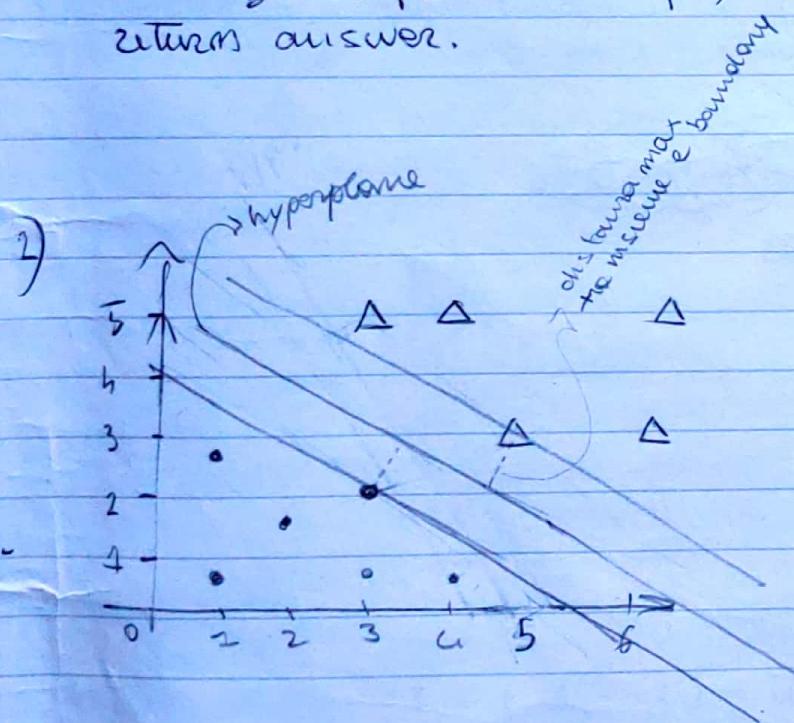


EXAM 11/02/19

1) INTERSECT (p_1, p_2) pastings lists T_1 AND T_2

answer $\leftarrow ()$
while $p_1 \neq \text{NIL}$ and $p_2 \neq \text{NIL}$
do if docID(p_1) = docID(p_2)
then ADD (answer, docID(p_1))
 $p_1 \leftarrow \text{next}(p_1)$
 $p_2 \leftarrow \text{next}(p_2)$
else ~~if~~ if docID(p_1) < docID(p_2)
then $p_1 \leftarrow \text{next}(p_1)$
~~else~~ $p_2 \leftarrow \text{next}(p_2)$
return answer.



Identify the linear maximum margin classifier. Draw 3 lines:
the two boundaries and the hyperplane
which of the vectors are support vector?

- Scegli gli elementi che vengono proiettati tra loro e creare un iperplano in modo da massimizzare la distanza.

devo massimizzare la distanza
tra i punti di una classe e il
decision boundary.

I support vector sono $(3, 3)$ e
 $(4, 5, 3)$

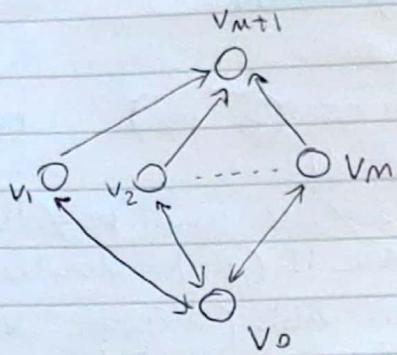
③ PageRank equations
teleporting probability α
personalization vector $\frac{1}{M}$

$$\sum_i \pi_i = 1 \quad \forall i=0 \dots M+1$$

$$\pi_0 = \frac{\alpha}{M} + \frac{(1-\alpha)}{2} \sum_{i=1}^M \pi_i$$

$$\pi_i = \frac{\alpha}{M} + \frac{(1-\alpha)}{M} \sum_{j \neq i} \pi_j \quad i=1 \dots M$$

$$\pi_{M+1} = \frac{\alpha}{M} + \frac{(1-\alpha)}{2} \sum_{i=1}^M \pi_i$$



④ personalise vector $p = (0 \dots 1)^T$ all components are 0 except for v^{M+1} . Compute pagerank π . Motivate the answer

$$\pi = \alpha p_i + (1-\alpha) \sum_i \frac{\pi_i}{d_i}$$

$$\begin{cases} \pi_0 = \alpha \cdot p_i + (1-\alpha) \sum_{l=1}^M \pi_l \end{cases}$$

$$\begin{cases} \pi_i = \alpha \cdot p_i + (1-\alpha) \pi_0 \end{cases}$$

$$\begin{cases} \pi_{M+1} = \alpha \cdot p_i + (1-\alpha) \sum_{l=1}^M \pi_l \end{cases}$$

$$\begin{cases} \pi_0 = (1-\alpha) \sum_{l=1}^M \pi_l \end{cases}$$

$$\begin{cases} \pi_i = \frac{(1-\alpha)}{2} \pi_0 \end{cases}$$

$$\pi_{M+1} = \alpha + \frac{(1-\alpha)}{2} \sum_{l=1}^M \pi_l$$

$$\text{formulae generali} \quad \pi = (1-\alpha)M + \alpha(I_p)$$

matrice identità

1) UNION (p_1, p_2) postings lists T₁ OR T₂

answer $\leftarrow ()$

while $p_1 \neq \text{NIL}$ and $p_2 \neq \text{NIL}$

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

if $\text{docID}(p_1) < \text{docID}(p_2)$

Then ADD (Answer, $\text{docID}(p_1)$)

$p_1 \leftarrow \text{next}(p_1)$

else if $\text{docID}(p_1) > \text{docID}(p_2)$

Then ADD (Answer, $\text{docID}(p_2)$)

$p_2 \leftarrow \text{next}(p_2)$

if $\text{docID}(p_1) = \text{docID}(p_2)$

Then ADD (Answer, $\text{docID}(p_1)$)

$p_1 \leftarrow \text{next}(p_1)$

$p_2 \leftarrow \text{next}(p_2)$

if $p_1 = \text{NIL}$

while ($p_2 \neq \text{NIL}$)

then ADD (Answer, $\text{docID}(p_2)$)

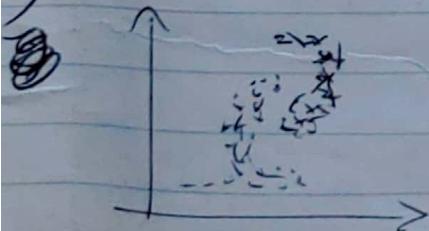
if $p_2 = \text{NIL}$

while ($p_1 \neq \text{NIL}$)

then ADD (Answer, $\text{docID}(p_1)$)

return answer

2) Better SVM or KNN ?



Linear SVM is not good and soft margin technique is not sufficient because dataset is not linearly separable.

So, in this case KNN is better.

However we can use kernelization to make this set linearly separable and then use SVM.

3) teleporting prob. &
personalization vector $\{1, 0, 0, 0\}$

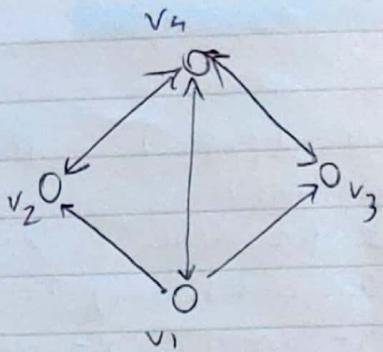
$$\sum_i \pi_i = 1 \quad i=1 \dots 4$$

$$\pi_2 = \pi_3 \quad \text{symmetry}$$

$$\pi_1 = \alpha + \frac{(1-\alpha)}{3} \pi_u$$

$$\pi_2 = \frac{(1-\alpha)}{3} \pi_1 + \frac{(1-\alpha)}{3} \pi_u$$

$$\pi_4 = \frac{(1-\alpha)}{3} \pi_1 + (1-\alpha) \pi_2 + (1-\alpha) \pi_3$$



(b) $\pi_1 \quad p_1 = \{1, 0, 0, 0\}^T$

$\pi_2 \quad p_2 = \{0, 1, 0, 0\}^T$

Explain in details how to calculate the personalized pagerank
with respect to the personalization vector $\{0.5, 0.5, 0.5, 0\}$

$$\begin{aligned}\pi_1 &= (1-\alpha) \sum_{j \in N_i} \frac{\pi_j}{d_j} + \alpha \cdot p_i \\ \pi_1 &= \frac{\alpha}{2} + \frac{(1-\alpha)}{3} \pi_u \\ \pi_2 &= \frac{\alpha}{2} + \frac{(1-\alpha)}{3} (\pi_1 + \pi_u) \\ \pi_u &= \frac{(1-\alpha)}{3} \pi_1 + (1-\alpha) (\pi_2 + \pi_3) \\ \pi_2 &= \pi_3\end{aligned}$$

$$\begin{aligned}p_3 &= \left\{ \frac{1}{2}, \frac{1}{2}, 0, 0 \right\} \\ \pi_1(p_3) &= \pi_1(p_1) - \frac{\alpha}{2} \\ \pi_2(p_3) &= \pi_2(p_2) - \frac{\alpha}{2} \\ \pi_3(p_3) &= \pi_3(p_1) = \frac{2}{\pi_3(p_2)} \\ \pi_u(p_3) &= \pi_u(p_1) = \pi_u(p_1)\end{aligned}$$

You have to suppose that you can jump with probability 1 to node v_1
or probability 2 to node v_2

If we have $p = \left\{ \frac{1}{2}, \frac{1}{2}, 0, 0 \right\}$ it means that we have to
~~subtract~~ subtract $\frac{\alpha}{2}$ to π_1 and π_2

$$\pi_2 = \pi_1 - \frac{\alpha}{2}$$

$$\pi_2 = \pi_2 - \frac{\alpha}{2}$$

thus, we have

$$v_{NB} = \arg \max P(v_3) \prod_i p_i(a_i/v_3)$$

assuming that the attribute values are conditionally independent given the target value, it is possible to notice that the probability of observing the conjunction $a_1 \dots a_m$ is just the product of the probability for the individual attribute. Substituting

2) $\rightarrow (\text{Yes}, \text{Yes}, \text{No})$

Which class (Fail or Pass) will be assigned to X by Naive Bayes?

$$v_{NB} = \underset{v_j \in V}{\operatorname{argmax}} P(v_j) \prod_i P(a_i/v_j)$$

attribute
probabilità di un evento

$$P(\text{Fail}) = 2/5$$

$$P(\text{Pass}) = 3/5$$

• probabilità ($\text{Yes}, \text{Yes}, \text{No} / P(\text{Fail})$)

$$P(\text{Confident} = \text{Yes} / P(\text{Fail})) = 1/2$$

$$P(\text{Student} = \text{Yes} / P(\text{Fail})) = 3/2$$

$$P(\text{Sick} = \text{No} / P(\text{Fail})) = 1/2$$

~~$$\Rightarrow P(\text{Fail}) \cdot \prod_i P(a_i / \text{Fail}) =$$~~

~~$$P(\text{Fail}) \cdot P(\text{Result} = \text{Fail}) \cdot P(\text{Confident} = \text{Yes} / \text{Result} = \text{Fail}) \cdot$$~~

~~$$P(\text{Student} = \text{Yes} / \text{Result} = \text{Fail}) \cdot$$~~

~~$$P(\text{Sick} = \text{No} / \text{Result} = \text{Fail}) =$$~~

$$= 2/5 \cdot 1/2 \cdot 1/2 \cdot 1/2 = \frac{1}{20} \approx 0,05$$

• probabilità ($\text{Yes}, \text{Yes}, \text{No} / P(\text{Pass})$)

$$P(\text{Confident} = \text{Yes} / \text{Result} = \text{Pass}) = 2/3$$

$$P(\text{Student} = \text{Yes} / \text{Result} = \text{Pass}) = 2/3$$

$$P(\text{Sick} = \text{No} / \text{Result} = \text{Pass}) = 1/3$$

$$P(\text{Result} = \text{Pass}) \cdot P(\text{Confident} = \text{Yes} / \text{Result} = \text{Pass}) \cdot P(\text{Student} = \text{Yes} / \text{Result} = \text{Pass})$$

$$P(\text{Sick} = \text{No} / \text{Result} = \text{Pass})$$

$$= 3/5 \cdot 2/3 \cdot 2/3 \cdot 1/3 = \frac{4}{45} \approx 0,08$$

Ora che dato $\text{Yes}, \text{Yes}, \text{No}$, la probabilità più alta è che passi poiché bisogna considerare $\boxed{\operatorname{argmax}}$

i valori degli attributi sono condizionalmente indipendenti.
 In questo modo le passate vederie per prob. delle attributi come il prob. delle singole attr.

