	1.532899	17.613393	5.251428e+07	2.613752e+08	1.761494e+08	1.039024e+08	103.773946	1.246853	3332.521354	1.522309	2.226630e+08	1647.220411
min	2014.000000	0.000000	5.000000e+03	5.470000e+02	5.470000e+02	1.377000e+03	0.600000	0.000000	0.000000	2015.000000	-1.510000e+08	-99.981875
25%	2016.000000	94.000000	6.000000e+06	2.084628e+06	1.177836e+06	5.622565e+06	13.550000	5.800000	266.000000	2016.000000	-3.898454e+06	-69.544079
50%	2017.000000	105.000000	1.900000e+07	2.935520e+07	1.424425e+07	2.740507e+07	22.168000	6.400000	1020.000000	2017.000000	8.197072e+06	63.263233
75%	2018.000000	118.000000	4.000000e+07	1.195200e+08	6.891399e+07	6.725403e+07	41.249000	7.000000	3038.000000	2018.000000	7.501105e+07	296.521358
max	2020.000000	209.000000	3.560000e+08	2.797801e+09	1.939128e+09	9.366622e+08	2103.518000	10.000000	24543.000000	2020.000000	2.441801e+09	42864.410011

if only want to focus on profitable movies

In [108]:	1# main_df = main_df[main_df.ROI>0]
-----------	-------------------------------------

▼

6 Exploratory data analysis

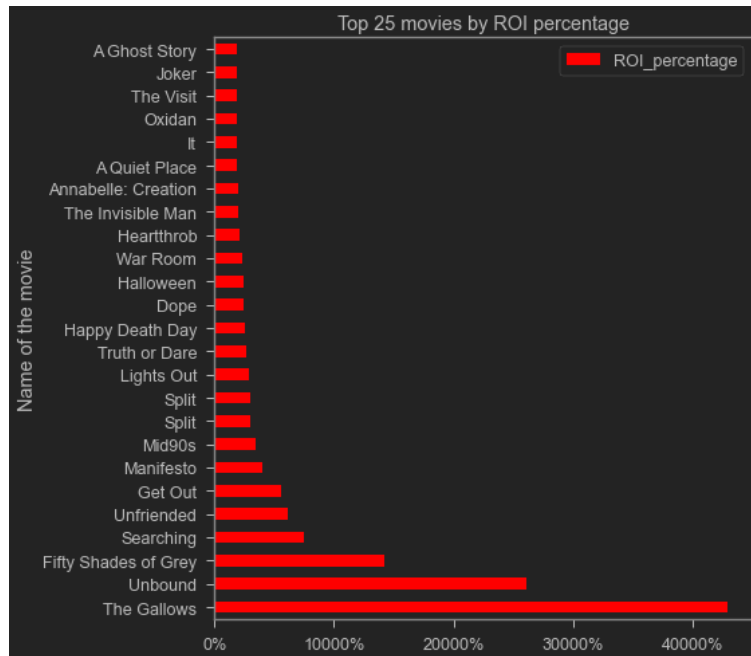
▼

6.1 EDA - top movie by return %

```

In [109]: 1ax = main_df.sort_values(by='ROI_percentage',
2                                ascending=False).head(25).plot(kind='barh',
3                                x='primaryTitle',
4                                y='ROI_percentage',
5                                color="red")
6plt.title('Top 25 movies by ROI percentage')
7plt.ylabel('Name of the movie')
8plt.tight_layout()
9ax.xaxis.set_major_formatter(format_add_percentage)
10plt.show()

```

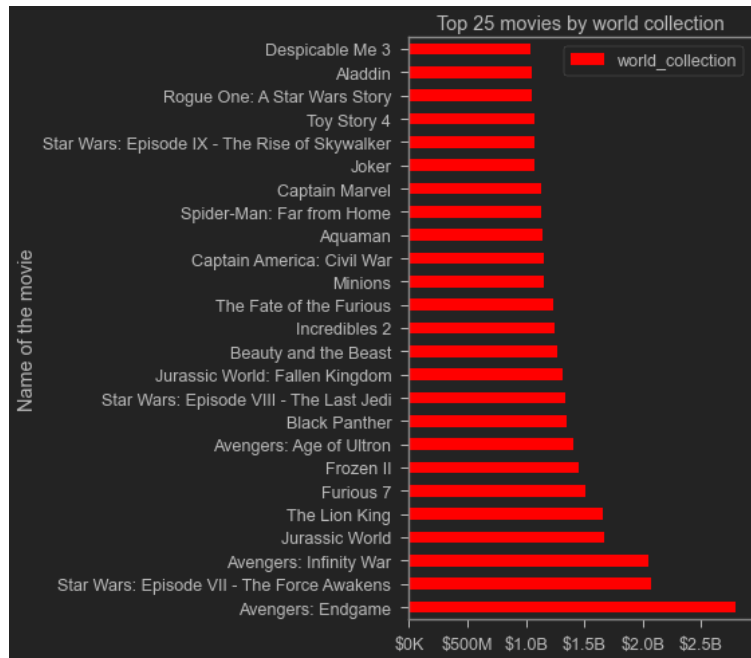


6.2 EDA - top movie by gross profit

```

In [110]:
1 ax = main_df.sort_values(by='world_collection',
2                           ascending=False).head(25).plot(kind='barh',
3                                                         x='primaryTitle',
4                                                         y='world_collection',
5                                                         color="red")
6 plt.title('Top 25 movies by world collection')
7 plt.ylabel('Name of the movie')
8 ax.xaxis.set_major_formatter(format_number)
9 plt.tight_layout()
10 plt.show()

```



6.3 EDA - profit by top 20 studio

```
In [111]: 1studio_df = main_df.copy()
```

```
In [112]: 1studio_df.loc[:, 'production_comp_exp'] = studio_df.production_comp.map(
2     lambda x: x.split(', '))
```

```
In [113]: 1studio_df_fig = studio_df.explode('production_comp_exp')
```

```
In [114]: 1studio_df_fig.head(3)
```

	imdb_id	primaryTitle	originalTitle	startYear	release_date	runtimeMinutes	budget	world_collection	int_collection	dom_collection	popularity	vote_average	vote
4	tt0369610	Jurassic World	Jurassic World	2015	2015-06-06	124	150000000	1.671713e+09	1.018131e+09	652385625.0	63.489	6.6	1659
4	tt0369610	Jurassic World	Jurassic World	2015	2015-06-06	124	150000000	1.671713e+09	1.018131e+09	652385625.0	63.489	6.6	1659
4	tt0369610	Jurassic World	Jurassic World	2015	2015-06-06	124	150000000	1.671713e+09	1.018131e+09	652385625.0	63.489	6.6	1659

```
In [115]: 1top_production_house_list = list(
2     studio_df_fig.production_comp_exp.value_counts().sort_values(
3     ascending=False)[:20].index)
```

```
In [116]: 1# to get Total worldwide $ collection by top 20 studios over the years
2studio_df_fig_0 = studio_df_fig[studio_df_fig['production_comp_exp'].isin(
3     top_production_house_list)]
```

```
In [117]: 1# Total worldwide $ collection by top 20 studios
2studio_df_fig_1 = studio_df_fig.groupby(
3     by='production_comp_exp').agg('sum').sort_values(by='world_collection',
4     ascending=False)[:20]
5# Total releases by top 20 studios
6studio_df_fig_2 = studio_df_fig.groupby(
7     by='production_comp_exp').agg('count').sort_values(by='world_collection',
8     ascending=False)[:20]
```

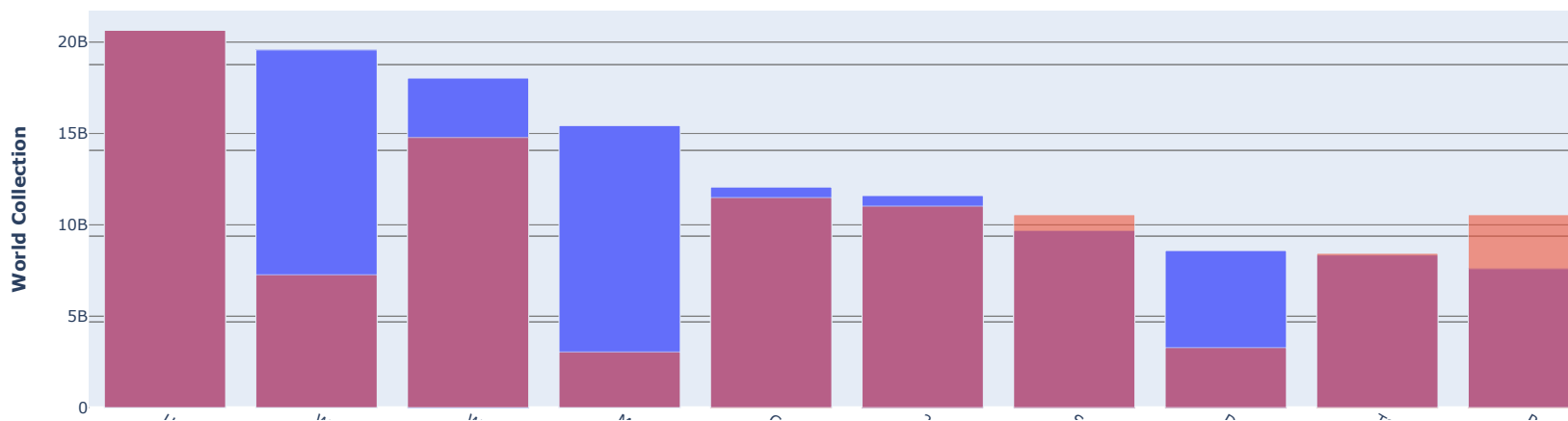
```
In [118]: 1# Collection Performance of top 10 movie studios
          2studio_df_fig_merged = pd.merge(
          3     studio_df_fig.groupby(by='production_comp_exp').agg('sum').sort_values(
          4         by='world_collection',
          5         ascending=False)['world_collection'].reset_index(),
          6     studio_df_fig.groupby(by='production_comp_exp').agg('count').sort_values(
          7         by='world_collection',
          8         ascending=False)['world_collection'].reset_index(),
          9     on='production_comp_exp')
          10# Budget Performance of top 10 movie studios
          11studio_df_fig_merged_1 = pd.merge(
          12     studio_df_fig.groupby(by='production_comp_exp').agg('sum').sort_values(
          13         by='budget',
          14         ascending=False)['budget'].reset_index(),
          15     studio_df_fig.groupby(by='production_comp_exp').agg('count').sort_values(
          16         by='budget',
          17         ascending=False)['budget'].reset_index(),
          18     on='production_comp_exp')
```

```

In [119]: 1## from https://plotly.com/python/multiple-axes/ ##official plotly how to instructions
2fig = make_subplots(specs=[[{"secondary_y": True}]]))
3# Add traces
4fig.add_trace(
5     go.Bar(x=studio_df_fig_merged.production_comp_exp[:10],
6             y=studio_df_fig_merged.world_collection_x[:10],
7             name="World Collection",
8             offset=True),
9     secondary_y=False,
10)
11fig.add_trace(
12     go.Bar(x=studio_df_fig_merged.production_comp_exp[:10],
13             y=studio_df_fig_merged.world_collection_y[:10],
14             name="Movie Released",
15             offset=True,
16             opacity=.6),
17     secondary_y=True,
18)
19# Add figure title
20fig.update_layout(title_text="Collection performance of top 10 movie studios")
21# Set x-axis title
22fig.update_xaxes(title_text="World Collection")
23# Set y-axes titles
24fig.update_yaxes(title_text="<b>World Collection</b>", secondary_y=False)
25fig.update_yaxes(title_text="<b>Number of Movie Released</b>",
26                 secondary_y=True)
27fig.show()

```

Collection performance of top 10 movie studios



Uni.

Wai.

Wai.

Mar.

Col.

<0%

50%

DEF.

ISr.

