

ČVUT, FAKULTA INFORMAČNÍCH TECHNOLOGIÍ
VYHLEDÁVÁNÍ NA WEBU A V MULTIMEDIÁLNÍCH DATABÁZÍCH
LETNÍ SEMESTR 2019/2020
ZÁVĚREČNÁ ZPRÁVA K PROJEKTU

LSI vektorový model

David Mašek a Kristýna Klesnilová

12. května 2020

OBSAH

1	Popis projektu	3
2	Způsob řešení	3
2.1	Preprocessing dokumentů	3
2.2	Výpočet vah termů	3
2.3	Implementace LSI	4
2.4	Vyhodnocení dotazu	5
3	Implementace	5
4	Příklad výstupu	6
5	Experimentální sekce	6
5.1	Určení optimálního počtu konceptů	6
5.2	Porovnání vlivu LSI na kvalitu výsledků vyhledávání s ohledem na výskyt synonym a homonym	7
5.2.1	Synonyma	7
5.2.2	Homonyma	8
5.3	Porovnání průchodu pomocí LSI vektorového modelu se sekvenčním průchodem databáze s ohledem na čas vykonání dotazu	9
5.4	Vliv různých vnitřních parametrů na výkon algoritmu (změna počtu konceptů, změna počtu extrahovaných termů, použití lemmatizace namísto stemmingu, odstranění číslovek při preprocessingu, použití jiného vzorce na výpočet vah termů...)	9
6	Diskuze	9
7	Závěr	9

1 POPIS PROJEKTU

V tomto projektu implementujeme *LSI vektorový model* sloužící k podobnostnímu vyhledávání v databázi anglických textových dokumentů. Tuto funkcionalitu následně vizualizujeme pomocí webového interface, který uživateli umožňuje procházet databázi článků na základě doporučení nejpodobnějších článků k právě čtenému.

V experimentální části projektu jsme se dále zaměřili na:

- Určení optimálního počtu konceptů
- Porovnání vlivu LSI na kvalitu výsledků vyhledávání s ohledem na výskyt synonym a homonym
- Porovnání průchodu pomocí LSI vektorového modelu se sekvenčním průchodem databáze s ohledem na čas vykonání dotazu
- Vliv různých vnitřních parametrů na výkon algoritmu (změna počtu konceptů, změna počtu extrahovaných termů, použití lemmatizace namísto stemmingu, odstranění číslovky při preprocesingu, použití jiného vzorce na výpočet vah termů...)

Celý náš projekt je volně dostupný k vyzkoušení na: <https://bi-vwm-lsi-demo.herokuapp.com/>.

2 ZPŮSOB ŘEŠENÍ

2.1 Preprocessing dokumentů

Jako první v naší aplikaci začínáme s preprocessingem dokumentů. Slova z jednotlivých dokumentů převedeme na malá písmena a odstraníme z nich nevýznamová slova a interpunkci. K identifikaci nevýznamových slov používáme seznam anglických nevýznamových slov. Jako parametr programu posíláme také, zda má z dokumentů odstranit i číslovky. Následně na zbylé termy aplikujeme *stemming* či *lemmatizaci*. Tím se snažíme slova, která mají stejný slovní základ, vyjádřit pouze jedním termem. Stemming to dělá pomocí algoritmu, kterým odsekává přípony a koncovky slova. Lemmatizace na to jde o něco chytřeji, podle kontextu slova se pokusí určit, o jaký slovní druh se jedná, a podle toho ho zkrátit.¹ Porovnání jejich použití v programu se dále věnujeme v experimentální části.

2.2 Výpočet vah termů

V aplikaci vytváříme matici M_w , která má v řádcích jednotlivé termy a ve sloupcích jejich váhy v jednotlivých dokumentech.