3) Indicate where the conditional independence assumption

The Bayesian approach to classify a new instance is to assign the most probable target value v_{map} given the attribute value $a_1 \dots a_m$ that describe the instance. $v_{map} = \operatorname{argmax} P(v_s)$
 by using the Bayes theorem, we can rewrite the previous equation as
 $v_{map} = \operatorname{argmax} P(v_s) \cdot P(a_1 \dots a_m / v_s)$

1)

$$TF = \frac{\text{numb. words}}{\text{numb. words in document}}$$

$$IDF = \log \frac{N}{df}$$

$$TF \cdot IDF$$

$\frac{N}{df} = \frac{\text{number of docs}}{\text{number of docs in which term occurs}}$

DOC	T	D2	TF = numb. words in document	IDF	TF/IDF
D1	anccaa		0	$\frac{N}{df} = \frac{5}{2}$	
	dush		0		
	pork		3		
	rabbit		0		
	recep		0		
	roast		0		
D2					

T	TF _{D1}	TF _{D2}	TF _{D3}	TF _{D4}	TF _{D5}
anccaa	0	0	3	0	0
dush					
pork					
rabbit					
recep					
roast					

Vegan
Dinner

$$Q = \begin{pmatrix} 1 & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{3} & 1 & \frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} & 1 \end{pmatrix}$$

	anua	pork	reape
D ₁	0	$\frac{15}{64}$	0
D ₂	$\frac{5}{42}$	$\frac{10}{84}$	0
D ₃	0	$\frac{10}{212}$	$\frac{5}{106}$
D ₄	0	0	$\frac{5}{30}$
D ₅	$\frac{5}{80}$	$\frac{5}{160}$	0

$$D_1 = \begin{pmatrix} 0 \\ \frac{15}{64} \\ 0 \end{pmatrix} \quad D_2 = \begin{pmatrix} 5/64 \\ 10/84 \\ 0 \end{pmatrix} \quad D_3 = \begin{pmatrix} 0 \\ 10/212 \\ 5/106 \end{pmatrix} \quad D_4 = \begin{pmatrix} 0 \\ 0 \\ 5/30 \end{pmatrix} \quad D_5 = \begin{pmatrix} 5/80 \\ 5/160 \\ 0 \end{pmatrix}$$

$$D_1 = \frac{15}{64} \quad D_2 = \frac{20}{84} \quad D_3 = \frac{10}{212} \quad D_4 = \frac{5}{30} \quad D_5 = \frac{15}{160}$$

$0,23 \quad 0,23 \quad 0,09 \quad 0,15 \quad 0,09$

$$\cos. \sinm (Q, D_1) = \frac{\frac{1}{3} \cdot 0 + \frac{1}{3} \cdot \frac{15}{64} + \frac{1}{3} \cdot 0}{\sqrt{\frac{15^2}{64}} \cdot \sqrt{\frac{1}{3}^2 + \frac{1}{3}^2 + \frac{1}{3}^2}} = \frac{\frac{5}{64}}{\frac{15\sqrt{3}}{64} \cdot \frac{\sqrt{3}}{3}} = \frac{5 \cdot \cancel{64}}{15 \cdot \cancel{64} \cdot \sqrt{3} \cdot \sqrt{3}} = \frac{1}{\sqrt{3}} \approx 0,36$$

$$(Q, D_2) = \frac{\frac{1}{3} \cdot \frac{5}{42} + \frac{1}{3} \cdot \frac{10}{84} + 0}{\frac{\sqrt{3}}{3} \cdot \sqrt{\left(\frac{5}{42}\right)^2 + \left(\frac{10}{84}\right)^2}} = \frac{\frac{\sqrt{3}}{3} \cdot \sqrt{0,02 + 0,02}}{\frac{\sqrt{3}}{3} \cdot \sqrt{0,02 + 0,02}} = \frac{\frac{10+10}{252}}{\frac{\sqrt{3} \cdot 0,02}{3}} = \frac{20}{\frac{\sqrt{3} \cdot 0,02}{3}} = \frac{20}{\frac{282}{84}} = \frac{20}{\frac{141}{42}} = \frac{20 \cdot 42}{141} = \frac{80}{23} \approx 3,48$$

$$0,24 \cdot 18,4 = 4,42$$

In grafo diretto è fortemente connesso se
esiste un percorso per andare in ogni nodo
(se ogni nodo del grafo è raggiungibile da ogni altro)

3) Unico stat. chsl? ?

⊗ Siccome il grafo bipartito è un grafo con un numero finito
di stati, allora esiste almeno una stationary distribution.
Tuttavia, siccome il grafo ^{e' non diretto} e' fortemente connesso, la
stationary distrib. non è unica

Non è fortemente connesso perché FGFR3 è collegato a
^(maestro se fosse stato orientato) ^{Bloodlet e Colm Me}
questi sono solo
collegati a FGFR3,
quindi non avranno
in nessun altro modo

b) Queste strutture è per nerals e anthoniles
Queste prendono alg. di questo genere

HUBS and AUTHORITIES (ESEMPIO INTERNET)

measure of importance similar to page rank.

Uses both outlinks and incoming links.

→ PUNTA

HUB: pagina che punta a un'altra pagina

AUTHORITY: pagina puntata da un'altra pagina

→ È PUNTATO

AUTHORITY SCORE: somma degli hub score delle pagine che puntano alla pagina.

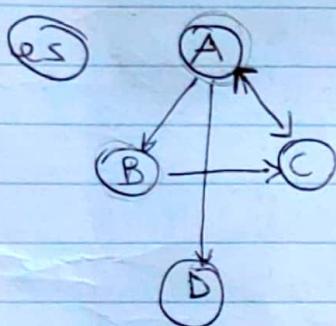
$$x_i^{(k)} = \sum_{j : e_{ji} \in E} y_j^{(k-1)}$$

HUB SCORE:

HUB SCORE: somma degli authority score delle pagine a cui punta

$$y_i^{(k)} = \sum_{j : e_{ij} \in E} x_j^{(k)}$$

↓
authority



A è hub di B D C hub

B è hub di C

C è hub di A

	A	B	C	D
A	0	1	1	1
B	0	0	1	0
C	1	0	0	0
D	0	0	0	0

A è authority di C

B è authority di A

C è authority di A e B

D è authority di A

$$i = 0$$

$$i = 1$$

$$y_i = 1$$

$$x_i^{(k)} = \sum_{j : e_{ji} \in E} y_j^{(k-1)}$$

authorities

$$\Rightarrow x_i = \sum_{j : e_{ij} \in E} y_j^0$$

quando conto gli authorities

$$x_A = \sum_{c : e_{CA} \in E} y_c^0 = 1$$

$$x_B = \sum_{a : e_{AB} \in E} y_a^0 = 1$$

$$x_C = \sum_{ab : e_{AC/BC}} y_a^0 + y_b^0 = 1 + 1 = 2$$

$$x_D = \sum_{a : e_{AD} \in E} y_a^0 = 1$$

normalizzo i valori al massimo valore

$$\frac{1}{2} = 0,5$$

$$\frac{1}{2} = 0,5$$

$$\frac{2}{2} = 1$$

$$\frac{1}{2} = 0,5$$

$$hub \gamma$$

$$y_i^{(k)} = \sum_{j : e_{ij} \in E} x_j^{(k)}$$

$$y_A = x_B^1 + x_D^1 + x_C^1 = 0,5 + 0,5 + 1 = 2$$

$$y_B = x_C^1 = 1$$

$$y_C = x_A^1 = 0,5$$

$$\frac{2}{2} = 1$$

$$\frac{1}{2} = 0,5$$

$$\frac{0,5}{2} = 0,25$$

$$\frac{0}{2} = 0$$

$$\boxed{[k=2]} \quad x_i^{(k)} = \sum y_j^{(k-1)} \rightarrow x_i^2 = \sum y_j^2$$

$$x \quad \left\{ \begin{array}{l} x_A = y_C = 0,25 \\ x_B = y_A = 1 \\ x_C = y_A + y_B = 1 + 0,5 = 1,5 \\ x_D = y_A = 1 \end{array} \right.$$

$$\left. \begin{array}{l} \text{normalizzazione} \\ x_A = \frac{0,25}{2,5} = 0,16 \\ x_B = \frac{1}{2,5} = 0,40 \\ x_C = \frac{1,5}{2,5} = 0,60 \\ x_D = \frac{1}{2,5} = 0,40 \end{array} \right\}$$

$$y \quad \left\{ \begin{array}{l} y_i^{(k)} = \sum x_j^{(k)} \rightarrow \text{moltre} \\ y_A = x_B^{(2)} + x_D^{(2)} + x_C^{(2)} = 0,66 + 0,66 + 1 = 2,32 \\ y_B = x_C^{(2)} = 1 \\ y_C = x_A^{(2)} = 0,16 \end{array} \right.$$

$$\left. \begin{array}{l} y_A = \frac{1,87}{2,32} = 0,80 \\ y_B = \frac{1}{2,32} = 0,43 \\ y_C = \frac{0,66}{2,32} = 0,28 \\ y_D = 0 \end{array} \right\}$$

→ questa operazione converge quando tutti i valori si trovano sotto un valore ~~fissato~~ fissato inizialmente
 (attualmente si continua)

hub = pagina che punta a un'altra pagina

authority = pagina che viene puntata

EXAM 23/01/20

- 3) HITS algorithm con adjacency matrix A
a(t) and h(t) are authorities and hubs

a) Indichiamo con E l'insieme di tutti i collegamenti diretti nel grafo che rappresenta la rete e chiamiamo e_{ij} il collegamento dal modo i al modo j .

A ciascuna pagina sono associati authority score $a_i^{(0)}$ e hub score $h_i^{(0)}$. In particolare $h_i^{(0)}$ viene inizializzato a 1.

Successivamente HITS effettua i punteggi calcolando

$$a_i^{(k)} = \sum_{j: e_{ij} \in E} h_j^{(k-1)} \quad h_i^{(k)} = \sum_{j: e_{ij} \in E} a_j^{(k)} \quad k = 1, 2, 3, \dots$$

possiamo scrivere queste formule in forma matriciale con la matrice di adiacenza L

$$L_{ij} = \begin{cases} 1 & \text{se c'è un link da pagina } i \text{ a } j \\ 0 & \text{altrimenti} \end{cases}$$

$$a^{(k)} = L^T h^{(k-1)} \quad h^{(k)} = L \cdot a^{(k)}$$

3) Inizialmente $h^{(0)}$ è un vettore colonna di tutti 1

3) Fino a convergenza ripetere:

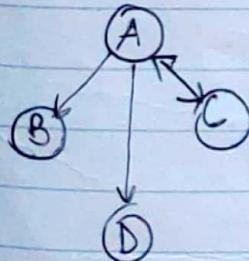
$$a^{(k)} = L^T h^{(k-1)}$$

$$h^{(k)} = L \cdot a^{(k)}$$

$$k = k + 1$$

Normalizzazione $a^{(k)}, h^{(k)}$ rispetto al valore più alto rispettivamente di $a^{(k)}$ e $h^{(k)}$

b) esempio network con vertici che hanno hub score = 0



In questo caso B e D hanno hub score 0 perché non puntano a niente

hub	A	B	C	D	
L_{ij}	A	0	1	1	1
	B	0	0	0	0
	C	1	0	0	0
	D	0	0	0	0

3) Discuss how paperrank could be used as an alternative measure of paper's scientific impact

EXAM 11/01/2011

1) general teleporting probability α

\rightarrow teleporter in own node

2) $\pi_2 = \pi_4$ (for symmetry)

$\pi_1 = \pi_5$ links ($1 \rightarrow$)

$$\pi_3 = \frac{\alpha}{7} + \frac{2}{2} \left(\frac{1-\alpha}{2} \pi_1 \right) + \frac{(1-\alpha)}{2} \pi_7$$

2 perché
 $\pi_2 = \pi_4$

link uscita da un nodo vicini

$$\pi_6 = \frac{\alpha}{7} + \frac{(1-\alpha)}{2} \pi_7 + \frac{(1-\alpha)}{3} \pi_3$$

$$\pi_7 = \frac{\alpha}{7} + \frac{(1-\alpha)}{2} \pi_6$$

$$\pi_2 = \frac{\alpha}{7} + \frac{(1-\alpha)}{2} \pi_1 + \frac{(1-\alpha)}{3} \pi_3 = \pi_4$$

$$\pi_1 = \frac{\alpha}{7} + \frac{(1-\alpha)}{1} \pi_2 \rightarrow \cancel{\pi_2} = \pi_5$$

$$\sum_i \pi_i = 1 \quad i=1 \dots 7$$

b) $\alpha = \frac{1}{2}$

$$\pi_1 = \frac{1}{14} + \left(1 - \frac{1}{2}\right) \pi_2 + \cancel{\pi_3}$$

$$\pi_2 = \frac{1}{14} + \left(1 - \frac{1}{2}\right) \pi_3 + \left(\frac{1}{2}\right) \pi_1$$

$$\pi_7 = \frac{1}{14} + \frac{1}{2} \pi_6$$

$$\pi_6 = \frac{1}{14} + \frac{1}{4} \pi_7 + \left(\frac{1}{2}\right) \pi_3$$

$$\pi_3 = \frac{1}{14} + \cancel{\pi_2} + \cancel{\pi_6} + \frac{1}{2} \pi_1 + \frac{1}{2} \pi_7$$

$$\pi_2 = \pi_6$$

$$\pi_1 = \pi_5$$

CALCO V' SINTI
LUNGHISSIMA
NON HO VOGLIA
DI FARLA
(LA SCIA PERDE)
PURPLE TV)

① What is the importance of teleporting probability with respect to the convergence of pagerank?