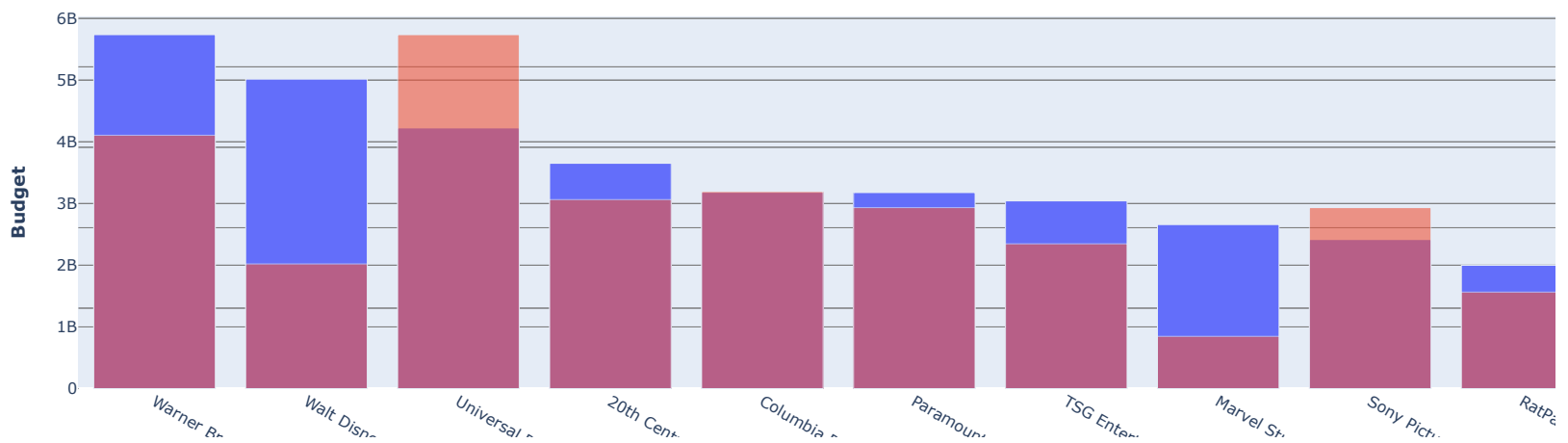
Far.


```

In [120]: 1## from https://plotly.com/python/multiple-axes/ ##official plotly how to instructions
2fig = make_subplots(specs=[[{"secondary_y": True}]])
3# Add traces
4fig.add_trace(
5     go.Bar(x=studio_df_fig_merged_1.production_comp_exp[:10],
6             y=studio_df_fig_merged_1.budget_x[:10],
7             name="Budget",
8             offset=True),
9     secondary_y=False,
10)
11fig.add_trace(
12     go.Bar(x=studio_df_fig_merged_1.production_comp_exp[:10],
13            y=studio_df_fig_merged_1.budget_y[:10],
14            name="Movie Released",
15            offset=True,
16            opacity=.6),
17     secondary_y=True,
18)
19# Add figure title
20fig.update_layout(title_text="Budget performance of top 10 movie studios")
21# Set x-axis title
22fig.update_xaxes(title_text="Budget")
23# Set y-axes titles
24fig.update_yaxes(title_text="<b>Budget</b>", secondary_y=False)
25fig.update_yaxes(title_text="<b>Number of Movie Released</b>",
26                 secondary_y=True)
27fig.show()

```

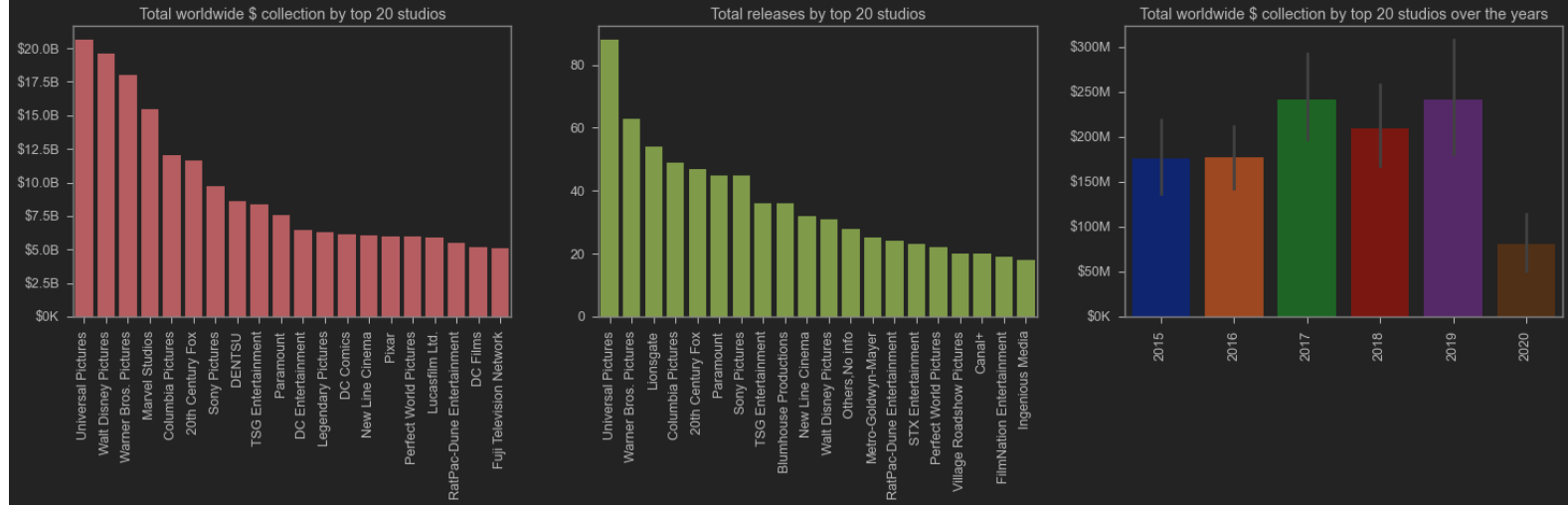
Budget performance of top 10 movie studios



	Marvel Studios and Walt Disney has the best release count to world collection ratio. It took Universal Pictures way more budget to achieve the top spot.

In [121]:

```
1plt.figure(figsize=(25, 5))
2plt.subplot(1, 3, 1)
3plt.xticks(rotation='vertical')
4
5plt.title('Total worldwide $ collection by top 20 studios')
6
7sns.barplot(y='world_collection',
8            x='production_comp_exp',
9            data=studio_df_fig_1.reset_index(),
10           color='r').yaxis.set_major_formatter(
11               format_number)
12plt.xlabel(None)
13plt.ylabel(None)
14
15
16plt.subplot(1, 3, 2)
17sns.barplot(y='world_collection',
18            x='production_comp_exp',
19            data=studio_df_fig_2.reset_index(),
20            color='g').set(xlabel=None, ylabel=None)
21
22plt.xticks(rotation='vertical')
23
24plt.title('Total releases by top 20 studios')
25
26plt.subplot(1, 3, 3)
27sns.barplot(data=studio_df_fig_0,
28            x='release_year',
29            y='world_collection',
30            palette='dark').yaxis.set_major_formatter(
31               format_number)
32plt.title('Total worldwide $ collection by top 20 studios over the years')
33plt.xticks(rotation='vertical')
34plt.xlabel(None)
35plt.ylabel(None)
36# plt.tight_layout()
37plt.show()
```



One caveat of this graph is that because of the nature of the data, if a movie has multiple studios attached to it then all earnings of it is counted as the studios sole earnings. This is the reason why the mismatch of metrics. All things set aside, from this graph a visual understanding can be achieved about top studios without trying to make sense of the numbers. Turning off ylabels of first two plots can help on that regard.

6.4 EDA - Relation between features

```
1corelation_filter_list = ['startYear',
2    'runtimeMinutes', 'budget', 'world_collection', 'int_collection',
3    'dom_collection', 'popularity', 'vote_average', 'vote_count']
```

```
1corr_df = main_df[corelation_filter_list]
```

```
1corr_df_matrix = corr_df.corr()
```

```
1corr_df_matrix.style.background_gradient()
```

	startYear	runtimeMinutes	budget	world_collection	int_collection	dom_collection	popularity	vote_average	vote_count
startYear	1.000000	-0.003422	0.053363	0.024102	0.035525	0.040447	0.260186	0.054454	-0.065041
runtimeMinutes	-0.003422	1.000000	0.426111	0.349988	0.351262	0.328371	0.159567	0.368076	0.439491
budget	0.053363	0.426111	1.000000	0.783042	0.787109	0.716023	0.321358	0.237517	0.675298
world_collection	0.024102	0.349988	0.783042	1.000000	0.986302	0.955519	0.236664	0.256312	0.791444
int_collection	0.035525	0.351262	0.787109	0.986302	1.000000	0.892850	0.237610	0.277242	0.763741
dom_collection	0.040447	0.328371	0.716023	0.955519	0.892850	1.000000	0.201020	0.267886	0.773461
popularity	0.260186	0.159567	0.321358	0.236664	0.237610	0.201020	1.000000	0.156417	0.241333
vote_average	0.054454	0.368076	0.237517	0.256312	0.277242	0.267886	0.156417	1.000000	0.366568
vote_count	-0.065041	0.439491	0.675298	0.791444	0.763741	0.773461	0.241333	0.366568	1.000000

```
In [126]: 1correlation_top_bottom(corr_df_matrix)

Positive correlations:
      index      feature_combo  correlation
0    31  int_collection and world_collection    0.986302
1    32  dom_collection and world_collection    0.955519
2    41  int_collection and dom_collection    0.892850
3    35  vote_count and world_collection    0.791444
4    22      int_collection and budget    0.787109
5    21      budget and world_collection    0.783042
6    53  dom_collection and vote_count    0.773461
7    44  int_collection and vote_count    0.763741
8    23      dom_collection and budget    0.716023
9    26      budget and vote_count    0.675298

-----
Negative correlations:
      index      feature_combo  correlation
0     8  startYear and vote_count   -0.065041
1     1  runtimeMinutes and startYear   -0.003422
2     3  startYear and world_collection    0.024102
3     4  int_collection and startYear    0.035525
4     5  dom_collection and startYear    0.040447
5     2      startYear and budget    0.053363
6    18      budget and startYear    0.053363
7     7  vote_average and startYear    0.054454
8    61  vote_average and popularity    0.156417
9    15  runtimeMinutes and popularity    0.159567
```

▼

6.4.1 Findings and observation

From those table it can be observed that world collection, international collection and domestic collection is highly correlated. It is expected as world collection is a dependent variable of the other two. And the later two are highly correlated. This is also expected, as this is indicator of a profitable versus flop movie. Better performing movies has higher popularity as explained by world collection versus vote count, vice versa. High budget movies perform better overseas.

Overall budget is the key for indicating performance both in international and domestic performance and feedback from movie consumers. No other standout correlation was found.

▼

7 Recommendations

▼

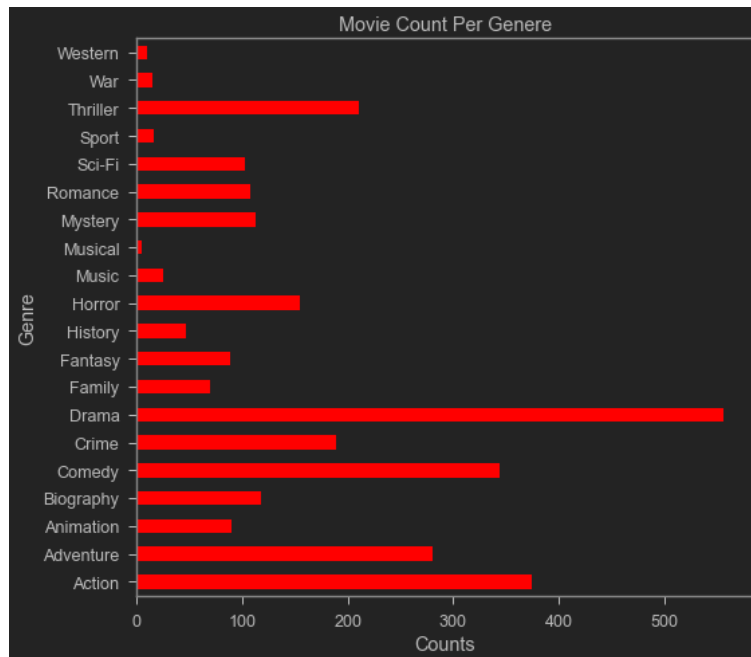
7.1 Which genre of movie to make, explained by top movie per genre

```
In [127]: 1genere_df = main_df.copy()
```

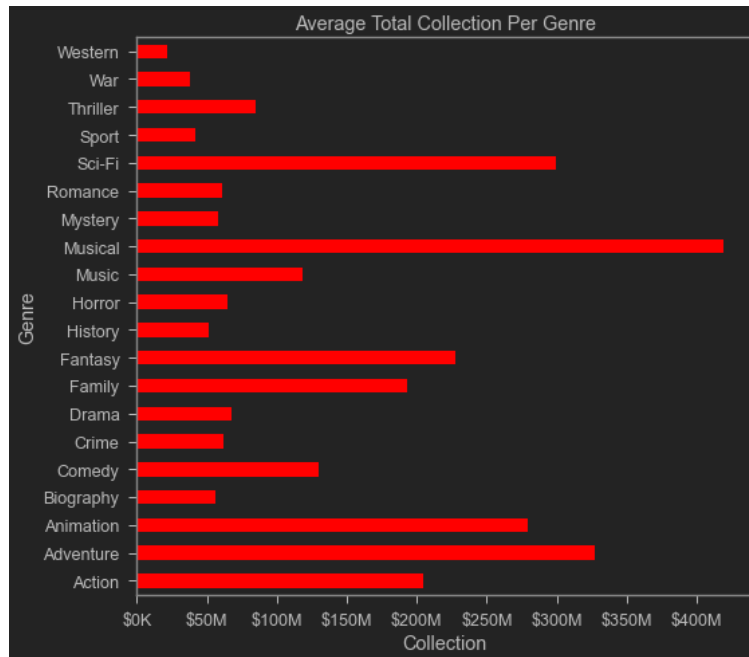
```
In [128]: 1genere_df.loc[:, 'genres_exp'] = genere_df.genres.map(lambda x: x.split(','))
```

```
In [129]: 1genere_df_fig = genere_df.explode('genres_exp')
```

```
In [130]: 1# Movie Count Per Genre
          2genre_df_fig.groupby('genres_exp').count()['imdb_id'].plot(kind='barh',
          3                                                    color="red")
          4plt.title('Movie Count Per Genre')
          5plt.ylabel('Genre')
          6plt.xlabel('Counts')
          7plt.tight_layout()
          8plt.show()
```