Začneme tím, že si vytvoříme matici počtu výskytů jednotlivých termů v jednotlivých dokumentech. Počet termů v této matici poté dále zredukujeme, abychom pracovali jen s těmi nejdůležitějšími. Funkci pro redukci termů posíláme následující parametry:

- *max_df* - termy nacházející se ve více % dokumentů, než udává číslo $100 * max_df$, z matice odstraníme
- *min_df* - termy nacházející se v méně nebo stejně dokumentech, než udává číslo *min_df*, z matice odstraníme
- *max_terms* - maximální počet termů, které si v aplikaci necháme

¹nlp.stanford.edu/IR-book/html/htmledition/stemming-and-lemmatization-1.html

- *keep_less_freq* - udává, zda si při výběru *max_terms* termů nechat ty nejméně či nejvíce často zastoupené v dokumentech

Zkoumání vlivu změny jednotlivých parametrů na výsledek LSI se dále podrobněji věnujeme v experimentální části. V programu vždy nastavujeme *min_df* alespoň na 1, abychom odstranili termy nacházející se pouze v 1 dokumentu, které nám do LSI nepřidávají žádné užitečné informace. (pravda?)

Z této zredukované matice poté již spočteme matici M_w . Pro výpočet vah jednotlivých termů používáme metodiku *tf-idf*. Váhu termu t_i v dokumentu d_j spočítáme podle vzorce:

$$w_{ij} = tf_{ij} \cdot idf_{ij} \quad (2.1)$$

kde tf_{ij} reprezentuje normalizovanou četnost termu t_i v dokumentu d_j a spočítáme ji podle vzorce²:

$$tf_{ij} = \frac{f_{ij}}{nt_j} \quad (2.2)$$

kde f_{ij} je četnost výskytu termu t_i v dokumentu d_j , kterou normalizujeme číslem nt_j vyjadřujícím celkový počet termů v dokumentu d_j . V přednášových slidech je použita normalizace jiná, tf_{ij} se tam počítá podle vzorce:

$$tf_{ij} = \frac{f_{ij}}{\max_i\{f_{ij}\}} \quad (2.3)$$

kde $\max_i\{f_{ij}\}$ vrací nejvyšší četnost termu t_i přes celou kolekci dokumentů. Tento způsob normalizace nám vrací spíše horší výsledky, jejich porovnáním se zabýváme v experimentální části. idf_{ij} reprezentuje převrácenou četnost t_i ve všech dokumentech a spočítá se podle vzorce:

$$idf_{ij} = \log_2\left(\frac{n}{df_i}\right) \quad (2.4)$$

kde n je celkový počet dokumentů a df_i reprezentuje celkový počet dokumentů obsahujících term t_i .

2.3 Implementace LSI

Jakmile máme vytvořenou matici vah termů M_w , můžeme přistoupit k samotné implementaci LSI. Princip LSI spočívá v tom, že s pomocí *singulárního rozkladu (SVD)* seskupíme tematicky podobné články do jednotlivých k konceptů. Vlivem počtu konceptů na kvalitu výsledků se dále zabýváme v experimentální sekci.

Singulární rozklad nám matici M_w rozloží následovně:

$$M_w = U \cdot S \cdot V^T \quad (2.5)$$

²<https://en.wikipedia.org/wiki/Tf%E2%80%93idf>

kde řádky matice U jsou obrazy řádků matice M_w , sloupce matice V jsou obrazy sloupců matice M_w a matice S obsahuje na diagonále *singulární hodnoty* (*absolutní hodnoty vlastních čísel*) matice M_w v sestupném pořadí. Z těchto matic získáme *concept-by-document* matici M_{cd} jako:

$$M_{cd} = S[k, k] \cdot V^T[k, :] \quad (2.6)$$

kde $S[k, k]$ značí prvních k řádků a sloupců matice S a $V^T[k, :]$ značí prvních k řádků matice V^T , kde k je počet konceptů. Nenásobíme tedy celou maticí M_{cd} , ale pouze její část podle počtu konceptů.

Matici projekce dotazu do prostoru konceptů M_q pak získáme jako:

$$M_q = U^T[k, :] \quad (2.7)$$

kde $U^T[k, :]$ značí prvních k řádků matice U^T .

2.4 Vyhodnocení dotazu

Při dotazu na nejpodobnější dokumenty k dokumentu d_j převedeme dotaz do prostoru konceptů na vektor V_c pomocí vzorce:

$$V_c = M_q \cdot M_{w:,j} \quad (2.8)$$

kde $M_{w:,j}$ značí j -tý sloupec matice M_w .

