

# Ekstrakcija strukturiranih podatkov s spleta

Seminarska naloga pri predmetu Iskanje in ekstrakcija podatkov s spleta

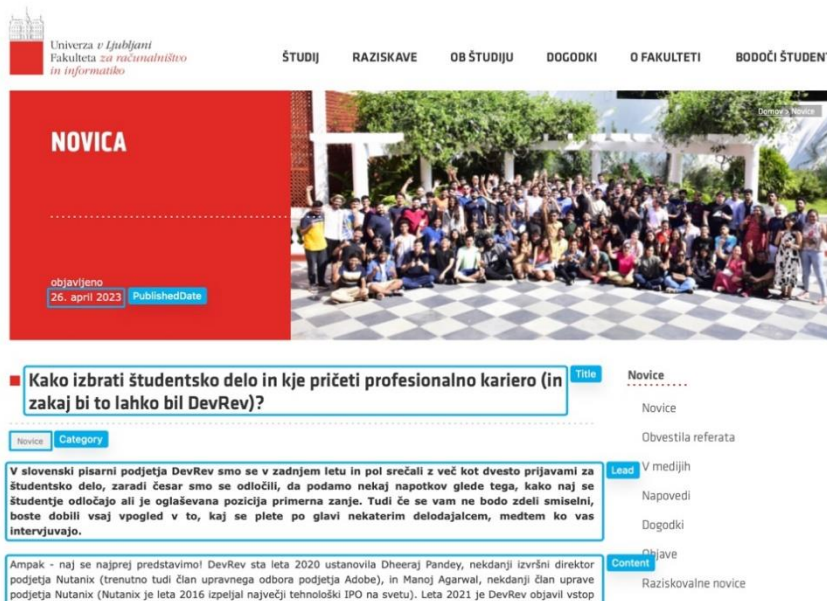
MENTOR: doc. dr. Slavko Žitnik

## 1. Uvod

Seminarska naloga opisuje implementacijo treh različnih algoritmov za ekstrakcijo strukturiranih podatkov s spleta.

## 2. Izbrane spletne strani

Kot lasten primer smo izbrali dve novici iz spletne strani univerze. Iz strani lahko razberemo naslov, datum objave, kategorijo, uvod in vsebino.



## 3. Regex

Overstock.com:

- **Title:** /<td valign="top">[\n\s]\*?<a\s.\*?><b>(.\*?)</b></a>/
- **List price:** /<td align="left"\s.\*?>[\n\s]\*?<s>(.\*?)</s>/
- **Price:** /<td align="left"\s.\*?>[\n\s]\*?<span class="bigred"><b>(.\*?)</b></span>/
- **Saving and saving percent:** /<td align="left"\s.\*?>[\n\s]\*?<span class="littleorange">(\\$.\*?)\s\((.\*?)\)</span>/
- **Content:** /<span class="normal">([\S\s]\*?)[\s]\*<br>/

Rtv.si:

- **Author:** /<div class="author-name">(.\*?)</div>/
- **Published time:** /<div class="publish-meta">[\n\s\t]\*(.\*?)<br>/

- **Title:** /<h1>(.\*?)</h1>/
- **Subtitle:** /<div class="subtitle">(.\*?)</div>/
- **Lead:** /<p class="lead">(.\*?)</p>/
- **Content:** /^(?:.\*)<div class="article-body">|<div class="gallery">(?:.\*)\$|<script\b[^\<]\*(?:(!\</script><[^\<]\*)\*\</script>|<[^\>]\*)\*>/s

#### Fri.uni-lj.si:

- **Published date:** /<div class="text">[\n\s\S]\*?<br>(.\*?)[\n\s]\*?</div>/
- **Title:** /<div class="heading-article\snovica-title">[\n\s]\*<ul>[\n\s]\*<li>s.\*?>(.\*</li>/
- **Category:** /<span class="kategorija">(.\*?)</span>/
- **Lead:** /<div class="novica-content">[\s\S]\*?<p\s?(?:class="rtejustify")?>(?:<.\*?>)\*(.\*?)(?:<.\*?>)\*</p>/
- **Content:** /^(?:.\*)<div class="novica-content">[\n\s\S]\*?<br>[\n\s]\*|</div>(?:.\*)\$|<[^\>]\*>|.\\n{2,}/s

## 4. XPath

#### Overstock.com:

XPath za korenski element seznama izdelkov:

/html/body/table[2]/tbody/tr[1]/td[5]/table/tbody/tr[2]/td/table/tbody/tr/td/table/tbody/tr[@bgcolor]

