

1.Cvičení

1.Príklad

Nájdite dotaz v tvare [vyraz_1 vyraz_2] (bez úvodzoviek), ktorý Google neinterpretuje ako súčin jednotlivých výrazov dotazu. Tzn. dotaz, keď bude položený Googlu, vráti aspoň jeden výsledok, ktorý obsahuje len jeden výraz z dotazu.

a) Prvých 20 výsledkov označte číslami 0, 1, 2 podľa toho, či výsledok obsahuje oba výrazy, len jeden z výrazov alebo ani jeden výraz.

b) Podľa tohto usúďte, či Google interpretuje všetky dotazy implicitne ako Booleovský prienik výrazov z dotazu.

Řešení: Řešení si každý musí vyzkoušet sám.

2. Príklad

Vytvorte invertovaný index, zostavený pre nasledujúcu kolekciu dokumentov:

Doc 1 new home sales top forecasts

Doc 2 home sales rise in july

Doc 3 increase in home sales in july

Doc 4 july new home sales rise

Řešení: napíšeme tabulku a za každé slovo napíšeme ve kterých dokumentech se vyskytuje.

new	1,4
home	1,2,3,4
sales	1,2,3,4
top	1
forecasts	1
rise	2,4
in	2,3
july	2,3,4
increase	3

3. Príklad

Nižšie je časť indexu s pozíciami v tvare term: doc1: <pos1, pos2, pos3, ..>; doc2<pos1, pos2, ..>

angels: 2: <36,174,252,651>; 4: <12,22,102,432>; 7: <17>;
fools: 2: <1,17,74,222>; 4: <8,78,108,458>; 7: <3,13,23,193>;
fear: 2: <87,704,722,901>; 4: <13,43,113,433>; 7: <18,328,528>;
in: 2: <3,37,76,444,851>; 4: <10,20,110,470,500>; 7: <5,15,25,195>;
rush: 2: <2,66,194,321,702>; 4: <9,69,149,429,569>; 7: <4,14,404>;
to: 2: <47,86,234,999>; 4: <14,24,774,944>; 7: <199,319,599,709>;
tread: 2: <57,94,333>; 4: <15,35,155>; 7: <20,320>;
where: 2: <67,124,393,1001>; 4: <11,41,101,421,431>; 7: <15,35,735>;

Ktoré dokumenty zodpovedajú nasledujúcim dotazom, kde každý výraz v uvozovkách je frázový dotaz (phrase query)?

- “fools rush in”
 - “fools rush in” AND “angels fear to tread”.
- Na ktorých pozíciách tieto dotazy matchujú?
- V uvedenom indexe je chyba, kde?

Řešení: Hledáme ve kterých dokumentech jdou slova za sebou. V závorkách jsou indexy pozice slov

a) 2 (1-3), 4 (8-10), 7 (3-5, 13-15)

b) 4 (8-10 AND 12-15)

c) V dokumentu 7 je na indexu 15 jak „where“ tak „in“

4. Príklad

Odporučte stratégiu spracovania dotazu (tangerine OR trees) AND (marmalade OR skies) AND (kaleidoscope OR eyes) vzhľadom na nasledujúce veľkosti postings zoznamov:

eyes 213312
kaleidoscope 87009
marmalade 107913
skies 271658
tangerine 46653
trees 316812

Řešení: Použil bych seznam s přeskokovači. Ty se nám vyplatí ale jen na AND, proto bych zpracovával AND jak nejdříve to bude možné. Přeskakuje se po \sqrt{n} .

5. Príklad

Máme dotaz zložený z dvoch výrazov. Postings zoznam jedného výrazu je zložený z nasledujúcich 16 položiek:

[4,6,10,12,14,16,18,20,22,32,47,81,120,122,157,180]

a druhý výraz ma postings zoznam len jednoprvkový:

[47].

Zistite koľko porovnaní a prečo je potreba vykonať na prienik týchto dvoch postings zoznamov s nasledujúcimi stratégiami:

a. použitie štandardných postings zoznamov

b. použitie postings zoznamov uložených s preskakujúcimi odkazmi, s dĺžkou preskoku \sqrt{P}

Řešení: Nejde vlastně o nic jiného než o procházení seznamu.

a) je potřeba projít 11 prvků seznamu, než se dojde ke shodě.

b) je potřeba 6 projití. Skok 4-12, skok 12-18, skok 18-32, skok 32-120, skok zpět 120-32, krok na 47

c)

6. Príklad

Pomocou matice vypočítajte Levenshteinovu vzdialenosť slov jablko a malina.

