

Matej Puk

IMDB vyhľadávač

Vyhľadávanie informácií

Študijný program: Softvérové inžinierstvo

Študijný odbor: 9.2.5. Softvérové inžinierstvo

Ročník: 2.

Cvičenie: Ut 11⁰⁰, miestnosť -1.65

Ak. rok: 2017/2018

Obsah

Zadanie.....	1
Crawler.....	2
Parser.....	3
Webová aplikácia.....	5
Zoradenie podľa hodnotenia.....	5
Výsledky vyhľadávania.....	6
Filtrovanie výsledkov na základe agregácií žánrov a hviezd.....	8
Vyhľadávanie od-do podľa hodnotenia.....	9
Histogram rokov vydania.....	10
Návrhy vyhľadávania.....	10
Vizualizácie.....	12
Vizualizácia 1.....	12
Vizualizácia 2.....	13
Vizualizácia 3.....	13
Mapovanie dokumentu movie v indexe imdb.....	15

Zadanie

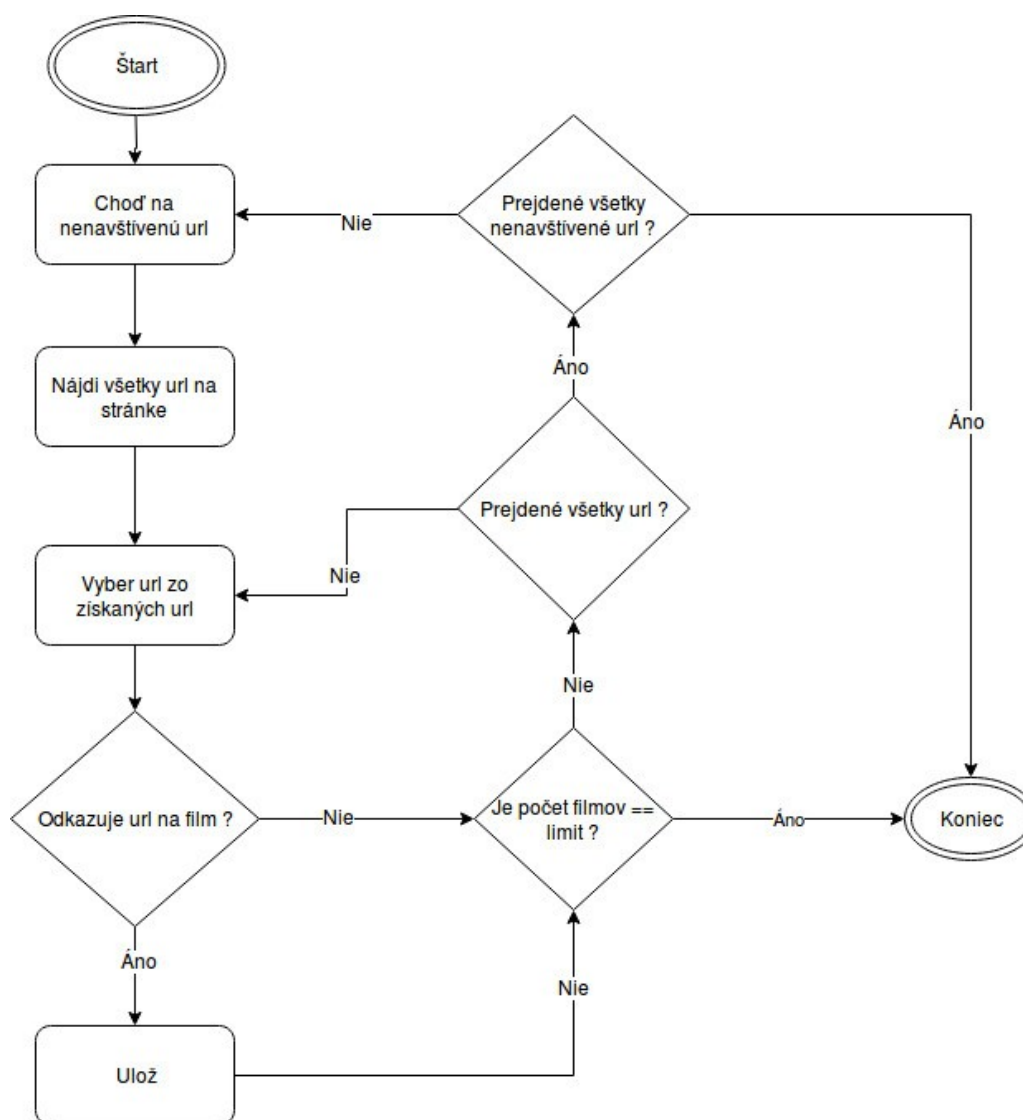
Cieľom zadania bolo vytvorenie crawlera, parsera pre získanie dát z vybratého internetového zdroja. Ako zdroj dát sme vybrali internetovú filmovú databázu **IMDB**. Následne po získaní dát sme vytvorili webovú aplikáciu, ktorá ponúka vyhľadávanie v získaných dátach na základe rôznych kritérií ako sú napríklad žáner filmu alebo jeho hodnotenie. Nakoniec sme vytvorili niekoľko vizualizácií získaných dát v nástroji **KIBANA**.

