# Coding assignment 2
## Web information extraction and retrieval

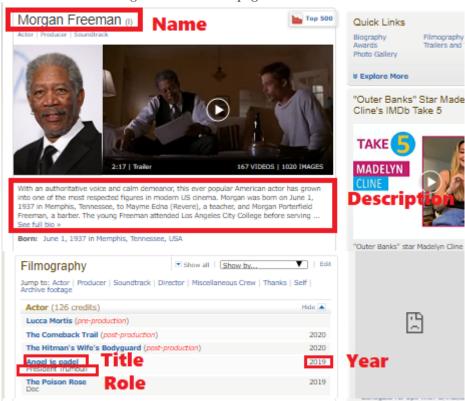Jošt Gombač, Marko Krajinović, Rok Petrovčič Vižintin

# 1   Page selection

Figure 1: imdb.com page data selection



We selected the name and description of the actors. The page also includes a list of films they have worked on (as an actor, producer or other role). For each film we extract its title, year and actors role in the film.

## 2 XPath

### 2.1 rtvslo.si

```
1  author = //div[@class='author-name']/text()
2  published_time = //div[@class='publish-meta']/text()
3  title = //h1/text()
4  subtitle = //div[@class='subtitle']/text()
5  lead = //p[@class='lead']/text()
6  content = //div[@class='article-body']//text()[not(ancestor::figure)
      ]
```

We additionally manipulate date to proper format.

### 2.2 overstock.com

```
1  rows = //table[@border='0' and @cellpadding='0' and @cellspacing='0'
2        and @width='100%']//tr[@bgcolor]
3  for each row do:
4    el = ./td[-1]
5    title = el.xpath(./a/b/text())
6    prices = el.xpath(./table/tbody/tr/td/table/tbody/tr/td//text())
7    list_price = prices[1]
8    price = prices[3]
9    saving = prices[-1].split()[0]
10   saving_percent = prices[-1].split()[-1]
11   content = el.xpath(./table/tbody/tr/td/span/text())
```

### 2.3 imdb.com

```
1  name = //td[@class='name-overview-widget__section']
2        /h1/span[@class='itemprop']/text()
3  description = //div[@class='name-trivia-bio-text']/div[@class='
      inline']
4        //text()[not(ancestor::span)]"
5  films = //div[contains(@class, 'filmo-row')]
6
7  for each film do:
8    title = film.xpath(./b//text())
9    year = film.xpath(./span[@class='year_column']/text())
10   role = film.xpath("./text()")
```

Each film role selection captured different number of elements as this field was not structured well.
We implemented some post-processing checks to insure the right data was captured.

## 3 Regular expressions

### 3.1 rtvslo.si

```
1  author = <div class="author-name">(.*?)</div>
2  published_time = <div class="publish-meta">(.*?)<br>
3  title = <h1>(.*?)</h1>
4  subtitle = <div class="subtitle">(.*?)</div>
5  lead = <p class="lead">(.*?)</p>
6  content = <article class="article">(.*?)</article>
7  remove_figures = <figure .*?>.*?</figure>
8  content = content.replace(remove_figures, "")
```

## 3.2 overstock.com

```
1  rows = <tr bgcolor=.*?>
2    <td valign="top" align="center">
3    <table><tbody><tr><td><a href=.*?><img src=.*? border="0"
4    width="80" height="80"></a></td></tr>
5    <tr><td><a href=.*?>More Info...</a></td></tr>
6    </tbody></table></td><td valign="top">
7    <a href=.*?><b>(.*?)</b></a><br>
8    <table><tbody><tr><td valign="top"><table>
9    <tbody><tr><td align="right" nowrap="nowrap"><b>List Price:</b></
     td>
10   <td align="left" nowrap="nowrap"><s>(.*?)</s></td></tr>
11   <tr><td align="right" nowrap="nowrap"><b>Price:</b></td>
12   <td align="left" nowrap="nowrap"><span
13    class="bigred"><b>(.*?)</b></span></td></tr>
14   <tr><td align="right" nowrap="nowrap"><b>You Save:</b></td>
15   <td align="left" nowrap="nowrap"><span
16    class="littleorange">(.*?)</span></td></tr>
17   </tbody></table>
18   </td><td valign="top"><span class="normal">(.*?)<br>
19   <a href=.*?><span class="tiny"><b>Click here to
20    purchase.</b></span></a></span><br>
21
22 for each row do:
23   construct JSON object
```

