



# Web Scraping W PHP

Katarzyna Mierzwa 269578

Mateusz Sęk 250487

Mateusz Kwiatkowski 274140

# SPIS TREŚCI

Co to jest Web Scraping?

Jak działa Web Scraping?

Wykorzystanie web scrapingu  
do naszego projektu

Jak stworzyć web scraping w  
języku PHP

Technologia

Bibliografia





# Co to jest web scraping?

Web Scraping to metoda, która wyciąga dane ze stron internetowych, co zastępuje odręcznie i powtarzające się kopiowanie i wklejanie danych. Najczęściej wykorzystywany jest:

- do porównywania cen przeróżnych produktów/ usług ([ceneo.pl](https://ceneo.pl)/[skyscanner.pl](https://skyscanner.pl))
- do pobierania danych o potencjalnych klientach pożądanym przez firmę
- do zbierania informacji o produktach/ usługach konkurencji
- do stron, które zbierają informacje z różnych źródeł na dany temat (np. oferty pracy)



# Jak działa Web Scraping?

Nasz kod wysyła żądanie do serwera, na którym znajduje się strona, którą chcemy scrapować. Pobiera on kod źródłowy strony, w takim sam sposób jak robi to przeglądarka. Następnie filtruje stronę i szuka elementów HTML, który wyodrębniliśmy w naszym kodzie, by następnie wyświetlić zawartość, którą chcemy wyodrębnić.

# Index

index.php > ...

```
1 <?php
2     require_once('simple_html_dom.php');//dołączamy bibliotekę php do tego projektu...
3     //biblioteka służy do przeszukiwania kodu html pod kątem wyszukiwanych wartości
4
5     require_once('ScraperItem.php');//dołączamy klasę na podstawie której tworzymy obiekty poszczególnych produktów
6
7     require_once('euroScraper.php');//pliki dołączamy za pomocą którego pobieramy dane z euro agd
8     require_once('mediaExpertScraper.php');//...z media expert
9
10    if(count($_GET) > 0){ //sprawdzamy czy ktoś użył wyszukiwarki
11        $mediaExpertURL = generateURLForMediaExpertScrapping($_GET); // jeżeli tak generujemy adres url media expert do pobrania produktów
12
13        if($mediaExpertURL !== false){
14            $mediaExpertProducts = getMediaExpertHTMLContent($mediaExpertURL); //jeżeli adres media expert jest różny od false
15            //będzie równy false tylko wtedy gdy użytkownik nie wybierze sklepu z wyszukiwarki
16        }
17        else{
18            $mediaExpertProducts = [];
19            //nie będzie żadnych produktów nie będzie niczego
20        }
21
22        $euroURL = generateURLForEuroScrapping($_GET); //ta sama historia
23
24        if($euroURL !== false){
25            $euroProducts = getEuroHTMLContent($euroURL);
26        }
27        else{
28            $euroProducts = [];
29        }
30    }
```