Řešení: Jednoduše řečeno, kolik písmen musíme změnit/přidat/odebrat, aby se ze slova A stalo slovo B. Pro jednoduchost můžeme počítat, že pro všechny operace je stejná váha.

jablko	jablko	jablko	jablko	jablko	jablko	jablko
malina	jalina	malina	malina	malina	malina	malina
0	1	1	2	3	4	5

7. Príklad

a) Nájdi dve rozdielne napísané podstatné mená (anglicky), ktoré majú rovnaký soundex kód.

b) Nájdi dve foneticky podobné podstatné mená (anglicky), ktoré majú rozdielny soundex kód

Řešení: Soundex index je čtyř místný kód. Vytváří se pomocí následujícího postupu.

Slovo „BOOKS“

1. Necháme první písmeno – „B“
2. Nahradíme skupiny písmen – „B0022“
3. Odstraníme duplicity – „B02“
4. Odstraníme nuly – „B2“
5. Změníme velikost na 4 (ořežeme nebo doplníme nuly) – „B200“

Skupiny písmen:

A, E, I, O, U, H, W, Y -> 0

B, F, P, V -> 1

C, G, J, K, Q, S, X, Z -> 2

D, T -> 3

L -> 4

M, N -> 5.

R -> 6

- a) Trash – T620
Thrash – T620
- b) Kill – K400
Bill – B400

8. Príklad

Vypíšte prvky slovníka permuterm indexu ktoré su generované slovom mama.

Řešení: Napíšeme za slovo znak „\$“ a napíšeme všechny permutace slova

mama\$

ama\$m

ma\$ma

a\$mam

\$mama

9. Príklad

Aké kľúče sú použiteľné na nájdenie termu s*ng v permuterm wildcard indexe.

Řešení: Doplňíme dolar a budeme permutovat tak dlouho, dokud nepřesuneme hvězdičku na konec.

s*ng\$
*ng\$\$
ng\$\$*

Hledá se tedy klíč začínající na „ng\$\$“

10. Príklad

Pre $n = 2$ a $1 \leq T \leq 30$, vykonajte krok za krokom simuláciu algoritmu 4.7 (Introduction to Information Retrieval). Vytvorte tabuľku, ktorá pre každý okamih v čase, v ktorom je spracovaných $T = 2 \cdot k$ tokenov ($1 \leq k \leq 15$), ukazuje ktoré zo štyroch indexov I_0, \dots, I_3 sú používané. Prvé tri riadky tabuľky sú uvedené nižšie:

	I_3	I_2	I_1	I_0
2	0	0	0	0
4	0	0	0	1
6	0	0	1	0

Řešení: /*chybí*/

2. Cvičení

1. Príklad

Kolekcia dokumentov obsahuje 4 slová: a, b, c, d. Vzájomná frekvencia slov je $a > b > c > d$. Celkový počet tokenov v kolekcii je 5000. Predpokladajte, že pre túto kolekciu presne platí Zipfov zákon. Aké sú frekvencie vyššie uvedených štyroch slov?

Řešení:

f = a > b > c > d

Aplikujeme Zipa

$$1/1k + 1/2k + 1/3k + 1/4k = 5000$$

Převédeme na společný jmenovatel

$$12/12k + 6/12k + 4/12k + 3/12k = 5000$$

Z toho vzejde že

$$25/12k = 5000, k = 2400$$

Po dosazení nám vyjde že

$$a = 2400, b = 1200, c = 800, d = 600$$

2. Príklad

γ -kódy je neefektívny pre veľké čísla (napr. 1000 alebo 10 000) pretože kódujú dĺžku offsetu v unárnom kóde. δ -kódy používajú gamma kód pre zakódovanie tejto dĺžky.

γ -kód je definovaný ako
unárny_kód(dĺžka(ofset(G))),ofset(G)

δ -kód je definovaný ako
 γ (dĺžka(ofset(G+1))),ofset(G+1)

Napríklad δ -kód pre $G=6$ je 10,0,11. 10,0 je γ -kód pre dĺžku (v tomto prípade 2). Kódovanie offsetu (11) je rovnaké ako v prípade γ -kódu pre $G = 7$.

Vypočítajte γ - a δ - kódy pre 1, 2, 3, 4, 31, 63, 127, 1023.

