

# **Vyhľadávanie informácií**

## **Zadanie 1**

### **Search**

Matúš Gáspár

2.ročník

2017/2018

# 1. Úvod

Úlohou bolo získať dáta z webu, oindexovať ich a vyhľadávať v nich podľa nami zvolených scenárov v Elasticsearchi.

Aplikácia pozostáva z vlastného rozhrania v Jave (Eclipse IDE), ktorá používa Maven dependency, kde v pom.xml súbore sú uložené dependencies. Na začiatku programu možné vybrať si z troch možností:

- A. Zapnúť crawler, ktorý bude prechádzať stránkou "hej.sk" a zbierať údaje z ich webu
- B. Zozbierané dáta vložiť do Elasticsearchu
- C. Samotné vyhľadávanie nad vloženými dátami

# 2. Crawler

Crawler je naprogramovaný v programovacom jazyku Java. Na začiatku (pri vybraní možnosti crawlovania webu hej.sk), sa do listu vloží počiatočný URL odkaz na stránku, na ktorej sa má začať crawlovať. Následne celé fungovanie je založené na tom, že sa prechádza tento list, až pokiaľ nebude prázdny a vyberá sa z neho vždy ďalší URL odkaz do neho vložený.

```
while (!urlsInProcess.isEmpty()) {  
    url = urlsInProcess.pop();.....
```

Pri prechádzaní každého jedného URL odkazu sa deje to, že sa okontroluje, či sme na správnej stránke a taktiež či sa nenachádzame na URL odkaze s košíkom alebo nejakým ".doc" súborom. V celom crawleri sa používa hashmapa už navštívených URL odkazov, aby sa neprechádzali znova. Preto sa vždy kontroluje, či už daný odkaz nie je v hasmape. Ak nie je, tak sa kontroluje kanonická URLka (s obsahom ? a #) a až následne sa na neho pripojí cez knižnicu "Jsoup.org" a vytvorí sa z neho HTML dokument, na ktorom sa pozbierajú všetky URL odkazy.

```
connection = Jsoup.connect(url).userAgent(browser);  
htmlDocument = connection.header("Accept-Encoding", "gzip, deflate").get();  
hrefsOnPage = htmlDocument.select("a[href]");
```

Zbierajú sa HTML stránky produktov, to znamená, že sa vždy kontroluje, či stránka, ktorá je práve prechádzaná, nie je stránka produktu. Ak nie, vloží sa len do hasmapy navštívených. Ak

áno, tak sa vyberú: **kategórie produktu, nadpisy, popisy, počet recenzií, prítomnosť dokumentácie k produktu, výrobca produktu, popisy z tabuliek a samotný URL odkaz.**

Príklad selektoru pre výber nadpisu produktu:

```
heading = htmlDocument.select(".product-summary .page-title");
```

Rovnako ako nadpis sa selektujú aj ostatné dáta zo stránky. Toto všetko sa vykonáva až kým v liste nie je žiadna URLka.

V niektorých prípadoch boli nutné modifikácie vyselektovaných dát ako napríklad počet recenzií:

```
review =  
Integer.parseInt(reviews.text().split("\\(")[1].split("\\)")[0].split("  
") [1]);
```

Všetky sparované dáta sa ukladali ako JSONy do súborov (použitá knižnica *JSON.org*), príklad:

```
object.put("Odkaz", url.toString());
```

### 3. Vkládanie dát do Elasticsearchu

Všetky sparované dáta sa musia uložiť do Elasticsearchu. Všetko sa to vykonáva pomocou *TransportClient* v Jave. Pripája sa na localhost, na ktorom beží Elasticsearch:

```
this.client = new PreBuiltTransportClient(Settings.EMPTY)  
.addTransportAddress(new InetSocketAddress(  
InetAddress.getByName("localhost"), 9300));
```

Následne sa prechádzajú všetky uložené súbory s JSONmi v kódovaní UTF-8 nakoľko sa v súboroch nachádzajú slovenské a české slová a vytvorí sa index pre všetky dáta:

```
String json = new String(Files.readAllBytes(Paths.get(path + index)),  
StandardCharsets.UTF_8);  
  
IndexResponse indexes = client.prepareIndex("produkty", "produkt",  
Integer.toString(index)).setSource(json,  
XContentType.JSON).get();
```