Celý projekt je implementovaný v jazyku **python**. Závislosti aplikácie sú spravované pomocou nástrojov **pip** a **bower**. Pre prepojenie python skriptov a webovej aplikácie s technológiou **Elasticsearch** bol použitý python modul **elasticsearch**.

Crawler

Prvou súčasťou zadania je **crawler**. Jeho implementácia sa nachádza v súbore **crawler.py**. Crawler je implementovaný pomocou triedy **Crawler**. Pri inicializácii je tejto triede zadáný štartovacia internetová stránka, z ktorej sa začína vyhľadávanie. V našom prípade je to stránka <http://www.imdb.com/genre/>. Následne je proces crawlovania spustený metódou **crawl()**, ktorej je zadáný argument **limit**, ktorý reprezentuje počet filmov, ktoré má crawler nájsť.

Na nasledujúcom diagrame je zobrazený algoritmus metódy **crawl()**.



Obrázok 1: Algoritmus metódy **crawl()**

Výstupom metódy **crawl()** je set jedinečných url odkazujúcich na filmy na stránke **IMDB**.

Parser

Druhou časťou zadania je vytvorenie parsera. Implementácia parsera sa nachádza v súbore `parser.py`. Parser je implementovaný triedou `Parser`. Pri inicializácii parsera je mu ako argument zadaný zoznam filmov, ktorý je výstupom crawlera. Následne je zoznam týchto filmov prejdený metódou `parse()`, ktorá prejde každý film v zozname uloží jeho html kód v podobe súboru na disk aby ho pri ďalšom spracovaní nebolo nutné znova sťahovať. Následne parser prejde celý html kód nachádzajúci sa v súbore a vyextrahuje z neho nasledovné dáta:

- hodnotenie filmu,
- počet ľudí, ktorý hodnotili film,
- názov filmu,
- rok vydania,
- hodnotenie obsahu,
- dĺžka,
- popis,
- žánre,
- režisér,
- dej filmu,
- tvorcovia,
- herci,
- hviezdy,
- ocenenia,
- kľúčové slová,
- krajina,
- jazyk,

