

# Introduction to Information Retrieval

معرفی بازیابی اطلاعات  
علی قنبری

## مرجع درس

---

- **An Introduction to Information Retrieval**, Christopher D. Manning, Prabhakar Raghavan, Hinrich Schütze. Cambridge University Press, Cambridge, England.

# ارزیابی

---

- میانترم ۵ نمره
- پایانترم ۱۲ نمره
- تمرین ۳ نمره
- پروژه ؟ نمره
- حضور و غیاب ۱ نمره

Ali.ghanbari289@gmail.com

# Introduction to Information Retrieval

فصل اول: بازیابی بولی  
علی قنبری

# تعریف بازیابی اطلاعات

Information retrieval (IR) is **finding** material (**usually documents**) of an **unstructured** nature (usually text) that satisfies an **information need** from within **large collections** (usually stored on computers).

بازیابی اطلاعات یافتن مواد (معمولا اسناد) از یک ماهیت بدون ساختار (معمولا متن) است که یک نیاز اطلاعاتی را داخل مجموعه‌های بزرگ (که معمولا در کامپیوتر ذخیره می‌شوند) برآورده می‌کند.

## مثالی از بازیابی اطلاعات

- مجموعه کتب شکسپیر شامل کتاب‌های قطوری است که بسیاری از افراد آن را دارند. فرض کنید می‌خواهید تعیین کنید که کدام نمایش شکسپیر شامل کلمات Brutus و Caesar است و شامل Calpurnia نیست.
- یک راه این است که از ابتدای کتاب شروع کنید و خط به خط متن را بخوانید و هر نمایشی که شامل اسم‌های Brutus و Caesar بود، ثبت کرده و نمایشی که اسم Calpurnia را داشت، مستثنی کنید. ساده‌ترین صورت بازیابی اسناد برای کامپیوتر، این نوع اسکن خطی در میان اسناد است.

# مثالی از بازیابی اطلاعات

- راه جلوگیری از اسکن خطی متون برای هر پرس و جو این است که اسناد را از پیش شاخص گذاری کنیم.
- فرض کنید که برای هر سند (در اینجا نمایش‌های شکسپیر) ثبت می‌کنیم که آیا شامل هر یک از کلماتی است که شکسپیر به کار برده است یا خیر (حدود ۳۲۰۰۰ کلمه مختلف شکسپیر به کار برده است).
- نتیجه عمل بالا یک ماتریس تلاقی دودویی عبارت - سند است.

# مثالی از ماتریس تلاقی

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth ...
ANTHONY	1	1	0	0	0	1
BRUTUS	1	1	0	1	0	0
CAESAR	1	1	0	1	1	1
CALPURNIA	0	1	0	0	0	0
CLEOPATRA	1	0	0	0	0	0
MERCY	1	0	1	1	1	1
WORSER	1	0	1	1	1	0
...						

یک به معنای وجود کلمه در نمایشنامه است و صفر به معنای عدم وجود کلمه در نمایشنامه است.



# بازیابی بولی

- مدل بولی، یکی از ساده ترین مدل های بازیابی اطلاعات است.
- Query ها عبارات بولی مانند CAESAR AND BRUTUS هستند.
- موتورهای جستجو همه اسنادی را برمیگرداند که عبارت بولین query برای آنها درست است.
- بسته به اینکه به سطر یا ستون ماتریس نگاه کنیم می توانیم برداری برای هر عبارت داشته باشیم که اسنادی را که در آن پدیدار می شوند، نشان می دهد. یا آنکه برداری برای هر سند داشته باشیم که عباراتی را که در آن ظاهر شده اند، نشان می دهد.

## بردار وقوع (Incidence vectors)

✓ برای پاسخ به پرسش BRUTUS AND CAESAR AND NOT CALPURNIA

- بردارهای متناظر با BRUTUS, CAESAR, CALPURNIA را برمی داریم.

- مکمل بردار CALPURNIA را محاسبه می کنیم.

- عملگر and بیتی روی سه بردار انجام می دهیم.

$$110100 \text{ AND } 110111 \text{ AND } 101111 = 100100$$

	Anthony and Cleopatra	Julius Caesar	The Tempest	Hamlet	Othello	Macbeth ...
ANTHONY	1	1	0	0	0	1
BRUTUS	1	1	0	1	0	0
CAESAR	1	1	0	1	1	1
CALPURNIA	0	1	0	0	0	0
CLEOPATRA	1	0	0	0	0	0
MERCY	1	0	1	1	1	1
WORSER	1	0	1	1	1	0
...						
result:	1	0	0	1	0	0

## پاسخ به پرسش

*Anthony and Cleopatra, Act III, Scene ii*

Agrippa [Aside to Domitius Enobarbus]: Why, Enobarbus,  
When Antony found Julius Caesar dead,  
He cried almost to roaring; and he wept  
When at Philippi he found Brutus slain.

