

Jakob Dekleva, Matjaž Bevc in Urban Poljšak

Poročilo za drugo domačo nalogo (Data extraction from the Web)

Domača naloga pri predmetu Iskanje in ekstrakcija podatkov s spleta

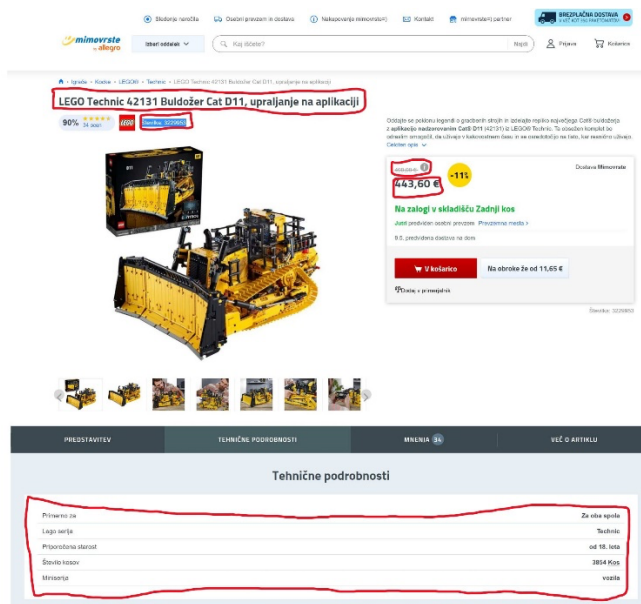
MENTORJA: prof. dr. Marko Bajec, doc. dr. Slavko Žitnik

1. Uvod

V tem poročilu bomo opisali izdelavo treh načinov ekstrakcije podatkov iz spleta: regularni izrazi, xpath in implementacija splošnega algoritma Webstemmer. Izbrane spletne strani so overstock.com, rtv.si in mimovrste.si.

2. Opis izbranih spletnih strani

Za dodatne dve strani smo izbrali produkte iz spletne strani mimovrste.si. Iz spletne stani smo ekstrahirali: naslov, številko izdelka, ceno, znižano ceno, dostopnost in tehnične lastnosti izdelka. Parametri so označeni na sliki 1.



Slika 1: ekstrahirani podatki

3. Implementacija regularnih izrazov in xpath-a

Regular Expressions (regex) omogočajo relativno učinkovito iskanje zelenih elementov v HTML kodi, čeprav orodje ni zasnovano za to uporabo. XPath je orodje, posebej zasnovano za iskanje elementov v XML in HTML strukturah. Omogoča enostavno definiranje lokacije elementa v drevesni strukturi in po potrebi iskanje njegovih pod-elementov. Iskanje običajno poteka tako, da definiramo tip in prepoznavno lastnost želenega elementa.

<pre>Overstock Title = (.*) ListPrice = (.*) Price = (.*) Saving = \$(.)*% SavingPercent = \$(.)*% Content = (.*) RTV title = <h1>(.*)</h1> subtitle = <div class="subtitle">(.*)</div> lead = <p class="lead">(.*)</p> author = <div class="author-name">(.*)</div> time = <div class="publish-meta">(.*)</div> content = <div class="article-body">(.*)</div> Minovrste title = <h1>(.*)</h1> number = (.*) availability = <div class="availability-box">(.*)</div> old_price = (.*) final_price = <div class="price_wrap">(.*)</div> table_params = <table class="product-parameters">(.*)</table> table_values = <tbody class="product-parameters">(.*)</tbody></pre>	<pre>Overstock title = //a[@href="http://www.overstock.com/cgi-bin/02.cgi?PAGE=PRODUCT&PROD_ID=4"] listPrice = //span[@class="bigred"] price = //span[@class="bigred"] saving = //span[@class="littleorange"] savingPercent = //span[@class="littleorange"] content = //span[@class="normal"] RTV title = //h1 subtitle = //div[@class="subtitle"] lead = //p[@class="lead"] author = //div[@class="author-name"] time = //div[@class="publish-meta"] content = //div[@class="article-body"] Minovrste title = //h1 number = //span[@class="detail_title"] availability = //div[@class="availability-box"] old_price = //span[@class="price_wrap"] final_price = //div[@class="price_wrap"] table_params = //table[@class="product-parameters"] table_values = //tbody[@class="product-parameters"]</pre>
---	---