We use teleporting probability in pagerank to avoid the interruption of random walk due to spider traps or dead ends. In fact teleporting probability is the possibility to jump to another node in the network, with a certain probability, to exit loops in the graphs.

512 128 64 32 16 8 4 2 1
 0 0 0 0 0 0 0 0 0

2) Show how can compress the list [2, 8, 17, 22, 30, 40, 52, 80]
using

- (a) variable byte encoding
- (b) γ encoding.

$\rightarrow 2, 6, 4, 5, 8, 10, 12, 28$
riservo i valori con gap

(a) 2 \rightarrow 000000010
 6 \rightarrow 000001000000000110
 9 \rightarrow 000100010001
 5 \rightarrow 000001001
 8 \rightarrow 00001000
 10 \rightarrow 00001010
 12 \rightarrow 00001100
 28 \rightarrow 00000000011100

quando cresce il codice ~~mettendo~~ separando ~~per~~ per ultimo 4 bit e mettendoci 0 1 come primo bit dell'ultimo byte.
Mentre come primo bit del byte prima metto uno 0.

Per decodificare con variable byte code leggiamo una sequenza di byte con combinazione bit 0 terminata da un byte con combinazione bit 1

Indica che è l'ultimo byte del numero

2 6 9 5
 0 100000010 100000110 10001001 10000101 0 8
 10001000 40001010 10001100 28 100011100

- 3) γ encoding

2 \rightarrow 00000010 offset \rightarrow 0 \rightarrow $m_{mano} = 10 \Rightarrow m_{mano} + offset$
 6 \rightarrow 0110 offset \rightarrow 10 \rightarrow $m_{mano} = 110 \rightarrow 11010$
 9 \rightarrow 0110 offset \rightarrow 001 \rightarrow $m_{mano} = 1110 \rightarrow 1110001$
 5 \rightarrow 0110 offset \rightarrow 01 \rightarrow $m_{mano} = 110 \rightarrow 11001$
 8 \rightarrow 0110 offset \rightarrow 000 \rightarrow $m_{mano} = 1110 \rightarrow 1110000$
 10 \rightarrow 0110 offset \rightarrow 010 \rightarrow $m_{mano} = 1110 \rightarrow 1110010$
 12 \rightarrow 0110 offset \rightarrow 100 \rightarrow $m_{mano} = 1110 \rightarrow 1110100$
 28 \rightarrow 0110 offset \rightarrow 1100 \rightarrow $m_{mano} = 11110 \rightarrow 111101100$

3) ② give a linear algorithm for web AND Information AND NOT Retrieval

INTERSECT ($p_1, p_2, \gamma p_3$)

```
result AND  $\leftarrow ()$  // gestisco AND  
answer  $\leftarrow ()$  // gestisco AND NOT  
while  $p_1 \neq \text{NULL}$  and  $p_2 \neq \text{NULL}$   
  if  $\text{docID}(p_1) = \text{docID}(p_2)$   
    ADD  $(p_1, \text{result AND})$   
     $p_1 \leftarrow \text{next}(p_1)$   
     $p_2 \leftarrow \text{next}(p_2)$   
  if else  $\text{docID}(p_1) < \text{docID}(p_2)$   
     $p_1 \leftarrow \text{next}(p_1)$   
  else ( $\text{docID}(p_2) < \text{docID}(p_1)$ )  
     $p_2 \leftarrow \text{next}(p_2)$ 
```

while $\text{result AND} \neq \text{NULL}$ and $p_3 \neq \text{NULL}$

```
if  $\text{docID}(\text{result AND}) = \text{docID}(p_3)$   
  result AND  $\leftarrow \text{next}(\text{result AND})$   
   $p_3 \leftarrow \text{next}(p_3)$ 
```

if else $\text{docID}(\text{result AND}) < \text{docID}(p_3)$

ADD $(\text{result AND}, \text{answer})$

result AND $\leftarrow \text{next}(\text{result AND})$

else ($\text{docID}(\text{result AND}) > \text{docID}(p_3)$)

~~ADD $(\text{result AND}, \text{answer})$~~

$p_3 \leftarrow \text{next}(p_3)$

~~return answer.~~

if $p_3 = \text{NULL}$

while $\text{result AND} \neq \text{NULL}$

ADD $(\text{answer}, \text{docID}(\text{result AND}))$

return answer

X AND NOT Y
prendo gli elementi
di X che non sono
in Y

Ex

$P_1: 3, 4, 5, 8, 10$

$P_2: 3, 4, 6, 8, 19$

result AND : 3, 4, 8

$P_3: 6, 8, 9$

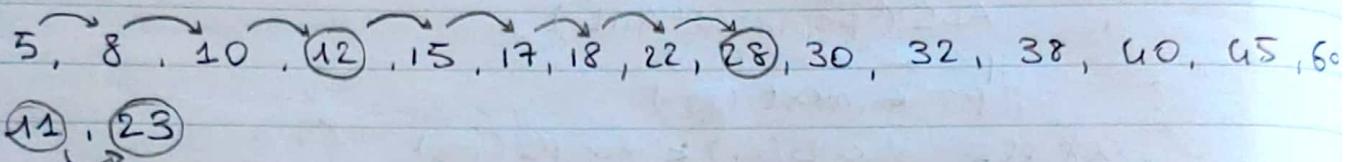
answer : 3, 4

(b) 2-word query. The postings lists are:

[5, 8, 10, 12, 15, 17, 18, 22, 28, 30, 32, 38, 40, 45, 60]
[11, 23]

how many comparisons for the intersection and justify the answer.

② Using standard posting lists



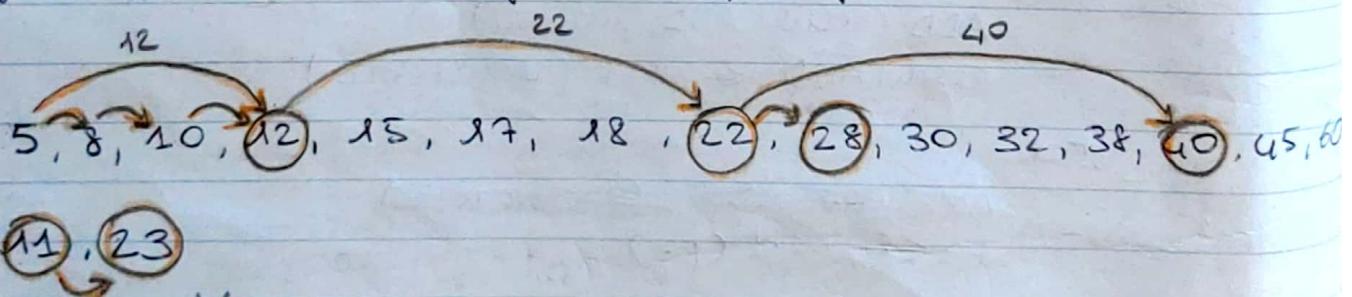
comincio con 5 e 11, $5 < 11$ quindi scendo fino a trovare un numero uguale o > 11 .

$12 > 11$ quindi prendo il successivo di 11 (23) e siccome $12 < 23$ scendo, successivi di 12 ecc.

In questo caso l'intersezione restituisce una lista vuota e ho effettuato 10 confronti

answer = {}

③ Using lists stored with skip pointers with skip length of \sqrt{P} where P is the length of the postings list.



answer = {}

(quasi)

Siccome $P = 15$, $\sqrt{P} = 3$ e qualcosa, quindi i salti saranno di $\frac{3}{2}$ elementi circa

comincio confrontando 5 e 11, $5 < 11$ quindi ~~prendo il successivo~~ dovrei prendere il successivo, ma siccome ho uno skip, provo a vedere se lo skip è ≤ 11 . $12 > 11$, quindi potrei avere l'11 tra gli elementi che ho saltato, allora riprendo il confronto da 5.

$8 < 11$, $10 < 11$. Siccome ~~non~~ a mom ho 11 nelle prime lista e $12 > 11$, allora prendo il successivo di 11, 23.

$23 > 12$ quindi provo a usare lo skip. $22 < 23$, quindi posso saltare i valori compresi tra 12 e 23 e tento lo skip a 40. $40 > 23$, quindi torno a 22 e comincio i successivi. $28 > 23$, quindi termine il confronto.

Ho effettuato 8 confronti, quindi è stato vantaggioso

STEMMING: vado a considerare solo le radice delle parole
COMPUTER → COMPUTER

LEMMAZATION: vado a considerare i lemmi delle parole

SONO → SONO
SEI → SEI → ESSERE
È → È → ESSERE

EXAM 26/01/2012

① What are AND and OR queries in Monotone's crawler URL frontier scheme? Explain how they work.

TRUE OR FALSE

① a) In a Boolean retrieval sys., stemming never lowers precision

FALSE: stemming can increase the retrieved set without increasing the number of relevant documents

② " . stemming never lowers recall

TRUE: stemming can only increase the retrieved set, which means increased or unchanged recall

③ Stemming increases the size of the vocabulary

FALSE: stemming decreases the size of the vocabulary because I reduce different words to the same radix.

④ Stemming should be invoked at indexing time but not while processing a query.

FALSE: the same processing should be applied to documents and queries to ensure matching terms.

⑤ Why skip pointers are not useful for (x OR y)?

Because in order to compute the union of the two, you have to run both the lists one element by one.

⑥ Binary index "New York" "York University"

⑦ Example of a document returned for a query New York University which is a false positive that should not be returned.

Document = "I like New York. ^{YORK} University ~~is awesome~~ is ^{awesome} -

Precision = $\frac{\# \text{relevant retrieved}}{\# \text{total retrieved}}$

Recall: $\frac{\# \text{relevant retrieved}}{\# \text{relevant total}}$

3) teleporting prob α

$$\sum_i \pi_i = 1 \quad i = 1, 2, 3, 4$$

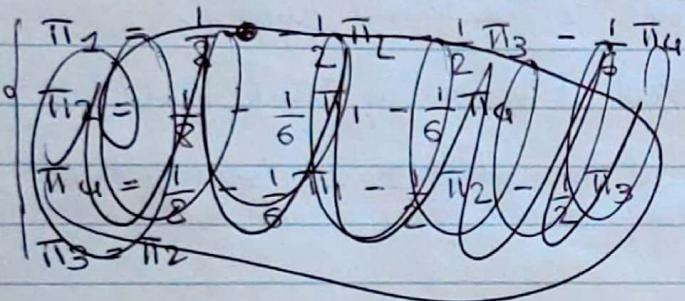
$$\pi_2 = \pi_3 \text{ symmetry}$$

$$\pi_1 = \frac{\alpha}{4} + \cancel{\frac{(1-\alpha)}{3}\pi_2} + \cancel{\frac{(1-\alpha)}{3}\pi_3} + \frac{(1-\alpha)}{3}\pi_4$$

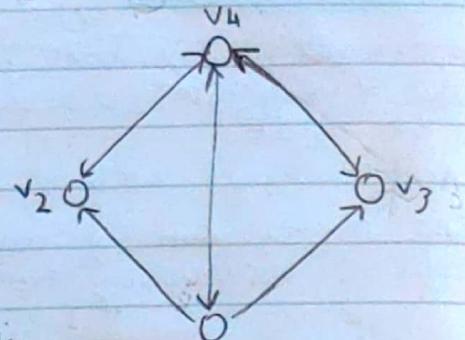
$$\pi_2 = \frac{\alpha}{4} + \frac{(1-\alpha)}{3}\pi_1 + \frac{(1-\alpha)}{3}\pi_4 \quad \begin{matrix} \text{numero di} \\ \text{archi uscenti dal} \\ \text{vicino} \end{matrix}$$

$$\pi_3 = \frac{\alpha}{4} + \frac{(1-\alpha)}{3}\pi_1 + \frac{(1-\alpha)}{3}\pi_4$$

$$(b) \alpha = \frac{1}{2}$$



modo vicini di
Modo vicini di
entranti
usciti



$$\left\{ \begin{array}{l} \pi_1 = \frac{1}{8} + \frac{1}{6}\pi_4 \\ \pi_2 = \frac{1}{8} + \frac{1}{6}\pi_1 + \frac{1}{6}\pi_4 \\ \pi_3 = \frac{1}{8} + \frac{1}{6}\pi_1 + \frac{1}{2}\pi_2 + \frac{1}{2}\pi_3 \\ \pi_4 = \frac{1}{4} + \frac{1}{6}\pi_1 + \frac{1}{6}\pi_4 \end{array} \right.$$

$$\pi_2 = \frac{1}{8} + \frac{1}{6}\left(\frac{1}{8} + \frac{1}{6}\pi_4\right) + \frac{1}{6}\pi_4$$

$$\pi_2 = \frac{1}{8} + \frac{1}{48} + \frac{1}{36}\pi_4 + \frac{1}{6}\pi_4 \rightarrow \frac{7}{48} + \frac{7}{36}\pi_4$$

$$\pi_4 = \frac{1}{8} + \frac{1}{6}\left(\frac{1}{8} + \frac{1}{6}\pi_4\right) + \pi_2 = \frac{1}{8} + \frac{1}{6}\left(\frac{1}{8} + \frac{1}{6}\pi_4\right) + \frac{7}{48} + \frac{7}{36}\pi_4$$

$$\text{Fraz } \frac{36-7-1}{36} \pi_4 = \frac{14}{48} \rightarrow \pi_4 = \frac{14}{48} \cdot \frac{36}{28} = \frac{3}{8}$$

$$\pi_2 = \frac{7}{48} + \frac{7}{36} \cdot \frac{3}{8} = \frac{14+7}{96} = \frac{21}{96} = \frac{7}{32} = \pi_3$$

$$\pi_1 = \frac{1}{8} + \frac{1}{6} \cdot \frac{3}{8} = \frac{9}{48} = \frac{3}{16}$$

$$\sum \pi_i = 1 \Rightarrow \frac{3}{16} + \frac{7}{32} + \frac{7}{32} + \frac{3}{8} = \frac{32}{32} = 1$$

OK

c) prove that for any graph the page rank of each node is at least $\frac{\alpha}{N}$

$$\forall v: \pi_v \geq \frac{\alpha}{N}$$

$$\pi_v = \frac{\alpha}{N} + (1-\alpha) \sum_{j \rightarrow v} \frac{\pi_j}{d_j}$$

where d_j is the number of out links

teleport following link \rightarrow always > 0 because the sum of rank in a graph = 1

whereas ~~and~~ $1 < \alpha < 0$ and $(1-\alpha) \sum \frac{\pi_j}{d_j} > 0$, $\Rightarrow \pi_v \geq \frac{\alpha}{N}$ sempre.

