

```
In [11]: 1 import numpy as np #linear algebra
2 import pandas as pd #data processing
3 # step 1 get top line
4 from bs4 import BeautifulSoup as bs
5 import bs4
6 import requests as rq # get url
7 import re
8 from datetime import date, time, datetime
9
10 pd.set_option('display.max_rows', None)
11 pd.set_option('display.max_columns', None)
12 pd.set_option('display.expand_frame_repr', False)
13 pd.set_option('max_colwidth', None)
14
15 # https://www.the-numbers.com/person/41500401-Robert-Downey-Jr#tab=acting
16 df_starsAct = pd.read_excel ('./Prj_Data/ImdbScrapingData/StarEarningsv1_NoM
17 df_starswriter = pd.read_excel ('./Prj_Data/ImdbScrapingData/StarEarningsv1_N
18 df_starsDir = pd.read_excel ('./Prj_Data/ImdbScrapingData/StarEarningsv1_NoM
19 df_starsProd = pd.read_excel ('./Prj_Data/ImdbScrapingData/StarEarningsv1_No
```

```
In [ ]:
```

```
1
```

```

In [74]: 1 dfmovies = pd.DataFrame(columns=['Actor', 'Born', 'ReleaseDate', 'Title', 'Ro
2
3 # df_starsAct_1050_1500 = df_starsAct.iloc[1051:]
4 # tab=technical
5 # "#tab=acting"
6
7 for index, row in df_starswriter.iterrows():
8     row["Link"]
9     currentname = row["Name"]
10    currentlink = row["Link"]
11    currentlink = currentlink + "#tab=technical"
12
13    # ['ReleaseDate', 'Title', 'Role', 'DomesticBox Office', 'InternationalBox O
14 #     the_getString = 'https://www.the-numbers.com/person/41500401-Robert-Do
15 the_getString = currentlink
16 #     print(the_getString)
17 Thecols = []
18 Themovies = []
19
20 r=rq.get(the_getString)
21 # 'html.parser'
22 # p=bs(r.text, 'lxml')
23 p=bs(r.text, 'html.parser')
24
25
26 # BornOn=p.find("table", id="all_acting_credits")
27 bornOns = p.find_all("a", href=re.compile("/on-this-day/"))
28 for bornon in bornOns:
29     #     print(bornon.get_text())
30     CurrentBornOn = str(bornon.get_text())
31
32
33
34 #     tbody=p.find("table", id="all_acting_credits")
35 tbody=p.find("table", id="all_technical_credits")
36 trs = tbody.find_all("tr")
37 # print(tbody)
38 for tr in trs:
39     #     print(tr)
40     #     print("loop")
41     #     if tr.find_all("th"):
42         ths = tr.find_all("th")
43         for th in ths:
44             Thecols.append(th.text)
45         #         print(th.text)
46         Thecols = Thecols
47
48
49     Themovies.append(currentname)
50     Themovies.append(CurrentBornOn)
51     if tr.find_all("td"):
52         tds = tr.find_all("td")
53         #         print(tds)
54
55     #         print(f' Lend of td {len(tds)}')
56     counter = 1

```

```

57
58     for td in tds:
59         if counter < len(tds):
60
61             #         print(f' the counter = {counter}')
62             #         print(td.text)
63             if not td.text:
64                 text = "NA"
65             elif td.text:
66                 text = td.text
67             #         print(text)
68             Themovies.append(text)
69             #         print(f' the counter = {counter}')
70
71         elif counter == len(tds):
72
73             #         print("inside >6")
74             counter = 0
75             #         print(f' the LENGTH {len(Themovies)}')
76
77
78         additionalColsToAdd = len(dfmovies.columns) - len(Themov
79
80
81         if additionalColsToAdd == 1:
82             #         print("appending1")
83             Themovies.append("NA")
84             #
85
86         if additionalColsToAdd == 2:
87             #         print("appending2")
88             Themovies.append("NA")
89             themovies.append("NA")
90
91
92         if len(Themovies) == len(dfmovies.columns):
93             #         print(Themovies)
94             #         dfmovies.loc[len(dfmovies)] = ['Dec 22, 2021', 'Si
95             #         themovies = ['Dec 22, 2021', 'Sing 2', 'Baxter', '
96             #         res = dict(zip(Thecols, Themovies))
97             x = pd.Series(Themovies, index = dfmovies.columns)
98             #         dfmovies = dfmovies.append(res, ignore_index=True)
99             #         dfmovies = dfmovies.append(x, ignore_index=True)
100            dfmovies.loc[len(dfmovies)] = Themovies
101            #         print(dfmovies)
102            Themovies = []
103            counter = counter +1
104            #         print("loop next")
105            #         dfmovies["Actor"] =
106            # print(dfmovies)

```

In [75]: 1 dfmovies.info()

```
In [72]: 1 dfmovies.tail(100)
```

...

```
In [ ]: 1
```

```
In [77]: 1 dfmovies.to_excel("starActorMovies_Writer.xlsx")
```

```
In [52]: 1 dfmovies.tail()
```

...

```
In [ ]: 1
```

```
In [ ]: 1
```