Řešení:

Napřed potřebujeme unární kód čísla, což je sled n jedniček ukončených nulou. Pak offset, což je binární zápis čísla, bez první jedničky. A length, což je počet bitů offsetu, reprezentováno unárně. γ -kód je pak „length, offset“. δ -kód vezme první složku (length) a zakóduje ji do gamma kódu, druhá složka zůstává stejná. S tím rozdílem že vstupní číslo se bere o jedno větší

1:

**Gamma = 0
Delta = 0**

2:

**Gamma = 10,0
Delta = 10,0,1**

3:

**Gamma = 10,1
Delta = 10,0,00**

4:

**Gamma = 110,00
Delta = 10,0,01**

a tak dále....

3. Príklad

Vypočítajte variabilný byte- a γ - kód pre postings zoznam $\langle 777, 17743, 294068, 31251336 \rangle$. Používajte medzery namiesto docID tam kde je to možné. Binárne kódy napíšte v 8 bitových blokoch.

Řešení: Zajímají nás rozdíly indexů a jejich γ - kódy.

777 = 00000110 10001001b

17743 = 00000001 00001010 11001111b

17743 – 777 = 00000001 00000100 11000110b

$\gamma = 111111111111110 00001001000110$

a tak dále...

4. Príklad

Posúďte tabuľku s frekvenciami slov troch dokumentov Doc1, Doc2, Doc3 nižšie. Vypočítajte tf-idf váhy termov *car*, *auto*, *insurance*, *best*, pre každý dokument. Idf hodnoty termov sú uvedené v tabuľke.

	Doc1	Doc2	Doc3	idf
car	27	4	24	1.65
auto	3	33	0	2.08
insurance	0	33	29	1.62
best	14	0	17	1.5

Řešení: Jednoduše se pronásobí řádky tabulky hodnotou IDf.

	Doc1	Doc2	Doc3
car	44,55	6,6	39,6
auto	6,24	68,64	0
insurance	0	53,46	46,98
best	21	0	25,5

5. Príklad

Vypočítajte normalizované Euklidovské vektory pre každý dokument z predchádzajúceho príkladu, kde každý vektor má štyri komponenty, jednu pre každý zo štyroch termov.

Řešení: Jde v podstatě jen o to, vytvořit normalizované vektory z dokumentů.

$$S = \sqrt{44,5^2 + 6,24^2 + 21^2} = 49,6$$

$$\text{Doc1}_{\text{car}} = 44,5 / 49,6 = 0,897$$

$$\text{Doc1}_{\text{auto}} = 6,24 / 49,6 = 0,126$$

$$\text{Doc1}_{\text{insurance}} = 0$$

$$\text{Doc1}_{\text{best}} = 21 / 49,6 = 0,423$$

$$\text{Doc1} = (0,897; 0,126; 0; 0,423)$$

a stejně pro další dokumenty...

6. Príklad

S váhami slov ako boli vypočítané v predchádzajúcom príklade, oznámkujte tri dokumenty podľa vypočítaného skóre pre dotaz car insurance, pre každý z nasledujúcich prípadov váženia slov:

a) váha termu je 1 ak sa v dotaze nachádza, inak 0

b) Euklidovské normalizované idf

Řešení:

a) $q = \text{„car insurance“}$

$$\text{Doc1: } (1+0+0+0) = 1$$

$$\text{Doc2: } (1+0+1+0) = 2$$

$$\text{Doc3: } (1+0+1+0) = 2$$

b) $q = \text{„car insurance“}$

$$\text{Doc1: } 0,897+0+0+0 = 0,897$$

$$\text{Doc2: } 0,0756 + 0 + 0,6127 + 0 = 0,6883$$

$$\text{Doc3: } 0,5953 + 0 + 0,7062 + 0 = 1,3015$$

7. Príklad

Vypočítajte vektor-space podobnosť medzi dotazom "digital cameras" a dokumentom "digital cameras and video cameras" doplnením prázdnych stĺpcov v tabuľke nižšie. Predpokladajte $N = 10\,000\,000$, logaritmicke váženie termov (stĺpce wf) pre dotaz aj dokumenty, idf váženie len pre dotaz a kosínovú normalizáciu len pre dokument.

"And" považujte za STOP slovo. Napíšte počty termov do tf stĺpca.