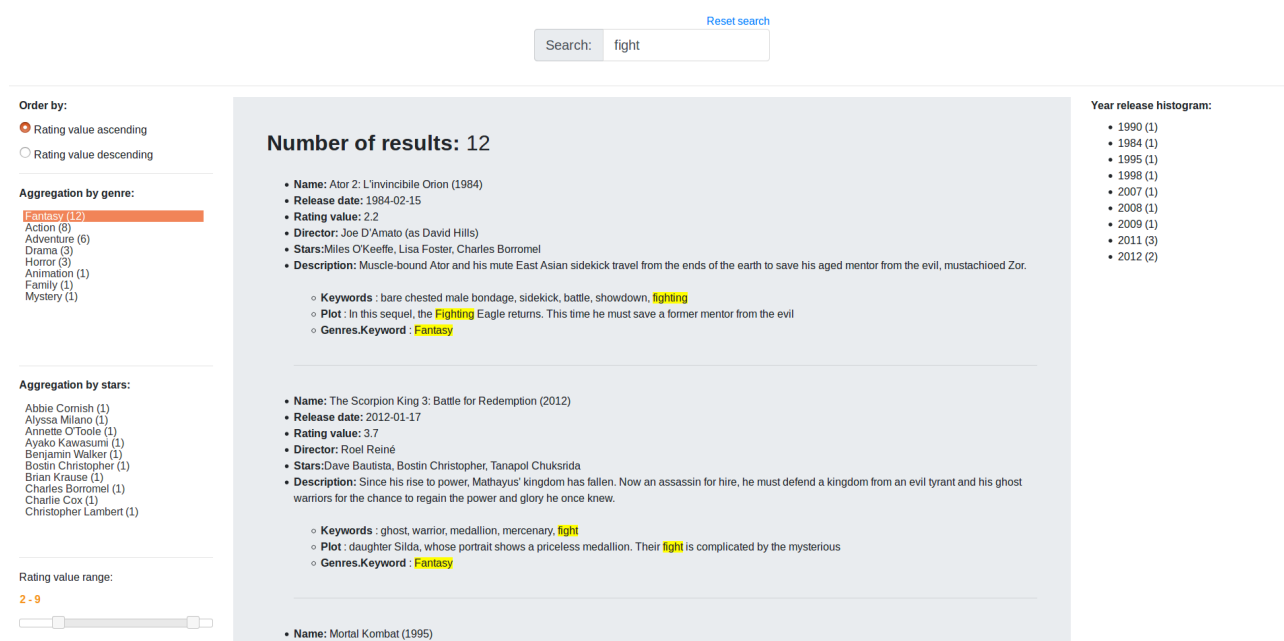
- dátum vydania.

Po vyextrahovaní dát je film uložený do databázy. Pre vyextrahovanie dát bol použitý modul **BeautifulSoup**.

Webová aplikácia

Pre implementáciu webovej aplikácie bol použitý framework **Flask**. Aplikácia je implementovaná v súbore `app.py`. Šablóny pre vykreslenie aplikácie sa nachádzajú v adresári **templates**. Aplikácia vykresľuje hlavnú šablónu **main.html** a následne do nej vykresľuje šablónu s výsledkami vyhľadávania **hits.html**. Aplikácia sa dopytuje do databázy cez metódu **search()**. V tejto metóde sa nachádza celé telo vyhľadávacieho dopytu a aj telo dopytu pre sugescie v prípade, vyhľadávací dopyt nenájde žiadne položky.

Na nasledujúcom obrázku je zobrazený celkový vzhľad aplikácie.



Obrázok 2: Webová aplikácie pre vyhľadávanie filmov

Na vrchu aplikácie v strede sa nachádza pole pre zadanie vyhľadávanej frázy. Zároveň je nad týmto poľom pridané tlačidlo pre celkové vynulovanie vyhľadávania.

Zoradenie podľa hodnotenia

Používateľ má možnosť zoradiť výsledky vyhľadávania podľa hodnotenia filmu a to vzostupne a zostupne.

Order by:

☒ Rating value ascending

☐ Rating value descending

Obrázok 3: Zoradenie výsledkov podľa hodnotenia filmu

Implementácia v dopyte:

```
"sort": [  
  {  
    "ratingValue": {  
      "order": o  
    }  
  },  
  "_score"  
]
```


Výsledky vyhľadávania

Ak dopyt nájde výsledky, tie sú zobrazené používateľovi uprostred stránky. Záznamy sú oddelené čiarou. Pre každý záznam je zobrazené meno, dátum vydania, hodnotenie, režisér, hviezdy a opis. Následne sú zobrazené všetky atribúty, v ktorých nastala zhoda a tá je zvýraznená žltou farbou.

