

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

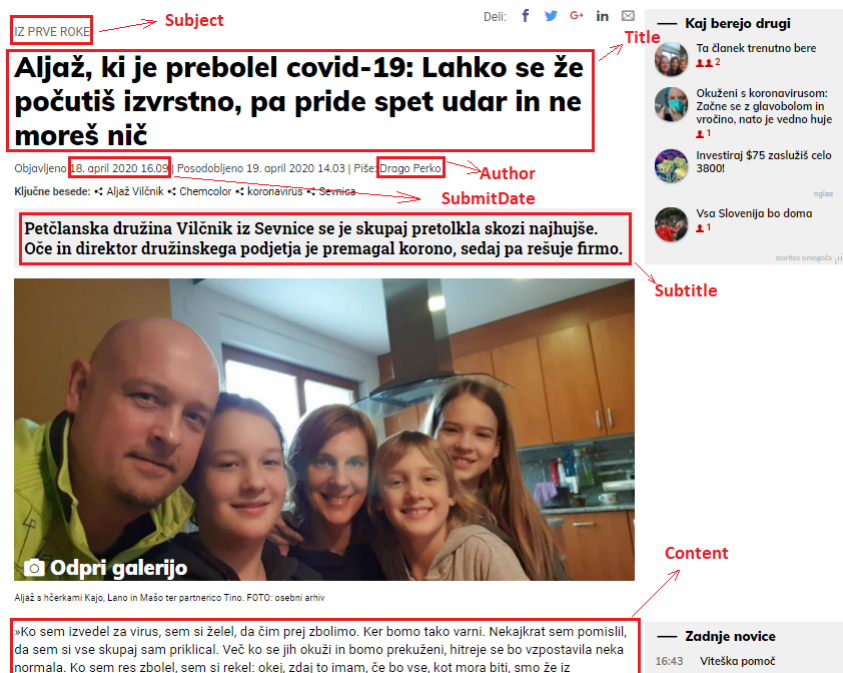


Figure 1: Definicija objektov za spletno stran slovenske novice

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]*(.)
Title	<header class="article-header">(.\n)*<h1>(.)</h1>
SubTitle	<header class="article-header">(.\n)* <div class="subtitle">(.)</div>
Lead	<header class="article-header">(.\n)* <p class="lead">[\n\s]*(.)</p>
Content	<div class="article-body">(\s \n .)*? <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.*(.)\n*
ListPrice	<td align="left" nowrap="nowrap"><s>(.)</s></td>
Price	(.)
Saving	(.) \(.*)
SavingPercent	\\$.* (.)
Content	<td valign="top">([\s\S]*?) <a href.*>(.)

3.3 slovenskenovice.si

Title	<h1 class="itemTitle">\s*([\s\S]*)</h1>
Subject	(.)
SubmitDate	\s*0bjavljeno (.)
Author	\s*Piše:\s+(.)
Subtitle	<h2 class="itemSubtitle">\s*(.)
Content	<div class="itemFullText" .*(.)</div><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> <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> <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> <code>/table/tbody/tr/td/table/tbody/tr/td[2]/table/tbody/tr/td[1]</code> <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> <code>/table/tbody/tr/td/table/tbody/tr/td[2]/table/tbody/tr/td[1]</code> <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> <code>/table/tbody/tr/td/table/tbody/tr/td[2]/table/tbody/tr/td[1]</code> <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> <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 /span[1]
Author	//*[@id="ocmContainer"]/div[1]/div/div[1]/div[1]/div/div[4]/span /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] //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);
RemoveUselessAttributs(soup);

```

5.2 Description

5.2.1 Difflib.compare()

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

- "-" line unikatna vrstica daototke 1
- "+" line unikatna vrstica daototke 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.

5.4 Izhod

5.4.1 Rtv1

```
<html><body class="first-page _www hide-submenu article user-logged-out"><div class=""
id="main-bar"><div class="top-container"><div class="row"><div class="col-12 d-flex
align-items-center"><ul id="main-menu"><li class="green"><a>#text</a></li></ul></div></div></div><div class="container article-container" id="main-container"><div class="section-heading
blue"><h3 class="section-title animated-circles-onhover"><a>#text</a></h3></div><div
class="edit-btn-container" id="article-edit-btn"><a class="edit-btn">#text</a></div>
<div class="news-container blue article-old article-type-1"><div class="row"><header
class="article-header"><div class="section-heading blue"><h3 class="section-title
animated-circles-onhover"><a>#text</a></h3></div><h1>#text</h1><div class="subtitle">#text</div>
<div class="article-meta-mobile"><div class="author-timestamp"><strong>#text</strong>
</div></div><div class="swiper-container swiper-header-gallery mobile-show-only
swiper-container-horizontal swiper-container-wp8-horizontal"><div class="swiper-wrapper"><figure cla
ss="photoswipe swiper-slide mobile-show-only"><figcaption>#text</figcaption></figure></div></div><p
class="lead">#text</p></header><div class="article-meta"><figure class="c-figure-emphasis"><ul
class="emphasis-list"><li>#text</li></ul></figure></div><div class="ar
ticle-body"><article class="article"><div class="gallery"><div class="swiper-container swiper-footer-gallery
swiper-container-horizontal swiper-container-wp8-horizontal"><div class="swiper-wrapper"><figure
class="photoswipe swiper-slide"><figcaption>#text</figcaption>
on</figure><figure class="photoswipe swiper-slide"><figcaption>#text</figcaption></figure><figure
class="photoswipe swiper-slide"><figcaption>#text</figcaption></figure><figure class="photoswipe
swiper-slide"><figcaption>#text</figcaption></figure><figure class="p
hotoswipe swiper-slide"><figcaption>#text</figcaption></figure><figure class="photoswipe
swiper-slide"><figcaption>#text</figcaption></figure><figure class="photoswipe
swiper-slide"><figcaption>#text</figcaption></figure><figure class="photoswipe
swiper-slide"><figcaption>#text</figcaption></figure><figure class="photoswipe
swiper-slide"><figcaption>#text</figcaption></figure><figure class="photoswipe
swiper-slide"><figcaption>#text</figcaption></figure><figure class="photoswipe
swiper-slide"><figcaption>#text</figcaption></figure><figure class="photoswipe
swiper-slide"><figcaption>#text</figcaption></figure><figure class="photos
wipe swiper-slide swiper-slide-prev"><figcaption>#text</figcaption></figure><figure class="photoswipe
swiper-slide swiper-slide-active"><figcaption>#text</figcaption></figure><figure class="photoswipe
swiper-slide swiper-slide-next"><figcaption>#text</figcaption></figure></div></div><div
class="gallery-bottom-thumbs"><figure class="photoswipe
gallery-bottom-thumb"><figcaption>#text</figcaption></figure><figure class="photoswipe
gallery-bottom-thum
more-photos"><figcaption>#text</figcaption></figure></div></div></article></div><div clas
s="article-column"><div class="sticky-block"><div class="right-block similar-articles"><h4 class="(block-title
.*)">#text</h4><div class="article-container"><div class="row"><div class="col-md-7 title"><a
class="list-title">#text</a></div></div></div><div class="ar
ticle-container"><div class="row"><div class="col-md-7 title"><a
class="list-title">#text</a></div></div></div></div></div><div class="article-footer"><div
```