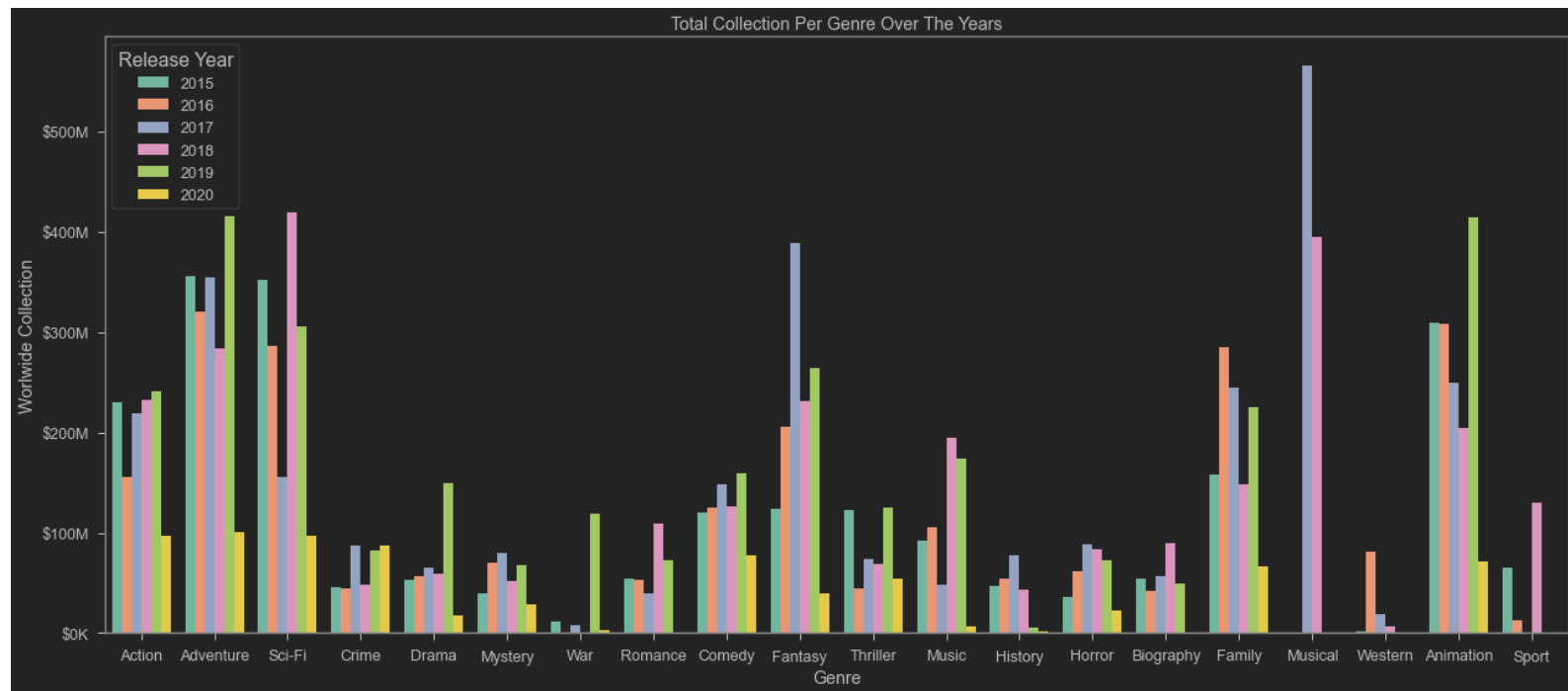
```
In [131]: 1# Average Total Collection Per Genre
2genre_df.groupby('genres_exp').mean()['world_collection'].plot(
3     kind='barh', color="red").xaxis.set_major_formatter(format_number)
4plt.title('Average Total Collection Per Genre')
5plt.ylabel('Genre')
6plt.xlabel('Collection')
7plt.tight_layout()
8plt.show()
```



```

In [132]:
1# styling
2# sns.set_style('ticks')
3fig, ax = plt.subplots()
4fig.set_size_inches(18, 8)
5# plotting
6sns.barplot(data=genere_df_fig,
7            x='genres_exp',
8            y='world_collection',
9            hue='release_year',
10           palette='Set2',
11           ci=None).yaxis.set_major_formatter(format_number)
12plt.title('Total Collection Per Genre Over The Years')
13plt.ylabel('Worldwide Collection')
14plt.xlabel('Genre')
15plt.legend(title='Release Year', title_fontsize='large')
16plt.tight_layout()
17plt.show()

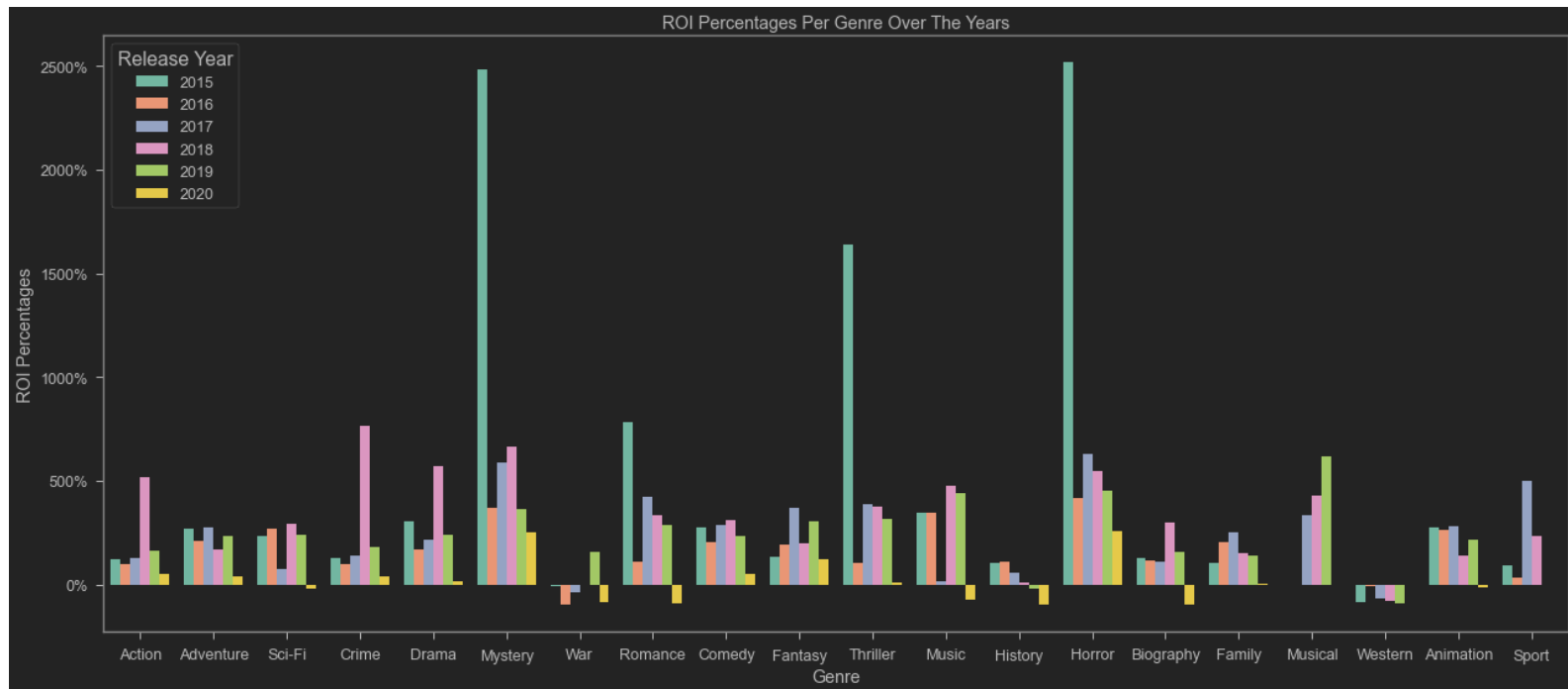
```




```

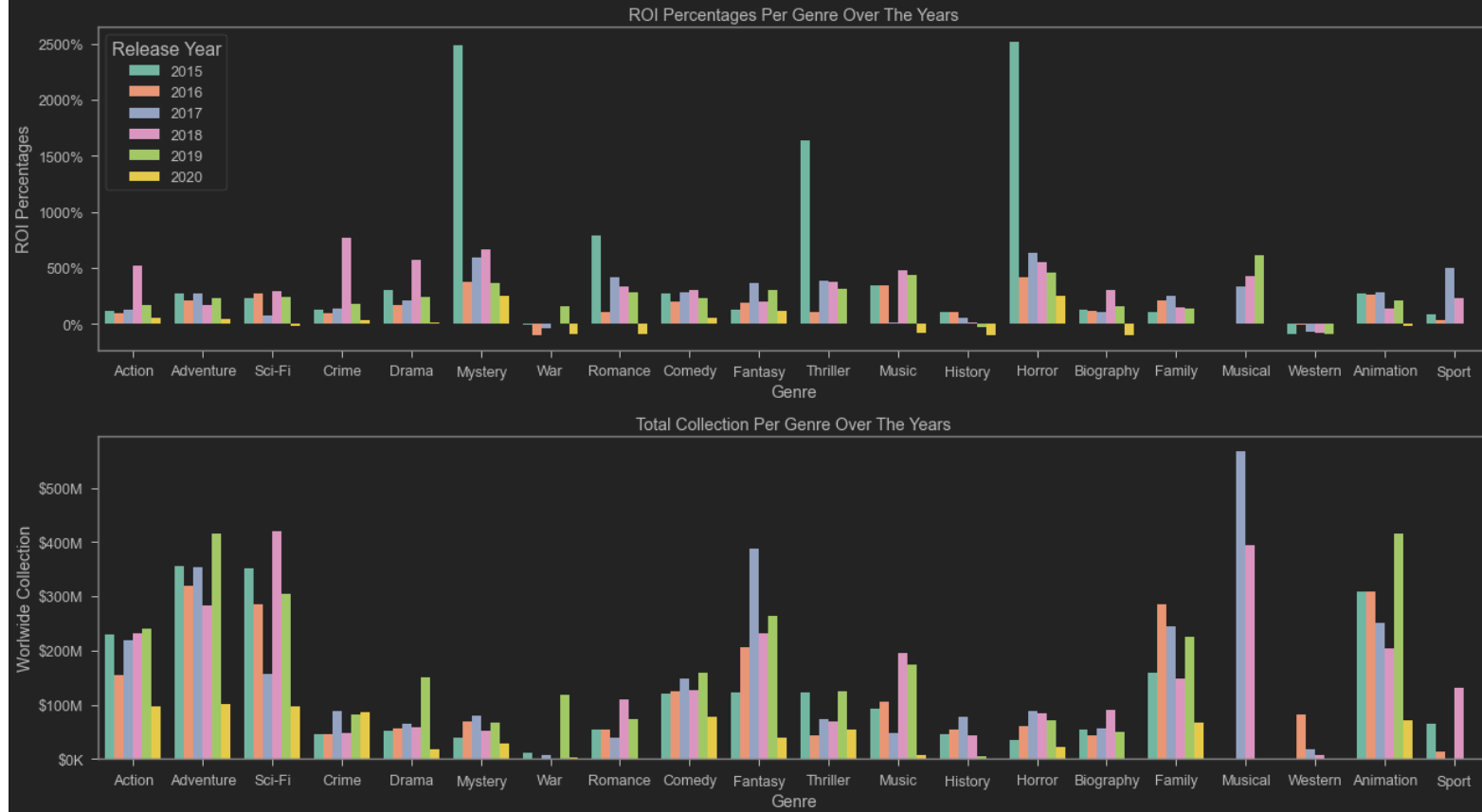
In [133]: 1# styling
          2# sns.set_style('ticks')
          3fig, ax = plt.subplots()
          4fig.set_size_inches(18, 8)
          5# plotting
          6sns.barplot(data=genere_df_fig,
          7             x='genres_exp',
          8             y='ROI_percentage',
          9             hue='release_year', palette='Set2',ci=None).yaxis.set_major_formatter(format_add_percentage)
10plt.title('ROI Percentages Per Genre Over The Years')
11plt.ylabel('ROI Percentages')
12plt.xlabel('Genre')
13plt.legend(title='Release Year', title_fontsize= 'large')
14plt.tight_layout()
15plt.show()

```



In [134]:

```
1plt.figure(figsize=(18, 10))
2# plotting
3plt.subplot(2, 1, 1)
4sns.barplot(data=genere_df_fig,
5            x='genres_exp',
6            y='ROI_percentage',
7            hue='release_year',
8            palette='Set2',
9            ci=None).yaxis.set_major_formatter(format_add_percentage)
10plt.title('ROI Percentages Per Genre Over The Years')
11plt.ylabel('ROI Percentages')
12plt.xlabel('Genre')
13plt.legend(title='Release Year', title_fontsize='large')
14plt.tight_layout()
15
16plt.subplot(2, 1, 2)
17sns.barplot(data=genere_df_fig,
18            x='genres_exp',
19            y='world_collection',
20            hue='release_year',
21            palette='Set2',
22            ci=None).yaxis.set_major_formatter(format_number)
23plt.title('Total Collection Per Genre Over The Years')
24plt.ylabel('Worldwide Collection')
25plt.xlabel('Genre')
26plt.legend().remove()
27plt.tight_layout()
28
29plt.show()
```



2015 was a good year for the industry. Animation has good performance but costly to make, hence lower percentage. Muscial had few good years then fell out of fashion. Action, Adventure, Family, Fantasy has been consistent performers. Horror and Mystery has high return percentage.

7.1.1 Action suggestion

Any one or combination of Action, Adventure, Animation is recommended. Animation and Action has 35% chance for occurring as genre combo. There is no landslide winner here, although this graphs can be used to figure out which one to avoid, for example western and war.

