# HW03 - Web Scrapping

```
1 # install gazpacho library
2 !pip install gazpacho
3
4 ## import function
5 from gazpacho import Soup
6 import requests
7
8 # IMDB website
9 url = "https://www.imdb.com/search/title/?groups=top_100&sort=user_rating,desc"
10 html = requests.get(url)
11
12 imdb = Soup(html.text)
```

    Requirement already satisfied: gazpacho in /usr/local/lib/python3.10/dist-packages (1.1)

```
1 # Movie Title
2 titles = imdb.find("h3", {"class": "lister-item-header"})
3 clean_titles = [title.strip() for title in titles]
```

```
1 # Movie Rate Type
2 rate_types = imdb.find("span", {"class": "certificate"})
3 clean_rate_types = [rate_type.strip() for rate_type in rate_types]
```

```
1 # Movie Length
2 lengths = imdb.find("span", {"class": "runtime"})
3 clean_lengths = [length.strip() for length in lengths]
```

```
1 # Movie Genre
2 genres = imdb.find("span", {"class": "genre"})
3 clean_genres = [genre.strip() for genre in genres]
```

```
1 # Movie Rating
2 ratings = imdb.find("div", {"class": "ratings-imdb-rating"})
3 clean_ratings = [float(rating.strip()) for rating in ratings]
```

```
1 # # Movie Release Year
2 years = imdb.find("span", {"class": "lister-item-year"})
3 clean_years = [int(year.strip().replace('(','').replace(')','')) for year in years]
```

```
1 # create DataFrame
2 import pandas as pd
3
4 # movie_database
5 movie_db = pd.DataFrame(data ={
6     "Title": clean_titles,
7     "Rate_type": clean_rate_types,
8     "Length": clean_lengths,
9     "Genre": clean_genres,
10    "Rating": clean_ratings,
11    "Released Year": clean_years
12    })
13
14 movie_db.head()
```

|   | Title | Rate_type | Length | Genre | Rating | Released Year |
|---|-------|-----------|--------|-------|--------|---------------|
| 0 | 1. The Shawshank Redemption (1994) | R | 142 min | Drama | 9.3 | 1994 |
| 1 | 2. The Godfather (1972) | R | 175 min | Crime, Drama | 9.2 | 1972 |
| 2 | 3. The Dark Knight (2008) | PG-13 | 152 min | Action, Crime, Drama | 9.0 | 2008 |
| 3 | 4. Schindler's List (1993) | R | 195 min | Biography, Drama, History | 9.0 | 1993 |
| 4 | 5. The Lord of the Rings: The Return of the Ki... | PG-13 | 201 min | Action, Adventure, Drama | 9.0 | 2003 |

```
1
```