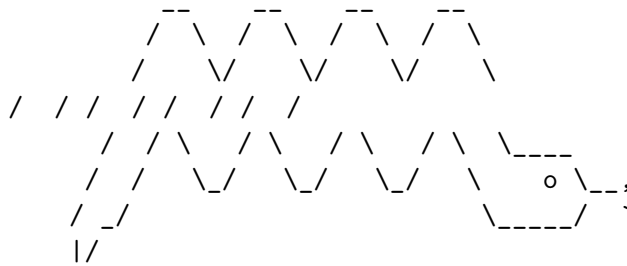# webScrappingImdb

March 9, 2022

## 0.1 Project 1

### 0.1.1 Web Scrapping IMDB

```
[1]: !mamba install bs4==4.10.0 -y
     !mamba install html5lib==1.1 -y
     !pip install lxml==4.6.4
```

```
              __     __     __     __
             /  \   /  \   /  \   /  \
            /    \/     \/     \/     \
       /  / /  / /  / /  /
          /  / \    / \    / \    / \   \____
         /  /   \_/   \_/   \_/   \    o  \__,
        /  _/                         \_____/ `
        |/

        mamba (0.15.3) supported by @QuantStack

        GitHub:  https://github.com/mamba-org/mamba
        Twitter: https://twitter.com/QuantStack

        Looking for: ['bs4==4.10.0']

        pkgs/r/linux-64          [<=>                     ] (00m:00s)
        pkgs/r/linux-64          [=>                      ] (00m:00s) 12 KB / ?? (78.72 KB/s)
        pkgs/r/linux-64          [=>                      ] (00m:00s) 12 KB / ?? (78.72 KB/s)
        pkgs/r/noarch            [>                       ] (--:--) Finalizing…
        pkgs/r/linux-64          [=>                      ] (00m:00s) 12 KB / ?? (78.72 KB/s)
```

```
pkgs/r/noarch            [>                    ] (--:--) Done
pkgs/r/noarch            [====================] (00m:00s) Done
pkgs/r/linux-64          [=>                   ] (00m:00s) 12 KB / ?? (78.72 KB/s)
pkgs/r/linux-64          [<=>                  ] (00m:00s) 12 KB / ?? (78.72 KB/s)
pkgs/r/linux-64          [<=>                  ] (00m:00s) 680 KB / ?? (2.19 MB/s)
pkgs/r/linux-64          [<=>                  ] (00m:00s) 680 KB / ?? (2.19 MB/s)
pkgs/main/linux-64       [<=>                  ] (00m:00s)
pkgs/r/linux-64          [<=>                  ] (00m:00s) 680 KB / ?? (2.19 MB/s)
pkgs/main/linux-64       [=>                   ] (00m:00s) 728 KB / ?? (2.34 MB/s)
pkgs/r/linux-64          [<=>                  ] (00m:00s) 680 KB / ?? (2.19 MB/s)
pkgs/main/linux-64       [=>                   ] (00m:00s) 728 KB / ?? (2.34 MB/s)
pkgs/main/noarch         [<=>                  ] (00m:00s)
pkgs/r/linux-64          [<=>                  ] (00m:00s) 680 KB / ?? (2.19 MB/s)
pkgs/main/linux-64       [=>                   ] (00m:00s) 728 KB / ?? (2.34 MB/s)
pkgs/main/noarch         [=>                   ] (00m:00s) 364 KB / ?? (1.17 MB/s)
pkgs/r/linux-64          [ <=>                 ] (00m:00s) Finalizing…
pkgs/main/linux-64       [=>                   ] (00m:00s) 728 KB / ?? (2.34 MB/s)
pkgs/main/noarch         [=>                   ] (00m:00s) 364 KB / ?? (1.17 MB/s)
pkgs/r/linux-64          [ <=>                 ] (00m:00s) Done
pkgs/main/linux-64       [=>                   ] (00m:00s) 728 KB / ?? (2.34 MB/s)
pkgs/main/noarch         [=>                   ] (00m:00s) 364 KB / ?? (1.17 MB/s)
pkgs/r/linux-64          [====================] (00m:00s) Done
pkgs/main/linux-64       [=>                   ] (00m:00s) 728 KB / ?? (2.34 MB/s)
pkgs/main/noarch         [=>                   ] (00m:00s) 364 KB / ?? (1.17 MB/s)
pkgs/main/linux-64       [=>                   ] (00m:00s) 728 KB / ?? (2.34 MB/s)
pkgs/main/noarch         [<=>                  ] (00m:00s) Finalizing…
pkgs/main/linux-64       [=>                   ] (00m:00s) 728 KB / ?? (2.34 MB/s)
pkgs/main/noarch         [<=>                  ] (00m:00s) Done
pkgs/main/noarch         [====================] (00m:00s) Done
pkgs/main/linux-64       [=>                   ] (00m:00s) 728 KB / ?? (2.34 MB/s)
