Nekyus 1:

Търсене чрез булеви заявен от Ключови думи. Обратен индекс.

stoyan @ lml.bas.bg

Оценяване на извлеченить документи

- Прецизност: Процентният дял на репевантните по отношение на потребността докименти спряно всички извлечени.
- Обхват: Процентот на извлечените релевантни докушенти спряно всиски релевантни докушенти в колекцията.

Матрица на срещания иногова длиа/докушент

	Антоний и Клеопатра		Бурята	Хамлет	Отело	Макбет
Антоний	1	1	0	0	0	1
Брут	1	1	0	1	0	0
Цезар	1	1	0	1	1	1
Калпурния	0	1	0	0	0	0
Клеопатра	1	0	0	0	0	0
милост	1	0	1	1	1	1
по-лош	1	0	1	1	1	0

1, ако пиесата съдържа съответната дума, 0 иначе.

Вектори на срещанията

На всяка клюхова дчиа съпоставяне вектор от 0/1 с размер броя на докушентите, който отразява срещанията За да отговорим на заявката *Брут и Цезар и не Калпурния*, извършваме побитови операции със съответните вектори:

•	110100 &		Антоний и Клеопатра		Бурята	Хамлет	Отело	Макбет
	110111 &	Антоний	1	1	0	0	0	1
		Брут	1	1	0	1	0	0
		Цезар	1	1	0	1	1	1 .
		Калпурния	0(1)	1 (0)	0(1)	0 (1)	0(1)	0 (1)
٠	~010000 = 101111	Клеопатра	1	0	0	0	0	0
		милост	1	0	1	1	1	1
	100100	по-лош	1	0	1	1	1	0

1 Камлет Антоний и Theonarpa

• Проблем: Представянето на могрицата на срещанията изисква

творде много памет • Наблюдение: Броят на единиците е малъх. Матрицата е шино

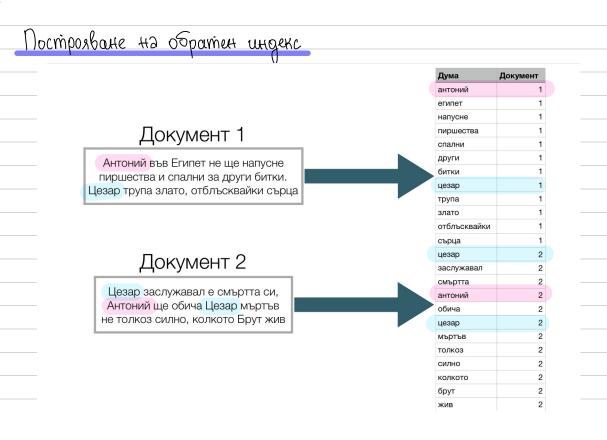
• Решение: Запаметяване само позициите на единиците

Етапи на предварителна обработка на текста

- Разбиване текста на единици (Tokenization)
 изрази като: 25-милиметров, полу-слято...
- · Hopmanusayus
 - мавни, махи бучви, съхращения, акроници: С.У. Текст
- Преобразуване до основни форши (Stemming)
 дунавски -> дунав, вървейки -> вървя
- · Tipenarbatte tha stop ginn
 - Съюзи, мендинетия, предлози...

Технологии и инструменти за предварителна обработка на текст

- Най-често базирани на експертни правила и речници
 - използване на регулярни изрази и релации
 - използване на речници за съкращения, акроними, ...
 - правила за преобразуване и замяна
 - използване на крайни автомати и преобразуватели



Сортираме спистка първо по дините, после по докунент.