Vektor V_c poté pomocí *kosinové podobnosti* porovnáme se sloupcovými vektory matice M_{cd} . Indexy nejpodobnějších sloupcových vektorů matice M_{cd} pak vrátíme jako indexy nejpodobnějších dokumentů k dokumentu dotazu d_j . Spolu s indexy vrátíme i samotnou hodnotu kosinové podobnosti.

3 IMPLEMENTACE

Celý projekt jsme programovali v jazyce *Python*. Práci nám velmi usnadnila jeho knihovna *NLTK*³ nabízející rozsáhlou funkcionalitu pro práci s přirozeným jazykem. Využili jsme například *WordNetLemmatizer* pro lemmatizaci či *SnowballStemmer* pro stemming. Dále jsme v programu hojně využívali Python knihovny *pandas*⁴ a *numpy*⁵.

Ukládání dat v projektu řešíme přes CSV soubory, ke kterým přistupujeme přes *pandas* funkce. V jednom souboru máme uložený dataset nad kterým provádíme LSI. V dalších souborech pak máme uložené matice M_w , M_{cd} a M_q , abychom je mohli cachovat a přepočítavat jen při změně LSI parametrů, které máme uložené v souboru *server/lisa_config.json*.

Procházení článků vizualizujeme v prohlížeči pomocí *Flask*⁶ web serveru. Jako dataset v našem projektu používáme anglicky psané novinové články stažené z *kaggle.com*⁷. Dataset nepoužíváme celý, vybrali jsme z něj pouze 996 článků.