7.2 Best time to release movie

```
In [135]: 1timing_df = main_df.copy()
```

```
In [136]: 1timing_df['release_month']=timing_df['release_date'].dt.month
```

- minor feature engineering

Release months are put in three bins based on market analysts opinion. The [dump months](https://en.wikipedia.org/wiki/Dump_months) (https://en.wikipedia.org/wiki/Dump_months) are what the film community calls the two periods of the year when there are lowered commercial and critical expectations for most new releases from American filmmakers and distributors.

1. January - May: Dump month
2. June - July: Summer
3. August - October: Dump month
4. November - December: Holidays

In [137]:

```
1 timing_df['release_timing'] = pd.cut(  
2     timing_df['release_month'],  
3     bins=[0, 6, 8, 10, 12],  
4     labels=['dump months', 'summer', 'dump months', 'holidays'],  
5     ordered=False)
```

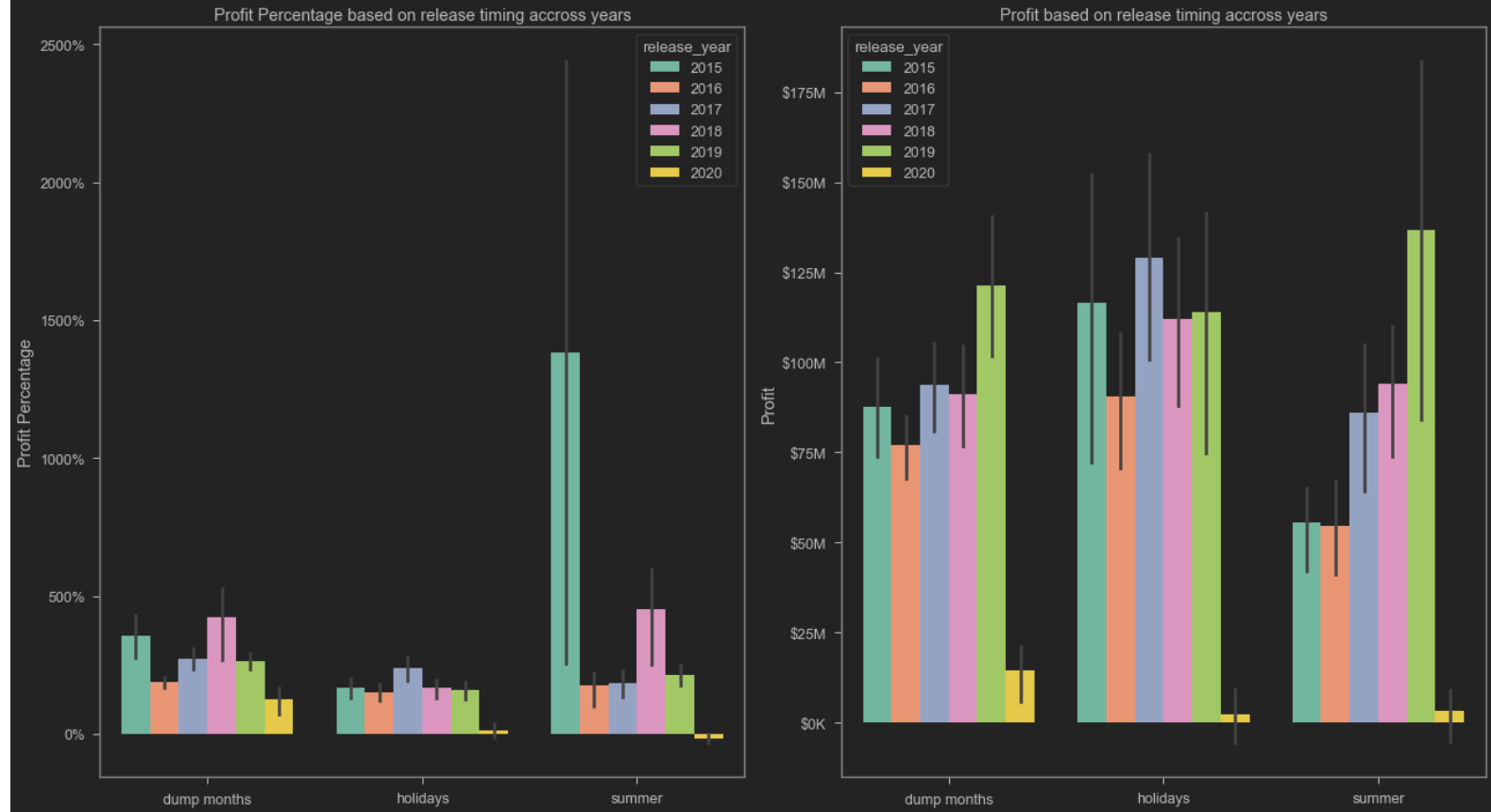
In [138]:

```
1 timing_df.head(3)
```

	imdb_id	primaryTitle	originalTitle	startYear	release_date	runtimeMinutes	budget	world_collection	int_collection	dom_collection	popularity	vote_average	vote
4	tt0369610	Jurassic World	Jurassic World	2015	2015-06-06	124	150000000	1.671713e+09	1.018131e+09	652385625.0	63.489	6.6	1659
6	tt0385887	Motherless Brooklyn	Motherless Brooklyn	2019	2019-10-31	144	26000000	1.847774e+07	9.200000e+06	9277736.0	75.020	6.8	842
11	tt0437086	Alita: Battle Angel	Alita: Battle Angel	2019	2019-01-31	122	170000000	4.049805e+08	3.191423e+08	85838210.0	175.798	7.2	6343

In [139]:

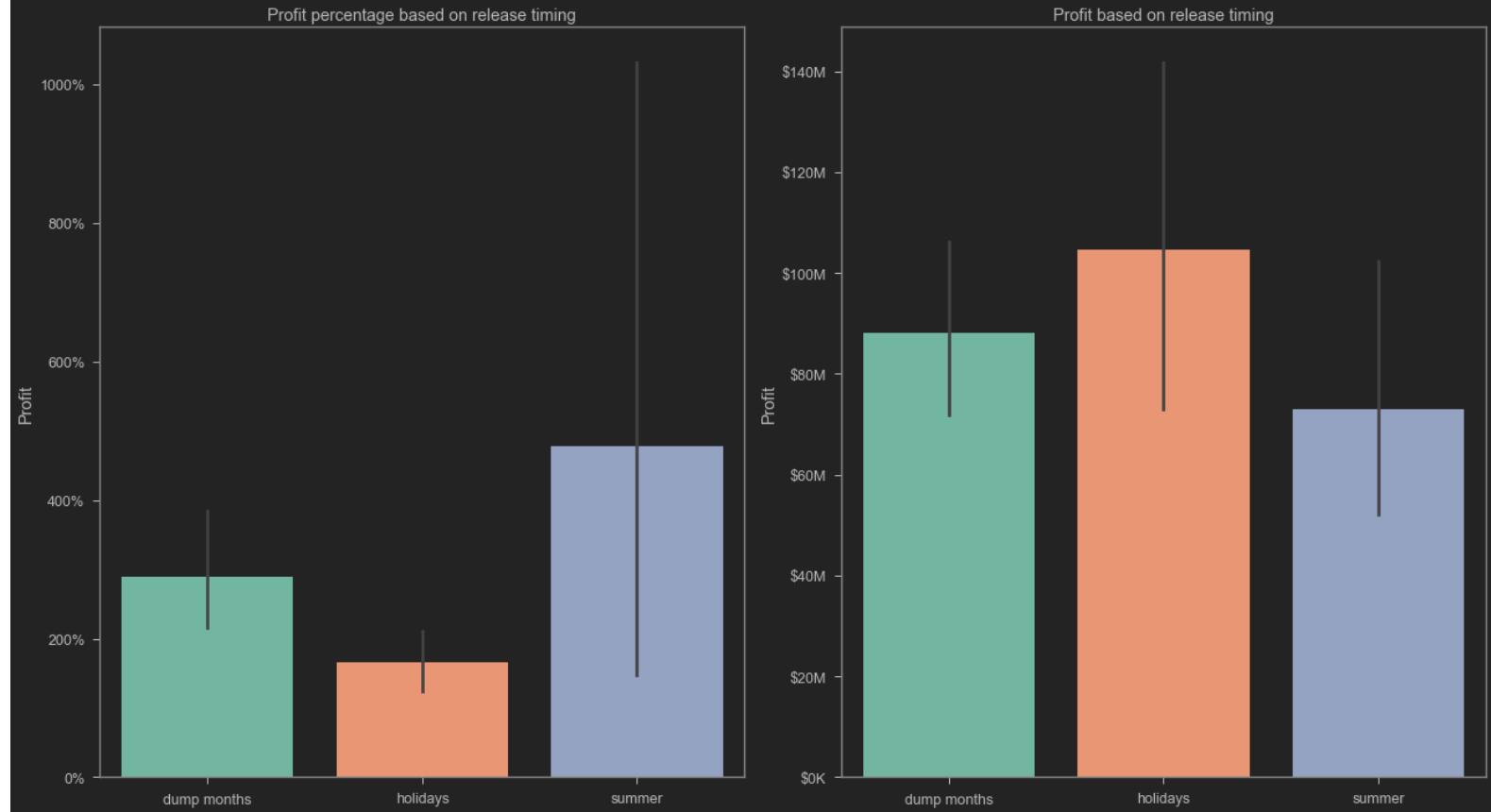
```
1plt.figure(figsize=(18, 10))
2plt.subplot(1, 2, 1)
3sns.barplot(data=timing_df,
4            x='release_timing',
5            y='ROI_percentage',
6            hue='release_year',palette='Set2',
7            ci=50).yaxis.set_major_formatter(format_add_percentage)
8plt.title('Profit Percentage based on release timing accross years')
9plt.ylabel('Profit Percentage')
10plt.xlabel("")
11
12plt.subplot(1, 2, 2)
13sns.barplot(data=timing_df,
14            x='release_timing',
15            y='ROI',
16            hue='release_year',palette='Set2',
17            ci=50).yaxis.set_major_formatter(format_number)
18plt.title('Profit based on release timing accross years')
19plt.ylabel('Profit')
20plt.xlabel("")
21plt.tight_layout()
22
23plt.show()
```



2015's Summer was good in terms of percentage return but weirdly did not generate much cash. Releasing movie in the holidays season is the safest bet. But summer is having a consistent raise, except for 2020. 2020's summer was not normal by any means, thus this is expected.

In [140]:

```
1plt.figure(figsize=(18, 10))
2plt.subplot(1, 2, 2)
3sns.barplot(data=timing_df, x='release_timing', y='ROI',
4            palette='Set2').yaxis.set_major_formatter(format_number)
5plt.title('Profit based on release timing')
6plt.ylabel('Profit')
7plt.xlabel("")
8plt.tight_layout()
9
10plt.subplot(1, 2, 1)
11sns.barplot(data=timing_df,
12            x='release_timing',
13            y='ROI_percentage',
14            palette='Set2').yaxis.set_major_formatter(format_add_percentage)
15plt.title('Profit percentage based on release timing')
16plt.ylabel('Profit')
17plt.xlabel("")
18plt.tight_layout()
19
20plt.show()
```