index.php > ...

```
30
31     }
32     else{
33         $euroProducts = getEuroHTMLContent(); //pobieramy domyslne produkty z euro
34         $mediaExpertProducts = getMediaExpertHTMLContent(); //pobieramy z media
35         // tylko wtedy gdy nic nie wybral nie wybral zadnego sklepu
36     }
37
38
39
40     $products = array_merge($euroProducts, $mediaExpertProducts); //laczymy w jedna tablice wszystkie produkty z tablic z euro i media
41     shuffle($products); // miwyszamy tablice zeby wyswietamy sie w lososowej kolejnosci
42 }>
43
44 <!DOCTYPE html>
45 <html lang="en">
46 <head>
47     <meta charset="UTF-8">
48     <meta http-equiv="X-UA-Compatible" content="IE=edge">
49     <meta name="viewport" content="width=device-width, initial-scale=1.0">
50     <title>Laptopy</title>
51     <link rel="stylesheet" href="https://cdn.jsdelivr.net/npm/bootstrap@5.1.3/dist/css/bootstrap.min.css" integrity="sha384-1BmE4kWbQ78iYhF
52 </head>
53 <body>
54 <nav class="navbar navbar-expand-lg navbar-light bg-light">
55     <div class="container-fluid">
56         <a class="navbar-brand" href="#">Wyszukiwanie laptopów</a>
57     </div>
58 </div>
59 </nav>
```

```

index.php > ...
64     <div class="form-group">
65         <label class="mb-2">Sklep</label>
66         <select multiple class="form-control" name="shop[]">
67             <option value="euro" <?php if(isset($_GET['shop'])) && in_array('euro', $_GET['shop'])): ?> selected <?php endif; ?>>Euro RTV AGD</option>
68             <option value="media" <?php if(isset($_GET['shop'])) && in_array('media', $_GET['shop'])): ?> selected <?php endif; ?>>Media Expert</option>
69         </select>
70     </div>
71 </div>
72 <div class="col-3">
73     <div class="form-group">
74         <label class="mb-2">Producent</label>
75         <select class="form-control" name="producent">
76             <option disabled selected>Wybierz</option>
77             <option <?php if(isset($_GET['producent'])) && $_GET['producent'] === 'hp'): ?> selected <?php endif; ?> value="hp">HP</option>
78             <option <?php if(isset($_GET['producent'])) && $_GET['producent'] === 'lenovo'): ?> selected <?php endif; ?> value="lenovo">Lenovo</option>
79             <option <?php if(isset($_GET['producent'])) && $_GET['producent'] === 'acer'): ?> selected <?php endif; ?> value="acer">Acer</option>
80             <option <?php if(isset($_GET['producent'])) && $_GET['producent'] === 'dell'): ?> selected <?php endif; ?> value="dell">Dell</option>
81             <option <?php if(isset($_GET['producent'])) && $_GET['producent'] === 'Apple'): ?> selected <?php endif; ?> value="Apple">Apple</option>
82         </select>
83     </div>
84 </div>
85 <div class="col-3">
86     <div class="form-group">
87         <label class="mb-2">Przekątna ekranu</label>
88         <select class="form-control" name="screen">
89             <option disabled selected>Wybierz</option>
90             <option <?php if(isset($_GET['screen'])) && $_GET['screen'] === '10'): ?> selected <?php endif; ?> value="10">10 cali</option>
91             <option <?php if(isset($_GET['screen'])) && $_GET['screen'] === '11'): ?> selected <?php endif; ?> value="11">11 cali</option>
92             <option <?php if(isset($_GET['screen'])) && $_GET['screen'] === '12'): ?> selected <?php endif; ?> value="12">12 cali</option>
93             <option <?php if(isset($_GET['screen'])) && $_GET['screen'] === '13'): ?> selected <?php endif; ?> value="13">13 cali</option>
94             <option <?php if(isset($_GET['screen'])) && $_GET['screen'] === '14'): ?> selected <?php endif; ?> value="14">14 cali</option>
95             <option <?php if(isset($_GET['screen'])) && $_GET['screen'] === '15'): ?> selected <?php endif; ?> value="15">15 cali</option>
96             <option <?php if(isset($_GET['screen'])) && $_GET['screen'] === '16'): ?> selected <?php endif; ?> value="16">16 cali</option>
97             <option <?php if(isset($_GET['screen'])) && $_GET['screen'] === '17'): ?> selected <?php endif; ?> value="17">17 cali</option>
98         </select>
99     </div>
100 </div>