Number of results: 12

- **Name:** Ator 2: L'invincibile Orion (1984)
 - **Release date:** 1984-02-15
 - **Rating value:** 2.2
 - **Director:** Joe D'Amato (as David Hills)
 - **Stars:** Miles O'Keeffe, Lisa Foster, Charles Borromel
 - **Description:** Muscle-bound Ator and his mute East Asian sidekick travel from the ends of the earth to save his aged mentor from the evil, mustachioed Zor.
 - **Keywords :** bare chested male bondage, sidekick, battle, showdown, fighting
 - **Plot :** In this sequel, the Fighting Eagle returns. This time he must save a former mentor from the evil
 - **Genres.Keyword :** Fantasy
-
- **Name:** The Scorpion King 3: Battle for Redemption (2012)
 - **Release date:** 2012-01-17
 - **Rating value:** 3.7
 - **Director:** Roel Reiné
 - **Stars:** Dave Bautista, Bostin Christopher, Tanapol Chuksrida
 - **Description:** Since his rise to power, Mathayus' kingdom has fallen. Now an assassin for hire, he must defend a kingdom from an evil tyrant and his ghost warriors for the chance to regain the power and glory he once knew.
 - **Keywords :** ghost, warrior, medallion, mercenary, fight
 - **Plot :** daughter Silda, whose portrait shows a priceless medallion. Their fight is complicated by the mysterious
 - **Genres.Keyword :** Fantasy
-

Obrázok 4: Výsledky vyhľadávania

Filtrovanie výsledkov na základe agregácií žánrov a hviezd

Používateľ má možnosť si nájdené výsledky ďalej vyfiltrovať podľa žánrov alebo hviezd, ktoré vo filme vystupujú. Naviac je pri každom výsledky agregácie zobrazená početnosť v zátvorke.

Aggregation by genre:

Fantasy (12)

Action (8)

Adventure (6)

Drama (3)

Horror (3)

Animation (1)

Family (1)

Mystery (1)

Aggregation by stars:

Abbie Cornish (1)

Alyssa Milano (1)

Annette O'Toole (1)

Ayako Kawasumi (1)

Benjamin Walker (1)

Bostin Christopher (1)

Brian Krause (1)

Charles Borromel (1)

Charlie Cox (1)

Christopher Lambert (1)

Obrázok 5: Agregácie podľa hviezd a žánrov

Implementácia v dopyte:

```
"aggs": {  
  "by_genre": {  
    "terms": {  
      "field": "genres.keyword",
```

```

    "size": 10
  },
  "by_stars": {
    "terms": {
      "field": "stars.keyword",
      "size": 10
    }
  }
}

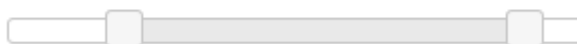
```

Vyhľadávanie od-do podľa hodnotenia

Používateľ má možnosť si výsledky vyfiltrovať podľa ich hodnotenia zadaním hodnoty od-do na **jquery-ui** prvku **slider**.

Rating value range:

2 - 9



Obrázok 6: Filter od-do podľa hodnotenia

Implementácia v dopyte:

```

"range": {
  "ratingValue": {
    "gte": int(rf),
    "lte": int(rt)
  }
}

```

Histogram rokov vydania

Na pravej strane v aplikácii má používateľ možnosť vidieť histogram ukazujúci početnosť filmov vydaných pre jednotlivé roky.

