

```
In [1]: 1 from selenium import webdriver
2 from selenium.webdriver.common.by import By
3 from selenium.webdriver.support.ui import WebDriverWait
4 from selenium.webdriver.support import expected_conditions as EC
5 from selenium.common.exceptions import TimeoutException
6 from IPython.display import clear_output
7 import pandas as pd
```

executed in 302ms, finished 15:44:36 2021-02-21

```
In [2]: 1 # set up selenium webscraping browser
2
3 option = webdriver.ChromeOptions()
4 option.add_argument('incognito')
5
6 browser = webdriver.Chrome(options=option)
```

executed in 1.17s, finished 15:44:37 2021-02-21



```

In [3]: 1 # scrape esrb.org for game information
2
3 pg = 1
4 games_list = []
5
6 while True:
7     browser.get("https://www.esrb.org/search/?searchKeyword=" \
8                 "&platform=Nintendo%20Switch%2CPlayStation%204%2CXB%20One" \
9                 "&rating=E%2CE%2B%2CT%2CM%2CA0&descriptor=All%20Content" \
10                "&pg={}&searchType=All".format(pg))
11     try:
12         # make selenium wait until games are loaded before moving on
13         element = WebDriverWait(browser, 10).until(
14             EC.presence_of_element_located((By.CLASS_NAME, 'game')))
15
16         # games on current page
17         results = browser.find_elements_by_xpath("//div[@class='game']")
18
19         titles, consoles, ratings, descriptors = [], [], [], []
20
21         # pull relevant content from results
22         for x in results:
23             # titles
24             title = x.find_element_by_css_selector('h2').text
25             titles.append(title)
26
27             # consoles
28             console = x.find_element_by_class_name('platforms').text
29             consoles.append(console)
30
31             # ratings are displayed on the page as an image,
32             # so we have to pull the rating from the image's url
33             xpath = x.find_element_by_css_selector('img')
34             rating = xpath.get_attribute('src')[58:-4]
35             ratings.append(rating)
36
37             # descriptors
38             desc = x.find_elements_by_css_selector('td')[1].text
39             descriptors.append(desc)
40
41         # collect the results as a list of dictionaries
42         for i in range(len(titles)):
43             game_dict = {}
44             game_dict['title'] = titles[i]
45             game_dict['consoles'] = consoles[i]
46             game_dict['rating'] = ratings[i]
47             game_dict['descriptors'] = descriptors[i]
48             games_list.append(game_dict)
49
50         # progress indicator
51         clear_output(wait=True)
52         print('Page: {}'.format(pg))
53         pg += 1
54
55         # stop running when no more games are found
56     except:
57         print('Stopped on page {}'.format(pg))
58         break

```

executed in 17m 59s, finished 16:02:35 2021-02-21

Page: 504

Stopped on page 505

In [4]:

```
1 # preview of games_list
2 games_list[:5]
```

executed in 13ms, finished 16:02:35 2021-02-21

Out[4]:

```
[{'title': 'Blizzard Arcade Collection',
  'consoles': 'Windows PC, PlayStation 4, Nintendo Switch, Xbox One',
  'rating': 'T',
  'descriptors': 'Blood, Fantasy Violence, Language, Use of Tobacco'},
 {'title': 'Rez Infinite',
  'consoles': 'PlayStation 4',
  'rating': 'E10plus',
  'descriptors': 'Fantasy Violence'},
 {'title': 'Hotshot Racing',
  'consoles': 'PlayStation 4, Nintendo Switch',
  'rating': 'E10plus',
  'descriptors': 'Alcohol Reference, Language, Mild Violence'},
 {'title': 'Sea of Solitude : The Director's Cut',
  'consoles': 'Nintendo Switch',
  'rating': 'T',
  'descriptors': 'Fantasy Violence, Language'},
 {'title': 'Ape Out',
  'consoles': 'Nintendo Switch',
  'rating': 'T',
  'descriptors': 'Blood and Gore, Violence'}]
```

