



NBA - prvá konzultácia

Analýza dát

Informácie o hráčovi z basketball-reference [https://www.basketball-reference.com/] budú pozostávať zo základných informácií o hráčovi zobrazené na obrázku tu:

The screenshot shows the Basketball Reference website for player Kareem Abdul-Jabbar. The page includes his photo, name, position (Center), height (7'2"), weight (238 lbs), birth date (April 16, 1947), college (UCLA), high school (Power Memorial Academy), draft (1st round, 1st pick), and career highlights (HOF in 1996). It also lists his NBA debut (1969-70) and career length (20 years).

V html tento element vyzerá nasledovne:

```
1 <div id="meta">
2   <div></div><div class="inner">
3     <div>
4       <img alt="Portrait of Kareem Abdul-Jabbar" data=""/>
5       <div>
6         <strong>Kareem Abdul-Jabbar</strong>
7         <br/>
8         Presently: Hall of Fame inducted July 2009
9         Kareem Abdul-Jabbar - Wikipedia
10        (Formerly known as Ferdinand Lewis Alcindor Jr.)
11        Kew, Cap, Mbaak, Big Fella, The TOWER of Power
12        Position: Center - Shoots: Right
13        Height: 7'2" (238cm, 180kg)
14        Born: April 16, 1947 (Age: 77-102d) in New York, New York, USA
15        College: UCLA
16        High School: Power Memorial in New York, New York
17        Draft: #1 Overall, 1st Round, 1st Pick, 1st overall, 1st overall Draft
18        NBA Debut: 1969-70
19        Hall of Fame: Inducted as Player in 1996 (Full Bio)
20        Career Length: 20 years
21      </div>
22    </div>
23  </div>
24  <div></div>
25  <div></div>
26  <div></div>
27  <div></div>
28  <div></div>
29  <div></div>
30  <div></div>
31  <div></div>
32  <div></div>
33  <div></div>
34  <div></div>
35  <div></div>
36  <div></div>
37  <div></div>
38  <div></div>
39  <div></div>
40  <div></div>
41  <div></div>
42  <div></div>
43  <div></div>
44  <div></div>
45  <div></div>
46  <div></div>
47  <div></div>
48  <div></div>
49  <div></div>
```

List informácií o hráčovi, ktoré budeme uchovávať:

- meno (text)
- výška (cm)
- váha (kg)
- pozícia (text)
- dominantná ruka (text)
- vek (v rokoch)
- rok narodenia (dátum)
- miesto narodenia (štát)
- stredná škola (text)
- college (text)
- draftovaný tímom (text)
- rok draftovania (rok)
- debut (dátum)
- dĺžka kariéry (v rokoch)
- sieň slávy (uvedený v roku pridania do siene)

Ďalej budú pozostávať z jeho vlastností na ihrisku (body, doskoky, asistencie, bloky, obrátené lopty, fauly, percentuálna úspešnosť,...) zobrazené na basketball-reference [https://www.basketball-reference.com/] v takomto formáte:

Season	Age	Foe	Lg	Pos	G	GS	MF	FG	FGA	3P%	3PA%	3PM%	FT	FTA	TRB	DRB	ORB	AST	STL	BLO	TOV	PF	PTS	PER	
1999-2000	24	1999-2000	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2000-2001	24	1999-2000	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2001-2002	24	2000-2001	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2002-2003	24	2001-2002	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2003-2004	24	2002-2003	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2004-2005	24	2003-2004	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2005-2006	24	2004-2005	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2006-2007	24	2005-2006	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2007-2008	24	2006-2007	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2008-2009	24	2007-2008	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2009-2010	24	2008-2009	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2010-2011	24	2009-2010	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2011-2012	24	2010-2011	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2012-2013	24	2011-2012	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2013-2014	24	2012-2013	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2014-2015	24	2013-2014	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2015-2016	24	2014-2015	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2016-2017	24	2015-2016	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2017-2018	24	2016-2017	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2018-2019	24	2017-2018	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2019-2020	24	2018-2019	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2020-2021	24	2019-2020	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2021-2022	24	2020-2021	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2022-2023	24	2021-2022	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2023-2024	24	2022-2023	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2024-2025	24	2023-2024	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2025-2026	24	2024-2025	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2026-2027	24	2025-2026	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2027-2028	24	2026-2027	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2028-2029	24	2027-2028	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2029-2030	24	2028-2029	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2030-2031	24	2029-2030	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2031-2032	24	2030-2031	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2032-2033	24	2031-2032	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2033-2034	24	2032-2033	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2034-2035	24	2033-2034	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2035-2036	24	2034-2035	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2036-2037	24	2035-2036	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2037-2038	24	2036-2037	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2038-2039	24	2037-2038	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2039-2040	24	2038-2039	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2040-2041	24	2039-2040	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2041-2042	24	2040-2041	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2042-2043	24	2041-2042	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2043-2044	24	2042-2043	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2044-2045	24	2043-2044	N	S	81	41.2	22.4	21.1	31.8	30.1	29.9	9.1	81.0	14.2	4.1	1.2	22.8	9.9	1.2	1.2	1.2	1.2	1.2	24.8	10.5
2045-2046	24	2044-2045	N	S	81	41.2	22.4	21																	

Linky všetkých hráčov sa budú zbierať veľmi jednoducho, keďže stránka ich triedi podľa začiatočného písmena v ich priezvisku, viď tu:

NBA & ABA Players with Last Names Starting with A

Index of Letters: A B C D E F G H I J K L M N O P Q R S T U V W X Y Z

177 Players

Player	Pos.	Yr.	Hgt.	Wt.	Birth Date	College		
Masa-Akiba, Aya	FWD	1999	F	6-0	240	June 21, 1999	None	
Dent-Ash, Aya	SG	1998	C/F	6-4	200	April 21, 1998	None	
Kurashiki-Moto, Ayaka*	FWD	1998	F	6-2	199	April 21, 1998	None	
Rahman-Mohd, Azizal	FWD	1993	2005	6-0	181	(March 3, 1993)	LNU	
Tara-Abed, Azrael	FWD	1998	2000	F	6-6	221	(November 3, 1998)	Holyoke, San Jose State
Garcia-Lukas, Azrael	FWD	1997	2008	F	6-6	204	(December 14, 1997)	California
Guerrero, Azucena	FWD	1997	1999	F	6-1	179	(May 2, 1997)	None

Podľa linku viem určiť, že ktoré začiatočné písmeno prezerám (napr. pre hráčov so začiatočným písmenom "D" v priezvisku je určená linka <https://www.basketball-reference.com/players/> [<https://www.basketball-reference.com/players/>]d]) a táto linka mi zobrazí všetky odkazy na daný list hráčov v tabuľke, čiže sa iba zoberie tabuľka hráčov a element po elemente sa budú zbierať linky hráčov. Takýmto spôsobom sa dajú vyzbierať všetci NBA hráči.

Zdroje dát

- NBA štatistiky hráčov [<https://www.basketball-reference.com/players/>]
- Wikipedia [<https://www.wikipedia.com>]

Poznamky: zoznam vlastností hraca, ukoncene

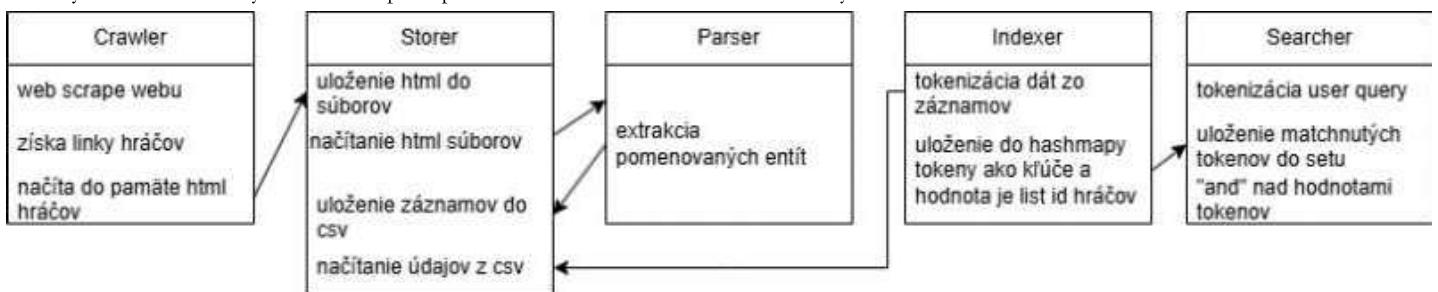
user/peter.bartos/konzultacia-1.txt · Last modified: 2024/10/07 12:37 by Administrator [/doku.php?id=user:admin]



NBA - druhá konzultácia

Architektúra

Crawler získava linky na hráčov, následne postahuje html hráčov z liniek a storer to uloží html do súborov. Tieto súbory sa načítajú v parser-i a ten vyextrahuje z nich elementy, ktoré sú zaujímavé. Tieto záznamy sa uložia do csv storera. Následne sa csv načítajú v indexer-i pomocou storera a jednotlivé záznamy sa začnú tokenizovať. Tieto tokeny sa uložia do hashmapy ako kľúče a vedú sa k nim listy id hráčov. V searcher-i sa zas táto hashmapu indexer-a zavolá, prebehne tokenizácia používateľského dopytu a jednotlivé kľúče z indexer-a sa matchujú a ukladajú do výsledku. Nad uloženými id hráčov vo výsledkoch sa spraví príenik a searcher mi dá vásledné záznamy.