Aké je konečné skóre podobnosti?

	df	Query				Document			Product
		tf	wf	idf	$q_i = wf \cdot idf$	tf	wf	$d_i = \text{normalized wf}$	
digital	10 000								
video	100 000								
cameras	50 000								

Řešení:

$q_1 = \text{„digital cameras“}$

$q_2 = \text{„digital cameras and video cameras“}$

tf = počet slov

wf = $\log(\text{tf}) + 1$

idf = $\log(N/\text{df})$

	df	Query				Document			Product
		tf	wf	idf	$q_i = wf \cdot idf$	tf	wf	$d_i = \text{normalized wf}$	
digital	10 000	1	1	3	3	1	1	0,52	1,36
video	100 000	0	0	2	0	1	1	0,52	0
cameras	50 000	1	1	2,3	2,3	2	1,3	0,68	1,56

8. Príklad

Ukážte, že pre dotaz *affection* je radenie skóre troch dokumentov z tabuľky nižšie v opačnom poradí ako pre dotaz *jealous gossip*. Dotaz je vážený normalizáciou tf.

	SaS	PaP	WH
affection	0.996	0.993	0.847
jealous	0.087	0.120	0.466
gossip	0.017	0	0.254

Řešení: /*chybí*/

3. Cvičení

1. Příklad

IR systém vrátí 8 relevantních dokumentov a 10 nerelevantních dokumentov. Dohromady je v kolekci 20 relevantních dokumentov. Aká je presnosť a úplnosť (precision, recall) systému pri tomto vyhľadávání?

Řešení: Je potřeba spočítat precision a recall. Tyto dva se počítají podle následujících vzorců.

$$\text{Precision} = \frac{\text{RELEVANTNI}}{\text{RELEVANTNI} + \text{NERELEVANTNI}}$$

$$\text{Recall} = \frac{\text{NALEZENE RELEVANTNI}}{\text{VSECHY RELEVANTNI}}$$

$$P = 8/18 = 4/9$$

$$R = 8/20 = 2/5$$

2. Príklad

Nasledujúci zoznam písmen R a písmen N reprezentuje relevantné (R) a nerelevantné (N) dokumenty vrátené v usporiadanom zozname 20 výsledkov ako odpoveď na dotaz z kolekcie 10 000 dokumentov. Prvý (najrelevantnejší) výsledok zoznamu je naľavo. Tento zoznam obsahuje 6 relevantných dokumentov.

Predpokladajte, že kolekcia obsahuje dohromady 8 relevantných dokumentov ku dotazu.

R R N N N N N N R N R N N N R N N N N R

- Aká je presnosť systému na prvých 20 výsledkoch?
- Aká je F_1 na prvých 20 výsledkoch?
- Aká je neinterpolovaná presnosť systému pri 25% pokrytí?
- Aká je interpolovaná presnosť systému pri 33% pokrytí?
- Predpokladajte, že týchto 20 dokumentov je úplný zoznam výsledkov systému. Aký je MAP systému pre tento dotaz?

Teraz predpokladajte, že systém vrátil všetkých 10 000 dokumentov v zoradenom zozname a hore je uvedených prvých 20 vrátených výsledkov.

- Aký najvyšší možný MAP môže tento systém dosiahnuť?
- Aký najnižší možný MAP môže tento systém dosiahnuť?
- Pri sade experimentov bolo vyhodnotených len prvých 20 výsledkov. Výsledok (e) bol použitý na na aproximovanie rozsahu (f)-(g). Aká veľká môže byť chyba pre výpočet MAP pri počítaní (e) namiesto (f) a (g) pre tento dotaz?

Řešení:

a) $P = 6/20 = 3/10$

b) $F_1 = \frac{2 \cdot P \cdot R}{P + R} = \frac{2 \cdot \frac{3}{10} \cdot \frac{3}{4}}{\frac{3}{10} + \frac{3}{4}} = \frac{9}{21} = \frac{3}{7}$

c) Hledáme množinu všech P, u kterých je $R = 25\% = 1/4$. Půjdeme postupně, od prvního.