pkgs/main/linux-64       [<=>                  ] (00m:00s) 728 KB / ?? (2.34 MB/s)
pkgs/main/linux-64       [ <=>                 ] (00m:00s) 1 MB / ?? (3.15 MB/s)
pkgs/main/linux-64       [  <=>                ] (00m:00s) 1 MB / ?? (3.15 MB/s)
pkgs/main/linux-64       [  <=>                ] (00m:00s) 2 MB / ?? (3.67 MB/s)
pkgs/main/linux-64       [   <=>               ] (00m:00s) 2 MB / ?? (3.67 MB/s)
pkgs/main/linux-64       [   <=>               ] (00m:00s) 3 MB / ?? (3.92 MB/s)
pkgs/main/linux-64       [    <=>              ] (00m:00s) 3 MB / ?? (3.92 MB/s)
pkgs/main/linux-64       [     <=>             ] (00m:00s) 4 MB / ?? (4.13 MB/s)
pkgs/main/linux-64       [     <=>             ] (00m:00s) Finalizing…
pkgs/main/linux-64       [     <=>             ] (00m:00s) Done
pkgs/main/linux-64       [====================] (00m:00s) Done

Pinned packages:
  - python 3.7.*


Transaction
```

```
Prefix: /home/jupyterlab/conda/envs/python

Updating specs:

 - bs4==4.10.0
 - ca-certificates
 - certifi
 - openssl


  Package                 Version  Build           Channel                      Size

  Install:


  + beautifulsoup4         4.10.0  pyh06a4308_0     pkgs/main/noarch
85 KB
  + bs4                    4.10.0  hd3eb1b0_0       pkgs/main/noarch
10 KB
  + soupsieve               2.3.1  pyhd3eb1b0_0     pkgs/main/noarch
34 KB

  Change:


  - certifi             2021.10.8  py37h89c1867_1   installed
  + certifi             2021.10.8  py37h06a4308_2   pkgs/main/linux-64
151 KB

  Upgrade:


  - ca-certificates     2021.10.8  ha878542_0       installed
  + ca-certificates      2022.2.1  h06a4308_0       pkgs/main/linux-64
122 KB
  - openssl                1.1.1l  h7f98852_0       installed
  + openssl                1.1.1m  h7f8727e_0       pkgs/main/linux-64
3 MB

  Summary:

  Install: 3 packages
  Change: 1 packages
  Upgrade: 2 packages

  Total download: 3 MB
```
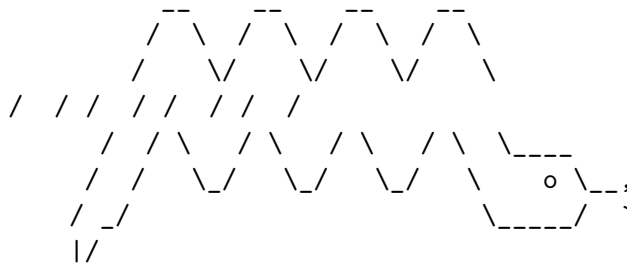
```
Downloading  [>                                        ] (00m:00s)    1.29 MB/s
Extracting   [>                                        ] (--:--)
Downloading  [==>                                      ] (00m:00s)    3.54 MB/s
Extracting   [>                                        ] (--:--)
Finished beautifulsoup4                  (00m:00s)            85
KB       1 MB/s
Downloading  [==>                                      ] (00m:00s)    3.54 MB/s
Extracting   [>                                        ] (--:--)
Downloading  [==>                                      ] (00m:00s)    3.54 MB/s
Extracting   [>                                        ] (--:--)
Downloading  [==>                                      ] (00m:00s)    3.54 MB/s
Extracting   [>                                        ] (--:--)
Finished certifi                         (00m:00s)           151
KB       2 MB/s
Downloading  [==>                                      ] (00m:00s)    3.54 MB/s
Extracting   [>                                        ] (--:--)
Downloading  [==>                                      ] (00m:00s)    3.54 MB/s
Extracting   [>                                        ] (--:--)
Downloading  [==>                                      ] (00m:00s)    3.54 MB/s
Extracting   [======>                                  ] (00m:00s)        1 / 6
Downloading  [===>                                     ] (00m:00s)    3.61 MB/s
Extracting   [======>                                  ] (00m:00s)        1 / 6
Downloading  [===>                                     ] (00m:00s)    3.61 MB/s