## 3.3 imdb.com

```
1  name = <td class="name-overview-widget__section">.*?
2      <h1 class="header"> <span class="itemprop">(.*?)</span>
3  description = <div class="name-trivia-bio-text">.*?<div
4          class="inline">(.*?)<span
5  films = <div class="filmo-row .*?>
6    <span class="year_column">
7     (.*?)
8    </span>
9    <b><a href=.*?
10   >(.*?)</a></b>(.*?)</?div
```

Each captured role was then post-processed to ensure proper data structure.

# 4 Wrapper

## 4.1 Description

The wrapper simultaneously iterates over the DOM of two web pages. The DOM of each page is cleaned before wrapper generation. We remove the following tags: *script*, *style*, *meta*, *link*, *noscript*, *svg*, *path*, *iframe* and *map*. We remove these, because they usually contain data irrelevant to the crawler. The user could remove tags they want to parse from the unwanted tags list if they deemed it necessary. We also remove empty strings and newlines. Each DOM element is then compared. If they match, they are added to the wrapper. If elements match in everything except their text content, that content is marked as interesting. Otherwise, the next match is found, and all unmatched elements in between get marked as optional. After a wrapper is generated, it is checked for repeating elements. Those are marked as such. the output is an almost HTML-like text file with a few differences. Optional elements are enclosed in ()?, repeating elements are enclosed in ()+ and potentially interesting text is replaced with **#Text**.

## 4.2 Pseudocode

```
1   GenerateWrapper(nodeA, nodeB, wrapper)
2     indexA = 0
3     indexB = 0
4     while nodeA.children.lenth > indexA and nodeB.children.length >
        indexB:
5      elementA = nodeA.children[indexA]
6      elementB = nodeB.children[indexB]
7      if elementA and elementB have matching tags and a matching id
        proeprty:
8         if elementA is text:
9           if elementA.textValue == elementB.textValue
10             wrapper.add(new textElement(elementA.textValue))
11           else
12             wrapper.add(new textElement("#text"))
13           indexA++, indexB++
14         else
15           tagElement = new tagElement(elementA.tag, elementA.
        attributes)
16           wrapper.add(tagElement)
17           if elementA.children.length > 0
18             GenerateWrapper(elementA, elementB, tagElement)
19           indexA++, indexB++
20       else:
21         matchingIndexA = FindIndexOfNextMatch(nodeA, nodeB, indexA,
        indexB)
22         matchingIndexB = FindIndexOfNextMatch(nodeB, nodeA, indexB,
        indexA)
23         if matchingIndexA and matchingIndexB are None:
24           add elementA and elementB as optional nodes to wrapper
25         if (matchingIndexA - indexA) < (matchingIndexB - indexB):
```

```
26          add elements in nodeA from indexA to matchingIndexA as
       optional nodes into the wrapper
27          indexA = matchingIndexA
28        else:
29          add elements in nodeB from indexB to matchingIndexB as
       optional nodes into the wrapper
30          indexB = matchingIndexB
31    add remaining elments from the node that still has unprocessed
       children as optional nodes into the wrapper
32    FindRepeatingNodes(wrapper)
33    return wrapper
34
35 FindRepeatingNodes(wrapper)
36    if wrapper has less than 2 children:
37      return
38    current = wrapper.children[0]
39    for child in wrapper.children[1:]:
40      if child and current are equal:
41        mark current as repeating
42        mark child as repeating
43        mark child as removed, so it is not printed when we stringify
       the wrapper
44       FindRepeatingNodes(child)
45      current = child
```

## 4.3   Page wrapper results

Due to the relatively naive implementation of our wrapper generation algorithm and wrapper pretty-printing, the wrappers tend to be rather long as they are similar to the actual HTML of the web page. Included are only parts of each wrapper which showcase the optional, list or text elements. The full outputs are present in the repository folder *results/C*.

