

Overview

This project is a response to the growing demand for Microsoft to venture into the movie making business. This move will enable Microsoft to tap into a new market while relaying a new service to their consumers. The success of Microsoft's new movie studio is dependent on making well-informed decisions that are data driven. The data driven decisions should propel the company forward with more sales and better customer engagement.

Business problem

Microsoft aims to establish a presence in the original video content industry, competing with well-established players like Netflix and Amazon. The primary challenge is to develop a winning strategy for creating content that not only competes effectively but also captures and retains audiences. To achieve this, Microsoft must set itself apart by making substantial investments in content development, talent acquisition, and marketing. Additionally, it needs to gain a deep understanding of audience preferences and emerging trends. Balancing the expenses of content creation with revenue sources, such as advertising or subscription models, is also critical.

The analysis is structured around three key factors:

Identifying the most prominent and appealing genres to produce content that resonates with viewers.

Investigating the correlation between the duration of a movie and its popularity, helping Microsoft make informed decisions about the length of its content.

Recognizing the best-performing studios in the movie box office, which can provide valuable insights into potential partnerships, acquisitions, or collaboration opportunities.

By addressing these aspects, Microsoft can enhance its ability to compete in the video industry, understand audience preferences, and establish a viable revenue model.

Data preparation

```
In [284]: #Loading needed libraries
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import matplotlib

Loading first data set(Box Office Mojo)

In [285]: movie_gross = pd.read_csv(r"C:\Users\ALLAN\Downloads\box_movie_gross.csv.gz")
print(movie_gross.head())

# movie_gross
#   tcconst      primary_title  original_title  start_year  runtime_minutes  genres  averagerating  numvotes
#   -----
#   0  tt0066767  One Day Before the Rainy Season  Ashad Ka Ek Din  2019      114.00000  Biography,Drama  7.2  43
#   1  tt0069049  The Other Side of the Wind  The Other Side of the Wind  2018      122.00000  Drama  6.9  4517
#   2  tt0069204  The Deathly Hallows Part 1  Saise Bada Sukh  2018      94.65404  Comedy,Drama  6.1  15
#   3  tt0069204  The Deathly Hallows Part 1  Saise Bada Sukh  2018      94.65404  Comedy,Drama  6.1  15
#   4  tt0186275  The Wandering Soap Opera  La Telenovela Errante  2017      80.0  Comedy,Drama,Fantasy  6.5  119

# movie_gross
#   foreign_gross  year
#   -----
#   0  652000000  2010
#   1  691300000  2010
#   2  664300000  2010
#   3  535700000  2010
#   4  513900000  2010

Loading second data set(Movie ratings)

In [286]: in_movie_ratings = pd.read_csv(r"C:\Users\ALLAN\Downloads\inmb_movie_ratings.csv.gz")
print(in_movie_ratings.head())

# in_movie_ratings
#   tcconst      primary_title  original_title  start_year  runtime_minutes  genres  averagerating  numvotes
#   -----
#   0  tt0066767  One Day Before the Rainy Season  Ashad Ka Ek Din  2019      114.00000  Biography,Drama  7.2  43
#   1  tt0069049  The Other Side of the Wind  The Other Side of the Wind  2018      122.00000  Drama  6.9  4517
#   2  tt0069204  The Deathly Hallows Part 1  Saise Bada Sukh  2018      94.65404  Comedy,Drama  6.1  15
#   3  tt0069204  The Deathly Hallows Part 1  Saise Bada Sukh  2018      94.65404  Comedy,Drama  6.1  15
#   4  tt0186275  The Wandering Soap Opera  La Telenovela Errante  2017      80.0  Comedy,Drama,Fantasy  6.5  119

# in_movie_ratings
#   foreign_gross  year
#   -----
#   0  652000000  2010
#   1  691300000  2010
#   2  664300000  2010
#   3  535700000  2010
#   4  513900000  2010

Loading third data set(Movie basics)

In [287]: in_movie_basics = pd.read_csv(r"C:\Users\ALLAN\Downloads\inmb_title_basics.csv.gz")
print(in_movie_basics.head())

# in_movie_basics
#   tcconst      primary_title  original_title  start_year  runtime_minutes  genres  averagerating  numvotes
#   -----
#   0  tt0066767  One Day Before the Rainy Season  Ashad Ka Ek Din  2019      114.00000  Biography,Drama  7.2  43
#   1  tt0069049  The Other Side of the Wind  The Other Side of the Wind  2018      122.00000  Drama  6.9  4517
#   2  tt0069204  The Deathly Hallows Part 1  Saise Bada Sukh  2018      94.65404  Comedy,Drama  6.1  15
#   3  tt0069204  The Deathly Hallows Part 1  Saise Bada Sukh  2018      94.65404  Comedy,Drama  6.1  15
#   4  tt0186275  The Wandering Soap Opera  La Telenovela Errante  2017      80.0  Comedy,Drama,Fantasy  6.5  119

# in_movie_basics
#   start_year  runtime_minutes  genres
#   -----
#   0  2019      114.0  Action,Crime,Drama
#   1  2019      114.0  Biography,Drama
#   2  2018      122.0  Drama
#   3  2018      NaN  Comedy,Drama
#   4  2017      80.0  Comedy,Drama,Fantasy

In [288]: # merge the two datasets (movie_ratings and movie_basics)
merged_data = pd.merge(in_movie_basics, in_movie_ratings, on='tcconst')
# The two datasets have been merged on a common column tcconst.
merged_dataset

Out[289]:
# merged_dataset
#   tcconst      primary_title  original_title  start_year  runtime_minutes  genres  averagerating  numvotes
#   -----
#   0  tt0066767  One Day Before the Rainy Season  Ashad Ka Ek Din  2019      114.00000  Biography,Drama  7.2  43
#   1  tt0069049  The Other Side of the Wind  The Other Side of the Wind  2018      122.00000  Drama  6.9  4517
#   2  tt0069204  The Deathly Hallows Part 1  Saise Bada Sukh  2018      94.65404  Comedy,Drama  6.1  15
#   3  tt0069204  The Deathly Hallows Part 1  Saise Bada Sukh  2018      94.65404  Comedy,Drama  6.1  15
#   4  tt0186275  The Wandering Soap Opera  La Telenovela Errante  2017      80.0  Comedy,Drama,Fantasy  6.5  119

# merged_dataset
#   foreign_gross  year
#   -----
#   0  652000000  2010
#   1  691300000  2010
#   2  664300000  2010
#   3  535700000  2010
#   4  513900000  2010

73856 rows x 8 columns
```