Crawler

Web scraping hráčov zo stránky basketball-reference [https://www.basketball-reference.com/] je možné spraviť veľmi jednoducho, keďže ak sa použije linka basketball-reference/players/a [https://www.basketball-reference.com/players/a/], tak sa zobrazia iba tí hráči, ktoré majú začiatocne písmeno "a" v priezvisku (čiže linka sa dá meniť /players/{písmeno} a budeme vedieť postupne získať všetky linky hráčov, keď použijeme všetky písmená v abecede). Linky hráčov sú zobrazené v tabuľke s ich menom a ďalšími dodatočnými informáciami:

The screenshot shows the Basketball Reference homepage with a search bar and navigation menu. Below the menu, a banner mentions they are hiring a Senior DevOps Engineer. The main content area is titled "NBA & ABA Players with Last Names Starting with A". Below this, there is an "Index of Letters" with links for each letter from A to Z. A yellow oval highlights a table titled "177 Players" containing player statistics and college information. The table includes columns for Player, From, To, Pos, Ht, Wt, Birth Date, and Colleges. Several rows of player data are listed, such as Alaa Abdelnaby, Zair Abdul-Aziz, Kareem Abdul-Jabbar*, Mahmoud Abdul-Rauf, Tariq Abdul-Wahad, Shareef Abdur-Rahim, and Tom Abernethy.

Player	From	To	Pos	Ht	Wt	Birth Date	Colleges
Alaa Abdelnaby	1991	1995	F-C	6-10	240	June 24, 1968	Duke
Zair Abdul-Aziz	1969	1978	C-F	6-9	235	April 7, 1946	Iowa State
Kareem Abdul-Jabbar*	1970	1989	C	7-2	225	April 16, 1947	UCLA
Mahmoud Abdul-Rauf	1991	2001	G	6-6	162	March 9, 1969	LSU
Tariq Abdul-Wahad	1998	2003	F	6-6	223	November 3, 1974	Michigan, San Jose State
Shareef Abdur-Rahim	1997	2008	F	6-9	225	December 11, 1976	California
Tom Abernethy	1977	1981	F	6-7	220	May 6, 1954	Indiana

Získavanie linky hráča

Linky na hráčov budú tiež uložené do csv súbora, aby vždy nemusel prebiehať web scraping. Pseudokód na získavanie všetkých liniek na hráčov vyzerá nasledovne:

Download

```
Function get_players_links:  
    Create a list "alphabet"  
    Create empty list "player_page_links"  
    For each "letter" in "alphabet":  
        "url" to "https://www.basketball-reference.com/players/<letter>"  
        "list_of_players_html" to the HTML content from making a GET request to "url"  
        "pattern" regex matching player links based on HTML  
        "player_links" finding all matches using "pattern" on "list_of_players_html"  
        For each "player_link" in "player_links":  
            Append "player_link" to "player_page_links"  
            Store "player_link" to player links csv file  
        Sleep for 3 seconds # comply with the website's crawler policy  
    Return "player_page_links"
```

Pre získanie liniek pre hráčov by sa využili iba natívne knižnice Pythonu:

- “re” - regexy
- “time” - spánok na určitý čas
- “requests” - vykonávanie dopytov na stránky

Získavanie informácií o hráčovi

Pseudokód na získavanie informácií o všetkých hráčoch vyzerá nasledovne:

Download

```
Function get_player_page_html_content(player_page_link: string):  
    "player_page_html" is the HTML result of making a GET request to "player_page_link"  
    Remove special characters from HTML in "player_page_html"  
    Store "player_page_html" content into a file  
    Return  
  
Function get_pages_of_players(player_page_links: list[string]):  
    For each "player_page_link" in "player_page_links":  
        Call get_player_page_html_content function with "player_page_link"  
        Sleep for 3 seconds # comply with the website's crawler policy  
    Return
```

Pre získanie html všetkých hráčov by sa tiež využili iba natívne knižnice Pythonu:

- “os” - ukladanie súborov (os api)
- “re” - regexy
- “time” - spánok na určitý čas
- “requests” - vykonávanie dopytov na stránky

Parser

Z jednotlivých html súboroch, ktoré majú informácie o danom hráčovi je potrebné vypарsovať pomocou regexov informácie, ktoré sme si určili v prvej konzultácii ([konzultácia 1](#)). Vyparsované informácie z html sa budú ukladať do csv súborov pre ľahkú manipuláciu s nimi neskôr. Html časti, ktoré budeme parsovať vyzerajú nasledovne:

```

1 ...
2 <div id="meta">
3     </div><!-- div.media-item --><div>
4         <h1>
5             <span>Kareem Abdul-Jabbar</span>
6         </h1>
7         <!-- all other pages -->
8         <p><strong>Pronunciation</strong>: \kuh-REEM ab-doo\ juh-BAR</p>
9     <p>
10        <strong><strong>Kareem Abdul-Jabbar</strong></strong>
11    </p>
12    <p><span class="desc">(Formerly known as Ferdinand Lewis Alcindor Jr.)</span></p>
13    <p>(Lew, Cap, Murdock, Big Fella, The Tower of Power)</p>
14    <p><strong>Position:</strong>Center<strong>Shoots:</strong>Right</p>
15    <p><span>7-2</span>,&nbsp;<span>225lb</span>&nbsp;[218cm,&nbsp;102kg]</p>
16    <p><strong>Born: </strong><span id="necro-birth" data-birth="1947-04-16"></span> <span><nobr>(Age:&nbsp;77-167d)</nobr></span><span>
17        in&nbsp;New York,&nbsp;<a href="/friv/birthplaces.fcgi?country=US&amp;state=NY">New York</a></span>
18        <span class="f-i f-us" style="">us</span>
19    </p>
20    <p><strong>College:</strong><a href="/friv/colleges.fcgi?college=ucla">UCLA</a></p>
21    <p><strong>High School:</strong>Power Memorial in New York, <a href="/friv/high_schools.fcgi?country=US&amp;state=NY">New York</a></p>
22    <p>
23        <strong>Draft:</strong>
24        <a href="/teams/MIL/draft.html">Milwaukee Bucks</a>, 1st round (1st pick, 1st overall), <a href="/draft/NBA_1969.html">1969 NBA Draft</a>
25    </p>
26    <p><strong>NBA Debut: </strong><a href="/boxscores/19691018MIL.html">October 18, 1969</a></p>
27    <p><strong>Hall of Fame:</strong> Inducted as Player in 1995 (<a href="/awards/hof.html">Full List</a>)</p>
28    <p><strong>Career Length:</strong>&nbsp;20 years</p>
29 </div>
30 ...

```

```

1 <tbody>
2     <tr id="per_game.1991" class="full_table" data-row="0">
3         <th scope="row" class="left " data-stat="season"><a href="/players/a/abdelal01/gamelog/1991">1990-91</a></th>
4         <td class="center " data-stat="age">22</td>
5         <td class="left " data-stat="team_id"><a href="/teams/POR/1991.html">POR</a></td>
6         <td class="left " data-stat="lg_id"><a href="/leagues/NBA_1991.html">NBA</a></td>
7         <td class="center " data-stat="pos">PF</td><td class="right " data-stat="g">43</td>
8         <td class="right iz" data-stat="gs">0</td><td class="right " data-stat="mp_per_g">6.7</td>
9         <td class="right " data-stat="fg_per_g">1.3</td><td class="right " data-stat="fga_per_g">2.7</td>
10        <td class="right " data-stat="fg3_per_g">0.0</td><td class="right iz" data-stat="fg3_pct">0.0</td>
11        <td class="right iz" data-stat="fg2_per_g">0.0</td><td class="right " data-stat="fg2a_per_g">2.7</td>
12        <td class="right " data-stat="fg2_pct">1.3</td><td class="right " data-stat="efg_pct">.474</td>
13        <td class="right " data-stat="ft_per_g">0.6</td><td class="right " data-stat="fta_per_g">1.0</td>
14        <td class="right " data-stat="ft_pct">.568</td><td class="right " data-stat="orb_per_g">0.6</td>
15        <td class="right " data-stat="drb_per_g">1.4</td><td class="right " data-stat="trb_per_g">2.1</td>
16        <td class="right " data-stat="ast_per_g">0.3</td><td class="right " data-stat="stl_per_g">0.1</td>
17        <td class="right " data-stat="blk_per_g">0.3</td><td class="right " data-stat="tov_per_g">0.5</td>
18        <td class="right " data-stat="pf_per_g">0.9</td><td class="right " data-stat="pts_per_g">3.1</td>
19        <td class="left iz" data-stat="award_summary"></td>
20     </tr>
21 </tbody>

```

Pre parsovanie informácií o hráčovi je použitý nasledujúci pseudokód:

[Download](#)

```

Function parse_player_info_from_html(player_page_html: string):
    'name' to the result of pattern '<strong>\s*<strong>(.*)</strong>\s*</strong>' in 'player_page_html'
    'height' to the result of pattern '\(\s*(\d+)cm' in 'player_page_html'
    'weight' to the result of pattern ',\s*(\d+)kg' in 'player_page_html'
    'position' to the result of pattern in 'player_page_html'
    'shoots' to the result of pattern '(?i)Shoots:\s*</strong>\s*(right|left)' in 'player_page_html'
    'born' to the result of pattern 'data-birth="(\d{4})(?:-(\d{2}))?(?:-(\d{2}))?"' in 'player_page_html'
    'died' to the result of pattern 'data-death="(\d{4})(?:-(\d{2}))?(?:-(\d{2}))?"' in 'player_page_html'
    'born_in_city' to the result of pattern '/friv/birthplaces\.fcgi\?country=(\[^&]+)' in 'player_page_html'
    'high_school' to the result of pattern 'High School[s]?:\s*</strong>\s*([^\<]+)\s*<' in 'player_page_html'
    'college' to the result of pattern 'College[s]?:\s*</strong>\s*<a\s*href=[\'"]([^\"]*)["]\s*>' in 'player_page_html'
    'nba_debut' to the result of pattern 'NBA Debut:\s*</strong>\s*<a href="[^"]*"(.*)</a>' in 'player_page_html'
    'career_length' to the result of pattern '(?:Career Length|Experience):\s*</strong>\s*(\d+ years)' in 'player_page_html'
    'hall_of_fame_year' to the result of pattern 'Inducted as Player in (\d+)' in 'player_page_html'
    Store variables into a player info csv file

```

Pre parsovanie vlastností hráča v jednotlivých sezónach je použitý nasledujúci pseudokód:

[Download](#)

```

Function parse_player_stats_from_html(player_page_html: str):
    pattern as '<tr id="per_game">\d+.*?>.*?</tr>' and store all occurrences in a list "per_game_rows_html"
    For each "per_game_row_html" in "per_game_rows_html":
        "season" to the result of pattern '"season"\s*>\s*<\s*a\s*href=[^"]+"\\s*>\s*([^\n]+)\s*<'
        "team" to the result of pattern '"team_id"\s*>\s*<\s*a\s*href=[^"]+"\\s*>\s*([^\n]+)\s*<'
        "league" to the result of pattern '"lg_id"\s*>\s*<\s*a\s*href=[^"]+"\\s*>\s*([^\n]+)\s*<'
        "pos" to the result of pattern '"pos"\s*>\s*([^\n]+)\s*<'
    Store variables into player stats csv file

```

Toto parsovanie sa vykoná pre každé html hráča a tým sa získa v csv súboroch databáza informácií, s ktorou sa bude ďalej pracovať.

