

# Regularni izrazi in Xpath - Iskanje in ekstrakcija podatkov iz spleta

Robert Košir, Aljoša Omejc

May 1, 2020

## 1 Uvod

V naslednjih nekaj poglavjih bomo na kratko opisali projekt pri predmetu Iskanje in ekstrakcija podatkov s spleta. Sam projekt se ukvarja z zajemom podatkov iz različnih strani z uporabo regularnih izrazov, xpathi in avtomatičnim algoritmom. Ob zagonu programa podate arugemnt, ki določi, katero metodo bo program uporabil.

## 2 Definicija spletne strani slovenskenovice.si

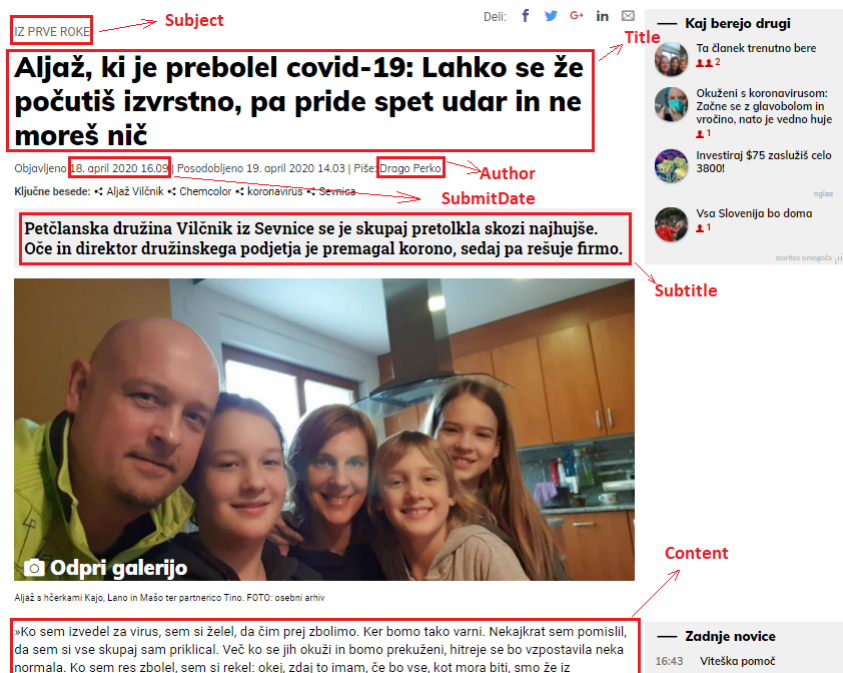


Figure 1: Definicija objektov za spletno stran slovenske novice

### 3 Regularni izrazi

V naslednjih podpoglavjih bom v obliki tabele podal regularne izraze za vsako stran posebj.

#### 3.1 rtvslo.si

|               |  |
|---------------|--|
| Author        | <div class="author-name">(.)</div>                                       |
| PublishedTime | <div class="publish-meta">[\n\s]*(.)<br>                                 |
| Title         | <header class="article-header">(.\n)*<h1>(.)</h1>                        |
| SubTitle      | <header class="article-header">(.\n)*<br><div class="subtitle">(.)</div> |
| Lead          | <header class="article-header">(.\n)*<br><p class="lead">[\n\s]*(.)</p>  |
| Content       | <div class="article-body">(\s \n .)*?<br><div class="article-column">    |

Sam content nam je v tem primeru vrnil cel html niz, kar pa smo dalje procesirali na naslednji način:

- Birsanje <script> tagov in njene vsebine  
re.sub(r"<script([\S\s]\*?)>([\S\s]\*?)</script>", "", content)
- Brisanje vseh HTML tagov  
re.sub(r"<[^>]\*>", "", content)
- Brisanje vseh \s  
re.sub(r"\s+", " ", content)

#### 3.2 overstock.com

|               |   |
|---------------|---|
| Title         | <a.*<b>(.)</b>\n*</a><br>   |
| ListPrice     | <td align="left" nowrap="nowrap"><s>(.)</s></td>  |
| Price         | <span class="bigred"><b>(.)</b></span>  |
| Saving        | <span class="littleorange">(.) \(.*)</span>   |
| SavingPercent | <span class="littleorange">\\$.* (.)</span>   |
| Content       | <td valign="top"><span class="normal">([\s\S]*?)<br><br><a href.*><span class="tiny"><b>(.)</b> |

#### 3.3 slovenskenovice.si

|            |   |
|------------|---|
| Title      | <h1 class="itemTitle">\s*([\s\S]*)</h1>                           |
| Subject    | <span class="itemSuperscript">(.)</span>                          |
| SubmitDate | <span class="itemDatePublished">\s*0bjavljeno (.)</span>          |
| Author     | <span class="itemAuthor">\s*Piše:\s+<span>(.)</span>              |
| Subtitle   | <h2 class="itemSubtitle">\s*<span>(.)</span>                      |
| Content    | <div class="itemFullText" .*([\s\S]*)<div class="itemInfoboxText" |

## 4 Xpath izrazi

V naslednjih podpoglavjih bom v obliki tabele podal xpath izraze, ki so iz vseh strani prebrale vnaprej določene objekte.