Year release histogram:

- 1990 (1)
- 1984 (1)
- 1995 (1)
- 1998 (1)
- 2007 (1)
- 2008 (1)
- 2009 (1)
- 2011 (3)
- 2012 (2)

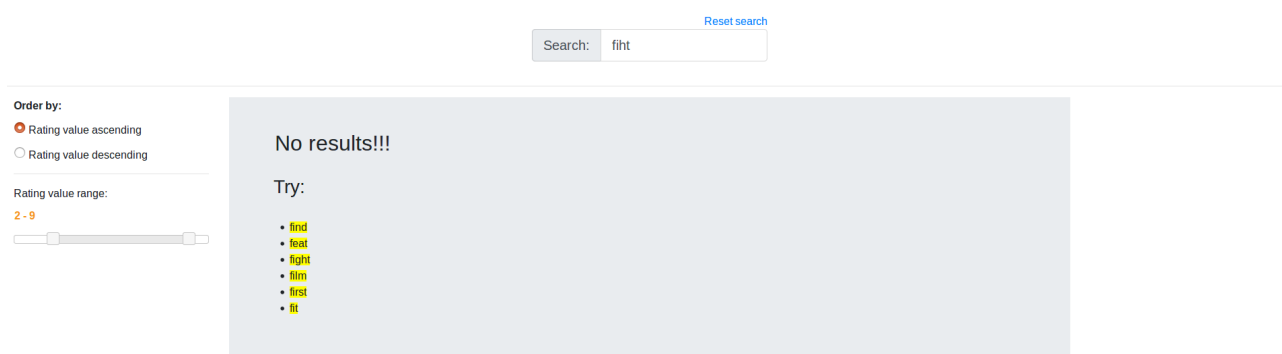
Obrázok 7: Histogram vydaných filmov
vzhľadom na rok

Implementácia v dopyte:

```
"by_releaseYears": {  
  "date_histogram": {  
    "field": "releaseDate",  
    "interval": "year"  
  }  
}
```

Návrhy vyhľadávania

Ak používateľ zadá frázu, ktorá vráti nulový počet výsledkov, aplikácia zavolá dopyt, ktorý navrhne používateľovi akú frázu podobnú tej, ktorú zadal má zadať aby získal výsledky. Na túto funkcionality bol využitý **phrase suggester**, ktorý bol mapovaním nastavený aby vyhľadával návrhy v atribútoch názov, popis a dej filmu.



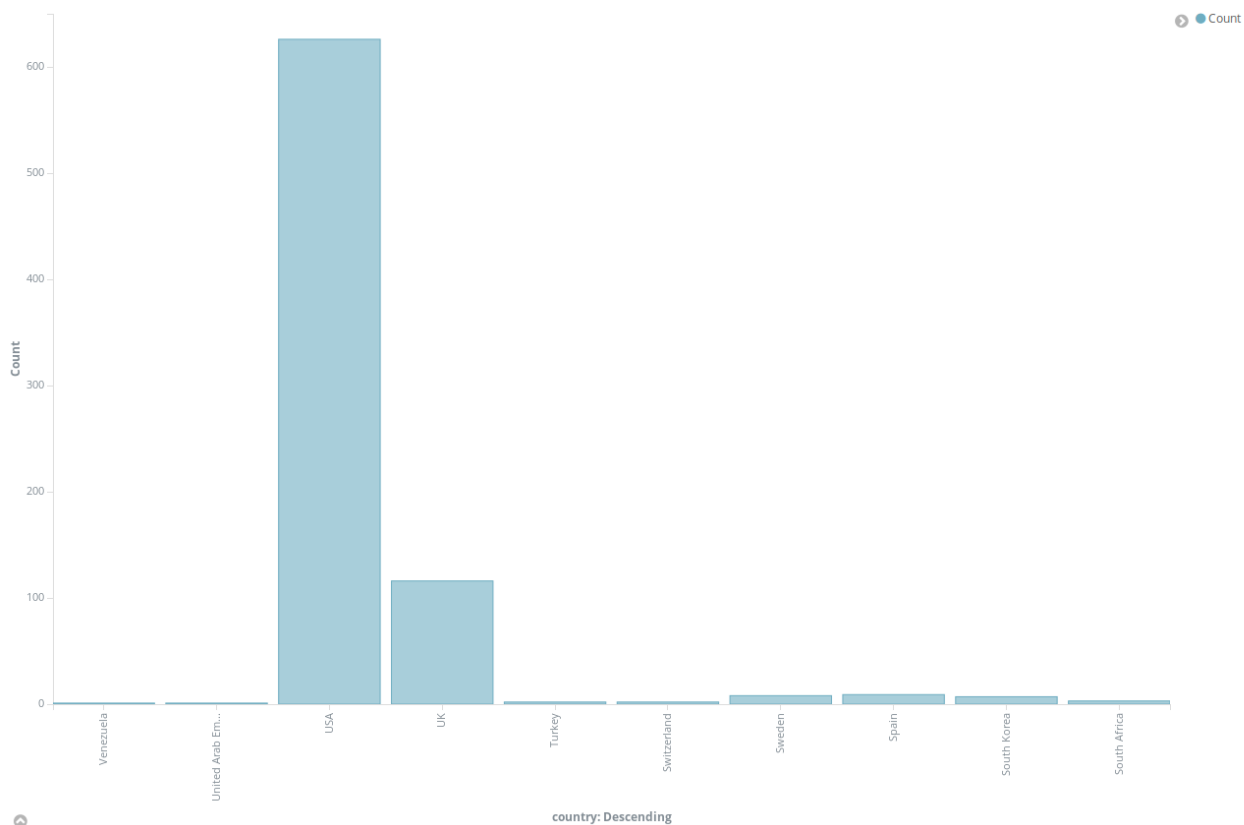
Obrázok 8: Návrhy pri nulovom počte výsledkov