Pre parsovanie html súborov by sa tiež využili iba natívne knižnice Pythonu:

- “os” - ukladanie súborov (os api)
- “re” - regexy
- “datetime” - dátumové operácie

Indexer

Terajšia verzia vlastného indexu sa robí iba nad atribútmi “name”, “position”, “shoots”, “college”, “high_school”, “life_status”, “season”, “team” a “born_in_country”. Ako index sa zadefinuje hash mapa a kľúče v nej sú jednotlivé tokeny (slová). Funguje to tak, že sa najprv prečítajú všetky záznamy z info o hráčoch csv súboru. Po načítaní sa ide záznam po zázname a tokenizujú sa jednotlivé atribúty (iba podľa medzery sa rozdelia slová do listu). Ak sa token nachádza v indexe, tak sa ako hodnota k tomuto kľúču pripíše do listu id hráča. Tým pádom keď vo vyhľadávači chceme hľadať podľa nejakého slova, tak sa mi skontroluje index pre dané slovo (skontroluje sa, či v indexe kľúč existuje) a keď sa dané slovo v indexe nájde, tak sa mi vráti list id-čiek všetkých hráčov, kde sa toto slovo nachádza. Následne sa to potom zoradí podľa vypočítaného tf-idf. TF a IDF sa počíta za behu indexácie a následne zbehne ešte ďalší cyklus cez tokeny, ktorý vypočíta TF-IDF a zapíše do indexu ako hodnotu (doc id, tf-idf-value).

Dáta v csv vyzerajú takto:

ID	Name	Height	Weight	Position	Shoots	Age	Born	Life_Status	Born_in_Country	High_School	College	NBA_Draft_Team	NBA_Draft_Year	NBA_Debut	Career_Length	Hall_of_Fam
0	Jaylen Tateque Ami	183cm	102kg	Point Guard	Righty	28yo	1990-05-04	alive	US	Mount St. Joseph	St. John's			2018-11-17	3 years	
1	Paul Carlyle Ami	180cm	77kg	Guard	Righty	64yo	1918-11-01	deceased	US	Central in Fort Wayne	Indiana			1948-11-03	3 years	
2	Mark Richard Ami	211cm	99kg	Center	Righty	61yo	1962-11-15	alive	US	Palo Verde in El Dorado	Dallas Mavericks	1985	1987-11-06	6 years	2010	
3	Chester J. Aubut	178cm	62kg	Guard	Righty	38yo	1916-05-18	deceased	US	Horace Mann in Michigan				1946-11-02	1 year	
4	James Randall Ami	203cm	99kg	Shooting G	Righty	38yo	1965-01-26	alive	US	Milton in Milton	Florida			1989-03-09	2 years	
5	Andrew Emil Ami	188cm	83kg	Shooting G	Righty	74yo	1845-07-06	deceased	US	Maryvale in Cheyenne	Catholic	Boston Celtics	1967		3 years	
6	Luke Robert Ami	205cm	102kg	Small Forward	Lefty	35yo	1989-06-20	alive	US	Galen in Reno	Nevada	Minnesota Timberwolves	2010	2010-11-01	8 years	
7	Kyle F. Anderson	206cm	104kg	Small Forward	Righty	31yo	1993-09-20	alive	US	St. Anthony in Los Angeles	UCLA	San Antonio Spurs	2014	2014-11-06	10 years	
8	LuMark Anthony	185cm	79kg	Point Guard	Righty	54yo	1969-11-11	alive	US	Paul Laurence Dunbar in Philadelphia				1999-03-06	1 year	
9	Thomas Craig Ami	201cm	99kg	Small Forward	Righty	71yo	1954-05-06	alive	US	Saint Joseph in Indianapolis	Indiana	Los Angeles Lakers	1976	1976-11-19	5 years	
10	Cole Hinton Ami	188cm	83kg	Point Guard	Righty	24yo	2000-05-15	alive	US	Archbishop Mitty	Orlando Magic	2020	2020-12-23	4 years		
11	Rafael Alderson	201cm	97kg	Small Forward	Righty	60yo	1964-07-22	alive	US	Snyder in Jersey	Syracuse	Phoenix Suns	1986	1986-11-01	6 years	
12	Brandon Simone Ami	196cm	85kg	Shooting G	Lefty	44yo	1980-06-16	alive	US	Vallejo in Vallejo	Pepperdine	Houston Rockets	2001	2001-11-02	3 years	
13	Alynn Austin Ami	183cm	79kg	Point Guard	Righty	37yo	1983-11-07	deceased	US	Weequahic in Newark	New Jersey	Philadelphia 76ers	1980	1980-10-28	11 years	

Pseudokód pre tento indexer vyzerá nasledovne:

[Download](#)

```

function build_player_info_index():
    "index" as hashmap containing key as token to value represented by list of player ids
    "tf" list for storage of term frequencies
    "idf" hashmap to store idf
    "players_doc_tokens" hashmap that stores all the unique tokens of each document
    open player info csv file
        read all lines of the csv and put into a list
        for each "row" in "rows":
            extract "player_id" from row's id column
            for each "player_property" in "PROPERTIES_IN_INDEX":
                if "player_property" in "row" is not empty:
                    "tokens" is set from splitting "player_property" in "row" by space and lowered
                    for each "token" in "tokens":
                        get the value from the "index" hashmap using "token" as key and append "player_id" to value
                    calculate term frequency and store in "tf"
                    calculate "idf" after aquiring all the tokens for current doc
            for each "row" in "rows":
                for each "token" in "players_doc_tokens"[ "doc id"]:
                    calculate "tf-idf" value
                    store the value in index for this "token"
    store "index" in csv file

```

Pre indexovanie by sa tiež využili iba natívne knižnice Pythonu:

- “collections” - pre hashmapu reprezentujúcu index
- “csv” - čítanie csv súboru

Searcher

Vlastný vyhľadávač je tiež veľmi jednoduchý. Rozdelí vyhľadávací dotaz na tokeny (čiže slová, teda rozdelí string podľa medzier do listu) a tieto tokeny hľadá v indexe (do indexu je poskytnutý token ako kľúč). Ak sa tento kľúč nachádza v indexe, tak sa uloží do výsledkov. Následne sa spraví prienik nad všetkými výsledkami a ten sa uloží do premennej. Teraz prejde index a hľadá podľa hodnôt z výsledkov a zapisuje si čo najvyššie tf-idf s doc id do listu s neusporiadanými výsledkami. Keď prejde index, tak neusporiadane výsledky usporiada podľa tf-idf a vráti výsledky usporiadane podľa tf-idf zostupne.

[Download](#)

```
function search(index, query):
    "query_tokens" is provided by tokenizing "query"
    "results" is initialized as a list
    for each "token" in "query_tokens":
        if "token" exists in "index":
            add "index"[ "token" ] to results
    if "results" is not empty:
        "unique_results" is intersection of all values in results
        gets "unsorted_results" by calculating score for each item in "unique_results"
        returns "sorted_results" that are sorted by highest score from "unsorted_results"
    else:
        return empty results
```

Pre vyhľadávač by sa tiež využili iba natívne knižnice Pythonu:

- “re” - regexy

Poznamky: crawler 5000 +/-, prerobiť reprezentáciu zaznamu/dokumentu, index na disku, odstrániť duplikáty v teams, skúsiť druhý index pre unikátne tímy v season_team, term id miesto termov

NBA - odovzdanie 1

Crawler & Parser

Dáta po použití crawlra:

html	
abdelal01.html	
abdulk01.html	
abdulma02.html	
abdulta01.html	
abdulza01.html	
abdursh01.html	
abernto01.html	
abjefo01.html	
abramjo01.html	
abrinalo01.html	
achiupry01.html	
ackeral01.html	
ackerdo01.html	
acresma01.html	
actonbu01.html	
acyqu01.html	
adamsal01.html	
adamsdo01.html	

Dáta po použití parsera:

Name	Height	Weight	Position	Shoots	Age*	Born	Deceased	Born_in	High_school	College	Nba_draft_team	Nba_draft_year	Nba_debut	Career_length	Hall_of_fame_Y	Seasons	Teams	Season_team
Malcolm	206	112	Center	Right	34	1970-1	False	US	Brainerd in ci galena	Phoenix Suns	1993	1994-01-13	1 year		1987	SDR	1967_SDR	
Sam Mac	201	99	Small Forw	Right	34	1970-1	False	US	Thornridge in iowa			1992-11-06	7 years		1956	SYR	1956_SYR	
Donald	208	106	Small Forw	Right	34	1970-1	False	US	Semi Valley in okla	Detroit Pistons	1992	1992-11-06	9 years		1946-1947_19_PIT_STB_BLB	1946_PIT_1947_STB		
Richard A.	211	114	Power Forw	Right	34	1970-1	False	US	Center in Ariz	Atlanta Hawks	1993	1996-01-27	2 years		2012-2013_20_TOR_TOR_SA	2012_TOR_2013_SA		
Robert W.	213	113	Center	Right	34	1969-1	False	US	Apple Valley in minnesot			1993-11-07	2 years		1953	NYK	1953_NYK	
Diva Lee	165	78	Point Guar	Right	34	1970-1	False	US	Will Rogers in arkansas	Milwaukee Bucks	1992	1992-11-06	7 years		1976-1977_19_LAL_LAL_GSN	1976_LAL_1977_LA		
Oliver J. R.	206	127	Center	Right	34	1970-1	False	US	Southwest in arkansas	Phoenix Suns	1992	1992-11-07	9 years		2006-2009	NIN_TOR	2006_NIN_2009_TOR	
Christopher	196	97	Small Forw	Right	34	1970-1	False	US	Fairfax in kentucky	Cleveland Cavalier	1993	1993-11-05	10 years		2016-2017_20_CMC_CMC_OH	2016_OHC_2017_O		
Alfonzo H.	208	100	Center	Right	34	1970-1	False	US	Indian River in georgia	Charlotte Hornets	1992	1992-11-13	15 years	2014	1987-1988_19_ICOS_BOS_CFR	1987_BOS_1988_IC		
Todd Max	213	113	Center	Right	34	1970-1	False	US	Central Merry memphis			1995-11-04	1 year		1995-1991_19_POR_POR_MI	1990_POR_1991_POR		
Terrell Jer	208	107	Power Forw	Right	34	1969-1	False	US	DeMatha Cat in maryland	New York Knicks	1990	1990-11-02	4 years		2020-2021_20_MIA_TOR_TO	2020_MIA_2021_TO		
Anthony J.	193	94	Shooting C	Left	34	1969-1	False	US	Paisley in kane missouri	Los Angeles Laker	1992	1992-11-06	13 years		1990-1991_19_DEN_DEN_OIL	1990_DEN_1991_OIL		
Eric Todd	201	97	Shooting C	Right	34	1970-1	False	US	Stevens in nebraska	Indiana Pacers	1994	1994-11-04	14 years		1968-1969_19_CIN_MIL_MIL	1968_OHC_1969_MIL		
David Wil	201	99	Small Forw	Right	34	1962-1	True	US	Menchville in northcarolina	Utah Jazz	1984	1984-12-09	2 years		1996-1997_19_VAN_VAN_WF	1996_VAN_WF_1997_VAN		
Raymond	193	92	Small Forw	Right	34	1928-	True	US	Washington in newstern	Rochester Royals	1951	1951-11-01	1 year		2005-2009_20_DET_DET_LAC	2005_DET_2008_LAC		
Richard Jr.	201	97	Power Forw	Right	34	1933-	True	US	Pottstown in duquesne	St. Louis Hawks	1955	1953-11-05	3 years		1987-1988_19_SAC_SAC_OR	1987_SAC_1988_OR		

Indexovanie (TF-IDF)

Termy sú poukladané podľa ich id a value:

```
1 id;term
...
3190 3188;greenville
3191 3189;greenway
3192 3190;greenwood
3193 3191;greer
3194 3192;greg
3195 3193;gregg
3196 3194;gregor
3197 3195;gregory
3198 3196;greig
3199 3197;greivis
3200 3198;grekin
3201 3199;grenadines
3202 3200;gresham
3203 3201;grevey
3204 3202;grey
3205 3203;griffin
3206 3204;griffith
3207 3205;grigsby
3208 3206;grimes
3209 3207;grimm
3210 3208;grimshaw
3211 3209;grimsley
3212 3210;grizzlies
3213 3211;groat
```

Prc **TF** som vyskúšal dva prístupy: prirodzený a logaritmický.

Prirodzený vyzerá nasledovne:

[Download](#)

```
...
{
  'joe': 0.037, 'depre': 0.037, 'shooting': 0.037, 'guard': 0.037, 'righty': 0.037,
  'stjohns': 0.037, 'westbury': 0.074, 'in': 0.037, 'alive': 0.037, '1970': 0.037,
  '1971': 0.037, '1972': 0.037, 'new': 0.11, 'york': 0.11, 'americans': 0.11,
  'united': 0.037, 'states': 0.037
},
```

Príklad pre logaritmický:

[Download](#)

```
...
{
  'mike': -2.89, 'theodore': -2.89, 'glenn': -2.89, 'shooting': -2.89, 'guard': -2.89,
  'righty': -2.89, 'sillinois': -2.89, 'coosa': -2.89, 'in': -2.89, 'rome': -2.89,
  'alive': -2.89, '1977': -2.89, '1978': -2.89, '1979': -2.89, '1980': -2.89,
  '1981': -2.89, '1982': -2.89, '1983': -2.89, '1984': -2.89, '1985': -2.89,
  '1986': -2.89, 'buffalo': -2.89, 'braves': -2.89, 'new': -1.79, 'york': -1.79,
  'knicks': -1.79, 'atlanta': -1.50, 'hawks': -1.50, 'milwaukee': -2.19,
  'bucks': -2.19, 'united': -2.89, 'states': -2.89
},
```

Prc **DF** som vyskúšal tiež dva prístupy: IDF a probabilistic IDF.

Príklad pre IDF:

[Download](#)

```
...
'mike': {
  1: 0.032, 24: 0.032, 37: 0.032, 128: 0.064, ...
},
```

Príklad pre probabilistic IDF vyzerá približne rovnako ako IDF len hodnoty sa menia:

[Download](#)

```
...
'mike': {
  1: 0.02, 24: 0.02, 37: 0.02, 128: 0.061, ...
},
```

Vypočítavam **TF-IDF** klasicky jednoduchým násobením:

[Download](#)

```
...
# calculate tf_idf value for each token
tf_idf[token] = tf[player_id].get(token, 0) * idf[token]
```

Následne to aj sortujem, keď TF-IDF vraciám z funkcie:

[Download](#)

```
...
# sorting the terms in the dict before returning
return dict(sorted(tf_idf.items()))
...
```

Skóre následne počítam **nasčítavaním**:

[Download](#)

```
...
# go through terms in query
for term in query_tokens:
    # if term is in terms dict
    if term in term_dict:
        # then add the score
        term_id = term_dict[term]
        score += index.get(term_id, 0)
...
```

Výsledky

Pri využití prirodzeného TF a IDF (v obr. naľavo) vyzerajú výsledky pre dopyt "michael chicago" trošku inak ako pri použití logaritmického TF a IDF (v obr. napravo):

```
1 1: | Michael Jeffrey Jordan | 61yo | alive |
2 2: | Michael Damien Sweetney | 41yo | alive |
3 3: | Alex Michael Caruso | 30yo | alive |
4 4: | Kristofer Michael Dunn | 30yo | alive |
5 5: | Aaron Michael Gray | 39yo | alive |
6 6: | Michael Edward Harper | 66yo | alive |
7 7: | Michael George Williams | 61yo | alive |
8 8: | Michael Donald Novak | 63yo | deceased |
9 9: | Michael David Ruffin | 47yo | alive |
10 10: | Andrew Michael Phillip | 79yo | deceased |
```

```
1 1: | Michael Jeffrey Jordan | 61yo | alive |
2 2: | Michael Damien Sweetney | 41yo | alive |
3 3: | Michael Edward Harper | 66yo | alive |
4 4: | Michael George Williams | 61yo | alive |
5 5: | Michael Donald Novak | 63yo | deceased |
6 6: | Alex Michael Caruso | 30yo | alive |
7 7: | Kristofer Michael Dunn | 30yo | alive |
8 8: | Aaron Michael Gray | 39yo | alive |
9 9: | Larry Michael Spriggs | 65yo | alive |
10 10: | Michael David Ruffin | 47yo | alive |
```

Ked' použijem prirodzený TF probabilistic IDF (v obr. napravo) tak pre dopyt "michael chicago" vyzerajú výsledky znova trošku inak oproti prirodzenému TF a IDF (v obr. naľavo):

```
1 1: | Michael Jeffrey Jordan | 61yo | alive |
2 2: | Michael Damien Sweetney | 41yo | alive |
3 3: | Alex Michael Caruso | 30yo | alive |
4 4: | Kristofer Michael Dunn | 30yo | alive |
5 5: | Aaron Michael Gray | 39yo | alive |
6 6: | Michael Edward Harper | 66yo | alive |
7 7: | Michael George Williams | 61yo | alive |
8 8: | Michael Donald Novak | 63yo | deceased |
9 9: | Michael David Ruffin | 47yo | alive |
10 10: | Andrew Michael Phillip | 79yo | deceased |
```

```
1 1: | Michael Jeffrey Jordan | 61yo | alive |
2 2: | Michael Damien Sweetney | 41yo | alive |
3 3: | Alex Michael Caruso | 30yo | alive |
4 4: | Kristofer Michael Dunn | 30yo | alive |
5 5: | Michael Edward Harper | 66yo | alive |
6 6: | Michael George Williams | 61yo | alive |
7 7: | Aaron Michael Gray | 39yo | alive |
8 8: | Michael Donald Novak | 63yo | deceased |
9 9: | Michael David Ruffin | 47yo | alive |
10 10: | Andrew Michael Phillip | 79yo | deceased |
```