### 4.3.1 rtvslo.si

```
1  (<div class="iAdserver adloaded" data-iadserver-zone="299" id="
      f4d0659e109b98fa56aa9d5c2e67bd1f" style="text-align: center;
      position: absolute;"></div>)+
2    <div id="main-container" class="container article-container" data-
        id="475392">
3     <div class="section-heading blue">
4       <h3 class="section-title animated-circles-onhover">
5         <a href="https://www.rtvslo.si/zabava/avtomobilnost/testi
            /">#Text</a>
6       </h3>
7     </div>
8     <div id="article-edit-btn" class="edit-btn-container" style="
        display:none">
9      <a href="https://admin.rtvslo.si/?c_mod=newsadmin&op=news&func
          =edit&id=475392" class="edit-btn" rel="nofollow" target="
          _blank">Uredi</a>
10    </div>
11    <div class="news-container blue article-old article-type-1">
12     <div class="row">
13       <header class="article-header">
14        <div class="section-heading blue">
15          <h3 class="section-title animated-circles-onhover">
16            <a href="https://www.rtvslo.si/zabava/avtomobilnost/
                testi/">#Text</a>
17          </h3>
18          <a class="comments-icon ml-auto" data-comments-number
                ="13" href="https://www.rtvslo.si/zabava/
                avtomobilnost/testi/audi-a6-50-tdi-quattro-nemir-v-
                premijskem-razredu/475392#article-comments-anchor"></
                a>
19        </div>
20        <h1>#Text</h1>
21        <div class="subtitle">#Text</div>
```

### 4.3.2 overstock.com

```
1  (<tr>
2    <td colspan="2" height="4">
3      <img src="jewelry01_files/ms.gif" width="1" height="4" border
      ="0"></img>
4    </td>
5  </tr>)?
6  (<tr bgcolor="#dddddd">
7    <td valign="top" align="center">
8      <table>
9        <tbody>
10         <tr>
11           <td>
12             <a href="http://www.overstock.com/cgi-bin/d2.cgi?PAGE=
      PROFRAME&PROD_ID=51595">
13               <img src="jewelry01_files/T917417.jpg" border="0"
      width="80" height="80"></img>
14             </a>
15           </td>
16         </tr>
17         <tr>
18           <td>
19             <a href="http://www.overstock.com/cgi-bin/d2.cgi?PAGE=
      PROFRAME&PROD_ID=51595">More Info...</a>
20           </td>
21         </tr>
22       </tbody>
23     </table>
24   </td>
25   <td valign="top">
26     <a href="http://www.overstock.com/cgi-bin/d2.cgi?PAGE=PROFRAME&
      PROD_ID=51595">
27       <b>14-kt. Diamond Cross Pendant (.06 TW)</b>
28     </a>
29     <br>
30     <table>
31       <tbody>
32         <tr>
33           <td valign="top">
34             <table>
35               <tbody>
36                 <tr>
37                   <td align="right" nowrap="nowrap">
38                     <b>List Price:</b>
39                   </td>
```

### 4.3.3 imdb.com

```
1  (<div class="filmo-row even" id="archive_footage-tt0480949">
2    <span class="year_column"> 2005 </span>
3    <b>
4      <a href="/title/tt0480949/?ref_=nm_flmg_arf_20">Venecia 2005:
          Cr nica de Carlos Boyero</a>
5    </b>
6    (TV Short)
7    <br>
8    Josef
9  </div>)?
10 (<div class="filmo-row odd" id="archive_footage-tt0410407">
11   <span class="year_column"> 2003 </span>
12   <b>
13     <a href="/title/tt0410407/?ref_=nm_flmg_arf_21">Orwell Rolls in
          His Grave</a>
14   </b>
15   (Documentary)
16   <br>
17   Self
18 </div>)?
19 (<div class="filmo-row even" id="archive_footage-tt0279641">
20   <span class="year_column"> 2002 </span>
21   <b>
22     <a href="/title/tt0279641/?ref_=nm_flmg_arf_22">Who Is Alan
          Smithee?</a>
23   </b>
24   (TV Movie documentary)
25   <br>
26   Self (uncredited)
27 </div>)?
28 (<div class="filmo-row odd" id="archive_footage-tt0296727">
29   <span class="year_column"> 2000 </span>
30   <b>
31     <a href="/title/tt0296727/?ref_=nm_flmg_arf_23">Lord Stanley's
          Cup: Hockey's Ultimate Prize</a>
32   </b>
33   (Video documentary)
34   <br>
35   Self
36 </div>)?
```