1. $P = 1/1 = 1$; $R = 1/8$

2. $P = 2/2 = 1$; $R = 2/8 = 1/4$ - toto nás zajímá

3. $P = 2/3$; $R = 2/8 = 1/4$ - takto to bude pokračovat až k dalšímu N

Pro $R = 25\% \Rightarrow P = \{1; 2/3; 1/8; 2/5; 1/3; 2/7; 2/8\}$

d) Stejný postup jako u „c“, s tím že hledáme pro $R = 1/3$.

e) Musíme spočítat MAP podle vzorce $MAP = \frac{\text{Suma } P \text{ u relevantních nálezů}}{\text{Počet } P}$.

$$MAP = \frac{1+1+\frac{3}{9}+\frac{4}{11}+\frac{5}{15}+\frac{6}{20}}{6} = 0,555$$

f) Předpokládáme, že seznam bude pokračovat ...RRNNN...

$$MAP = \frac{1+1+\frac{3}{9}+\frac{4}{11}+\frac{5}{15}+\frac{6}{20}+\frac{7}{21}+\frac{8}{22}}{6} = 0,503$$

g) Předpokládáme, že seznam bude končit ...NNRRR

$$MAP = \frac{1+1+\frac{3}{9}+\frac{4}{11}+\frac{5}{15}+\frac{6}{20}+\frac{7}{9999}+\frac{8}{10000}}{6} = 0,417$$

3. Príklad

Nižšie je tabuľka ukazujúca ako dvaja znalci ohodnotili relevanciu množiny 12 dokumentov k nejakej informačnej potrebe (0=nerlevantné, 1=relevantné).

Predpokladajme, že ste vyvinuli IR systém, ktorý pre tento dotaz vráti dokumenty {4, 5, 6, 7, 8}.

DocID	Judge 1	Judge 2
1	0	0
2	0	0
3	1	1
4	1	1
5	1	0
6	1	0
7	1	0
8	1	0
9	0	1
10	0	1
11	0	1
10	0	1

- Vypočítajte Kappa mieru zhody medzi týmito znalcami.
- Vypočítajte presnosť, pokrytie a F_1 vášho systému, ak je dokument relevantný len ak sa na ňom zhodli obaja znalci.
- Vypočítajte presnosť, pokrytie a F_1 vášho systému, ak je dokument relevantný ak si to myslí aspoň jeden zo znalcov.

Řešení:

a) Kappa míra se počítá pomocí vzorce $k = \frac{P(A)-P(E)}{(1-P(E))}$; kde $P(A) = \frac{\text{shoda všech sudích}}{\text{počet prvků}}$ a

$P(E) = p(\text{rel})^2 + p(\text{notrel})^2$; kde $p(\text{rel}) = \frac{\text{suma všech "ano" všech sudích}}{\text{počet položek} \cdot \text{počet sudích}}$ a

$p(\text{notrel}) = \frac{\text{suma všech "ne" všech sudích}}{\text{počet položek} \cdot \text{počet sudích}}$

$$k = k = \frac{\frac{1}{3} - \frac{1}{2}}{1 - \frac{1}{2}} = -\frac{1}{3}; P(A) = \frac{4}{12} = \frac{1}{3}; P(E) = \left(\frac{6+6}{24}\right)^2 + \left(\frac{6+6}{24}\right)^2 = \frac{1}{2}$$

$$\text{b) } F_1 = \frac{2 \cdot \frac{2}{12} \cdot 1}{\frac{2}{12} + 1} = \frac{1/3}{7/6} = \frac{6}{21}$$

$$\text{c) } F_1 = \frac{2 \cdot \frac{10}{12} \cdot 1}{\frac{10}{12} + 1} = \frac{5/3}{11/6} = \frac{30}{33} = \frac{10}{11}$$

4. Príklad

Užívateľov prvotný dotaz je "cheap CDs cheap DVDs extremely cheap CDs". Užívateľ preskúma dva dokumenty d1 a d2. Ohodnotí dokument d1 "CDs cheap software cheap CDs" ako relevantný a d2 "cheap thrills DVDs" nerelevantný. Predpokladajme, že používame jednoduchú tf bez dĺžkovej normalizácie vektorov. Použitím Rocchio relevance feedbacku aký by bol prepracovaný vektor dotazu po zvážení relevance feedbacku? $\alpha=1$, $\beta=0.75$, $\gamma=0.25$

Řešení: Napíšeme tabuľku, kde budú rozepsány četnosti slov dotazu a dokumentů.

	Cheap	CDs	DVDs	extremly	software	thrill
dotaz	3	2	1	1	0	0
d1	2	2	0	0	1	0
d2	1	0	1	0	0	1

Pak zkonstruujeme vektor s následujícím postupem: $(\alpha * \text{vektor dotazu}) + \frac{(\beta * \text{suma vektor relevantní})}{\text{počet relevantních}} - \frac{(\gamma * \text{suma vektor nerelevantní})}{\text{počet nerelevantních}} = \text{vysledna relevance}$

$$1*(3,2,1,1,0,0)+(0,75*(2,2,0,0,1,0))/1-(0,25*(1,0,1,0,0,1))/1=(4,25; 3,5; 0,75; 1; 0,75; -0,25)$$

5. Príklad

Prečo je pozitívny feedback pravdepodobne lepší ako negatívny feedback pre IR systém? Prečo je možno lepšie použiť na feedback len jeden nerelevantný dokument ako ich použiť viac?