Extracting   [======>                                  ] (00m:00s)        1 / 6
Finished soupsieve                       (00m:00s)            34
KB     468 KB/s
Downloading  [===>                                     ] (00m:00s)    3.61 MB/s
Extracting   [======>                                  ] (00m:00s)        1 / 6
Downloading  [===>                                     ] (00m:00s)    3.61 MB/s
Extracting   [======>                                  ] (00m:00s)        1 / 6
Downloading  [===>                                     ] (00m:00s)    3.25 MB/s
Extracting   [======>                                  ] (00m:00s)        1 / 6
Finished bs4                             (00m:00s)            10
KB     119 KB/s
Downloading  [===>                                     ] (00m:00s)    3.25 MB/s
Extracting   [======>                                  ] (00m:00s)        1 / 6
Downloading  [===>                                     ] (00m:00s)    3.25 MB/s
Extracting   [======>                                  ] (00m:00s)        1 / 6
Downloading  [====>                                    ] (00m:00s)    4.57 MB/s
Extracting   [======>                                  ] (00m:00s)        1 / 6
Finished ca-certificates                 (00m:00s)           122
KB       1 MB/s
Downloading  [=====>                                   ] (00m:00s)    4.57 MB/s
Extracting   [======>                                  ] (00m:00s)        1 / 6
Downloading  [=====>                                   ] (00m:00s)    4.57 MB/s
Extracting   [======>                                  ] (00m:00s)        1 / 6
Downloading  [=====>                                   ] (00m:00s)    4.57 MB/s
```

```
Extracting      [=============>                                  ] (00m:00s)          2 / 6
Downloading     [=====>                                          ] (00m:00s)      4.57 MB/s
Extracting      [=============>                                  ] (00m:00s)          2 / 6
Downloading     [=====>                                          ] (00m:00s)      4.57 MB/s
Extracting      [===================>                            ] (00m:00s)          3 / 6
Downloading     [===========================================] (00m:00s)     26.67 MB/s
Extracting      [===================>                            ] (00m:00s)          3 / 6
Downloading     [===========================================] (00m:00s)     26.67 MB/s
Extracting      [===================>                            ] (00m:00s)          3 / 6
Downloading     [===========================================] (00m:00s)     26.67 MB/s
Extracting      [=========================>                      ] (00m:00s)          4 / 6
Downloading     [===========================================] (00m:00s)     26.67 MB/s
Extracting      [==========================>                     ] (00m:00s)          4 / 6
Downloading     [===========================================] (00m:00s)     26.67 MB/s
Extracting      [==============================>                 ] (00m:00s)          5 / 6
Finished openssl                              (00m:00s)                    3
MB      23 MB/s
Downloading     [===========================================] (00m:00s)     26.67 MB/s
Extracting      [================================>               ] (00m:00s)          5 / 6
Downloading     [===========================================] (00m:00s)     26.67 MB/s
Extracting      [================================>               ] (00m:00s)          5 / 6