In [5]:

```
1 # dataframe of collected data
2 df = pd.DataFrame(games_list)
3 df.head()
```

executed in 30ms, finished 16:02:35 2021-02-21

Out[5]:

	title	consoles	rating	descriptors
0	Blizzard Arcade Collection	Windows PC, PlayStation 4, Nintendo Switch, Xb...	T	Blood, Fantasy Violence, Language, Use of Tobacco
1	Rez Infinite	PlayStation 4	E10plus	Fantasy Violence
2	Hotshot Racing	PlayStation 4, Nintendo Switch	E10plus	Alcohol Reference, Language, Mild Violence
3	Sea of Solitude : The Director's Cut	Nintendo Switch	T	Fantasy Violence, Language
4	Ape Out	Nintendo Switch	T	Blood and Gore, Violence

In [6]:

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

executed in 13ms, finished 16:02:35 2021-02-21

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5034 entries, 0 to 5033
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  -
0   title       5034 non-null   object
1   consoles    5034 non-null   object
2   rating      5034 non-null   object
3   descriptors 5034 non-null   object
dtypes: object(4)
memory usage: 157.4+ KB
```

```
In [7]: 1 # split consoles and descriptors columns
2 df.descriptors = df.descriptors.map(lambda x: x.split(', '))
3 df.consoles = df.consoles.map(lambda x: x.split(', '))
4
5 df.head()
```

executed in 15ms, finished 16:02:35 2021-02-21

Out[7]:

	title	consoles	rating	descriptors
0	Blizzard Arcade Collection	[Windows PC, PlayStation 4, Nintendo Switch, X...	T	[Blood, Fantasy Violence, Language, Use of Tob...
1	Rez Infinite	[PlayStation 4]	E10plus	[Fantasy Violence]
2	Hotshot Racing	[PlayStation 4, Nintendo Switch]	E10plus	[Alcohol Reference, Language, Mild Violence]
3	Sea of Solitude : The Director's Cut	[Nintendo Switch]	T	[Fantasy Violence, Language]
4	Ape Out	[Nintendo Switch]	T	[Blood and Gore, Violence]

```
In [8]: 1 # choose consoles
2 consoles_list = ['PlayStation 4', 'Xbox One', 'Nintendo Switch']
3 for i in range(len(df)):
4     df.consoles[i] = [x for x in df.consoles[i] if x in consoles_list]
5
6 df.head()
```

executed in 424ms, finished 16:02:36 2021-02-21

Out[8]:

	title	consoles	rating	descriptors
0	Blizzard Arcade Collection	[PlayStation 4, Nintendo Switch, Xbox One]	T	[Blood, Fantasy Violence, Language, Use of Tob...
1	Rez Infinite	[PlayStation 4]	E10plus	[Fantasy Violence]
2	Hotshot Racing	[PlayStation 4, Nintendo Switch]	E10plus	[Alcohol Reference, Language, Mild Violence]
3	Sea of Solitude : The Director's Cut	[Nintendo Switch]	T	[Fantasy Violence, Language]
4	Ape Out	[Nintendo Switch]	T	[Blood and Gore, Violence]

```
In [9]: 1 # rearrange dataframe to put target variable at the end
2 df = df[['title', 'consoles', 'descriptors', 'rating']]
3 df.head()
```

executed in 13ms, finished 16:02:36 2021-02-21

Out[9]:

	title	consoles	descriptors	rating
0	Blizzard Arcade Collection	[PlayStation 4, Nintendo Switch, Xbox One]	[Blood, Fantasy Violence, Language, Use of Tob...	T
1	Rez Infinite	[PlayStation 4]	[Fantasy Violence]	E10plus
2	Hotshot Racing	[PlayStation 4, Nintendo Switch]	[Alcohol Reference, Language, Mild Violence]	E10plus
3	Sea of Solitude : The Director's Cut	[Nintendo Switch]	[Fantasy Violence, Language]	T
4	Ape Out	[Nintendo Switch]	[Blood and Gore, Violence]	T

```
In [11]: 1 # pickle it
2
3 df.to_pickle('esrb_ratings.pkl')
```

executed in 20ms, finished 16:08:50 2021-02-21