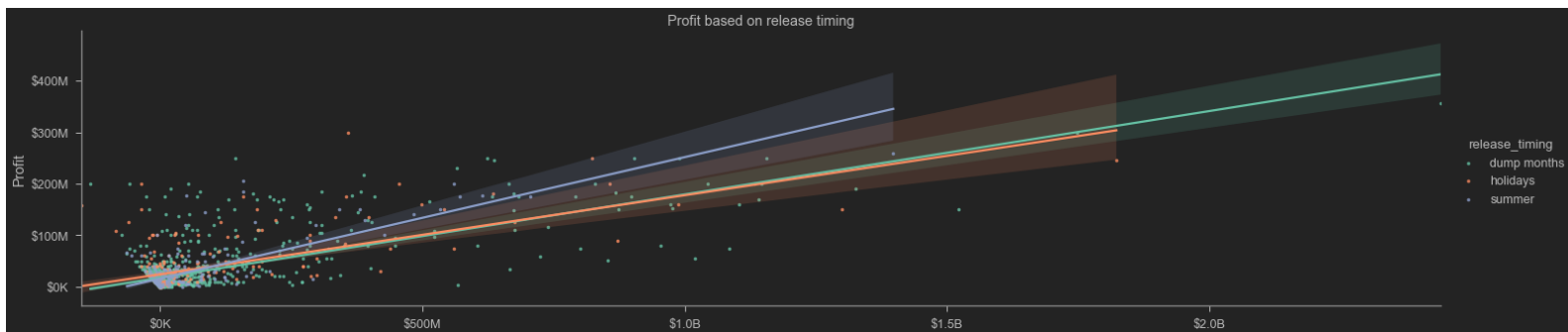


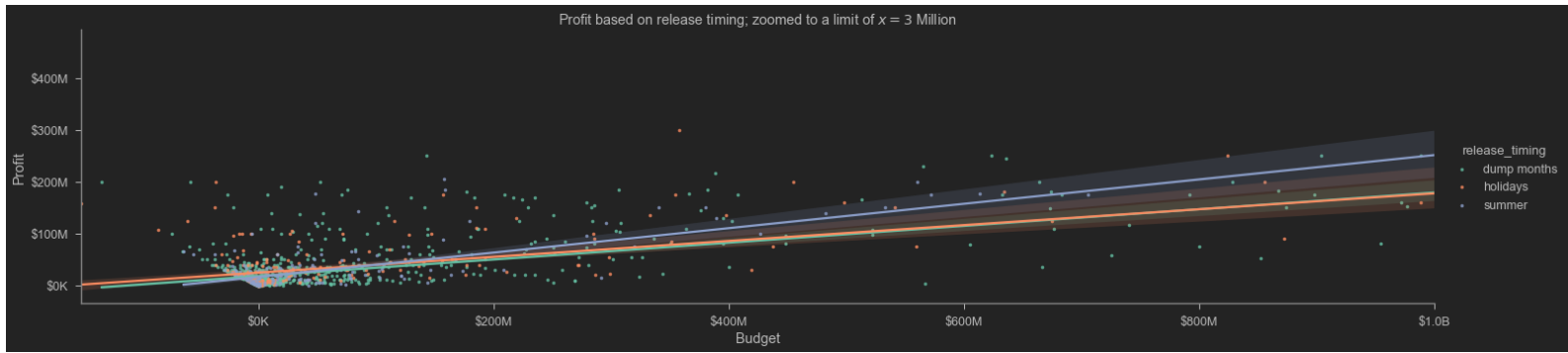
Movies released in holidays earn consistent returns but costs more. Summer is more dollar generating and volatile in a good way, on a uptrend.


```

In [141]:
1# Profit based on release timing
2g = sns.lmplot(data=timing_df,
3               x='ROI',
4               y='budget',
5               hue='release_timing',
6               fit_reg=True,
7               markers='.',
8               aspect=4,palette='Set2',
9               robust=True)
10for ax in g.axes.flat:
11    ax.yaxis.set_major_formatter(format_number)
12    ax.xaxis.set_major_formatter(format_number)
13plt.title('Profit based on release timing')
14plt.ylabel('Profit')
15plt.xlabel("")
16
17g = sns.lmplot(data=timing_df,
18               x='ROI',
19               y='budget',
20               hue='release_timing',
21               fit_reg=True,
22               markers='.',
23               aspect=4,palette='Set2',
24               robust=True)
25plt.xlim(right=1000000000)
26for ax in g.axes.flat:
27    ax.yaxis.set_major_formatter(format_number)
28    ax.xaxis.set_major_formatter(format_number)
29plt.title('Profit based on release timing; zoomed to a limit of $x=3$ Million')
30plt.ylabel('Profit')
31plt.xlabel("Budget")
32
33plt.show()

```



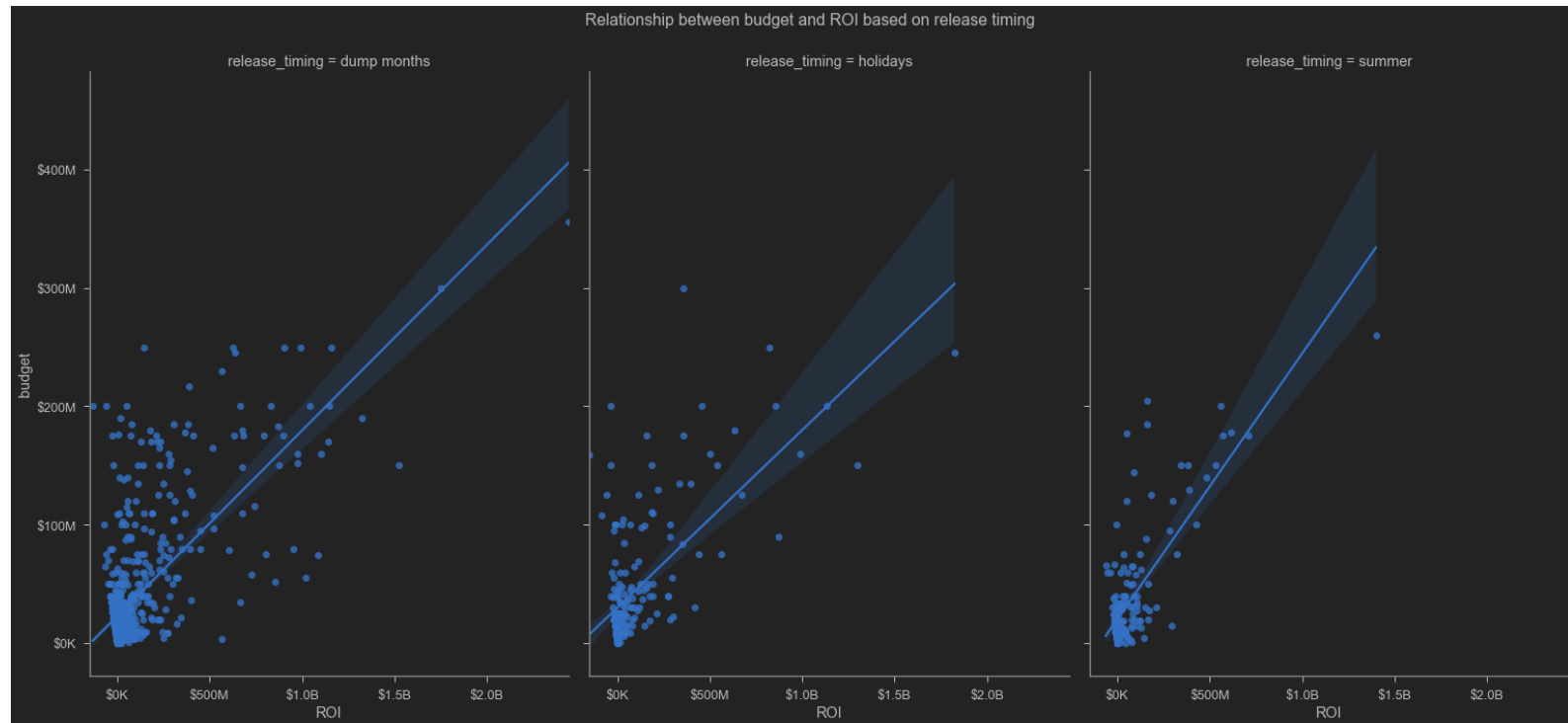


```

In [142]:
1# Relationship between budget and ROI based on release timing
2g = sns.FacetGrid(
3    timing_df, col='release_timing',
4    height=10, aspect=.7, palette='Set2')
5g.map(sns.regplot, 'ROI', 'budget')
6for ax in g.axes.flat:
7    ax.yaxis.set_major_formatter(format_number)
8    ax.xaxis.set_major_formatter(format_number)
9g.fig.subplots_adjust(top=0.9)
10g.fig.suptitle('Relationship between budget and ROI based on release timing')

```

```
Text(0.5, 0.98, 'Relationship between budget and ROI based on release timing')
```



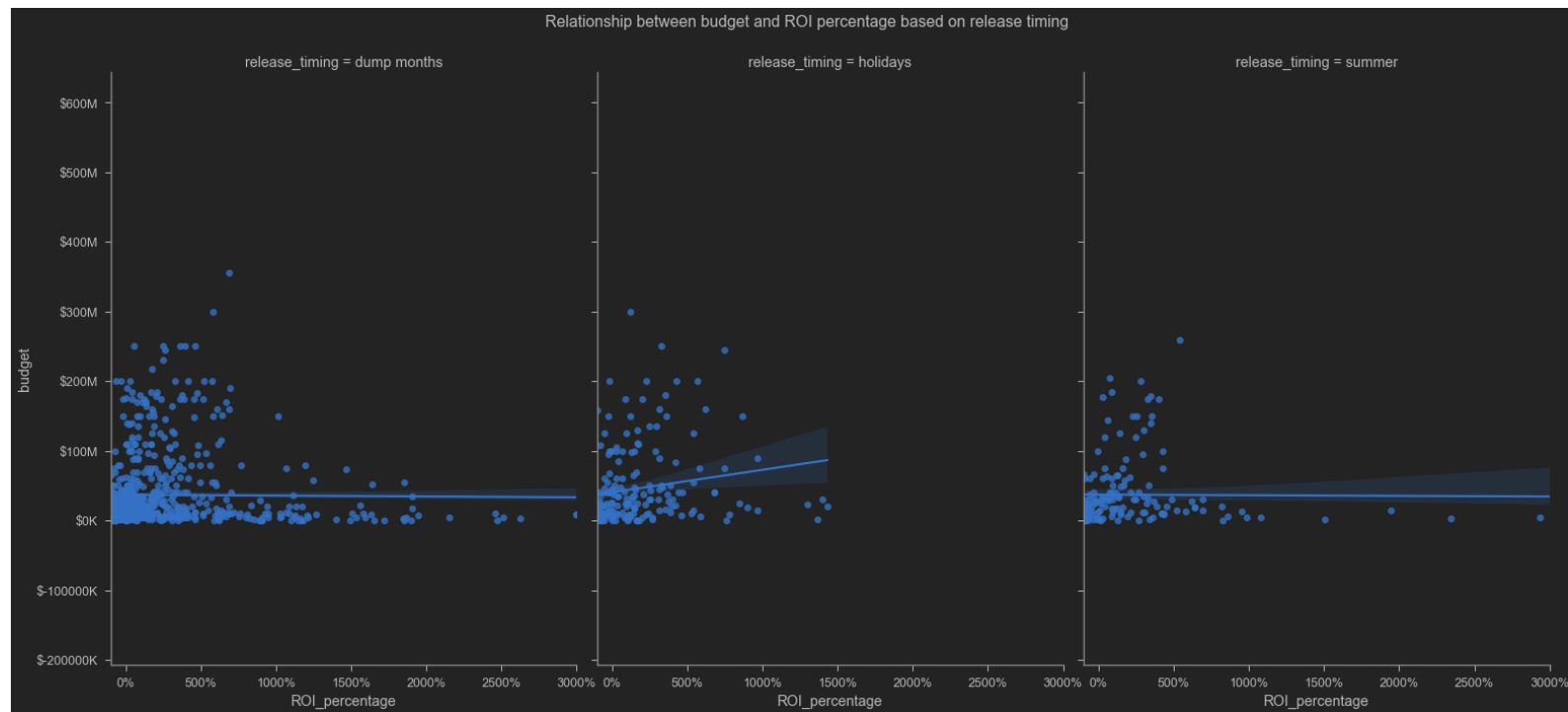
Producing movies for summer release is more costly, but return is steeper. Number of movies beyond 500 million is more frequent as well as observation counts are higher for holidays release, and the line is flatter meaning less costly to produce. Holidays season is the better option.

```

In [143]: 1# Relationship between budget and ROI based on release timing
          2g = sns.FacetGrid(
          3     timing_df, col='release_timing',
          4     height=10, aspect=.7, palette='Set2')
          5g.map(sns.regplot, 'ROI_percentage', 'budget')
          6plt.xlim(right=3000)
          7for ax in g.axes.flat:
          8     ax.yaxis.set_major_formatter(format_number)
          9     ax.xaxis.set_major_formatter(format_add_percentage)
         10g.fig.subplots_adjust(top=0.9)
         11g.fig.suptitle('Relationship between budget and ROI percentage based on release timing')

```

```
Text(0.5, 0.98, 'Relationship between budget and ROI percentage based on release timing')
```



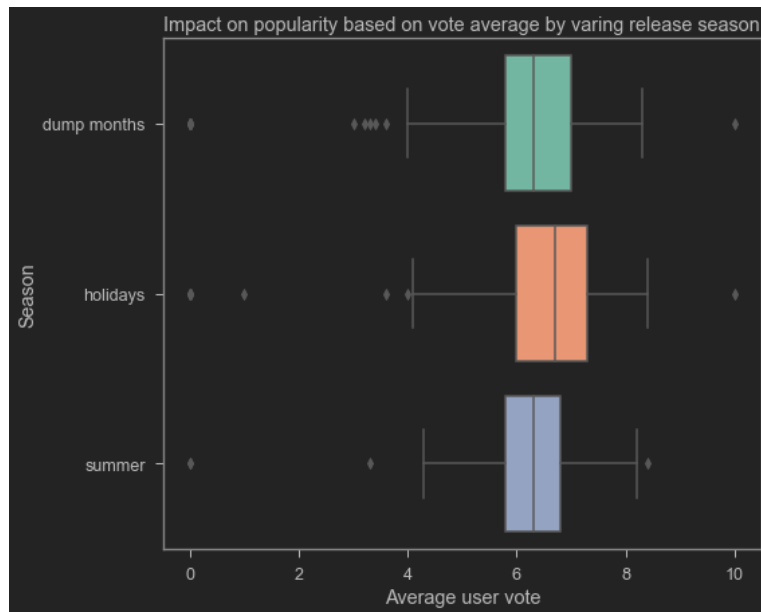
summer has the best earning potential. But line is steeper for holidays, confirming the point made on the previous graph.

```

In [144]: 1# impact on popularity based on vote average by varing release season
2sns.boxplot(x='vote_average', y='release_timing',data=timing_df, palette='Set2')
3plt.title('Impact on popularity based on vote average by varing release season')
4plt.ylabel('Season')
5plt.xlabel("Average user vote")

Text(0.5, 0, 'Average user vote')

```



Holidays movies are more popular and catch people on good mood maybe? Or content is less experimental. Reasoning can not be drawn from this figure but it can be said that holidays movies are more popular, which is good for entering the market with a more favorable impression on people.

▼ 7.2.1 Action suggestion

My recommendation is to focus for release schedule in the holidays season. There is higher probability of financial and critical success for movies released in that time frame. It is relatively cheaper to make than the next best option; i.e., Summer.

