

# Locality sensitive hashing

• Binary vectors

$K, L \rightarrow$  # repetitions

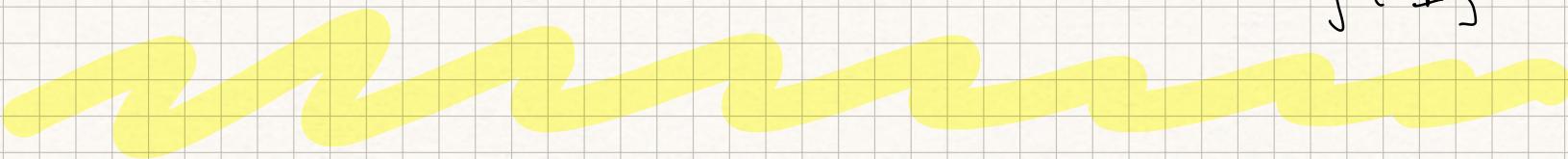
↳ projection

$\begin{matrix} I_1 \\ I_2 \\ \vdots \\ I_L \end{matrix}$

sets of  $K$  points each

$$h_{i,j}(P) = p_{i_1} p_{i_2} \dots p_{i_K}$$

$$\{i_j \in I_j\}$$

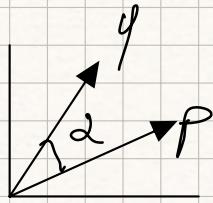


## COSINE SIMILARITY

Given vectors  $p, q \in \mathbb{R}_{\geq 0}^m$ , we define cosine similarity as

$$\cos \angle = \frac{\underline{p \cdot q}}{\|p\| \cdot \|q\|} \rightarrow p \cdot q = \sum_{i=1}^m p_i q_i$$

$\|q\| = \sqrt{\sum_{i=1}^m q_i^2}$

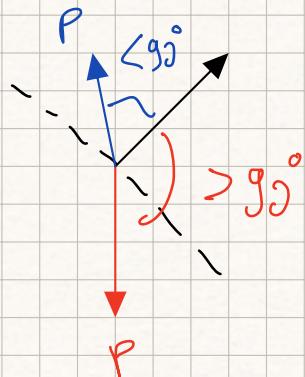


## BIG VECTORS

## SKETCHING FOR COSINE SIMILARITY

pick random  $\tau \in \mathbb{R}^m$ , define  $h_\tau(p) = \text{sign}(p \cdot \tau)$

$$P(h_\tau(p) = h_\tau(q)) = 1 - \frac{1}{\pi}$$



"sketch vector"

number vectors  $r_1, r_2, \dots, r_k$

$\forall p \in P: \text{SKETCH}(p) = \langle h_{r_1}(p), \dots, h_{r_k}(p) \rangle \in \{-1, +1\}^k$

$$\in \mathbb{R}_{\geq 0}^k$$

$$\frac{\text{SKETCH}(p) \setminus \text{sketch components}}{k} = \frac{\sqrt{P(h_{r_1}(p) > h_{r_2}(q))}}{\sqrt{P(h_{r_1}(p) > h_{r_2}(q))}} = 1 - \frac{2}{\pi}$$

## MIN HASHING

$A = \{1, 4, 6, 8, 9\} \rightarrow \text{SHINGLING + CAPPABIN}$

$B = \{2, 3, 4, 6, 7, 8\}$

$C = \{2, 5, 11, 12\}$

use MIN-MAPPING to estimate jaccard similarity by using  
a sketch of size 2, mapping  $V = \{0, \dots, 13\}$

permutation  $\pi_1(x) = 2x+1 \bmod 13$

permutation  $\pi_2(x) = x+2 \bmod 13$

$$S(A, B) = \frac{|S(A) \cap S(B)|}{|S(A) \cup S(B)|}$$

	$\pi_1$	$\pi_2$	
$A$	3, 9, 0, 4, 6	3, 6, 8, 10, 14	$\rightarrow \text{SKETCH}(A) = \{0, 3\}$
$B$	5, 7, 9, 0, 2, 4	4, 5, 6, 8, 9, 10	$\rightarrow \text{SKETCH}(B) = \{0, 4\}$

$$c \quad S_1, 11, 10, 12 \quad | \quad 4, 7, 0 \quad \rightarrow \text{SkorCor}(c) = \langle S_1 \rangle$$