Downloading     [===========================================] (00m:00s)     26.67 MB/s
Extracting      [================================>               ] (00m:00s)          5 / 6
Downloading     [===========================================] (00m:00s)     26.67 MB/s
Extracting      [===========================================] (00m:00s)          6 / 6
Preparing transaction: done
Verifying transaction: done
Executing transaction: done


                    __    __    __    __
                   /  \  /  \  /  \  /  \
                  /    \/    \/    \/    \
        / / /  / /  / /  /
                / / \   / \   / \   / \  \____
               / /   \_/   \_/   \_/   \    o \__,
              / _/                      \_____/  `
            |/

           mamba (0.15.3) supported by @QuantStack

           GitHub:  https://github.com/mamba-org/mamba
           Twitter: https://twitter.com/QuantStack
```

```
Looking for: ['html5lib==1.1']

pkgs/main/linux-64       Using cache
pkgs/main/noarch         Using cache
pkgs/r/linux-64          Using cache
pkgs/r/noarch            Using cache


Pinned packages:
  - python 3.7.*



Transaction

  Prefix: /home/jupyterlab/conda/envs/python

  Updating specs:

   - html5lib==1.1
   - ca-certificates
   - certifi
   - openssl


  Package           Version  Build            Channel                 Size

  Install:


  + html5lib           1.1  pyhd3eb1b0_0  pkgs/main/noarch         91 KB
  + webencodings     0.5.1  py37_1        pkgs/main/linux-64       19 KB

  Summary:

  Install: 2 packages

  Total download: 110 KB



Downloading  [======>                                  ] (00m:00s)  134.47 KB/s
Extracting   [>                                        ] (--:--)
Finished webencodings                        (00m:00s)              19
KB     134 KB/s
Downloading  [======>                                  ] (00m:00s)  134.47 KB/s
```

```
Extracting     [>                                    ] (--:--)
Downloading    [======>                    ] (00m:00s)  134.47 KB/s
Extracting     [>                                    ] (--:--)
Downloading    [======>                    ] (00m:00s)  134.47 KB/s
Extracting     [>                                    ] (--:--)
Downloading    [=====================================] (00m:00s)  752.86 KB/s
Extracting     [>                                    ] (--:--)
Finished html5lib              (00m:00s)          91
KB      622 KB/s
Downloading    [=====================================] (00m:00s)  752.86 KB/s
Extracting     [>                                    ] (--:--)
Downloading    [=====================================] (00m:00s)  752.86 KB/s
Extracting     [>                                    ] (--:--)
Downloading    [=====================================] (00m:00s)  752.86 KB/s
Extracting     [====================>         ] (00m:00s)      1 / 2
Downloading    [=====================================] (00m:00s)  752.86 KB/s
Extracting     [====================>         ] (00m:00s)      1 / 2
Downloading    [=====================================] (00m:00s)  752.86 KB/s
Extracting     [=====================================] (00m:00s)      2 / 2
Preparing transaction: done
Verifying transaction: done
Executing transaction: done
Requirement already satisfied: lxml==4.6.4 in
/home/jupyterlab/conda/envs/python/lib/python3.7/site-packages (4.6.4)
```

```python
[4]: import requests
     from bs4 import BeautifulSoup
     import pandas
```