```



# EuroScrapper

euroScrapper.php > PHP Intelephense > getEuroHTMLContent

```
1  <?php
2
3  function getEuroHTMLContent($url = 'https://www.euro.com.pl/laptopy-i-netbooki.bhtml'){ //funkcja majaca za zadanie pobranie danych ze strony euro i przetworzenie ich do
4      $euroHTML = file_get_html($url); //pobranie danych ze strony euro
5
6      $products = []; //deklarujemy pusta tablice produktow
7
8      foreach($euroHTML->find('.product-box') as $product){ //przechodzimy petlo po wszystkich elemetach html ktora zwrocila funkcja
9          $attributes = []; //pusta tablica w ktorej znajduje sie poszczegone czesci opisu produktow
10
11          if(!count($product->find('.product-name')) || !count($product->find('.product-photo .photo-hover')) || !count($product->find('.price-normal'))){
12              continue; //sprawdzamy czy istnieje nazwa zdjecie i cena , jezeli nie istnieje ktorych z tych elementow to przechodzmy do kolejnej iteracji
13          }
14
15          $title = trim(preg_replace('/\s+/',' ', $product->find('.product-name')[0]->plaintext)); //pobieramy tytul usuwamy z niego spacje przed i po
16          $image = 'https://www.euro.com.pl' . $product->find('.product-photo .photo-hover')[0]->{'data-hover'}; //pobieramy link do zdjecia prduktu
17          $price = trim(preg_replace('/\s+/',' ', $product->find('.price-normal')[0]->plaintext)); //pobieramy cene i usuwamy spacje z pocztaku i konca
18
19          foreach($product->find('.product-attributes .attributes-row') as $attribute){ //iterujemy po poszczegolnych fragmentach opisu produktow
20              array_push($attributes, trim(preg_replace('/\s+/',' ', $attribute->plaintext))); //i wkladamy do tablicy i usuwamy biale znaki
21          }
22
23          $shopUrl = 'https://www.euro.com.pl' . $product->find('.product-name a')[0]->href; //pobieramy link do produkty
24          //dodajemy do tablicy nowy obiekt
25          array_push($products, new ScrapperItem($title, $image, $price, $attributes, $shopUrl, 'https://f00.esfr.pl/img/desktop/euro/logo.png'));
26      }
27
28      return $products; //zwracamy cala tablice produktow
29  }
30
31  function getScreenSize($size){ //funkcja ktora przeksztalca wielkosc ekranu na taka ktora jest akceptowalna przez wyszukiwarke euro
32      $sizes = [
33          '10' => '11',
34          '11' => '11-2-13-1',
35          '12' => '11-2-13-1',
36          '13' => '13-2-14-1',
37      ]
```



```
57
58
59 function generateURLforEuroScrapping($params){ //funkcja ktora generuje adres url
60     if(isset($_GET['shop']) && !in_array('euro', $_GET['shop'])) { //jesli w wyszukiwarce nie zostal wybrany sklep euro to funkcja zwroci false
61         return false;
62     }
63
64     $url = 'https://www.euro.com.pl/laptopy-i-netbooki{producent}{przekatna}{ram}{dysk}{matryca}.bhtml'; //bazowy adres wyszukiwarki
65
66     //sprawdzamy czy w wyszukiwarce zostal wybrany producent jezeli tak to dadajemy go do linku z pozostalymi wartosciami tak samo czyli ram itd
67     if(isset($_GET['producent'])) {
68         $url = str_replace('{producent}', $_GET['producent'], $url);
69     }
70     else {
71         $url = str_replace('{producent}', '', $url);
72     }
73
74     if(isset($_GET['screen'])) {
75         $url = str_replace('{przekatna}', "przekatna-ekranu-cale-" . getScreenSize($_GET['screen']), $url);
76     }
77     else {
78         $url = str_replace('{przekatna}', '', $url);
79     }
80
81     if(isset($_GET['ram'])) {
82         $url = str_replace('{ram}', "pamiec-ram-2!" . $_GET['ram'] . "-gb", $url);
83     }
84     else {
85         $url = str_replace('{ram}', "", $url);
86     }
87     if(isset($_GET['disc'])) {
88         $url = str_replace('{dysk}', "dysk-ssd!" . getDiskCapacity($_GET['disc']), $url);
89     }
90     else {
91         $url = str_replace('{dysk}', "", $url);
92     }
93
94     if(isset($_GET['matrix'])) {
```

# ScrapperItem

ScrapperItem.php > PHP Intelephense > ScrapperItem

```
1  <?php
2
3  class ScrapperItem{
4      public $title;
5      public $image;
6      public $price;
7      public $description;
8      public $shop_url;
9      public $shop;
10
11  public function __construct($title = '', $image = '', $price = '', $description = [], $shop_url = '', $shop = ''){
12      $this->title = $title;
13      $this->image = $image;
14      $this->price = $price;
15      $this->description = $description;
16      $this->shop_url = $shop_url;
17      $this->shop = $shop;
18  }
19 }
```



# Biblioteka Simple HTML DOM Parser

Jest to skrypt który potrafi pobrać zawartość innej strony oraz przefiltrować ją w poszukiwaniu interesującego nas elementu. PHP Simple HTML DOM potrafi dokładnie skopiować wybrany fragment strony. Działaniem przypomina bibliotekę jQuery bo właśnie w podobny sposób poszukujemy wybranych elementów, po klasie, id, po nazwie tagu HTML. Możemy używać skrypt w pętli co daje potężne możliwości analizy całych serwisów.



# Technologia

Stos technologiczny:

- Visual Studio Code
- PHP
- Bootstrap
- Xampp Apache



# Bibliografia

1. [Link](#)
2. [Link](#)
3. [Link](#)
4. [Link](#)

