# Ekstrakcija podatkov

Enes Fejzoski, Marko Kofol, Elian Mugerli 10. Maj 2024

To poročilo opisuje implementacijo treh različnih pristopov za ekstrakcijo strukturiranih podatkov iz spleta: regularne izraze (regex), XPath in algoritem za samodejno ekstrakcijo vsebine (RoadRunner). Poudarek je na pridobivanju podatkov iz štirih primerjalnih spletnih mest: rtvslo.si, bolha.com, overstock.com in sporttv.si.

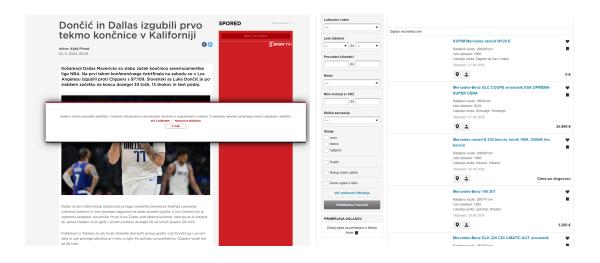
## 1. Uvod

Ekstrakcija strukturiranih podatkov iz spleta je ključna naloga za različne aplikacije, vključno z web scrapingom, pridobivanjem informacij in združevanjem podatkov. To poročilo raziskuje tri pogoste pristope za doseganje tega: regularne izraze, XPath in algoritme za samodejno ekstrakcijo vsebine.

Regularni izrazi (regex) so zmogljiva orodja za iskanje in manipulacijo besedila na podlagi vzorcev. So jedrnati in učinkoviti pri pridobivanju določenih formatov podatkov iz HTML kode. XPath je poizvedovalni jezik, posebej zasnovan za navigacijo in ekstrakcijo podatkov iz XML dokumentov, ki se lahko uporabijo za predstavitev HTML strukture. Algoritmi za samodejno ekstrakcijo vsebine, kot je RoadRunner, se osredotočajo na učenje pravil ekstrakcije iz nabora spletnih strani, s čimer avtomatizirajo postopek identifikacije in pridobivanja relevantnih podatkovnih elementov.

# 2. Izbira dodatnih spletnih strani

Poleg že dodeljenih spletnih domen smo za širitev analize izbrali tudi strani sporttv.si in bolha.com.



## 3. Implementacija

## 3.1. Regularni izrazi

Regularni izrazi (angl. "regular expressions" ali "regex") so izrazi, ki se uporabljajo za iskanje in ujemanje vzorcev v besedilu. Regularni izrazi so implementirani v Pythonu z uporabo knjižnice re. Tukaj je primer za pridobivanje podatkov iz novic na rtvslo.si:

```
import re
def rtv(html):
   title_regex = re.compile(r'<h1>(.*?)</h1>')
   subtitle regex = re.compile(r'<div class="subtitle">(.*?)</div>')
   lead_regex = re.compile(r'(.*?)')
   content_regex = re.compile(
       r'<div *?class="article-body">(.*?)<div class="gallery">',
       re.DOTALL)
   author_regex = re.compile(r'<div class="author-name">(.*?)</div>')
   published_time_regex = re.compile(r'<div</pre>
class="publish-meta">\n\t\t(.*?)<br>')
   data = {
        'Naslov': re.findall(title_regex, html)[0],
        'Podnaslov': re.findall(subtitle_regex, html)[0],
        'Uvod': re.findall(lead regex, html)[0],
        'Vsebina': re.findall(content_regex, html)[0],
        'Avtor': re.findall(author_regex, html)[0],
        'Cas Objavve': re.findall(published_time_regex, html)[0]
   }
   return data
```

#### Razlaga kode:

- 1. **import** re: Uvozi knjižnico re za delo z regularnimi izrazi.
- 2. **def** rtv(html): Definira funkcijo rtv, ki sprejme HTML vsebino kot vhod.

- 3. **title\_regex = re.compile(r'<h1>(.\*?)</h1>')**: Ustvari objekt regularnega izraza title\_regex za ujemanje naslova znotraj oznake <h1>. (.\*?) zajame vse znake znotraj začetne in končne oznake <h1>.
- 4. Podobne vrstice sledijo za subtitle\_regex, lead\_regex, author\_regex in published\_time\_regex, vsaka cilja na določene HTML elemente in zajame njihovo vsebino znotraj oklepajev.
- 5. **re.DOTALL**: Ta zastavica omogoča regexu, da se ujema z znaki nove vrstice (\n), kar je potrebno za zajetje vsebine, ki se lahko razteza na več vrstic.