# 1 Using requests method to get the html code of website

### 1.0.1 Use of try and expect

### 1.0.2 Use of raise_for_status() method

```python
[26]: source = requests.get('https://www.imdb.com/chart/top/')
      source.raise_for_status()
```

```python
[27]: soup = BeautifulSoup(source.text,'html.parser')
```

```
[41]: ### using movies find tr from tbody and movies length using print(len(movies))
      ## now lets open a loop

      movies = soup.find('tbody',class_= 'lister-list').find_all('tr')



      for movie in movies :
          #only text

          name = movie.find('td',class_='titleColumn').a.text
          rank = movie.find('td',class_='titleColumn').get_text(strip =True).split('.
       ↪')[0]
          year = movie.find('td',class_='titleColumn').span.text.strip('()')
          rating = movie.find('td',class_ ='ratingColumn imdbRating').strong.text
          print(rank,name,year,rating)
```

```
1 The Shawshank Redemption 1994 9.2
2 The Godfather 1972 9.2
3 The Dark Knight 2008 9.0
4 The Godfather: Part II 1974 9.0
5 12 Angry Men 1957 9.0
6 Schindler's List 1993 8.9
7 The Lord of the Rings: The Return of the King 2003 8.9
8 Pulp Fiction 1994 8.9
9 The Lord of the Rings: The Fellowship of the Ring 2001 8.8
10 The Good, the Bad and the Ugly 1966 8.8
11 Forrest Gump 1994 8.8
12 Fight Club 1999 8.8
13 Inception 2010 8.7
14 The Lord of the Rings: The Two Towers 2002 8.7
15 Star Wars: Episode V - The Empire Strikes Back 1980 8.7
16 The Matrix 1999 8.7
17 Goodfellas 1990 8.7
18 One Flew Over the Cuckoo's Nest 1975 8.6
19 Se7en 1995 8.6
20 Seven Samurai 1954 8.6
21 It's a Wonderful Life 1946 8.6
22 The Silence of the Lambs 1991 8.6
23 Saving Private Ryan 1998 8.6
24 City of God 2002 8.6
25 Life Is Beautiful 1997 8.6
26 The Green Mile 1999 8.6
27 Star Wars 1977 8.6
28 Interstellar 2014 8.6
```

29 Terminator 2: Judgment Day 1991 8.5
30 Back to the Future 1985 8.5
31 Spirited Away 2001 8.5
32 Psycho 1960 8.5
33 Léon: The Professional 1994 8.5
34 The Pianist 2002 8.5
35 Parasite 2019 8.5
36 The Lion King 1994 8.5
37 Gladiator 2000 8.5
38 American History X 1998 8.5
39 The Usual Suspects 1995 8.5
40 The Departed 2006 8.5
41 Spider-Man: No Way Home 2021 8.5
42 The Prestige 2006 8.5
43 Casablanca 1942 8.5
44 Whiplash 2014 8.5
45 The Intouchables 2011 8.5
46 Modern Times 1936 8.5
47 Once Upon a Time in the West 1968 8.5
48 Hara-Kiri 1962 8.5
49 Grave of the Fireflies 1988 8.5
50 Alien 1979 8.4
51 Rear Window 1954 8.4
52 City Lights 1931 8.4
53 Memento 2000 8.4
54 Apocalypse Now 1979 8.4
55 Cinema Paradiso 1988 8.4
56 Indiana Jones and the Raiders of the Lost Ark 1981 8.4
57 Django Unchained 2012 8.4
58 WALL·E 2008 8.4
59 The Lives of Others 2006 8.4
60 Sunset Blvd. 1950 8.4
61 The Shining 1980 8.4
62 Paths of Glory 1957 8.4
63 The Great Dictator 1940 8.4
64 Avengers: Infinity War 2018 8.4
65 Witness for the Prosecution 1957 8.4
66 Aliens 1986 8.4
67 American Beauty 1999 8.3
68 The Dark Knight Rises 2012 8.3
69 Dr. Strangelove or: How I Learned to Stop Worrying and Love the Bomb 1964 8.3
70 Joker 2019 8.3
71 Spider-Man: Into the Spider-Verse 2018 8.3
72 The Batman 2022 8.3
73 Old Boy 2003 8.3
74 Braveheart 1995 8.3
75 Toy Story 1995 8.3
76 Amadeus 1984 8.3