*Hamlet, Act III, Scene ii*

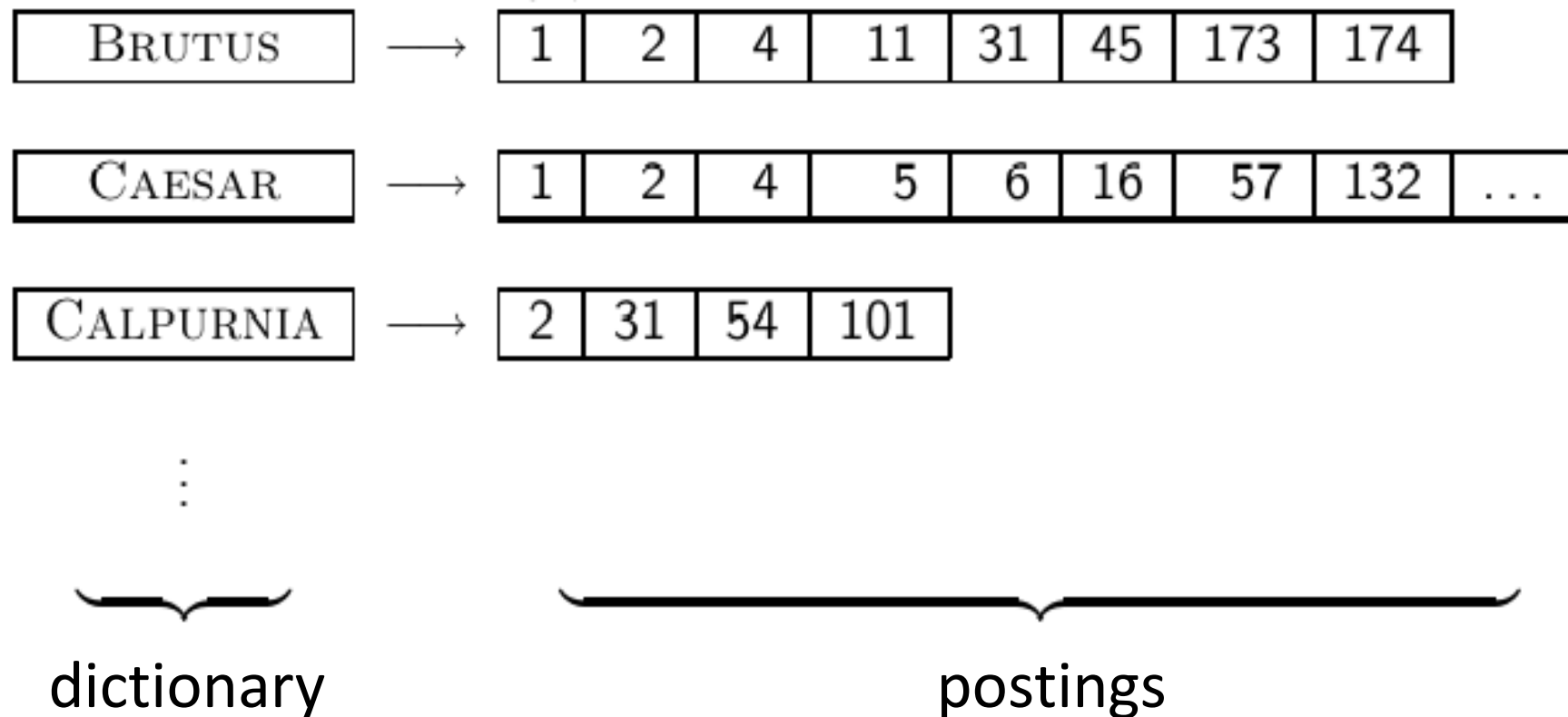
Lord Polonius: I did enact Julius Caesar: I was killed i'  
the Capitol; Brutus killed me.

# مجموعه‌های بزرگ

- تعداد  $N = 10^6$  متن با حدود 1000 کلمه (۲ تا ۳ صفحه کتاب) را در نظر بگیرید.
- اگر هر کلمه ۶ بایت حافظه بخواهد، مجموعه اسناد حدود ۶ گیگابایت خواهد بود.
- فرض کنیم 500'000 کلمه مجزا در این اسناد وجود داشته باشد.
- ماتریس تلاقی شامل  $10^6 * 500'000$  عنصر خواهد بود.
- می‌توان نشان داد که حداقل 99.8% عناصر این ماتریس صفر است.
- این ماتریس **تنگ** است.
- نمایش بهتر از ماتریس تلاقی آن است که فقط عناصر یک را ذخیره کنیم.