Checking the data and datatypes

```
In [289]: # Inspect a few random rows of the DataFrame to get a sense of the data's structure and content
merged_dataset.sample(5)

Out[289]:
# merged_dataset.sample(5)
#   tcconst      primary_title  original_title  start_year  runtime_minutes  genres  averagerating  numvotes
#   -----
#   0  tt0066767  One Day Before the Rainy Season  Ashad Ka Ek Din  2019      114.00000  Biography,Drama  7.2  43
#   1  tt0069049  The Other Side of the Wind  The Other Side of the Wind  2018      122.00000  Drama  6.9  4517
#   2  tt0069204  The Deathly Hallows Part 1  Saise Bada Sukh  2018      94.65404  Comedy,Drama  6.1  15
#   3  tt0069204  The Deathly Hallows Part 1  Saise Bada Sukh  2018      94.65404  Comedy,Drama  6.1  15
#   4  tt0186275  The Wandering Soap Opera  La Telenovela Errante  2017      80.0  Comedy,Drama,Fantasy  6.5  119

# merged_dataset.sample(5)
#   foreign_gross  year
#   -----
#   0  652000000  2010
#   1  691300000  2010
#   2  664300000  2010
#   3  535700000  2010
#   4  513900000  2010

73856 rows x 8 columns

In [290]: merged_dataset.shape # checks for the number of rows and columns
(73856, 8)

In [291]: merged_dataset.shape # checks for the number of rows and columns
(73856, 8)

In [292]: movie_gross.shape # checks for the number of rows and columns
(3387, 5)

In [293]: merged_dataset.columns # prints out the column names
Out[293]:
Index(['tcconst', 'primary_title', 'original_title', 'start_year', 'runtime_minutes', 'genres', 'averagerating', 'numvotes'], dtype='object')

In [294]: movie_gross.columns # prints out the column names
Out[294]:
Index(['title', 'studio', 'domestic_gross', 'foreign_gross', 'year'], dtype='object')

In [295]: merged_dataset.info() # checks for the overview of the data
Out[295]:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 73856 entries, 0 to 73855
Data columns (total 8 columns):
#   Column      Non-Null Count  Dtype
---  ---
0  tcconst      73856 non-null  object
1  primary_title  73856 non-null  object
2  original_title  73856 non-null  object
3  start_year     73856 non-null  int64
4  runtime_minutes  66236 non-null  float64
5  genres        73856 non-null  object
6  averagerating  73856 non-null  float64
7  numvotes      73856 non-null  int64
dtypes: float64(2), int64(2), object(4)
memory usage: 4.5+ MB

In [296]: movie_gross.info() # checks for the overview of the data
Out[296]:
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3387 entries, 0 to 3386
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  ---
0  title        3387 non-null  object
1  studio       3387 non-null  object
2  domestic_gross  3387 non-null  float64
3  foreign_gross  3387 non-null  object
4  year         3387 non-null  int64
dtypes: float64(1), int64(1), object(3)
memory usage: 132.4+ KB

put description here

In [297]: merged_dataset.dtypes # dtype attribute
Out[297]:
tcconst      object
primary_title  object
original_title  object
start_year     int64
runtime_minutes  float64
genres         object
averagerating  float64
numvotes       int64
dtype: object

In [298]: movie_gross.dtypes # dtype attributes
Out[298]:
title        object
studio       object
domestic_gross  float64
foreign_gross  object
year         int64
dtype: object

In [299]: movie_gross.describe() # check for a statistical summary of the data
Out[299]:
# movie_gross.describe()
#   domestic_gross  year
#   -----
#   count  3387000000  3387000000
#   min  6.68265e+07  2010.000000
#   25%  1.200000e+05  2012.000000
#   50%  1.400000e+06  2014.000000
#   75%  2.790000e+07  2016.000000
#   max  9.367000e+08  2018.000000
```

