

Ekstrakcija podatkov

Skupina **DOMACI-NJOKI**: Niki Bizjak, Bojan Vrangeloski, Uroš Škrjanc

Povzetek—Cilj prve seminarske naloge je bil napisati pajka, ki zna s spleta prenašati spletne vsebine in jih shranjevati v podatkovno bazo. Pred iskanjem podatkov po taki podatkovni bazi, pa je treba prenesene vsebine najprej prečistiti in iz njih izluščiti pomembne informacije. V drugem seminarju si bomo ogledali tri različne načine za ekstrakcijo podatkov iz HTML vsebin.

I. UVOD

Druga seminarska naloga pri predmetu Iskanje in ekstrakcija podatkov s spleta je namenjena pridobivanju informacij iz vsebin, ki jih je s spleta prenesel pajek. V seminarju smo si ogledali tri različne načine ekstrakcije - dva taka, ki zahtevata veliko uporabnikovega poseganja in razumevanja strukture strani in tretjega, ki zna informacije prepoznati in izluščiti avtomatično, s primerjavo podobnih spletnih strani.

Ekstrakcijo podatkov smo izvajali na štirih različnih spletnih straneh, ki jih glede na strukturo razdelimo na dve skupini:

- **list** - stran vsebuje *seznam* izdelkov, artiklov, člankov, ipd.
- **detail** - stran prikazuje podatke za en izdelek, članek, ipd.

Tabela I, prikazuje naše testne strani, razvrščene glede na tipa strani, opisana zgoraj.

Stran	Tip
overstock.com	list
rtvslo.si	detail
bolha.com	list
avto.net	detail

Tabela I: Tip strani za posamezno domeno, nad katero smo izvajali ekstrakcijo podatkov

II. REGULARNI IZRAZI

Regularni izrazi nam omogočajo učinkovito iskanje informacij v nizih z uporabo končnih avtomatov. Regularni izrazi se lahko uporabijo za preverjanje, če posamezen niz vsebuje iskan vzorec. Pri implementaciji ekstrakcije z regularnimi izrazi, smo si pomagali s vgrajeno Python knjižnico `re`.

Za ekstrakcijo podatkov s strani Overstock, smo uporabili regularne izraze, prikazane na sliki 1. Najprej smo za vsako iskano polje na strani definirali svoj regularni izraz, nato pa smo jih z ustreznimi vmesnimi regularnimi izrazi, povezali v skupen izraz, ki zna hkrati poiskati vse podatke. Na podoben način smo ekstrahirali podatke tudi iz ostalih spletnih strani. Regularne izraze za spletno stran RTV SLO lahko vidimo na sliki 2, slika 3, pa prikazuje regularne izraze za spletno stran Avto.net.

[illegible]

Slika 1: Regularni izrazi, uporabljeni za ekstrakcijo podatkov s strani Overstock

[illegible]

Slika 2: Regularni izrazi, uporabljeni za ekstrakcijo podatkov s strani RTV SLO

III. XPATH

XPath je poizvedovalni jezik, ki je bil razvit za namene ekstrakcije podatkov iz XML podatkovnih struktur. Izrazi so sestavljeni v obliki datotečnih poti in lahko vsebujejo pogoje za spust po določeni veji drevesa, klice funkcij nad vejami, ipd. Slike 4, 5, 6 in 7 prikazujejo XPath izraze, ki smo jih uporabili za ekstrakcijo podatkov iz izbranih spletnih strani.

Največje težave sta predstavljali strani `Auto.net` in `Overstock`, saj elementi niso bili semantično označeni. Ostale strani imajo namreč z razredi ali identifikatorji označene elemente na strani, kar olajša iskanje, saj lahko take elemente hitro najdemo s filtriranjem. Spletna stran `Overstock` je sestavljena iz vgnezenih tabel, kar še posebej oteži iskanje, saj je potrebno poznati indekse elementov, ki jih iščemo. Če se taka

```

1 link_regex = "<a class='stretched-link' href='{({})}'><sup>1</sup></a>"
2 title_regex = "<tbody class='GO-Results-Naviz {^}^>{<+></span>{<+>/div>"
3 data_regex = "<td class='{s}'+tr{<+>td class='w-25 d-none d-md-block pl-3'>{<+></td>{<+>td class='w-75 pl-3'>{<+>{<+>}"
4 "<tr>{<+>tr{<+>td{<+>td class='d-none d-md-block pl-3'>{<+></td>{<+>td{<+>td class='pl-3'>{<+></td>{<+>}"
5 "<tr>{<+>tr{<+>td class='d-none d-md-block pl-3'>Gorivoc</td>{<+>td class='pl-3'>{<+></td>{<+>td{<+>tr{<+>tr{<+>}"
6 "<td class='d-none d-md-block pl-3'>Menjalikan</td>{<+>td class='pl-3 text-truncate'>{<+></td>{<+>}"
7 "<tr>tr class='d-none d-md-table-row'>{<+>td class='d-none d-md-block pl-3'>Motor</td>{<+>td class='pl-3 text-truncate'>{<+>}"
8 "{<+></td>{<+></tr>{<+>tr{<+>tbody>{<+>table>"
9
10 price_regex = "<div class='GO-Results-Top-)>Price-TXT-[{BbsTt\^}]^>{<+>/div>"

```

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

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

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.