

# 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.

## 1. 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.

### III. XPATH

XPath je proizvedovalni jezik, ki je bil razvit za namene ekstrakcije podatkov iz XML podatkovnih struktur.

[illegible]

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

```
1 title_regex = r"<[^\>+]+[^\>]*><div class=\"subtitle\">[^\>+]+</div>[^\>+]*</div>"
2 author_regex = r"<div class=\"author-timestamp\">[^\>+]+<strong>[^\>+]*</strong>[^\>+]/div>"
3 lead_regex = r"<p class=\"lead\">[^\>+]+</p>"
4 content_regex = r"<div>[^\>+]*?<figure>[^\>+]*?<[\{|\$|}%*><div lass=\"gallery\">[^\>+]*</div>[^\>+]*?</div>"
5 regex = title_regex + author_regex + r"<[^\>+]*>" + lead_regex + r"<[^\>+]*>" + content_regex
```

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

#### 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 dostopan 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.

```

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

```

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

## 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.