Data preprocessing and cleaning

We now need to identify and correct or remove incorrect, incomplete, incorrectly formatted, corrupted, duplicate or irrelevant data within the provided data.

```
In [300]: # Display the count of null values in each column
null_values = movie_gross.isnull()
print(null_values.sum())

# movie_gross
#   title  studio  domestic_gross  foreign_gross  year
#   -----
#   0  0  5  28  1350
#   1  0  5  28  1350
#   2  0  5  28  1350
#   3  0  5  28  1350
#   4  0  5  28  1350
#   5  0  5  28  1350
#   6  0  5  28  1350
#   7  0  5  28  1350
#   8  0  5  28  1350
#   9  0  5  28  1350
#   10 0  5  28  1350
#   11 0  5  28  1350
#   12 0  5  28  1350
#   13 0  5  28  1350
#   14 0  5  28  1350
#   15 0  5  28  1350
#   16 0  5  28  1350
#   17 0  5  28  1350
#   18 0  5  28  1350
#   19 0  5  28  1350
#   20 0  5  28  1350
#   21 0  5  28  1350
#   22 0  5  28  1350
#   23 0  5  28  1350
#   24 0  5  28  1350
#   25 0  5  28  1350
#   26 0  5  28  1350
#   27 0  5  28  1350
#   28 0  5  28  1350
#   29 0  5  28  1350
#   30 0  5  28  1350
#   31 0  5  28  1350
#   32 0  5  28  1350
#   33 0  5  28  1350
#   34 0  5  28  1350
#   35 0  5  28  1350
#   36 0  5  28  1350
#   37 0  5  28  1350
#   38 0  5  28  1350
#   39 0  5  28  1350
#   40 0  5  28  1350
#   41 0  5  28  1350
#   42 0  5  28  1350
#   43 0  5  28  1350
#   44 0  5  28  1350
#   45 0  5  28  1350
#   46 0  5  28  1350
#   47 0  5  28  1350
#   48 0  5  28  1350
#   49 0  5  28  1350
#   50 0  5  28  1350
#   51 0  5  28  1350
#   52 0  5  28  1350
#   53 0  5  28  1350
#   54 0  5  28  1350
#   55 0  5  28  1350
#   56 0  5  28  1350
#   57 0  5  28  1350
#   58 0  5  28  1350
#   59 0  5  28  1350
#   60 0  5  28  1350
#   61 0  5  28  1350
#   62 0  5  28  1350
#   63 0  5  28  1350
#   64 0  5  28  1350
#   65 0  5  28  1350
#   66 0  5  28  1350
#   67 0  5  28  1350
#   68 0  5  28  1350
#   69 0  5  28  1350
#   70 0  5  28  1350
#   71 0  5  28  1350
#   72 0  5  28  1350
#   73 0  5  28  1350
#   74 0  5  28  1350
#   75 0  5  28  1350
#   76 0  5  28  1350
#   77 0  5  28  1350
#   78 0  5  28  1350
#   79 0  5  28  1350
#   80 0  5  28  1350
#   81 0  5  28  1350
#   82 0  5  28  1350
#   83 0  5  28  1350
#   84 0  5  28  1350
#   85 0  5  28  1350
#   86 0  5  28  1350
#   87 0  5  28  1350
#   88 0  5  28  1350
#   89 0  5  28  1350
#   90 0  5  28  1350
#   91 0  5  28  1350
#   92 0  5  28  1350
#   93 0  5  28  1350
#   94 0  5  28  1350
#   95 0  5  28  1350
#   96 0  5  28  1350
#   97 0  5  28  1350
#   98 0  5  28  1350
#   99 0  5  28  1350
#   100 0  5  28  1350
#   101 0  5  28  1350
#   102 0  5  28  1350
#   103 0  5  28  1350
#   104 0  5  28  1350
#   105 0  5  28  1350
#   106 0  5  28  1350
#   107 0  5  28  1350
#   108 0  5  28  1350
#   109 0  5  28  1350
#   110 0  5  28  1350
#   111 0  5  28  1350
#   112 0  5  28  