$$\tilde{S}(A, B) = \frac{1}{2}$$

$$\tilde{S}(A, C) = 0/2 = 0$$

$$\tilde{S}(B, C) = 0/2 = 0$$

$$S(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{3}{8}$$

$T_1$  = white xmes

compute IL by encoding

$T_2$  = Xnew xmes happy

the porting list with  $\text{yopt}$

$T_3$  = happy white

$T_4$  = red red

happy  $\rightarrow 2, 3$

GAP encoding

red  $\rightarrow 4, 2$

white  $\rightarrow 1, 3$

Xnew  $\rightarrow 1, 2$

2, 1

4

1, 2

1, 1

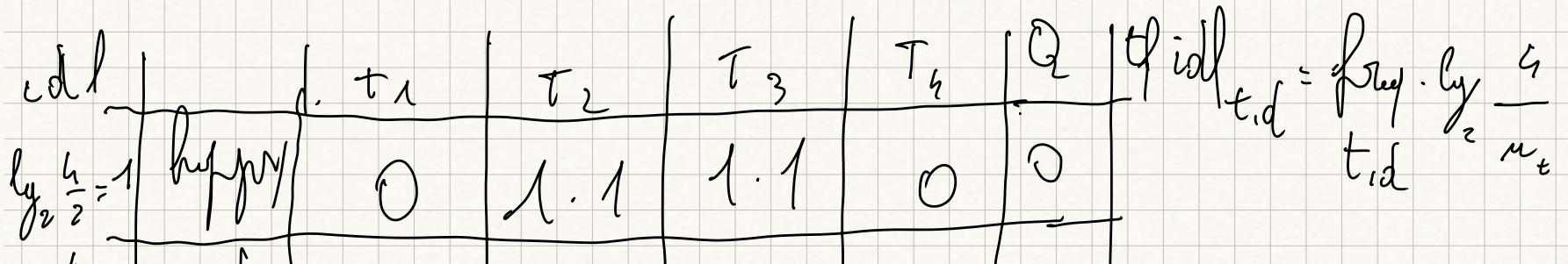
y of GAP

1-1 "0"  
2-1 "0" 2  
0101

00100

1010

11



$\text{by}_2 \frac{n}{1} = 2$	red	0	0	0	2 · 2	1 · 2
1	white	1 · 1	0	1 · 1	0	0
1	X <sub>mes</sub>	1 · 1	2 · 1	0	0	1 · 1

$$Q = \text{red} \cdot X_{\text{mes}}$$

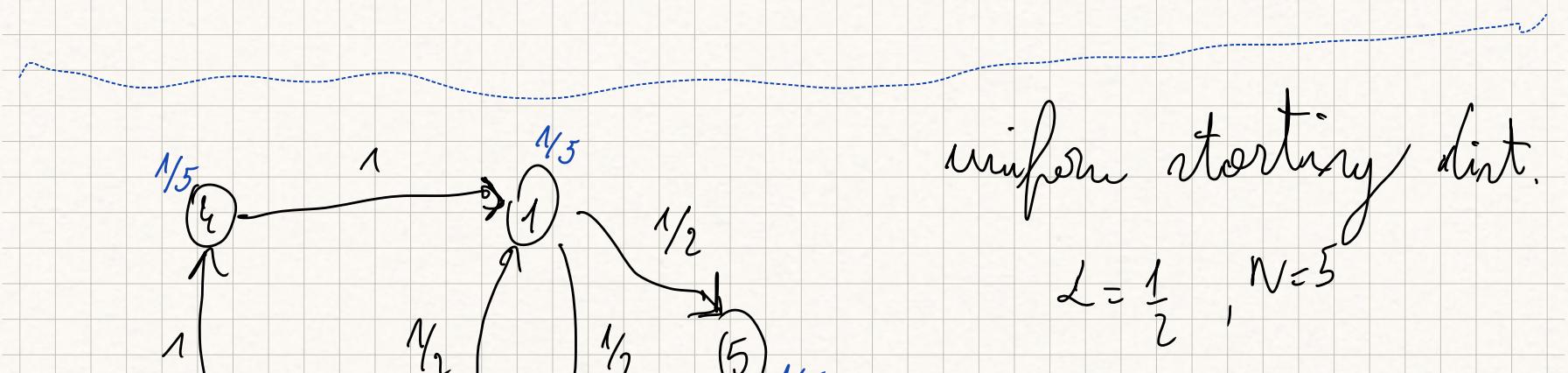
only regular products

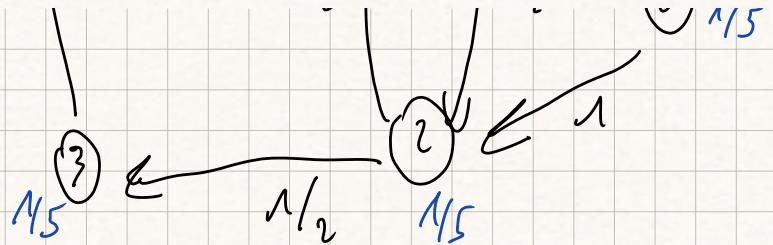
$$\cos(\tau_1, \alpha) = 1$$

$$\cos(\tau_2, \alpha) = 2$$

$$\cos(\tau_3, \alpha) = 0$$

$$\cos(\tau_u, \alpha) = 8 \rightarrow \text{most similar}$$





$$P_{r_2}[1] = \frac{1}{2} \left( P_{r_2}[4] \cdot 1 + P_{r_2}[2] \cdot \frac{1}{2} \right) + \frac{1}{2} \cdot \frac{1}{5}$$