³<https://www.nltk.org/>

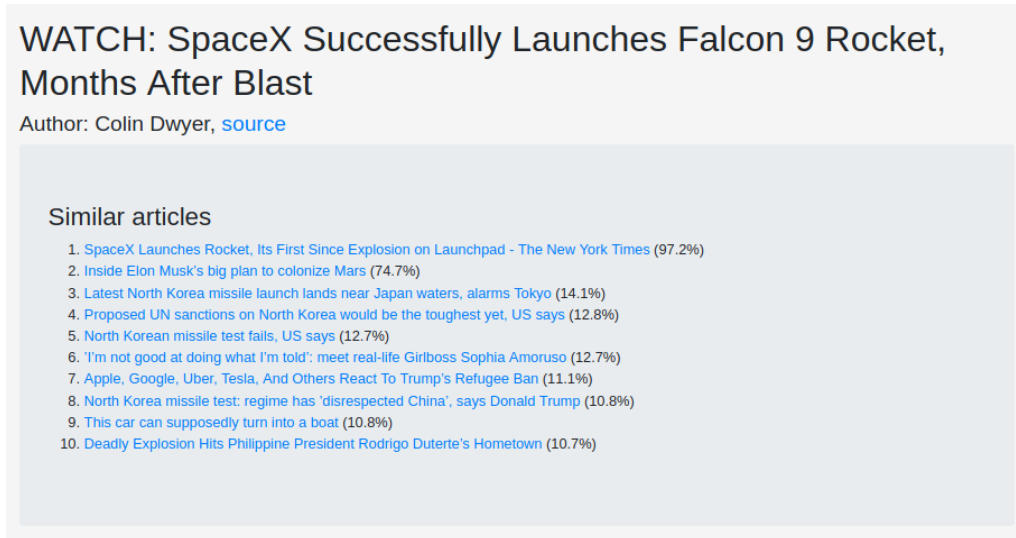
⁴<https://pandas.pydata.org/>

⁵<https://numpy.org/>

⁶<https://flask.palletsprojects.com/en/1.1.x/>

⁷<https://www.kaggle.com/snapcrack/all-the-news>

4 PŘÍKLAD VÝSTUPU

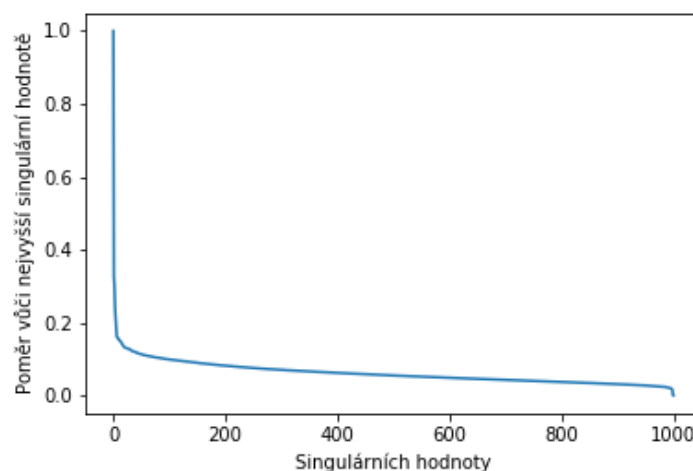


Obrázek 4.1: Příklad výstupu aplikace

Na obrázku 4.1 je vidět konkrétní vstup a výstup naší aplikace. Zobrazí se název článku dotazu, jméno jeho autora a také samotný text článku. Naše aplikace dále uživateli nabídne seznam 10 nejpodobnějších článků i s určenou kosinovou podobností v procentech. Je vidět, že aplikace vrací víceméně přesně to, co bychom čekali. K článku o tom, že firma SpaceX vypustila do vesmíru raketu, vrátí články týkající se firmy SpaceX či raket.

5 EXPERIMENTÁLNÍ SEKCE

5.1 Určení optimálního počtu konceptů



Obrázek 5.1: Důležitost konceptů

Vezmeme si singulární hodnoty, které nám vrátil singulární rozklad v rovnici 2.5. Vizualizujeme-li si v grafu 5.1, jak klesá poměr nejvyšší singulární hodnoty vůči zbylým singulárním hodno-

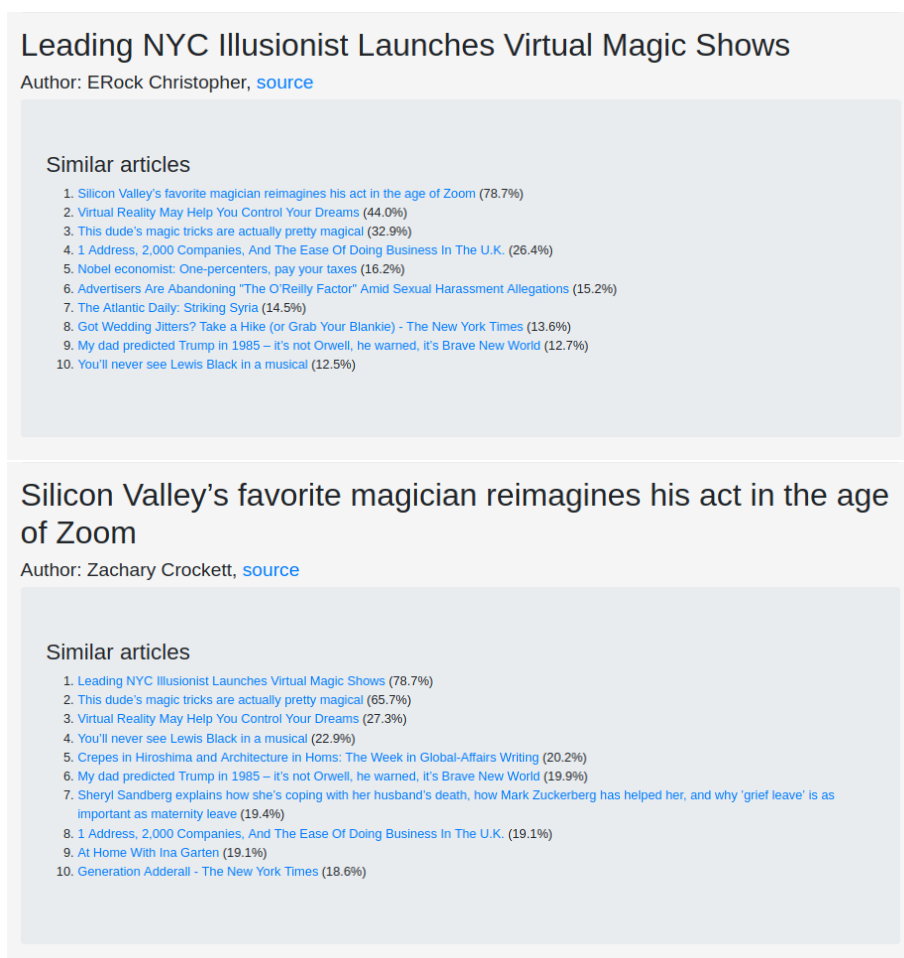
tám, vidíme z toho také, jak klesá důležitost konceptů v datasetu. Počet konceptů k v naší aplikaci tedy podle tohoto grafu určíme jako 400.