### Príklad JSONu v elasticu:

```
{
  "Výrobca": "Siemens",
  "Popis": "Přáli byste si, aby vám příprava jídla přinášela stejnou radost jako stolování? Ke splnění tohoto snu vám pomůže kval",
  "Kategória": ["Velké domácí spotřebiče", "Mikrovlnné rúry"],
  "Odkaz": "https://www.hej.sk/mikrovlнна-rura-siemens-hf24g541/",
  "Cena": 117,
  "Názov": "Mikrovlnná rúra Siemens HF24G541 nerez",
  "Informácie o produkte": {
    "Hlavné informácie": ["Záruka: 24 mesiacov", "Typ: voľne stojacia", "Mikrovlnný výkon: 900 W", "Gril - výkon: 1200 W", "Gri",
    "Ďalšie informácie": ["Antibakteriálny vnútorný priestor: Nie", "Ovládanie: mechanické-gombíkové, tlačidlóvé - zatlačacie"],
  },
  "Dokumentácia k produktu": "Áno",
  "Počet recenzií": 0
}
```

## 4. Aplikácia vyhľadávania v dátach

Hlavné okno aplikácie vyhľadávania vyzerá takto:

Zadaj query

mikrovlína rúra

Hľadať

☐ S dokumentáciou

Cena od:Cena do:

☐ V popise☐ V kategóriách☐ Vo výrobkoch☒ V názvoch

Zoraď vzostupne podľa:

☐ Ceny

Zoraď zostupne podľa:

☐ Počtu recenzií

Myslíš:

mikrovlnné

mikrovlnná

mikrovlnný

mikrovlínka

mikrovlinní

-----ďalšie slovo myslíš: -----

rúka  
rúra  
rury  
roura  
ruda

-----

Priemerná cena produktov bez filtrov: 195.00584795321637

Názov	Cena	Výrobca	Dokumentácia k produ.	Počet recenzií
Mikrovlnná rúra Electrolux EMS28201OW	129	Electrolux	Ano	2
Mikrovlnná rúra Gorenje MMO20DBII čierna	81	Gorenje	Ano	1
Mikrovlnná rúra Bosch CMG633BS1 nerez	787	Uvedený v názve	Ano	0
Mikrovlnná rúra Siemens BF634LGS1 nerez	480	Siemens	Ano	0
Mikrovlnná rúra Bosch BFR634GS1 nerez	503	Uvedený v názve	Ano	0
Mikrovlnná rúra Gorenje BM171E2XG čierna	218	Gorenje	Ano	0
Mikrovlnná rúra Siemens HF15M564 nerez	232	Siemens	Ano	1
Mikrovlnná rúra Siemens HF15G541 nerez	104	Siemens	Ano	12
Mikrovlnná rúra Electrolux EMT25207OX nerez	279	Electrolux	Ano	8
Mikrovlnná rúra Zanussi ZSM17100XA nerez	169	Zanussi	Ano	0
Mikrovlnná rúra Siemens CM656GBS1 nerez	980	Siemens	Ano	0
Mikrovlnná rúra Electrolux EMT25207OW biela	289	Electrolux	Ano	0
Mikrovlnná rúra Siemens HF24GS41 nerez	117	Siemens	Ano	0
Mikrovlnná rúra Candy CMG2071DS strieborná	65	Candy	Ano	0
Mikrovlnná rúra Siemens BF634RGS1 nerez	502	Siemens	Ano	0
Mikrovlnná rúra Electrolux EMS20107OX nerez	189	Electrolux	Ano	5
Mikrovlnná rúra Electrolux EMS30400OX nerez	139	Electrolux	Ano	9
Mikrovlnná rúra Bosch BFL634GS1 nerez	498	Uvedený v názve	Ano	0
Mikrovlnná rúra Concept KTV4444 nerez	489	Concept	Ano	0
Mikrovlnná rúra Electrolux EMT25207OB čierna	289	Electrolux	Ano	0
Mikrovlnná rúra Electrolux EMS26004OK čierna	409	Electrolux	Ano	0
Mikrovlnná rúra Concept MTV3020 nerez	184	Concept	Ano	0
Mikrovlnná rúra AEG Mastery KME721000M nerez	519	AEG	Ano	0
Mikrovlnná rúra AEG Mastery MSB2547D-M nerez	299	AEG	Ano	0
Mikrovlnná rúra Electrolux EMS 21400 W biela	95	Electrolux	Ano	38
Mikrovlnná rúra Gallet FMOM 420W biela	48	Gallet	Ano	11
Mikrovlnná rúra Gorenje BM 300 X nerez	123	Gorenje	Ano	2
Mikrovlnná rúra Whirlpool MW 4200 IX nerez	261	Uvedený v názve	Ano	2
Mikrovlnná rúra AEG Mastery MBE2658D-M nerez	479	AEG	Ano	0
Mikrovlnná rúra Gorenje MO 4250 CLB čierna	108	Gorenje	Ano	0
Mikrovlnná rúra Mora MT 120 W biela	55	Mora	Ano	19
Mikrovlnná rúra ETA MORELO 0209 90000 biela	62	ETA	Ano	0
Mikrovlnná rúra Gorenje MO 17 DW biela	61	Gorenje	Ano	2
Mikrovlnná rúra Whirlpool MCP 349 BL čierna	210	Uvedený v názve	Ano	0
Mikrovlnná rúra Gallet FMOM 171W biela	51	Gallet	Ano	0
Mikrovlnná rúra Sencor SMW 6420DS čierna/strieborná	69	Sencor	Nie	0
Mikrovlnná rúra Amica TMI 20 AFB čierna	159	Amica	Ano	0
Mikrovlnná rúra AEG MFC3026S-M nerez/ocel	159	AEG	Ano	2
Mikrovlnná rúra ETA Klasiko 1208 90000 biela	69	ETA	Ano	9
Mikrovlnná rúra ETA Klasiko 2208 90000 strieborná	79	ETA	Ano	8
Mikrovlnná rúra Bosch HMT 75M451 nerez	96	Uvedený v názve	Ano	4
Mikrovlnná rúra Samsung Muse3 MS23FE301FAW/EO	102	Uvedený v názve	Ano	6

