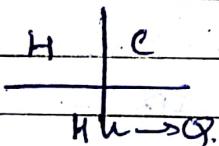


Web - Collection of docs.  
 Internet - Interconnected Network.  
 HTTP

WI → AI + Advanced IT.

Turing Test → In, AI, Turing Test is a method to determine whether a computer is Human Being.



Reverse Turing Test →

Computer asks question e.g. CAPTCHA

(1) Intelligent Web Page System

→ BI → CRM → Web Marketing / Services  
 (Cust. Rel. / Publishing)  
 (Marketly)

(2) Knowledge Network & Management

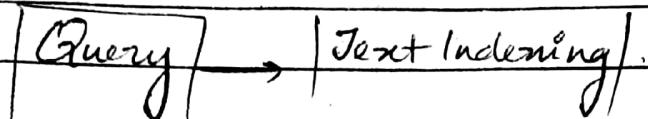
→ Electronic library  
 → Ontology and indexing  
 → Semantic Web.

(3) Web Mining

→ Multimedia mining → Web log mining.  
 → Text Mining, → Web log structure.  
 → Web based Ontology → Web warehousing

Index →

list of words / names / terms



- 1) Documents
- 2) Tokenizing. ~ Breaking Down Sentence.
- 3) Normalizing
- 4) Stemming. ~ Root Word.
- 5) Stop Words. ~ Common Word.

Serial no.	term	Doc ID
1	bag	$1 \rightarrow 2$
2	ball	$3 \rightarrow 12$
3	ball	$1 \rightarrow 12$

Class Index      ↙ list of docs.

```

create(D)      (populates list of documents)
for d in D
  (each document in D)
    for w in d
      (each word in every document d)
        i = index.lookup(w)
        (check whether word w is present in index)
        (if i is not present in index,
         lookup function returns -ve(i))
  
```

if  $i < 0$

$j = \text{index.add}(w)$

( $j$  will return to which position  
(or serial no.) the new word  $w$   
is added to)

$\text{index.append}(j, \text{do d.id});$

(add the document id to the  
list of doc id

4 cat 2

else

$\text{index.append(d.id)}$

D1:

In the end it  
doesn't even matter

D2:

It might matter  
who knows

1)

Tokenize

Terms	DocID
In	1
the	1
end	1
it	1
doesn't	1
even	1
matter	1

Terms      DocID

It	1
might	1
matter	1
who	1
knows.	1

2) Sort ~~Terms~~ & merge both lists.

Doesn't |  
 and |  
 even |  
 In |  
 It |  
 It 2  
 knows. 2.  
 matter. 1.

### Inverted Index:

Now, frequency comes into picture. Posting  
 Dictionary.

Sl. no.	term	Docfq.	DocID list
1	Doesn't	1	1
2	and	1	1
3	In	1	1
4	It	2	1 → 2
5	knows.	1	2

Document frequency  
 Not  
 Term frequency

'It' twice in doc1

It	1	1
----	---	---

'It' twice in doc1 & 'It' once in doc2.

It	2	1 → 2
----	---	-------

# # Complexity

Date / / 20

'n' documents. 'm' words.  
(Total docs.) in index. 'w' words per doc.

(1) complexity of reading every word

$$O(nw)$$

(2) for lookup function, at most 'm' words  
 $O(\log m)$

Assump: Structure used to store Term index is balanced Binary Tree

(3) One more function to append to docID the new document.

Considering all of these, the Total Time complexity is :  $O(nw \log m)$

I love pets  
pets are good

I don't like  
pets.

S1: term DocID

I	1
love	1
pets	1
pets	1
are	1
good	1

term DocID

I	2
don't	2
like	2
pets.	2

S2: Term ID

are	1
don't	2
good	1
I	1
I	2
like	2
love	1
pets	1
pets	2

S3: Sl.no. term DocID

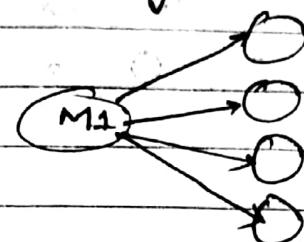
1	are	1	1
2	don't	1	2
3	good	1	1
4	I	2	1 → 2
5	like	1	2
6	love	1	1
7	pets	2	1 → 2

## Web Intelligence :

### Page Rank, Hyperlinks, link Analysis.

(Q) How to tell if a mail account is a spam account; without looking into contents.

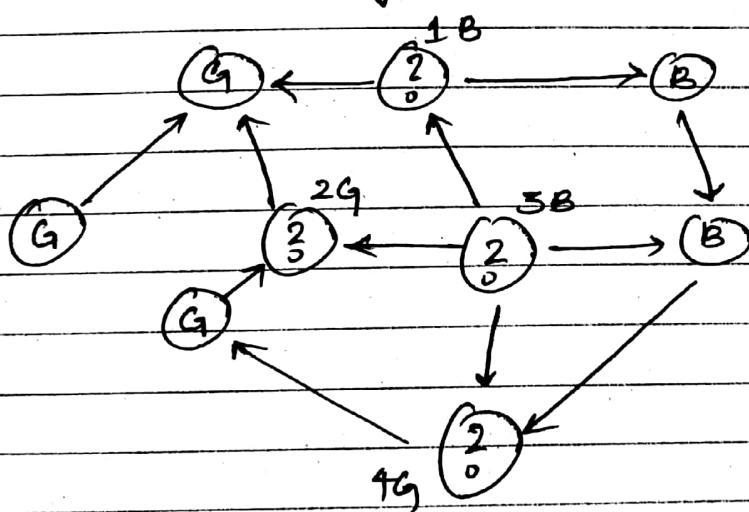
(A) No In-link.