1350
#   113 0  5  28  1350
#   114 0  5  28  1350
#   115 0  5  28  1350
#   116 0  5  28  1350
#   117 0  5  28  1350
#   118 0  5  28  1350
#   119 0  5  28  1350
#   120 0  5  28  1350
#   121 0  5  28  1350
#   122 0  5  28  1350
#   123 0  5  28  1350
#   124 0  5  28  1350
#   125 0  5  28  1350
#   126 0  5  28  1350
#   127 0  5  28  1350
#   128 0  5  28  1350
#   129 0  5  28  1350
#   130 0  5  28  1350
#   131 0  5  28  1350
#   132 0  5  28  1350
#   133 0  5  28  1350
#   134 0  5  28  1350
#   135 0  5  28  1350
#   136 0  5  28  1350
#   137 0  5  28  1350
#   138 0  5  28  1350
#   139 0  5  28  1350
#   140 0  5  28  1350
#   141 0  5  28  1350
#   142 0  5  28  1350
#   143 0  5  28  1350
#   144 0  5  28  1350
#   145 0  5  28  1350
#   146 0  5  28  1350
#   147 0  5  28  1350
#   148 0  5  28  1350
#   149 0  5  28  1350
#   150 0  5  28  1350
#   151 0  5  28  1350
#   152 0  5  28  1350
#   153 0  5  28  1350
#   154 0  5  28  1350
#   155 0  5  28  1350
#   156 0  5  28  1350
#   157 0  5  28  1350
#   158 0  5  28  1350
#   159 0  5  28  1350
#   160 0  5  28  1350
#   161 0  5  28  1350
#   162 0  5  28  1350
#   163 0  5  28  1350
#   164 0  5  28  1350
#   165 0  5  28  1350
#   166 0  5  28  1350
#   167 0  5  28  1350
#   168 0  5  28  1350
#   169 0  5  28  1350
#   170 0  5  28  1350
#   171 0  5  28  1350
#   172 0  5  28  1350
#   173 0  5  28  1350
#   174 0  5  28  1350
#   175 0  5  28  1350
#   176 0  5  28  1350
#   177 0  5  28  1350
#   178 0  5  28  1350
#   179 0  5  28  1350
#   180 0  5  28  1350
#   181 0  5  28  1350
#   182 0  5  28  1350
#   183 0  5  28  1350
#   184 0  5  28  1350
#   185 0  5  28  1350
#   186 0  5  28  1350
#   187 0  5  28  1350
#   188 0  5  28  1350
#   189 0  5  28  1350
#   190 0  5  28  1350
#   191 0  5  28  1350
#   192 0  5  28  1350
#   193 0  5  28  1350
#   194 0  5  28  1350
#   195 0  5  28  1350
#   196 0  5  28  1350
#   197 0  5  28  1350
#   198 0  5  28  1350
#   199 0  5  28  1350
#   200 0  5  28  1350
#   201 0  5  28  1350
#   202 0  5  28  1350
#   203 0  5  28  1350
#   204 0  5  28  1350
#   205 0  5  28  1350
#   206 0  5  28  1350
#   207 0  5  28  1350
#   208 0  5  28  1350
#   209 0  5  28  1350
#   210 0  5  28  1350
#   211 0  5  28  1350
#   212 0  5  28  1350
#   213 0  5  28  1350
#   214 0  5  28  1350
#   215 0  5  28  1350
#   216 0  5  28  1350
#   217 0  5  28  1350
#   218 0  5  28  1350
#   219 0  5  28  1350
#   220 0  5  28  1350
#   221 0  5  28  1350
#   222 0  5  28  1350
#   223 0  5  28  1350
#   224 0  5  28  1350
#   225 0  5  28  1350
#   226 0  5  28  1350
#   227 0  5  28  1350
#   228 0  5  28  1350
#   229 0  5  28  1350
#   230 0  5  28  1350
#   231 0  5  28  1350
#   232 0  5  28  1350
#   233 0  5  28  1350
#   234 0  5  28  1350
#   235 0  5  28  1350
#   236 0  5  28  1350
#   237 0  5  28  1350
#   238 0  5  