4)

NNN e K-means, | NNN con
inverted ~~list~~ list)

2) Show that for normalized vectors, Euclidean distance gives the same proximity ordering as the cosine measure

A)

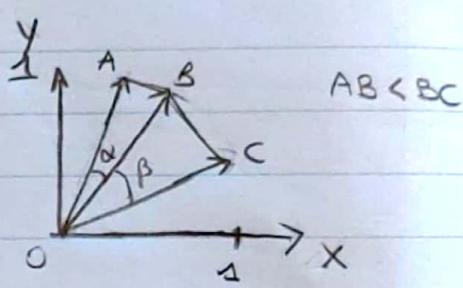
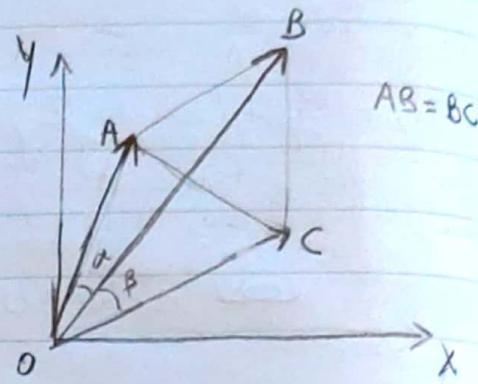
Normalized vectors have module = 1.

Let's consider this example, where $\alpha < \beta$,

~~this means that we know that OA is better than OC more to OB because cosine is greater, so is~~ ~~more similar than OC.~~

However, the triangle ABC is equilateral, so the euclidean distance in this case would lead to misleading results because it will consider OA and OC equidistant from OB.

To solve this problem we normalize vectors, making them of module = 1. In this way the angle will remain unchanged, while the euclidean distance will be coherent with the cosine similarity because it will show that the two vector with smaller angle ~~are~~ have smaller Euclidean distance. (same proximity ordering)



③ Is this true for non normalized vectors?

The previous statement is not true for non normalized vectors and the proof is in the previous exercise

1) for which queries skip pointers are ~~not~~ useful and for which are useful. Briefly explain your answer.

(a) $x \text{ AND } y$, where x frequent term and y rare
in this case it could be useful,

ex $x: 1, 2, 5, 15, 19, 20, 32, 35$
 $y: 19$

facile pochi passaggi

(b) $x \text{ AND } y$ where both x and y are frequent
USEFUL because it save times

(c) $x \text{ OR } y$
(d) $x \text{ OR } y$ } USELESS because in OR we have to ~~sum~~ all elements of both lists.

(B) pseudocode merging two postings lists term1 AND term2
using Skip Pointers

MERGE SKIPPOINTERS (p_1, p_2)

Answer $\leftarrow ()$

while $p_1 \neq \text{NULL}$ and $p_2 \neq \text{NULL}$

if $\text{docID}(p_1) = \text{docID}(p_2)$

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

$p_1 \leftarrow \text{next}(p_1)$

$p_2 \leftarrow \text{next}(p_2)$

if $\text{docID}(p_1) < \text{docID}(p_2)$

if hasSkip(p_1) and $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$

while hasSkip(p_1) and $\text{docID}(\text{skip}(p_1)) \leq \text{docID}(p_2)$

$p_1 \leftarrow \text{skip}(p_1)$

else $p_1 \leftarrow \text{next}(p_1)$

if $\text{docID}(p_2) < \text{docID}(p_1)$

if hasSkip(p_2) and $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$

while hasSkip(p_2) and $\text{docID}(\text{skip}(p_2)) \leq \text{docID}(p_1)$

$p_2 \leftarrow \text{skip}(p_2)$

else $p_2 \leftarrow \text{next}(p_2)$

return answer.

2) teleporting prob α

$$M = 4$$

$$\alpha = \frac{1}{2}$$

$$\sum_{i=0}^{M+1} \pi_i = 1 \quad i = 0 \dots M+1$$

per simmetria $\pi_1 = \pi_2 = \dots = \pi_M$

$$\pi_0 = \frac{\alpha}{M+2} + \frac{(1-\alpha)}{2} \left(\sum_{i=1}^M \pi_i \right) \quad \forall i \in \{1 \dots M\}$$

$$\pi_i = \frac{\alpha}{M+2} + \frac{(1-\alpha)}{2} \pi_0 + \frac{1-\alpha}{M} \pi_0$$

$$\pi_{M+1} = \frac{\alpha}{M+2} + \frac{(1-\alpha)}{2} \left(\sum_{i=1}^M \pi_i \right)$$

$$\pi_0 = \frac{\alpha}{6} + \frac{(1-\alpha)}{2} \pi_1 + \frac{(1-\alpha)}{2} \pi_2 + \frac{(1-\alpha)}{2} \pi_3 + \frac{(1-\alpha)}{2} \pi_4$$

$$\pi_5 = \frac{\alpha}{6} + \frac{1-\alpha}{2} \pi_1 + \frac{1-\alpha}{2} \pi_2 + \frac{1-\alpha}{2} \pi_3 + \frac{1-\alpha}{2} \pi_4$$

$$\pi_1 = \frac{\alpha}{6} + \frac{1-\alpha}{2} \pi_0$$

$$\pi_1 = \pi_2 = \pi_3 = \pi_4$$

$$\pi_0 = \frac{1}{12} + \frac{1}{4} (\pi_1 + \pi_2 + \pi_3 + \pi_4) = \frac{1}{12} + \frac{1}{4} \left(\frac{1}{12} + \frac{1}{8} \pi_0 \right)$$

$$\pi_0 - \frac{1}{8} \pi_0 = \frac{1}{12} + \frac{1}{48} \rightarrow \left(1 - \frac{1}{8}\right) \pi_0 = \frac{1}{12} + \frac{1}{48}$$

$$\pi_0 = \frac{1}{6} \cdot \frac{8}{7} = \frac{4}{21} \pi_0 = \frac{4}{21} \cdot \frac{16}{15} = \frac{64}{315} = 0,19$$

$$\begin{cases} \pi_0 = 0,19 \\ \pi_5 = 0,19 \end{cases}$$

$$\pi_0 = 0,19$$

$$\pi_5 = 0,19$$

$$\pi_1 = \frac{1}{12} + \frac{1}{8} (0,19) = \frac{1}{12} + 0,02375 = 0,107 = \pi_2 = \pi_3 = \pi_4$$

$$\sum = 0,19 + 0,19 + (0,107 \times 4) = 0,38 + 0,428 = 0,81$$

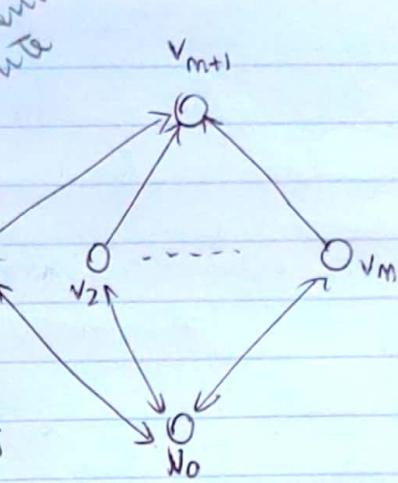
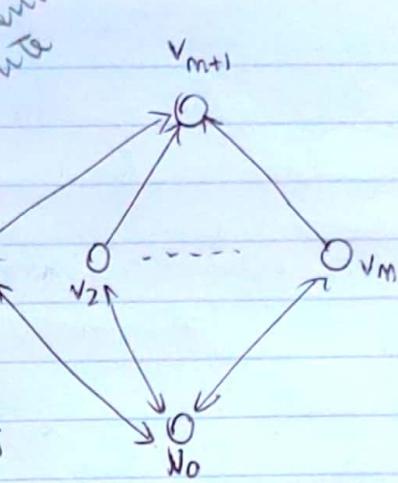
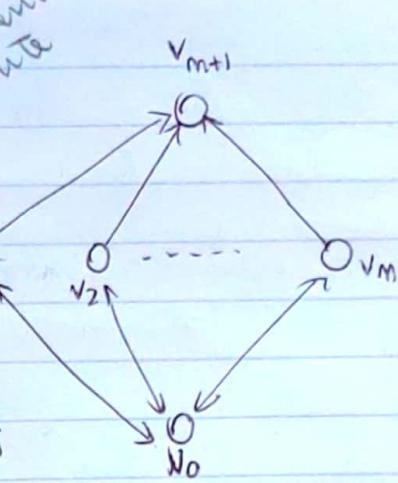
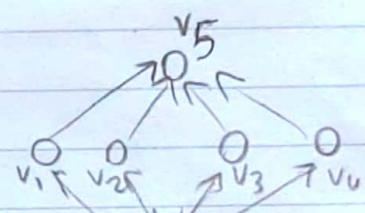
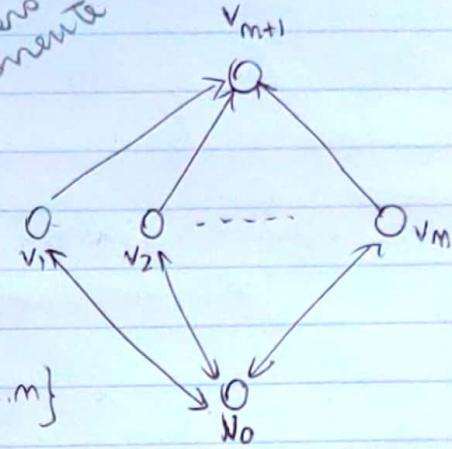
dove deve essere
giusto ma non
viene 1

IMPORTANTE PER COSINE SIM & EUCLIDEAN DIST.

N.B. ~~un~~ un documento viene rappresentato in uno spazio vettoriale come un vettore. I componenti del vettore sono gli elementi del dizionario che forma il documento (ovvero i termini senza ripetizioni).

Più è lungo il documento, più è lungo il vettore.

Modi archi uscenti
Modi vicini al modo
corrente, ovvero verso
nel modo corrente



Additivum β UD

Mantova Mem

}

e

⑤ interpolation 33%?

Vedo nella struttura a cosa corrisponde il 33% delle recall.

Siccome ho 10 elementi rilevanti totali, il 33% sarà intorno ai 3 elementi (33% di 10). Annulli sconsiglio la lista fino a che non ottengo al 33%.

es) $\frac{1}{10} = 0,1$ $\frac{2}{10} = 0,2$ $\frac{3}{10} = 0,3$ → siccome questo è il 30%, potrei usare questo

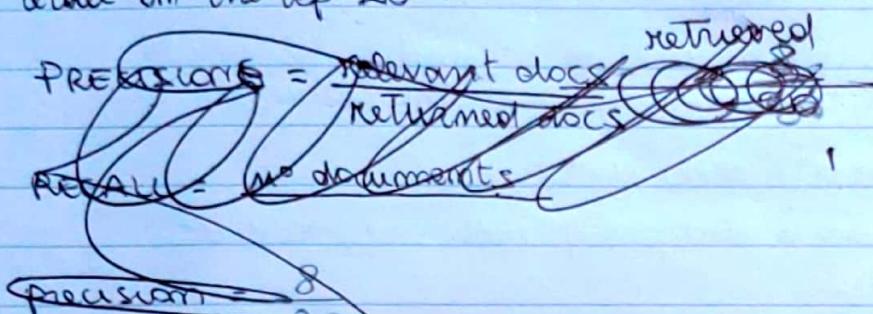
$$\text{precision} = \frac{3}{4} \approx 75\%$$

Annulli all 33% di recall corrisponde circa al 75% di precision

④ R e N R = relevant returned
 N = non relevant returned
 30 documents

RRNRN NRRNN NRNNN NRNNR NNNNN NRNRN

- ① precision of the system on the top 20?
 ② recall on the top 20?



$$\text{PRECISION} = \frac{\# \text{ relevant docs retrieved}}{\# \text{ total retrieved}} = \frac{8}{20} = \frac{2}{5} = 0,4$$

$$\text{RECALL} = \frac{\# \text{ relevant docs retrieved}}{\# \text{ total relevant docs}} = \frac{8}{10} = \frac{4}{5} = 0,8$$

- ③ F_1 measure on top 20?

$$F_1\text{-SCORE} = \frac{2 \cdot \text{PRECISION} \cdot \text{RECALL}}{\text{PRECISION} + \text{RECALL}} = \frac{2 \cdot \frac{2}{5} \cdot \frac{4}{5}}{\frac{2}{5} + \frac{4}{5}} = \frac{8}{15}$$

- ④ show the precision-recall curve

precision y recall x

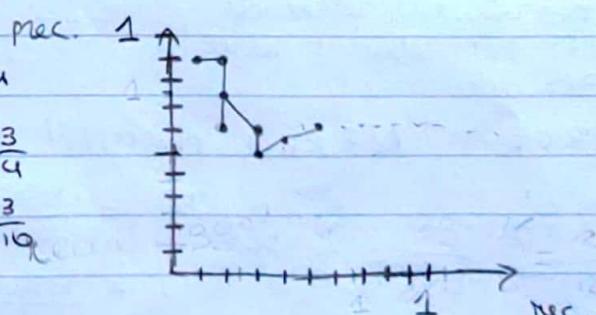
tacca pure prec. e recall sul primo
 elemento, poi sul primo e secondo, poi sul
 primo secondo e terzo... ecc.

$$pr_1 = \frac{1}{1} \quad pr_2 = \frac{2}{2} \quad pr_3 = \frac{2}{3} \quad pr_4 = \frac{3}{4}$$

$$rec_1 = \frac{1}{20} \quad rec_2 = \frac{2}{10} \quad rec_3 = \frac{2}{10} \quad rec_4 = \frac{3}{10}$$

$$pr_5 = \frac{3}{5} \quad pr_6 = \frac{3}{6} \quad pr_7 = \frac{4}{7} \quad pr_8 = \frac{5}{8}$$

$$rec_5 = \frac{3}{10} \quad rec_6 = \frac{3}{10} \quad rec_7 = \frac{4}{10} \quad rec_8 = \frac{5}{10}$$



- ⑤ interpolation 33% recall

Precision : $\frac{\text{relevant documents retrieved}}{\text{retrieved documents}}$

Recall = $\frac{\text{relevant retrieved documents}}{\text{total relevant}}$

Exam 20/07/2015

1) True or False?

(A) (a) Stemming always increases precision

FALSE: because stemming increases the retrieved set
increasing the number of relevant documents does not always increase
So stemming ~~decreases~~ precision (it can be unchanged or lower)

(b) Stemming increases recall

FALSE: stemming can increase or unchanged recall
because it can only increase the retrieved set

(c) Stemming reduces the size of the dictionary

TRUE: because we have a dictionary with only radices

(B) Are skip pointers useful for $X \text{ AND NOT } Y$

No because we have to consider all elements of X ~~and~~ that are not my
and discard those that are equals. There could be some useful
skip pointers ~~more~~ in general no.

(C) assume a biword index. Example with New York University which
is a false positive

Document: "I like New York. New University, here is great."

↳ they are both in the same document ~~s~~ but not
next to each other

2) $R = \text{relevant retrieved}$

$N = \text{non relevant retrieved}$

30 document

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

① precision. top 20? = $\frac{\text{relevant doc retrieved}}{\text{total retrieved}} = \frac{9}{20}$

② recall top 20? = $\frac{\text{relevant doc retrieved}}{\text{total rel. doc}} = \frac{9}{11}$

③ $F_1\text{-score} = 2 \cdot \frac{\text{prec. rec}}{\text{prec} + \text{rec}} = 2 \cdot \frac{\frac{9}{20} \cdot \frac{9}{11}}{\frac{9}{20} + \frac{9}{11}} =$

④ draw precision-recall curve

④ draw precision and recall curve.

precision on y axis

recall on x axis

from 0 to 1

to calculate it we calculate

precision and recall ~~for each~~ firstly for
the first element, then for the first two,
then the three and so on ...

$$\text{prec}_1 = \frac{0}{\cancel{2}} = 0$$

$$\text{prec}_{1,2} = \frac{1}{2}$$

$$\text{rec}_1 = \frac{0}{\cancel{10}} = 0$$

$$\text{rec}_{1,2} = \frac{1}{11} = 0,09$$

$$\text{prec}_{1,2,3} = \frac{2}{3} = 0,66$$

$$\text{rec}_{1,2,3} = \frac{2}{11} = 0,18$$

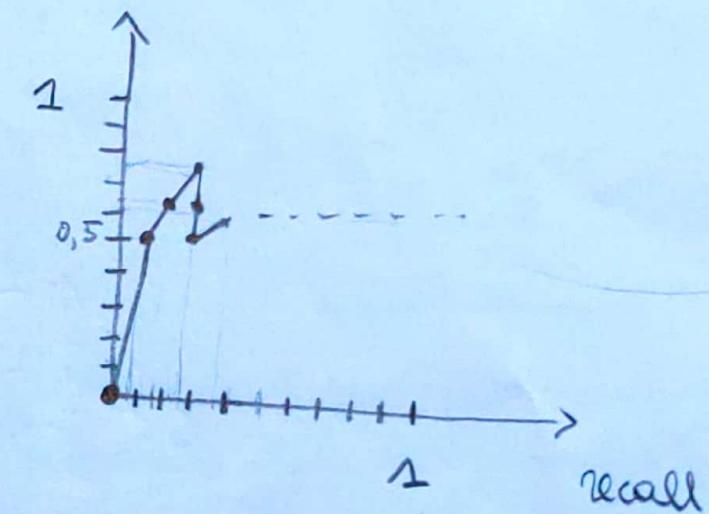
$$\text{prec}_7 = \frac{\cancel{4}}{7} = 0,57$$

$$\text{rec}_7 = \frac{4}{11} = 0,36$$

$$\text{prec}_{1,2,3,4} = \frac{3}{4} = 0,75$$

$$\text{rec}_{1,2,3,4} = \frac{3}{11} = 0,27$$

prec.



$$\text{prec}_5 = \frac{3}{5} = 0,6$$

$$\text{rec}_5 = \frac{3}{11} =$$

$$\text{prec}_6 = \frac{3}{6} = 0,5$$

$$\text{rec}_6 = \frac{3}{11} =$$

- 3) teleporting probability α , calculate personalized pagerank personalization vector $\{1, 0, 0, 0\}$

$$\sum \pi_i = 1 \quad i=1..4$$

$$\pi_2 = \pi_3 \text{ per simmetria}$$

$$\pi_2 = \frac{\alpha}{4} + \frac{(1-\alpha)}{3} \pi_4 \quad \begin{array}{l} \text{e viene è quello utente} \\ \text{nel modo corrente} \\ \text{numero di modi uscire dal vicino} \end{array}$$

$$\pi_2 = \frac{\alpha}{4} + \frac{(1-\alpha)}{3} \pi_4 + \frac{(1-\alpha)}{3} \pi_1$$

$$\pi_u = \frac{\alpha}{4} + \frac{(1-\alpha)}{3} \pi_1 + \frac{(1-\alpha)}{3} \pi_2 + \frac{(1-\alpha)}{2} \pi_3$$

② $p_1 = \{1, 0, 0, 0\}$

$$\pi_1 = \frac{\alpha}{4} + \frac{(1-\alpha)}{3} \pi_4$$

$$\pi_2 = \frac{(1-\alpha)}{3} \pi_4 + \frac{(1-\alpha)}{3} \pi_1 = \pi_3$$

$$\pi_4 = \frac{(1-\alpha)}{3} \pi_1 + (1-\alpha)(\pi_2 + \pi_3)$$

③ $p_2 = \{0, 1, 0, 0\}$

$$\pi_1 = \frac{(1-\alpha)}{3} \pi_u$$

$$\pi_2 = \alpha + \frac{(1-\alpha)}{3} (\pi_u + \pi_1)$$

$$\pi_3 = \frac{(1-\alpha)}{3} (\pi_u + \pi_1)$$

$$\pi_u = \frac{(1-\alpha)}{3} \pi_1 + (1-\alpha)(\pi_1 + \pi_3)$$

- ④ explain in detail how to calculate the personalized pagerank with respect to the personalization vector $p_3 = \{0, 5, 0, 5, 0, 0\}$ without solving the system.

$$p_3 = \{\frac{1}{2}, \frac{1}{2}, 0, 0\}$$

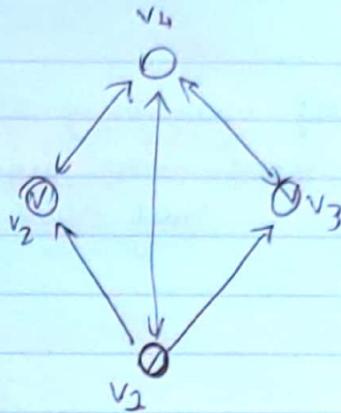
We know that the positions in p correspond to π_1 and π_2 , this means that the probability is shared between those two

$$\pi_1(p_3) = \pi_1(p_1) - \frac{\alpha}{2}$$

$$\pi_2(p_3) = \pi_2(p_2) - \frac{\alpha}{2}$$

$$\pi_3(p_3) = \pi_3(p_1) = \pi_3(p_2)$$

$$\pi_u(p_3) = \pi_u(p_1) = \pi_u(p_2)$$



4) $\text{tf} \times \text{idf}$ weighting scheme

- (a) two docs only frequent words (the, a, an, of etc) in common
- (b) " " that have no words in common
- (c) " " many rare words in common

1,3	0
0	0,3
1	0
1	0,3
0	0,3
1,3	0
2	0
1	0

512 256 128 64 32 16 8 4 2 1
0 0 0 0 0 0 0 0 0 0

EXAM 09/2017

2) ③ compress using byte encoding [5, 7, 18, 19, 28, 40, 52, 20]
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
5 2 11 1 9 12 12 28

5 → 00000101 ⇒ 10000101

2 → 000000010 ⇒ 10000010

11 → 00001011 ⇒ 10001011

1 → 00000001 ⇒ 10000001

9 → 00001001 ⇒ 10001001

12 → 00001100 ⇒ 10001100

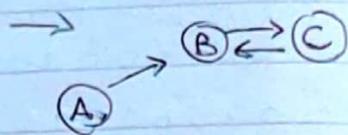
12 → " "

28 → 00011100 ⇒ ~~1~~0011100

3) Importance of teleporting probability with the convergence of pagerank.

Teleporting probability is important in pagerank because it avoids spider traps which are when the random walk can't escape from two nodes.

~~With teleporting probability~~, we can both follow a certain rule or we can jump to another node



~~$\pi_i = (1-\alpha) \sum_{j \rightarrow i} \frac{\pi_j}{d_j} + \alpha p_i$~~

where p_i is the personalized vector that shows with what probability you can jump on a node rather than another.

② equations for probability α

$$\sum \pi_i = 1 \quad i=1 \dots 4$$

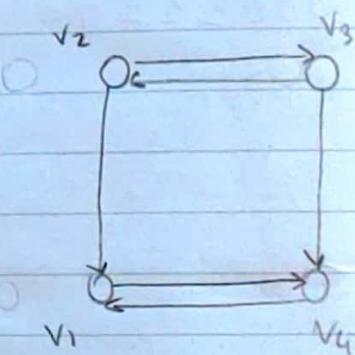
$$\pi_1 = \frac{\alpha}{4} + \frac{(1-\alpha)}{2} \pi_2 + \frac{(1-\alpha)}{2} \pi_4$$

$$\pi_2 = \frac{\alpha}{4} + \frac{(1-\alpha)}{2} \pi_3$$

$$\pi_3 = \frac{\alpha}{4} + \frac{(1-\alpha)}{2} \pi_2$$

$$\pi_4 = \frac{\alpha}{4} + \frac{(1-\alpha)}{2} \pi_1 + \frac{(1-\alpha)}{2} \pi_3$$

$$\begin{cases} \pi_2 = \pi_3 \\ \pi_1 = \pi_4 \end{cases} \quad \left. \begin{matrix} \text{symmetry} \\ \text{ } \end{matrix} \right\}$$



③ $\alpha = \frac{1}{2}$

$$\begin{cases} \pi_1 = \frac{1}{8} + \frac{1}{4} \pi_2 + \frac{1}{2} \pi_4 \\ \pi_2 = \frac{1}{8} + \frac{1}{4} \pi_3 \\ \pi_3 = \frac{1}{8} + \frac{1}{4} \pi_2 \\ \pi_4 = \frac{1}{8} + \frac{1}{2} \pi_1 + \frac{1}{2} \pi_3 \end{cases}$$

$$\begin{cases} \pi_2 = \frac{1}{8} + \frac{1}{4} \left(\frac{1}{8} + \frac{1}{4} \pi_2 \right) \rightarrow \frac{1}{8} + \frac{1}{32} + \frac{1}{16} \pi_2 \\ \pi_2 = \frac{8}{32} \cdot \frac{8}{153} = \frac{1}{6} \\ \pi_2 = \pi_3 = \frac{1}{6} \end{cases}$$

$$\pi_1 = \frac{1}{8} + \frac{1}{24} + \frac{1}{2} \pi_1 \rightarrow \frac{1}{2} \pi_1 = \frac{4}{24} \rightarrow \pi_1 = \frac{1}{3} = \pi_4$$

$$\sum = 1 \rightarrow \frac{1}{6} + \frac{1}{6} + \frac{1}{3} + \frac{1}{3} = \frac{1}{3} + \frac{2}{3} = 1 \quad \text{OR}$$

→ Unsupervised learning.

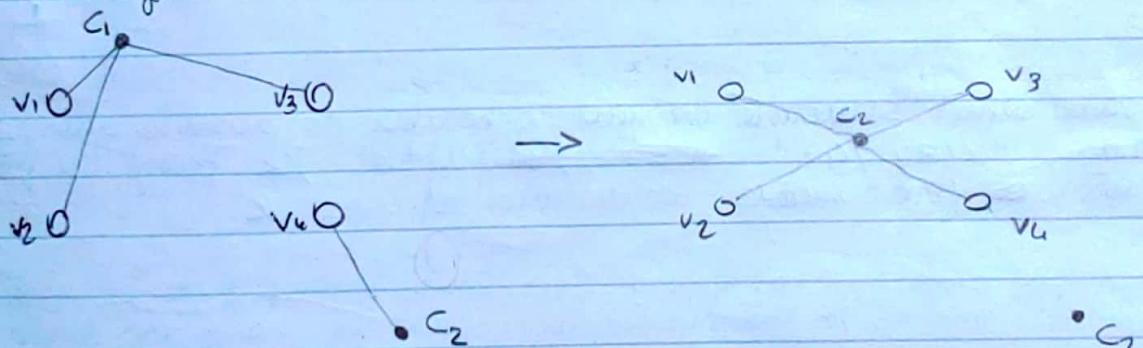
4) Explain how the K-MEANS works. Write the algorithm

① The K-means algorithm is used for classification of ~~problems~~ problems. Initially we take centroids ~~at~~ random and we take all points in a cluster that are nearest to the centroid. After that, the centroid is recalculated positioning it in a more central position of the cluster. This procedure is iterated until the convergence is reached and all elements of the clusters are correctly classified.

Algorithm

random positions for the centroid
until convergence: ~~do~~ → convergence when the old clusters are the same as the new.
do:
assign each point in the distribution to ~~the~~ the nearest centroid.
do the mean to recalculate the new centroid positions

② Show that if the initial assignment is unlucky the K-means solution might be bad.



③ Why K-means converges?

K-means converges when the centroids do not update their position so the last position is the less expensive one in terms of costs. ~~it converges because at every iteration the position of the centroids is adjusted until the center of the cluster is reached, that is when is less costly.~~

It converges because we minimize the distortion measure through EM algorithm.

minimize on distortion measure
(derivate su libri machine learning (4.24))

- 1) How to handle $X \text{ AND NOT } Y$? Why naive implementation expensive?
 2) Write algorithm merge.

To handle $X \text{ AND NOT } Y$ we take elements of X that are not in Y
 (so all elements of X except for those in Y too)

The naive implementation is expensive because we have to run all
 elements in the list X and see if there are elements in Y too and
 then discard them. We would ~~do~~ have to use two cycles immediately so
~~we can obtain this result so the cost will be $\Theta(n^2M)$~~

~~MERGE OPERATION (p_1, p_2) In fact we can if the smaller list is in
 answer of p_1 and the biggest one
 until $p_1 = p_2$ and~~

In a naive implementation, we can consider an adjacency matrix
 with one entry for each element of the dictionary of the documents.
 However, if we have big collections of documents (like thousands),
 the adjacency matrix will be too big to be computed.

So this is not the best way to proceed and we use posting lists
 and even skip pointers when possible to obtain ~~the best operation~~
~~operations has~~ less costly operations.

- b) describe what structure we need to be able to answer queries
 such as " $x/3y/4z$ " ~~with~~ with $/k$ being the proximity
 operator ~~is~~ that means at distance at most k .

To handle this type of query (proximity queries) we use the
 positional index. In the positional index we have the posting
 list storing, for each document, the position in which the
 terms appear.

ex $t_1 \rightarrow \text{docs: } 2, 42, 45, 80; \text{ doc2: } 7, 21 \dots$

We can use proximity search to find ~~the~~ matching
 documents

3)

Describe an external memory algorithm for the implementations of the power iteration method for pagerank computation

?

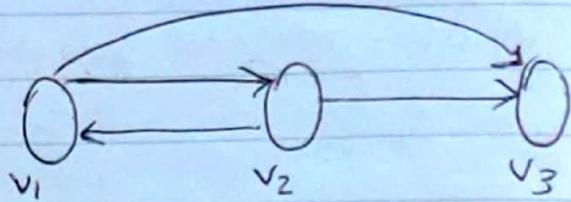
Le) teleporting probability α

$$\sum \pi_i = 1 \quad i = 1, 2, 3$$

$$\pi_1 = \frac{\alpha}{3} + \frac{(1-\alpha)}{2} \pi_2$$

$$\pi_2 = \frac{\alpha}{3} + \frac{(1-\alpha)}{2} \pi_1$$

$$\pi_3 = \frac{\alpha}{3} + \frac{(1-\alpha)}{2} \pi_2 + \frac{(1-\alpha)}{2} \pi_1$$



$$\alpha = \frac{1}{2}$$

$$\pi_1 = \frac{1}{6} + \frac{1}{4} \pi_2$$

$$\pi_2 = \frac{1}{6} + \frac{1}{4} \pi_1$$

$$\pi_3 = \frac{1}{6} + \frac{1}{4} \pi_2 + \frac{1}{4} \pi_1$$

$$\pi_1 = \frac{1}{6} + \frac{1}{4} \left(\frac{1}{6} + \frac{1}{4} \pi_1 \right) = \frac{1}{6} + \frac{1}{4} \cdot \frac{1}{6} + \frac{1}{4} \cdot \frac{1}{4} \pi_1$$

$$\frac{1}{8} \pi_1 = \frac{5}{24} \Rightarrow \pi_1 = \frac{5}{24} \cdot \frac{2}{1} = \frac{5}{12}$$

$$\pi_2 = \frac{1}{6} + \frac{1}{4} \left(\frac{1}{6} + \frac{1}{4} \cdot \frac{5}{12} \right) = \frac{1}{6} + \frac{1}{4} \cdot \frac{1}{6} + \frac{1}{4} \cdot \frac{5}{12}$$

$$\frac{15}{16} \pi_2 = \frac{5}{24} \Rightarrow \pi_2 = \frac{5}{24} \cdot \frac{16}{15} = \frac{2}{9}$$

$$\pi_3 = \frac{1}{6} + \frac{1}{4} \cdot \frac{2}{9} + \frac{1}{4} \cdot \frac{2}{9} = \frac{1}{6} + \frac{1}{18} + \frac{5}{18} = \frac{11}{18}$$

$$2 = \frac{1}{2}$$

$$\left\{ \begin{array}{l} \pi_1 = \frac{1}{6} + \frac{1}{4} \pi_2 \\ \pi_2 = \frac{1}{6} + \frac{1}{4} \pi_1 \\ \pi_3 = \frac{1}{6} + \frac{1}{4} \pi_2 + \frac{1}{4} \pi_1 \end{array} \right.$$

$$\left\{ \begin{array}{l} \pi_2 = \frac{1}{6} + \frac{1}{4} \left(\frac{1}{6} + \frac{1}{4} \pi_2 \right) = \frac{1}{6} + \frac{1}{24} + \frac{1}{16} \pi_2 \\ \frac{15}{16} \pi_2 = \frac{5}{24} \Rightarrow \pi_2 = \frac{5}{24} \cdot \frac{16}{15} = \frac{2}{9} \\ \pi_1 = \frac{1}{6} + \frac{1}{24} \cdot \frac{2}{9} = \frac{1}{6} + \frac{1}{18} = \frac{14}{18} = \frac{7}{9} \end{array} \right.$$

$$\pi_3 = \frac{1}{6} + \frac{1}{24} \cdot \frac{2}{9} + \frac{1}{24} \cdot \frac{2}{9} = \frac{1}{6} + \frac{1}{18} + \frac{5}{18} = \frac{11}{18}$$

$$\frac{2}{9} + \frac{2}{9} + \frac{5}{18} = \frac{4+4+5}{18} = \frac{16}{18}$$

$$\frac{16}{18}$$

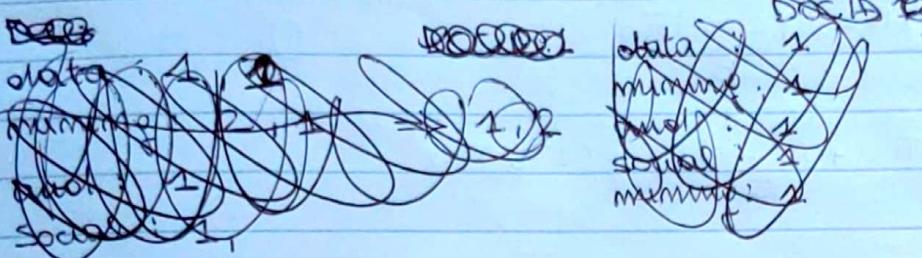
$$\frac{2}{9} + \frac{5}{18} + \frac{23}{84} = \frac{16+20+23}{84}$$

$$\frac{1}{6} + \frac{1}{4} \left(\frac{1}{6} + \frac{1}{4} \right)$$

1) three textual documents (D_1, D_2, D_3)

- D_1 : data mining and social mining
- D_2 : social network analysis
- D_3 : data mining

① write down the posting lists corresponding to the above document



term	DOC ID
data :	1
mining :	1
and :	1
social :	1
network :	1

term	DOC ID
social :	2
network :	2
analysis:	2

term	DOC ID
data :	3
mining:	3

term	doc frequency
data :	2
mining :	3
and :	1
social :	2
network :	1
analysis :	1

posting list	
data :	[1, 3]
mining :	[1, 3]
and :	[1]
social :	[1, 2]
network :	[2]
analysis :	[2]

② write down document frequency (df)

③ write down the term frequency (tf) for document D_1

④ write down the formula that relates the $tf \cdot idf$ weight for a term given its tf and df .

term	df
data	2
mining	2
and	1
social	2
network	1
analysis	1

$tf \cdot df$
2/5
2/5
1/5
2/5
pesante rispetto alle lunghezze dei documenti

$tf \cdot df$
1
2
1
1
pesante rispetto alle lunghezze dei documenti

↳ tf assente, non rispetto alle lunghezze dei documenti

the $tf \cdot idf$ formula is: $tf_{td} \times idf_t = (1 + \log_{10} tf_{td}) \log_{10} \frac{N}{df_t}$
where N is the number of documents in the collection.

teleportation probability α

3) personalized vectors:

$$p_1 = (0, 1, 0, 0)^T$$

$$p_2 = (0, 0, 1, 0)^T$$

4)

$$\sum \pi_i = 1 \quad i=1 \dots 4$$

$$\pi_1 = \alpha p_1 + \frac{1-\alpha}{3} \pi_u$$

$$\pi_2 = \alpha p_2 + \frac{1-\alpha}{3} \pi_1 + \frac{1-\alpha}{3} \pi_u$$

$$\pi_3 = \pi_2 \text{ symmetry}$$

$$\pi_4 = \alpha p_1 + \frac{1-\alpha}{3} \pi_1 + \frac{1-\alpha}{3} \pi_2 + \frac{1-\alpha}{3} \pi_3$$

$$p_1 = (0, 1, 0, 0)$$

$$\pi_1 = \frac{1-\alpha}{3} \pi_u$$

$$\pi_2 = \alpha + \frac{1-\alpha}{3} (\pi_1 + \pi_u)$$

$$\pi_3 = \cancel{\alpha} \frac{1-\alpha}{3} (\pi_1 + \pi_u)$$

$$\pi_u = \frac{1-\alpha}{3} (\pi_1 + (1-\alpha)(\pi_2 + \pi_3))$$

$$p_2 = (0, 0, 1, 0)$$

$$\pi_1 = \frac{1-\alpha}{3} \pi_u$$

$$\pi_2 = \frac{1-\alpha}{3} (\pi_1 + \pi_u)$$

$$\pi_3 = \cancel{\alpha} + \frac{1-\alpha}{3} (\pi_1 + \pi_u)$$

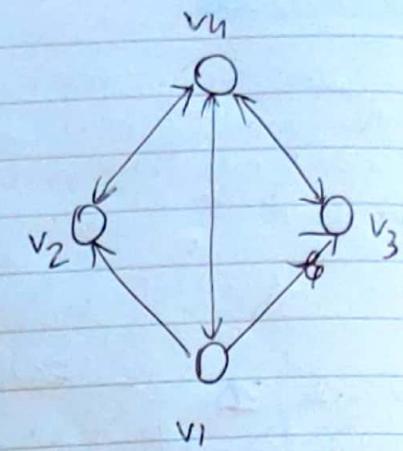
$$\pi_4 = \frac{1-\alpha}{3} \pi_1 + (1-\alpha)(\pi_2 + \pi_3)$$

2) Show how any personalized pagerank vector π corresponding to personalization vector $\alpha_1 p_1 + \alpha_2 p_2$ ($\alpha_1 + \alpha_2 = 1$) can be computed from π_1 and π_2 .
 If you can rigorously show why this works

$$\pi = \alpha_1 p_1 + \alpha_2 p_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \\ 0 \end{pmatrix} \alpha_1 + \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} \alpha_2 = (0 \alpha_1 \alpha_2 0)$$

$$\text{with } \alpha_1 + \alpha_2 = 1$$

$$\pi = \begin{pmatrix} \pi_1 \\ \pi_2 \\ \pi_3 \\ \pi_u \end{pmatrix} \Rightarrow \begin{aligned} \pi_1 &= (1-\alpha) \frac{\pi_u}{3} \\ \pi_2 &= \alpha_1 + (1-\alpha) \frac{(\pi_u + \pi_1)}{3} \\ \pi_3 &= \alpha_2 + (1-\alpha) \frac{\pi_u + \pi_1}{3} \\ \pi_u &= (1-\alpha) \left(\pi_2 + \pi_3 + \frac{3}{3} \pi_1 \right) \end{aligned}$$



If we have $\alpha_1, \alpha_2, \pi_u, \pi_i$, we can calculate π_u and π_i
 So we will calculate any personalised pagerank vector π .

~~We know that~~ We know that π_3 and π_2 have similar equations
 by symmetry.

$$\left. \begin{array}{l} \pi_2 = \alpha_1 + \underbrace{\frac{(1-\alpha)}{3}(\pi_u + \pi_i)}_{\text{fattore uguale}} \\ \pi_3 = \alpha_2 + \underbrace{\frac{(1-\alpha)}{3}(\pi_u + \pi_i)}_{\text{fattore uguale}} \end{array} \right\}$$

$$\pi_2 - \alpha_1 = \text{fattore}$$

$$\pi_3 - \alpha_2 = \text{fattore}$$

$$\Rightarrow \pi_2 - \alpha_1 = \pi_3 - \alpha_2$$

$$\pi_2 = \pi_3 - \alpha_2 + \alpha_1$$

$$\pi_3 = \alpha_2 - \alpha_1 + \pi_2$$

whereas $\alpha_1 + \alpha_2 = 1 \Rightarrow \pi_3 = \pi_2 + \alpha_2 - 1 + \alpha_2$

$$= \pi_2 - 1 + 2\alpha_2$$

$$\begin{aligned} \pi_u &= (1-\alpha) \left[\pi_2 + \underbrace{(\pi_2 - 1 + 2\alpha_2)}_{\pi_3} + \pi_{1/3} \right] \\ &= (1-\alpha) \left[2\pi_2 - 1 + 2\alpha_2 + \pi_{1/3} \right] \end{aligned}$$

So π_3 and $\pi_{2/1}$ can be expressed in terms of π_i and π_2
 therefore, we can calculate any $\pi = (\pi_1, \pi_2, \pi_3, \pi_u)$

EXAM 20-10-15

1) ① per le query $a \text{ AND NOT } b$ con x e y elementi in a e b rispettivamente, in quanto tempo $O(x+y)$ si può eseguire l'intersezione?

Sì perché abbiamo bisogno di leggere entrambi i set, nel caso peggiore, prendendo solo gli elementi in a ma non in b .

② stessa domanda per $a \text{ OR NOT } b$

$a \text{ OR NOT } b$ → significa che devo prendere tutto tranne gli elementi non in a ma in b .

Il costo in questo modo nel caso peggiore è $O(x+y)$

Algorithm $a \text{ OR NOT } b$ (p_1, p_2)

answer $\leftarrow ()$

while $p_1 \neq \text{NULL}$ and $p_2 \neq \text{NULL}$

if ~~$p_1 = p_2$~~ $\text{docID}(p_1) = \text{docID}(p_2)$

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

$p_1 \leftarrow \text{next}(p_1)$

$p_2 \leftarrow \text{next}(p_2)$

else if $\text{docID}(p_1) < \text{docID}(p_2)$

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

$p_1 \leftarrow \text{next}(p_1)$

else $\text{docID}(p_2) < \text{docID}(p_1)$

$p_2 \leftarrow \text{next}(p_2)$

while $p_1 \neq \text{NULL}$

if ~~$p_2 = \text{NULL}$~~

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

$p_1 \leftarrow \text{next}(p_1)$

return answer

Ex) $p_1 : 1, 3, 4, 7, 9, 12, 15$

$p_2 : 4, 5, 8, 9, 10$

answer : 1, 3, 4, 7, 9, 12, 15

prendo tutto tranne gli elementi che stanno solo in p_2



③ MERGE(p_1, p_2)

answer $\leftarrow ()$

while $p_1 \neq \text{NULL}$ and $p_2 \neq \text{NULL}$

if $\text{docID}(p_1) = \text{docID}(p_2)$

$p_1 \leftarrow \text{next}(p_1)$

$p_2 \leftarrow \text{next}(p_2)$

else if $\text{docID}(p_1) < \text{docID}(p_2)$

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

$p_1 \leftarrow \text{next}(p_1)$

else $\text{docID}(p_1) > \text{docID}(p_2)$

$p_2 \leftarrow \text{next}(p_2)$

while $p_1 \neq \text{NULL}$

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

$p_1 \leftarrow \text{next}(p_1)$

return answer

ex

$p_1: 1, 3, 5, 6, 15$

$p_2: 3, 6, 14, 18$

result: 1, 5, 15

a AND NOT b

elementi di a che non sono in b

SLIDE: CLASSIFICATION NAÏVE BAYES

	doc ID	words in document	in $C = \text{China}$?
training set	1	chim, Beij, chum	yes
	2	chum, chum, Shamp.	yes
	3	chim, Macao	yes
	4	Tokyo, Japan, chum	no
Test set	5	chum, chum, chum, Tokyo, Jap	?

FACCIO PRIMA SENZA SMOOTHING

quindi non dovo aggiungere 1 al numeratore e al numero di elementi del vocabolario al denominatore

~~$P(\text{China}/\text{yes}) = \frac{3}{4}$~~ $P(C = \text{yes}) = \frac{3}{4}$ $P(C = \text{no}) = \frac{1}{4}$

~~$P(\text{chum}/\text{yes}) = \frac{5}{8}$~~ $P(\text{Beij}/\text{yes}) = \frac{1}{8}$ $P(\text{Shamp}/\text{yes}) = \frac{1}{8}$

~~$P(\text{chum}/\text{no}) = 0$~~ $P(\text{Beij}/\text{no}) = 0$ $P(\text{Shamp}/\text{no}) = 0$

~~$P(\text{chum}/\text{no}) = \frac{1}{3}$~~

~~$P(\text{Macao}/\text{yes}) = \frac{1}{8}$~~

~~$P(\text{Macao}/\text{no}) = 0$~~

~~$P(\text{Tokyo}/\text{yes}) = 0$~~

~~$P(\text{Tokyo}/\text{no}) = \frac{1}{3}$~~

~~$P(\text{Jap}/\text{yes}) = 0$~~

~~$P(\text{Jap}/\text{no}) = \frac{1}{3}$~~

~~$P(\text{chum}/\text{yes}) = \frac{5}{8}$~~

~~$P(\text{chum}/\text{no}) = \frac{3}{8}$~~

adesso calcolo le probabilità del test set (ds) rispetto a yes o no

~~$P(\text{yes}/\text{ds}) = P(C = \text{yes}) \cdot P(\text{chum}/\text{yes})^3 \cdot P(\text{Tokyo}/\text{yes}) \cdot P(\text{Jap}/\text{yes}) =$~~

$$= \frac{3}{4} \cdot \frac{125}{512} \cdot 0 \cdot 0 = 0$$

~~$P(\text{no}/\text{ds}) = P(C = \text{no}) \cdot P(\text{chum}/\text{no})^3 \cdot P(\text{Tokyo}/\text{no}) \cdot P(\text{Jap}/\text{no}) =$~~

$$= \frac{1}{4} \cdot \frac{1}{27} \cdot \frac{1}{3} \cdot \frac{1}{3} = 0,001$$

quindi senza smoothing ds è classificato come no

lo smoothing serve per evitare che delle probabilità vengano 0 come succede per $P(\text{ds}/\text{yes})$ in questo caso.



• FACCIO con SHOOTING : $\frac{\text{numero di elementi/prob} + 1}{\text{numero totale/prob} + \text{numero totale del dominio}}$

$$P(\text{chim}/\text{yes}) = \frac{5+1}{8+6} = \frac{6}{14} = \frac{3}{7}$$

$$P(\text{chim}/\text{no}) = \frac{1+1}{3+6} = \frac{2}{9}$$

$$P(\text{falso}/\text{yes}) = \frac{0+1}{8+6} = \frac{1}{14}$$

$$P(\text{gap}/\text{no}) = \frac{1+1}{3+6} = \frac{2}{9}$$

$$P(\text{YES/ds}) = \frac{3}{4} \cdot \left(\frac{3}{7}\right)^3 \cdot \frac{1}{14} \cdot \frac{1}{14} = 0,0003$$

$$P(\text{NO/ds}) = \frac{1}{4} \cdot \left(\frac{2}{9}\right)^3 \cdot \frac{2}{9} \cdot \frac{2}{9} = 0,0001$$

In questo caso viene classificato come yes ~~False~~ ~~True~~

EXAM - 2015 - 16

2) RRNRR NRRNN N RNNNN NRNNNR NNNNN NRNRN

④ precision and recall 10

$$\text{precision} = \frac{\# \text{ relevant retrieved documents}}{\# \text{ retrieved documents}} = \frac{6}{10}$$

$$\text{recall} = \frac{\# \text{ relevant retrieved documents}}{\# \text{ relevant documents total}} = \frac{6}{11}$$

⑤ precision-recall curve

$$\text{prec}_1 = \frac{1}{1} = 1 \quad \text{prec}_{1,2} = \frac{2}{2} = 1$$

$$\text{rec}_1 = \frac{1}{11} = 0,09 \quad \text{rec}_{1,2} = \frac{2}{11} = 0,18$$

$$\text{prec}_{1,2,3} = \frac{2}{3} = 0,66$$

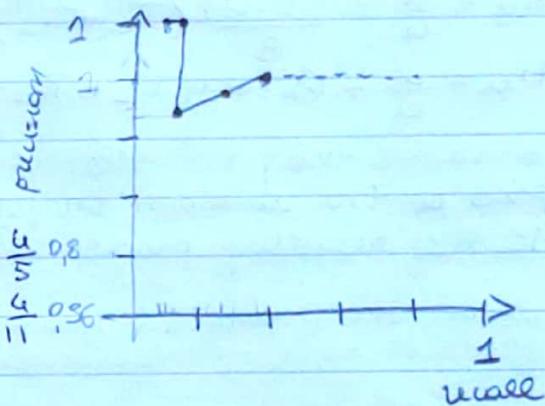
$$\text{prec}_4 = \frac{3}{4} = 0,75$$

$$\text{prec}_5 = \frac{4}{5} = 0,8$$

$$\text{rec}_{1,2,3} = \frac{2}{11}$$

$$\text{rec}_4 = \frac{3}{11} = 0,27$$

$$\text{rec}_5 = \frac{4}{11} = 0,36$$

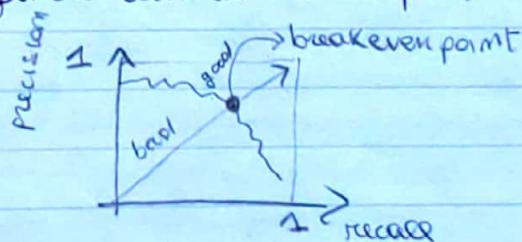


⑥ must there always be a break even point between precision and recall? Either show there must be or give a counter-example.

COSA E' IL BREAK EVEN POINT :

il break even point e' il punto che
mostra che prima ci sono solo risultati
peggiori e dopo ci sono solo
risultati migliori.

Si ottiene dall'intersezione tra la
bisettrice e la curva precision-recall.



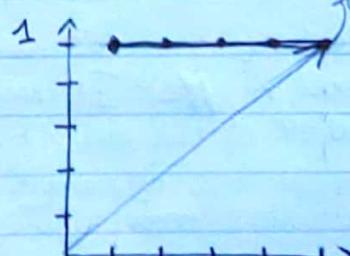
break even point

To solve the problem, consider a collection
RRRRR, then the graph will be this \Rightarrow

This means that the break even point
will coincide with the maximum point
~~possible~~ in the graph $(1, 1)$.

This means means that there is no
improvement in the result after it.