(1) No Good Node points to a bad node.

(2) If a node is pointing to a bad Node, then, it is bad.

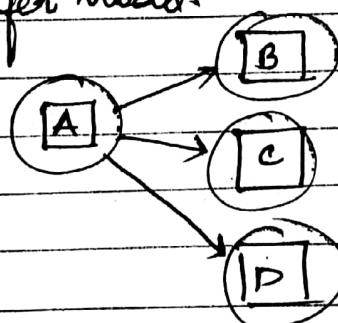
(3) If good nodes points to a node, then, it is good.



### # Page Rank

#### Random Surfer Model:

Page Rank:



$$\text{Prob} = 1/n$$

$n \rightarrow$   
No. of  
outlinks

Once you reach a node that has no outlink then, we use teleport operation

- (1) Navigate to a different page
- (2) Go Back

Teleport used when:

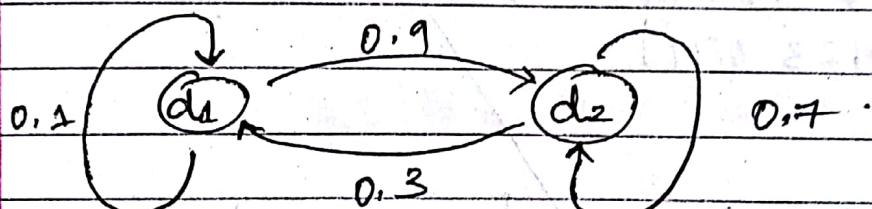
- (1) Node has no outlink
- (2)  $0 < \alpha < 1$ .

# To convert Adjacency Matrix to TPM

(1) Divide each '1' in A by total no. of 1's in the rows.

(2) Multiply resulting matrix by  $1-\alpha$ .

(3) Add.  $\alpha/N$  to every entry to the resulting matrix, to obtain



$$\begin{matrix} x_1 & x_2 \end{matrix}$$

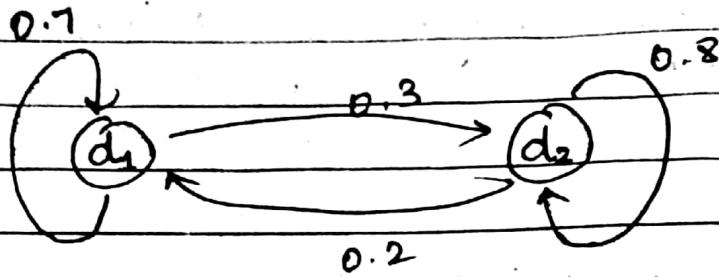
$$P_t(d_1) \quad P_t(d_2)$$

$$\begin{bmatrix} 0.1 \\ 0.3 \end{bmatrix} \begin{bmatrix} 0.1 & 0.9 \\ 0.3 & 0.7 \end{bmatrix}$$

$$\begin{aligned} P_{11} &= 0.1 & P_{12} &= 0.9 \\ P_{21} &= 0.3 & P_{22} &= 0.7 \end{aligned}$$

$$\begin{array}{ccccc} 0 & 1 & 0.3 & 0.7 & \vec{x}P \\ 0.3 & 0.7 & 0.24 & 0.76 & \vec{x}P^2 \\ 0.24 & 0.76 & 0.252 & 0.748 & \vec{x}P^3 \\ 0.252 & 0.748 & 0.2496 & 0.7504 & \vec{x}P^4 \\ 0.25 & 0.75 & 0.25 & 0.75 & \vec{x}P^5 \end{array}$$

# Power Method.



I

$$\begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 0.1 & 0.3 \\ 0.2 & 0.8 \end{bmatrix}$$

Ans.  $\begin{bmatrix} 0.4 & 0.6 \end{bmatrix}$

II

$$= \begin{bmatrix} 0.3 & 0.7 \end{bmatrix} \begin{bmatrix} 0.3 & 0.7 \\ 0.2 & 0.8 \end{bmatrix}$$

$$= 0.09 + 0.14 \quad 0.21 + 0.56$$

$$= \begin{bmatrix} 0.23 & 0.77 \end{bmatrix}$$

III

$$\begin{bmatrix} 0.23 & 0.77 \end{bmatrix} \begin{bmatrix} 0.3 & 0.7 \\ 0.2 & 0.8 \end{bmatrix}$$

$$= \begin{bmatrix} 0.2223 & 0.7777 \end{bmatrix}$$

$$= \begin{bmatrix} 0.2223 & 0.7777 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 0.1 & 0.3 \\ 0.2 & 0.8 \end{bmatrix}$$

Ans.  $\begin{bmatrix} 0.4 & 0.6 \end{bmatrix}$

PageRank Vector =  $\vec{\pi} = (\pi_1, \pi_2)$   
 $= (0.4, 0.6)$

Date / / 120

$$\Pi = \Pi \begin{bmatrix} 0.7 & 0.3 \\ 0.2 & 0.8 \end{bmatrix}$$

$$\Pi = \begin{bmatrix} x & y \end{bmatrix}$$

$$\begin{bmatrix} x & y \end{bmatrix} = \begin{bmatrix} x & y \end{bmatrix} \begin{bmatrix} 0.7 & 0.3 \\ 0.2 & 0.8 \end{bmatrix}$$

$$\begin{bmatrix} x & y \end{bmatrix} = \begin{bmatrix} 0.7x + 0.2y & 0.3x + 0.8y \end{bmatrix}$$

$$x = 0.7x + 0.2y$$

$$x + y = 1$$

②

$$0.3x = 0.2y$$

$$y = \frac{3}{2}x$$

$$x + \frac{3}{2}x = 1$$

$$\frac{5x}{2} = 1$$

$$x = \frac{2}{5} = 0.4$$

(1)  $\alpha = 0.4$  consider Graph

$1 \rightarrow 2 \quad 2 \rightarrow 1 \quad 2 \rightarrow 3 \quad 3 \rightarrow 2$