77 Coco 2017 8.3
78 Inglourious Basterds 2009 8.3
79 The Boat 1981 8.3
80 Avengers: Endgame 2019 8.3
81 Princess Mononoke 1997 8.3
82 Once Upon a Time in America 1984 8.3
83 Good Will Hunting 1997 8.3
84 Toy Story 3 2010 8.3
85 Requiem for a Dream 2000 8.3
86 3 Idiots 2009 8.3
87 Singin' in the Rain 1952 8.3
88 Your Name. 2016 8.3
89 Star Wars: Episode VI - Return of the Jedi 1983 8.3
90 Reservoir Dogs 1992 8.3
91 Eternal Sunshine of the Spotless Mind 2004 8.3
92 2001: A Space Odyssey 1968 8.3
93 Citizen Kane 1941 8.3
94 High and Low 1963 8.3
95 Lawrence of Arabia 1962 8.3
96 M 1931 8.3
97 Capernaum 2018 8.3
98 North by Northwest 1959 8.3
99 The Hunt 2012 8.3
100 Vertigo 1958 8.3
101 Amélie 2001 8.3
102 A Clockwork Orange 1971 8.3
103 Full Metal Jacket 1987 8.2
104 Scarface 1983 8.2
105 Double Indemnity 1944 8.2
106 Come and See 1985 8.2
107 The Apartment 1960 8.2
108 Taxi Driver 1976 8.2
109 To Kill a Mockingbird 1962 8.2
110 Hamilton 2020 8.2
111 The Sting 1973 8.2
112 L.A. Confidential 1997 8.2
113 Up 2009 8.2
114 Heat 1995 8.2
115 Ikiru 1952 8.2
116 Snatch 2000 8.2
117 Die Hard 1988 8.2
118 Indiana Jones and the Last Crusade 1989 8.2
119 A Separation 2011 8.2
120 Metropolis 1927 8.2
121 Bicycle Thieves 1948 8.2
122 Incendies 2010 8.2
123 1917 2019 8.2
124 Like Stars on Earth 2007 8.2

```
125 Batman Begins 2005 8.2
126 For a Few Dollars More 1965 8.2
127 Dangal 2016 8.2
128 Downfall 2004 8.2
129 The Kid 1921 8.2
130 Some Like It Hot 1959 8.2
131 The Father 2020 8.2
132 All About Eve 1950 8.2
133 Green Book 2018 8.2
134 The Wolf of Wall Street 2013 8.2
135 Unforgiven 1992 8.2
136 Casino 1995 8.2
137 Pan's Labyrinth 2006 8.2
138 Judgment at Nuremberg 1961 8.2
139 Ran 1985 8.2
140 A Beautiful Mind 2001 8.2
141 The Sixth Sense 1999 8.2
142 Monty Python and the Holy Grail 1975 8.2
143 There Will Be Blood 2007 8.2
144 The Truman Show 1998 8.2
145 Yojimbo 1961 8.2
146 The Treasure of the Sierra Madre 1948 8.2
147 Shutter Island 2010 8.2
148 The Great Escape 1963 8.2
149 Rashomon 1950 8.1
150 Jurassic Park 1993 8.1
151 Kill Bill: Vol. 1 2003 8.1
152 Finding Nemo 2003 8.1
153 No Country for Old Men 2007 8.1
154 Raging Bull 1980 8.1
155 The Elephant Man 1980 8.1
156 V for Vendetta 2005 8.1
157 Gone with the Wind 1939 8.1
158 Chinatown 1974 8.1
159 Inside Out 2015 8.1
160 Lock, Stock and Two Smoking Barrels 1998 8.1
161 The Thing 1982 8.1
162 Dial M for Murder 1954 8.1
163 The Secret in Their Eyes 2009 8.1
164 Howl's Moving Castle 2004 8.1
165 The Bridge on the River Kwai 1957 8.1
166 Trainspotting 1996 8.1
167 Three Billboards Outside Ebbing, Missouri 2017 8.1
168 Warrior 2011 8.1
169 Gran Torino 2008 8.1
170 Fargo 1996 8.1
171 My Neighbor Totoro 1988 8.1
172 Prisoners 2013 8.1
```