28  1350
#   239 0  5  28  1350
#   240 0  5  28  1350
#   241 0  5  28  1350
#   242 0  5  28  1350
#   243 0  5  28  1350
#   244 0  5  28  1350
#   245 0  5  28  1350
#   246 0  5  28  1350
#   247 0  5  28  1350
#   248 0  5  28  1350
#   249 0  5  28  1350
#   250 0  5  28  1350
#   251 0  5  28  1350
#   252 0  5  28  1350
#   253 0  5  28  1350
#   254 0  5  28  1350
#   255 0  5  28  1350
#   256 0  5  28  1350
#   257 0  5  28  1350
#   258 0  5  28  1350
#   259 0  5  28  1350
#   260 0  5  28  1350
#   261 0  5  28  1350
#   262 0  5  28  1350
#   263 0  5  28  1350
#   264 0  5  28  1350
#   265 0  5  28  1350
#   266 0  5  28  1350
#   267 0  5  28  1350
#   268 0  5  28  1350
#   269 0  5  28  1350
#   270 0  5  28  1350
#   271 0  5  28  1350
#   272 0  5  28  1350
#   273 0  5  28  1350
#   274 0  5  28  1350
#   275 0  5  28  1350
#   276 0  5  28  1350
#   277 0  5  28  1350
#   278 0  5  28  1350
#   279 0  5  28  1350
#   280 0  5  28  1350
#   281 0  5  28  1350
#   282 0  5  28  1350
#   283 0  5  28  1350
#   284 0  5  28  1350
#   285 0  5  28  1350
#   286 0  5  28  1350
#   287 0  5  28  1350
#   288 0  5  28  1350
#   289 0  5  28  1350
#   290 0  5  28  1350
#   291 0  5  28  1350
#   292 0  5  28  1350
#   293 0  5  28  1350
#   294 0  5  28  1350
#   295 0  5  28  1350
#   296 0  5  28  1350
#   297 0  5  28  1350
#   298 0  5  28  1350
#   299 0  5  28  1350
#   300 0  5  28  1350
#   301 0  5  28  1350
#   302 0  5  28  1350
#   303 0  5  28  1350
#   304 0  5  28  1350
#   305 0  5  28  1350
#   306 0  5  28  1350
#   307 0  5  28  1350
#   308 0  5  28  1350
#   309 0  5  28  1350
#   310 0  5  28  1350
#   311 0  5  28  1350
#   312 0  5  28  1350
#   313 0  5  28  1350
#   314 0  5  28  1350
#   315 0  5  28  1350
#   316 0  5  28  1350
#   317 0  5  28  1350
#   318 0  5  28  1350
#   319 0  5  28  1350
#   320 0  5  28  1350
#   321 0  5  28  1350
#   322 0  5  28  1350
#   323 0  5  28  1350
#   324 0  5  28  1350
#   325 0  5  28  1350
#   326 0  5  28  1350
#   327 0  5  28  1350
#   328 0  5  28  1350
#   329 0  5  28  1350
#   330 0  5  28  1350
#   331 0  5  28  1350
#   332 0  5  28  1350
#   333 0  5  28  1350
#   334 0  5  28  1350
#   335 0  5  28  1350
#   336 0  5  28  1350
#   337 0  5  28  1350
#   338 0  5  28  1350
#   339 0  5  28  1350
#   340 0  5  28  1350
#   341 0  5  28  1350
#   342 0  5  28  1350
#   343 0  5  28  1350
#   344 0  5  28  1350
#   345 0  5  28  1350
#   346 0  5  28  1350
#   347 0  5  28  1350
#   348 0  5  28  1350
#   349 0  5  28  1350
#   350 0  5  28  1350
#   351 0  5  28  1350
#   352 0  5  28  1350
#   353 0  5  28  1350
#   354 0  5  28  1350
#   355 0  5  28  1350
#   356 0  5  28  1350
#   357 0  5  28  1350
#   358 0  5  28  1350
#   359 0  5  28  1350
#   360 0  5  28  1350
#   361 0  5  28  1350
#   362 0  5  28  1350
#   363 0  5  28  1350
#   3
```