$4 \rightarrow 3 \quad 4 \rightarrow 2 \quad 4 \rightarrow 1 \quad 1 \rightarrow 3$

	1	2	3	4.
1	0	1	1	0
2	1	0	1	0
3	0	1	0	0
4	1	1	1	0

Step 1: Divide by total no. of 1's in row

$$\begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & \frac{1}{3} & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 \end{bmatrix}$$

$$\alpha = 0.4 \Rightarrow 1 - \alpha = 0.6.$$

Step 2: Multiply the resulting matrix by  $1 - \alpha$

$$[0.6] \begin{bmatrix} 0 & 0.5 & 0.5 & 0 \\ 0.5 & 0 & 0.5 & 0 \\ 0 & 1 & 0 & 0 \\ 0.33 & 0.33 & 0.33 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0.3 & 0.3 & 0 \\ 0.3 & 0 & 0.3 & 0 \\ 0 & 0.6 & 0 & 0 \\ 0.198 & 0.198 & 0.198 & 0 \end{bmatrix}$$

Step (3) : Add  $\alpha/n$  to everything :

$$\alpha = 0.4 \quad n = 4$$

$n \rightarrow$  Total no. of nodes.

$$\begin{bmatrix} 0.1 & 0.4 & 0.4 & 0.1 \\ 0.4 & 0.1 & 0.4 & 0.1 \\ 0.1 & 0.7 & 0.1 & 0.1 \\ 0.3 & 0.3 & 0.3 & 0.1 \end{bmatrix}$$

↑ TPM from Adjacency Matrix

Now, to get PageRank

$$\begin{bmatrix} 1 & 0 \end{bmatrix} \begin{bmatrix} 0.1 & 0.4 & 0.4 & 0.1 \\ 0.4 & 0.1 & 0.4 & 0.1 \\ 0.1 & 0.7 & 0.1 & 0.1 \\ 0.3 & 0.3 & 0.3 & 0.1 \end{bmatrix}$$

Page Rank  $\rightarrow$  human memory

Word-word Association

E.g. when 'Apple' is given people think 'Orange' as the next word

Experiment were done & made a semantic network which works better than web for searching

Letter-word Association

'a'  $\rightarrow$  apple  
 $\rightarrow$  aeroplane

Also  $\rightarrow$  ends up relearning a bunch of words of words

Google indexing & human memory is not same.

Human memory is poor in recalling facts but Google does it easily  
context - human needs context to recall

fuzzy - each & every detail not possible

sparse - a small hint can also trigger to recall a whole incident

LSH  $\rightarrow$  Locality Sensitive Hashing

Shingling  $\rightarrow$  divide doc. into sets where no element is repeated

Given 2 new assignments, how to find if they have been plagiarized.  
doc. similarity

## # PageRank

more  $\rightarrow$  human memory

"Google & the mind"

Psychological.

### (1) Word-Word Association

formed semantic network

web-graphs  $\Rightarrow$  hyperlink

### (2) Algorithm:

Ends up returning a bunch of words of words

Search works better for a semantic network as compared to a web graph

### (Q) Is human memory same as Google indexing

(A) No,

- Human mind is poor at recalling facts.
- " " needs context
- fuzzy
- sparse

Searching: Structured Data: complex in nature.

LSH: Locality Sensitive Hashing

(Q) Given a network assignments, how to find if they have been plagiarized.

(A) Document Similarity

Shingling : A common technique of representing sets.

LSH → Locality Sensitive Hashing.

General idea of LSH →  
Hashing items such that similar items fall into the same bin.

Starting point = Similar docs → Similar sets

hard part → Arranging similar items such that they fall into the same bin

Starting point = Similar docs → Similar sets.

# Applications of set similarities :

(1) Grouping of docs, pages with similar content  
→ topic

(2) Netflix → Recommending.

(3) Movie fans.

(4) Entropy Matching.

## #1 Applications of Document Similarity!

(1) Mirror sites

(2) Plagiarism.

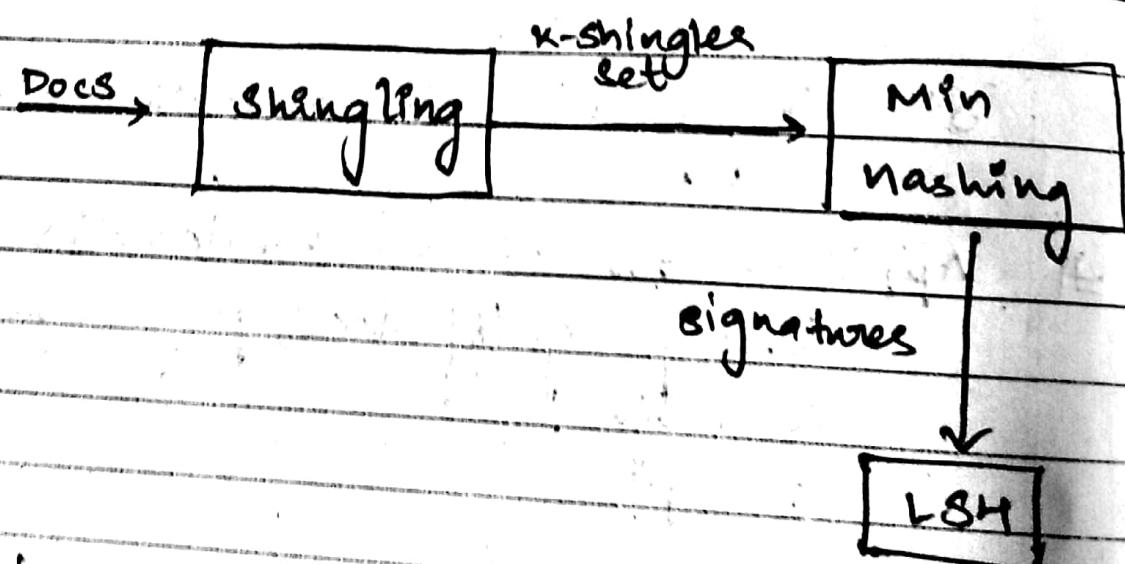
(3) e-news.

## # Calculate Similarity

1) Shingling.

(2) Min-Mashing.

(3) LSH.



Shingling: set of k-sensitive words

2K abcab → {ab, bc, ca, }  
Set

## Shingler & Doc Similarity

→ most of the shingles.

→ 2k shingles (other inter-changing paragraph)

Min Hashing:

→ Jaccard similarity.

→ Signature Matrix

\*

$$JC \rightarrow \text{sim}(C_1, C_2) = |C_1 \cap C_2| / |C_1 \cup C_2|$$

Rows → universal elements:

cols → sets

$$JC = A / (A + B + C)$$

C1	C2		C1	C2
0	1	2/5	A	1
1	0		B	1
1	1	*	C	0
0	0		D	0
1	1	*		
0	1			

0  
1  
1  
0

Maturity level of an organization w.r.t. to the ability of an org for continuous improvement in a particular

### # Min Hashing:

Element	$s_1$	$s_2$	$s_3$	$s_4$
b	0	0	1	0
e	0	0	1	0
a	1	0	0	1
d	1	0	1	1
c	0	1	0	1
;				

$$h(s_1) = a$$

$$h(s_2) = e$$

$$h(s_3) = b$$

$$h(s_4) = a$$

Row	$c_1$	$c_2$	$c_3$	$h(x) = x \bmod 5$	$g(x) = 2x + 1 \bmod 5$
1	1	0			
2	0	1			
3	1	1			
4	1	0			
5	0	1			

$$JC = 1/5$$

Date / / 20

	Sig 1	Sig 2
$h(x)$	$\infty$	$\infty$
$g(x)$	$\infty$	$\infty$

$$h(1) = 1$$

$$g(1) = 3$$

$$\begin{array}{cc} \text{sig 1} & \text{sig 2} \\ h(1) = 1 & 1 \\ g(1) = 3 & 3 \end{array}$$

$$\begin{array}{cc} h(2) = 2 & 1 \\ g(2) = 0 & 3 \end{array}$$

$$\begin{array}{cc} h(3) = 3 & 1 \\ g(3) = 2 & 2 \end{array}$$

$$\begin{array}{cc} h(4) = 4 & 1 \\ g(4) = 4 & 2 \end{array}$$

$$\begin{array}{cc} h(5) = 0 & 1 \\ g(5) = 1 & 2 \end{array}$$

Row	$S_1$	$S_2$	$S_3$	$S_4$	
0	1	0	0	1	$h_1(x) = (x+1) \bmod 5$
1	0	0	1	0	$h_2(x) = (3x+1) \bmod 5$
2	0	1	0	1	
3	1	0	1	1	
4	0	0	1	0	

$S_1$	$S_2$	$S_3$	$S_4$	<del>if</del>
$\infty$	$\infty$	$\infty$	$\infty$	$h_1(0) = 1$
$\infty$	$\infty$	$\infty$	$\infty$	$h_2(0) = 1$

$$h_1(0) = 1 \quad | \quad 1 \quad 00 \quad 00 \quad 1$$

$$h_2(0) = 2 \quad | \quad 1 \quad 00 \quad \infty \quad 1$$

$$h_1(1) = 2 \quad | \quad 1 \quad \infty \quad 2 \quad 1$$

$$h_2(1) = 4 \quad | \quad 1 \quad \infty \quad 4 \quad 1$$

$$h_1(2) = 3 \quad | \quad 1 \quad 3 \quad 2 \quad 1$$

$$h_2(2) = 2 \quad | \quad 1 \quad 2 \quad 4 \quad 1$$

$$h_1(3) = 4 \quad | \quad 1 \quad 3 \quad 2 \quad 1$$

$$h_2(3) = 0 \quad | \quad 0 \quad 2 \quad \text{X} \quad 0 \\ 0$$

$$h_1(4) = 0 \quad | \quad 1 \quad 3 \quad 0 \quad 1$$

$$h_2(4) = 3 \quad | \quad 0 \quad 2 \quad 0 \quad 0$$

$$\text{Sim}(S_1, S_4) =$$

0	0	1	0	1
1	0	1	0	0
2	1	0	0	1
3	0	0	1	0
4	0	0	1	1
5	1	0	0	0

$$h_1(x) = (2x+1) \bmod 6$$

$$h_2(x) = (3x+2) \bmod 6$$

$$h_3(x) = (3x+2) \bmod 6$$

S1	S2	S3	S4
----	----	----	----

$$h_1(0) = 1$$

$\infty$	1	$\infty$	1
----------	---	----------	---

$$h_2(0) = 2$$

$\infty$	2	$\infty$	2
----------	---	----------	---

$$h_3(0) = 2$$

$\infty$	2	$\infty$	2
----------	---	----------	---

$$h_1(1) = 3$$

$\infty$	1	$\infty$	1
----------	---	----------	---

$$h_2(1) = 5$$

$\infty$	2	$\infty$	2
----------	---	----------	---

$$h_3(1) = 1$$

$\infty$	1	$\infty$	2
----------	---	----------	---

$$h_1(2) = 5$$

5	1	$\infty$	1
---	---	----------	---

$$h_2(2) = 2$$

2	2	$\infty$	2
---	---	----------	---

$$h_3(2) = 0$$

0	1	$\infty$	0
---	---	----------	---

$$h_1(3) = 1$$

5	1	1	1
---	---	---	---

$$h_2(3) = 5$$

2	2	5	2
---	---	---	---

$$h_3(3) = 5$$

0	1	5	0
---	---	---	---

$$h_1(4) = 3$$

5	1	1	1
---	---	---	---

$$h_2(4) = 2$$

2	2	2	2
---	---	---	---

$$h_3(4) = 4$$

0	1	4	0
---	---	---	---

$$h_1(5) = 5$$

5	1	1	1
---	---	---	---

$$h_2(5) = 5$$

2	2	2	2
---	---	---	---

$$h_3(5) = 3$$

0	1	4	0
---	---	---	---

$$2/3 = 66\frac{2}{3}\%$$

kind. of similar

## # Adsense, Keywords & mutual information

transmitted signal =  $\xrightarrow{\text{web-page content}}$  AdSense received signal → sequence of messages  
mutual information.

### Advertising Messages

## # Inverse Search.

pages to keywords.

query words to pages.

Tf - Idf.

- How important a word is.
- Tf-idf increases proportionality to the number of times.
- Rarer words make better keywords

2 ways to search (or show ads)

(1) Search based on keyword

(2) ⚡ Inverse Search

(Based on ~~freq~~ <sup>imp</sup> word in website)

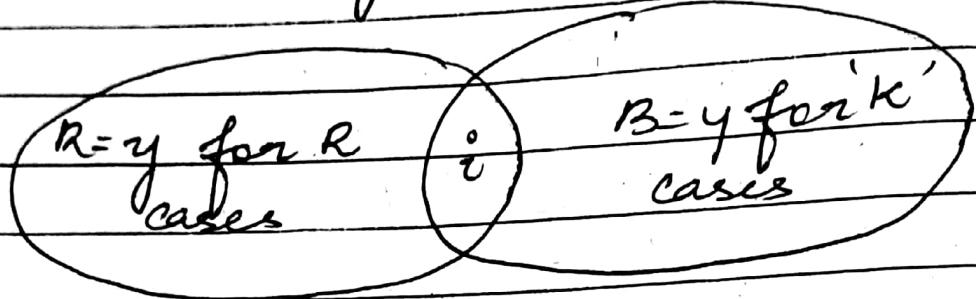
$$P(B|R) = \frac{i}{n}$$

$$P(R) = n/n$$

$$P(R \text{ and } B) = i/n = (i/n) * (n/n)$$

$$P(B, R) = P(B|R) \cdot P(R)$$

This is Baye's rule.



$$P(B, R) = P(B|R) \cdot P(R)$$

$$= P(R|B) \cdot P(B)$$

⇒ 8

Types of prob:

→ Conditional.  $P(x_1 | x_2)$

→ Posterior  $P(x | x_1)$

→ A Priori  $P(x_1) P(x_2) P(x)$

Independence:

$$\frac{i}{c} = \frac{n}{n} \equiv \frac{1}{n} = \frac{c}{n}$$

$$P(R) = n/n \quad P(C) = c/n$$

$$P(R|C) = i/c \quad P(C|R) = i/n$$

R & C are independent

(Q) Machine A & B produce produce 10% & 90% effectively. of B the pass of component intended for motor industry. From the experience, it is known that prob. that A produces a defective comp is 0.01  
 $B \rightarrow$  defective component  $\rightarrow 0.05$   
 Pick random comp  
 Find the prob. that it was made by A or by B.

Prob (A) = (produced by A)

$$P(A) = 0.1$$

$$P(D|A) = 0.01$$

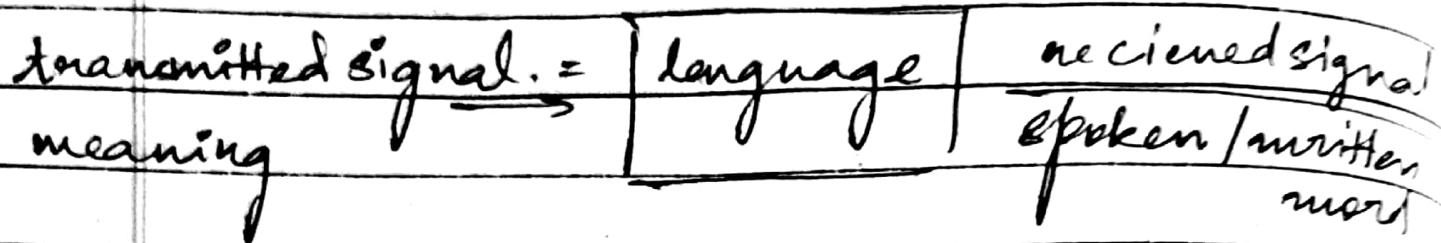
$$P(B) = 0.9$$

$$P(D|B) = 0.05$$

$$P(A|D) = \frac{P(D|A) \cdot P(A)}{P(D|A) \cdot P(A) + P(D|B) \cdot P(B)}$$

$$= 0.02$$

# Language & Information



language is redundant.  
language tries to maintain "uniform information density"

- Besides TF-idf.
    - keyword is present in how many docs.
    - semantics of a word.
    - bipartite graph.
  - Machine learning.