# شاخص معکوس

برای هر کلمه، لیستی از همه اسنادی که شامل آن کلمه هستند را ذخیره می‌کنیم.



# ساخت شاخص معکوس

1 جمع‌آوری اسنادی که باید شاخص گذاری شوند.

Friends, Romans, countrymen. So let it be with Caesar ...

2 تبدیل هر سند به لیستی از نشانه‌ها

Friends Romans countrymen So ...

3 پیش پردازش، تولید لیستی از نشانه‌های نرمال شده

friend roman countryman so ...

4 شاخص گذاری اسنادی که هر عبارت در آن رخ می‌دهد. (به کمک ایجاد شاخص وارونه شامل لغت نامه و پست‌ها)

## مرحله دوم و سوم

**Doc 1.** I did enact Julius Caesar: I was killed i' the Capitol; Brutus killed me.

**Doc 2.** So let it be with Caesar. The noble Brutus hath told you Caesar was ambitious:



**Doc 1.** i did enact julius caesar i was killed i' the capitol brutus killed me

**Doc 2.** so let it be with caesar the noble brutus hath told you caesar was ambitious



**Doc 1.** i did enact julius caesar i was  
killed i' the capitol brutus killed me  
**Doc 2.** so let it be with caesar the  
noble brutus hath told you caesar was  
ambitious



term	docID
i	1
did	1
enact	1
julius	1
caesar	1
i	1
was	1
killed	1
i'	1
the	1
capitol	1
brutus	1
killed	1
me	1
so	2
let	2
it	2
be	2
with	2
caesar	2
the	2
noble	2
brutus	2
hath	2
told	2
you	2
caesar	2
was	2
ambitious	2

# مرحله چهارم – گام ۱

- درون یک مجموعه سند، فرض می‌کنیم که هر سند یک شماره سریال منحصر به فرد به نام شناسه سند دارد.
- در طول ساخت شاخص به راحتی می‌توانیم یک عدد صحیح متوالی را برای هر سند جدید، هنگامی که اولین بار با آن مواجه شدیم، اختصاص دهیم.

# مرحله چهارم – گام ۲ (مرتب سازی بر اساس حروف الفبا)

term	docID		term	docID
i	1		ambitious	2
did	1		be	2
enact	1		brutus	1
julius	1		brutus	2
caesar	1		capitol	1
i	1		caesar	1
was	1		caesar	2
killed	1		caesar	2
i'	1		did	1
the	1		enact	1
capitol	1		hath	1
brutus	1		i	1
killed	1		i	1
me	1	⇒	i'	1
so	2		it	2
let	2		julius	1
it	2		killed	1
be	2		killed	1
with	2		let	2
caesar	2		me	1
the	2		noble	2
noble	2		so	2
brutus	2		the	1
hath	2		the	2
told	2		told	2
you	2		you	2
caesar	2		was	1
was	2		was	2
ambitious	2		with	2

# مرحله چهارم - گام ۳ (ساخت لیست پست‌ها و تعداد تکرار آن‌ها)

term	docID		term	doc. freq.	→	postings lists
ambitious	2		ambitious	1	→	2
be	2		be	1	→	2
brutus	1		brutus	2	→	1 → 2
brutus	2		capitol	1	→	1
capitol	1		caesar	2	→	1 → 2
caesar	1		did	1	→	1
caesar	2		enact	1	→	1
caesar	2		hath	1	→	2
did	1		i	1	→	1
enact	1		i'	1	→	1
hath	1		it	1	→	2
i	1		julius	1	→	1
i	1		killed	1	→	1
i'	1		let	1	→	2
it	2		me	1	→	1
julius	1		noble	1	→	2
killed	1		so	1	→	2
killed	1		the	2	→	1 → 2
let	2		told	1	→	2
me	1		you	1	→	2
noble	2		was	2	→	1 → 2
so	2		with	1	→	2
the	1					
the	2					
told	2					
you	2					
was	1					
was	2					
with	2					

# مثال

- شاخص وارونه‌ای برای مجموعه اسناد زیر رسم کنید.
- سند ۱:

new home sales top forecasts

- سند ۲:

home sales rise in july

- سند ۳:

increase in home sales in july

- سند ۴:

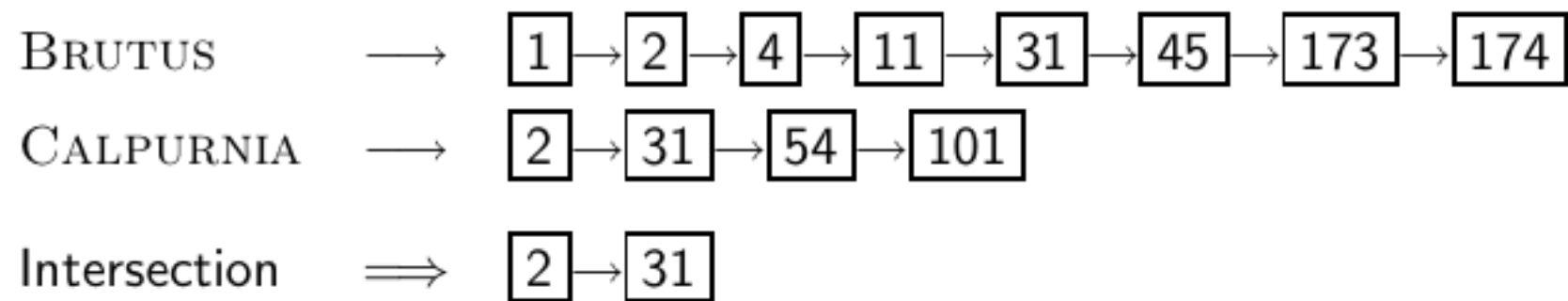
july new home sales rise

فرض کنید همه نشانه‌های به کار رفته در اسناد، نرمال هستند.

## یک مثال از پرس و جوی عطفی ساده به کمک شاخص معکوس

- پرس و جوی BRUTUS AND CALPURNIA را در نظر بگیرید.
- برای پیدا کردن پاسخ پرس و جوی بالا به کمک شاخص معکوس:
  ۱. محل BRUTUS را در دیکشنری پیدا می کنیم.
  ۲. پست‌های آنرا بازیابی می کنیم.
  ۳. محل CALPURNIA را در دیکشنری پیدا می کنیم.
  ۴. پست‌های آنرا بازیابی می کنیم.
  ۵. اشتراک دو لیست پست را محاسبه می کنیم.

## یک مثال از پرس و جوی عطفی ساده به کمک شاخص معکوس (ادامه)

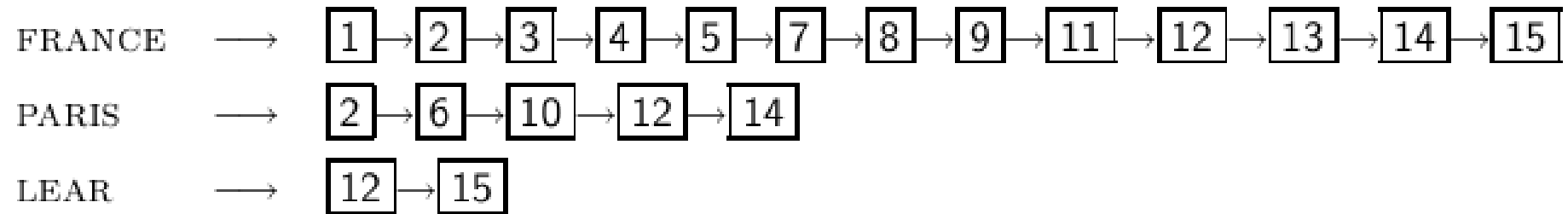


# الگوریتم اشتراک دو لیست پست

```
INTERSECT( $p_1, p_2$ )  
1   $answer \leftarrow \langle \rangle$   
2  while  $p_1 \neq \text{NIL}$  and  $p_2 \neq \text{NIL}$   
3  do if  $\text{docID}(p_1) = \text{docID}(p_2)$   
4      then  $\text{ADD}(answer, \text{docID}(p_1))$   
5           $p_1 \leftarrow \text{next}(p_1)$   
6           $p_2 \leftarrow \text{next}(p_2)$   
7  else if  $\text{docID}(p_1) < \text{docID}(p_2)$   
8      then  $p_1 \leftarrow \text{next}(p_1)$   
9      else  $p_2 \leftarrow \text{next}(p_2)$   
10 return  $answer$ 
```

این الگوریتم فقط زمانی درست است که لیست پست‌ها مرتب باشند.

## مثال دیگر



نتیجه پرس و جوی ((paris AND NOT france) OR lear را محاسبه نمایید.