▼ 7.3 Franchise performance analysis leading to recommendation

```

In [145]: 1# getting a copy of main_df
2franchise_df_main = main_df.copy()

```

```

In [146]: 1# getting all movies that are part of a franchise
2franchise_df = franchise_df_main[~main_df['belongs_to_collection.name'].isna(
3)]

```

In [147]:	<div>▼</div> <pre>1# renaming column for use later 2franchise_df = franchise_df.rename(3 columns={"belongs_to_collection.name": "belongs_to_collection"})</pre>	
In [148]:	<div>▼</div> <pre>1# getting all movies that are not part of a franchise, yet! 2non_franchise_df = franchise_df_main[3 main_df['belongs_to_collection.name'].isna()].copy()</pre>	
▼	<h3>7.3.1 Franchise info</h3>	
	By franchise I mean serialization of movies either based on a related intellectual property or sharing same cinematic universe.	
In [149]:	<div>▼</div> <pre>1# list of unique franchise names 2list_of_franchise = franchise_df['belongs_to_collection'].unique()</pre>	
In [150]:	<pre>1franchise_df_ = franchise_df.groupby('belongs_to_collection').mean(2).ROI_percentage.sort_values(ascending=False).reset_index()</pre>	

```
In [151]:  
▼ 1# formatting  
▼ 2format_dict = {  
3     'ROI': '${0:,.0f}',  
4     'budget': '${0:,.0f}',  
5     'ROI_percentage': '{:.2f}%'  
6}  
7# performance of movies that are part of a franchise  
▼ 8franchise_df.groupby('belongs_to_collection').mean()[[  
9     'ROI_percentage', 'ROI', 'budget'  
▼ 10]].sort_values(  
11     by='ROI_percentage',  
▼ 12     ascending=False)[:20].style.format(format_dict).background_gradient(  
13         cmap='afmhot')
```

	ROI_percentage	ROI	budget
belongs_to_collection			
The Gallows Collection	42864.41%	\$42,864,410	\$100,000
Searching Collection	7446.20%	\$74,462,037	\$1,000,000
Fifty Shades Collection	5115.27%	\$403,622,827	\$38,000,000
Unfriended Collection	3845.35%	\$38,453,538	\$1,000,000
Lights Out Collection	2938.14%	\$143,968,835	\$4,900,000
Halloween Collection	2456.15%	\$245,614,941	\$10,000,000
A Quiet Place Collection	1905.61%	\$323,952,971	\$17,000,000
Children of the Corn Collection	1721.12%	\$13,768,989	\$800,000
Escape Room - Collection	1630.13%	\$146,712,077	\$9,000,000
Happy Death Day Collection	1565.97%	\$88,139,709	\$6,900,000
屏住呼吸（系列）	1506.54%	\$149,147,649	\$9,900,000
Minions Collection	1466.82%	\$1,085,444,662	\$74,000,000
It Collection	1335.45%	\$635,875,764	\$57,000,000
Annabelle Collection	1307.14%	\$246,384,238	\$22,500,000
Insidious Collection	1304.35%	\$130,434,738	\$10,000,000
Despicable Me Collection	1193.50%	\$954,800,131	\$80,000,000
The Purge Collection	1020.08%	\$116,322,071	\$11,500,000
Call Me by Your Name Collection	947.22%	\$37,888,660	\$4,000,000
Deadpool Collection	932.53%	\$700,732,819	\$84,000,000
Saw Collection	929.53%	\$92,952,888	\$10,000,000

	Most franchise earn a lot on their investment. This is expected as there is a reason for film makers to visit same universe several times. More often than not it is because of their proven success record and popularity among movie consumers.
▼	7.3.1.1 which genre to franchise

```
In [152]: 1print('On an average films that are part of a franchise earn {:.2f}% return.'.
2         format(franchise_df.ROI_percentage.mean()))

On an average films that are part of a franchise earn 727.47% return.
```

```
In [153]: 1# joining and filtering using SQL statements
2list_of_franchise_df0 = sqldf("""SELECT
3     DISTINCT belongs_to_collection,
4     a.ROI AS 'ROI', b.world_collection,
5     b.genres
6     FROM franchise_df AS a
7     JOIN franchise_df AS b
8     USING(belongs_to_collection);""")
9
10format_dict = {'ROI': '${0:,.0f}', 'world_collection': '${0:,.0f}'}
11
12list_of_franchise_df0[~list_of_franchise_df0.belongs_to_collection.
13                      duplicated()].sort_values(
14                      by='ROI',
15                      ascending=False)[:15].style.background_gradient(
16                      cmap='bwr').hide_index().format(format_dict)
```

belongs_to_collection	ROI	world_collection	genres
Star Wars Collection	\$1,823,455,919	\$1,074,144,248	Action,Adventure,Fantasy
Jurassic Park Collection	\$1,521,713,208	\$1,310,464,680	Action,Adventure,Sci-Fi
The Lion King (Reboot) Collection	\$1,397,870,986	\$1,657,870,986	Adventure,Animation,Drama
The Fast and the Furious Collection	\$1,325,255,622	\$1,238,764,765	Action,Adventure,Crime
Frozen Collection	\$1,300,026,933	\$1,450,026,933	Adventure,Animation,Comedy
The Avengers Collection	\$1,155,403,694	\$1,405,403,694	Action,Adventure,Sci-Fi
Black Panther Collection	\$1,147,597,973	\$1,347,597,973	Action,Adventure,Sci-Fi
Minions Collection	\$1,085,444,662	\$1,159,444,662	Adventure,Animation,Comedy
The Incredibles Collection	\$1,043,089,244	\$1,243,089,244	Action,Adventure,Animation
Aquaman Collection	\$988,485,886	\$1,148,485,886	Action,Adventure,Fantasy
Captain Marvel Collection	\$976,462,972	\$1,128,462,972	Action,Adventure,Sci-Fi
Despicable Me Collection	\$954,800,131	\$1,034,800,131	Adventure,Animation,Comedy
Captain America Collection	\$903,561,649	\$1,153,561,649	Action,Adventure,Sci-Fi
Toy Story Collection	\$898,394,593	\$1,073,394,593	Adventure,Animation,Comedy
Jumanji Collection	\$872,174,450	\$800,059,707	Action,Adventure,Comedy

Observation: None of them fall into a single genre.


```
In [154]: 1# joining and filtering using SQL statements
2list_of_franchise_df = sqldf(
3     """SELECT
4         belongs_to_collection,
5         a.ROI_percentage AS 'ROI%',
6         b.genres
7     FROM franchise_df AS a
8     JOIN franchise_df AS b
9     USING(belongs_to_collection);""")
```

```
In [155]: 1# most often produced genre for serialization of movies
2list_of_franchise_df.loc[:, 'genres_exp'] = list_of_franchise_df.genres.map(
3     lambda x: x.split(','))
4
5franchise_genre = list_of_franchise_df.explode('genres_exp').groupby(
6     'genres_exp').agg(['count', 'mean']).sort_values(by=('ROI%', 'count'),
7     ascending=False)
8franchise_genre.columns = [
9     " ".join(pair) for pair in franchise_genre.columns
10]
11franchise_genre=franchise_genre.reset_index()
12franchise_genre.style.background_gradient(cmap='PRGn')
```

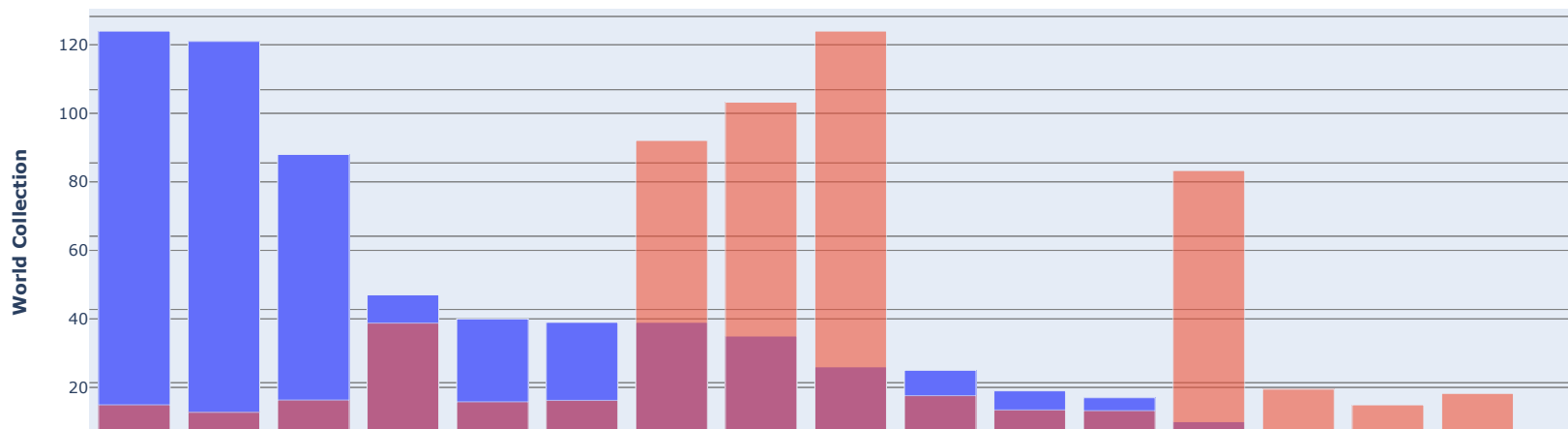
	genres_exp	ROI% count	ROI% mean
0	Adventure	124	349.019417
1	Action	121	297.557482
2	Comedy	88	381.703050
3	Drama	47	906.836685
4	Animation	40	369.577146
5	Sci-Fi	39	379.531594
6	Horror	39	2153.693791
7	Thriller	35	2414.741535
8	Mystery	26	2900.197396
9	Crime	25	411.506582
10	Fantasy	19	314.229739
11	Family	17	309.154043
12	Romance	10	1948.204978
13	Music	4	456.125992
14	Sport	2	348.766674
15	Musical	1	426.755804
16	Biography	1	100.663038

```

In [156]: 1## from https://plotly.com/python/multiple-axes/ ##official plotly how to instructions
2fig = make_subplots(specs=[[{"secondary_y": True}]]))
3# Add traces
4fig.add_trace(
5     go.Bar(x=franchise_genre['genres_exp'],
6             y=franchise_genre['ROI% count'],
7             name="Movies released",
8             offset=True),
9     secondary_y=False,
10)
11fig.add_trace(
12     go.Bar(x=franchise_genre['genres_exp'],
13             y=franchise_genre['ROI% mean'],
14             name="ROI% mean",
15             offset=True,
16             opacity=.6),
17     secondary_y=True,
18)
19# Add figure title
20fig.update_layout(title_text="Most often produced genre for serialized movies ")
21# Set x-axis title
22fig.update_xaxes(title_text="Genre")
23# Set y-axes titles
24fig.update_yaxes(title_text="<b>World Collection</b>", secondary_y=False)
25fig.update_yaxes(title_text="<b>Number of Movie Released</b>",
26                 secondary_y=True)
27fig.show()

```

Most often produced genre for serialized movies



Adventure, Action, Comedy market is saturated. Horror, Thriller, Mystery release count is lower with higher mean return percentage. This recommendation will alter if we look at collection instead of ROI% because those genre requires less budget, so the return percentage is generally higher.

7.3.2 non franchise info

```
In [157]: 1print(  
2     'On an average films that are not part of a franchise earn {:.2f}% return.'  
3     .format(non_franchise_df.ROI_percentage.mean()))
```

On an average films that are not part of a franchise earn 186.17% return.

