```python
1 # web scraping
2 !pip install gazpacho
3
```

```python
1 from gazpacho import Soup
2 import requests
3 import pandas as pd
4 import numpy as np
```

```python
1 from gazpacho import Soup
2 import requests
3 import pandas as pd
4 import numpy as np
5 url = "https://www.animenewsnetwork.com/encyclopedia/ratings-anime.php?top50=popular&n=100"
6 html = requests.get(url)
7 top_100 = Soup(html.text)
```

```python
1 # Find anime name
2 titles = top_100.find("a", {"href" : "/encyclopedia/anime.php"})
3 titles[0:5]
```

```
    [<a href="/encyclopedia/anime.php">Anime</a>,
     <a href="/encyclopedia/anime.php?id=2960">Fullmetal Alchemist (TV)</a>,
     <a href="/encyclopedia/anime.php?id=6592">Death Note (TV)</a>,
     <a href="/encyclopedia/anime.php?id=13">Cowboy Bebop (TV)</a>,
     <a href="/encyclopedia/anime.php?id=377">Spirited Away (movie)</a>]
```

```python
1 # Clean titles
2 # Delete non anime name
3 clean_titles =[title.strip() for title in titles ]
4 clean_titles.pop(0)
5 clean_titles[0:5]
```

```
    ['Fullmetal Alchemist (TV)',
     'Death Note (TV)',
     'Cowboy Bebop (TV)',
     'Spirited Away (movie)',
     'Princess Mononoke (movie)']
```

```python
1 # Find rating
2 ratings = top_100.find("td", {"class" : "r"})
3 ratings[0:6]
```

```
    [<td class="r">rating</td>,
     <td class="r">nb. votes</td>,
     <td class="r">8.72</td>,
     <td class="r">13523</td>,
     <td class="r">8.86</td>,
     <td class="r">13107</td>]
```

```python
1 # Clean rating
2 clean_ratings = [rating.strip() for rating in ratings]
3 clean_ratings.pop(0)
4 clean_ratings.pop(0)
5 clean_ratings[0:10]
```

```
    ['8.72',
     '13523',
     '8.86',
     '13107',
     '8.93',
     '12444',
     '8.96',
     '10706',
     '8.93',
     '9974']
```

```python
1 # In the list, it contain rating and number vote.
2 # So, we have to extract and create new list for that.
3 # Extract only rating
4 float_numbers = [float(number) for number in clean_ratings]
5 filtered_ratings = [rating for rating in float_numbers if rating <= 10]
6
```

```
7 # Extract only number vote
8 int_numbers = [int(float(number)) for number in clean_ratings]
9 filtered_number_vote = [number for number in int_numbers if number >10]
10
11
```

```
1 df = pd.DataFrame(data = {
2     "titles" : clean_titles,
3     "rating" : filtered_ratings,
4     "number_vote" : filtered_number_vote
5 })
6
7 df.head()
```

|   | titles | rating | number_vote |
|---|---|---|---|
| 0 | Fullmetal Alchemist (TV) | 8.72 | 13523 |
| 1 | Death Note (TV) | 8.86 | 13107 |
| 2 | Cowboy Bebop (TV) | 8.93 | 12444 |
| 3 | Spirited Away (movie) | 8.96 | 10706 |
| 4 | Princess Mononoke (movie) | 8.93 | 9974 |

```
1 # Further more, we want to extract type of anime in parentheses
2 # Extract type into the new column
3 df["Type"] = df["titles"].str.extract(r'\((.*)\)', expand =False)
4
5 # Remove () in tiltes
6 df['titles'] = df['titles'].str.replace(r'\(.*?\)', '').str.strip()
7
8 df.head(10)
```

```
<ipython-input-17-8598cdf67f3e>:6: FutureWarning: The default value of regex will change from True to F
  df['titles'] = df['titles'].str.replace(r'\(.*?\)', '').str.strip()
```

|   | titles | rating | number_vote | Type |
|---|---|---|---|---|
| 0 | Fullmetal Alchemist | 8.72 | 13523 | TV |
| 1 | Death Note | 8.86 | 13107 | TV |
| 2 | Cowboy Bebop | 8.93 | 12444 | TV |
| 3 | Spirited Away | 8.96 | 10706 | movie |
| 4 | Princess Mononoke | 8.93 | 9974 | movie |
| 5 | Melancholy of Haruhi Suzumiya | 8.56 | 10225 | The) Melancholy of Haruhi Suzumiya (TV |
| 6 | Elfen Lied | 8.29 | 10530 | TV |
| 7 | Neon Genesis Evangelion | 8.32 | 10372 | TV |
| 8 | Code Geass: Lelouch of the Rebellion | 8.85 | 9308 | TV |
| 9 | Bleach | 7.94 | 9242 | TV |

```
1 # There are still wrong values in index = 5
2 # Because, the original name is (The) Melancholy of Haruhi Suzumiya (TV)
3 # It was extract by regex
4 # So, just replace it.
5 df.iloc[5,3] = "TV"
6 df.iloc[5,0] = "(The) Melancholy of Haruhi Suzumiya"
7
8 df.head(10)
```

| | titles | rating | number_vote | Type |
|---|---|---|---|---|
| **0** | Fullmetal Alchemist | 8.72 | 13523 | TV |
| **1** | Death Note | 8.86 | 13107 | TV |
| **2** | Cowboy Bebop | 8.93 | 12444 | TV |
| **3** | Spirited Away | 8.96 | 10706 | movie |
| **4** | Princess Mononoke | 8.93 | 9974 | movie |
| **5** | (The) Melancholy of Haruhi Suzumiya | 8.56 | 10225 | TV |
| **6** | Elfen Lied | 8.29 | 10530 | TV |
| **7** | Neon Genesis Evangelion | 8.32 | 10372 | TV |
| **8** | Code Geass: Lelouch of the Rebellion | 8.85 | 9308 | TV |
| **9** | Bleach | 7.94 | 9242 | TV |

Colab paid products  -  Cancel contracts here

✓　0s　completed at 10:10 PM

| | titles | rating | number_vote | Type |
|---|---|---|---|---|