$$P_{r_2}[2] = \frac{1}{2} \left( \frac{1}{5} \cdot \frac{1}{2} + \frac{1}{5} \cdot 1 \right) + \frac{1}{10} = \frac{1}{4}$$

$$= \frac{5}{20}$$

$$= \frac{1}{2} \left( \frac{1}{5} \cdot 1 + \frac{1}{5} \cdot \frac{1}{2} \right) + \frac{1}{2} \cdot \frac{1}{5} =$$

$$= \frac{1}{2} \cdot \frac{3}{10} + \frac{1}{10} = \frac{5}{20} = \frac{1}{4} = \frac{5}{20}$$

$$P_{r_2}[3] = \frac{1}{2} \left( \frac{1}{5} \cdot \frac{1}{2} \right) + \frac{1}{10} = \frac{3}{20}$$

$$P_{r_2}[4] = \frac{1}{2} \left( \frac{1}{5} \cdot 1 \right) + \frac{1}{10} = \frac{2}{10} = \frac{1}{5} = \frac{2}{20}$$

$$P_{r_2}[5] = \frac{1}{2} \left( \frac{1}{5} \cdot \frac{1}{2} \right) + \frac{1}{10} = \frac{3}{20}$$

Compute the Personalized PR with respect to  $\{1, 5\} \div 5$

$$Pr(u) = 2 \cdot \cancel{f(u)} + (1-2) \cdot \begin{cases} 1 & u \in S \\ 0 & u \notin S \end{cases}$$

$$Pr[1] = \frac{3}{20} + \frac{1}{2} \cdot \frac{1}{2} = \frac{8}{20}$$

$$Pr[2] = \frac{3}{20} + \frac{1}{2} \cdot 0 = \frac{3}{20}$$

$$Pr[3] = \frac{1}{20} + \frac{1}{2} \cdot 0 = \frac{1}{20}$$

$$Pr[4] = \frac{1}{10} + \frac{1}{2} \cdot 0 = \frac{1}{10}$$

$$Pr[5] = \frac{1}{20} + \frac{1}{2} \cdot \frac{1}{2} = \frac{6}{20}$$

- compute the similarity of node ② wrt node ① & node ⑤

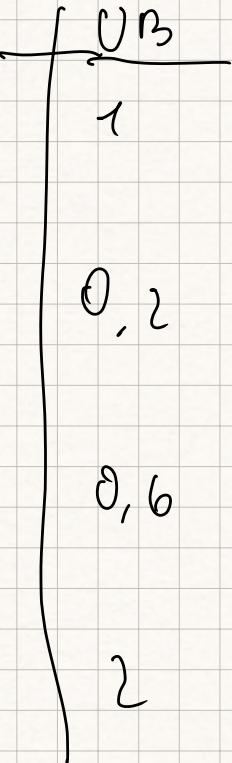
- Compute the PPR of G wrt S and then return the most similar node to S

$E_1 \rightarrow 1, 5, 6, 7, 8, 11$

$t_2 \rightarrow 4, 13, 15$

$t_3 \rightarrow 5, 6, 8, 9$

$t_4 \rightarrow 2, 3, 5, 6$



current  $\theta = 3.7$

- what's next pivot?
- is it more computed?

not

$t_1 \rightarrow 1, 5, 6 \dots$

$t_4 \rightarrow 2, 3, 5 \dots$

$$e.g. 3, 6 > 3.5 = \theta$$

$t_2 \rightarrow a, 13 \dots$	$ $	0,2
$t_3 \rightarrow 5, 6, 8 \dots$	$ $	0,5

pivot

5 is a potential candidate  
to enter in the heap,  
thus we compute its  
full score

D = {<sup>1</sup>shed, <sup>2</sup>mean, <sup>3</sup>stan, <sup>4</sup>too, <sup>5</sup>ome}

