In [2]:
```python
#importing libraries
from urllib.request import Request, urlopen
from bs4 import BeautifulSoup as soup
import pandas as pd
import matplotlib as plt
import seaborn as sns
import time
```

# Scrapping 150 movies data from given url

```
In [61]:    1  #Scrapping 150 movies data from following url
            2  movies_name = []
            3  movies_rating = []
            4  movies_genre = []
            5  movies_release_date = []
            6  movies_runtime = []
            7  movies_director = []
            8  movies_links = []
            9  movies_budget = []
           10  movies_revenue = []
           11  for i in range(1,151):
           12      main_url = 'https://www.themoviedb.org/movie?page=' + str(i)
           13      req = Request(main_url , headers={'User-Agent': 'Mozilla/5.0'})
           14      webpage = urlopen(req).read()
           15      page_soup = soup(webpage, "html.parser")
           16      time.sleep(0.2)
           17      soup_body = page_soup.body
           18      print(i,main_url)
           19      for j in range(0,1):
           20          a_tag_movie_link = soup_body.find_all('a',class_='image')
           21          href_data = a_tag_movie_link[j].get('href')
           22          movie_title = a_tag_movie_link[j].get('title')
           23          movie_url = 'https://www.themoviedb.org/' + str(href_data)
           24          print(movie_url)
           25          req_jloop = Request(movie_url , headers={'User-Agent': 'Mozilla/5.0'
           26          webpage_jloop = urlopen(req_jloop).read()
           27          page_soup_jloop = soup(webpage_jloop, "html.parser")
           28          soup_body_jloop = page_soup_jloop.body
           29          page_wrap_class = soup_body_jloop.find_all('div',class_='page_wrap m
           30          try :
           31              release_span = page_wrap_class[j].find_all('span',class_='releas
           32              release_text = release_span[0].get_text()
           33              genres_span = page_wrap_class[j].find_all('span',class_='genres'
           34              genres_text = genres_span[0].get_text()
           35              runtime_span = page_wrap_class[j].find_all('span',class_='runtim
           36              runtime_text = runtime_span[0].get_text()
           37              li_profile = page_wrap_class[0].find_all('div',class_='user_scor
           38              rating = li_profile[0].get('data-percent')
           39              money_data = page_wrap_class[0].find_all('section',class_='facts
           40              money_text = money_data[0].find_all('p')
           41              budget_value=money_text[2].text
           42              budget = budget_value.split()[1]
           43              #print(budget)
           44              revenue_value =money_text[3].text
           45              revenue = revenue_value.split()[1]
           46              #print(revenue)
           47              for k in range(0,1):
           48                  li_profile = page_wrap_class[k].find_all('li',class_='profil
           49                  dr = li_profile[k]
           50                  director_text = (dr.find('a').text)
           51          except :
           52              pass
           53          print(j,movie_url,movie_title,release_text,genres_text,runtime_text,
           54          movies_name.append(movie_title)
           55          movies_rating.append(rating)
           56          movies_genre.append(genres_text.strip())
```

```python
57          movies_release_date.append(release_text.strip())
58          movies_runtime.append(runtime_text.strip())
59          movies_director.append(director_text)
60          movies_links.append(movie_url)
61          movies_budget.append(budget)
62          movies_revenue.append(revenue)
63          time.sleep(0.1)
64  #creating dataFrame
65  df = pd.DataFrame({
66      'Name' : movies_name,
67      'Rating' : movies_rating,
68      'Genre' : movies_genre,
69      'Release date' : movies_release_date,
70      'Runtime' : movies_runtime,
71      'Director' : movies_director,
72      'Budget ($)' : movies_budget,
73      'Revenue ($) ' : movies_revenue,
74      'Url' : movies_links
75  })
76  #removing special character
77  df['Genre'] = df['Genre'].map(str).apply(lambda x: x.encode('utf-8').decode(
78  #print(df)
79  #converting dataframe into CSV
80  time.sleep(0.1)
81  df.to_csv('Movies scrapped data.csv',index = False)
82  print('CSV Successfully Created')
83
```

In [ ]:    1