Metóda *queryMaker*(*final TransportClient client*) v triede *Elastic*, je zodpovedná za všetky query, ktoré je možné v aplikácii spúšťať. Query je možné napísať do textovej poľa a kliknúť na tlačidlo "Hľadať", vtedy sa hľadajú záznamy (by default) v názvoch produktov. Pre vyhľadávania v názvoch je použité *boolQuery*, cez JAVA API, ktorý používa *QueryBuilder*:

```
BoolQueryBuilder query = QueryBuilders.boolQuery();
```

Pokiaľ nie je zaškrtnutý checkbox "V názvoch" tak sa vykonáva *boolQuery* s *occurrence-om must*:

```
query.must(QueryBuilders
    .matchQuery("Názov", window.queryField.getText()).minimumShouldMatch("30%")
    .fuzziness(Fuzziness.AUTO).operator(Operator.AND).slop(5));
```

Inak sa používa *boolQuery* s *occurrence-om should* nakoľko chceme zobrazit' výsledky, ktoré sú relevantné aj pre záznamy, kde sa slovo nachádza v názvoch. Používam *fuzziness* na to, aby bolo možné sa do určitej miery pomýliť v písaní slov, použitý je aj *operator AND*, aby daný názov obsahol všetky hľadané slová (nielen niektoré z nich), následne je použitý aj *slop* s hodnotou 5, na to, že hľadané slová nemusia byť bezprostredne za sebou, ale môže byť aj o 5 miest ďalej:

```
query.should(QueryBuilders
    .matchQuery("Názov", window.queryField.getText())
    .minimumShouldMatch("30%")
    .fuzziness(Fuzziness.AUTO).operator(Operator.AND).slop(5));
```

Vyhľadávať je možné rovnako aj kategóriách produktov, v popise produktov a aj vo výrobcach. Možno je aj označiť, hľadať tie produkty, ktoré majú dokumentáciu. Tieto checkboxy je možné zaškrtnúť všetky naraz. Pri jednotlivých checkboxoch, okrem "V názvoch" sa používa *boolQuery* s *occurrence-om must*.