.) build a 2-gram index

.) show the candidate

dictionary strings for

$Q = \text{mean with } e=1$

\$ shed

\$ d → 1

d o → 1

o e → 1

\$ m → 2

m e → 2

{ mean

$$L - K \cdot e = 3 - 2 \cdot 1 = 1$$

$\emptyset m \rightarrow 2, 3$

$m \circ \rightarrow 2$

$\$ 0 \rightarrow 3$

$\emptyset t \rightarrow 3$

$t 0 \rightarrow 3$

$\emptyset m \rightarrow 3, 5$

$\$ 2 \rightarrow 4$

$z 0 \rightarrow 4$

$\emptyset o \rightarrow 4$

$\$ o \rightarrow 5$

$\$ m o m$

$\$ m \rightarrow 2$

$m o \rightarrow 2$

$\emptyset m \rightarrow 3, 5$



$D = \{ \text{abe}, \text{bae}, \text{bab} \}$

1

2

3

$Q = \text{bab}$

1-error match

candidates

$D_1 = D$

$D_2^- = \{ \text{be}[1], e\text{e}[1], \text{b}\text{b}[1]$   
 $\text{b}\text{e}[2], \text{e}\text{b}[2]$   
 $\text{abe}[3], \text{bbe}[3], \text{bae}[3], \text{bab}[3] \}$

