

Slika 3: Regularni izrazi, uporabljeni za ekstrakcijo podatkov s strani Ayto.net

```

1 title = '/tr/td[2]/a/b/text()'
2 list_price = '/tr/td[2]/table/table/tbody/tr[1]/td[2]/s/text()'
3 price = '/tr/td[2]/table/table/tbody/tr[2]/td[2]/span/b/text()'
4 saving = '/tr/td[2]/table/table/tbody/tr[3]/td[2]/span/text()'
5 content = '/tr//span[@class="normal"]//text()'

```

Slika 4: XPath izrazi, uporabljeni za ekstrakcijo podatkov s strani Overstock

```

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

```

Slika 5: XPath izrazi, uporabljeni za ekstrakcijo podatkov s strani RTV SLO

stran le malo spremeni, bi naši XPath izrazi ne delovali več.

IV. AVTOMATIČNA EKSTRAKCIJA PODATKOV

Na spletu se od uvedbe programskega jezika PHP naprej pojavlja dosti dinamičnih spletnih strani. Take strani hranijo podatke v podatkovni bazi in glede na dostopen URL naslov izvedejo poizvedbo v podatkovni bazi in na podlagi rezultata prikažejo podatke v uporabniku prijazni obliki. Pri avtomatični ekstrakciji podatkov, bi radi na podlagi dveh strani, ki imata enako strukturo, a drugačno vsebino, zgradili program, ki zna samodejno pridobiti podatke.

V našem primeru smo za avtomatično ekstrakcijo podatkov uporabili algoritem ROADRUNNER [1], [2]. Ta sicer v originalni implementaciji deluje na seznamu žetonov (angl. token), v našem primeru pa smo za predstavitev HTML vsebin uporabili kar DOM drevo. Algoritem deluje tako, da v pomnilnik prebere dve datoteki - ovojnico (angl. wrapper) in drugi dokument. S pomočjo funkcije GENERALIZE, skuša nato ovojnico ustrezno uskladiti z dokumentom tako, da išče istoležne elemente, v njih išče polja (angl. fields) in opsijske elemente (angl. optionals).

Implementiran algoritem zna na strani prepoznati polja, tj. elemente, ki se pojavijo na obeh straneh, a imajo drugačne vrednosti. Prav tako zna program prepoznati opsijske elemente, tj. elemente, ki se pojavijo v enem izmed dokumentov, v drugem pa ne. Našo implementacijo bi lahko izboljšali tako,

```

1 title = '//h3/text()'
2 price = '//div[@id="FINANCE"]/following-sibling::div/p/text()'
3 date_published = '//i[contains(@class, "fa-calendar")]/parent::div/text()'
4 number_of_views = '//i[contains(@class, "fa-bar-chart")]/parent::div/text()'
5 seller = '//div[contains(text(), "Prodajalec")]/following-sibling::div/following-sibling::div/strong/text()'

```

Slika 6: XPath izrazi, uporabljeni za ekstrakcijo podatkov s strani Avto.net

```

1 title = '/li//h3[@class="entity-title"]/a/text()'
2 price = '/li//li[@class="price-item"]/strong/text()'
3 location = '/li//div[@class="entity-description-main"]/text()'
4 date_posted = '/li//time/text()'

```

Slika 7: XPath izrazi, uporabljeni za ekstrakcijo podatkov s strani Bolha

```

<html>
  Books of: <b>PCDATA</b>
  
  <ul>
    <li>
      <i>Title:</i> #PCDATA
    </li>
    <li>
      <i>Title:</i> #PCDATA
    </li>
    <li optional="True">
      <i>Title:</i> Javascript
    </li>
  </ul>
</html>

```

Slika 8: Rezultat izvajanja avtomatične ekstrakcije podatkov

da bi dodali prepoznavanje iteratorjev, tj. elementov, ki se ponavljajo. Tega dela nam zaradi pomanjkanja časa, ni uspelo implementirati.

Slika 8 prikazuje rezultat izvajanja avtomatične ekstrakcije podatkov nad primerom iz članka [1]. Naš program opsijske elemente označi z atributom `optional="True"`. Kot lahko vidimo, je naša implementacija ustrezno prepoznala polja in opsijske elemente.

V. ZAKLJUČEK

Pisanje regularnih in XPath izrazov za pridobivanje podatkov je zamudno in od razvijalca zahteva poznavanje oziroma razumevanje strukture spletne strani. Kljub temu, pa daje tak način ekstrakcije podatkov boljše rezultate, saj podatke označimo oziroma umestimo v nek kontekst. Razvijalec lahko npr. na mestih kjer se pojavijo cela števila, le-ta iz nizov pretvori v številčno predstavitev in s tem še dodatno opiše podatke. Avtomatična ekstrakcija podatkov je princip, pri katerem sistem na strani sam najde polja, ki vsebujejo podatke. Rezultati, tj. neoznačene izluščene informacije pa so v tem primeru pogosto pomanjkljivi oziroma težko razumljivi.

Prva metoda je uporabna predvsem v primerih, ko je potrebno izluščiti podatke iz ene same spletne strani. Pisanje izrazov je časovno potratno in je smiselno, kadar vemo, da se bo struktura spletne strani spreminjala zelo redko.

Druga metoda je zelo uporabna v primeru, ko imamo za ekstrahirati podatke iz ogromne količine različnih strani in smo za prihranek časa, pripravljeni žrtvovati natančnost.

LITERATURA

- [1] V. Crescenzi, G. Mecca, and P. Merialdo, "The roadrunner project: Towards automatic extraction of web data," 01 2001.
- [2] —, "Automatic web information extraction in the road runner system," in *International Conference on Conceptual Modeling*. Springer, 2001, pp. 264–277.