Řešení: Protože závěry je-li nebo není-li lepší systém s negativním feedbackem jsou dosti nepřesné, a pozitivní feedback nám stačí. Jako negativní je lepší použít jen jeden nejvyšší.

6. Príklad

Prečo je prírastková relevancia viacej realistické merítka užívateľskej spokojnosti? Udajte príklad kde neprírastková metrika ako napríklad presnosť alebo úplnosť je zavádzajúce merítka užívateľskej spokojnosti a naopak prírastková je lepšia?

Řešení: Řekl bych, že je to tehdy, kdy uživatel hledá spíše podle toho, který dokument vidělo více uživatelů.

4. Cvičení

1. Příklad

Každý z dvou webových vyhledávacích systémů A a B zo svojich indexov generujú veľké množstvo stránok rovnomerne náhodne. 30% stránok z A sa nachádza v indexe B a 50% stránok z B sa nachádza v indexe A. Aký je pomer stránok medzi systémami A a B?

Řešení:

$$\begin{aligned}x|A| &= y|B| \\ 0,3|A| &= 0,5|B| \\ \frac{|A|}{|B|} &= \frac{5}{3}\end{aligned}$$

2. Příklad

Každý z dvou webových vyhledávacích systémů A a B zbierajú (crawl) náhodnú, ale rovnako veľkú podmnožinu Webu. Niektoré zozbierané stránky sú duplikáty - presné textové kópie na rôznych URL.

Predpokladajte, že sú duplikáty distribuované rovnomerne medzi stránkami zozbierané systémom A aj B. Ďalej predpokladajte, že duplikát má presne dve kópie - žiadne stránky nemajú viac ako dve kópie. A indexuje stránky bez eliminácie duplikátov, kdežto B indexuje len jednu kópiu duplikovaných stránok. Tieto dve náhodné podmnožiny majú rovnakú veľkosť pred odstránením duplikátov.

Ak sa 45% stránok z A nachádza v indexe B, a 50% stránok z B v indexe A, aká veľká časť Webu sa skladá zo stránok, ktoré nemajú duplikáty?

Řešení:

$$\begin{aligned}x|A| &= y\left(|B| - \frac{D}{2}\right) \\ 0,45|A| &= 0,5\left(|B| - \frac{D}{2}\right) \\ 0,9|A| &= |B| - \frac{D}{2} \\ 1,8|A| &= 2|B| - D \\ -0,2|A| &= D // |A| \text{ a } |B| \text{ jsou si rovny} \\ D &= 0,2 * 0,5 = 0,1 // |A| \text{ je polovina} \\ D &= 10\% \text{ webu}\end{aligned}$$

3. Príklad

Daný je nasledujúci web graf.

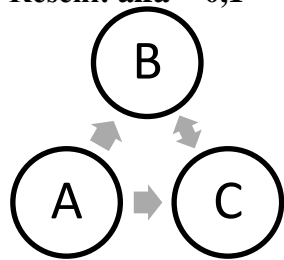
$V = \{a, b, c\}$, $E = \{a \rightarrow b, a \rightarrow c, b \rightarrow c, c \rightarrow b\}$

Vypočítajte PageRank, hub skóre a autoritatívne skóre pre každú z troch stránok. Zoradíte stránky podľa jednotlivých skóre a pozorujte prípadné väzby.

Pre výpočet PageRank môžete predpokladať, že sa v každom kroku náhodnej prechádzky teleportujeme na náhodnú stránku s pravdepodobnosťou 0.1 a s rovnomernou distribúciou stránok, na ktoré sa teleportujeme.

Pre huby a autority normalizujte skóre tak, aby maximum bolo 1.