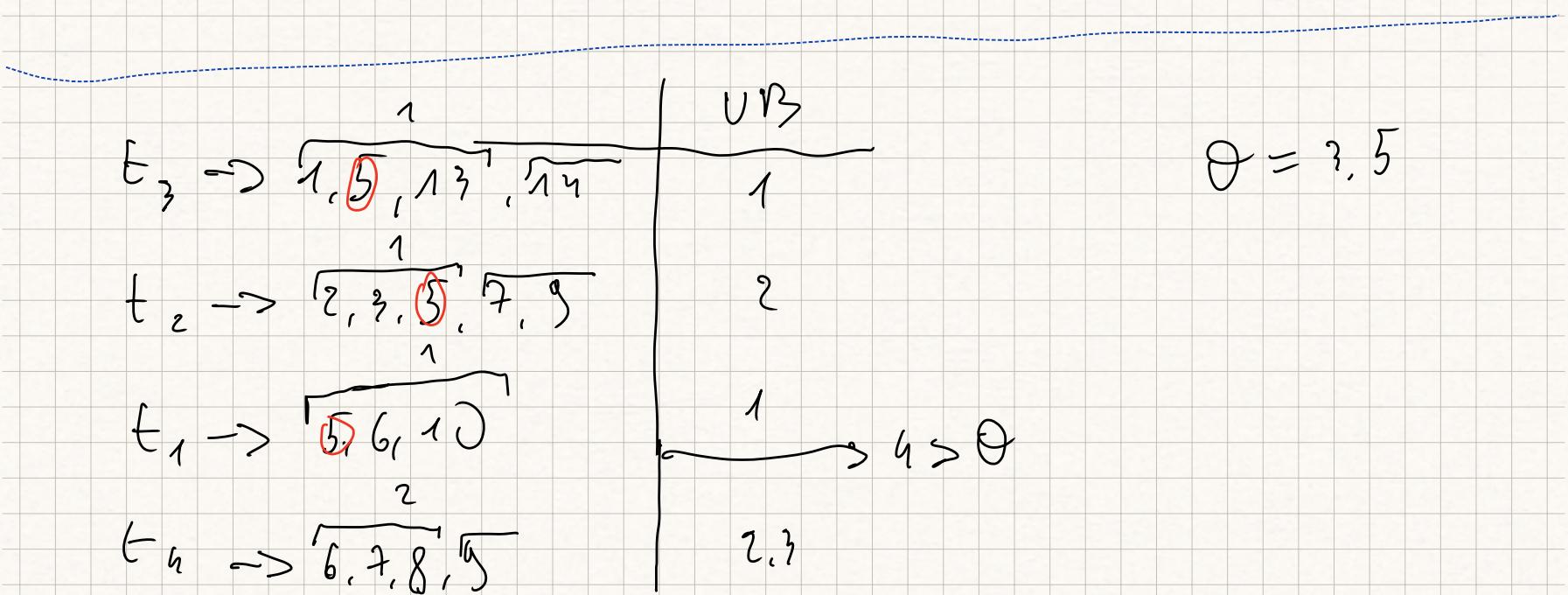
$Q = \text{bab}$

•  $Q$  in  $D_1$  = exact

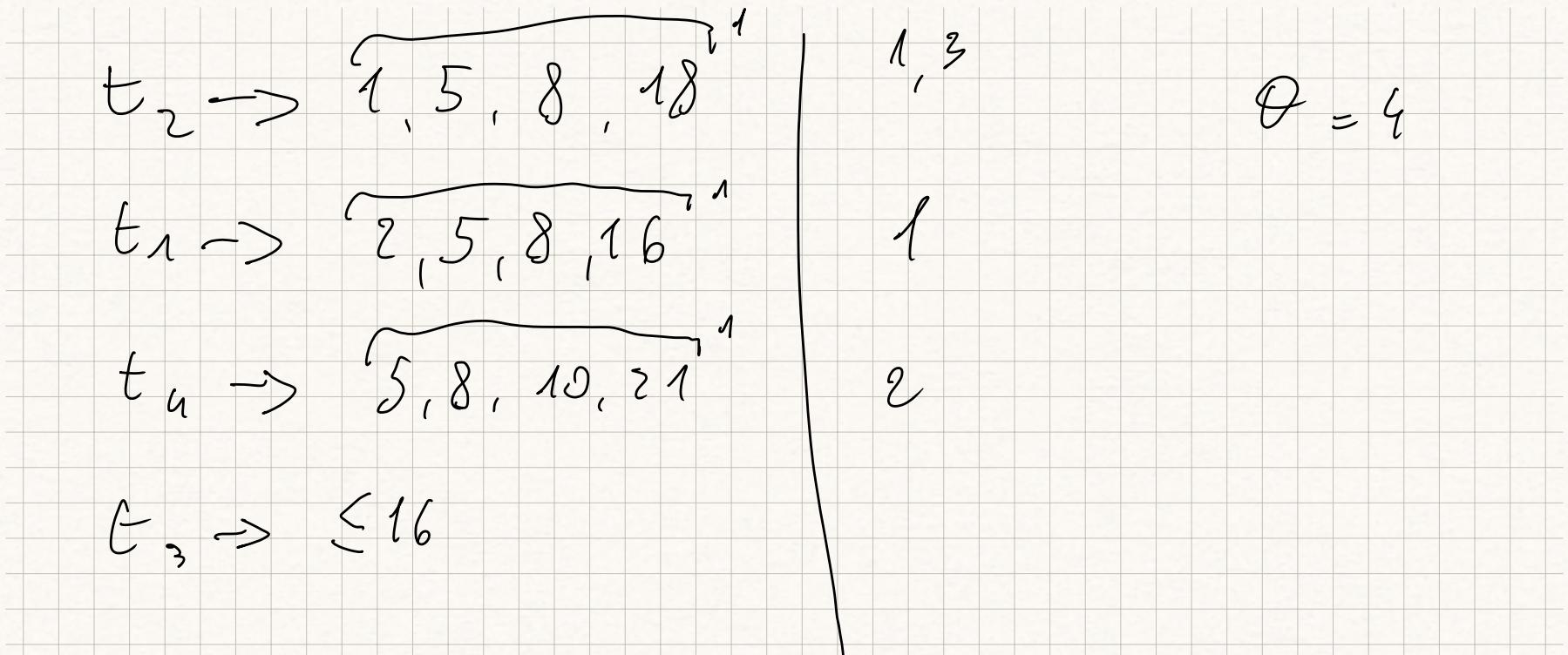
•  $Q^-$  in  $D_1 = \{ \text{ab}, \text{bb}, \text{ba} \}$  in  $D_1$

$\bar{Q}$  in  $D_2 = \{ab, bb, ba\}$  in  $D_2 = S_1, S_2 = ab, ba$

$Q$  in  $D_2 = S_3 = baba$



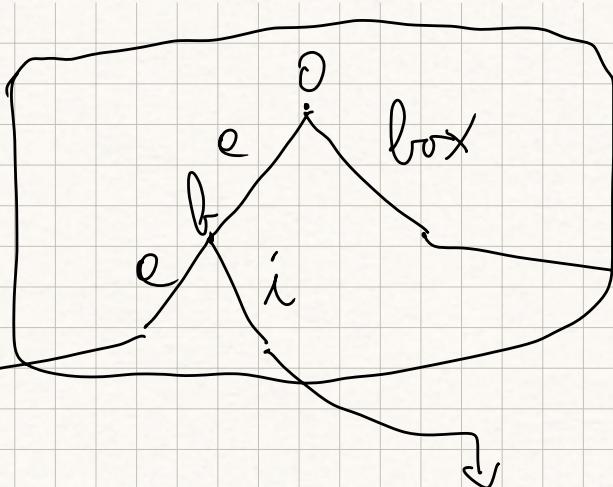
(5) is diverted because  $3 < 3, 5$



$D = \{ \text{ele, elte, elert, eli, boss, bot, box, brue, cat} \}$

- i) Design a tree level index which uses a block of 3 strings and front coding on each block

