

# A web scrapper case

Rafael Campos Nunes

April 16, 2019

## Contents

<b>1</b>	<b>Requirements</b>	<b>1</b>
<b>2</b>	<b>Inspection</b>	<b>1</b>

## 1 Requirements

The objective of the web scrapper is to list every product with the following characteristics encoded as a JSON:

1. Name of the product
2. Link to one or more of its images
3. Price of the product

The page from which all of this information will be scrapped is: <https://nerdstore.com.br/categoria/especiais/game-of-throne/s>

## 2 Inspection

When inspecting the page I found that the elements are enclosed in the following structure:

---

```
1 <ul class="products ...">
2   <li class="... product ...">
3     <a class="woocommerce-LoopProduct-link woocommerce-loop-product__link">
4       <img class="" src="">
```

```

5      <img class="" src="">
6      <h2 class="woocommerce-loop-product__title">
7          <span class="ellip">text ellipsed?</span>
8          <span class="ellip-line">text ellipsed 2?</span>
9      </h2>
10     <span class="price">
11         <span class="woocommerce-Price-amount amount">
12             <span class="woocommerce-Price-currencySymbol">R$</span>
13             "49,90"
14         </span>
15     </span>
16 </a>
17 </li>
18 ...
19 </ul>

```

---

By inspection I see that a product may or may not have more than one image and that are specific classes to grab the specified element, though this shouldn't be enough I think because they (the developers) might change on the long run. I should return a JSON that looks somewhat like this:

```

1  o = {
2      'name'   : '',
3      'img'    : ['', ..., ''],
4      'price'  : '44,90'
5  }

```

---

The price comes with a currency tag and that should be important to have but it's not asked by the head hunters.

What I should do now is to take the big block of products that are contained inside the *ul* tag with the *products* class and work from there.