173 Million Dollar Baby 2004 8.1
174 Blade Runner 1982 8.1
175 The Gold Rush 1925 8.1
176 Catch Me If You Can 2002 8.1
177 On the Waterfront 1954 8.1
178 Children of Heaven 1997 8.1
179 Harry Potter and the Deathly Hallows: Part 2 2011 8.1
180 The Third Man 1949 8.1
181 Gone Girl 2014 8.1
182 Ben-Hur 1959 8.1
183 12 Years a Slave 2013 8.1
184 The General 1926 8.1
185 The Deer Hunter 1978 8.1
186 Wild Strawberries 1957 8.1
187 Before Sunrise 1995 8.1
188 In the Name of the Father 1993 8.1
189 Pather Panchali 1955 8.1
190 Mr. Smith Goes to Washington 1939 8.1
191 The Grand Budapest Hotel 2014 8.1
192 Room 2015 8.1
193 Sherlock Jr. 1924 8.1
194 Hacksaw Ridge 2016 8.1
195 How to Train Your Dragon 2010 8.1
196 The Wages of Fear 1953 8.1
197 Memories of Murder 2003 8.1
198 The Seventh Seal 1957 8.1
199 Barry Lyndon 1975 8.1
200 The Big Lebowski 1998 8.1
201 Klaus 2019 8.1
202 Mad Max: Fury Road 2015 8.1
203 Wild Tales 2014 8.1
204 Monsters, Inc. 2001 8.1
205 Mary and Max 2009 8.1
206 Jaws 1975 8.1
207 The Passion of Joan of Arc 1928 8.1
208 Hotel Rwanda 2004 8.1
209 Rocky 1976 8.1
210 Dead Poets Society 1989 8.1
211 Tokyo Story 1953 8.1
212 Platoon 1986 8.1
213 Ford v Ferrari 2019 8.1
214 The Terminator 1984 8.1
215 Stand by Me 1986 8.1
216 Rush 2013 8.0
217 Into the Wild 2007 8.0
218 The Wizard of Oz 1939 8.0
219 Logan 2017 8.0
220 Spotlight 2015 8.0

```
221 Network 1976 8.0
222 Groundhog Day 1993 8.0
223 The Exorcist 1973 8.0
224 Ratatouille 2007 8.0
225 Hachi: A Dog's Tale 2009 8.0
226 The Incredibles 2004 8.0
227 Dersu Uzala 1975 8.0
228 The Best Years of Our Lives 1946 8.0
229 Before Sunset 2004 8.0
230 Rebecca 1940 8.0
231 Dune 2021 8.0
232 The Grapes of Wrath 1940 8.0
233 My Father and My Son 2005 8.0
234 Cool Hand Luke 1967 8.0
235 To Be or Not to Be 1942 8.0
236 Amores perros 2000 8.0
237 The Battle of Algiers 1966 8.0
238 Pirates of the Caribbean: The Curse of the Black Pearl 2003 8.0
239 The Sound of Music 1965 8.0
240 Life of Brian 1979 8.0
241 The 400 Blows 1959 8.0
242 Persona 1966 8.0
243 It Happened One Night 1934 8.0
244 La Haine 1995 8.0
245 Aladdin 1992 8.0
246 Beauty and the Beast 1991 8.0
247 Jai Bhim 2021 8.0
248 Gandhi 1982 8.0
249 The Help 2011 8.0
250 The Handmaiden 2016 8.0
```

### 1.0.3 using movies find tr from tbody and movies length using print(len(movies))

## 1.1 now lets open a loop

movies = soup.find('tbody',class_= 'lister-list').find_all('tr')

for movie in movies :

```
first tr tag
print(movie)
break


 #print first td


name = movie.find('td',class_='titleColumn')


print(name)
break
```

```python
    #only a

name = movie.find('td',class_='titleColumn').a
print(name)
break
 #only text

name = movie.find('td',class_='titleColumn').a.text
print(name)
break

#print rank
 rank = movie.find('td',class_='titleColumn').text
 print(rank)
 break
```

[48]:

```
      File "/tmp/ipykernel_342/3443480531.py", line 1
    install openpyxl
               ^
  SyntaxError: invalid syntax
```

[ ]: