

Recuperação de Informação – modelo vectorial

Ideia – cada documento é um vector

- Numa colecção de documentos existem t termos distintos
 - identificados depois de construído o dicionário
- Os termos formam um espaço de vectores
 - dimensão do espaço = t = | dicionário |
 - ... com 2 termos é bidimensional, ... com n termos n -dimensional
- A cada termo i do documento j é atribuído um peso
 - que é um valor real ω_{ij}
- Assim, um documento d_j representa-se pelo vector
 - $d_j \equiv \langle \omega_{1j}, \omega_{2j}, \dots, \omega_{tj} \rangle$
 - ... que tem t dimensões

... cada interrogação também é um vector

- Uma interrogação j pode ser vista como um pequeno documento!
 - ou seja, como um vector de pesos...
 - $\langle \omega_{1j}, \omega_{2j}, \dots, \omega_{tj} \rangle$ (com um peso por cada uma das t dimensões)
- ... o que permite reformular a noção de interrogação
 - já não se pergunta “quais os documentos que contêm estes termos?”
- ... então, “dado um documento que outros lhe são semelhantes?”
 - passa a ser a pergunta para a qual se pretende resposta
- O modelo vectorial foi desenvolvido no sistema SMART
 - por Salton, 1970; depois explorado na recuperação de informação Web

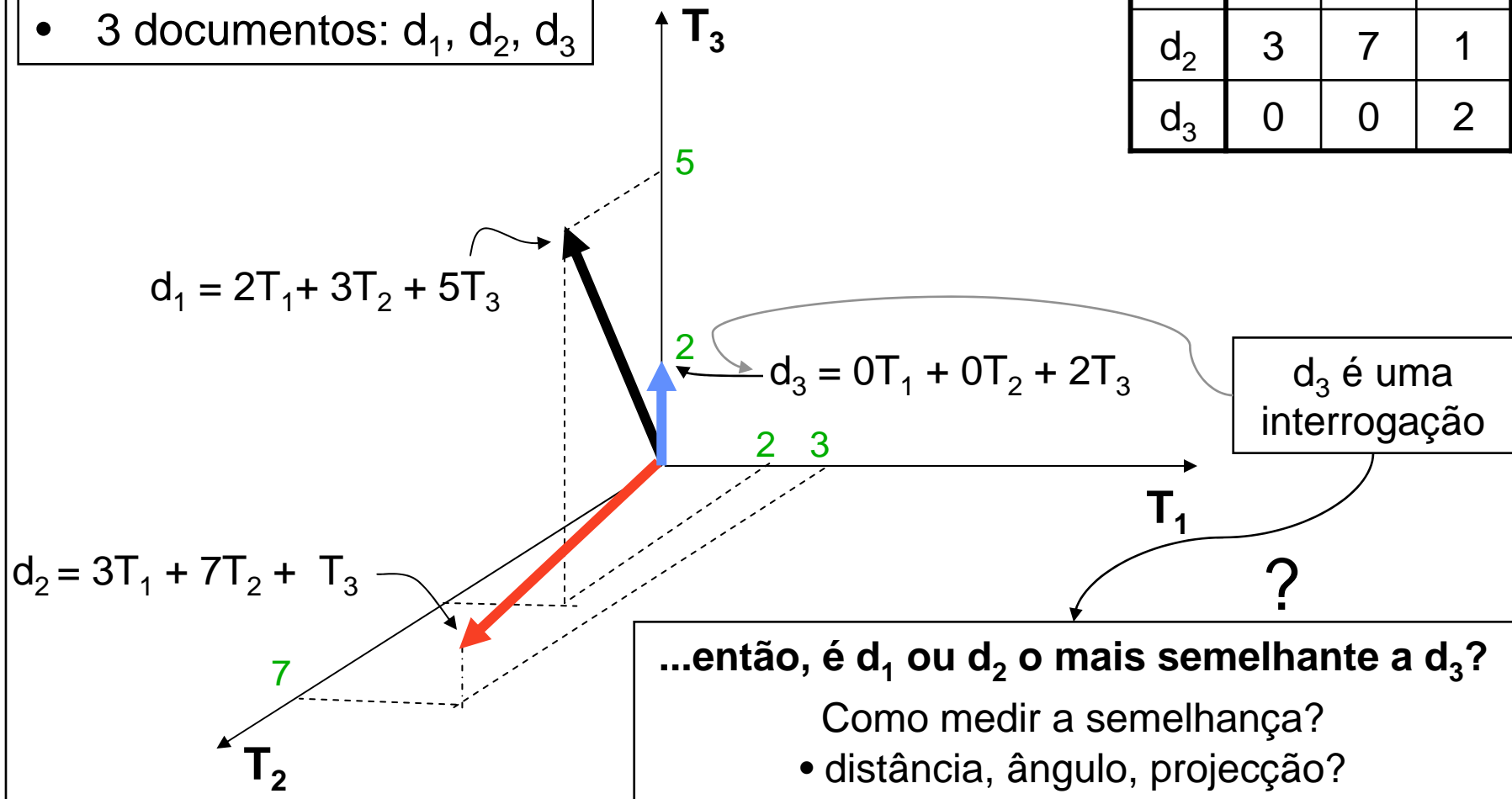
Representação gráfica

Exemplo

- 3 termos: T_1 , T_2 , T_3
- 3 documentos: d_1 , d_2 , d_3

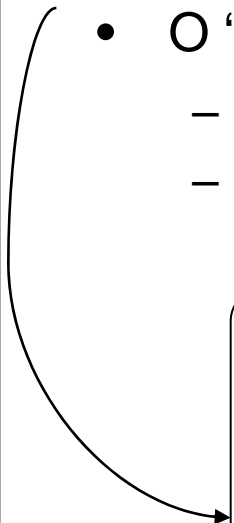
matriz de pesos

ω_{ij}	T_1	T_2	T_3
d_1	2	3	5
d_2	3	7	1
d_3	0	0	2



Representação da colecção de documentos

- A colecção de documentos representa-se no espaço de vectores
 - como uma matriz de “termos \times documentos”
- O “peso do termo no documento”
 - é representado por cada elemento da matriz
 - ... zero indica termo sem significado ou inexistente (no documento)



	T_1	T_2	\dots	T_t
d_1	ω_{11}	ω_{21}	\dots	ω_{t1}
d_2	ω_{12}	ω_{22}	\dots	ω_{t2}
\vdots	\vdots	\vdots		\vdots
\vdots	\vdots	\vdots		\vdots
d_n	ω_{1n}	ω_{2n}	\dots	ω_{tn}

Recordar:
matriz do modelo Booleano,
onde $\omega_{ij} \in \{0, 1\}$

Peso de cada termo – perspectiva da frequência

Atenção:

A literatura da recuperação de informação (RI, ou IR) usa o conceito de “frequência” para significar “quantidade”.

... i.e. não se divide pelo número total de termos no documento (o que tornaria a “quantidade” numa “frequência”)

Assim adoptaremos a ideia de que,
“frequência do termo no documento”

≡

“número de ocorrências do termo no documento”

Peso de cada termo – simplesmente a sua frequência?

- A forma mais simples de definir o peso de cada termo
 - é considerar que corresponde simplesmente à frequência
- Ou seja, o peso, $\omega_{t,d}$, do termo t no documento d , seria dado por
 - $\omega_{t,d} = f_{t,d}$ (onde $f_{t,d}$ é a frequência do termo t no documento d)
- É importante manter uma perspectiva global do peso de cada termo
 - e para isso é preciso normalizar a frequência
- ... normalizar a frequência do termo corresponde a considerar
 - $\omega_{t,d} = \text{tf}_{t,d} = f_{t,d} / \max_{d \in C} \{ f_{t,d} \}$ (onde C é a colecção de documentos)
- ... ou seja, divide-se a frequência do termo no documento pelo valor máximo da frequência daquele termo na colecção de documentos
 - obtendo um valor entre 0 e 1 que é comparável com os restantes

Exemplo – frequência normalizada do termo

Que filme! Um filme sobre como realizar um filme acerca de um gato!

docA.txt

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

docB.txt

$C \equiv$ colecção de documentos

$$tf_{t,d} = f_{t,d} / \max_{d \in C} \{ f_{t,d} \}$$

termo t	$f_{t,docA}$	$f_{t,docB}$
gato	1	2
filme	3	1

**Qual é a matriz de pesos $tf_{t,d}$
i.e., qual a matriz com as
frequências normalizadas?**

Exemplo – frequência normalizada do termo (cont.)

Que filme! Um filme sobre como realizar um filme acerca de um gato!

docA.txt

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

docB.txt

$C \equiv$ colecção de documentos

frequência máxima do termo na colecção

termo t	$f_{t,\text{docA}}$	$f_{t,\text{docB}}$	$\max_{d \in C} \{f_{t,d}\}$
gato	1	2	2
filme	3	1	3

$$tf_{t,d} = f_{t,d} / \max_{d \in C} \{f_{t,d}\}$$

termo t	$tf_{t,\text{docA}}$	$tf_{t,\text{docB}}$
gato	0.5	1
filme	1	0.3

... frequência do termo – limitação e ideia para a “atenuar”

- A principal limitação da perspectiva da frequência do termo é que
 - todos os termos se consideram igualmente importantes
 - ... na avaliação da relevância de um documento face a uma interrogação
- Por exemplo, numa colecção sobre “apólices de seguros”
 - é natural que o termo “seguro” ocorra em todos os documentos
- ... é preciso “atenuar” o efeito dos termos que “surgem demasiado”
 - na avaliação da relevância de um documento

Ideia:

- factor para reduzir $tf_{t,d}$ (frequência do termo t no documento d), e
- redução aumenta quando aumenta o número de termos t na colecção.

... “atenuar” efeito dos termos que “surgem demasiado”

Ideia:

- factor para reduzir $tf_{t,d}$ (frequência do termo t no documento d), e
- redução aumenta quando aumenta o número de termos t na colecção.

Para construir esse factor de redução há duas medidas possíveis:

- $cf_t \equiv$ número de ocorrências do termo t na colecção, e
- $df_t \equiv$ número de documentos que têm o termo t na colecção.

Aquelas medidas são conhecidas como:

- $cf_t \equiv$ frequência na colecção (*collection frequency*), e
- $df_t \equiv$ frequência do documento (*document frequency*).

Reduzir frequência do termo

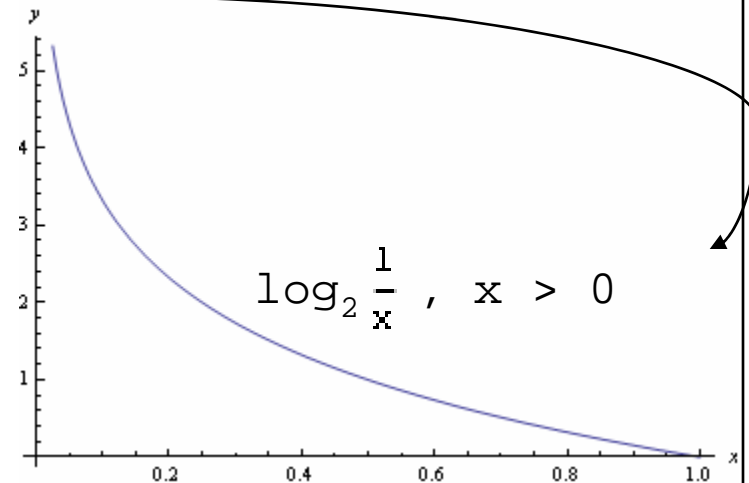
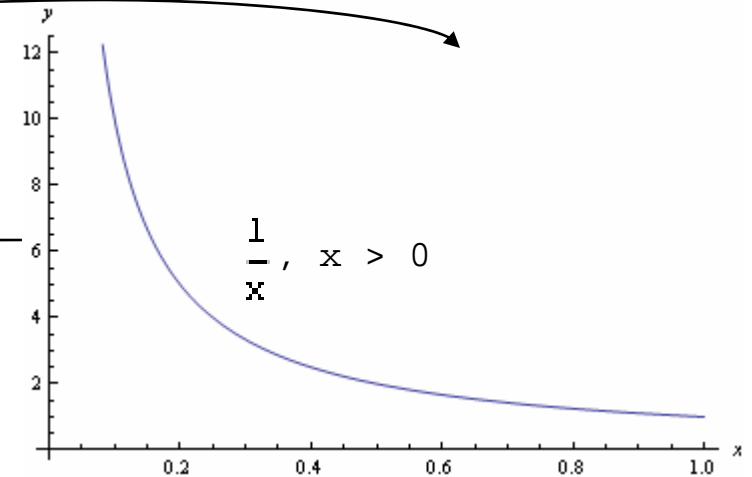
- Pode reduzir-se $tf_{t,d}$ multiplicando-o pelo inverso de uma das medidas

- $tf_{t,d} \times 1 / cf_t$, ou
- $tf_{t,d} \times 1 / df_t$

... mas a função $1/x$ é demasiado abrupta!

O seu logaritmo (de $1/x$) permite:

- decrescer de modo “menos abrupto”, e
- “melhor comportamento” perto de zero



... mas que medida, cf ou df , usar?

Exemplo ilustrativo

termo t	cf_t	df_t
ferrari	10442	23
seguro	10440	3997

- A frequência da colecção (cf) e a frequência do documento (df)
 - podem ter comportamento bastante diferente (um do outro)
- ... `ferrari` e `seguro` têm valor de **cf** idêntico, mas
 - o valor de **df** é bastante diferente entre eles!
- Intuitivamente, pretendemos que,
 - os poucos documentos com `ferrari` tenham maior peso (numa interrogação sobre `ferrari`) do que os que contém `seguro`.
- ... assim, é usual usar-se **df** como medida para reduzir $tf_{t,d}$.

Função “inversa da frequência do documento”

Como usar a frequência do documento (**df**) para reduzir o peso de um termo?

Dado uma colecção com **N** documentos, considera-se:

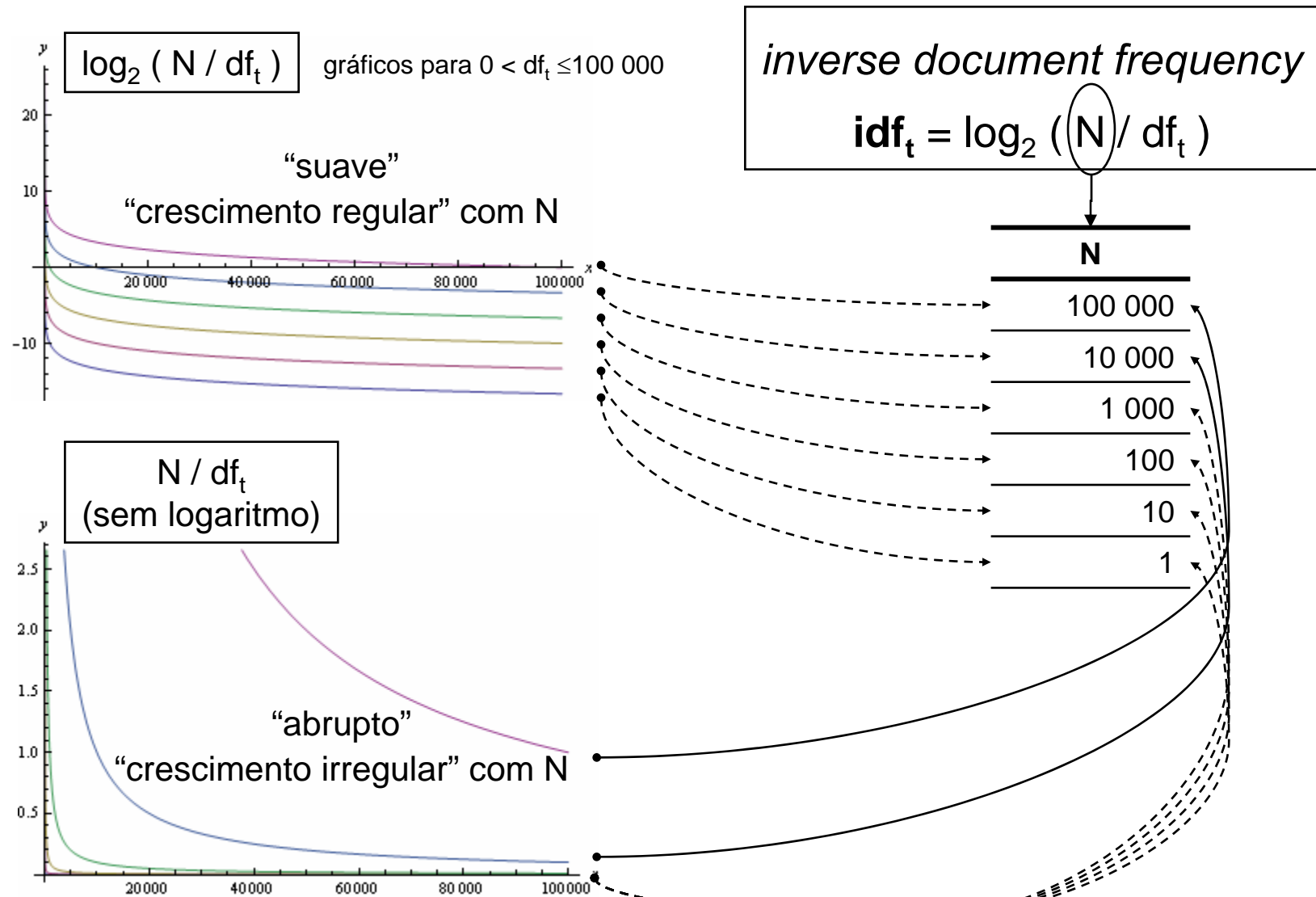
$$\mathbf{idf}_t = \log_2 (N / df_t), \text{ para } 0 < df_t \leq N$$

Recordar:
 $\log (1) = 0$

idf_t \equiv inversa da frequência do documento
(*inverse document frequency*)

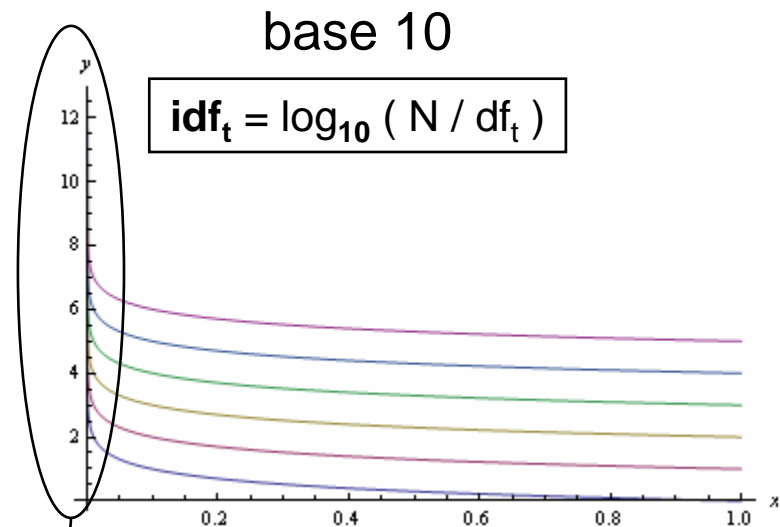
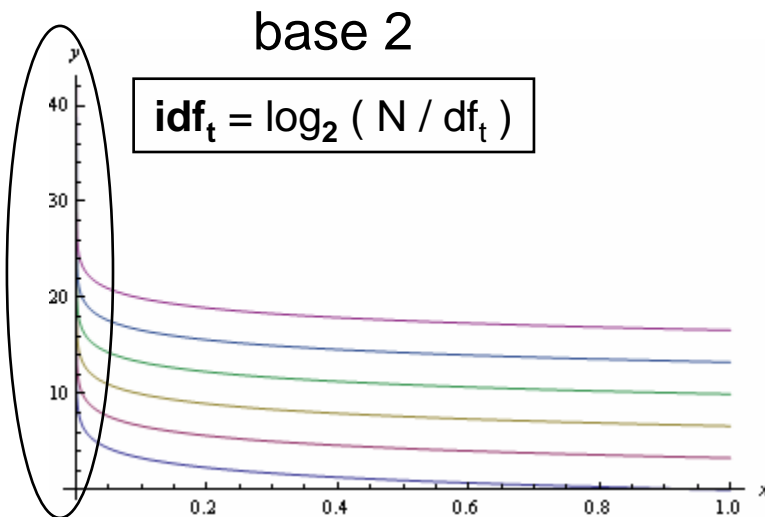
O logaritmo torna a função inversa menos abrupta e com “crescimento regular” face ao número de documentos.

... o efeito da dimensão da colecção (e papel do logaritmo)



... que base usar?

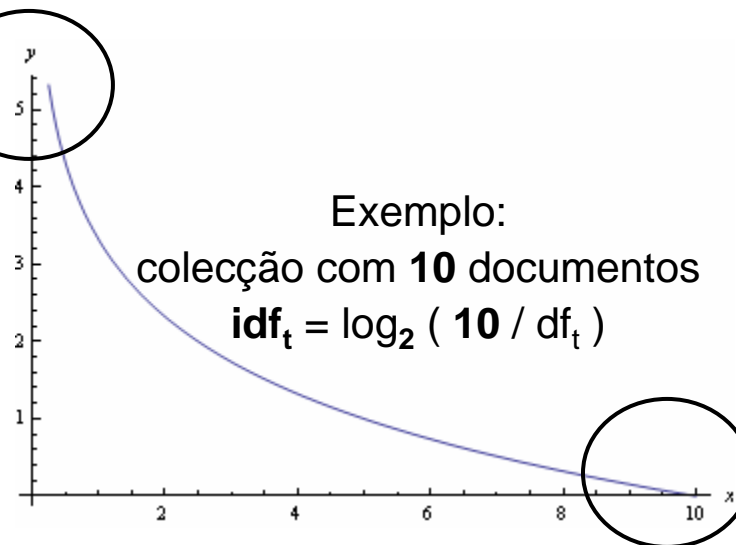
O logaritmo pode usar qualquer outra base, para além da base 2 (e.g. 10), que não afecta o comportamento global da função.



... aumentar a base reduz os valores de **idf** .

O que oferece o “idf” (*inverse document frequency*)?

O idf_t dá ideia da capacidade discriminatória do termo t .
É uma medida da raridade (“do quanto de raro”) do termo na colecção.



Exemplo:
colecção com **10** documentos
 $idf_t = \log_2 (10 / df_t)$

notar que o
máximo valor
de df_t é 10

O que é um termo, t , raro?

- é um que ocorre em pouco documentos, i.e. cujo valor de df_t é baixo!

→ O idf_t de um termo, t , raro é alto.

O idf_t de um termo, t , frequente é baixo.

“atenuar” efeito dos termos que “surgem demasiado”

... recordar a ideia inicial:

- factor para reduzir $tf_{t,d}$ (frequência do termo t no documento d), e
- redução aumenta quando aumenta o número de termos t na colecção.

$$tf-idf_{t,d} = tf_{t,d} \times idf_t$$

$$f_{t,d} / \max_{d \in C} \{f_{t,d}\}$$

$$\log_2 (N / df_t)$$

“frequência do termo t no documento d ”

≡

“número de ocorrências do termo t no documento d ”

número de documentos que têm o termo t na colecção

Análise teórica em: Kishore Papineni, NAACL 2, 2002.

... em síntese – o peso “tf-idf”

- Atribuir um peso, $\omega_{t,d}$, tf-idf a cada termo, t , em cada documento d

$$\omega_{t,d} = \text{tf}_{t,d} \times \log_2 (N / \text{df}_t)$$

- $\text{tf}_{t,d} \equiv$ número de ocorrências do termo t no documento d
 - ◊ ...este valor pode estar normalizado
- $N \equiv$ número total de documentos na colecção
- $\text{df}_t \equiv$ número de documentos que, na colecção, têm o termo t
- Aumenta com o número de ocorrências do termo num documento
- Aumenta com a “raridade” do termo na colecção de documentos

Qual o peso de um termo que ocorre em todos os documentos?

$$\text{df}_t = N \Rightarrow \omega_{t,d} = 0$$

Em quantos documentos deve um termo ocorrer para que o seu peso tenha valor máximo?

$$\text{df}_t = 1 \text{ (domínio dos naturais)} \Rightarrow \omega_{t,d} \text{ máximo}$$

Exemplo – peso “tf-idf”

Dado uma colecção com 23456 documentos:

- sejam os termos A, B e C, e
- sejam 2 documentos d1 e d2.

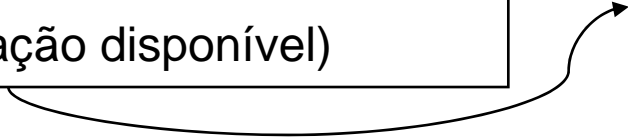
Considere-se que os termos A, B e C ocorrem:

- no documento d1, respectivamente 5, 5 e 2 vezes, e
- no documento d2, respectivamente 15, 1 e 6 vezes.

Considere-se que os termos A, B e C ocorrem na colecção:

- em, respectivamente 2000, 5 e 2 documentos.

Construa a matriz de “termos \times documentos”
(para a informação disponível)



	T_1	T_2	\dots	T_t
d_1	ω_{11}	ω_{21}	\dots	ω_{t1}
d_2	ω_{12}	ω_{22}	\dots	ω_{t2}
\vdots	\vdots	\vdots		\vdots
\vdots	\vdots	\vdots		\vdots
d_n	ω_{1n}	ω_{2n}	\dots	ω_{tn}

... exemplo – peso “tf-idf”

Cálculos

$N = 23456$

tf	A	B	C
d1	5	5	2
d2	15	1	6

df	2000	5	2
----	------	---	---

$\log_2 (N/df)$	3.55	12.20	13.52
-----------------	------	-------	-------

tf x $\log_2 (N/df)$	A	B	C
d1	17.76	60.98	27.04
d2	53.28	12.20	81.11

$\omega_{t,d}$	A	B	C
d1	17.76	60.98	27.04
d2	53.28	12.20	81.11
...
d23456

$$\omega_{t,d} = \text{tf}_{t,d} \times \log_2 (N / \text{df}_t)$$

... não depende de d, pelo que $\omega_{t,d}$ mantém a proporção dos respectivos $\text{tf}_{t,d}$.

Por exemplo:

$$\text{tf}_{C,d2} / \text{tf}_{C,d1} = 6 / 2 = 3$$

$$\omega_{C,d2} / \omega_{C,d1} = 81.11 / 27.04 = 3$$

Mas, para um documento, d, altera-se muito a relação entre os pesos, $\omega_{t,d}$, dos seus termos t.

Por exemplo:

$$\text{tf}_{B,d1} / \text{tf}_{A,d1} = 5 / 5 = 1$$

$$\omega_{B,d1} / \omega_{A,d1} = 60.98 / 17.76 = 3.43$$

A noção de proximidade

- Temos técnica para construir uma matriz de pesos
 - calcula-se $\text{tf-idf}_{t,d}$ e fica-se com um espaço de vectores
 - ... representados como uma matriz “termos \times documentos”
- Mas, como usar esse espaço para escolher documentos?
 - cada documento é um vector de pesos, portanto
 - ... como comparar um vector (documento) com outro vector?
- Como decidir quanto à proximidade entre vectores (documentos)?
 - baseada na distância de Manhattan?
 - baseada na distância euclidiana?
 - baseada no produto interno entre vectores?
 - baseada no valor do co-seno do ângulo entre vectores?
 - ... as 2 anteriores podem ver-se como uma única (o co-seno pode ver-se como produto interno normalizado)

$$\begin{pmatrix} & T_1 & T_2 & \dots & T_t \\ d_1 & \omega_{11} & \omega_{21} & \dots & \omega_{t1} \\ d_2 & \omega_{12} & \omega_{22} & \dots & \omega_{t2} \\ \vdots & \vdots & \vdots & & \vdots \\ d_n & \omega_{1n} & \omega_{2n} & \dots & \omega_{tn} \end{pmatrix}$$

Abordagem 1 – distância de Manhattan

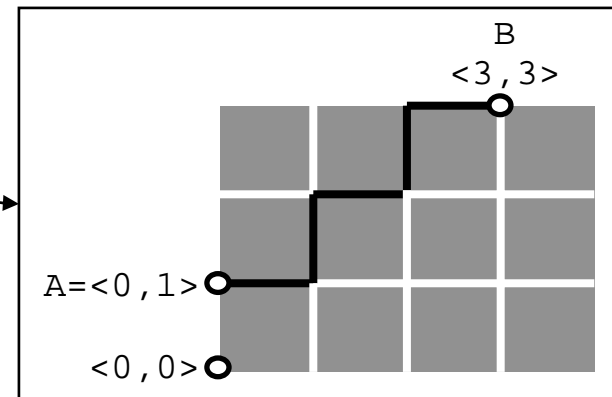
Distância de Manhattan, ou distância de blocos.

Inspira-se na ideia de que as cidades Americanas têm um formato em grelha.

Distância de Manhattan entre

$A = \langle 0, 1 \rangle$ e $B = \langle 3, 3 \rangle$

$$3 - 0 + 3 - 1 = 3 + 2 = 5$$



Distância de Manhattan

$$ManhDist(X, Y) = \sum_{i=1}^n |x_i - y_i|$$

Calcular distância de Manhattan entre

$A = \langle 0, 3, 2, 1, 10 \rangle$ e $B = \langle 2, 7, 1, 0, 0 \rangle$

$$|0-2| + |3-7| + |2-1| + |1-0| + |10-0| = 18$$

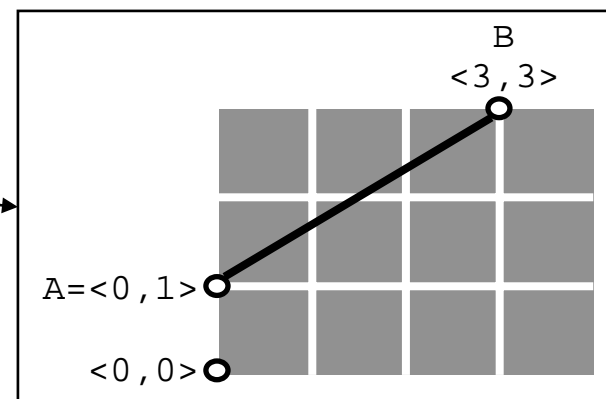
Abordagem 2 – distância euclidiana

Distância em linha recta entre dois pontos.

Distância euclidiana entre

$A = \langle 0, 1 \rangle$ e $B = \langle 3, 3 \rangle$

$$[(3 - 0)^2 + (3 - 1)^2]^{1/2} = [9 + 4]^{1/2} = 3.6$$



Distância euclidiana entre os documentos (vectores) d_j e d_k

$$|d_j - d_k| = \sqrt{\sum_{i=1}^n (d_{i,j} - d_{i,k})^2}$$

Calcular distância euclidiana entre

$d_1 = \langle a, b, c \rangle$ e $d_2 = \langle x, y, z \rangle$

$$\sqrt{\text{Abs}[a - x]^2 + \text{Abs}[b - y]^2 + \text{Abs}[c - z]^2}$$

... limitações – abordagem 1 e 2

- As métricas sobre distância sofrem grande influência da
 - da dimensão do documento
- Documentos pequenos tendem a ser próximos não pelo conteúdo
 - mas pela sua dimensão
- Para comparar distâncias é preciso normalizar
 - dividindo cada componente pelo módulo (tamanho) do vector
 - ◇ ... i.e. para um documento d e dados k pesos, $\omega_{1,d}, \dots, \omega_{k,d}$, fazer
 - ◇ $\omega_{i,d} / (\omega_{1,d}^2 + \dots + \omega_{k,d}^2)^{1/2}$, para cada $1 \leq i \leq k$
 - ao normalizar fica-se, para cada documento d_j , com
 - ◇ $|d_j| = (\omega_{1,d}^2 + \dots + \omega_{k,d}^2)^{1/2} = 1$ (vector projectado em esfera de raio 1)
- ... mas, olhando apenas para os ângulos (em vez das distâncias)
 - é também possível obter uma perspectiva normalizada

$$|\vec{d}_j| = \sqrt{\sum_{i=1}^n w_{i,j}^2}$$

Abordagem 3 – produto interno

O produto interno (ou escalar) entre dois vectores dá uma medida do “peso” (ou força) daqueles vectores quando projectados numa mesma direcção.

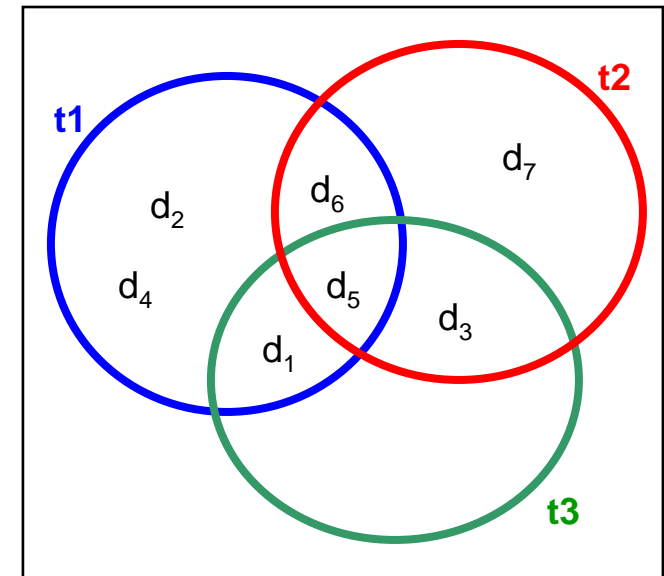
- A semelhança (similaridade) entre documentos pode medir-se pelo
 - produto interno dos seus vectores
- ... dado um documento d_j e interrogação q (também é documento)

$$\text{sim}(d_j, q) = d_j \bullet q = \sum_{i=1}^t \omega_{ij} \times \omega_{iq}$$

- onde ω_{iq} é o peso do termo i na interrogação q
- Para vectores binários dá o número de termos comuns em d_j e q
 - corresponde à dimensão da intersecção dos conjuntos termos
- Para vectores com pesos reais dá a soma dos produtos dos pesos
 - dos termos que existem em ambos os documentos

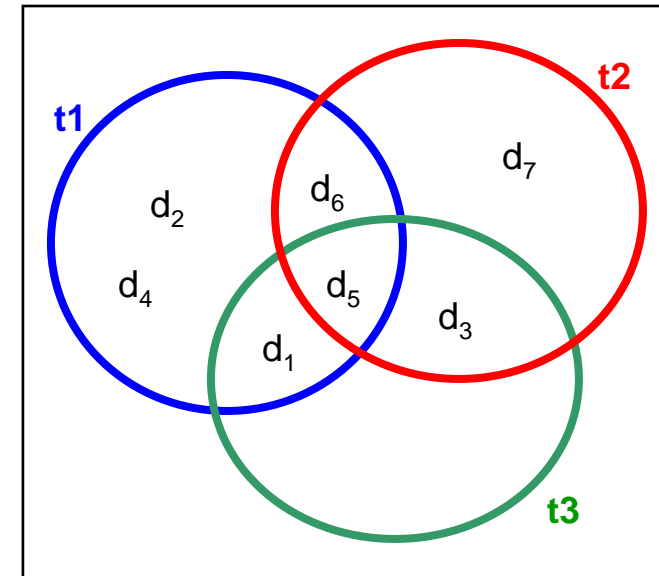
... exemplo – produto interno (vectores binários)

	t1	t2	t3	$q \bullet d_j$
d_1	1	0	1	2
d_2	1	0	0	1
d_3	0	1	1	2
d_4	1	0	0	1
d_5	1	1	1	3
d_6	1	1	0	2
d_7	0	1	0	1
q	1	1	1	



... exemplo – produto interno (limitações)

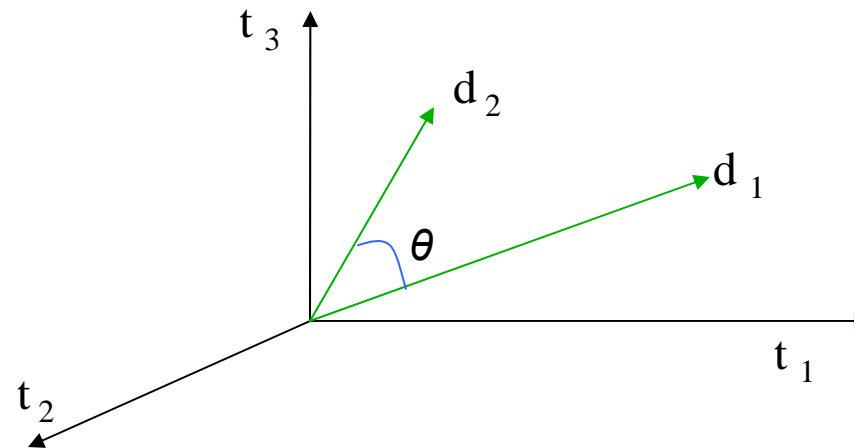
	t1	t2	t3	$q \bullet d_j$
d_1	1	0	1	4
d_2	1	0	0	1
d_3	0	1	1	5
d_4	1	0	0	1
d_5	1	1	1	6
d_6	1	1	0	3
d_7	0	1	0	2
q	1	2	3	



Favorece documentos longos com grande número de termos únicos repetidos.

Favorece o número de termos comuns nos documentos (d_j e q), mas não contempla o número de termos que não são comuns.

.. termos, documentos e ângulos



t_1 , t_2 e t_3 são termos;

d_1 e d_2 são documentos e portanto têm um peso para cada termo;

d_1 e d_2 são portanto vectores de pesos,

e.g., $d_1 = 15t_1 + 4t_2 + 0t_3 = \langle 15, 4, 0 \rangle$

e.g., $d_2 = 2t_1 + 6t_2 + 3t_3 = \langle 2, 6, 3 \rangle$

$d_1 \bullet d_1 = 15 \times 15 + 4 \times 4 + 0 \times 0 = 229$

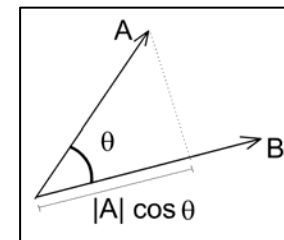
... o produto interno depende da dimensão dos vectores.

Mas o ângulo formado pelos vectores não depende de nenhuma dimensão!

... mais sobre o produto interno

- O produto interno (ou escalar) de vectores usa-se para determinar
 - a componente escalar de um vector numa determinada direcção
 - ... também designado por “ projecção ” do vector nessa direcção
- O produto interno entre os vectores A e B é dado por

$$A \bullet B = |A| |B| \cos(\theta) \quad \theta \equiv \text{ângulo entre A e B}$$



- $A \bullet B$ não é um vector; é um escalar
 - $A \bullet B = 0$, para vectores perpendiculares
 - $A \bullet B = |A| |B|$, para A e B colineares com mesmo sentido ($\cos(0) = 1$)
 - $A \bullet B = -|A| |B|$, para A e B colineares de sentido inverso ($\cos(\pi) = -1$)
 - ... ou seja, o produto interno varia em: $-|A| |B| \leq A \bullet B \leq |A| |B|$

ângulo nulo \rightarrow co-seno = 1 \rightarrow **semelhança máxima**

ângulo recto \rightarrow co-seno = 0 \rightarrow **diferença máxima**

Abordagem 4 – ângulo entre vectores

- A partir do valor do produto interno é simples calcular
 - o valor do co-seno do ângulo formado pelos vectores, e.g. A e B

$$A \bullet B = |A| |B| \cos(\theta)$$

\Leftrightarrow

$$\cos(\theta) = (A \bullet B) / (|A| |B|)$$

Recordar, módulo do vector $d_j = \langle \omega_{1j}, \dots, \omega_{nj} \rangle$

$$|\vec{d}_j| = \sqrt{\sum_{i=1}^n w_{i,j}^2}$$

- Ou seja, o co-seno pode ver-se como o produto interno normalizado
 - ... $A / |A|$ transforma A num vector de módulo 1 (e o mesmo para B)
- Em síntese, a similaridade (semelhança) entre documentos, d_j e d_k ,
 - pode calcular-se como $\text{sim}(d_j, d_k)$ sendo o valor do co-seno entre d_j e d_k ,

$$\text{sim}(d_j, d_k) = \frac{\vec{d}_j \cdot \vec{d}_k}{|\vec{d}_j| |\vec{d}_k|} = \frac{\sum_{i=1}^n w_{i,j} w_{i,k}}{\sqrt{\sum_{i=1}^n w_{i,j}^2} \sqrt{\sum_{i=1}^n w_{i,k}^2}}$$

$\omega_{t,d}$ é o peso do termo t no documento d

Exemplo – dois documentos e uma interrogação

$$\text{sim}(d_j, q) = \frac{\vec{d}_j \cdot \vec{q}}{|\vec{d}_j| \cdot |\vec{q}|} = \frac{\sum_{i=1}^t (w_{ij} \cdot w_{iq})}{\sqrt{\sum_{i=1}^t w_{ij}^2 \cdot \sum_{i=1}^t w_{iq}^2}}$$

$$d_1 = 2t_1 + 3t_2 + 5t_3$$

$$d_2 = 3t_1 + 7t_2 + t_3$$

$$q = 0t_1 + 0t_2 + 2t_3$$

$$d_1 \bullet q = 2 \times 0 + 3 \times 0 + 5 \times 2 = \mathbf{10}$$

$$d_2 \bullet q = 3 \times 0 + 0 \times 0 + 1 \times 2 = \mathbf{2}$$

$$\text{sim}(d_1, q) = 10 / [(4 + 9 + 25)(0 + 0 + 4)]^{1/2} = \mathbf{0.81}$$

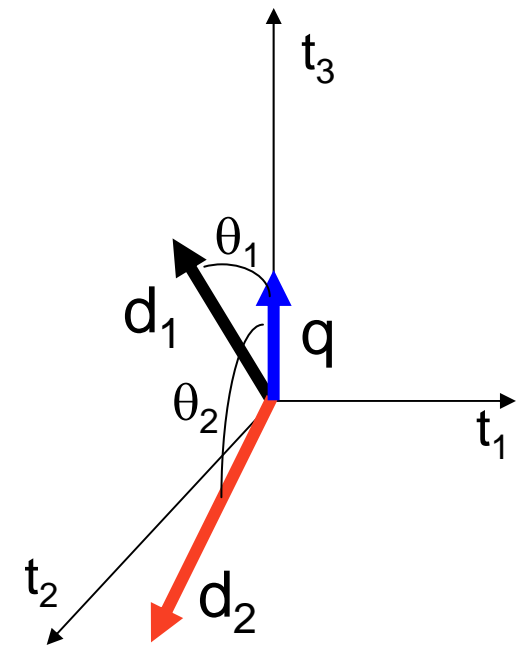
$$\text{sim}(d_2, q) = 2 / [(9 + 49 + 1)(0 + 0 + 4)]^{1/2} = \mathbf{0.13}$$

Similaridade dos documentos em relação à interrogação q:

- usando co-seno: d_1 é 6.23 vezes melhor que d_2
- usando produto interno: d_1 é apenas 5 vezes melhor que d_2

Exemplo

- 3 termos: t_1, t_2, t_3
- 3 documentos: d_1, d_2, q



... outros exemplos

Ordenar por ordem crescente de similaridade (medida pelo co-seno):

- (A) dois documentos que só têm palavras usuais em comum
 - e.g. “de”, “para”, “um”, ...
- (B) dois documentos que não têm palavras em comum
- (C) dois documentos que têm muitas palavras raras em comum,
 - e.g. “esferoidal”, “decangular”, ... (tem forma esférica e tem 10 ângulos)

$$\text{sim}(d_j, d_k) = \frac{\vec{d}_j \cdot \vec{d}_k}{\|\vec{d}_j\| \|\vec{d}_k\|} = \frac{\sum_{i=1}^n w_{i,j} w_{i,k}}{\sqrt{\sum_{i=1}^n w_{i,j}^2} \sqrt{\sum_{i=1}^n w_{i,k}^2}}$$

A relação de ordem seria:

- **(C)** numerador cresce com as raras; provavelmente também terá usuais.
- **(A)** numerador só cresce com usuais.
- **(B)** numerador é zero.

Exemplo – co-seno com vectores normalizados

Se os vectores estiverem normalizados, então:

$$\cos(\vec{d}_j, \vec{d}_k) = \vec{d}_j \cdot \vec{d}_k$$

Os pesos, ω_{ij} , são simplesmente $tf_{t,d}$ (número de ocorrências do termo t em d)

		affection	jealous	gossip	$ d_j $
Sense and Sensibility	d1	115	10	2	115.451
Pride and Prejudice	d2	58	7	0	58.421
Wuthering Heights	d3	20	11	6	23.601

		affection	jealous	gossip
		$ d_j $	$ d_j $	$ d_j $
Sense and Sensibility	d1	0.996	0.087	0.017
Pride and Prejudice	d2	0.993	0.120	0.000
Wuthering Heights	d3	0.847	0.466	0.254

Recordar, módulo do vector $d_j = \langle \omega_{1j}, \dots, \omega_{nj} \rangle$

$$|\vec{d}_j| = \sqrt{\sum_{i=1}^n w_{i,j}^2}$$

Recordar, normalizar corresponde a dividir cada componente, ω_{ij} , de d_j pelo módulo de d_j .

$$\cos(d1, d2) = .996 \times .993 + .087 \times .120 + .017 \times 0.0 = \mathbf{0.999}$$

$$\cos(d1, d3) = .996 \times .847 + .087 \times .466 + .017 \times .254 = \mathbf{0.889}$$

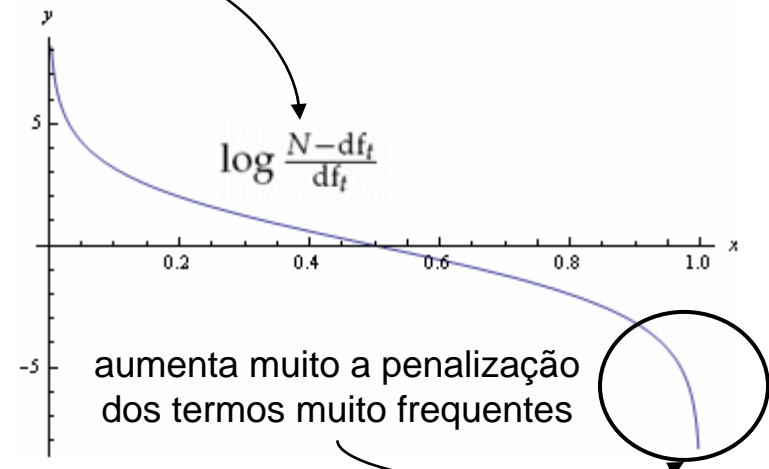
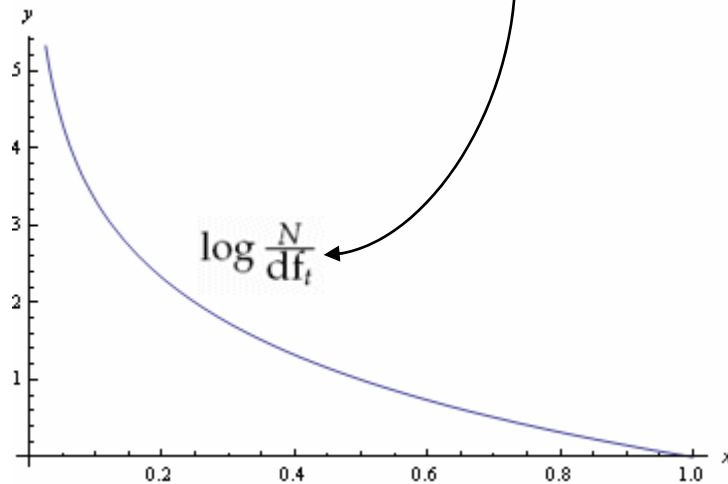
d1 e d2 ambos de Jane Austin;
d3 é de Emile Bronte;
d1 e d2 estão muito próximos!

Síntese – “td”, “df”, “tf-idf” e normalização (várias métricas)

$tf_{t,d} \equiv$ número de ocorrências do termo t em d

$df_t \equiv$ número de documentos que, na colecção, têm o termo t

Term frequency		Document frequency		Normalization	
n (natural)	$tf_{t,d}$	n (no)	1	n (none)	1
l (logarithm)	$1 + \log(tf_{t,d})$	t (idf)	$\log \frac{N}{df_t}$	c (cosine)	$\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_M^2}}$
a (augmented)	$0.5 + \frac{0.5 \times tf_{t,d}}{\max_t(tf_{t,d})}$	p (prob idf)	$\log \frac{N - df_t}{df_t}$	b (byte size)	$1 / charLeng^\alpha, \alpha < 1$
L (log ave)	$\frac{1 + \log(tf_{t,d})}{1 + \log(\text{ave}_{t \in d}(tf_{t,d}))}$				
b (boolean)	$tf_{t,d} > 0$				

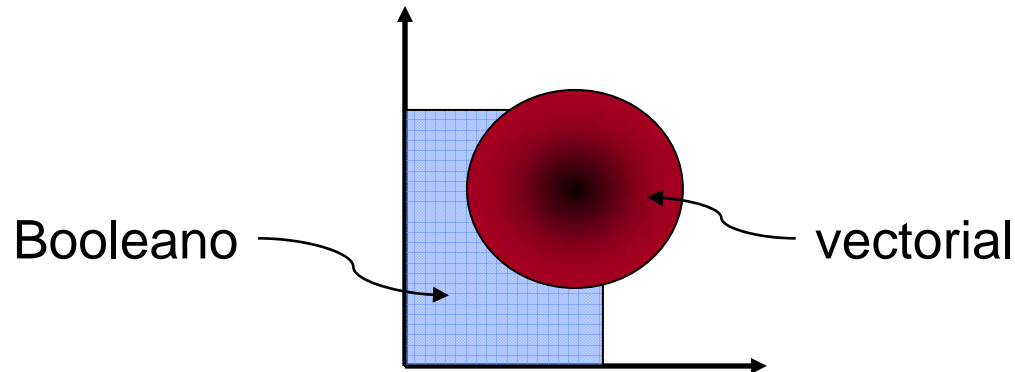


Comentários ao modelo vectorial

- Abordagem baseada em conceitos formais (matemáticos) simples
 - mas com muito adequados ao problema em causa
- Considera perspectiva local (tf) e global (idf) da ocorrência de termos
 - a combinação das perspectivas (tf-idf) reduz limitação de ambas
- Permite obter respostas parciais
 - pode desprezar documentos com proximidade quase nula
- Permite ordenar as respostas
 - de acordo com o valor de proximidade face a determinada interrogação
- ... na prática permite obter bons resultados
 - quanto à cobertura e relevância das respostas (a detalhar mais tarde)

... interrogações no modelo Booleano ‘versus’ vectorial

- Os modelos vectorial e Booleano não funcionam muito bem juntos!
 - a noção de similaridade não “combina bem” com a de presente/ausente
- No espaço dos termos a similaridade de vectores define esferas
 - e.g. “todos os documentos com co-seno ≥ 0.5 ” face à interrogação
- A interrogação Booleana devolve hiper-rectângulos
 - e suas intersecções e uniões
- ... “entalhe rectangular” versus “cavilha redonda”



Implementação do modelo vectorial – na versão ‘naïve’

Input: uma colecção C de documentos e uma interrogação q
Output: uma lista ordenada de documentos

- para cada documento $d_j \in C$
 - converter d_j num vector de pesos tf-idf
- converter a interrogação q num vector de pesos tf-idf
- para cada documento $d \in C$
 - calcular $\text{score}_j = \text{sim}(d_j, q)$
- ordenar documentos por ordem decrescente de score_j
- apresentar, ao utilizador, os documentos de topo

Complexidade temporal: $O(|T| \cdot |C|)$; mau para grandes T & C !

$|T| = 10,000$; $|C| = 100,000$; $|T| \cdot |C| = 1,000,000,000$

Implementação do modelo vectorial – na prática

- Os documentos que não têm qualquer termos da interrogação
 - não afectam o resultado da interrogação!
- Tentar identificar os documentos que contêm
 - pelo menos 1 termo da interrogação
- Usar índices invertidos para suportar a pesquisa

Implementar funções de pré-processamento

- “tokenization” (eliminar pontuações, espaços, etc)
- remoção de “stop words” (“de”, “e”, “com”, etc)
- “stemming” (passar termos para forma canónica)
 - e.g. “casa”, “casinha”, “casarão” → “casa”

Objectivo:

**Reduzir o
espaço da
pesquisa!**

Tornar eficiente o cálculo dos co-seno

- Encontrar K documentos na colecção
 - aqueles que estão mais próximos da interrogação
- ... implica calcular os K maiores co-senos entre
 - a interrogação e K documentos
- Problema:
 - calcular, ou garantir a similaridade, do co-seno de forma eficiente
 - escolher os k co-seno de forma eficiente

Objectivo:

Escolher os K maiores co-senos
sem ter que os calcular todos!

Garantir similaridade do co-seno de forma eficiente

- Numa interrogação sem termos repetidos
 - o valor normalizado dos termos é igual!
- ... na interrogação pode considerar-se
 - que os seus termos têm valor 1 (i.e. não se normaliza)
- Para quaisquer dois documentos d_1 e d_2 tem-se,

q	gato	rato	rei	$ q $
	1	1	1	1.73

gato	rato	rei
$ q $	$ q $	$ q $
0.58	0.58	0.58

$$\vec{V}(q) \cdot \vec{v}(d_1) > \vec{V}(q) \cdot \vec{v}(d_2) \Leftrightarrow \vec{v}(q) \cdot \vec{v}(d_1) > \vec{v}(q) \cdot \vec{v}(d_2).$$

vector de pesos
vector apenas com 1s e 0s

para provar basta
dividir ambos os termos
da desigualdade por $|q|$
($q / |q|$ é q normalizado)

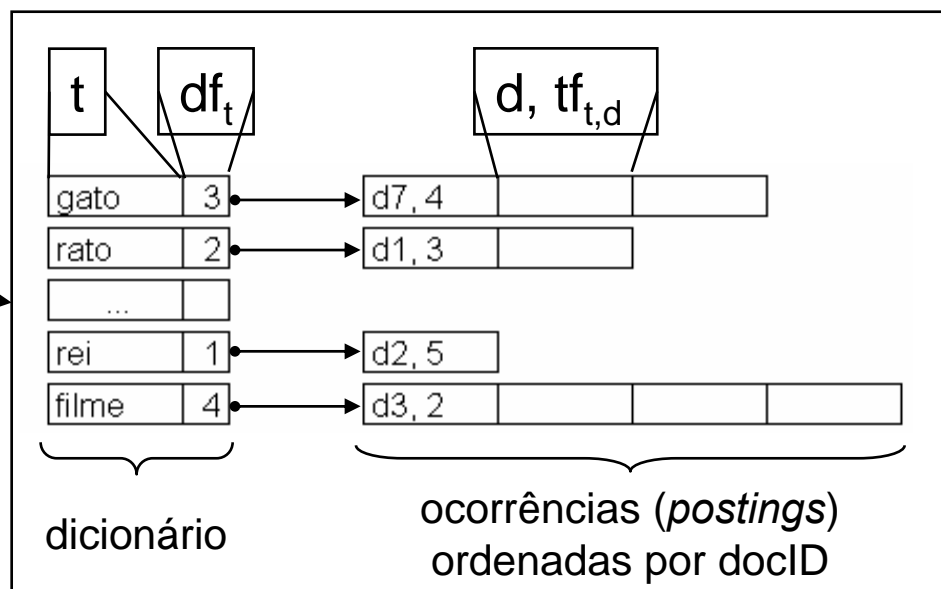
- Para qualquer documento d , calcular $\vec{V}(q) \cdot \vec{v}(d)$ corresponde
 - apenas a somar os pesos de d que constam da interrogação q
- Assim, garante-se relação de ordem sem calcular o co-seno!

... calcular similaridade de forma eficiente

- Manter uma lista de índices invertidos
 - com informação sobre o valor de “df” e “tf”
- Dado uma interrogação q
 - calcular $\vec{V}(q) \cdot \vec{v}(d)$ percorrendo as ocorrências de cada termo, t , em q
 - acumulando, para cada documento d , o valor “tf-idf” de cada termo t

Recordar:

$$\text{tf-idf}_{t,d} = \omega_{t,d} = \text{tf}_{t,d} \times \log_2 (N / \text{df}_t)$$

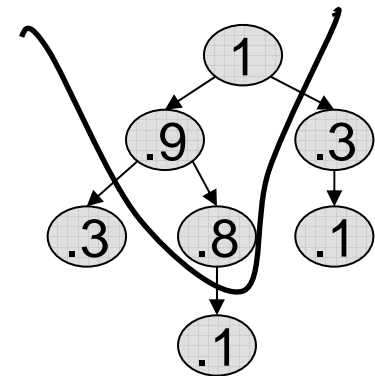


Processo idêntico ao do modelo Booleano!

No entanto, aqui é atribuída uma pontuação (score) a cada documento.

Escolher os K co-seno de forma eficiente

- Uma abordagem consiste em ordenar todos os co-seno
 - e escolher os K maiores
 - ... melhor algoritmo de ordenação tem custo pior caso $O(N \log N)$
- Outra abordagem mais eficiente é a de manter uma estrutura
 - cujo custo de manutenção compense depois na procura dos K co-seno
- Manter uma árvore binária onde
 - o valor de cada nó é maior do que a dos seus filhos (“heap”)
- ... dados N documentos,
 - construir a “heap” tem custo pior caso $O(2N)$
 - obter os K nós tem custo pior caso $O(2 \log N)$



... comparar ordenação com manutenção de “heap”

Com ordenação:

ordenar: $O(N \log N)$

obter K co-senos: $O(1)$

Total: $(N \log N)$

Com árvores (“heap”):

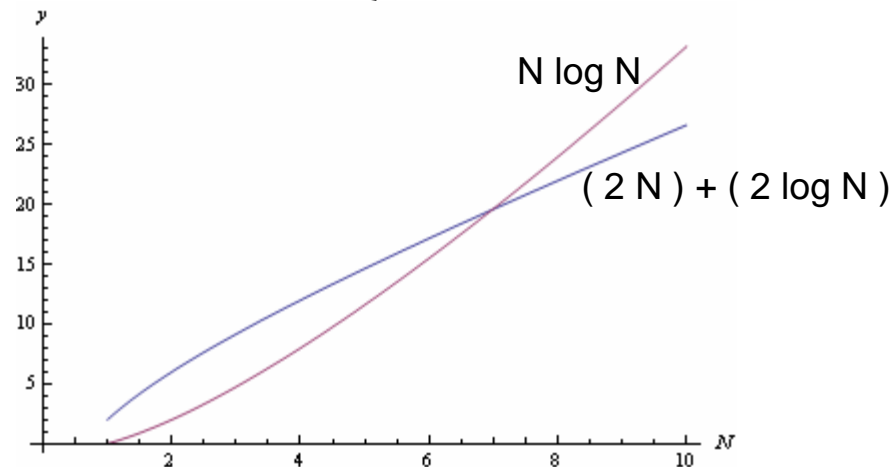
ordenar: $O(2N)$

obter K co-senos: $O(2 \log N)$

Total: $(2N) + (2 \log N)$

$(N \log N)$ cresce mais depressa, a partir de certa altura, do que $(2N) + (2 \log N)$

Para valores grandes de N (dimensão da colecção) compensa manter o “heap”



Os “mais altos” ou os “bons candidatos a mais alto”?

Duas alternativas (face a uma interrogação):

1. procurar os K co-seno mais altos na colecção
(esta foi a abordagem seguida até agora)
2. procurar K bons candidatos a ter o maior co-seno
(esta será a próxima abordagem)

Em geral, passar de uma pesquisa global para diversas pesquisas locais:

- reduz complexidade, mas
- não garante que se encontre a solução óptima.
- ... apesar disso as soluções sub-óptimas podem ser suficientemente boas!

Obter os K co-seno “provavelmente” mais altos

- Queremos os melhores K candidatos a ter co-seno mais alto
 - mas, podem surgir “infiltrados”!
 - ... i.e., um candidato que afinal não tem um dos K co-seno mais altos
- O risco tem o objectivo de reduzir custo de computação
 - tentando não afectar grandemente a percepção do utilizador quanto à relevância dos K documentos apresentados
- ... a própria medida de similaridade (co-seno) já é uma estimativa
 - da noção de relevância do utilizador
 - portanto uma estimativa próxima dessa pode ser suficiente
 - ... de facto a “efectiva relevância” é uma medida de cada utilizador!
- Assim, obter os K co-seno “provavelmente” mais altos
 - não é necessariamente uma abordagem pior para o utilizador
 - do que a de obter exactamente os K co-seno mais altos!

Reduzir complexidade

- Ao procurar os K melhores candidatos reduzindo complexidade
 - deixamos de fazer uma pesquisa exaustiva no espaço global
 - passamos a fazer pesquisas em espaços locais mais pequenos
- ... portanto pode acontecer que nas pesquisas locais
 - “escapem” documentos por estarem fora dos espaços pesquisados
 - “surjam” documentos com alto valor local mas baixo valor global
- Algumas das técnicas para reduzir complexidade incluem
 - heurísticas para descartar índices (“index elimination”)
 - lista campeã (“champion list”, ou “fancy list”)
 - corte por agrupamento (“cluster pruning”)

— Reduzir complexidade – heurísticas para descartar índices

- Heurística 1 – limiar (“threshold”) ε para valor de “idf”
 - apenas considerar documentos com termos cujo “idf” exceda ε .
 - ... os termos com baixo “idf” passam a ser vistos como “stop words”
 - ... em geral origina grande redução no cálculo dos co-seno
- Heurística 2 – interrogação conjuntiva
 - só considerar documentos que tenham todos (ou maioria) dos termos
 - apenas para esses documentos será efectuado o cálculo do co-seno
 - ... a interrogação passa a ser vista como uma conjunção
- ... na interrogação conjuntiva podem ter que se repetir pesquisas
 - quando há menos que K candidatos para os termos considerados
 - e.g., não há K candidatos para 4 termos, repetir pesquisa para 3 termos

Reduzir complexidade – lista campeã

Técnica

lista campeã, ou
“champion list”, ou
“fancy list”.

Pré-Processamento:

1. para cada termo, t , construir o conjunto dos r documentos com maior $tf_{t,d}$;
o valor de r é escolhido ‘a priori’;
aqueles são os documentos da “lista campeã”.

Interrogação q :

1. construir conjunto A com a união das listas campeãs para os termos em q .
2. calcular co-seno apenas dos documentos em A .

O valor de r deve ser alto, em relação a K , para se obterem bons resultados.
Não é necessário ter o mesmo valor de r para todos os termos do dicionário.

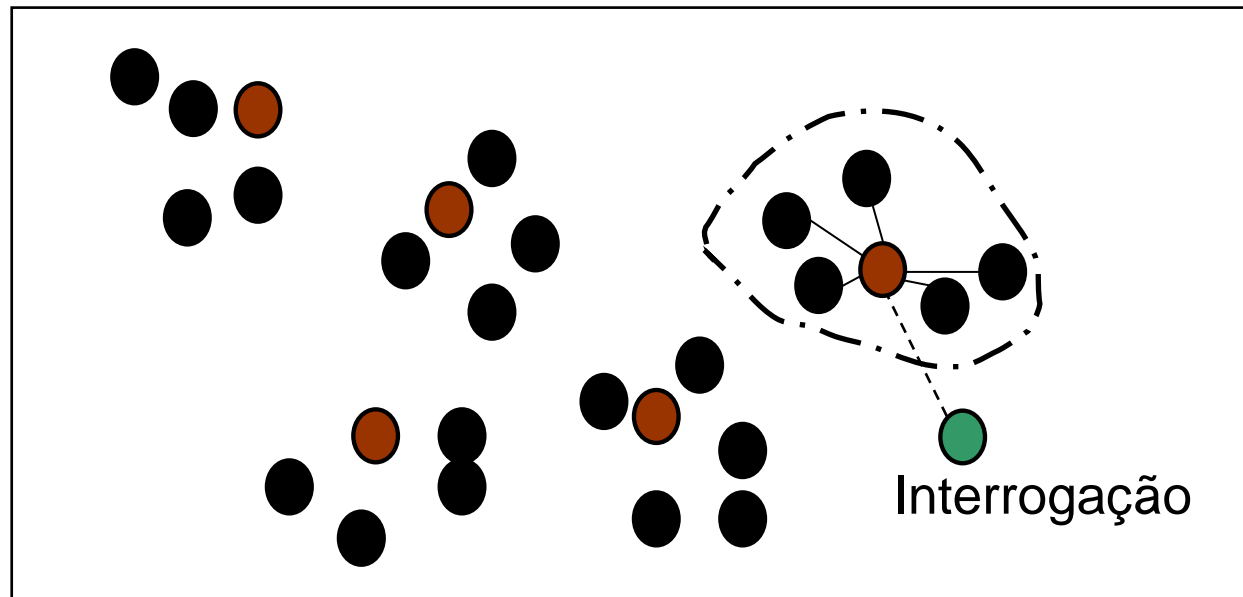
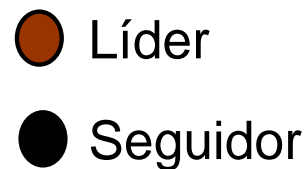
Reduzir complexidade – corte por agrupamento

Pré-Processamento:

1. escolher aleatoriamente \sqrt{N} documentos da colecção;
chamar a cada um desses documentos líder.
2. para cada um dos restantes documentos calcular o seu líder mais próximo
chamar a cada uma desses documento seguidor.

Técnica

corte por agrupamento,
ou “cluster pruning”.



... corte por agrupamento (as partições e a interrogação)

Interrogação q:

1. encontrar o líder (documento) L mais próximo de q (calcular \sqrt{N} co-senos).
2. os documentos candidatos são L e os seus seguidores;
calcular co-seno dos documentos candidatos.

- O espaço de documentos é fraccionado \sqrt{N} em grupos
 - cada grupo é representado pelo seu líder
- ... na partição de seguidores induzida pelo \sqrt{N} líderes
 - o número esperado de seguidores para cada líder será
 - $\approx (N / \sqrt{N}) = \sqrt{N}$

... corte por agrupamento (líderes aleatórios e variação)

- Usar líderes aleatórios é rápido e pode também
 - reflectir a distribuição do espaço de vectores (documentos)
- ... é provável que uma região densa em documentos
 - origine múltiplos líderes e assim partições finas dessa região
- Variação do corte por agrupamento: sejam b_1 e b_2 inteiros
 - pré-processamento:
 - ◊ para cada seguidor identificar os b_1 líderes mais próximos.
 - interrogação q :
 - ◊ considerar os b_2 líderes mais próximos de q .
 - ... note-se que a técnica base corresponde a ter $b_1 = b_2 = 1$
- ... aumentar b_1 e b_2 aumenta a possibilidade de encontrar os K
 - documentos que de facto têm o valor de co-seno mais alto
 - ... à custa de mais processamento!

Métricas de avaliação do resultado da pesquisa

		avaliação do utilizador (humano)	
		Relevante	Não Relevante
resposta da pesquisa (máquina)	Recuperado	positivo (p)	falso positivo (fp)
	Não Recuperado	falso negativo (fn)	negativo (n)

Métricas:

Precisão (“Precision”) \equiv fracção dos documentos recuperados que é relevante

$P(\text{Relevante} \mid \text{Recuperado}) =$

$= \#(\text{documentos relevantes e recuperados}) / \#(\text{documentos } \underline{\text{recuperados}})$

Cobertura (“Recall”) \equiv fracção os documentos relevantes que são recuperados

$P(\text{Recuperado} \mid \text{Relevante}) =$

$= \#(\text{documentos relevantes e recuperados}) / \#(\text{documentos } \underline{\text{relevantes}})$

... métricas de avaliação – cálculo

		avaliação do utilizador (humano)	
		Relevante	Não Relevante
resposta da pesquisa (máquina)	Recuperado	positivo (p)	falso positivo (fp)
	Não Recuperado	falso negativo (fn)	negativo (n)

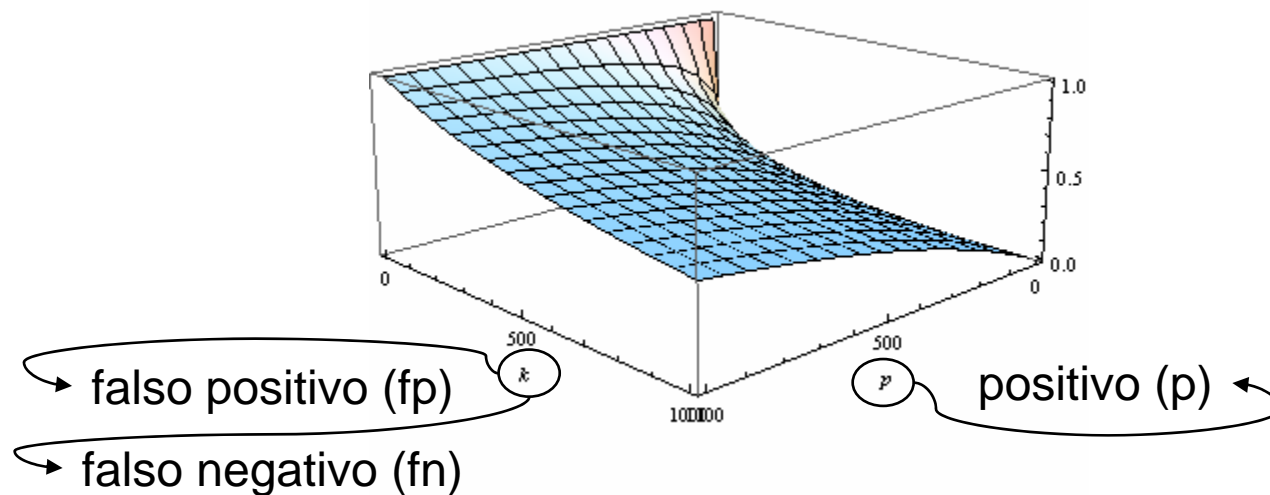
Precisão (“Precision”) \equiv fracção dos documentos recuperados que é relevante
 $= \#(\text{documentos relevantes e recuperados}) / \#(\text{documentos } \underline{\text{recuperados}})$
 $= p / (p + fp)$

Cobertura (“Recall”) \equiv fracção os documentos relevantes que são recuperados
 $P(\text{Recuperado} \mid \text{Relevante}) =$
 $= \#(\text{documentos relevantes e recuperados}) / \#(\text{documentos } \underline{\text{relevantes}})$
 $= p / (p + fn)$

Perspectiva de variação – precisão e cobertura

Precisão (“Precision”) $\equiv p / (p + fp)$

Cobertura (“Recall”) $\equiv p / (p + fn)$



A **precisão** aumenta quando diminuem os falsos positivos.

A **cobertura** aumenta diminuem os falsos negativos.

A **precisão** e a **cobertura** diminuem quando diminuem os positivos.

... precisão e cobertura – “tensão” entre ambas

Para obter óptima cobertura basta recuperar todos os documentos da colecção!
No entanto teria baixa precisão (muitos não seriam relevantes).

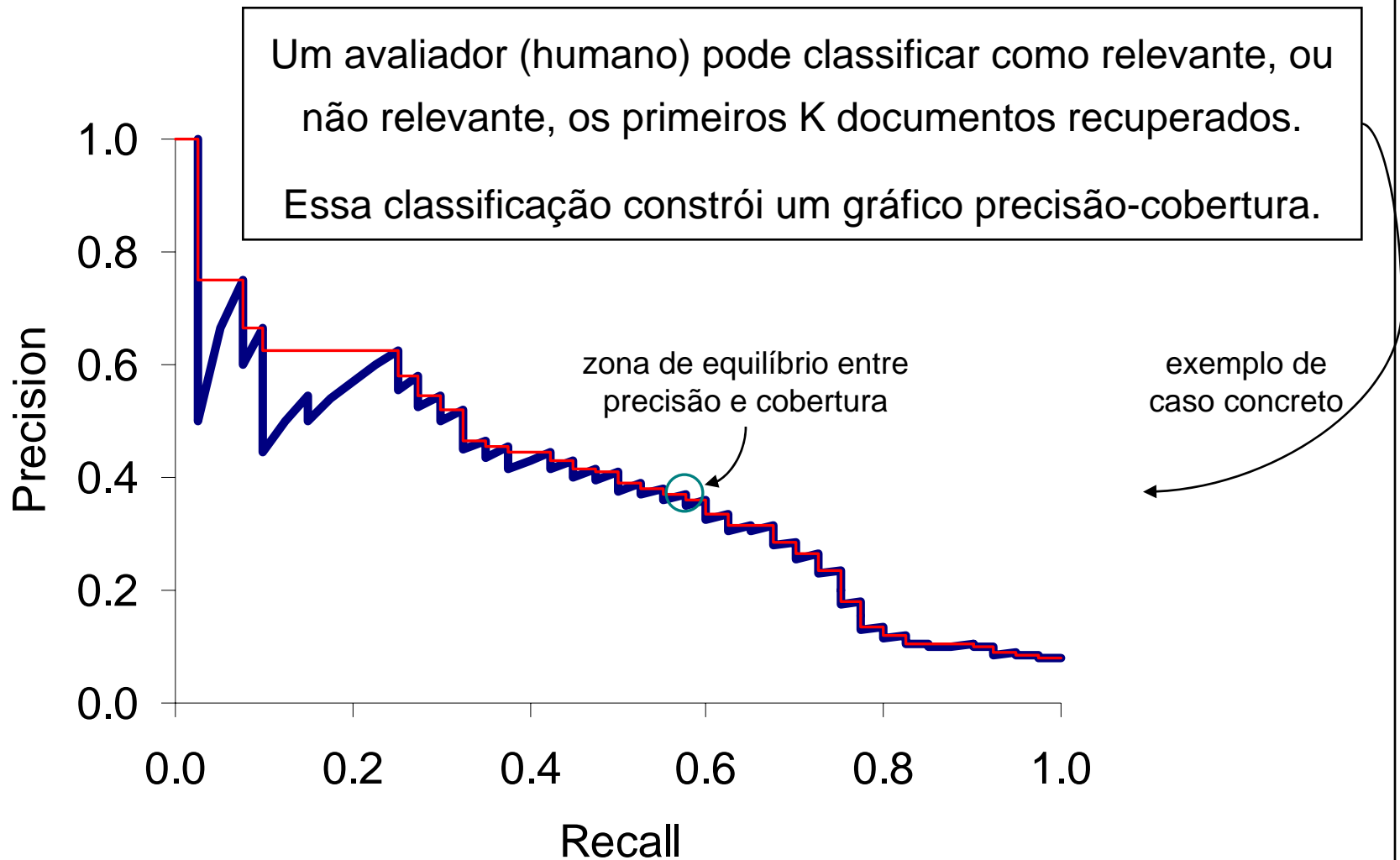
A cobertura não decresce com o número de documentos recuperados.
A precisão pode decrescer com o número de documentos recuperados.

A precisão decresce, em sistemas reais (e.g. comerciais) usualmente com:

- o aumento de documentos recuperados, ou
- o aumento da cobertura.

... empiricamente pode comprovar-se num qualquer um motor de busca.

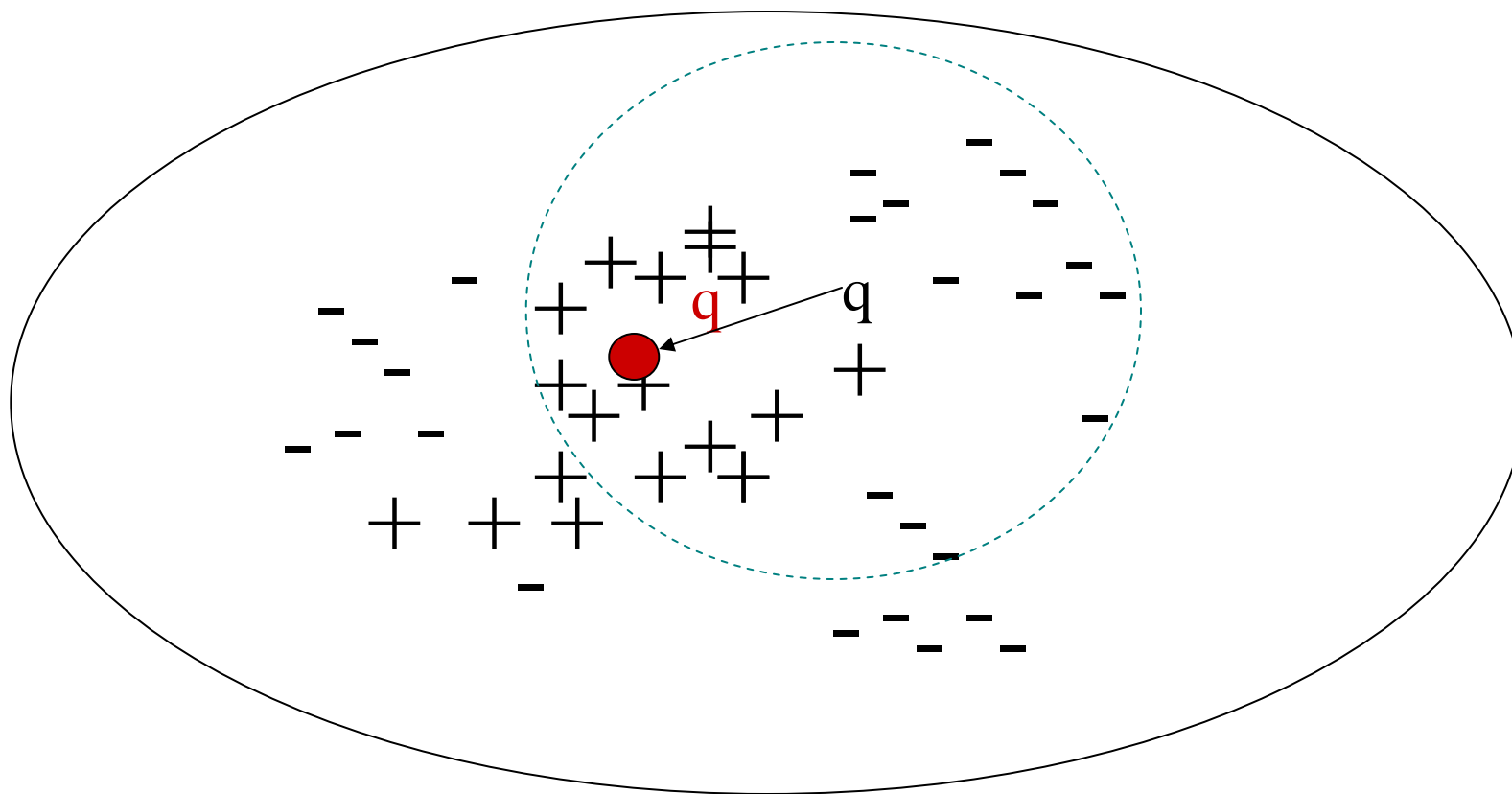
Gráfico precisão-cobertura



Retroacção

- Aprender com exemplos
 - neste caso um exemplo é “a avaliação do utilizador a uma resposta”
- ... os exemplos podem ser positivos ou negativos
 - positivos: documentos considerados relevantes (pelo utilizador)
 - negativos: documentos considerados não relevantes (pelo utilizador)
- Como utilizar os exemplos para melhorar o desempenho?
 - “modificando a interrogação”
 - ... o método mais conhecido é o da “retroacção de Rocchio”
- Como modificar a interrogação?
 - adicionar novos termos
 - ajustar pesos de termos antigos
 - combinar as aproximações anteriores

Retroacção de Rocchio – ideia



Retroacção de Rocchio – formulação

**Nova
Interrogação**

Parâmetros

$$\vec{q}_m = \alpha \vec{q} + \frac{\beta}{|D_r|} \sum_{\forall \vec{d}_j \in D_r} \vec{d}_j - \frac{\gamma}{|D_n|} \sum_{\forall \vec{d}_j \in D_n} \vec{d}_j$$

**Interrogação
Inicial**

**Documentos
Relevantes**

**Documentos
Não Relevantes**

Características do Modelo Vectorial

- Bom desempenho
 - obteve melhores resultados na TREC
- Intuitivo e simples de implementar
 - mais estudado e avaliado
 - implementado no sistema SMART; desenvolvido Cornell 1960-1999
 - ... ainda usado
- Assume independência dos termos
- Assume que documentos e interrogações são o mesmo
- Dificuldade em “afinar” os diversos parâmetros

A reter...

- Modelo Vectorial pertence a uma família de modelos heurísticos
 - ... para recuperação de informação
- Normalização dos pesos (tf-idf) conduz a bons resultados
- Retroacção de Rocchio é modelo eficiente de alteração de pesos
- Modelo Vectorial é usado num grande número de aplicações
- Dificuldade em “afinar” os diversos parâmetros
- Generalização do Modelo Vectorial
 - ... tem maior fundamentação teórica mas é pouco usado na prática

Exercício

Documentos:

Austen; *Sense and Sensibility* (SaS), *Pride and Prejudice* (PaP);
Bronte; *Wuthering Heights* (WH)

Calcular a semelhança
do documento SaS
com os documentos
Pap e WH

	SaS	PaP	WH
<i>affection</i>	115	58	20
<i>jealous</i>	10	7	11
<i>gossip</i>	2	0	6

	SaS	PaP	WH
<i>affection</i>	0,996	0,993	0,847
<i>jealous</i>	0,087	0,120	0,466
<i>gossip</i>	0,017	0,000	0,254

$$\cos(\text{SaS}, \text{PaP}) = .996 \times .993 + .087 \times .120 + .017 \times 0.0 = 0.999$$
$$\cos(\text{SaS}, \text{WH}) = .996 \times .847 + .087 \times .466 + .017 \times .254 = 0.929$$