- **Title:** ./td[2]/a/b/text()
- **List price:** ./td[2]/table/tbody/tr/td[1]/table/tbody/tr[1]/td[2]/s/text()
- **Price:** ./td[2]/table/tbody/tr/td[1]/table/tbody/tr[2]/td[2]/span/b/text()
- **Saving and saving percent:** ./td[2]/table/tbody/tr/td[1]/table/tbody/tr[3]/td[2]/span/text()
- **Content:** ./td[2]/table/tbody/tr/td[2]/span/text()

#### Rtv.si:

- **Author:** //\*[@id="main-container"]/div[3]/div/div[1]/div[1]/div
- **Published time:** //\*[@id="main-container"]/div[3]/div/div[1]/div[2]/text()[1]
- **Title:** //\*[@id="main-container"]/div[3]/div/header/h1
- **Subtitle:** //\*[@id="main-container"]/div[3]/div/header/div[2]
- **Lead:** //\*[@id="main-container"]/div[3]/div/header/p
- **Content:** //\*[@id="main-container"]/div[3]/div/div[2]/descendant::\*[not(name() = "script")]/text()

#### Fri.uni-lj.si:

- **Published date:** //\*[@id="banner-header"]/div[2]/div/div[3]/text()[2]
- **Title:** //\*[@id="katedre-container"]/div[2]/span
- **Category:** //\*[@id="katedre-container"]/div[1]/div[1]/ul/li
- **Lead:** //\*[@id="katedre-container"]/div[3]/p[1]/strong/span/span/span **ali** //\*[@id="katedre-container"]/div[3]/p[1]/strong
- **Content:** //\*[@id="katedre-container"]/div[3]/\*[position() > 2]

## 5. RoadRunner

Pri implementaciji algoritma RoadRunner smo si pomagali s predlaganim člankom [1]. Pred zagonom algoritma smo iz vhodnih HTML dokumentov odstranili odvečne značke in komentarje, ter iz HTML elementov ustvarili žetone. Generirane ovojnice vseh treh strani so smiselne. Delovanje smo preverili tudi z dvema avtomatskima testoma, ki primerjata naš rezultat s primeroma iz članka.

### 5.1 Psevdokoda delovanja algoritma

If the length of both pages is the same or one of the indexes is greater or equal length of the page:

Return current wrapper

If the first tokens of both pages are equal:

Add it to the wrapper, increment both indexes and recursively call the function again

If both tokens are database field:

Add "#PCDATA" to the wrapper, increment both indexes and recursively call the function again

Check for iterators:

Fix the wrapper, add iterator tags , increment page index and recursively call the function again

Check for optionals:

Add optional tags , increment page index and recursively call the function again

Return empty wrapper

## 5.2 Ovojnica Overstock

```
<html><head><title>Overstock.com, save up to 80% every day!</title><link></link>  
<link></link></head><body><input></input><table><tbody><tr><td></td><td><table>  
<tbody><tr><td><table><tbody><tr><td><a><img></a></td><td><img></td><td><a>  
<img></a></td></tr></tbody></table><map><area></area><area></area><area></area>  
<area></area><area></area></map></td></tr><tr><td><a><img></a><a><img></a>  
<a><img></a><a><img></a><a><img></a><a><img></a><a><img></a></td></tr>  
</tbody></table></td></tr><tr><td><img></td></tr><form></form><tr><td><table>  
<tbody><tr><td><table><tbody><tr><td><span>Search:</span><select><option>All  
Stores</option><option>Home & Garden</option><option>Electronics & Computers  
</option><option>Books, Movies, CDs, Games</option><option>Jewelry, Watches & Gifts  
</option><option>Sports, Travel & Toys</option><option>Worldstock</option><option>  
Apparel, Shoes & Access.</option></select></td><td></td><td><input></input></td>  
<td><input></input></td></tr></tbody></table></td><td><img><br></br><a><img>  
</a></td><td><img></td></tr></tbody></table></td></tr><tr><td><img></td></tr>  
</tbody></table><table><tbody><tr><td><img></td><td><br></br><table><tbody>  
<tr><td><table><tbody>( <tr><td></td><td><a>#PCDATA</a></td></tr>)+ </tbody></table>  
</td></tr></tbody></table><br></br><span><b>Stores</b></span><br></br><table>  
<tbody><tr><td><a>Apparel, Shoes & Access.</a></td></tr><tr><td><a>Books,  
Movies, CDs, Games</a></td></tr><tr><td><a>Electronics & Computers</a></td></tr>  
<tr><td><a>Home & Garden</a></td></tr><tr><td><a><b>Jewelry, Watches & Gifts  
</b></a></td></tr><tr><td><a>Sports, Travel & Toys</a></td></tr><tr><td><a>  
Worldstock</a></td></tr></tbody></table><br></br><span><b>New Stock</b></span>  
<br></br><table><tbody><tr><td><a>Ralph Lauren $29.95</a></td></tr><tr><td><a>  
Ben Sherman 53% off</a></td></tr><tr><td><a>Pre-order Harry Potter DVD</a></td>  
</tr><tr><td><a>HP 2GHz System $499</a></td></tr><tr><td><a>New Items within 7  
Days</a></td></tr></tbody></table><br></br><span><b>Customer Service</b></span>  
<br></br><table><tbody><tr><td><a>Shopping Cart & Checkout</a></td></tr><tr>  
<td><a>Track Your Order</a></td></tr><tr><td><a>Your Account</a></td></tr><tr>  
<td><a>Help & FAQ</a></td></tr><tr><td><a>Best Price Guarantee</a></td></tr>  
</tbody></table><br></br><span><b>About Us</b></span><br></br><table><tbody>  
<tr><td><a>About Us</a></td></tr><tr><td><a>Privacy & Security</a></td></tr>  
<tr><td><a>Terms & Conditions</a></td></tr><tr><td><a>Become An Affiliate</a>  
</td></tr><tr><td><a>Business Purchases</a></td></tr><tr><td><a>Have Products  
to Sell?</a></td></tr><tr><td><a>Investor Relations</a></td></tr></tbody>  
</table><img><br></br></td><td><img></td><td><img></td><td><br></br><b><a>  
Jewelry, Watches & Gifts</a>><a>Jewelry</a>><a>#PCDATA</a>>View All</b><br>  
</br><table><tbody><tr><td><br></br><table><tbody><tr><td><table><tbody>  
(<tr><td><span><b>#PCDATA</b></span></td><td><b>#PCDATA</b>#PCDATA(<b>30<b>of<b>296<b>|>)?  
<a>#PCDATA</a>|<a>#PCDATA</a>|<a>#PCDATA</a>(|<a>Quantity</a>|<a>Markdowns</a>)?)  
</td></tr>)+ </tbody></table></td></tr></tbody></table></td></tr><tr><td><table>  
<tbody><tr><td><table><tbody><tr>( <td><span>#PCDATA(<br><br></span><b>Click here  
to purchase.<b>)</span><a><span><br><br></span></td>+ </tr></tbody></table></td>  
</tr><tr><td><img></td></tr><tr><td><table><tbody><tr><td><a><img></a></td>  
</tr><tr><td><a>More Info...</a></td></tr></tbody></table></td><td><a><b>  
#PCDATA</b></a><br></br><table><tbody><tr><td><table><tbody><tr><td><b>List  
Price:</b></td><td><s> #PCDATA </s></td></tr><tr><td><b>Price:</b></td><td>  
<span><b> #PCDATA</b></span></td></tr><tr><td><b>You Save:</b></td><td><span>  
#PCDATA</span></td></tr></tbody></table></td><td><span> #PCDATA<br></br><a>  
<span><b>Click here to purchase.</b></span><a></span><br></br></td></tr>  
</tbody></table></td></tr><tr><td><img></td></tr><tr><td><table><tbody><tr>  
<td><a><img></a></td></tr><tr>( <td><span>#PCDATA(<br><br></span><b>Click here to  
purchase.<b>)</span><a><span><br><br></span></td>+ </tr></tbody></table></td></tr>  
<tr><td><img></td></tr><tr><td><table><tbody><tr><td><a><img></a></td></tr>  
<tr>()+ </tr><tr>()+ </tr><tr>()+ </tr></tbody></table>  
(<td></tr></tbody></table/><td></tr></tbody></table/><td></tr></tbody></table/><td></tr></tbody></table/>  
<td></tr></tbody></table/><td></tr></tbody></table/><td></tr></tbody></table/>  
<td></tr></tbody></table/><td></tr></tbody></table/><td></tr></tbody></table/>  
<td></tr></tbody></table/><td></tr></tbody></table/><td></tr></tbody></table/><td></tr></tbody></table/  
<td></tr></tbody></table/><td></tr></tbody></table/> )? <map><area></area><area></area><area>
```



[illegible]

## 5 **Literatura**

- [1] Crescenzi, Valter & Mecca, Giansalvatore & Merialdo, Paolo. (2001). RoadRunner: Towards Automatic Data Extraction from Large Web Sites. <https://www.vldb.org/conf/2001/P109.pdf>