

Program No:16

Aim:-Implement a simple web crawler

Program Code

```
import requests

import lxml

from bs4 import BeautifulSoup

url = "https://www.rottentomatoes.com/top/bestofrt/"

headers = {

    'User-Agent' : 'Mozilla/5.0(Windows NT 6.1; WOW64) AppleWebKit/537.36

(KHTML, like Gecko) Chrome/63.0.3239.132 Safari/537.36 QIHU 360SE'

}

f = requests.get(url, headers = headers)

movies_lst = []

soup = BeautifulSoup(f.content, 'html.parser')

movies = soup.find('table', {

    'class': 'table'

}).find_all('a')

print(movies)

num = 0

for anchor in movies:

    urls = 'https://www.rottentomatoes.com' + anchor['href']

    movies_lst.append(urls)

print(movies_lst)

num += 1

movie_url = urls

movie_f = requests.get(movie_url, headers = headers)

movie_soup = BeautifulSoup(movie_f.content, 'lxml')

movie_content = movie_soup.find('div', {
```

```

'class': 'Movies_synopsis clamp clamp-6 js-clamp'
))
print(num,urls,'\n','Movie:' + anchor.string.strip())
print('Movie info:' + movie_content.string.strip())

```

Output

```

C:\Users\ajcemca\PycharmProjects\pythonProject1\venv\Scripts\python.exe C:/Users/ajcemca/PycharmProjects/pythonProject1/WDM.py
[<a class="unstyled articleLink" href="/m/it_happened_one_night">
    It Happened One Night (1934)</a>, <a class="unstyled articleLink" href="/m/citizen_kane">
    Citizen Kane (1941)</a>, <a class="unstyled articleLink" href="/m/the_wizard_of_oz_1939">
    The Wizard of Oz (1939)</a>, <a class="unstyled articleLink" href="/m/modern_times">
    Modern Times (1936)</a>, <a class="unstyled articleLink" href="/m/black_panther_2018">
    Black Panther (2018)</a>, <a class="unstyled articleLink" href="/m/parasite_2019">
    Parasite (Gisaengchung) (2019)</a>, <a class="unstyled articleLink" href="/m/avengers_endgame">
    Avengers: Endgame (2019)</a>, <a class="unstyled articleLink" href="/m/1003707-casablanca">
    Casablanca (1942)</a>, <a class="unstyled articleLink" href="/m/knives_out">
    Knives Out (2019)</a>, <a class="unstyled articleLink" href="/m/us_2019">
    Us (2019)</a>, <a class="unstyled articleLink" href="/m/toy_story_4">
    Toy Story 4 (2019)</a>, <a class="unstyled articleLink" href="/m/lady_bird">
    Lady Bird (2017)</a>, <a class="unstyled articleLink" href="/m/mission_impossible_fallout">
    Mission: Impossible - Fallout (2018)</a>, <a class="unstyled articleLink" href="/m/blackkkkiansman">
    Blackkkkiansman (2018)</a>, <a class="unstyled articleLink" href="/m/get_out">
    Get Out (2017)</a>, <a class="unstyled articleLink" href="/m/the_irishman">

```

```

    Baby Driver (2017)</a>, <a class="unstyled articleLink" href="/m/spider_man_homecoming">
    Spider-Man: Homecoming (2017)</a>, <a class="unstyled articleLink" href="/m/godfather_part_ii">
    The Godfather, Part II (1974)</a>, <a class="unstyled articleLink" href="/m/the_battle_of_algiers">
    The Battle of Algiers (La Battaglia di Algeri) (1967)</a>]
'https://www.rottentomatoes.com/m/it_happened_one_night', 'https://www.rottentomatoes.com/m/citizen_kane', 'https://www.rottentomatoes.com/m/the_wizard_of_oz_1939',
'https://www.rottentomatoes.com/m/modern_times', 'https://www.rottentomatoes.com/m/black_panther_2018', 'https://www.rottentomatoes.com/m/parasite_2019',
'https://www.rottentomatoes.com/m/avengers_endgame', 'https://www.rottentomatoes.com/m/1003707-casablanca', 'https://www.rottentomatoes.com/m/knives_out',
'https://www.rottentomatoes.com/m/us_2019', 'https://www.rottentomatoes.com/m/toy_story_4', 'https://www.rottentomatoes.com/m/lady_bird', 'https://www.rottentomatoes.com/m/mission_impossible_fallout',
'https://www.rottentomatoes.com/m/mission_impossible_fallout', 'https://www.rottentomatoes.com/m/blackkkkiansman', 'https://www.rottentomatoes.com/m/get_out', 'https://www.rottentomatoes.com/m/the_irishman'
Movie:The Battle of Algiers (La Battaglia di Algeri) (1967)

```