Vyhľadávanie v popisoch je podobné ako v názvoch s pridaním *boostu*, aby sa zvýšilo skóre výsledkov nájdených v popisoch:

```
query.must(QueryBuilders.matchQuery("Popis", window.queryField.getText())
    .fuzziness(Fuzziness.AUTO).slop(50).boost(5));
```

Ostatné query ako hľadanie vo výrobcach, hľadanie produktov s dokumentáciou a hľadanie v kategóriách sú takéto, všetky s *occurrence-om must*:

```
query.must(QueryBuilders.matchQuery("Výrobca", window.queryField.getText())
    .fuzziness(Fuzziness.AUTO));
```

```
query.must(QueryBuilders.matchQuery("Dokumentácia k produktu", "Áno"));
```

```
query.must(QueryBuilders.matchQuery("Kategória",
    window.queryField.getText()).fuzziness(Fuzziness.AUTO));
```

Nakoniec sa vyskladá výsledná query cez `.prepareSearch()`, kde sa nastaví default filter pre cenu od 0 po 999999999 (táto cena môže byť taktiež zmenená používateľom a následne sú záznamy filtrované podľa zadanej ceny) a vypíše sa prvých 1000 nájdených záznamov do tabuľky v GUIku:

```
responses.add(client.prepareSearch().setQuery(query).setSize(1000).setPostFilter(QueryBuilders.rangeQuery("Cena").from(cena_od).to(cena_do)));
```

V aplikácii sa nachádzajú aj tlačidlá pre zoradenie výsledkov zostupne podľa ceny a vzostupne podľa recenzií. Zoradenie je urobené spôsobom, že nájdené výsledky sa neustále ukladajú do premennej, ku ktorej potom zoradovanie pristupuje:

```
Iterator<String> hits = getHits(responses.get(responses.size() - 1)
.addSort("Cena", SortOrder.ASC).execute().actionGet(), urls,
model).iterator()
```

Jednoducho je len ku získaným *hits* (nájdenej matchov) pridané `.addSort` podľa toho či to chceme vzostupne alebo zostupne `SortOrder.ASC` a `SortOrder.DES`.

V JSONe môžu byť tieto query vyskladané nasledovne, príklad ako je to možné:

```
{
  "query": {
    "bool": {
      "must": {
        "match": {
          "Názov": {
            "query": "mikrovlnn rura electorlux",
            "fuzziness": "AUTO",
            "operator": "and",
            "slop": 5
          }
        }
      },
    },
    "filter": {
      "range": {
        "Cena": { "gte": 0, "lte": 900 }
      }
    },
    "minimum_should_match": "30%",
    "boost": 1.0
  },
  "sort": [{"Cena": {"order": "asc"}}]
}
```

Vyššie napísaný JSON je možné upraviť na každú z daných requestov, ktoré boli v JAVE napísané cez QueryBuilder. V JSONe sa nachádza *sort* (ktorý je zapnutý pri zoradení cien, je tam možné pridať zoradenie podľa recenzií rovnakým spôsobom), nachádza sa tam *fuzziness* (ktorý používam cez queryBuilder), *boost*, *range*, *operator*, a to všetko je v boolQuery, teda presne tak ako som to robil v JAVE cez queryBuilder.

V projekte sa taktiež nachádza *suggest*(autocomplete), *highlighting* (vyznačovanie nájdených slov v názvoch) a agregácia typu priemerná cena.

### **Suggest(autocomplete)**

Suggest je vykonávaný pomocou POST request priamo na localhost kde je spustený Elasticsearch. JSON, ktorý je vykonávaný pri autocomplete:

```
"{\r\n" +
"  \"suggest\": {\r\n" +
"    \"autocomplete\": {\r\n" +
"      \"text\": \"\"+query+\"\", \r\n" +
"      \"term\": {\r\n" +
"        \"field\": \"\"+field+\"\" \r\n" +
"      } \r\n" +
"    } \r\n" +
"  } \r\n" +
"}";
```

Hodnotou textu v JSONe je napísané slovo pri hľadaní daného záznamu. Hodnotou vo *fielde* je daný stĺpec v tabuľke (ja to vykonávam nad stĺpcom "Názov"). Všetko sa to vykonáva hneď ako používateľ napíše aspoň 5 písmen. Ak nemá čo navrhnúť tak, nenavrhone nič. V programe sa vykonáva to, že z daného *responzu*, získam objekt JSONu zo *suggest* a z neho zoberiem pole JSONov z *autocomplete*. Takto dostanem JSON, v ktorom sú *options*, kde sa nachádzajú dané návrhy. Tieto options si iba zoberiem a vypíšem.