### 4.1 rtvslo.si

|               |   |
|---------------|---|
| Author        | <code>//*[@id="main-container"]/div[3]/div/div[1]/div[1]/div</code>                                     |
| PublishedTime | <code>//*[@id="main-container"]/div[3]/div/div[1]/div[2]/text()[1]</code>                               |
| Title         | <code>//*[@id="main-container"]/div[3]/div/header/h1</code>   |
| SubTitle      | <code>//*[@id="main-container"]/div[3]/div/header/div[2]</code>   |
| Lead          | <code>//*[@id="main-container"]/div[3]/div/header/p</code>  |
| Content       | <code>//*[@id="main-container"]/div[3]/div/div[2]</code><br><code>//*[@not(self::script)]/text()</code> |

### 4.2 overstock.com

|                 |  |
|-----------------|--|
| Title           | <code>/html/body/table[2]/tbody/tr[1]/td[5]/table/tbody/tr[2]/td</code><br><code>/table/tbody/tr/td/table/tbody/tr/td/a/b</code>   |
| ListPrice       | <code>/html/body/table[2]/tbody/tr[1]/td[5]/table/tbody/tr[2]/td</code><br><code>/table/tbody/tr/td/table/tbody/tr/td[2]/table/tbody/tr/td[1]</code><br><code>/table/tbody/tr[1]/td[2]/s</code>      |
| Price           | <code>/html/body/table[2]/tbody/tr[1]/td[5]/table/tbody/tr[2]/td</code><br><code>/table/tbody/tr/td/table/tbody/tr/td[2]/table/tbody/tr/td[1]</code><br><code>/table/tbody/tr[2]/td[2]/span/b</code> |
| Saving(Percent) | <code>/html/body/table[2]/tbody/tr[1]/td[5]/table/tbody/tr[2]/td</code><br><code>/table/tbody/tr/td/table/tbody/tr/td[2]/table/tbody/tr/td[1]</code><br><code>/table/tbody/tr[3]/td[2]/span</code>   |
| Content         | <code>/html/body/table[2]/tbody/tr[1]/td[5]/table/tbody/tr[2]/td</code><br><code>/table/tbody/tr/td/table/tbody/tr/td[2]/table/tbody/tr/td[2]/span</code>  |

Tu smo z uporabo enega xpatha pridobili Saving in SavingPercent, ki pa smo ju kasneje razcepili.

### 4.3 slovenskenovice.si

|            |   |
|------------|---|
| Title      | //*[@id="ocmContainer"]/div[1]/div/div[1]/div[1]/div/div[3]/h1                                  |
| Subject    | //*[@id="ocmContainer"]/div[1]/div/div[1]/div[1]/div/div[2]/span                                |
| SubmitDate | //*[@id="ocmContainer"]/div[1]/div/div[1]/div[1]/div/div[4]/span<br>/span[1]                    |
| Author     | //*[@id="ocmContainer"]/div[1]/div/div[1]/div[1]/div/div[4]/span<br>/span[5]/span               |
| Subtitle   | //*[@id="ocmContainer"]/div[1]/div/div[1]/div[1]/div/div[6]/h2/span                             |
| Content    | //*[@id="ocmContainer"]/div[1]/div/div[2]/div[1]/div[1]/div[1]<br>//text()[not(parent::script)] |

## 5 Automatska ekstrakcija

### 5.1 Psevdokoda

---

#### Algorithm 1: Get Page Wrapper

---

```

Result: wrapper
page1,page2 = removeNewLinesOnContent();
diffLines = diffLib.compare(page1, page2);
pairs = GetPairs(diffLines) ;
pairsAttr = ShortenAttributes(pairs);
result = replaceDynamicText(pairsAttr, page1);
soup = BeautifulSoup(result);
RemoveUselessTags(soup);
RemoveUselessAttributes(soup);

```

---

### 5.2 Description

#### 5.2.1 DiffLib.compare()

s primerjanjem dobimo skupaj združeni datoteki, kjer se vrstice z začnejo z [" ", "- ", "+ ", "?"]

- "-" line unikatna vrstica datotke 1
- "+" line unikatna vrstica datotke 2
- " " line enak v obeh
- "?" šum

#### 5.2.2 GetPairs()

Pare pridobimo tako, da gledamo zaporedne vrstice izbiramo po metrikah: - enaka začetna oznaka - dovolj podobni atributi

#### 5.2.3 ShortenAttributes()

Skrajšamo vse attribute tako, da jih primerjamo in vzamemo začetni del, kjer sta enaka.

```
Vhod:
    page1: <h4 class="block-title blue">
    page2: <h4 class="block-title green">
Izhod:
    <h4 class="(block-title .*)"
```

#### 5.2.4 ReplaceDynamicText()

Za vsak par pogledava, če ima različno vsebino in v primeru, da ima potem delava podobno, kot pri roadrunnerju, da zamenjava vsebino z oznako "#text"

```
Vhod:
    -<div class="subtitle">Test novega modela</div>
    +<div class="subtitle">Test nove generacije</div>
Izhod:
    <div class="subtitle">#text</div>
```

#### 5.2.5 RemoveUselessTags()

Odstranila sva oznake [iframe, img, br, footer, nav, script], ker se nama je zdelo, da večinoma ne držijo nobenih uporabnih informacij.

Obdrži vse oznake, ki same ali pa njihovi nasledniki vsebujejo #text. Na ta način dobimo drevo, kjer listi vsebujejo koristne informacije.

#### 5.2.6 RemoveUselessAttributs()

Obdržimo samo atributa class in id, ki pomagata navigirati po drevesu. V primeru, da bi imela več prostora bi obdržala tudi nekatere druge kot so title, src.

### 5.3 Opcijske značke

Za implementacijo opsijskih značk nama je zmanjkalo časa. Vendar bi jo implemenitrjala tako, da vzameva vse različne vrstice iz datoteke1 in datoteke2, ki nimajo svojega para.

## 6 Zaključek

Sama seminarska naloga je bila zelo zanimiva. Njaveč problem nama je povzročala 3 naloga, a vendar mislim, da sva uspešno opravila večino dela.