stored alphabetically  
sorted



internal  
memory

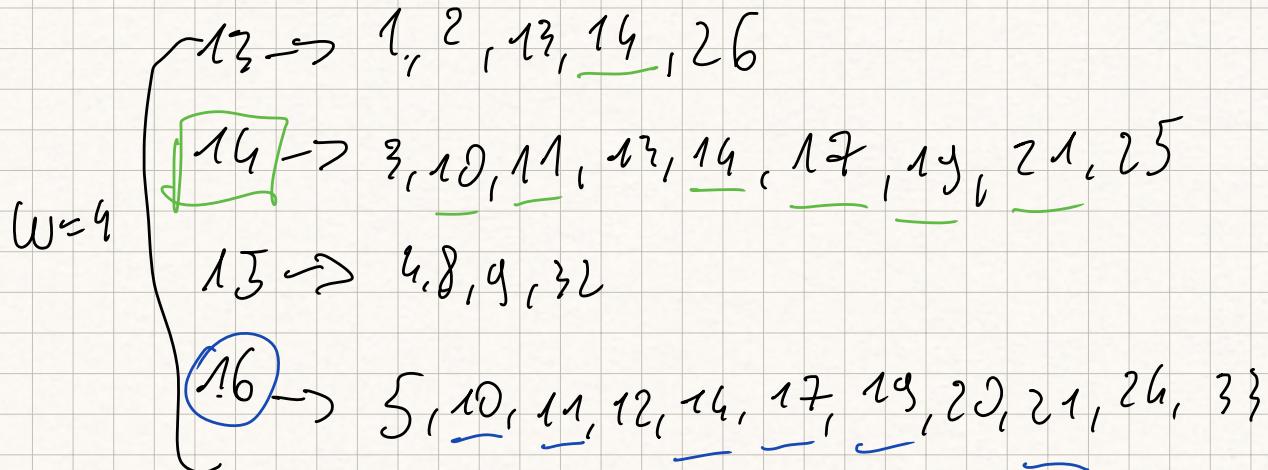
disk: [ $\langle 0, \text{ele} \rangle, \langle 3, \text{te} \rangle, \langle 3, \text{rt} \rangle$ ] [ $\langle 0, \text{obj} \rangle, \langle 0, \text{box} \rangle, \langle 2, \text{f} \rangle$ ]

[ $\langle 0, \text{box} \rangle, \langle 1, \text{nm} \rangle, \langle 0, \text{act} \rangle$ ]

2) Search for  $Q = \begin{cases} \text{ecc} & | \\ \text{ea} & | \\ \text{obj} & | \\ \text{t0} & : \end{cases}$  1 I/O  $\rightarrow$  access to disk  
exact  
no disk access  
: occurs to block three

3) Prefix search for  $Q = e$  : occur block 1 and 2

Web Graph



compress 16 in the best way

id	ref	outd	copy blocks	extra memory
14	--	--		
15	--	--		
16	2	11	0 1 1 0 1 1 1 1 0	5, 12, 20, 26, 33

↙ ↘  
g

# blocks	copy block
5	<del>0, 1, 2, 1, 1, 1</del> $\hookrightarrow 0, 0, 1, 0, 1$

$$S = (1, 4, 6, 10, 12, 14)$$

Elias - Fano

1 | 0 0 | 0 1

$$w=4$$

$$l = \lceil \log_2 \frac{m}{n} \rceil = \left\lceil \log_2 \frac{15}{6} \right\rceil = 2$$

4	0 1 0 0
6	0 1 1 0
10	1 0 1 0
12	1 1, 0 0
14	1 1, 1 0
	h : l

$$m = (\max S) + 1 = 15$$

$$L = \underbrace{0100101000010}_{l \cdot m}$$

$$H = 1011010110$$

$$\#1 = n$$

$$\#0 = 2^h$$

now let's do the reverse

$$L = 100010101000$$

$$\underbrace{00}_{00} \quad \underbrace{01}_{01} \quad \underbrace{10}_{10} \quad \underbrace{11}_{11}$$

$$H = 1011010110$$

	$h=2$	$l=2$
2	0 0	1 0
4	0 1	0 0

$$S ?$$

$h=2$  cut the #0 in H

$$m = 5$$

$$l = |L| = 2$$

6	01	10
10	10	10
12	11	00

$n$

$$S = 2, 4, 6, 10, 12$$

$$r(i) = \alpha \sum_{j \rightarrow i} r(j) + (1-\alpha) \frac{1}{N}$$

*Contribution depends on G's structure*

*Teleportation step*

$r(i)$   $i \in S$

Personalized / topic-based

$$r(i) =$$

for  $u$

$$\begin{cases} 1 & \text{if } u = i \\ (1-\alpha) \cdot \frac{1}{\text{volume of paths}} & \text{otherwise} \end{cases}$$

$r(i) = \text{similarity}$

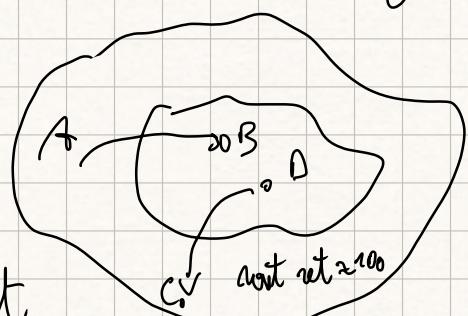
between nodes  $i$  and  $u$

volume of paths

Iteration  $u \sim i$

$H(\star)$  (query dependent)

- Compute "root ret" = list of pages answering  $Q$
- expand to base ret  $\Rightarrow$  root ret
- random walk over the base ret

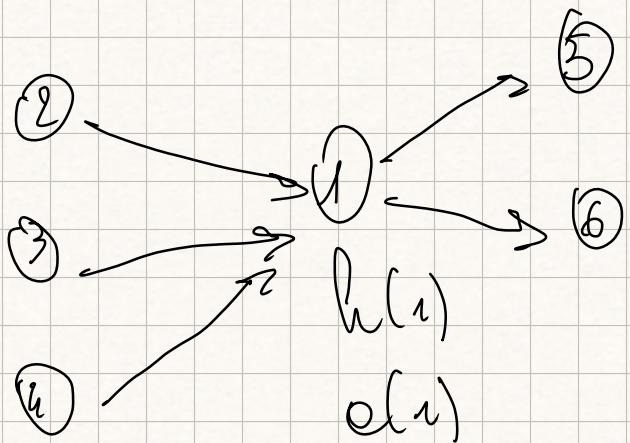


base net 1000

Authority score

hubness score

- good authority is a page pointed by good hub pages
- " hub page is a page that points to good authority pages  
noway



$$h(1) = \sqrt{w(1,5)} + \sqrt{w(1,6)}$$

$$h(1) = h(2) + h(3) + h(4)$$



$\vec{h}$  = vector of all hubness score

$\vec{e} = \vec{1}$   $\leftrightarrow$  authority score

eigenvalue

$$\lambda \cdot \vec{e} = (\vec{A}^T \cdot \vec{A}) \cdot \vec{e} \leftarrow \text{eigenvector}$$

$$\lambda \cdot \vec{h} = (\vec{A} \cdot \vec{A}^T) \vec{h}$$

$$\begin{cases} \vec{e} = \vec{A}^T \cdot \vec{h} = \vec{A}^T \cdot \vec{A} \cdot \vec{e} \\ \vec{h} = \vec{A} \cdot \vec{e} = \vec{A} \cdot \vec{A}^T \cdot \vec{h} \end{cases}$$

adjacency matrix  
of the base net

. TF

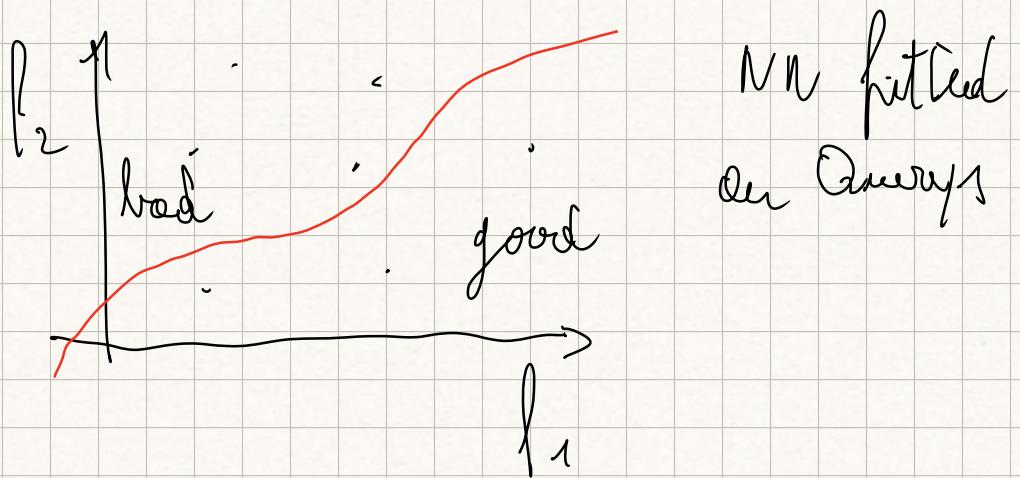
. Page Rank (HITS): ...