### **Highlighting(vyznačovanie slov v názvoch)**

Highlighting je vykonávaný pomocou GET requestu na localhost so spusteným Elasticsearchom. JSON na highlighting:

```
"{\r\n" +
"  \"query\": {\r\n" +
"    \"match\": {\"\"+field+\"\": {\r\n" +
"      \"query\": \"\"+query+\"\", \r\n" +
"      \"fuzziness\": \"AUTO\", \r\n" +
"      \"operator\": \"and\" \r\n" +
"    } \r\n" +
"  } \r\n" +
"  , \r\n" +
"  \"highlight\": {\r\n" +
"    \"fields\": {\"r\" +
```

```
"      \"Názov\" : {}\\r\\n\" +
"    }\\r\\n\" +
"  }\\r\\n\" +
"}";
```

Vykonáva sa to tak, že response ktorý príde vráti JSON *highlight*, v ktorom sa nachádzajú otagované slová v názvoch produktov (nakolko robím highlight pre názvy). Príklad: <em>rúra </em>. Využívam automatický fuzziness na čiastočný *misspelling* a taktiež používam operátor AND, aby to v názvoch obsahovalo všetky slová, ktoré používateľ vyhľadáva.

To znamená, že je potrebné získať všetky slová medzi týmito tagmi. Tieto slová získam a následne prechádzam tabuľku s výsledkami, kde aplikujem html tag(<b>) na zvýraznenie daných slov v názvoch.

### Agregácia(priemerná cena)

Query pre vypočítanie priemernej ceny produktov je možné vidieť iba pri hľadaní slov v názvoch produktov (t.j. pri zaškrtnutej položky "v názvoch") . Vytvára sa POST request na localhost s Elasticsearchom. JSON pre agregáciu je nasledovný:

```
{"\\r\\n\" +
"  \"query\": {\\r\\n\" +
"    \"match\" : {\" + \"\"+field+\"\": {\\r\\n\" +
"      \"query\" : \"\"+query+\"\",\\r\\n\" +
"      \"fuzziness\": \"AUTO\",\" +
"      \"operator\": \"and\"\\r\\n\" +
"    }\\r\\n\" +
"  }\\r\\n\" +
"},\\r\\n\" +
"  \"aggs\": {\\r\\n\" +
"    \"avg_grade\" : {\\r\\n\" +
"      \"avg\" : {\\r\\n\" +
"        \"field\" : \"Cena\"\\r\\n\" +
"      }\\r\\n\" +
"    }\\r\\n\" +
"  }\\r\\n\" +
"}";
```

Priemerná cena je vypočítaná z cien vrátených produktov (hits). Používam fuzziness, operator AND. Objekt JSONu, ktorý sa vráti po úspešnom vypočítaní priemernej ceny je *aggregations* a z neho je potrebné si zobráť *avg\_grade* a *value*. Tu sa potom nachádza výsledná vypočítaná priemerná cena.

Pri pracovaní z JSONom v Jave je vždy nutné nastaviť na aký localhost pristupovať, aký typ requestu a následne mu poslať JSON ako string:

```
String url = "http://localhost:9200/produkty/produkt/_search";
URL obj = new URL(url);
URLConnection connection = (URLConnection) obj.openConnection();

connection.setRequestMethod("GET");
```



***Zaujímavosť v projekte:***

Aplikácia obsahuje aj jednu zaujímavú vec a tou je, že na zobrazené produkty v tabuľke je možné kliknúť a tým zobrazíť daný produkt vo webovom prehliadači:

```
Desktop.getDesktop().browse(uri);
```