So there is no always be a break even point.



$$\begin{array}{lll} \text{prec}_1 = 1 & \text{prec}_2 = 1 & \text{prec}_3 = 1 \\ \text{rec}_1 = \frac{1}{5} & \text{rec}_2 = \frac{2}{5} & \text{rec}_3 = \frac{3}{5} \\ \text{prec}_4 = 1 & \text{prec}_5 = 1 & \text{prec}_6 = 1 \\ \text{rec}_4 = \frac{4}{5} & \text{rec}_5 = \frac{5}{5} = 1 & \text{rec}_6 = 1 \end{array}$$

3) α telep. prob

$$\sum_i \pi_i = 1 \quad i=1 \dots 4$$

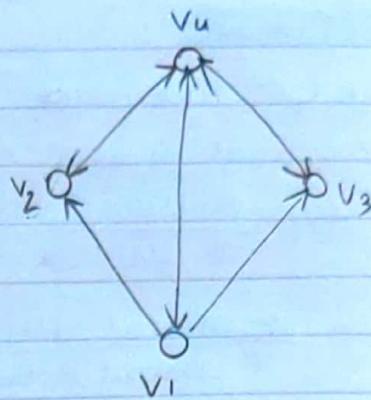
$$\pi_1 = \frac{\alpha}{4} + \frac{(1-\alpha)}{3} \pi_u$$

$\pi_2 = \pi_3$ by symmetry

$$\pi_2 = \frac{\alpha}{4} + (1-\alpha) \left(\frac{\pi_1}{3} + \frac{\pi_u}{3} \right)$$

$$\pi_3 = \frac{\alpha}{4} + \frac{(1-\alpha)}{3} (\pi_1 + \pi_u)$$

$$\pi_u = \frac{\alpha}{4} + (1-\alpha) \left(\pi_2 + \pi_3 + \frac{\pi_1}{3} \right)$$



- (b) assume that the teleport. prob is modified: at every step, with prob α the random surfer jumps to v_1 . Is the resulting process still a Markov chain? (justify the answer)

In a markov chain the probability to go on a certain node depends only on the current state

Quindi è markov chain o no?

Siccome è una modifica al personalization vector, che passa da $\{\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}\}$ a $\{1, 0, 0, 0\}$, quindi significa che la struttura del graph è invariata lo stessa e non ha violato nessuna proprietà delle markov chain. (perché le probabilità di finire su un nodo continuo a essere indipendente degli stati precedenti ma dipende da quello corrente).

i) Describe the assumptions of a Naive Bayes in the bag of words model

a) The bag of words model is one of the simplest language models used in NLP. It makes a unigram model of the text by keeping track of the number of occurrences of each word.

In the bag of words model you only take individual words into account and give each word a specific subjectivity score

The Naive Bayes model is also based on the bag of words model.

In Naive Bayes, the assumption is that the attribute values are conditionally independent.

b)

- ⑥ compute the coefficient of a boolean classifier without smoothing
- (a) browsing tiger safari. | apple
 - (b) afice video exam | not apple
 - (c) eum mountain osx. | apple
 - (d) mountain safari browsing tiger | not apple

Senza smoothing significa che non devo considerare i fattori di connivenza a numeratore e denominatore

$$P(\text{apple}) = \frac{1}{2} \quad P(\text{not apple}) = \frac{1}{2}$$

$$P(\text{browsing}/\text{apple}) = \frac{1}{6} \quad P(\text{tiger}/\text{apple}) = \frac{1}{6} \quad P(\text{safari}/\text{apple}) = \frac{1}{6}$$

$$P(\text{eum}/\text{apple}) = \frac{1}{6} \quad P(\text{mountain}/\text{apple}) = \frac{1}{6} \quad P(\text{osx}/\text{apple}) = \frac{1}{6}$$

$$P(\text{afice}/\text{not apple}) = \frac{1}{7} \quad P(\text{video}/\text{not apple}) = \frac{1}{7} \quad P(\text{eum}/\text{not apple}) = \frac{1}{7}$$

$$P(\text{mountain}/\text{not apple}) = \frac{1}{7} \quad P(\text{safari}/\text{not apple}) = \frac{1}{7} \quad P(\text{tiger}/\text{not apple}) = \frac{1}{7}$$

$$P(\text{browsing}/\text{not apple}) = \frac{1}{7}$$

- ⑦ classify the query document: eum mountam safari

$$P(\text{apple}/\text{testdocument}) = P(\text{apple}) \cdot P(\text{eum}/\text{apple}) \cdot P(\text{mount}/\text{apple}) \cdot$$

$$P(\text{safari}/\text{apple}) = \frac{1}{2} \cdot \frac{1}{6} \cdot \frac{1}{6} \cdot \frac{1}{6} = 0,027$$

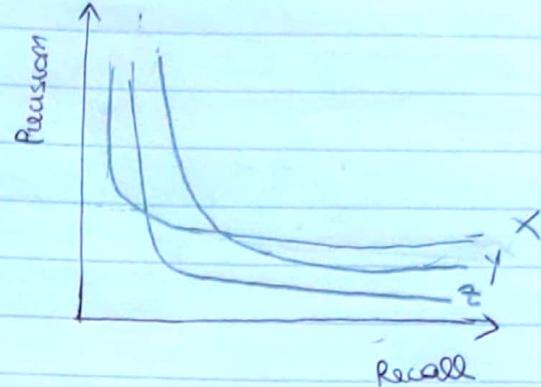
$$P(\text{not apple}/\text{testdocument}) = P(\text{not apple}) \cdot P(\text{eum}/\text{not apple}) \cdot P(\text{mount}/\text{not apple}) \cdot$$

$$P(\text{safari}/\text{not apple}) = \frac{1}{2} \cdot \frac{1}{7} \cdot \frac{1}{7} \cdot \frac{1}{7} = \frac{1}{49} = 0,0204$$

$$P(\text{apple}/\text{testdocument}) = \arg \max \Rightarrow 0,027, \text{ quindi la query viene classificata come apple}$$

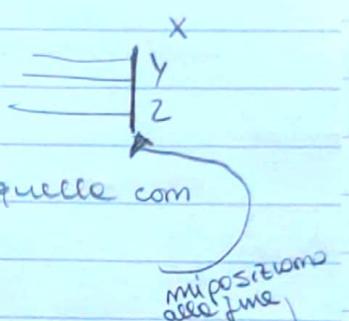
EXAM 21/04/2021

- 1) Describe to obtain the picture on the right comparing search engines
X (blue) Y (red) Z (green)



TUTTI I DOCUMENTI RILEVANTI \rightarrow RECALL.

Ci interessa le recall



- ② recall è uguale per tutti, ma bisogna prendere quelle con precision più alta per evitare noise.

Ora la scelta giusta è le X

Se vuole prendere all the relevant docs we care about recall.

- ③ see of the systems return 100% of relevant document because recall \rightarrow 1

The difference is that X retrieves all relevant docs but X has retrieved less ~~more~~ non-relevant documents perché la precision è più alta.

- ④ providing only relevant documents.
aumenta le precise, perché tra quelli rilevanti vuole maggioranza di relevant.

quando tutti hanno preso tutti i documenti rilevanti ma X ha preso anche pochi non rilevanti (meno noise).

[NOISE]: amount of not relevant documents

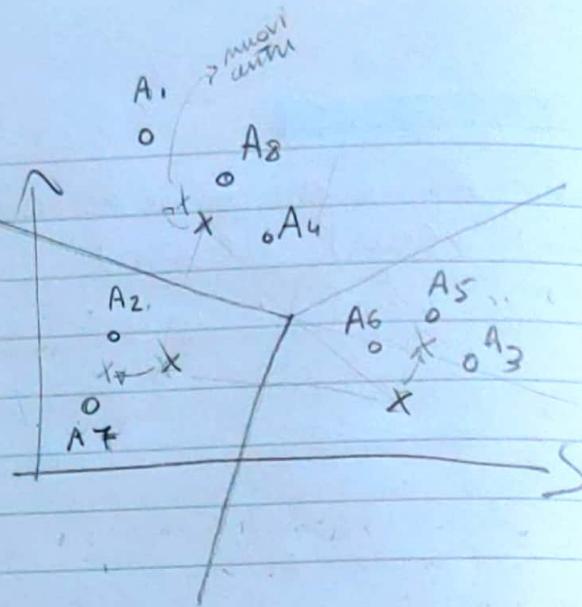
If you want to retrieve at most 50% of ~~documents~~ relevant docs Y \rightarrow perché ha precisione alta

If you want ~~at most 10%~~ more than 50%, best solution is X because precision is higher.

Mi interessa il rapporto tra doc rilevanti e doc totali trovati, quindi uso precisione (non mi interessa quanti rilevanti prendo ma mi interessa solo che siano rilevanti).

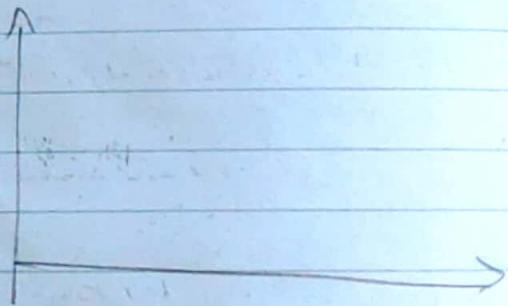
- 2) K-means per 3 classi
 passo per bisezione
 1 distanza perpendicolare
 alle distanze da 2 centroidi

Initially the centroids are given randomly, then we consider the distance between them and find the ~~perpendicolo~~ perpendicular bisect between them. Then the centroids are repositioned at the center of the new cluster and proceed until convergence.



- 22) A_1, A_2, A_4 rossi
 A_5, A_6 blu
 2 classi

(In questo caso uso kNN e classifico quelle che restano M base a quanti elementi di una certa classe ci sono nel vicinato dell'elemento)



- 3) D_1 = removed from domain Ontario to domain England
 D_2 = " " " England " "
 D_3 = " " England to ~~into~~ domain Ontario

- ① which docs have the same representation in the bag of words model? Bernoulli model

You represent docs as sets of words.

$$\tilde{D}_1 = \{\text{He, moved, from, domain} \dots\}$$

we do not count the occurrences

$$\tilde{D}_2 = \tilde{D}_3 \text{ non conto le occorrenze}$$

$$\tilde{D}_1 = \tilde{D}_2 = \tilde{D}_3$$

(perché con bernoulli voglio solo vedere se un dato appartiene a una certa classe, non quante volte compare)

EXAM 2012-09-10

② Vector Space Model = Multinomial

$\tilde{D}_1 = \{He, moved, from, London\}$ ma abbiamo London con 2



$\tilde{D}_1 = \{He, moved, from, London, \dots\}$ coordinate CONTA LE RIPETIZIONI

$\tilde{D}_1 = \{(He, 1), (moved, 1), (from, 1), (London, 2), \dots\} = \tilde{D}_2$

D_2 è uguale

mentre D_3 è diverso perché London compare una volta

③ suppose D_1, D_2, D_3 all belong to the same class "c"

~~P~~ $p(t/c)$ according to the Bernoulli model

$$p(t/c) = \frac{P(t) \cdot P(c/t)}{P(c)}$$

⇒ BAYES THEOREM

assumptions: all docs belong to same class, quindi $p(c) = 1$

$$p(c/t) = 1$$
 perche' ^{so} che il termine compare in tutti i documenti e non solo in $t \in c$

$$\text{So } p(t/c) = \frac{P(t) \cdot 1}{1} = P(t)$$

In the bag of word model you have to ~~count~~ e quindi per il numero di elementi del set,

$$p(t/c) = \frac{1}{f}$$

take the word and ~~the document~~

per il numero di elementi del set,
(perche' in Bernoulli non conta le occorrenze)

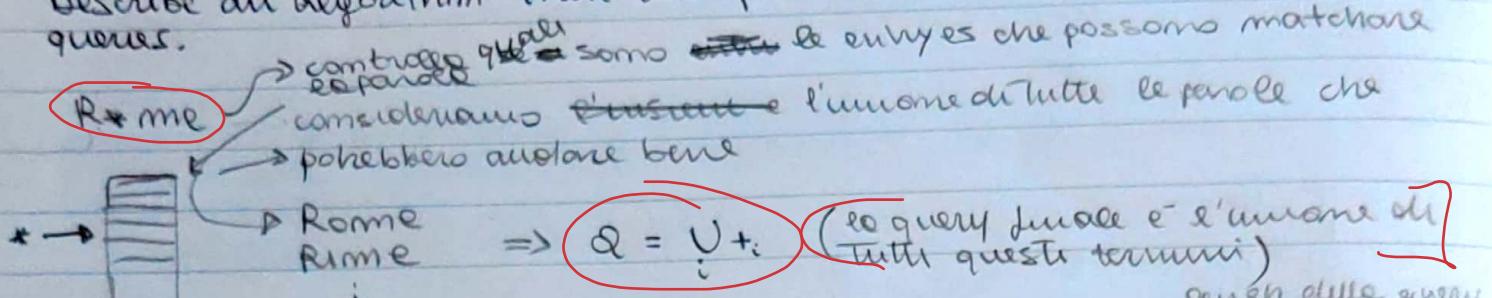
e assume ~~normalization~~

ESEMPIO PROF

Suppose that you want to answer queries containing wildcards

like: R*ame, R*I

Describe an algorithm that will process and solve those wildcards queries.



We have to compute the complexity



$$O(|V| \cdot |q|)$$

cardinalità di vettore

MENO DISPENDIOSO

Rom*

Suppose the vocabulary is not sorted, first of all we sort it to do by many search

$$O(|q| \log |V| + \# \text{ occurrences having Rom as prefix})$$

~~Suppose the vocabulary is not sorted~~

*imp → albero binario inverso
 $O(|q| \log |V| + \# \text{ occ})$

Rom* imp → just Rom*
S T them *imp

$$\Rightarrow Q = S \cap T$$

$$O(|q| \log |V| + \# \text{ occ})$$

Margine

SVM: la più grande distanza
tra il decisom boundary e il più
vicino punto dell'insieme

NAIVE BAYES

la classificazione e' troppo difficile da
risolvere se ci sono tanti attributi, quindi
uso la semplice versione del naive bayes