Vizualizácie

V nasledujúcej kapitole sú zobrazené tri vizualizácie uskutočnené nad vyextrahovanou databázou.

Vizualizácia 1

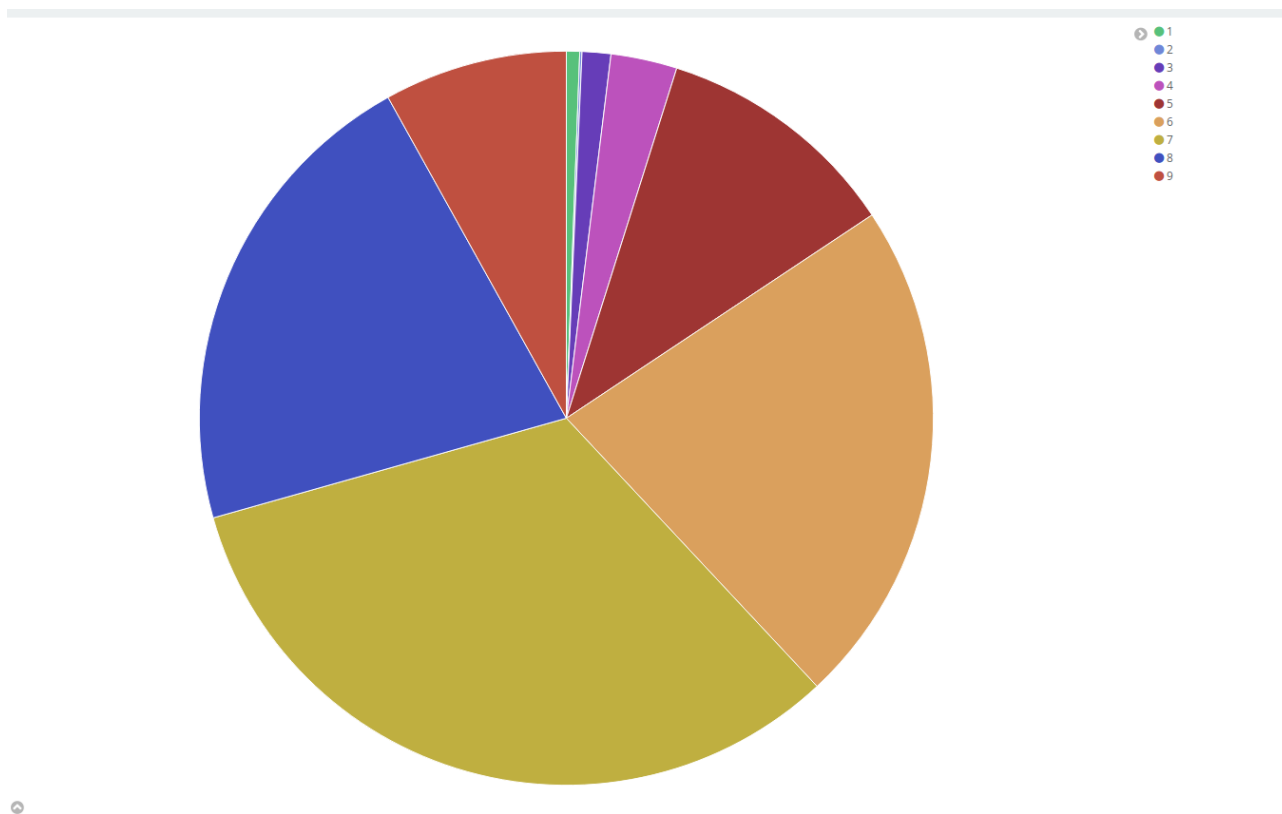
Nasledujúca vizualizácia zobrazuje počet filmov vzhľadom na krajinu za obdobie od roku 2012 do roku 2017.