## Prediction using conditional probability

## # Sentiment Analysis:

lib

Date 1/120

limiting physical.

External control

Hardware limitation

Mutual Information:

double summation

$$I(F, B) = \sum_{f, b} p(f, b) \log \frac{p(f, b)}{p(f)p(b)}$$

Limitations of MI:

## ASSIGNMENT - 1

$$P(A) = .1, P(B) = .9$$

$$P(D/A) = .01, P(D/B) = .05$$

$$P(A/D) = \frac{P(D/A) \cdot P(A)}{P(D/A) \cdot P(A) + P(D/B) \cdot P(B)} = \frac{.01 \times .1}{(.01 \times .1) + (.05 \times .9)} = \frac{.001}{.046} = \underline{\underline{.022}}$$

$$P(B/D) = \frac{P(D/B) \cdot P(B)}{P(D/A) \cdot P(B) + P(D/B) \cdot P(A)} = \frac{.05 \times .9}{.046} = \underline{\underline{.978}}$$

$$h_1(x) = (2x + 4) \% 5$$

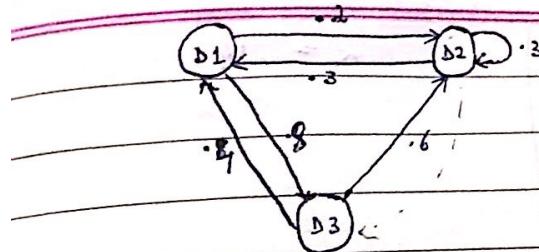
$$h_2(x) = (3x - 1) \% 5$$

## Minhash SM

	$s_1$	$s_2$	$s_3$	$s_4$	$h_1(x)$	$\text{sig}_1$	$\text{sig}_2$	$\text{sig}_3$	$\text{sig}_4$
Row									
0	1	0	0	1	$h_2(x)$	$\infty$	$\infty$	$\infty$	$\infty$
1	0	0	1	0		$\infty$	$\infty$	$\infty$	$\infty$
2	0	1	0	1	$h_1(0)$	4	$\infty$	$\infty$	4
3	1	0	1	1	$h_2(0)$	4	$\infty$	$\infty$	4
4	0	0	1	0	$h_1(1)$	4	$\infty$	$\infty$	4
					$h_2(1)$	4	$\infty$	<u>1</u>	4
					$h_1(2)$	4	$\infty$	2	4
					$h_2(2)$	4	0	<u>1</u>	3
					$h_1(3)$	0	<u>3</u>	0	0
					$h_2(3)$	3	0	2	0
					$h_1(4)$	0	3	0	0
					$h_2(4)$	3	0	1	0

0	3	0	0
3	0	1	0

||



$$P = \begin{bmatrix} 0 & .2 & -.8 \\ -.3 & .3 & 0 \\ -.8 & -.6 & 0 \end{bmatrix}$$

$$t = [0 \ 0 \ 1]$$

Power Method

$$\begin{array}{ccc} x(D_1) & x(D_2) & x(D_3) \\ \hline 0 & 0 & 1 \end{array}$$

$$\begin{array}{cccccc} & & & .4 & .6 & 0 \\ 0 & 0 & 1 & -18 & .26 & .32 \\ \cdot 4 & \cdot 6 & 0 & -206 & .306 & -144 \\ \cdot 18 & \cdot 26 & \cdot 32 & \cdot 1494 & \cdot 2194 & \cdot 1648 \\ \cdot 206 & \cdot 306 & \cdot 144 & \cdot 2342 & \cdot 3442 & \cdot 4216 \end{array} \rightarrow X P^2$$

$$P = \begin{bmatrix} 0 & .2 & .8 \\ -.3 & .3 & -.4 \\ -.4 & -.6 & 0 \end{bmatrix}$$

Power Method

$$\begin{array}{ccc} x(D_1) & x(D_2) & x(D_3) \\ \hline 0 & 0 & 1 \end{array}$$

$$\begin{array}{cccccc} & & & .4 & .6 & 0 \\ 0 & 0 & 1 & -18 & .26 & .56 \\ \cdot 4 & \cdot 6 & 0 & -302 & .45 & -248 \\ \cdot 18 & \cdot 26 & \cdot 56 & -2342 & .3442 & \cdot 4216 \\ \cdot 302 & \cdot 45 & \cdot 248 & -2719 & -403 & \cdot 325 \\ \cdot 2342 & \cdot 3442 & \cdot 4216 & -2509 & \cdot 3702 & \cdot 3787 \end{array} \rightarrow X P^5$$