5.2 Porovnání vlivu LSI na kvalitu výsledků vyhledávání s ohledem na výskyt synonym a homonym

Pro zjištění kvality výsledků vyhledávání s ohledem na výskyt synonym a homonym jsme do datasetu přidali 4 články.

5.2.1 Synonyma

Jak si náš model poradí se synonymy jsme testovali na článcích *"Silicon Valley's favorite magician reimagines his act in the age of Zoom"* a *"Leading NYC Illusionist Launches Virtual Magic Shows"*. Testovali jsme tedy anglická synonyma *illusionist* a *magician*.

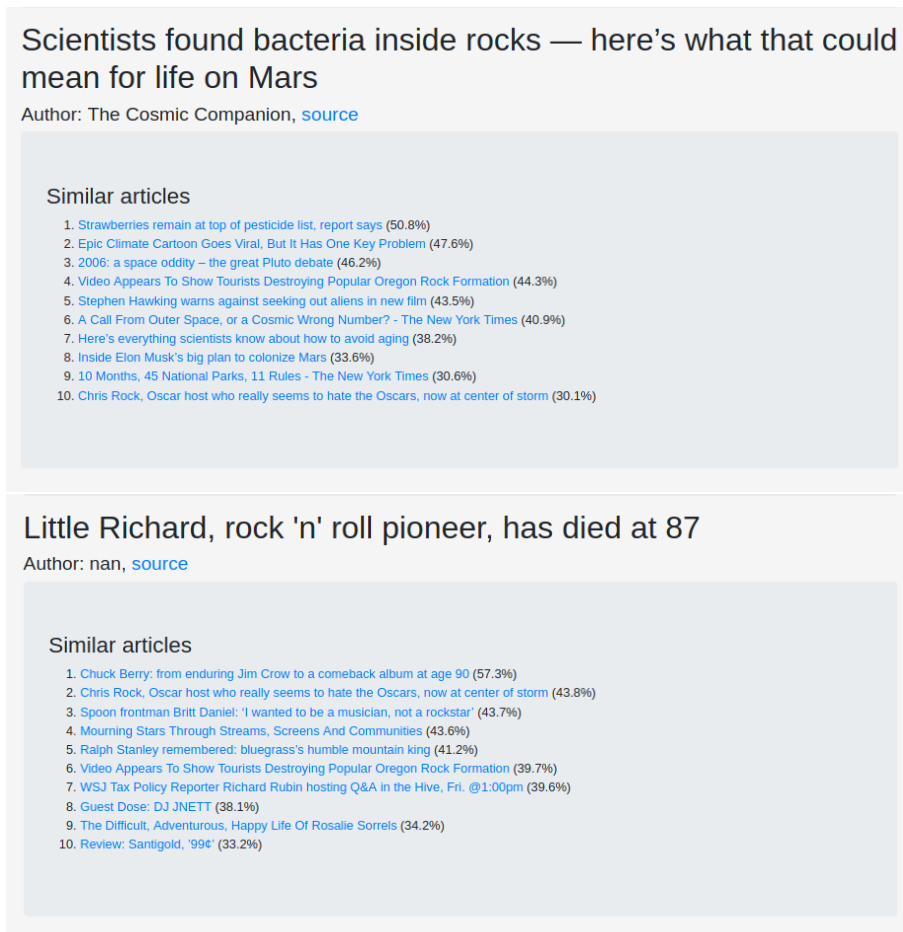


Obrázek 5.2: Synonyma

Na obrázku 5.2 je vidět, že náš model synonyma v článcích identifikoval správně. V prvním článku se sice ani jednou přímo nevyskytuje slovo *magician* a v druhém ani jednou slovo *illusionist*, ale v obou se vyskytují slova se slovním základem *magic* nebo například slova *show* a *audience* a model tak články správně identifikoval jako podobné.

5.2.2 Homonyma

K testování jak si náš článek poradí s homonymy jsme použili články *"Little Richard, rock 'n' roll pioneer, has died at 87"* a *"Scientists found bacteria inside rocks — here's what that could mean for life on Mars"*, které oba obsahují slovo *rock* ovšem jednou ve významu skály a jednou ve významu rockového muzikanta.