Obrázok 9: Vizualizácia 1

Vizualizácia 2

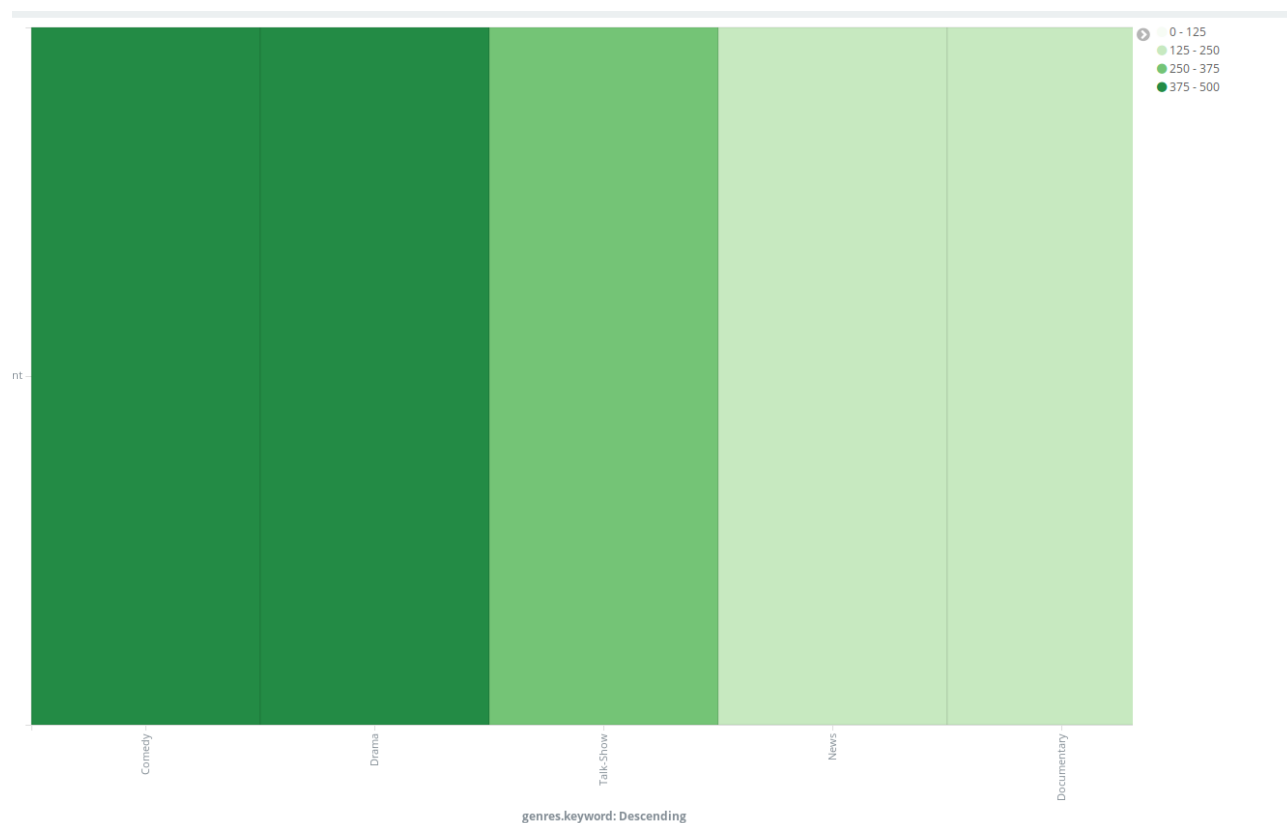
Nasledujúca vizualizácia zobrazuje počet filmov vzhľadom na ich hodnotenie od 1 do 10 po intervaly 1 za obdobie od roku 2012 do roku 2017.