Řešení: $\alpha = 0,1$



$$\begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \text{ (každý riádek musí dávať súčet 1) } = \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} * (1 - \alpha) =$$
$$\begin{pmatrix} 0 & \frac{9}{20} & \frac{9}{20} \\ 0 & 0 & \frac{9}{10} \\ 0 & \frac{9}{10} & 0 \end{pmatrix} \text{ nuly nahradíme } \left(\frac{\alpha}{\text{počet uzlů}} \right) = \begin{pmatrix} \frac{1}{30} & \frac{9}{20} & \frac{9}{20} \\ \frac{1}{30} & \frac{1}{30} & \frac{9}{10} \\ \frac{1}{30} & \frac{9}{10} & \frac{1}{30} \end{pmatrix} * 30 =$$
$$= \begin{pmatrix} 1 & 90 & 90 \\ 1 & 1 & 180 \\ 1 & 180 & 1 \end{pmatrix}$$

Dále vypočítám vlastní vektor matice

$$\begin{bmatrix} 1 - \lambda & 90 & 90 \\ 1 & 1 - \lambda & 180 \\ 1 & 180 & 1 - \lambda \end{bmatrix} = \dots \text{ Vlastní vektor je } (1/30, 29/60, 29/60)$$

4. Príklad

Priemerný vstupný stupeň všetkých uzlov vybraného grafu webu je 9. Čo môžeme povedať o priemernom výstupnom stupni všetkých uzlov tohto grafu?

Řešení:

$$A = \begin{pmatrix} 0 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \quad A^T = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{pmatrix}$$

$$\vec{h} = A * A^T * \vec{h} = \begin{pmatrix} 2 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix} * \vec{h} = (4, 2, 2) - \text{Hub, odkazy na authority}$$

$$\vec{a} = A^T * A * \vec{a} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 2 & 1 \\ 0 & 1 & 2 \end{pmatrix} * \vec{a} = (0, 3, 3) - \text{Authority, obsah jako takový}$$

5. Cvičení

1. Příklad

Pre XML dokument uvedený nižšie napíšte XPath výrazy.

- Vráťte všetky názvy (title elementy), ako kurzu, tak oddelenia.
- Vráťte názvy kurzov, ktoré majú v názve výraz "programming"
- Vráťte priezviská inštruktorov učiacich aspon jeden kurz, ktorý má vo svojom popise slovo "software".
- Vráťte priezviská profesorov učiacich aspon jeden kurz, ktorý má vo svojom popise slovo "software".

```
<Course_Catalog>
  <Department Code="CS">
    <Title>Computer Science</Title>
    <Chair>
      <Professor>
        <First_Name>Jennifer</First_Name>
        <Last_Name>Widom</Last_Name>
      </Professor>
    </Chair>
    <Course Number="CS106A" Enrollment="1070">
      <Title>Programming Methodology</Title>
      <Description>Introduction to the engineering of computer applications
emphasizing modern software engineering principles.</Description>
      <Instructors>
        <Lecturer>
          <First_Name>Jerry</First_Name>
          <Middle_Initial>R.</Middle_Initial>
          <Last_Name>Cain</Last_Name>
        </Lecturer>
        <Professor>
          <First_Name>Eric</First_Name>
          <Last_Name>Roberts</Last_Name>
        </Professor>
        <Professor>
          <First_Name>Mehran</First_Name>
          <Last_Name>Sahami</Last_Name>
        </Professor>
      </Instructors>
    </Course>
    <Course Number="CS106B" Enrollment="620">
      <Title>Programming Abstractions</Title>
      <Description>Abstraction and its relation to programming.</Description>
      <Instructors>
        <Professor>
          <First_Name>Eric</First_Name>
          <Last_Name>Roberts</Last_Name>
        </Professor>
        <Lecturer>
          <First_Name>Jerry</First_Name>
          <Middle_Initial>R.</Middle_Initial>
          <Last_Name>Cain</Last_Name>
        </Lecturer>
      </Instructors>
      <Prerequisites>
        <Prereq>CS106A</Prereq>
      </Prerequisites>
    </Course>
  </Department>
</Course_Catalog>
```

Řešení: základní příkazy XPath – „//“ znamená cokoliv, contains(kde, co)

- a) //Title
- b) //Course[contains(Title,'Programming')]/Title
- c) //Course[contains(Description,'Software')]/Instructors/Last_Name
- d) //Course[contains(Description,'Software')]/Instructors/Professor/Last_Name

2. Příklad

Vypočítajte podobnosť medzi dotazmi a im zodpovedajúcimi cestami v dokumente z Příkladu 1.