```

more-photos"><figcaption>#text</figcaption></figure></div></div></article></div><div class="
s="article-column"><div class="sticky-block"><div class="right-block similar-articles"><h4 class="(block-title
.*)>#text</h4><div class="article-container"><div class="row"><div class="col-md-7 title"><a
class="list-title">#text</a></div></div></div><div class="ar
ticle-container"><div class="row"><div class="col-md-7 title"><a
class="list-title">#text</a></div></div></div></div></div><div class="article-footer"><div
class="article-tags"><a class="tag tag-red-dark">#text</a><a class="tag tag-red-dark">#text</a></div><d
iv class="share"><a class="error-report">#text</a></div></div></div></div><div class="news-block
blue"><div class="row"><div class="col-lg-3 col-md-6"><div class="md-news"><h3><span
class="news-cat"><a>#text</a></span></h3></div></div><div class="col-lg-3 col-md-6"
"><div class="md-news"><h3><span class="news-cat"><a>#text</a></span></h3></div></div><div
class="col-lg-3 col-md-6"><div class="md-news"><h3><span
class="news-cat"><a>#text</a></span></h3></div></div></div></div><div
id="article-comments-anchor"><div class="section
-heading blue"><h3 class="section-title animated-circles-onhover"><a
class="btn-show-comments">#text</a></h3></div><div class="edit-btn-container"
id="article-comments-edit-btn"><a class="edit-btn">#text</a></div><div class="news-block
article-comments article-comm
ents-toggle hide-comments"><p
class="hidden-comments-notice">#text</p></div></div></div></body></html>

```


5.4.2 Overstock1

[illegible]

5.4.3 Slovenske novice

[illegible]

```

v></div><div class="container item_break_02_cont outer_cont"><div class="row col_grid"><div
class="col-xs-12"><div><div class="custom"><div
class="OUTBRAIN" id="outbrain_widget_1"><div class="ob-widget ob-strip-layout AR_1"><div
class="ob-widget-section ob-first"><ul class="ob-widget-items-container"><li
class="ob-dynamic-rec-container ob-recldx-0 ob-o"><a class="ob-dynamic-rec-link"><span class="ob
it ob-rec-text">#text</span></a></li><li class="ob-dynamic-rec-container ob-recldx-1 ob-o"><a
class="ob-dynamic-rec-link"><span class="ob-unit ob-rec-text">#text</span></a></li></ul><ul
class="ob-widget-items-container ob-multi-row ob-row-1"><li class="ob-dynamic-r
ec-container ob-recldx-7 ob-o"><a class="ob-dynamic-rec-link"><span class="ob-unit
ob-rec-text">#text</span><span class="ob-unit
ob-rec-source">#text</span></a></li></ul></div></div></div><div class="OUTBRAIN"
id="outbrain_widget_2"><div class="ob-widget ob-strip-l
ayout AR_2"><div class="ob-widget-section ob-first"><ul class="ob-widget-items-container"><li
class="ob-dynamic-rec-container ob-recldx-2 ob-o"><a class="ob-dynamic-rec-link"><span class="ob
ob-rec-text">#text</span></a></li><li class="ob-dynamic-rec-container
ob-recldx-3 ob-o"><a class="ob-dynamic-rec-link"><span class="ob-unit
ob-rec-text">#text</span></a></li></ul><ul class="ob-widget-items-container ob-multi-row ob-row-1">
class="ob-dynamic-rec-container ob-recldx-4 ob-o"><a class="ob-dynamic-rec-link"><span clas
s="ob-unit ob-rec-text">#text</span></a></li><li class="ob-dynamic-rec-container ob-recldx-5 ob-o">
class="ob-dynamic-rec-link"><span class="ob-unit
ob-rec-text">#text</span></a></li></ul></div></div></div></div></div></div></div></div></div></div></div>
</body></

```

html>

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.