. Nervenre Q's terms in URL

. // " in titles of pages

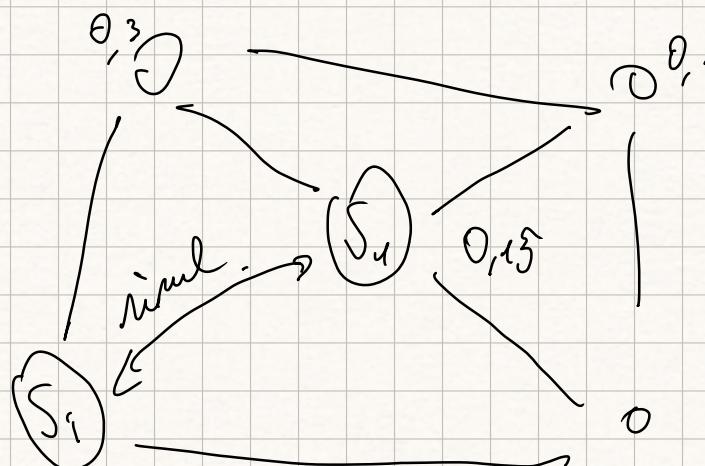
. // " in the first lines

. Terms of Q are close to each other



$S_i$  = sentence

$$\text{similarity}(S_i) = \sum_{w \in S_i} \frac{\text{TF-IDF}(w)}{|S_i|}$$



Text rank

$$\text{Sim}(S_i, S_j) =$$

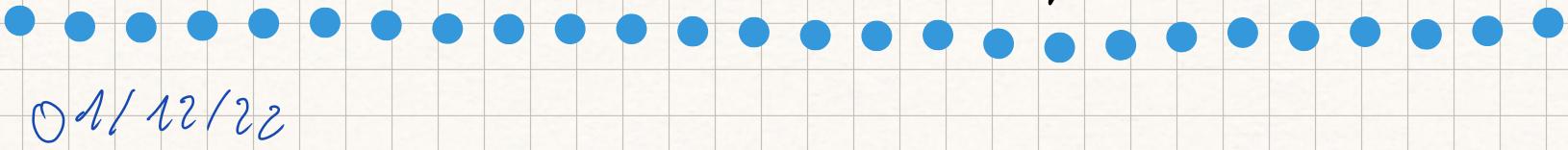
$$= \underbrace{(S_i \cap S_j)}_{\text{ly}(S_i) + \text{ly}(S_j)}$$

$$(1) \text{Sim}(S_i, S_j) = \text{cosine sim}$$

1-UV

$$TF-IDF \quad (S_i, S_j)$$

## ⑦ Pruning



04/12/22

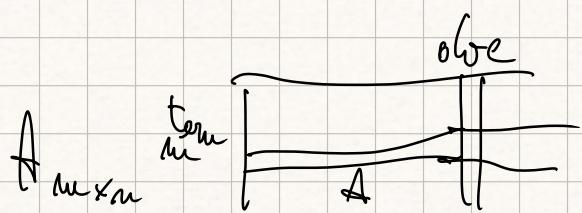
Vector Space model

m dim space  $\rightarrow$  k dim space

- Latent Semantic Indexing (LSI) - data dependent
- Random Projections - data independent (with high prob.)

preserve the distances b/w pairs of vectors in the reduced space

Singular-value decomposition



• define  $T = A \cdot A^T$  term-term matrix

$$T_{m \times m} = \begin{matrix} j \\ | \\ \vdots \\ i \\ | \\ s \end{matrix} = \begin{matrix} t_i \\ | \\ A \end{matrix} \cdot \begin{matrix} t_j \\ | \\ A^T \end{matrix}$$

$\Rightarrow \min(t_i, t_j)$

• Define  $D \geq A^T \cdot A$  a ~~be~~-olve matrix

$$D_{m \times m} = \begin{matrix} & \\ & \end{matrix} = m \begin{matrix} & \\ & D^T \\ & \end{matrix} \cdot \begin{matrix} & \\ & A \end{matrix} m$$

$\min(d_1, d_2)$

$$A = U \cdot \Sigma \cdot V^T$$

$U, V$  orthogonal,  $\Sigma$  diagonal  $\rightarrow$  singular values

latent concepts  
latent dimensions

$$\begin{matrix} & \\ & \end{matrix} = \begin{matrix} & \\ & \end{matrix} \cdot \begin{matrix} & \\ & \Sigma \end{matrix} \cdot \begin{matrix} & \\ & V \end{matrix}$$

concept term zeros u element concept

II eigvalues

decreasing value of S.V.

$U$  = matrix of  $n$  eigenvectors of  $T$

$V$  = " " of  $D$

per confrontare che due generiche colonne di  $V^T$

$$A = U \cdot \Sigma \cdot V^T$$

~~$$U^T \cdot A = U^T \cdot U \cdot \Sigma \cdot V^T$$~~

$$U^T \quad | \quad A \quad | \quad V^T$$

q<sub>new</sub>

$$q' = U^T \cdot q$$

q<sub>new</sub>

EJ:

	d <sub>1</sub>	d <sub>2</sub>	m
day	5	0	
cat	0	0	
PC	0	3	
m			

$$A = \begin{matrix} & d_1 & d_2 \\ \text{day} & 5 & 0 \\ \text{cat} & 0 & 0 \\ \text{PC} & 0 & 3 \end{matrix}$$

U

animal computer science

$k = 3$

$V^A$