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

classmate  
Date \_\_\_\_\_  
Page \_\_\_\_\_

d<sub>0</sub> d<sub>1</sub> d<sub>2</sub> d<sub>3</sub> d<sub>4</sub> d<sub>5</sub> d<sub>6</sub>

d <sub>0</sub>	0	0	1	0	0	0	0
d <sub>1</sub>	0	1	1	0	0	0	0
d <sub>2</sub>	1	0	1	0	0	0	0
d <sub>3</sub>	0	0	0	1	1	0	0
d <sub>4</sub>	0	0	0	0	0	0	1
d <sub>5</sub>	0	0	0	0	0	1	1
d <sub>6</sub>	0	0	0	1	1	0	1

$$\alpha = .142857$$

divide by no. 3's in each row

0	0	1	0	0	0	0
0	$\frac{1}{2}$	$\frac{1}{2}$	0	0	0	0
$\frac{9}{2}$	$\frac{1}{2}$	$\frac{1}{2}$	0	$\frac{1}{2}$	0	$\frac{1}{2}$
$\frac{9}{2}$	0	$\frac{1}{2}$	0	$\frac{1}{2}$	0	$\frac{1}{2}$
$\frac{9}{2}$	0	0	$\frac{1}{2}$	0	$\frac{1}{2}$	0
$\frac{9}{2}$	0	0	$\frac{1}{2}$	0	$\frac{1}{2}$	$\frac{1}{2}$
0	0	0	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

initial result

$$(.4) \times \quad (.4) \times \quad (.4) \times$$

$$-.4 \quad -.4 \quad -.4$$

$$0 \quad 0 \quad 0$$

$$.8 \quad .8 \quad .8$$

$$.8 \quad .8 \quad .8$$

$$.8 \quad .8 \quad .8$$

multiply by  $(1-\alpha)$

0	0	-86	0	0	0	0
0	.43	.43	0	0	0	0
.43	0	.43	0	0	0	0
0	0	0	.43	.43	0	0
.08	0	0	0	0	.86	0
.08	0	0	.36	0	.08	.43
.08	0	0	.287	.287	0	.287

$$8.888\ldots = 9$$

$$.888\ldots = 1$$

$$.888\ldots = 1$$

$$(.4) \times \quad (.4) \times \quad (.4) \times$$

$$.01 \quad .01 \quad .01$$

$$.01 \quad .01 \quad .01$$

$$.01 \quad .01 \quad .01$$

add  $\alpha \times \text{d/N} = \text{new value}$

.02	.02	.888	.02	.02	.02	.02
.02	.02	.45	.02	.02	.02	.02
.45	.02	.45	.02	.02	.02	.02
.02	.02	.02	.45	.45	.02	.02
.02	.02	.02	.02	.02	.02	.88
.02	.02	.02	.02	.02	.45	.45

$$.2188 \quad .3488 \quad .4788$$

$$.2188 \quad .202 \quad .0188$$

$$.2188 \quad .202 \quad .0188$$

5) Doc 1 → 1 did not

Doc 2 → 6 did not

Team Doc #

J did

exact

Julius

Cesar

J was

killed

J'

The

capital

brother

killed

me

to

let

it

be

with

Cesar

the

whole

brother

both

told

you

call

we

and

thanked

+

ea

l

meals  
start  
again

classmate  
Date  
Page

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

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

Term	Doc #	Term	Doc #
I	1	ambitious	2
did	1	be	2
enact	1	brutus	1
Julius	1	brutus	2
Caesar	1	Capitol	1
I'	1	Caesar	1
was	1	Caesar	2
killed	1	Caesar	2
I'	1	did	1
the	1	enact	1
Capitol	1	both	2
Brutus	1	I	1
killed	1	i'	2
me	1	it	2
so	2	Julius	1
let	2	killed	1
it	2	killed	1
be	2	so let	2
with	2	me	1
Caesar	2	noble	2
the	2	so	2
noble	2	the	2
Brutus	2	told	2
hath	2	you	2
told	2	was	1
you	2	was	2
Caesar	2	with	2
was	2		
ambitious	2		

Inverted Index :-

term doc-freq. postings list

ambitious | 1 → [ ] → 2

be | 1 → [ ] → 2

Brutus | 2 → [ ] → 1 → 2

## Notes.

### Shingles:

- Define a  $k$ -shingle for a doc to be any substring of length  $k$  found in the doc.
- Then we associate with each doc the set of  $n$  shingles that appear one or more times in that doc.

How to pick  $k$  (Shingle size).

- $k$  should be picked large enough so that probability of any given shingle appearing in any given doc is low.
- For large docs, such as research articles,  $k=9$  is considered safe.

Hashing Shingles: we can repeat all the  $9$ -shingles and hash them down to  $4$  by  $1$ s.

Shingles Built from Words. Using say, a stop word followed by any two words as a shingle.

Ex 3.3.2 Minhash Signatures.

Ele	S.	S.	S,	S,
0	0	1	0	1
			0	0
2	1	0	0	1
3	0	0	1	0
4	0	0	1	1
5	0	0	0	0

$h_1(2)$	3	2	0	3
$h_2(2)$	0	2	0	0
$h_1(3)$	3	2	0	3
$h_2(3)$	0	2	3	0
$h_1(4)$	3	2	0	2
$h_2(4)$	0	2	1	0
$h_1(5)$	3	2	0	2
$h_2(5)$	0	2	1	0

Signature matrix -  $\begin{bmatrix} 3 & 2 & 0 & 2 \\ 0 & 2 & 1 & 0 \end{bmatrix}$