Obrázek 5.3: Homonyma

Na obrázku 5.3 je vidět, že u prvního článku, týkajícího se vlivu objevu bakterie v podmořských skalách na možnost života na Marsu, si model poradil opravdu dobře a vrátil nám články týkající se vesmíru nebo přírody. Jediné co model spletlo je článek o Chrise Rockovi, jehož příjmení, které se ve článku vyskytuje velmi často, jsme při předzpracování převedli na malá písmena do tvaru *rock* a model tak neví, že se jedná o vlastní jméno. I přes tento problém jsme se v naší aplikaci rozhodli všechna slova převádět na malá písmena. Kdybychom se rozhodli převádět na malá písmena pouze slova na začátcích vět a ostatní vlastní jména nechávali s velkým písmenem na začátku, identifikoval by nám model například slovo *school* a slovo *School* v názvu nějaké školy jako 2 rozdílná slova. Je tedy na zvážení autorů aplikace, jak k problému přistoupí.

U druhého článku týkajícího se smrti rockového hudebníka nám jako první vrátil velmi relevantní článek o jiném stárnoucím hudebníkovi. Jako druhý článek nám model opět vrátil článek o Chrise Rockovi, jež je teď pro model velmi relevantní, jelikož je to článek o celebritě a obsahuje spoustu výskytů slova *rock*. Třetí, čtvrtý a pátý článek jsou opět relevantní články o hudebnících, ale jako šestý článek už model vrátí článek o skále namísto o hudebnících. Je tedy

vidět, že velmi relevantní články o hudebnících model identifikuje správně jako nejpodobnější, u méně relevantních článků týkajících se hudebníků, hudby nebo celebrit už má model problém identifikovat, jestli jsou relevantnější tyto články nebo články o skalách. Velmi záleží na tom, kolik podobných slov se ve člancích vyskytuje.

5.3 Porovnání průchodu pomocí LSI vektorového modelu se sekvenčním průchodem databáze s ohledem na čas vykonání dotazu

Pro odsimulování sekvenčního průchodu databází nastavíme počet konceptů stejný jako počet dokumentů databáze. Toho docílíme tak, že v souboru *server/lsa_config.json* nastavíme parametr *k* na 0.

```
In [111]: df_concept_by_doc, df_query_projection = transform_to_concept_space(df_tf_idf, k=0, customSVD=False)

start = time.time()
best_match = get_n_nearest(df_tf_idf, df_concept_by_doc, df_query_projection, 999, n=10)
end = time.time()
time_seq = end - start
print("Čas vykonání dotazu při sekvenčním průchodu databází: " + str(time_seq) + " sekund")

Čas vykonání dotazu při sekvenčním průchodu databází: 0.18472576141357422 sekund

In [112]: df_concept_by_doc, df_query_projection = transform_to_concept_space(df_tf_idf, k=400, customSVD=False)

start = time.time()
best_match = get_n_nearest(df_tf_idf, df_concept_by_doc, df_query_projection, 999, n=10)
end = time.time()
time_norm = end - start
print("Čas vykonání dotazu při průchodu databází pomocí LSI vektorového modelu: " + str(time_norm) + " sekund")

Čas vykonání dotazu při průchodu databází pomocí LSI vektorového modelu: 0.13696837425231934 sekund

In [113]: print("Rozdíl: " + str(time_seq - time_norm))

Rozdíl: 0.04775738716125488
```

Obrázek 5.4: Porovnání času vykonání dotazu při sekvenčním průchodu databáze

Na obrázku 5.4 je vidět, že při sekvenčním průchodu databází se čas vykonání dotazu zpomalí asi o 0.05 sekund. Při spuštění Flask serveru a testování v prohlížeči je však rozdíl v rychlosti načítání stránky zobrazující obsah článku a 10 jemu nejpodobnějších článků mnohem výraznější, kolem 3 sekund.

5.4 Vliv různých vnitřních parametrů na výkon algoritmu (změna počtu konceptů, změna počtu extrahovaných termů, použití lemmatizace namísto stemmingu, odstranění číslovek při preprocesingu, použití jiného vzorce na výpočet vah termů...)

6 DISKUZE

7 ZÁVĚR