Obrázok 10: Vizualizácia 2

Vizualizácia 3

Nasledujúca vizualizácia zobrazuje početnosť filmov vzhľadom na ich žáner za obdobie od roku 2012 do roku 2017.



Obrázok 11: Vizualizácia 3

Mapovanie dokumentu movie v indexe imdb

V nasledujúcej ukážke kódu je zobrazené nastavenia mapovania pre dokument typu **movie** v indexe **imdb**.

```
PUT imdb
{
  "settings": {
    "index": {
      "number_of_shards": 1,
      "analysis": {
        "analyzer": {
          "suggest": {
            "type": "custom",
            "tokenizer": "standard",
            "filter": ["standard", "shingle"]
          }
        },
        "filter": {
          "shingle": {
            "type": "shingle",
            "min_shingle_size": 2,
            "max_shingle_size": 5
          }
        }
      }
    },
    "mappings": {
      "movie": {
        "properties": {
          "name": {
            "type": "text",
            "analyzer": "english",
          },
          "fields": {
            "suggest": {
              "type": "text",
              "analyzer": "suggest"
            }
          }
        },
        "plot": {
          "type": "text",
          "analyzer": "english",
        }
      }
    }
  }
}
```

```

        "fields": {
            "suggest": {
                "type": "text",
                "analyzer": "suggest"
            }
        },
    },
    "genres": {
        "type": "text",
        "analyzer": "english",
    },
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    },
    },
    "director": {
        "type": "text",
        "analyzer": "english"
    },
    },
    "keywords": {
        "type": "text",
        "analyzer": "english"
    },
    },
    "awards": {
        "type": "text",
        "analyzer": "english"
    },
    },
    "stars": {
        "type": "text",
        "analyzer": "english",
    },
    "fields": {
        "keyword": {
            "type": "keyword",
            "ignore_above": 256
        }
    },
    },
    "duration": {
        "type": "text",
        "analyzer": "english"
    },
    },
    "actors": {
        "type": "text",
        "analyzer": "english"
    },
    },
    "creators": {
        "type": "text",
        "analyzer": "english"
    },

```

```

    },
    "description": {
        "type": "text",
        "analyzer" : "english",
    "fields": {
        "suggest": {
            "type": "text",
            "analyzer": "suggest"
        }
    }
    },

    "ratingValue": {
        "type": "float"
    },
    "ratingCount": {
        "type": "integer"
    },

    "language": {
        "type": "keyword"
    },
    "country": {
        "type": "keyword"
    },

    "releaseDate": {
        "type": "date"
    }
}
}
}

```