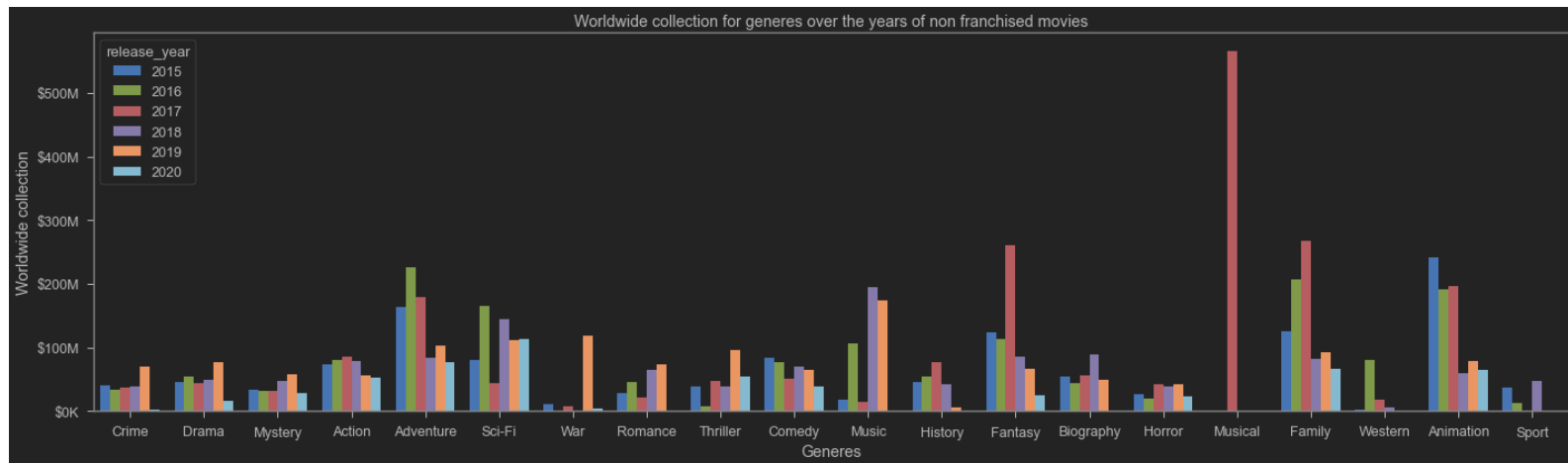
```
In [158]: 1non_franchise_df.loc[:, 'genres_exp'] = non_franchise_df.genres.map(  
2     lambda x: x.split(','))
```

```
In [159]: 1non_franchise_df = non_franchise_df.explode('genres_exp')
```

```

In [160]:
1# Worldwide collection for genres over the years of non franchised movies
2fig, ax = plt.subplots()
3fig.set_size_inches(20, 6)
4# plotting
5sns.barplot(data=non_franchise_df,
6            x='genres_exp',
7            y='world_collection',
8            hue='release_year',ci=None).yaxis.set_major_formatter(format_number)
9plt.title('Worldwide collection for genres over the years of non franchised movies')
10plt.ylabel('Worldwide collection')
11plt.xlabel("Genres")
12plt.tight_layout()
13
14plt.show()

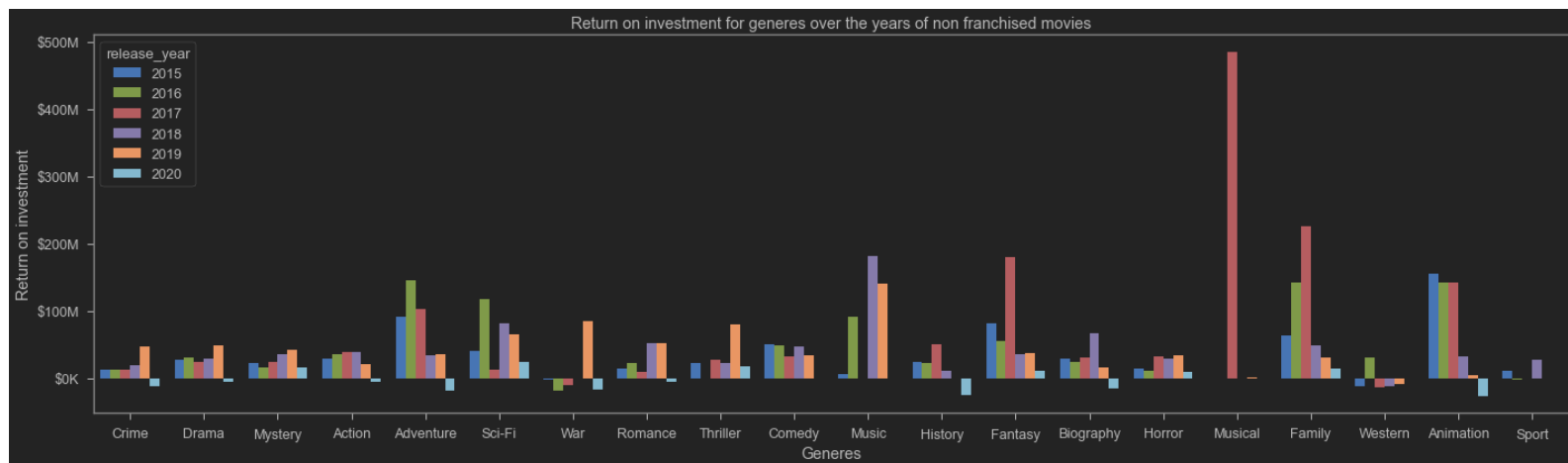
```



```

In [161]:
1# Return on investment for genres over the years of non franchised movies
2fig, ax = plt.subplots()
3fig.set_size_inches(20, 6)
4# plotting
5sns.barplot(data=non_franchise_df,
6            x='genres_exp',
7            y='ROI',
8            hue='release_year',ci=None).yaxis.set_major_formatter(format_number)
9plt.title('Return on investment for genres over the years of non franchised movies')
10plt.ylabel('Return on investment')
11plt.xlabel("Genres")
12plt.tight_layout()
13
14plt.show()

```



Non franchised movies are experiencing a hard time in the box office. The general trend is downwards across the board except Crime and Drama and Mystery. Mystery, Sci-Fi and Horror did well in 2020. Those three genres have high correlation.

7.3.3 Side by side comparison

Converting franchise info in to a boolean array

```

In [162]:
1franchise_df_main.loc[~main_df['belongs_to_collection.name'].isna(),
2                    'franchise'] = True

```

```

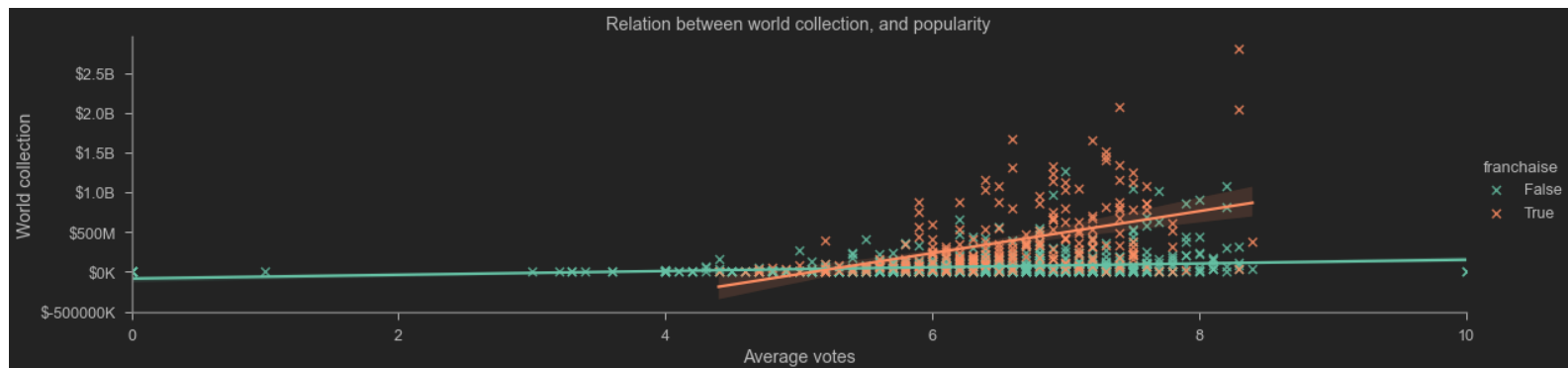
In [163]:
1franchise_df_main.loc[main_df['belongs_to_collection.name'].isna(),
2                    'franchise'] = False

```

```

In [164]:
1 g = sns.lmplot(data=franchise_df_main,
2               x='vote_average',
3               y='world_collection',
4               hue='franchise',
5               height=4,
6               aspect=4,
7               palette='Set2',
8               markers='x')
9 for ax in g.axes.flat:
10     ax.yaxis.set_major_formatter(format_number)
11
12 plt.title('Relation between world collection, and popularity')
13 plt.ylabel('World collection')
14 plt.xlabel("Average votes")
15
16 g = sns.lmplot(data=franchise_df_main,
17               x='dom_collection',
18               y='int_collection',
19               hue='franchise',
20               height=4,
21               aspect=4,
22               palette='Set2',
23               markers='x')
24 for ax in g.axes.flat:
25     ax.yaxis.set_major_formatter(format_number)
26     ax.xaxis.set_major_formatter(format_number)
27 plt.title(
28     'Relation between world collection, and budget. (Cost to sales ratio)')
29 plt.ylabel('World collection')
30 plt.xlabel("Budget")
31
32 plt.show()

```



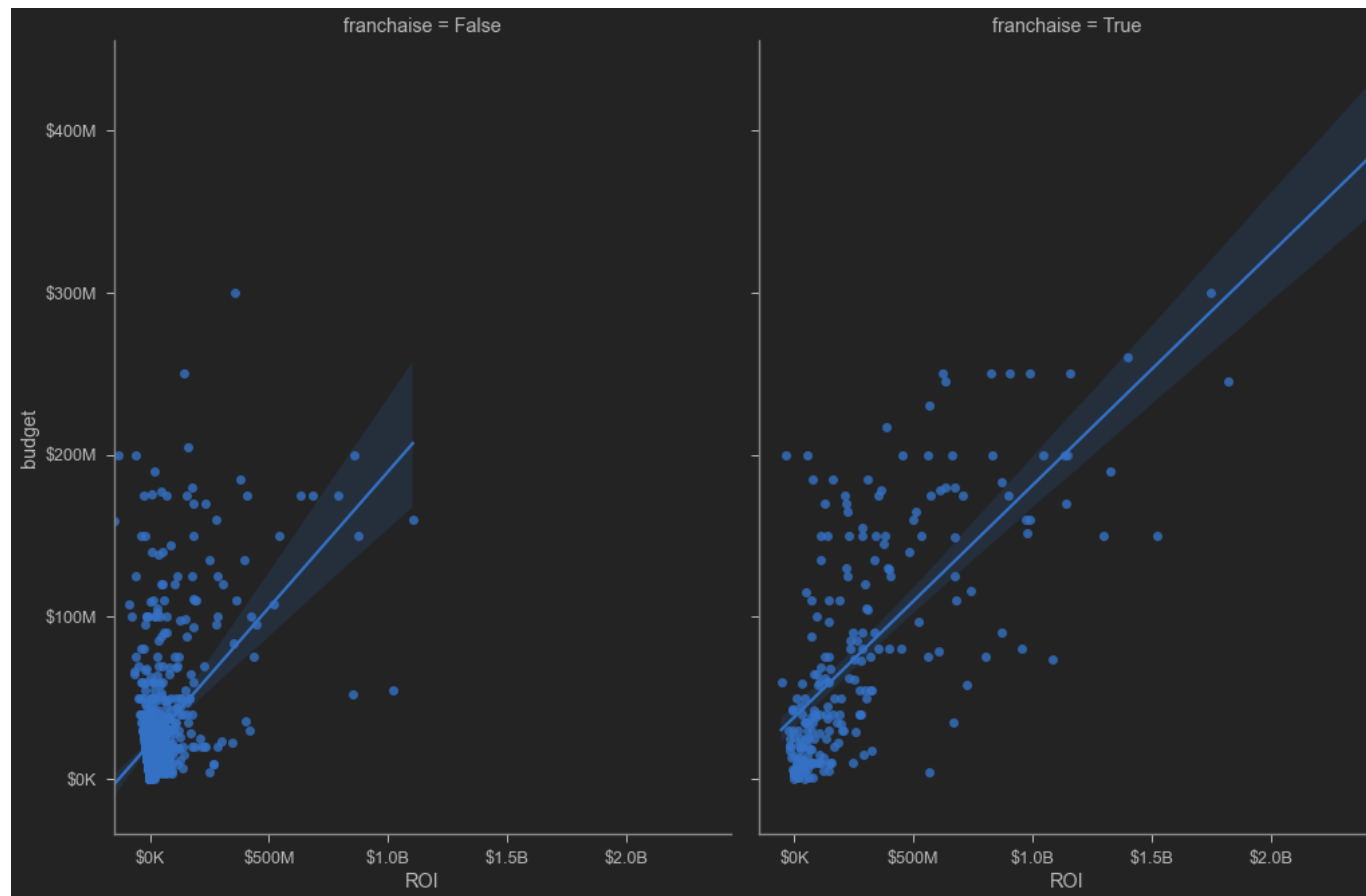


Franchised movies are often more popular with greater success in international market.

```

In [165]:
1g = sns.FacetGrid(
2    franchise_df_main, col='franchise',
3    height=10, aspect=.7, palette='Set2')
4g.map(sns.regplot, 'ROI', 'budget')
5for ax in g.axes.flat:
6    ax.yaxis.set_major_formatter(format_number)
7    ax.xaxis.set_major_formatter(format_number)
8g.fig.subplots_adjust(top=0.9)
9# g.fig.suptitle('Relationship between budget and ROI based on release timing')