Po použití logaritmického TF a probabilistic IDF (v obr. napravo) tak pre dopyt "michael chicago" vyzerajú výsledky úplne inak oproti prirodzenému TF a IDF (v obr. naľavo):

```
1 1: | Michael Jeffrey Jordan | 61yo | alive |
2 2: | Michael Damien Sweetney | 41yo | alive |
3 3: | Alex Michael Caruso | 30yo | alive |
4 4: | Kristofer Michael Dunn | 30yo | alive |
5 5: | Aaron Michael Gray | 39yo | alive |
6 6: | Michael Edward Harper | 66yo | alive |
7 7: | Michael George Williams | 61yo | alive |
8 8: | Michael Donald Novak | 63yo | deceased |
9 9: | Michael David Ruffin | 47yo | alive |
10 10: | Andrew Michael Phillip | 79yo | deceased |
```

```
1 1: | Michael Damien Sweetney | 41yo | alive |
2 2: | Michael Jeffrey Jordan | 61yo | alive |
3 3: | Michael Edward Harper | 66yo | alive |
4 4: | Michael George Williams | 61yo | alive |
5 5: | Michael Donald Novak | 63yo | deceased |
6 6: | Alex Michael Caruso | 30yo | alive |
7 7: | Kristofer Michael Dunn | 30yo | alive |
8 8: | Aaron Michael Gray | 39yo | alive |
9 9: | Larry Michael Spriggs | 65yo | alive |
10 10: | Michael David Ruffin | 47yo | alive |
```

Tiež som pre tento dopyt "michael chicago" skúšal rátať vypísaný tím pre danú sezónu ako slovo iba raz, aby to nefavorizovalo hráčom, ktorý odohral za daný tím viac sezón, preto v predošlých obrázkoch vidime medzi prvými pozíciami vždy "Michael Jeffrey Jordan", lebo za "Chicago Bulls" odohral veľa sezón. Po vyskúšaní výsledky dopytu vyzerajú nasledovne:

```
1 1: | Michael Jeffrey Jordan | 61yo | alive |
2 2: | Michael Damien Sweetney | 41yo | alive |
3 3: | Alex Michael Caruso | 30yo | alive |
4 4: | Kristofer Michael Dunn | 30yo | alive |
5 5: | Aaron Michael Gray | 39yo | alive |
6 6: | Michael Edward Harper | 66yo | alive |
7 7: | Michael George Williams | 61yo | alive |
8 8: | Michael Donald Novak | 63yo | deceased |
9 9: | Michael David Ruffin | 47yo | alive |
10 10: | Andrew Michael Phillip | 79yo | deceased |
```

```
1 1: | Michael Edward Harper | 66yo | alive |
2 2: | Michael George Williams | 61yo | alive |
3 3: | Michael Donald Novak | 63yo | deceased |
4 4: | Michael Damien Sweetney | 41yo | alive |
5 5: | Larry Michael Spriggs | 65yo | alive |
6 6: | Alex Michael Caruso | 30yo | alive |
7 7: | Michael David Holton | 63yo | alive |
8 8: | Aaron Michael Gray | 39yo | alive |
9 9: | Kristofer Michael Dunn | 30yo | alive |
10 10: | Michael Jeffrey Jordan | 61yo | alive |
```

Zdrojový kód

[xbartosp2_odovzdanie_1.zip](#)

NBA - štvrtá konzultácia

PySpark

Celý wikidump treba načítať' postupne, aby program nespadol s nedostatkom pamäte. To program robí tak, že číta súbor pomocou natívnej knižnice pre python bz2 ešte komprimovaný. Ďalej si zadefinuje pole "pages" a postupne číta súbor riadok po riadku a hľadá "<page>" a "</page>" tagy. Keď nájdete "<page>", tak vie, že začína stránka, tak to appenduje do stringu "page" každý prečítaný riadok. Keď nájde tag "</page>", tak končí appendovanie stringu, uloží "page" string do poľa "pages", premaže "page" a ide znova. Teraz sa nazbiera určitý počet stránok v poli "pages" a ten počet je definovaný premennou "pages_size", čo predstavuje chunk stránok. Pseudokód vyzerá nasledovne:

[Download](#)

```
Function read_wikidump_in_chunks(wikidump_path: string, pages_size=5000):
    Open file reading with bz2 as "FILE":
        Init "pages" list
        Init "page" string
        Init "is_in_page" boolean
        For "line" in "FILE":
            If length of "pages" == "pages_size":
                Process pages further
                Reset "pages" list
            If "<page>" in "line" and "is_in_page" is False:
                Reset "page" string
                Set "is_in_page" to True
            Elif "</page>" not in "line" and "is_in_page" is True:
                Append "line" to "page" string
            Elif "</page>" in "line" and "is_in_page" is True:
                Append "page" string to "pages" list
                Set "is_in_page" to False
```

Ked' sa teda nazbiera určitý počet stránok v poli, tak sa ide najprv určiť, že či je stránka vôbec relevantná. Ak je relevantná, tak sa spustia regexy na extrakciu nových entít ("Rookie of the year", "NBA Most Valuable Player", "NBA Champion", "Hall of Famer"), ktoré doplnia predošlé dátá, viď obrázok:



```
1930 | {{Infobox basketball biography
1931 |   | name = Kobe Bryant
1932 |   | image = Kobe Bryant 2014.jpg
1933 |   | image_size = 240
1934 |   | caption = Bryant with the [[Los Angeles Lakers]] in 2014
1935 |   | alt = Bryant handling the basketball
1936 |   | birth_date = {{birth date|1978|8|23}}
1937 |   | birth_place = [[Philadelphia, Pennsylvania]], U.S.
1938 |   | death_date = {{death date and age|2020|1|26|1978|8|23}}
1939 |   | death_place = [[Calabasas, California]], U.S.
1940 |   | resting_place = [[Pacific View Memorial Park]]
1941 |   | height_ft = 6
1942 |   | height_in = 6
1943 |   | height_footnote = <!--SEE DISCUSSION AT Talk:Kobe Bryant#Height edit-->{{#tag:ref|In 2017, Vanessa Bry
1944 |   | weight_lb = 212
1945 |   | high_school = [[Lower Merion High School|Lower Merion]]<br />([[Ardmore, Pennsylvania]])
1946 |   | draft_year = 1996
1947 |   | draft_round = 1
1948 |   | draft_pick = 13
1949 |   | draft_team = [[Charlotte Hornets]]
1950 |   | career_start = 1996
1951 |   | career_end = 2016
1952 |   | career_position = [[Shooting guard]]<!--Primarily a SG. He only played SF one season out of 20.-->
1953 |   | career_number = 8, 24
1954 |   | years1 = {{nbay|1996|start}}-{{nbay|2015|end}}
1955 |   | team1 = [[Los Angeles Lakers]]
1956 |   | highlights = <!-- See talk page at [[Talk:Kobe Bryant#Infobox highlights]] re: Academy Award -->
1957 |   | * 5x [[List of NBA champions|NBA champion]] ({{nbafy|2000}}-{{nbafy|2002}}, {{nbafy|2009}}, {{nbafy|2010}})
1958 |   | * 2x [[NBA Finals Most Valuable Player Award|NBA Finals MVP]] ({{nbafy|2009}}, {{nbafy|2010}})
1959 |   | * [[NBA Most Valuable Player Award|NBA Most Valuable Player]] ({{nbay|2007|end}})
1960 |   | bbr = bryanko01
```

Najprv sa ale toto pole stránok Sparkom paralelne spracuje a pomocou flatMap funkcie sa regexom vyfiltrujú iba relevantné stránky. Ďalej sa flatMapom a filtrom z tohto RDD vyextrahujú relevantné informácie pre hráčov NBA a tie sa následne uložia do CSV. Pseudokód vyzerá nasledovne:

[Download](#)

```
Function extract_relevant_records(text: str):
    Gets "matches" from all regex matches "{{Infobox[\s\S]basketball[\s\S]biography[\s\S]*?}}[ \t\r\n]*''"
    Return empty list if no "matches", otherwise return "matches"

Function extract_additional_info(relevant_page: str, players_names: list[str]):
    Init "player_name" from searching via regex "\|s*name\s*=s*\|[^|\n|]+"
    If "player_name" from "relevant_page" is in "players_names":
        Extract "roty" via regex "\*[\s\S]?(\d+)?x?[\s\S]?[\[\].*?NBA Rookie of the Year\]\]"
        Extract "mvp" via regex "\*[\s\S]?(\d+)?x?[\s\S]?[\[\].*?NBA Most Valuable Player\]\]"
        Extract "champ" via regex "\*[\s\S]?(\d+)?x?[\s\S]?[\[\].*?NBA champion\]\]"
        Extract "hall_of_fame" via regex "\| HOF_player"

Function wikidump_parse_pages(pages: list[str], players_names: list[str]):
    Get "wiki_rdd" that is initialized through sparks parallelize function on "pages" list
    Get "relevant_pages" using sparks flatMap on "wiki_rdd" and function "extract_relevant_records"
    Get "additional_information" using sparks flatMap and filter on "relevant_pages" using "extract_additional_info"
    Write "additional_information" into CSV dump
```

V tejto časti boli použité knižnice Pythonu:

- "re" - regexy (natívna)

- “os” - na narábanie s OS (natívna)
- “bz2” - na čítanie komprimovaných súborov (natívna)
- “pyspark”

Lucene

Indexer

Najprv sa začne čítať CSV s uloženými dátami z wiki a použije sa knižnica “csv” pre načítavanie riadkov rovno do dictionary, aby sa nemusela robíť logika na parsovanie fieldov. Teraz pôjde cyklus záznam po zázname a na každý záznam sa vytvorí lucene dokument a k tomu dokumentu sa budú pridávať jednotlivé fieldy. Textové fieldy sú zadefinované pre “name”, “position”, “shoots”, “born_country”, “highschool”, “college”, “nba_roty”, “nba_mvp”, “nba_champion”, “seasons”, “teams” a “hall_of_fame”. String fieldy sú zadefinované pre “weight”, “height”, “age”, “life_status”. Pseudokód pre indexer vyzerá nasledovne:

[Download](#)

```
Function index_wiki_data():
    Init "WRITER" as lucene's index writer
    With open read file with raw wiki player info as "FILE"
        Init "READER" using "csv" lib dict reader
        For "player" in "READER":
            Init "doc" as lucene document for player
            Add "name" field as TextField to "doc"
            Add "height" field as StringField to "doc"
            ...
            Add "teams" field as TextField to "doc"
        Using "WRITER" add "doc" to index
    Commit with "WRITER"
    Close "WRITER"
```

V tejto časti boli použité knižnice Pythonu:

- “os” - na narábanie s OS (natívna)
- “csv” - na narábanie s CSV súborom (natívna)
- “lucene”

Searcher

Používateľ poskytne query na vyhľadávanie. Najprv sa inicializuje lucene query parser, kde sa mu poskytnú všetky fieldy a štandardný analyzér. Ďalej z query parsera sa zavolá funkcia parse, do ktorej ako vstupné parametre ide používateľská query, list fieldov, list flagov pre jednotlivé fieldy a analyzér. Potom nad lucene index searcherom sa zavolá funkcia search a ako vstupné parametre ide používateľova query a číslo reprezentujúce aký počet výsledkov sa vráti. Z výsledkov sa zavolá atribút scoreDocs, ktorý predstavuje jednotlivé hity. Následne v cykle sa z týchto výsledkov vytahuje doc id a vytiahne sa celý dokument z lucene search indexera. Nakoniec sa tieto výsledky zobrazia používateľovi. Pseudokód pre searcher vyzerá nasledovne:

[Download](#)

```
Function search(query: str, fields: list[str], flags: list, result_num=10):
    Init "SEARCHER" as lucene's index searcher
    Init "ANALYZER" as lucene's standard analyzer
    Init "query_parser" as lucene's query parser with params "fields" and "ANALYZER"
    Init "query" with "query_parser" parse function with params "query", "fields", "flags" and "ANALYZER"
    Init "hits" with "SEARCHER" search function with params "query" and "result_num"
    Init "results" list to store result documents
    For "hit" in "hits":
        Init "doc" with "SEARCHER" document function with param "hit" using attribute "doc"
        Append "doc" into "results" list
    Return "results" list
```

Vzorka výstupu vyhľadávača:

```
SEARCHER search center champion
(+)
Listing found results
| NAME | AGE | VITAL |
1: | Miloš Babíč | 55yo | alive |
2: | Radislav Čurčić | 59yo | alive |
3: | Darko Miličić | 39yo | alive |
4: | Nikola Jokić | 29yo | alive |
5: | Rastko Cvetković | 54yo | alive |
6: | Vlade Divac | 56yo | alive |
7: | Nenad Krstić | 41yo | alive |
8: | Boban Marjanović | 36yo | alive |
9: | Filip Petrušev | 24yo | alive |
10: | Miroslav Raduljica | 38yo | alive |
(-) Write number to lookup from results: 5
(+)
Detailed player result
| NAME | AGE | POSITION | SHOOTS | HEIGHT | WEIGHT | HIGHSCHOOL | COLLEGE | BORN | VITAL | ROTY | MVP | CHAMPION | HOF | SEASONS |
5: | Nikola Jokić | 29yo | Center | Righty | 211cm | 128kg | Unknown | Unknown | Serbia | alive | No | Yes | Yes | Not a member | 2015_DEN 2016_DEN 2017_DEN 2018_DEN 2019_DEN 2020_DEN 2021_DEN 2022_DEN 2023_DEN |
(-) Write number to lookup from results:
```

V tejto časti boli použité knižnice Pythonu:

- “os” - na narábanie s OS (natívna)
- “csv” - na narábanie s CSV súborom (natívna)
- “lucene”

Zapracovanie nedostatkov

4/5. Konzultacia - body na zapracovanie:

- treba spraviť fuzzy merge záznamov
- extrakcia z textu wikipedie pre bud' HoF alebo rank (vybrať som si HoF)
- extrakcia nickname-u z wikipedie
- porovnanie s 1. riešením → 2. odovzdanie 02.12.2024 (pretože sa nedá porovnať ani s google a ani s basketball-reference)

Fuzzy merge

Je potrebné pri merge-ovaní dát z wikipédie a dát z basketball-reference počítať skóre podobnosti. To bolo počítané podľa mena hráča a ak bolo nad 80%, tak sa záznam o hráčovi z wikipédie spojil zo záznamom o hráčovi z basketball-reference. Pseudokód vyzerá nasledovne:

[Download](#)

```
Function fuzzy_merge_based_on_names(csv_name: str, wiki_name: str):
    Init "similarity_score" as integer
    Use function "fuzz.ratio" with parameters "csv_name" and "wiki_name" to calculate "similarity_score"
    if "similarity_score" is higher than 80:
        return True
    return False
```

V tejto časti bola použitá knižnica Pythonu:

- “fuzzywuzzy” - na počítanie similarity score medzi ret’azcami (externá)

Extrakcia vety o uvedení hráča v sieni slávy (HOF)

Miesto nedostačujúceho regexu v pôvodnej implementácii bolo potrebné nahradit takým, aby v úvodnej časti odseku o hráčovi bola vyregexovaná celá veta o jeho uvedení v HOF pre obohatenie Lucene indexu textom. Vytiahnutie tejto vety bolo spravené tak, že najprv sa vyregexoval celá úvodná kapitola o hráčovi. Odtiaľto sa spomedzi všetkých viet podľa vhodne určených slov (member|inductee|induct|ed,vote|ed,enshrine|ed,..., hall of fame,...) táto veta vytiahla. Vďaka starému regexu na tento field, ktorý sa t’ahal z infoboxu, bolo možné určiť úspešnosť tohto regexu a je to 83.5% (66 z 79). Pseudokód extrakcie vyzerá nasledovne:

[Download](#)

```
def extract_first_section_and_sentence_about_hof(page_text: str, player_name: str):
    Init regex "''{player_name}'''.*?\n\s*==\s*' to extract whole introduction text
    Init "section_match" from using regex on "page_text"
    If "section_match" was found:
        Init "intro_text" from the found match
        Init regex "(?:[A-Z][^.?!]*?s)(?:induct(?:ee|ed)?|vote(?:d)?|enshrine(?:d)?|member|elect(?:ed)?|name(?:d)?))\s+?(?:\[.\]?)?([:\[Nn]aismith\s+?)?(?:.*?\s+)?[Hh]all\s+o
            to match the sentence
        Init "sentence_match" from using regex on "intro_text"
        if "sentence_match":
            return "hof_sentence" string from "sentence_match"
    return ""
```

Extrakcia prezývky hráča

Pri poslednej konzultácii sa prišlo na to, že ako redirecty sa na wikipédií ukladajú prezývky hráčov, napr. je stránka “Black Mamba” a jediný jej kontent je redirect na Kobe Bryanta (#REDIRECT \[[Kobe Bryant]\]). Takýto redirect vyzerá nasledovne:

The screenshot shows a Wikipedia edit page for the article 'Black Mamba (basketball player)'. The page title is 'Editing Black Mamba (basketball player)'. At the top, there is a warning message: '⚠ You are not logged in. Your IP address will be publicly visible if you make any edits. If you log in or create an account, your edits will be attributed to a user-name, among other benefits.' Below this, there is a note: 'Content that violates any copyrights will be deleted. Encyclopedic content must be verifiable through citations to reliable sources.' The edit summary is '#REDIRECT [[Kobe Bryant]]'. The preview window shows the redirected content: 'This page is a redirect. The following categories are used to track and monitor this redirect:

- From an alternative name: This is a redirect from a title that is another name or identity such as an alter ego, a nickname, or a synonym of the target, or of a name associated with the target.
 - This redirect leads to the title in accordance with the naming conventions for common names to aid searches and writing. It is not necessary to replace these redirected links with a piped link.
 - If this redirect is an incorrect name for the target, then {{R from incorrect name}} should be used instead.

When appropriate, protection levels are automatically sensed, described and categorized.'

Ciž implementácia pre hľadanie prezývek vyzerala tak, že všade sa iba po stránkach kontrolujú redirecty, kde do regexu sa natlačí meno hráča. Zoberie sa zoznam mien hráčov z basketball-reference a na každej stránke sa skúšajú takéto regexy a ak sa nájde, tak sa uloží prezývka. Pseudokód vyzerá nasledovne:

[Download](#)

```
def try_to_extractNickname_from_page(page_text: str, players_names: list[str]):
    for player_name in players_names:
        Init "redirect_regex" as "#REDIRECT \[\{{player_name}\]\]"
        Init "redirect_match" on "page_text" using pattern "redirect_regex"
        If "redirect_match":
            # nickname is in the page
            Init "nickname_regex" as "<title>(.+?)</title>"
            Init "nickname_match" on "page_text" using pattern "nickname_regex"
            If "nickname_match":
                return "nickname" from "nickname_match"
```

NBA - odovzdanie 2

Spark

Na celej wikipédii sa nachádzalo 3196 NBA hráčov z 5204 NBA hráčov, čo tvorí 61.4%. Pre týchto hráčov sa podarilo pridať nové informácie ohľadom trofejí hráčov, ako "Rookie of the year", "Most Valuable Player", "Champion" a "Hall of Fame". Z celkového wiki datasetu iba 39 NBA hráčov má trofej "Rookie of the year" (1.2% z 3196). Ďalej "Most Valuable Player" má iba 23 NBA hráčov (0.7% z 3196). Pohár NBA vyhralo 343 NBA hráčov (10.7% z 3196). Nakoniec v "Hall of Fame" je 79 NBA hráčov (2.4% z 3196). Po zapracovaní nedostatkov z poslednej konzultácií bola pridaná ďalšia entita "nickname" z wikipédie, kde 347 hráčov malo nejakú prezývku (10.8% z 3196).