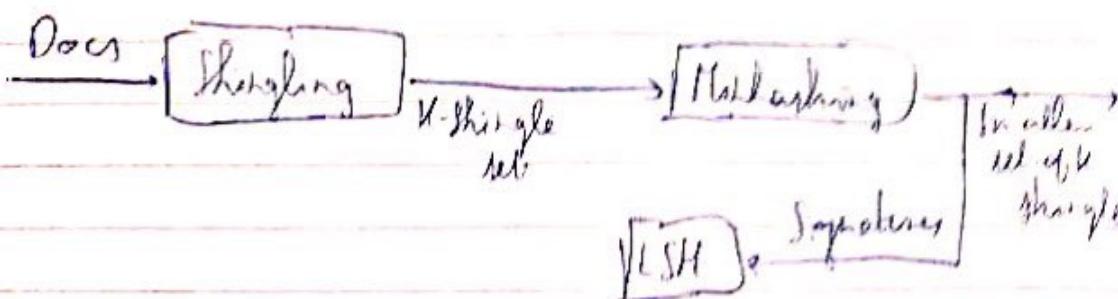
Locality Sensitive Hashing (LSH): Instead of finding similarity of every pair (tedious), often we want to find only the most similar pairs or all pairs that are above some lower bound in similarity. For this we use LSH.

Approach: "Hash" items several times so that similar items are more likely to be hashed to the same bucket than dissimilar items are.

→ Then, pairs within the same bucket are compared (candidate pairs). Only these are checked for similarity.

Note: Pairs that hash to the same bucket are called

## False Positives



## Cheat

Cheat 3.4.3 from the slides.

Associative memory: A type of computer memory from which items may be retrieved by matching some part of their content, rather than by specifying their address.  
→ AM is much slower than RAM.  
→ Used in multilevel memory systems.  
→ To retrieve a word from AM, a search key (or digest) is compared with locality tag bits of all stored words, and corresponding matches are found.  
→ AM is expensive to implement.

Space distributed memory: Generalized RAM for long binary words

→ Main attribute is that it is sensitive to similarity: A word can be read back by not only giving its address but also by giving one close word, as measured by Hamming distance.



### Jaccard Similarity:

We shall focus initially on a particular notion of “similarity”: the similarity of sets by looking at the relative size of their intersection. This notion of similarity is called “Jaccard similarity”.

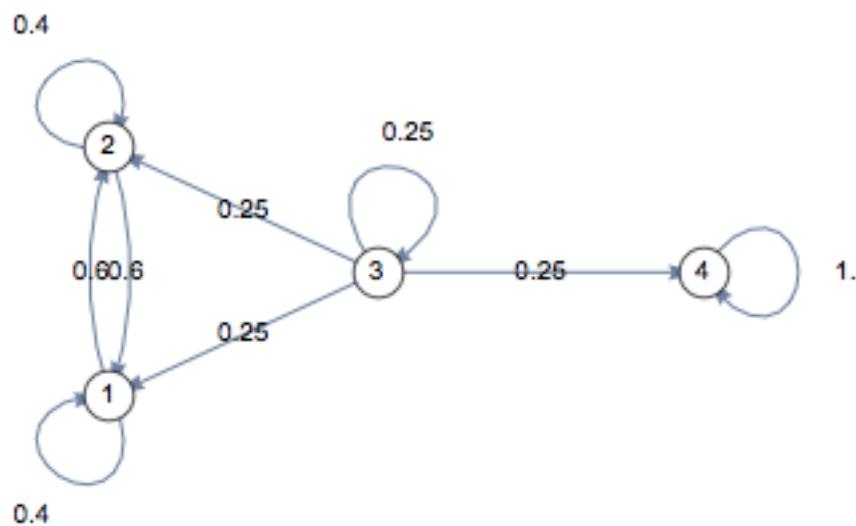
**Formula:**  $|S \cap T|/|S \cup T|$

### Markov Chain:

A Markov chain is a mathematical system that experiences transitions from one state to another according to certain probabilistic rules. The defining characteristic of a Markov chain is that no matter how the process arrived at its present state, the possible future states are fixed. In other words, the probability of transitioning to any particular state is dependent solely on the current state and time elapsed. The state space, or set of all possible states, can be anything: letters, numbers, weather conditions, baseball scores, or stock performances.

Markov chains may be modelled by finite state machines, and random walks provide a prolific example of their usefulness in mathematics. They arise broadly in statistical and information-theoretical contexts and are widely employed in economics, game theory, queueing (communication) theory, genetics, and finance. While it is possible to discuss Markov chains with any size of state space, the initial theory and most applications are focused on cases with a finite (or countably infinite) number of states.

Markov chain Monte Carlo (MCMC) is a technique for estimating by simulation the expectation of a statistic in a complex model. Successive random selections form a Markov chain, the stationary distribution of which is the target distribution. It is particularly useful for the evaluation of posterior distributions in complex Bayesian models.



### Semantic Graphs:

A **semantic graph** is a network of heterogeneous nodes and links annotated with a domain ontology. In information analysis, **semantic graphs** are generated and applied in a visual analysis approach known as link analysis.

Example: RDF.

### **Conditional Probability:**

$$P(a|y) = (P(y|a)*P(a))/P(y).$$

### **Shannon Coding:**

In the field of **data compression**, **Shannon coding**, named after its creator is a **lossless data compression** technique for constructing a **prefix code** based on a set of symbols and their probabilities (estimated or measured). It is suboptimal in the sense that it does not achieve the lowest possible expected code word length like **Huffman coding** does, and never better but sometime equal to the **Shannon-Fano coding**.

Easy Video of 3 minutes: <https://www.youtube.com/watch?v=dJCck1OgsIA>