- a) //Instructors/Last_Name#Cain
- b) //Course/Instructors/Lecturer/Last_Name#Cain

Řešení: Potřebujeme určit kontextovou podobnost, což je funkce

$CR = \frac{1+\text{počet uzlů v dotazu}}{1+\text{počet uzlů v dokumentu, pro které to sedí}}$ nebo pokud je dotaz špatný tak je $CR = 0$.

a) $\frac{1+2}{1+6} = \frac{3}{7}$

b) $\frac{1+4}{1+6} = \frac{5}{7}$

3. Příklad

Spočítajte, koľko štruktúrnych termov (structural terms, dvojíc kontext/term <c,t>) je v XML strome na nižšie.

```
<Course>
  <Title>Programming Abstractions</Title>
  <Description>Abstraction and its relation to programming</Description>
  <Instructors>
    <Professor>
      <First_Name>Eric</First_Name>
      <Last_Name>Roberts</Last_Name>
    </Professor>
  </Instructors>
</Course>
```

Řešení: Musí se počítat pro každé slovo, všechny možné podstromy.

Programming Abstractions – **dvě slova**

<Title>Programming Abstractions</Title> – **první podstrom**

<Course>

<Title>Programming Abstractions</Title> – **druhý podstrom**

</Course>

2*(1+1) = 4 – A dále pro každé další slova

6*(1+1) = 12

1*(1+1+1+1) = 4

1*(1+1+1+1) = 4

4+12+4+4=24

4. Příklad

Ktorý z dokumentov uvedených nižšie má rovnaké alebo rozdielne bag of words reprezentácie pre Bernouliho a multinomický model? Aké sú rozdiely

Doc1: He moved from London, Ontario, to London, England.

Doc2: He moved from London, England, to London, Ontario.

Doc3: He moved from England to London, Ontario.

Řešení:

Podle Bernouliho jsou všechny stejné, protože Bernouli nezahrnuje duplikáty.

{London, London} = {London}

Podle multinomického modelu jsou různé, protože zahrnuje duplikáty

{London, London} ≠ {London}

5. Príklad

Na základe dát z tabuľky nižšie

- odhadnite multinomické Naive Bayes klasifikátory,
- aplikujte ich na testovací dokument,
- odhadnite Bernoulli Naive Bayes klasifikátor,
- aplikujte ho na testovací dokumnet

Nemusíte odhadovať parametre, ktoré na klasifikáciu dokumnetu nie sú potrebné.

	docID	obsah dokumentu	je v c = China?
trénovacia množ.	1	Taipei Taiwan	áno
	2	Macao Taiwan Shanghai	áno
	3	Japan Sapporo	nie
	4	Sapporo Osaka Taiwan	nie
testovacia množ.	5	Taiwan Taiwan Sapporo	?

Řešení:

Klasifikátory se určují pro každé slovo podle

$$P(\text{slovo}|+/-) = \frac{(\text{počet výskytů } v + \text{nebo}-) + 1}{(\text{počet slov } v + \text{nebo}-) + (\text{počet různých slov v sadě})}$$

$$P(+/-) = \frac{(\text{počet} + \text{nebo} - \text{dokumentu})}{\text{počet všech dokumentu}}$$

Pro konkrétní aplikaci pak

$$P(+/-|\text{test}) = P(+/-) * (\text{všechny váhy slov v testu pronásobené})$$

a)

$$P(\text{Taiwan}|+) = (2+1)/(5+7) = 1/4$$

$$P(\text{Saporro}|+) = (0+1)/(5+7) = 1/12$$

$$P(+)= 2/4 = 1/2$$

$$P(\text{Taiwan}|-) = (1+1)/(5+7) = 1/6$$

$$P(\text{Saporro}|-) = (2+1)/(5+7) = 1/4$$

$$P(-) = 2/4 = 1/2$$

b)

$$P(+|\text{test}) = 1/2 * 1/4 * 1/4 * 1/12 = 0,0026$$

$$P(-|\text{test}) = 1/2 * 1/6 * 1/6 * 1/4 = 0,0035$$

c)

$$P(\text{Taiwan}|+) = (2+1)/(2+2) = 3/4$$

$$P(\text{Saporro}|+) = (0+1)/(2+2) = 1/4$$

$$P(+)= 2/4 = 1/2$$

$$P(\text{Taiwan}|-) = (1+1)/(2+2) = 1/2$$

$$P(\text{Saporro}|-) = (2+1)/(2+2) = 3/4$$

$$P(-) = 2/4 = 1/2$$

d)

$$P(+|\text{test}) = 1/2 * 3/4 * 1/4 = 0,0937$$

$$P(-|\text{test}) = 1/2 * 1/2 * 3/4 = 0,1875$$