Dáta pred spracovaním Sparkom vyzerajú nasledovne:

Career_length	Seasons	Teams	Season_team	Nba_roty	Nba_mvp	Nba_champiōn	Hall_of_fame
5 years	2010 2011 2012	GSW HOU CHA	2010_GSW 2011				
7 years	2008 2009 2010	CHA CHA DAL T	2008_CHA 2009				
1 year	2020	PHO	2020_PHO				
11 years	2007 2008 2009	DET DET DEN DE	2007_DET 2008				
7 years	1968 1969 1970	SDR SDR POR PC	1968_SDR 1969				
3 years	2016 2017 2018	OKC OKC OKC	2016_OKC 2017				
13 years	1975 1976 1977	PHO PHO PHO F	1975_PHO 1976				
2 years	1959 1960	DET LAL	1959_DET 1960				
11 years	1985 1986 1987	SAC WSB DEN D	1985_SAC 1986				
20 years	1969 1970 1971	MIL MIL MIL MIL	1969_MIL 1970				
8 years	2010 2011 2012	OKC OKC HOU S	2010_OKC 2011				
7 years	2017 2018 2019	BRK BRK BRK BR	2017_BRK 2018				
1 year	1984	GSW	1984_GSW				
5 years	1990 1991 1992	POR POR MIL BC	1990_POR 1991				
5 years	1986 1987 1988	DEN WSB WSB V	1986_DEN 1987				
1 year	2018	DEN	2018_DEN				
2 years	2005 2008 2008	DET DET LAC	2005_DET 2008				

Dáta po spracovaní Sparkom vyzerajú nasledovne:

Team	Season_team	Nba_roty	Nba_mvp	Nba_champiōn	Hall_of_fame	Nickname	Hof_text
MEM N 2021	MEM 20				Hall of Famer		
MIL MIL 1969	MIL 1970	NBA Rookie	NBA Most Val	NBA Champion	Hall of Famer		
MIL MIL 1996	MIL 1997			NBA Champion	Hall of Famer		He played 18 seasons in the National Basketball Association (NBA) and was
PHI PHI 1984	PHI 1985		NBA Most Val		Hall of Famer	Sir Charles	Barkley is a two-time inductee into the Naismith Memorial Basketball Hall of
STL STL 1962	STL 1963				Hall of Famer		A three-time ABA All-Star and two-time NBA All-Star, Beaty was inducted i
BOS TC 1997	BOS 1998			NBA Champion	Hall of Famer		In 2024, it was announced that Billups would be inducted into the Naismith
BOS BL 1979	BOS 1980	NBA Rookie	NBA Most Val	NBA Champion	Hall of Famer		He was inducted into the Naismith Memorial Basketball Hall of Fame twice
CIN CIN 1958	CIN 1959				Hall of Famer		
VVK N 1947	VVK 1948			NBA Champion	Hall of Famer		
LAL LAL 1996	LAL 1997		NBA Most Val	NBA Champion	Hall of Famer	Black Mamba	Bryant was posthumously voted into the Naismith Memorial Basketball Hal
CHI CHI 1984	CHI 1985	NBA Rookie	NBA Most Val	NBA Champion	Hall of Famer	MI	Jordan was twice inducted into the Naismith Memorial Basketball Hall of F
MEM N 2004	MEM 20				Hall of Famer		

Implementácia čítania veľkého súboru:

[Download](#)

```
def read_wikidump_in_chunks(file_path: str, pages_size=5000):
    # opening compressed wikidump
    with bz2.open(file_path, "rt", encoding="utf-8") as FILE:
        pages = []
        page = ""
        in_page = False
        # page extraction
        for line in FILE:
            # if its time to process the pages chunk in spark
            if len(pages) == pages_size:
                # parse the pages
                wikidump_parse(pages)
                # reset the pages chunk
                pages = []
            # if page is in line, loop is in page, so reset page string
            if START_PAGE in line and not in_page:
                in_page = True
                page = ""
            # if loop is in page, then append the line to page string
            elif END_PAGE not in line and in_page:
                page += line
            # if at the end of the page, then return this specific page and continue then
            elif END_PAGE in line and in_page:
                pages.append(page)
                in_page = False
```

Lucene - indexer

Implementácia lucene indexera:

[Download](#)

```
def index_wiki_data():
    # opening file to read
    with open(constants.WIKI_PLAYER_INFO_PATH, mode='r', encoding='utf-8') as FILE:
        # reading file rows into dicts
        READER = csv.DictReader(FILE, delimiter=';')
        # going through each player record
        for player in READER:
            # doc record for player
            doc = Document()
            # adding fields for indexing
            doc.add(TextField("name", player["name"], Field.Store.YES))
            doc.add(StringField("weight", player["weight"], Field.Store.YES))
            doc.add(StringField("height", player["height"], Field.Store.YES))
            doc.add(TextField("position", player["position"], Field.Store.YES))
            country_of_origin = constants.COUNTRY_CODES[player["born_in_country"]]
            ...
            doc.add(TextField("nba_roty", player["nba_roty"], Field.Store.YES))
            doc.add(TextField("nba_mvp", player["nba_mvp"], Field.Store.YES))
            doc.add(TextField("nba_champion", player["nba_champion"], Field.Store.YES))
            doc.add(TextField("hall_of_fame", player["hall_of_fame"], Field.Store.YES))
            ...
            doc.add(StoredField("season_team", player["season_team"]))
            # write the document of the newly populated player
            WRITER.addDocument(doc)
    WRITER.commit()
    WRITER.close()
```

Lucene - searcher

Default fields:

- name
- position
- shoots
- height
- age
- college
- highschool
- country
- alive

Implementácia lucene searchera:

[Download](#)

```
def full_text_search(query_string: str, searcher, analyzer, fields: list[str], flags: list, result_num=10):
    # init query parser
    query_parser = MultiFieldQueryParser(fields, analyzer)
    query = query_parser.parse(query_string, fields, flags, analyzer)
    # retrieving top n results
    hits = searcher.search(query, result_num).scoreDocs
    # now retrieving the results
    results = []
    for hit in hits:
        # getting the document
        doc = searcher.storedFields().document(hit.doc)
        # getting the result into a dict
        result = {field: doc.get(field) for field in fields if doc.get(field)}
        results.append(result)
    # returning results
    return results
```

Výsledky

Naďavo pylucene napravo vlastné riešenie (výsledky ohraničené na 10 podľa konzultácie):

1. Pre dopyt zobrazený na obrázku dole je potrebné po vyhľadávaní získať iba hráčov zo Srbska, ktorí popri tom majú pozíciu pivota. Obidve riešenia tu spravili všetko správne, lišia sa iba usporiadanie. Každý jeden z nich je aj v jednom aj druhom riešení je zo Srbska a hral/hráva pozíciu pivota. P = 10/10+0 = 100%, R = 10/10+0 = 100%. Zobrazenie vykonaného dopytu:

```
1 -> search serbia center
2 (+) Listing found results
3 | NAME | AGE | VITAL |
4 1: | Miloš Babić | 55yo | alive |
5 2: | Radisav Čurčić | 59yo | alive |
6 3: | Darko Miličić | 39yo | alive |
7 4: | Rastko Cvetković | 54yo | alive |
8 5: | Vlade Divać | 56yo | alive |
9 6: | Nikola Jokić | 29yo | alive |
10 7: | Nenad Krstić | 41yo | alive |
11 8: | Boban Marjanović | 36yo | alive |
12 9: | Filip Petrušev | 24yo | alive |
13 10: | Miroslav Raduljica | 36yo | alive |
```

```
2 -> search serbia center
3 (+) Listing found results
4 | NAME | AGE | VITAL |
5 1: | Radisav Čurčić | 59yo | alive |
6 2: | Miloš Babić | 55yo | alive |
7 3: | Rastko Cvetković | 54yo | alive |
8 4: | Dragomir Tarlać | 51yo | alive |
9 5: | Vladimir Vraneš | 41yo | alive |
10 6: | Miroslav Raduljica | 36yo | alive |
11 7: | Filip Petrušev | 24yo | alive |
12 8: | Željko Rebrača | 52yo | alive |
13 9: | Darko Miličić | 39yo | alive |
14 10: | Nikola Jokić | 29yo | alive |
```

2. Ďalší dopyt je ohľadom nájdenia ľavákov, ktorí hrali pozíciu rozohrávača v sezóne 1994. Oba dopyty sú znova správne, všetci sú ľaváci, rozohrávači a zo sezóny 1994. Tentokrát sa lišia v zoradení iba 3 hráčov. Pre Lucene sa ale musí query lišiť od query vlastného riešenia. P = 10/10+0 = 100%, R = 10/10+0 = 100%. Zobrazenie druhého dopytu:

```

1 -> search position:"point guard" AND shoots:lefty AND seasons:1994
2 (+) Listing found results
3 | NAME | AGE | VITAL | SHOOTS | POSITION |
4 1: | Anthony Guy Bennett | 55yo | alive | Lefty | Point Guard |
5 2: | Brooks James Thompson | 45yo | deceased | Lefty | Point Guard |
6 3: | Johnny Earl Dawkins Jr. | 61yo | alive | Lefty | Point Guard |
7 4: | Gregory Carleton Anthony | 57yo | alive | Lefty | Point Guard |
8 5: | John Kevin Crotty | 55yo | alive | Lefty | Point Guard |
9 6: | Elliot Lamonte Perry | 55yo | alive | Lefty | Point Guard |
10 7: | Darrick David Martin | 53yo | alive | Lefty | Point Guard |
11 8: | Nickey Maxwell Van Exel | 53yo | alive | Lefty | Point Guard |
12 9: | Kenneth Anderson | 54yo | alive | Lefty | Point Guard |
13 10: | Avery DeWitt Johnson | 59yo | alive | Lefty | Point Guard |