Regularni izrazi

Xpath-a

Slika 2: Implementacija regularnih izrazov in xpath

4. Implementacija algoritma Webstemmer

Funkcija webstemmer deluje tako, da sprejme dve HTML strani ter pragove za; merjenje podobnosti, odkrivanje naslova in glavnega besedila. Razčleni strani, izračuna matriko podobnosti in združi podobne strani v grozde. Nato odstrani statične bloke in odkrije naslov ter glavno besedilo znotraj grozdov. Rezultate pretvori v JSON objekt in jih vrne kot končni rezultat.

Psevdokoda rešitve

- webStemmer(html_page1, html_page2, sim_threshold, diff_threshold, title_threshold, main_text_threshold)
 - wrapper ← []
- Razčlenjevanje strani:
 - parsed_pages ← parse_page(html) za vsako html stran
- Izračun podobnosti:
 - sim_matrix ← compute_similarity_matrix(parsed_pages)
- Grupiranje in ustvarjanje vzorca postavitev:
 - clusters ← cluster_pages(sim_matrix, sim_threshold)
- Obdelava vsakega grozda:
 - Za vsak grozd:
 - Izračunaj diff_scores za bloke v grozdu
 - Odstrani statične bloke (diff_scores > diff_threshold)
 - Poišči naslov in glavno besedilo (upoštevaj title_threshold, main_text_threshold)
 - Dodaj rezultate v wrapper

6. Vrni wrapper kot JSON niz

```
[
  "a",
  "span"
],
"Click here to purchase."
],
[
  [
    "tr",
    "td",
    "tr",
    "td",
    "table",
    "tbody",
    "tr",
    "td",
    "a",
    "tr",
    "td"
  ],
  "More Info..."
],
[
  [
    "td",
    "a"
  ],
  "14-kt. Diamond 7.5-8 mm Pearl Pendant"
],
[
  [
    "table",
    "tbody",
    "tr",
    "td",
    "table",
    "tbody",
    "tr",
    "td"
  ],
  "List Price:"
],
[
  [
    "td"
  ],
  "$196.99"
],
[
  [
    "tr",
    "td"
  ],
  "Price:"
],
[
```

Overstock

```
"Temni na\u010din BETA"
],
[
  [
    "BETA"
  ],
  [
    [
      "label",
      "div",
      "div",
      "div",
      "div"
    ],
    "Prijavljen"
  ],
  [
    [
      "button"
    ],
    "\u00d77"
  ],
  [
    [
      "div"
    ],
    "Odjava"
  ],
  [
    [
      [
        "Uporabni\u0161ki na\u010dun"
      ],
      [
        "div",
        "div"
      ],
      "Temni na\u010din BETA"
    ],
    [
      [
        "BETA"
      ],
      [
        [
          "label",
          "div",
          "div",
          "div"
        ],
        [
          "\u010dasovno obdobje po meri"
        ],
        [
          "button"
        ]
      ]
    ]
  ]
],
```

RTV

```
],
"Partner"
],
[
  [
    "g"
  ],
  [
    "Pridru\u017eite se mimovrste=)Partner programu"
  ],
  [
    [
      [
        "Item 1 of 5"
      ],
      [
        [
          "div",
          "div",
          "div",
          "div",
          "ul"
        ],
        [
          "Nakupovanje na mimovrste=)"
        ],
        [
          "li",
          "a",
          "li",
          "a",
          "li",
          "a",
          "li",
          "a",
          "li",
          "a",
          "li",
          "a",
          "div",
          "div",
          "ul"
        ],
        "Kontakt in pomo\u010d"
      ],
      [
        "li",
        "a",
        "li",
        "a"
      ]
    ]
  ]
],
```

Mimovrste

Slika 3: izhodni wrapper za Overstock, RTV in Mimovrste