```



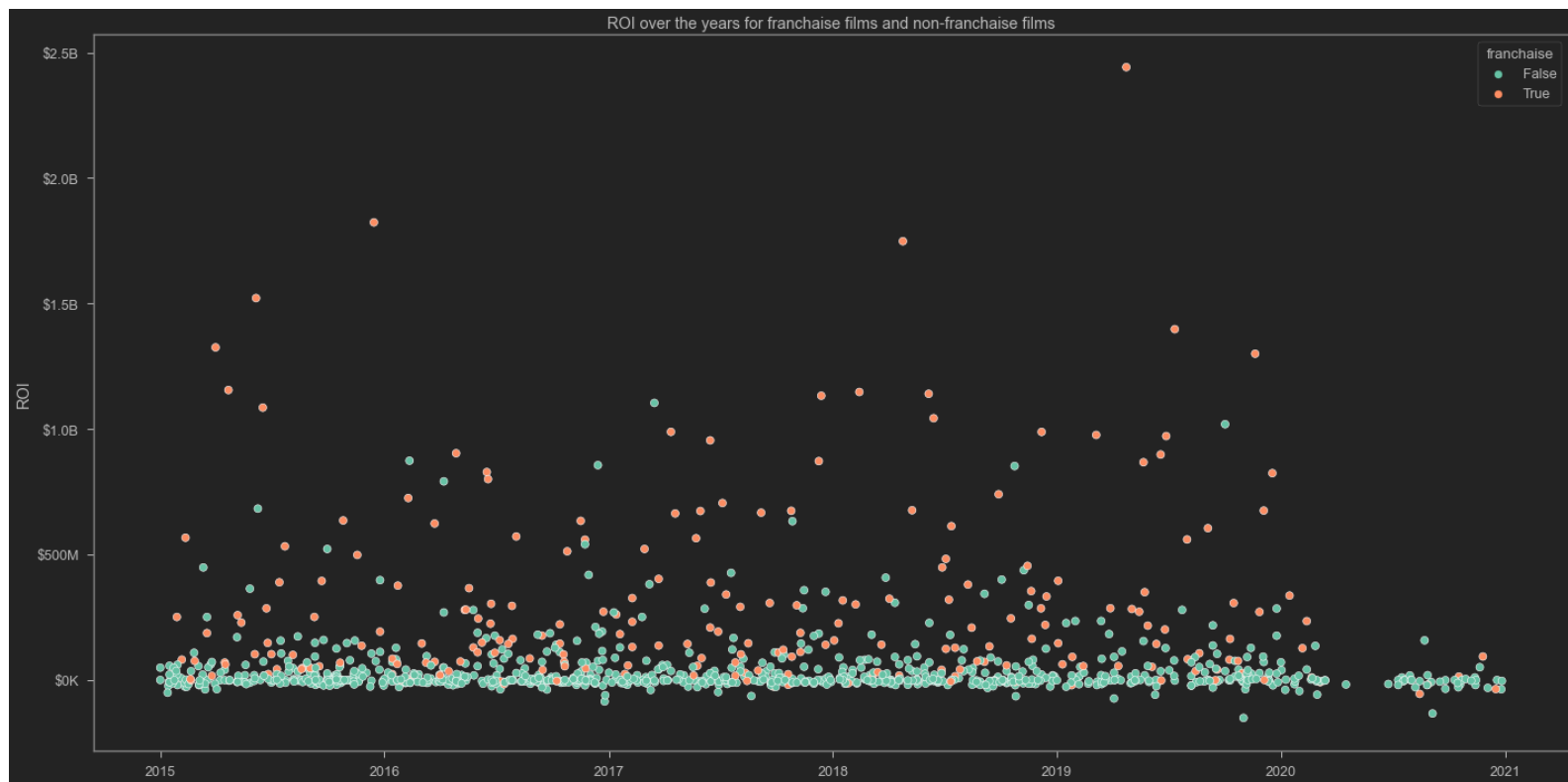


	Franchised movies require bigger budget bit their return is also significantly higher.
--	--

```

In [166]:
1# ROI over the years for franchise films and non-franchise films
2plt.figure(figsize=(20, 10))
3sns.scatterplot(x='release_date',
4               y='ROI',
5               hue='franchise',
6               data=franchise_df_main,
7               palette='Set2',
8               legend='brief').yaxis.set_major_formatter(format_number)
9plt.title('ROI over the years for franchise films and non-franchise films')
10plt.ylabel('ROI')
11plt.xlabel("")
12plt.tight_layout()

```



This straightforward time series of box office gross profit of the two categories is the simplest but most layman friendly chart that demonstrate the stark difference between them. Franchised movies are consistently outperforming the other category.

7.3.4 Action suggestion

All the analysis leads towards starting a movie franchise in a shared movie universe. This must be be priority when selecting genre, director and other crew and cast. There must be option for serialization in the future. And for this Horror, Thriller, Mystery and Adventure, Action, Comedy genre should be prioritized. It very rare that a movie falls in only one genre this days.

8 Conclusion

	<p>Lets summarize and reiterate:</p> <hr/> <ol style="list-style-type: none">1. My recommendation is to focus for release schedule in the <i>holidays season</i>. There is higher probability of financial and critical success for movies released in that time frame. It is relatively cheaper to make than the next best option; i.e., Summer.2. Any one or combination of <i>Action, Adventure, Animation</i> is recommended. Animation and Action has 35% chance for occurring as genre combo. There is no landslide winner here, although this graphs can be used to figure out which one to avoid, for example western and war.3. All the analysis leads towards starting a <i>franchise</i> in a shared movie universe. This must be be priority when selecting genre, director and other crew and cast. There must be option for serialization in the future. And for this Horror, Thriller, Mystery or Adventure, Action, Comedy genre combination should be prioritized. <p>It very rare that a movie falls in only one genre this days.</p>	
▼	<h2>9 Next Steps</h2> <p>Further analyses could yield additional insights to further improve considerations for creating a new movie:</p> <hr/> <ul style="list-style-type: none">• Performance of other language movies and markets.• Focusing on low budget movies versus high budget movies performance and rational.• Movies performance in home and international market.• Recommending lead director.• Recommending movie cast classified on genre.• Focus only on 2020 data and find pattern and trend.	
▼	<h2>10 For More Information</h2> <p>See the full analysis in the Jupyter Notebook (./student.ipynb) or review this presentation (./presentation.pdf).</p>	
▼	<h2>11 Appendix</h2>	
▼	<h3>11.1 Most produced genre combo</h3>	
In [167]:	<pre>1combo_genre = main_df_raw.iloc[:,18:-1].copy()</pre>	

In [168]:

```
1 combo_genre = combo_genre.corr()  
2 combo_genre.style.background_gradient(cmap='PuRd')
```

	Action	Adventure	Biography	Drama	Fantasy	Comedy	War	Crime	Romance	Family	History	Sci-Fi	Thriller	Western	Sport	Mystery	Horror
Action	1.000000	0.234557	-0.039696	-0.155409	0.021305	-0.134239	0.007438	0.184287	-0.124091	-0.082017	0.012833	0.119460	0.087795	-0.003914	-0.014039	-0.065094	-0.052483
Adventure	0.234557	1.000000	-0.032004	-0.191267	0.096487	-0.008579	-0.034706	-0.047212	-0.109738	0.157636	-0.041994	0.055559	-0.091786	0.000736	-0.030389	-0.057172	-0.071896
Biography	-0.039696	-0.032004	1.000000	0.132731	-0.040755	-0.121452	0.024546	0.013829	-0.048016	-0.034802	0.203732	-0.038157	-0.069462	-0.008145	0.076512	-0.053194	-0.065236
Drama	-0.155409	-0.191267	0.132731	1.000000	-0.066611	-0.290930	0.063189	0.023217	0.038453	-0.103478	0.117032	-0.076056	-0.111885	0.019962	0.037100	-0.009435	-0.184330
Fantasy	0.021305	0.096487	-0.040755	-0.066611	1.000000	-0.016738	-0.020963	-0.058588	-0.019053	0.072728	-0.040150	0.008284	-0.063046	-0.007915	-0.023789	0.013609	0.054360
Comedy	-0.134239	-0.008579	-0.121452	-0.290930	-0.016738	1.000000	-0.080062	-0.109071	0.107414	0.008391	-0.114813	-0.090470	-0.248184	-0.034028	-0.027722	-0.148315	-0.150491
War	0.007438	-0.034706	0.024546	0.063189	-0.020963	-0.080062	1.000000	-0.039289	-0.019791	-0.023206	0.118629	-0.017206	-0.025188	0.002384	-0.014776	-0.022964	-0.030458
Crime	0.184287	-0.047212	0.013829	0.023217	-0.058588	-0.109071	-0.039289	1.000000	-0.114093	-0.073954	-0.046635	-0.048057	0.136428	0.006373	-0.037173	0.082603	-0.059221
Romance	-0.124091	-0.109738	-0.048016	0.038453	-0.019053	0.107414	-0.019791	-0.114093	1.000000	-0.063814	-0.031822	-0.042308	-0.115350	-0.017523	-0.024942	-0.068354	-0.105189
Family	-0.082017	0.157636	-0.034802	-0.103478	0.072728	0.008391	-0.023206	-0.073954	-0.063814	1.000000	-0.031871	-0.023233	-0.094161	-0.010673	0.010173	-0.058735	-0.076032
History	0.012833	-0.041994	0.203732	0.117032	-0.040150	-0.114813	0.118629	-0.046635	-0.031822	-0.031871	1.000000	-0.028010	-0.053683	-0.005363	0.006565	-0.036475	0.075418
Sci-Fi	0.119460	0.055559	-0.038157	-0.076056	0.008284	-0.090470	-0.017206	-0.048057	-0.042308	-0.023233	-0.028010	1.000000	0.043575	-0.010579	-0.017144	0.040065	-0.029989
Thriller	0.087795	-0.091786	-0.069462	-0.111885	-0.063046	-0.248184	-0.025188	0.136428	-0.115350	-0.094161	-0.053683	0.043575	1.000000	0.011551	-0.043187	0.176246	0.199766
Western	-0.003914	0.000736	-0.008145	0.019962	-0.007915	-0.034028	0.002384	0.006373	-0.017523	-0.010673	-0.005363	-0.010579	0.011551	1.000000	0.002428	-0.015553	-0.005981
Sport	-0.014039	-0.030389	0.076512	0.037100	-0.023789	-0.027722	-0.014776	-0.037173	-0.024942	0.010173	0.006565	-0.017144	-0.043187	0.002428	1.000000	-0.030591	-0.017649
Mystery	-0.065094	-0.057172	-0.053194	-0.009435	0.013609	-0.148315	-0.022964	0.082603	-0.068354	-0.058735	-0.036475	0.040065	0.176246	-0.015553	-0.030591	1.000000	-0.033203
Horror	-0.052483	-0.071896	-0.065236	-0.184330	0.054360	-0.150491	-0.030458	-0.059221	-0.105189	-0.076032	-0.056715	0.075418	0.199766	-0.005981	-0.036922	0.242051	1.000000
Music	-0.067461	-0.041102	0.054277	0.025775	-0.005686	-0.007963	-0.010543	-0.052291	0.022317	-0.004684	-0.014798	-0.029989	-0.057237	-0.003808	-0.017649	-0.033203	-0.000000
Animation	0.060473	0.353164	-0.036639	-0.203516	0.028996	-0.012312	-0.027168	-0.071642	-0.087790	0.140764	-0.039954	0.014958	-0.090670	-0.015114	-0.005959	-0.057700	-0.000000
Musical	-0.033499	-0.017786	-0.014353	-0.028121	0.000515	-0.007671	-0.005422	-0.030310	0.015053	0.010627	-0.001970	-0.016579	-0.032767	0.006416	-0.005377	-0.021159	-0.000000

```
In [169]: 1correlation_top_bottom(combo_genre)
```

Positive correlations:

	index	feature_combo	correlation
0	359	Music and Musical	0.552813
1	38	Animation and Adventure	0.353164
2	316	Mystery and Horror	0.242051
3	1	Action and Adventure	0.234557
4	50	History and Biography	0.203732
5	256	Thriller and Horror	0.199766
6	7	Crime and Action	0.184287
7	255	Mystery and Thriller	0.176246
8	29	Family and Adventure	0.157636
9	198	Animation and Family	0.140764

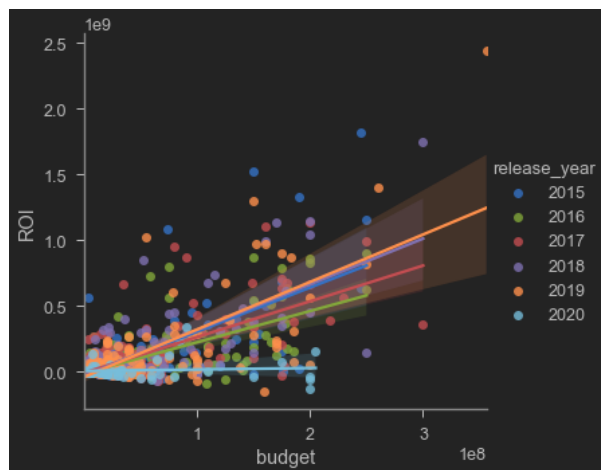
Negative correlations:

	index	feature_combo	correlation
0	65	Comedy and Drama	-0.290930
1	112	Comedy and Thriller	-0.248184
2	78	Animation and Drama	-0.203516
3	23	Drama and Adventure	-0.191267
4	76	Drama and Horror	-0.184330
5	3	Action and Drama	-0.155409
6	116	Comedy and Horror	-0.150491
7	115	Mystery and Comedy	-0.148315
8	5	Comedy and Action	-0.134239
9	8	Romance and Action	-0.124091

▼ 11.2 Variability of profitability on different metrics

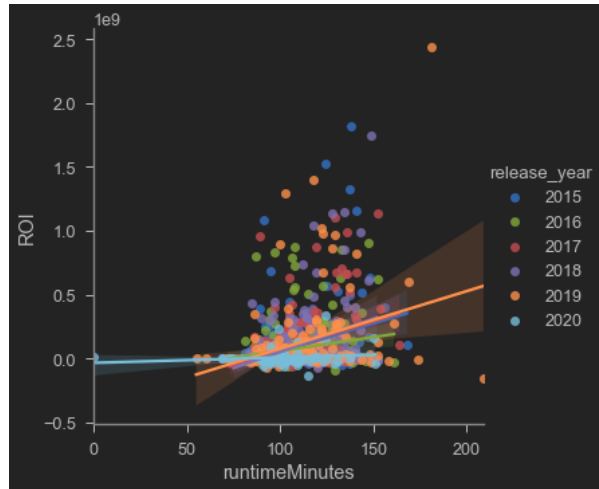
▼ 11.2.1 budget vs profitability

```
In [170]: 1sns.lmplot(data=main_df, x='budget', y='ROI',hue='release_year')
2plt.show()
```



▼ 11.2.2 runtime on profitability

```
In [171]: 1sns.lmplot(data=main_df, x='runtimeMinutes', y='ROI',hue='release_year')
          2plt.show()
```



11.2.3 user rating on profitability

```
In [172]: 1sns.lmplot(data=main_df, x='vote_average', y='ROI',hue='release_year')
          2plt.show()
```

