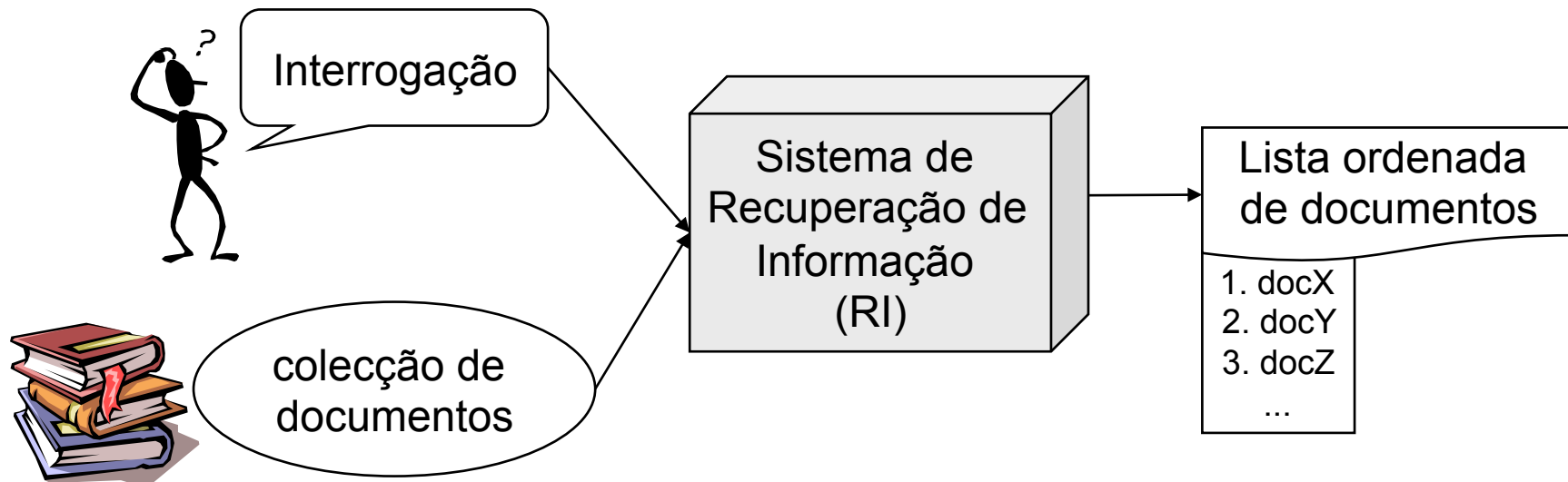


# Recuperação de Informação – modelo Booleano

## Modelo usual de funcionamento

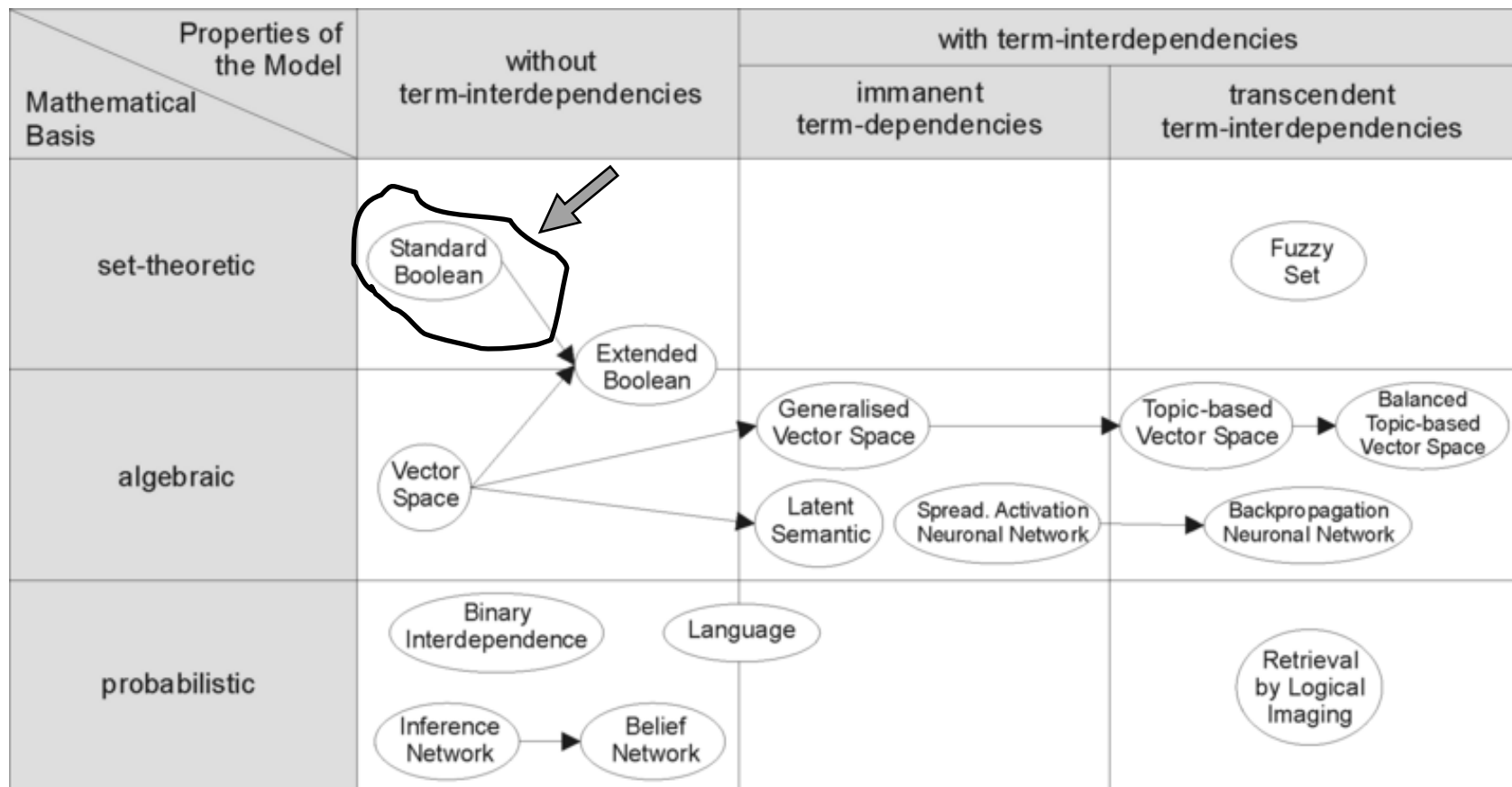
- Entrada:
  - uma colecção de documentos com texto escrito em língua natural
  - uma interrogação do utilizador escrita como sequência de caracteres
- Saída:
  - uma lista ordenada de documentos relevantes para a interrogação



## ... recuperação de informação em documentos de texto

- Processos
  - indexar documentos
  - extrair os documentos com determinada informação
- Preocupações
  - 1. identificar os documentos relevantes
  - 2. extrair eficazmente grandes quantidades de documentos
- ... existem diferentes modelos para implementar os processos
  - contemplando as preocupações!

# Classificação dos modelos de recuperação de informação



## Modelo Booleano – sem interdependência de termos

- A interrogação tem a forma de uma expressão Booleana de termos
  - e.g. gato AND rato AND NOT Rússia

O rato roeu a rolha do garrafão do rei da Rússia; o rei ficou zangado!

docA.txt

O gato e o rato vivem no palácio do rei.

docB.txt

"Gato branco, gato preto" é um filme (de Kusturika) com imagens surrealistas.

docC.txt

colecção de documentos



docB.txt

gato AND rato AND NOT Rússia



## ... funcionamento geral do processo de indexação

Considerando uma colecção de documentos:

- Percorrer todos os documentos e, para cada documento,
  - construir o conjunto de todos os termos desse documento
- Fazer a união de todos os conjuntos de termos
- Construir matriz de incidência  $\omega$ 
  - $\omega_{ik} = 1$ , se o termo  $i$  ocorre no documento  $k$
  - $\omega_{ik} = 0$ , se o termo  $i$  não ocorre no documento  $k$

## Exercício – construir matriz de incidência

O rato roeu a rolha do  
garrafão do rei da  
Rússia; o rei ficou  
zangado!

docA.txt

O gato e o rato vivem  
no palácio do rei.

docB.txt

"Gato branco, gato  
preto" é um filme (de  
Kusturika) com  
imagens surrealistas.

docC.txt

coleção de documentos

## ... exemplo – conjunto de termos e matriz de incidência

docA {  
o  
rato  
roeu  
a  
rolha  
do  
garrafão  
do  
rei  
da  
rússia  
;  
o  
rei  
ficou  
zangado  
!  
- - - - -  
o  
gato  
e  
o  
rato  
vivem  
no  
palácio  
...

docB {

alguns termos	documentos		
	docA	docB	docC
...	...	...	...
rato	1	1	0
rússia	1	0	0
gato	1	1	1
palácio	0	1	0
filme	0	0	1
...	...	...	...



gato AND rato AND NOT Rússia

**Como usar a matriz de incidência  
para obter a resposta?**

?



## ... recuperação de informação com o modelo Booleano



gato AND rato AND NOT Rússia

	docA	docB	docC
...	...	...	...

111 AND 110 AND (NOT 100) =  
111 AND 110 AND 011 =  
110 AND 011 =  
**010**

### Operadores:

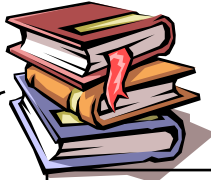
AND, OR e NOT

bit a bit (*bitwise*) sobre  
cada linha, da matriz  
de incidência, relativa  
a cada um dos termos  
da interrogação.

docB.txt



## Exemplo – vamos aumentar a dimensão da colecção



Admita-se uma colecção  
com a seguinte dimensão:

nº de documentos = 1M ( $\approx 1.000.000$ )  
nº de palavras por documento = 1K ( $\approx 1.000$ )  
dimensão de cada palavra = 6 byte  
Total de informação =  $1M \times 1K \times 6 =$   
 $= 6G$  de informação nos documentos

Atenção!  
aqui estão palavras e não termos

o  
rato  
roeu  
a  
rolha  
...

Admita-se um conjunto de termos  
com a seguinte dimensão:

nº de termos distintos = 500K ( $\approx 500.000$ )

E a matriz de incidência

- qual a sua dimensão total?
- qual a proporção de 0s e 1s?

## Exemplo – matriz de incidência muito esparsa

- Dimensão total da matriz de incidência
  - 500K termos distintos x 1M documentos
  - $\approx$  meio trilião ( $\frac{1}{2} 10^{12}$ ) de células
  - $\approx$   $\frac{1}{2}$  terabyte de memória (só para 1M de documentos)!
- Número máximo de 1s na matriz de incidência
  - 1M documentos, cada um com 1K termos, ou seja,
  - ... há no máximo  $1M \times 1K = 1G$  (1 bilhão,  $10^9$ ) de 1s na matriz
  - ... i.e. caso em que todas as palavras de cada documento são termos
- Proporção de 0s e 1s
  - $(10^9 \div \frac{1}{2} 10^{12}) \times 100 = (2 \times 10^{-3}) \times 100 = 0.002 \times 100 =$
  - **no máximo, 0.2% de 1s**, e portanto
  - **no mínimo, 99.8% de 0s**

**Demasiado esparsa!**  
Então porque não  
representar apenas os 1s?

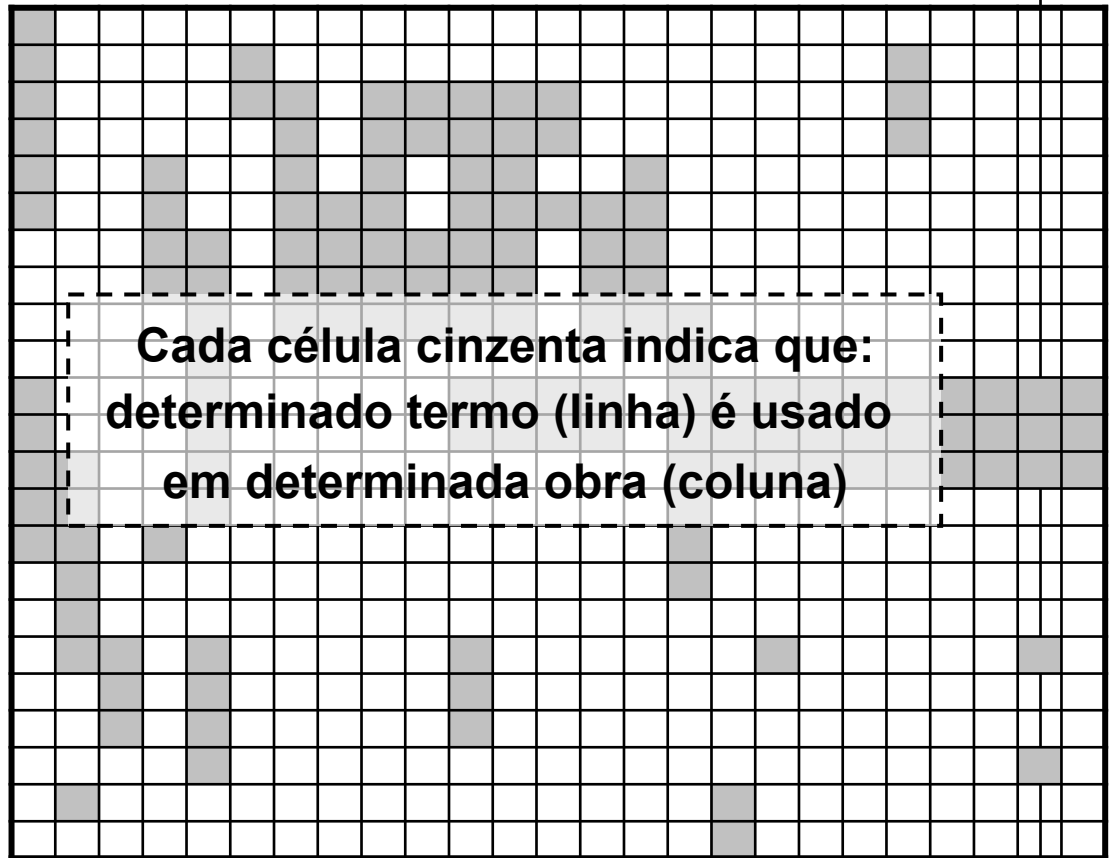
... matriz de incidência muito esparsa (a intuição)

Em geral, cada documento só contém um pequeno subconjunto do conjunto total dos termos!  
Que obra literária incluirá todo o léxico de uma Língua?

Algumas obras literárias;  
e.g., “Os Lusíadas”, “Os Maias”, etc

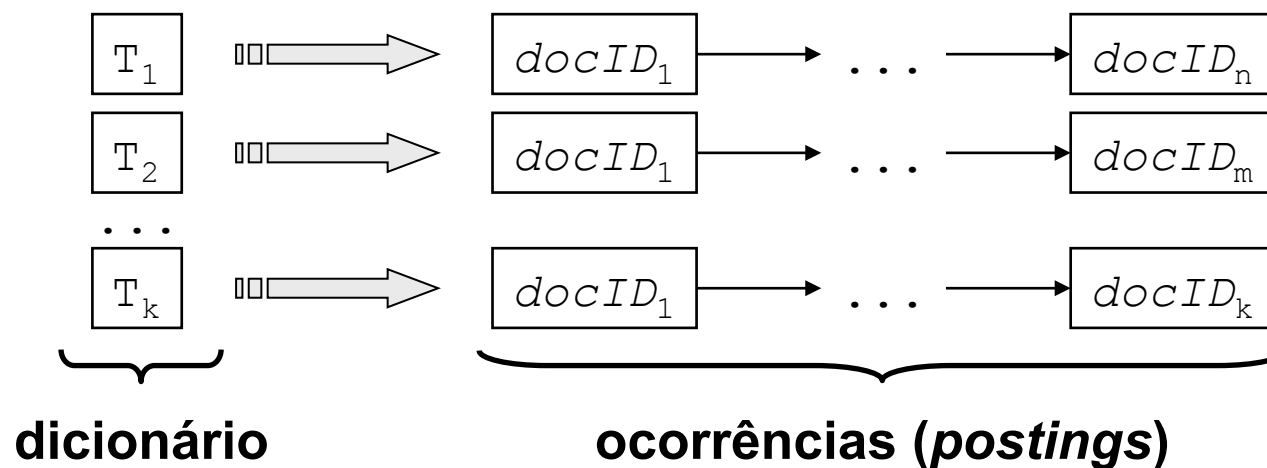
Lista de termos  
e.g., todo o léxico do Português  
i.e. “o nosso dicionário”

**Cada célula cinzenta indica que:  
determinado termo (linha) é usado  
em determinada obra (coluna)**



## Processo de indexação – usando índices invertidos

- Atribuir a cada documento um valor inteiro único: *docID*
  - e.g. atribuído sequencialmente a cada novo documento analisado
- Para cada termo T
  - armazenar a lista de todos os *docID* que contêm T
- ... usar listas dinâmicas em vez de *arrays* estáticos
  - pois o processo evolui por incrementos sucessivos



## Exercício – indexar a colecção de documentos

O rato roeu a rolha do  
garrafão do rei da  
Rússia; o rei ficou  
zangado!

docA.txt

O gato e o rato vivem  
no palácio do rei.

docB.txt

"Gato branco, gato  
preto" é um filme (de  
Kusturika) com  
imagens surrealistas.

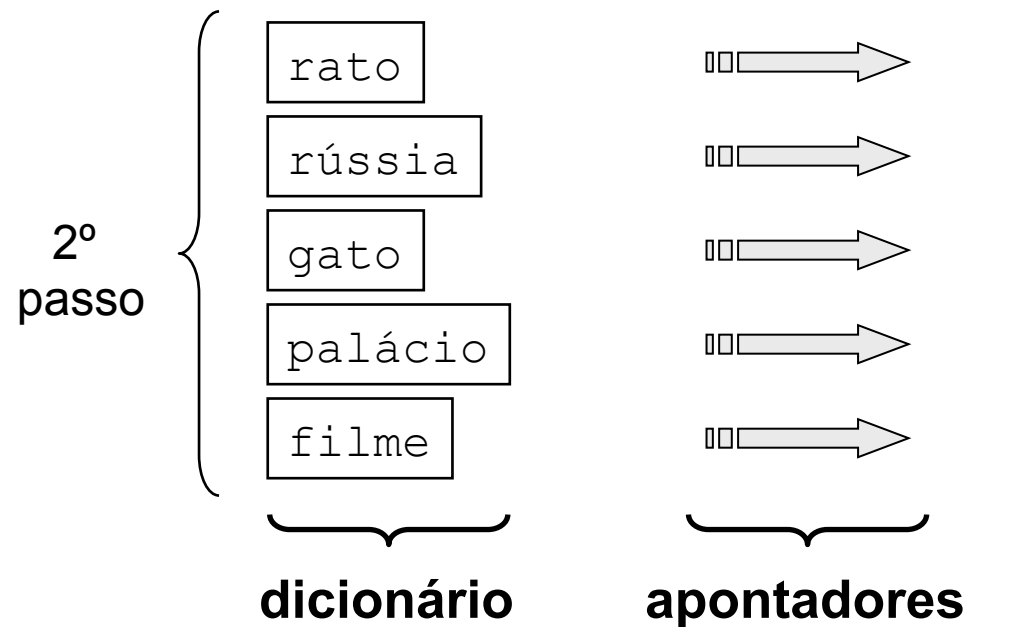
docC.txt

colecção de documentos

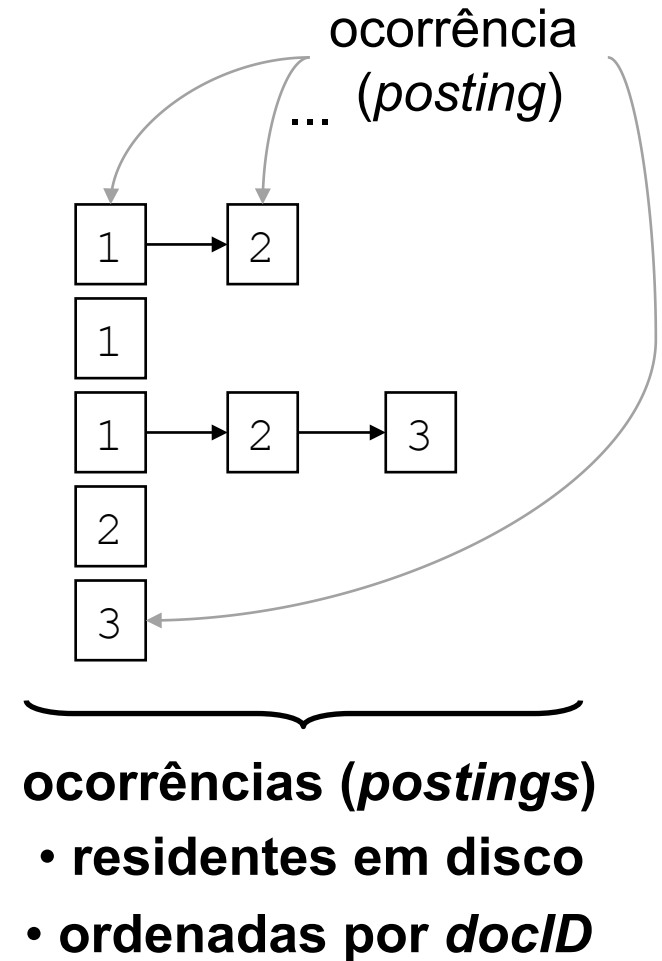
## ... exemplo – índices invertidos

O mesmo que o índice remissivo de um livro!

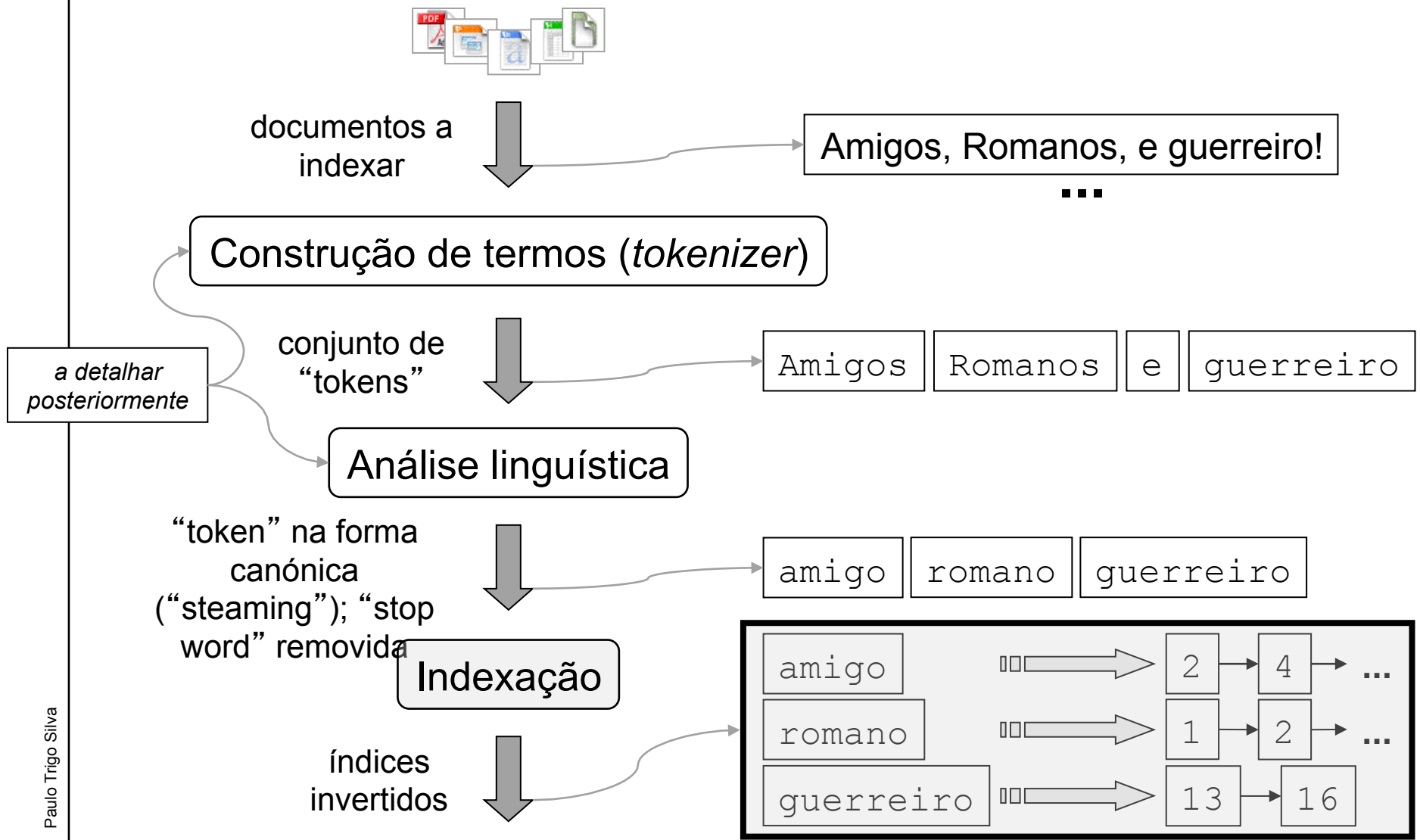
1º passo {  
docA ≡ 1  
docB ≡ 2  
docC ≡ 3



em geral mantidos em memória



# Construção de índices invertidos – processo geral





## Indexação – passo 1

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

julius.txt  
**docID = 1**

So let it be with  
Caesar. The noble  
Brutus hath told you  
Caesar was ambitious

brutus.txt  
**docID = 2**

construir  
sequência de pares  
<termo, docID>

termo	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

## Indexação – passo 2

termo	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

ordenar por  
<termo>

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

## Indexação – passo 3

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

- agregar múltiplos termos de um mesmo documento;
- registrar frequência (número de ocorrências) nesse documento

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



# Como processar uma interrogação?

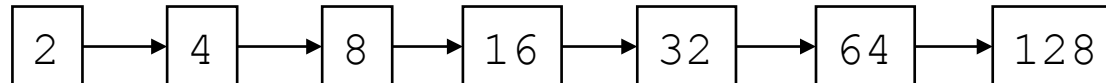


Brutus AND Caesar

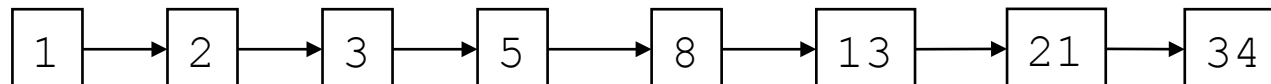
1. localizar **Brutus** no dicionário
2. obter as sua lista ocorrências
3. localizar **Caesar** no dicionário
4. obter as sua lista ocorrências
5. **intersectar** as duas listas de ocorrências

vamos  
admitir  
esta  
saída

**Brutus**



**Caesar**

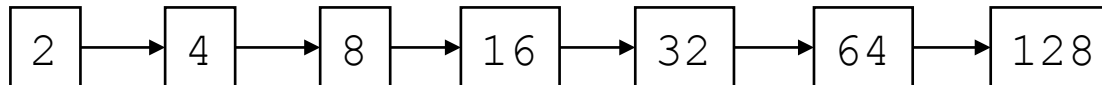


Como calcular a intersecção das duas listas?

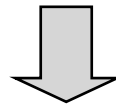
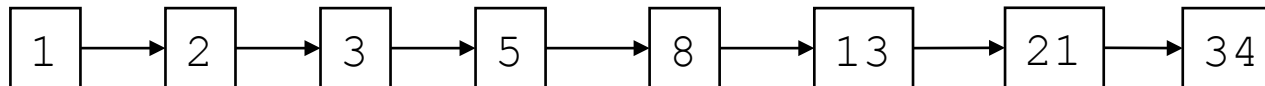
Como calcular a intersecção das duas listas?

Percorrendo as duas listas simultaneamente,  
pois estão ordenadas por `docID`

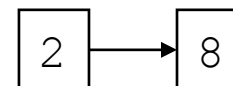
**Brutus**



**Caesar**



**Brutus AND Caesar**



Qual a relação entre o tempo para realizar a operação  
de intersecção e a dimensão das listas de ocorrências?

## Operador de conjunção (AND) – propriedades

Tempo aumenta **linearmente** com dimensão das listas.  
Pressuposto essencial: ordenação, por `docID`, das listas.

i.e. se as listas tiverem dimensão  $N$  e  $K$ , a intersecção tem  $O(N+K)$  operações.

$O(x)$   $\equiv$  tem  
ordem de  
grandeza  $x$

```
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  $answer$ 
```

pressuposto  
da ordenação

## Optimizar processamento das interrogações



Brutus AND Caesar AND Calpurnia

- Como otimizar o processamento de uma interrogação?
  - escolhendo a sequência de operações que origina o menor trabalho!
- Heurística para otimizar o processamento
  - “processar termos por ordem crescente de documentos na coleção”
- ... se começar por intersectar as duas listas de menor dimensão,
  - então os resultados intermédios serão sempre inferiores à menor lista!
  - ... e assim faremos o menor trabalho (comparações, afectações, etc)

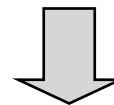
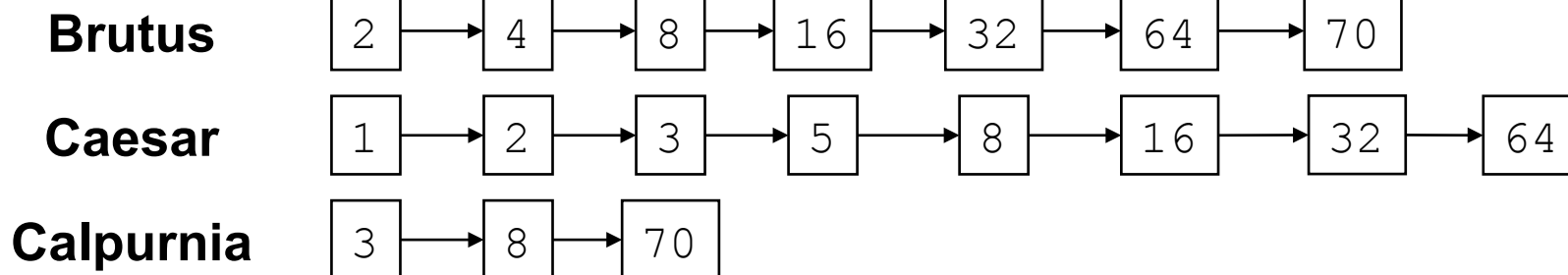
... é essencial ter, no dicionário e por termo, o número de documentos  
A escolha é feita sem acesso a disco (dicionário está em memória)



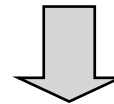
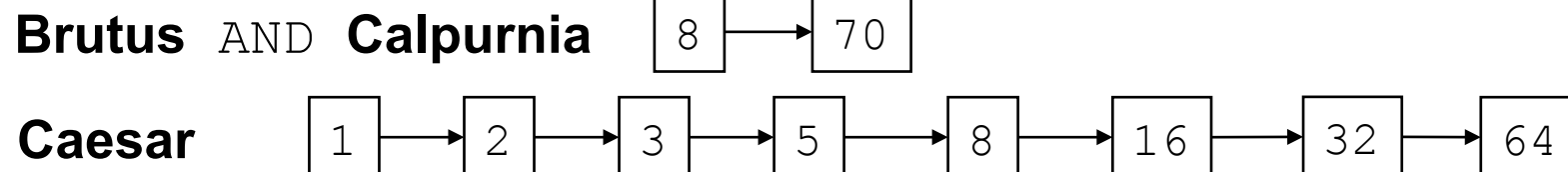
## Exemplo – otimizar processamento



Brutus AND Caesar AND Calpurnia



menores listas = Brutus; Calpurnia



( **Brutus AND Calpurnia** ) AND **Caesar** 8

## ... expressão de interrogações conjuntivas – algoritmo



$t_1$  AND  $t_2$  AND  $t_3$  AND ... AND  $t_n$

INTERSECT( $\langle t_1, \dots, t_n \rangle$ )

```
1  terms ← SORTBYINCREASINGFREQUENCY( $\langle t_1, \dots, t_n \rangle$ )
2  result ← POSTINGS(FIRST(terms))
3  terms ← REST(terms)
4  while terms  $\neq$  NIL and result  $\neq$  NIL
5  do result ← INTERSECT(result, POSTINGS(FIRST(terms)))
6     terms ← REST(terms)
7  return result
```

... ordenar

... algoritmo anterior

## Outras interrogações Booleanas

... uma curiosidade:  
Calpurnia Pisonis (nasceu por volta de 77 a.C.), filha de Lucius Calpurnius Piso Caesoninus, foi uma nobre romana, a terceira e última esposa de Júlio César,



(Brutus OR Caesar) AND (Calpurnia OR Lucius)

1. obter por termo o número de documentos em que ocorre
2. calcular soma das frequências de cada disjunção (OR)
3. processar por ordem crescente da dimensão de cada OR

perspectiva “conservadora” da  
estimativa da dimensão resultante

## Operador de disjunção (OR)

... substituir, no algoritmo da conjunção (AND), por

~~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 ADD( $\text{answer}, \text{docID}(p_1)$ )
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  $\text{answer}$ 
```

**Union** ( $p_1, p_2$ )

```
...
Add( $\text{answer}, \text{docID}(p_2)$ )
do if  $\text{docID}(p_1) \neq \text{docID}(p_2)$ 
...
```

o resto do algoritmo não se altera

Importante:

admite-se que Add só insere caso o último elemento inserido seja diferente do que se pretende inserir.

## Operador de subtracção (composição de AND e NOT)



Brutus AND NOT Caesar

... substituir, no algoritmo da conjunção (AND), por

~~INTERSECT( $p_1, p_2$ )~~

```
1 answer ← {}  
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 ADD(answer, 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
```

**Minus** ( $p_1, p_2$ )

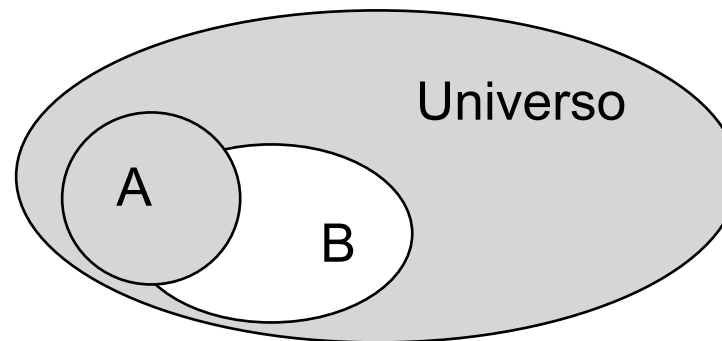
```
...  
do if  $\text{docID}(p_1) \neq \text{docID}(p_2)$   
...
```

o resto do algoritmo não se altera

## Operador para composição de OR e NOT



Brutus OR NOT Caesar



$$\begin{aligned} A \text{ OR NOT } B \\ = \\ \text{Universo} - (B - A) \end{aligned}$$

resolver invocando as funções já definidas  
ou  
construir novo algoritmo para otimizar esta operação...

## Exercício (ordem de processamento)



(tangerina OR árvore) AND  
(goiabada OR nuvem) AND  
(caleidoscópio OR olho)

termo	frequência na coleção
olho	213312
caleidoscópio	97009
goiabada	107913
nuvem	271658
tangerina	46653
árvore	316812

Recomende uma ordem de  
processamento (de operadores)  
para executar a interrogação

## Exercício (criticar motor de busca)

- Utilize o motor de busca sobre textos de Shakespeare, em
  - <http://www.rhymezone.com/shakespeare/>
- Experimente a pesquisa por “keyword” para encontrar:
  - Brutus AND Caesar
  - Brutus AND Caesar AND NOT Calpurnia
  - Brutus OR Caesar
  - Brutus OR NOT Caesar
  - ... aceda a alguns dos documentos encontrados.



Proponha 5 aspectos a melhorar  
neste motor de pesquisa.



... pesquisar textos de Shakespeare

RhymeZone Shakespeare Search: Brutus Caesar - Windows Internet Explorer

http://www.rhymezone.com/r/ss.cgi?q=Brutus+Caesar&mode=kw

Google

IPLNet: ... DEETC Rhym... Quick Re...

**RhymeZone**  
Shakespeare Search

Browse: [Comedies](#), [Tragedies](#), [Histories](#), [Poetry](#), [Help](#), [Coined words](#), [Most popular lines](#)

Find word or phrase: Brutus Caesar Search

☐ Word or phrase ☒ Keywords ☐ Start a line

Limit to: All, [Comedies](#), [Tragedies](#), [Histories](#), [Poetry](#)

Keyword search results:

Why *brutus* rose against *caesar*, this is my answer: ➔ [Julius Caesar: III, ii](#)  
So let it be with *caesar*. the noble *brutus* ➔ [Julius Caesar: III, ii](#)  
For *brutus*, as you know, was *caesar*'s angel: ➔ [Julius Caesar: III, ii](#)  
Except immortal *caesar*, speaking of *brutus* ➔ [Julius Caesar: I, ii](#)  
*caesar* than you shall do to *brutus*. the question of ➔ [Julius Caesar: III, ii](#)  
*brutus* will start a spirit as soon as *caesar*. ➔ [Julius Caesar: I, ii](#)  
Soothsayer. flourish. enter *caesar*, *brutus*, ➔ [Julius Caesar: III, i](#)  
*caesar*'s, to him I say, that *brutus*' love to *caesar* ➔ [Julius Caesar: III, ii](#)

8 results returned.

Internet 100%

Brutus AND Caesar

... um dos textos encontrados

Except immortal *caesar*, speaking of *brutus* ➔ [Julius Caesar: I, ii](#)

Forgets the shows of love to other men.

CASSIUS: Then, Brutus, I have much mistook your passion;  
By means whereof this breast of mine hath buried  
Thoughts of great value, worthy cogitations.  
Tell me, good Brutus, can you see your face?

BRUTUS: No, Cassius; for the eye sees not itself,  
But by reflection, by some other things.

CASSIUS: 'Tis just:  
And it is very much lamented, Brutus,  
That you have no such mirrors as will turn  
Your hidden worthiness into your eye,  
That you might see your shadow. I have heard,  
Where many of the best respect in Rome,  
Except immortal Caesar, speaking of Brutus  
And groaning underneath this age's yoke,  
Have wish'd that noble Brutus had his eyes.

BRUTUS: Into what dangers would you lead me, Cassius.

## Outras pesquisas (para além dos termos)

### Pesquisar:

o rato roeu a rolha

Instituto Superior de Engenharia de Lisboa

frases

Torvalds *PERTO-DE* Linux

Gates *PERTO-DE* Microsoft

proximidade  
(capturar posição  
no documento)

autor = Ullman AND  
( texto *CONTÉM* autómatos )

zonas do  
documento



## Como apresentar o resultado da pesquisa?

O modelo Booleano responde a uma interrogação com o conjunto dos documentos que satisfazem a condição de pesquisa.

- Mas é preciso oferecer primeiro aquilo que é mais relevante!
  - ... implica impor relação de ordem ao conjunto obtido (passar a lista)
  - ... implica agrupar documentos que cobrem vários aspectos da interrogação
- Para definir uma relação de ordem
  - é necessário “medir a proximidade” da interrogação a cada documento
- Para agrupar documentos
  - é necessário identificar os documentos que “melhor cobrem” a interrogação

## Outros problemas – agrupamento e classificação

### O problema do agrupamento (“clustering”):

Dado um conjunto de documentos caracterizar subconjuntos (i.e., fazer grupos ou “clusters”) com conteúdo semelhante.

... algoritmo muito conhecido: **K-means**

constrói K grupos otimizando critério de partição (semelhança).

### O problema da classificação (“classification”):

Dado um conjunto de tópicos (classes) e um novo documento D (objecto), decidir a que tópico (classe) pertence esse documento.

... algoritmo muito conhecido: “**K Nearest Neighbors**” (KNN)

atribui o objecto à classe mais comum nos seus K vizinhos.

... existem diversos outros algoritmos para atacar estes problemas!

## Sistemas de RI – código fonte aberto (“open source”)

- Lucene (“The Apache Software Foundation”)
  - <http://lucene.apache.org/>
- Lemur (“CMU, University of Massachusetts”)
  - <http://www.lemurproject.org/>
- Zettair (“RMIT, Australia”)
  - <http://www.seg.rmit.edu.au/zettair/>
- MG (“RMIT & Melbourne, Australia; Waikato, New Zealand”)
  - <http://www.ncsi.iisc.ernet.in/raja/netlis/wise/mg/mainmg.html#Features>
- ... e outros!

## Exercício

Identifique sistemas de recuperação de informação.  
Faça uma análise comparativa de 3 sistemas de recuperação de informação.