```

3. Ten istý dopyt ako druhý, ale tentokrát je nechaná rovnaká query v Lucene vyhľadávači ako vo vlastnom riešení. Usporadanie výsledkov je o dosť iné, ale hlavne sú získané odlišné záznamy. Lucene kvôli jeho default-nému OR berie výsledky, ktoré nie sú relevantné. Tým pádom vznikajú horšie metriky: P = 5/5+5 = 50%, R = 5/5+5 = 50%. Zobrazenie tretieho dopytu:

```

1 -> search point guard lefty 1994
2 (+) Listing found results
3 | NAME | AGE | VITAL | SHOOTS | POSITION |
4 1: | Melvin Jermaine Booker | 52yo | alive | Righty | Point Guard |
5 2: | Litterial Green | 56yo | alive | Righty | Point Guard |
6 3: | Devin Armani Booker | 28yo | alive | Righty | Shooting Guard |
7 4: | Anthony Guy Bennett | 55yo | alive | Lefty | Point Guard |
8 5: | Brooks James Thompson | 45yo | deceased | Lefty | Point Guard |
9 6: | Samuel David Hauser | 26yo | alive | Righty | Small Forward |
10 7: | Ryan Lorthridge | 52yo | alive | Lefty | Shooting Guard |
11 8: | Johnny Earl Dawkins Jr. | 61yo | alive | Lefty | Point Guard |
12 9: | Gregory Carleton Anthony | 57yo | alive | Lefty | Point Guard |
13 10: | John Kevin Crotty | 55yo | alive | Lefty | Point Guard |

```

4. Štvrtý dopyt sa zameriava na hráčov, ktorí žijú, hrali pozíciu kridla a na VŠ chodili do UCLA. Výsledky pre obidve riešenia sú všetky relevantné, aj keď záznamy sa líšia, či už usporiadaní alebo celkovo. Sú iba 3 záznamy rovnaké (Jaquez, Henderson, Fields). P = 10/10+0 = 100%, R = 10/10+0 = 100%. Zobrazenie štvrtého dopytu:

```

1 -> search players from ucla playing small forward who are still alive
2 (+) Listing found results
3 | NAME | AGE | VITAL | POSITION | COLLEGE |
4 1: | Kyle F. Anderson | 31yo | alive | Small Forward | Ucla |
5 2: | Trevor Anthony Ariza | 59yo | alive | Small Forward | Ucla |
6 3: | Matt Kelly Barnes | 44yo | alive | Small Forward | Ucla |
7 4: | Darren Keefe Daye | 64yo | alive | Small Forward | Ucla |
8 5: | Keith Raymond Erickson | 80yo | alive | Small Forward | Ucla |
9 6: | Kenneth Henry Fields | 62yo | alive | Small Forward | Ucla |
10 7: | Milton Henderson Jr. | 48yo | alive | Small Forward | Ucla |
11 8: | Jaime Jaquez Jr. | 23yo | alive | Small Forward | Ucla |
12 9: | Marques Kevin Johnson | 68yo | alive | Small Forward | Ucla |
13 10: | Jason Alan Kapono | 43yo | alive | Small Forward | Ucla |

```

5. Piaty dopyt mal vyhľadať ľavákov hráčov z Kansasu, ktorí sú strelníci a hrali v Los Angeles Lakers. Takýto boli v NBA iba dvaja. Vlastné riešenie tým pádom zobrazí iba 2 záznamy. Naopak Lucene prinesie aj ďalších 8 záznamov okrem týchto dvoch, ktoré sú podobné. P = 2/2+8 = 20%, R = 2/2+0 = 100%. Zobrazenie piatého dopytu:

```

1 -> search lefty players from Kansas who played or are playing shooting guard in los angeles lakers
2 (+) Listing found results
3 | NAME | SHOOTS | POSITION | HIGHSCCHOOL |
4 1: | Gail Charles Goodrich Jr. | Lefty | Shooting Guard | Polytechnic in Los Angeles |
5 2: | Robert Antell Sims Jr. | Righty | Shooting Guard | Jordan in Los Angeles |
6 3: | Kareem Lamar Rush | Lefty | Shooting Guard | Pembroke Hill in Kansas City |
7 4: | Nicholas Aaron Young | Righty | Shooting Guard | Cleveland in Los Angeles |
8 5: | Anthony Eugene Peeler | Lefty | Shooting Guard | Paseo in Kansas City |
9 6: | Jerome D. Henderson | Righty | Power Forward | Jefferson in Los Angeles |
10 7: | Darius Aaron Morris | Righty | Point Guard | Windward in Los Angeles |
11 8: | Larry Donnell Drew Sr. | Righty | Point Guard | Wyandotte in Kansas City |
12 9: | Owayne L. Poole Sr. | Righty | Shooting Guard | Manual Arts in Los Angeles |
13 10: | Lucius Oliver Allen Jr. | Righty | Point Guard | Wyandotte in Kansas City |

```

```

1 -> search point guard lefty 1994
2 (+) Listing found results
3 | NAME | AGE | VITAL | SHOOTS | POSITION |
4 1: | Anthony Guy Bennett | 55yo | alive | Lefty | Point Guard |
5 2: | Brooks James Thompson | 45yo | deceased | Lefty | Point Guard |
6 3: | Johnny Earl Dawkins Jr. | 61yo | alive | Lefty | Point Guard |
7 4: | John Kevin Crotty | 55yo | alive | Lefty | Point Guard |
8 5: | Elliot Lamonte Perry | 55yo | alive | Lefty | Point Guard |
9 6: | Gregory Carleton Anthony | 57yo | alive | Lefty | Point Guard |
10 7: | Darrick David Martin | 53yo | alive | Lefty | Point Guard |
11 8: | Nickey Maxwell Van Exel | 53yo | alive | Lefty | Point Guard |
12 9: | Kenneth Anderson | 54yo | alive | Lefty | Point Guard |
13 10: | Avery DeWitt Johnson | 59yo | alive | Lefty | Point Guard |

```

```

1 -> search point guard lefty 1994
2 (+) Listing found results
3 | NAME | AGE | VITAL | SHOOTS | POSITION |
4 1: | Anthony Guy Bennett | 55yo | alive | Lefty | Point Guard |
5 2: | Brooks James Thompson | 45yo | deceased | Lefty | Point Guard |
6 3: | Johnny Earl Dawkins Jr. | 61yo | alive | Lefty | Point Guard |
7 4: | John Kevin Crotty | 55yo | alive | Lefty | Point Guard |
8 5: | Elliot Lamonte Perry | 55yo | alive | Lefty | Point Guard |
9 6: | Gregory Carleton Anthony | 57yo | alive | Lefty | Point Guard |
10 7: | Darrick David Martin | 53yo | alive | Lefty | Point Guard |
11 8: | Nickey Maxwell Van Exel | 53yo | alive | Lefty | Point Guard |
12 9: | Kenneth Anderson | 54yo | alive | Lefty | Point Guard |
13 10: | Avery DeWitt Johnson | 59yo | alive | Lefty | Point Guard |

```

```

1 -> search players from ucla playing small forward who are still alive
2 (+) Listing found results
3 | NAME | AGE | VITAL | POSITION | COLLEGE |
4 1: | Ray Lynn Shackelford | 77yo | alive | Small Forward | Ucla |
5 2: | Milton Henderson Jr. | 47yo | alive | Small Forward | Ucla |
6 3: | Jaime Jaquez Jr. | 23yo | alive | Small Forward | Ucla |
7 4: | Keith Kensi Owens | 55yo | alive | Small Forward | Ucla |
8 5: | Michael Edward Lynn | 78yo | alive | Small Forward | Ucla |
9 6: | Dijon Lynn Thompson | 41yo | alive | Small Forward | Ucla |
10 7: | James Robert Wilkes | 66yo | alive | Small Forward | Ucla |
11 8: | Edward Charles O'Bannon Jr. | 52yo | alive | Small Forward | Ucla |
12 9: | Kenneth Henry Fields | 62yo | alive | Small Forward | Ucla |
13 10: | Shabazz Naige Muhammad | 31yo | alive | Small Forward | Ucla |

```

```

1 -> search lefty players from Kansas who played or are playing shooting guard in los angeles lakers
2 (+) Listing found results
3 | NAME | SHOOTS | POSITION | HIGHSCCHOOL |
4 1: | Karem Lamar Rush | Lefty | Shooting Guard | Pembroke Hill in Kansas City |
5 2: | Anthony Eugene Peeler | Lefty | Shooting Guard | Paseo in Kansas City |

```

Priemerný precision: 74%. Priemerný recall: 90%.

Výhodou vlastného riešenia je to, že je presné a vždy vráti iba to, čo si používateľ vypýtal. Na druhú stranu ale vracia niekedy menej než dohodnutých 10 záznamov a nesnaží sa vôbec pri tom vraciať aj podobné záznamy. Presné vyhľadávanie je tiež niekedy nevhod, keď používateľ presne nevie, že čo hľadá a systém mu nevie ani trochu napovedať. Zatiaľ čo Lucene vracia oproti presnému výsledkom aj omáčky, ale na druhú stranu pri lepšie špecifikovanej query vie vraciať rovnako presné výsledky ako vracať vlastné riešenie. Nevýhodou tohto je, že používateľ musí byť aspoň trochu zaškolený ako to vyhľadávanie funguje na to, aby vedel toto presné vyhľadávanie využívať.

Zdrojový kód

[xbartosp2_odovzdanie_2.zip](#)