Дума	Документ
тоний	1
египет	1
напусне	1
пиршества	1
спални	1
други	1
битки	1
цезар	1
трупа	1
злато	1
отблъсквайки	1
сърца	1
цезар	2
аслужавал	2
мъртта	2
нтоний	2
обича	2
цезар	2
мъртъв	2
толкоз	2
силно	2
колкото	2
брут	2
жив	2

Създаване на обратен индекс

- Обединяваме повторенията
- Разделяме речника от срещанията
- Добавяме поле за честота (броя) на срещанията

•			
	Дума	Документ	
	антоний	1	
	антоний	2	
	битки	1	
	брут	2	
	други	1	
	египет	1	
	жив	2	
	заслужавал	2	
	злато	1	
	колкото	2	
	мъртъв	2	
	напусне	1	
	обича	2	
	отблъсквайки	1	
	пиршества	1	
	силно	2	
	смъртта	2	
	спални	1	
	сърца	1	
	толкоз	2	
	трупа	1	
	цезар	1	
	цезар	2	

	Речни	Спис	ък с	реща	пия	
	Дума	Честота	0.00000			
	антоний	(2)	†	1	→	2
	битки	1	1	1		
	брут	1	1	2		
	други	1	→	1		
	египет	1	1	1		
	жив	1	→	2		
	заслужавал	1	+	2		
	злато	1	+	1		
	колкото	1	→	2		
>	мъртъв	1	→	2		
	напусне	1	→	1		
	обича	1	→	2		
	отблъсквайки	1	→	1		
	пиршества	1	→	1		
	силно	1	→	2		
	смъртта	1	→	2		
	спални	1	+	1		
	сърца	1	→	1		
	толкоз	1	→	2		
	трупа	1	→	1		
	цезар	2	→	1	→	2
		_		_	_	_

Булево търсене грез обратен индекс

Респизиране на коннонкция: сливане на срещенията

Разпендаме заявката Бриг и Цегар

ВАЖНО: Списъците трябва да са сортирани по номер на документ

Epyr
$$2 \rightarrow 4 \rightarrow 8 \rightarrow 16 \rightarrow 32 \rightarrow 64 \rightarrow 128$$

4esap $1 \rightarrow 2 \rightarrow 3 \rightarrow 5 \rightarrow 8 \rightarrow 13 \rightarrow 21 \rightarrow 34$

Altopum BM 3a Chubahe Ha copthpathy chucbyn

Алгоритом за произволни булеви заявки

Au B -> C

Сливаме всеки от подизразите и получаваме списък от срещания, след което сливаме получените списъци до получаването на срещанията за цялата заявка:

Конюнкция на п терия: Оптинизация на изпълнението на заявкато: запосване с най-късите списъци и сливаме в нарастващ ред

Позиционен обрамен индекс

Съставни имена, фрази и изрази
Примери: Иван Иванов, Стара Захора, операционня системв...

- Първо решение добавяне на двойки от думи към речника
- Огромно нарастване на речника и списъците от срещания
 - необходимо разбиване на двойки и проверка дали са последователниИма смисъл само за често срещани съставни имена и

При нужда от намиране на тройки или четворки е

_ изрази

Алтернативно рещение: използване на позиционен индекс

Яби списъкы от срещания на всяко срещане на теры в документ добавяме списък на позициите в документа на съответните срещания

Mpunep: ga voge un ga he voge:

да: бъде: 2 -> 1,17,74,222,551; 1 -> 17,19; 4 -> 8,16,190,429,433; 4 -> 17,191,291,430,434; 7 -> 13,23,191; 5 -> 14,19,101;

ymetit

- Позволява в булевата заявка да се добави ограничение за разстоянието между термовете:
 - да и/1 бъде

компютърна и/3 мрежа

- Изисква съществено допълване на алгоритмите за сливане
- Изисква между 2 и 4 пъти повече памет за представяне на индексите.

Толериране на форми на дадена дума

- · Чрез използване на stemming:
 - позиционен -> позиция, дунавската -> дунав, ...
 - втори речник върху "основите" на думите, за всяка основа се посочва списък от ключови думи със съответната основа
- Води до повече резултати, в някои случай нерелевантни:
 - хлебарка -> хляб, националистически -> национален -> национал, ...
 - Не се толерират грешки в изписването