#### 3.2. XPath

XPath je drug pristop za pridobivanje podatkov iz HTML dokumentov. Uporablja potenco osnovano sintakso za lociranje določenih elementov in atributov znotraj drevesne strukture HTML. Tukaj je primer za pridobivanje podatkov iz novic na rtvslo.si:

```
from lxml import html
def rtv(html_content):
   tree = html.fromstring(html_content)
   title_xpath = "//h1/text()"
   subtitle_xpath = "//div[@class='subtitle']/text()"
   lead_xpath = "//p[@class='lead']/text()"
   content_xpath = "//div[@class='article-body']//text()"
   author_xpath = "//div[@class='author-name']/text()"
   published time xpath = "//div[@class='publish-meta']/text()"
   title = tree.xpath(title_xpath)[0]
   subtitle = tree.xpath(subtitle xpath)[0]
   lead = tree.xpath(lead_xpath)[0]
   content = ' '.join(tree.xpath(content_xpath))
   author = tree.xpath(author xpath)[0]
   published time = tree.xpath(published time xpath)[0]
   data = {
        'Naslov': title,
        'Podnaslov': subtitle,
        'Uvod': lead,
        'Vsebina': content,
        'Avtor': author,
```

```
'Cas Objavve': published_time
}
return data
```

Razlaga XPath nekaterih izrazov:

- **1.** //h1/text(): Ta izraz izbere vse elemente <h1> kjerkoli v dokumentu (označeno z //) in izvleče njihovo besedilno vsebino.
- 2. //div[@class='subtitle']/text(): Izbere vse elemente div s pripisanim razredom "subtitle" in izvleče njegovo besedilno vsebino.
- 3. Podobni izrazi sledijo za lead, content, author in published\_time\_xpath, ki ciljajo na določene elemente glede na njihova imena oznak ali razredov in izvlečejo njihovo besedilno vsebino.
- 4. content = ' '.join(tree.xpath(content\_xpath)): Za vsebino uporabimo tree.xpath za pridobitev seznama tekstovnih elementov znotraj elementa "article-body". Nato jih združimo z presledkom (" "), da ustvarimo eno samo besedo, ki predstavlja celotno vsebino.

XPath ponuja bolj strukturiran in deklarativen način za navigacijo po drevesu HTML v primerjavi z regularnimi izrazi. Bolje se spopade s kompleksnimi strukturami HTML, vendar zahteva poznavanje XPath sintakse in specifične hierarhije elementov HTML.

### 3.3. RoadRunner

Tretji pristop za pridobivanje podatkov iz HTML dokumentov je implementacija algoritma RoadRunner. Tega smo razvili s pomočjo grajenja DOM dreves in pripadajočih elementov.

Implementacijo sestavljajo naslednje funkcije:

- 1. **build\_dom\_tree**: Ta rekurzivno sestavi DOM drevo iz razčlenjene HTML strani s BeautifulSoup. Hkrati pa preskoči oznake, ki niso prisotne na seznamu allowed\_tags.
- 2. **Build\_generalized\_tree**: Funkcija primerja dve HTML drevesi, predstavljeni z objekti *Node*. Ugotavljanje skladja oznak, neskladja besedil in ponavljajoča se vozlišča. Glede na te primerjave nastavi vozlišča kot izbirna ali ponovljiva.
- Pomožne funkcije: Skripta vključuje več pomožnih funkcij, kot je check\_repeated\_nodes za preverjanje, ali so vozlišča ponovljena, check\_string\_mismatch za preverjanje neskladij besedil in print\_dom\_tree za izpis generaliziranega DOM drevesa z oznakami za izbirnost in ponovitev.
- 4. **Print\_dom\_tree**: Funkcija za izpis DOM drevesa rekurzivno izpiše generalizirano DOM drevo z oznakami za izbirnost in ponovitev, zamaknjeno glede na globino vozlišča.

Spodaj so podani izseki implementacije.

```
class Node:
    def __init__(self, tag, text, classes, children, optional):
        self.tag = tag
        self.text = text
        self.classes = classes
        self.children = children
        self.optional = optional
        self.repeatable = False
```

```
def road_runner(pages):
    trees = []
    for page in pages:
        soup = BeautifulSoup(page, 'lxml')
        tree = build_dom_tree(soup.find("html"))

        trees.append(tree)

    wrapper = build_generalized_tree(trees[0], trees[1])
    print_dom_tree(None, wrapper)

def build_dom_tree(tag):
    node = Node(tag.name, tag.string if tag.string else None,
    tag.attrs.get('class', []), [], False)

if allowed_tags is None or tag.name in allowed_tags:
    for child in tag.children:
        if isinstance(child, str):
            continue

    if child.name in allowed_tags:
            node.children.append(build_dom_tree(child))

return node
```

```
def build_generalized_tree(tree1, tree2):
    if tree1.tag != tree2.tag:
        tree1.optional = True
        return tree1

if check_string_mismatch(tree1, tree2):
        tree1.text = "#PCDATA"

else:
    if tree1.children:
```

```
return False

return True

def check_string_mismatch(node1, node2):
   if node1.text and node2.text:
     return node1.text != node2.text
   return False
```

Spodaj bodo podani se primeri tesnih izhodov:

#### 1. RTV

```
</div><div></div></div></div><div><div><div><<div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div>
2></h2><
v><div></div><div></div></div><div><div><div>lskanje</div></div><div><div><div><div
></
div><div><div><div><div><div></div></div></div></div></div></div><div><div><div><div><div></div>
v><div><div><div><div><div></div><div><div><div><div>
><div>
</div></div></div>Zaradi testiranja je dodanih umetnih 2 sekunde
></div></div><div><d
iv><div><div></div></div></div><div><div><div></div></div></div></div></div></div
/div
></div><div><div></div><div><div><div><div><div><div><div>
Merljak</div></div><div></div><div><div><div></div><div></div></div></div></div></div>
</div></div><div><div><di
v></div></div><div><div><div></div></div></div></div></div></div></div></div></div></div></div></div>
</div></div></div><div><div></div></div></div></div></div></div><div><div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div>
div><
/div></div><div><div><div><div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div>
iv><div></div><div></div></div></div></div></div></div></div><div><div><div><div>Test
noveg
a modela</div></div><div><div>Test nove
generacije</div></div><div><div>Test nove
generacije</div></div><div>Preizkus novega modela</div></div></div
><div><div><div></div></div></div><div><div><div><div><div><div><giv><div><div><div>
</div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></di>
div><
```

/div></div><div></div><div><div><div><div>

v><div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></di>

v></div><div></div><div><div><div>Razvoj avtomobilske

varnosti</div></div><div><qiv>Ekipa Avtomobilnosti na avtomobilskem salonu v

Żenevi</div></div><div><div>Obu

ditev legende</div></div><div><div>Ogrožena naravna

><div></div></div><div><div><div><div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></di>

div><div></div><div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></di>

#### 2. Overstock

>

td>
#PCDATA
#PCDATA

#PCDATA#PCDATA

#PCDATA#PCDATA

NecklacesPendant

sPins/Slides

Earrings/Pendants<Solitaire Rings</td>

Wedding/Anniversar

```
v Bands View
AllApparel, Shoes &
Access.Apparel, Shoes & Access.Books
. Movies, CDs, GamesBooks, Movies, CDs, GamesElectronics &
ComputersElectronics & ComputersHome & GardenHome &
GardenJewelry, Watc
hes & GiftsJewelry, Watches & GiftsSports, Travel & ToysSports, Travel
ToysWorldstockWorldstock<td
>
able>Shopping Cart & CheckoutShopping Cart &
CheckoutT
rack Your OrderTrack Your OrderYour AccountYour
AccountHelp & FAQHelp & FAQBest Price GuaranteeBest
Price Guarantee<
/tbody>About UsAbout UsPrivacy &
SecurityPrivacy & SecurityTerms & ConditionsTerms &
ConditionsBecome An
AffiliateBecome An AffiliateBusiness PurchasesBusiness
PurchasesHave Products to Sell?Have Products to
Sell?Investor Relations<td
>Investor
Relations
>
<tb
body>More Info...More
Info...</tbod
y>
Price:#PCDATAPrice:#PCDATAYou
Save:#PCDATA<
/tbody>
Info...<t
d>More
Info...
Price:#PCDATAYou Sa
ve:#PCDATA
More Info...
```

>tody></

ble>List

Price:#PCDATAPrice:#PCDATAYou

Save:#PCDATA</tbody

>

>More

Info...More

Price:#PCDATAPrice:#PCDATAYou

Save:</t

d>#PCDATA

>

></td

>More Info...More

Info...List

Price:#PCDATAPrice:#PCDATA

You

Save:#PCDATA

tr>Kr>

o...More

Info...List

Price:#PCDATAPrice:#PCDATA

You

Save:#PCDATA

Info...

Price:#PCDATAPrice:

>#PCDATAYou

Save:#PCDATA

tr>

More Info...More

Info...List

Price:\$639.99Price:\$199.99

/tr>You Save:\$440.00

(68%)

>More Inf

o...More

Info...

Price:\$305.00Price:\$119.99

You Save:\$185.01

(60%)

Brenda of Falmouth MA USA says..Brenda of Falmouth M

A USA says..Yet elegantA USA says..It is subtle yet elegant

and has such a beautiful shine to the gold.""The quality of this piece is absoul

```
tly amazing, it is subtle yet elegant and has such a beautiful shine to the
qold."</t
d>More
Info...More
Info...
Price:$999.99
You Save:$640.00
(64%)
>More Info...
More
Info...><<td>
Price:$508.99Price:$179.99You Save
:$329.00
(64%)
>
List
Price:$196.99Price:$69.99You
Save:$12
7.00
(64%)
>
body>List
Price:$1,369.99Price:$409.99You
Save:$960.00 (70%)
More Info...</tab
le>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table>table><t
Price:$1,635.00Price:$609.99You
Save:$1,025.01 (62%)</t
r>
able><
td
>
/tr></body></html>
```

S pridobljenimi rezultati sicer nismo bili najbolj zadovoljni in ima tako podana implementacija veliko prostora za nadgradnjo.