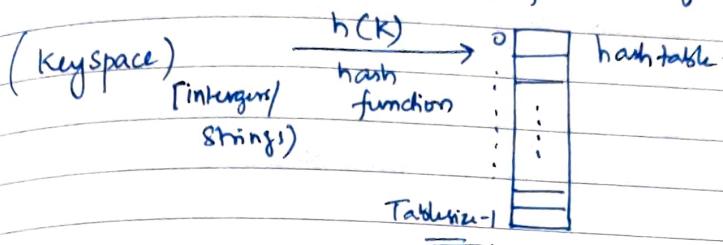


## ① Hashing Techniques:-

A hash table is an array of some fixed size (usually prime no).



Hash func :-

1. Simple & fast to compute
2. Avoid collisions.
3. Keys distributed evenly among the Cells.

Collision - When 2 keys map to the same location in hash table.

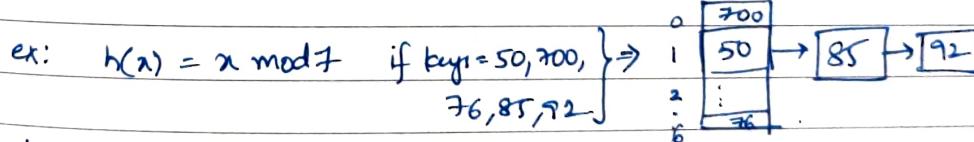
② 2 ways to resolve Collisions:-

- (i) Separate Hashing Chaining.
- (ii) Open Addressing

→ Separate Chaining :- All the keys that map to same hash value are kept in a (list/bucket).

Note: 'hashing' is <sup>technique</sup> used to index & retrieve items in a database in a reasonable amount of complexity).

→ Main idea is to make each cell of hash table point to linked list of records that have same hash function value.



disadv :- Cache performance is not good as keys are stored using linked list.

→ Wastage of space, Extra space for links ( $\because$  we are storing address).

→ If chain becomes long, search space can become  $O(n)$  in worst case.

Time Complexity :  $O(1+\alpha)$   $\alpha = \text{Avg. keys per slot (or)} \Rightarrow \frac{n}{m} \rightarrow \# \text{ keys stored in table}$   
 $\text{load factor} \rightarrow \# \text{ slot in table}$

## Open Addressing :-

- All items (keys) are stored in table itself. ( $\text{size of table} \geq \#(\text{keys})$ ).
- Hash func. specifies order of slots to probe (try) for a key (for insert / search, ~~only~~ not just one slot).

(i) Linear Probing :-  $h_i(x) = [(Hash(x) + i) \% \text{HashTable size}]$

(try for next empty slot)

ie If  $h_0 = (Hash(x) + 0) \% \text{HashTable size}$  if it is full, then try  
 $h_1 = (Hash(x) + 1) \% \text{HashTable size}$  " try  $h_2$  for  $h_3$   
 and so on.

(Keep probing until empty slot is found).

(ii) Quadratic Probing:  $h_i(x) = [Hash(x) + i^2 \% \text{HashTable size}]$ .

(iii) Double Hashing :  $h_i(x) = [Hash(x) + i * Hash(x)^2 \% \text{HashTable size}]$ .  
 (here diff is 2 diff hash func. are used).

Comparison :-

<u>Linear probing</u>	<u>Quadratic probing</u>	<u>Double hashing</u>
→ Easy to implement	→ Avg cache performance	→ poor cache performance
→ Best cache performance	→ Suffers less clustering than linear probing $(\because$ easier to find empty slot here)	→ No clustering
→ Suffers from clustering		→ Requires more computation time.

Complexity : ' $\alpha$ ' → load factor  $\Leftrightarrow O\left(\frac{1}{1-\alpha}\right)$  time.  
 $= \frac{n}{m}$  → # keys to be inserted in hash table.

## Cuckoo Hashing :-

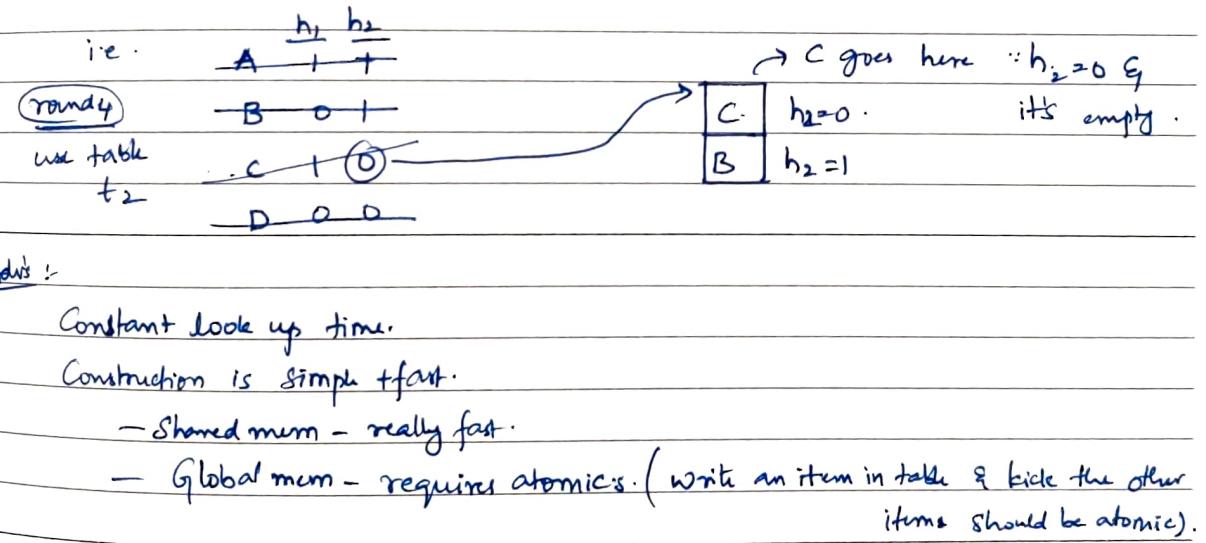
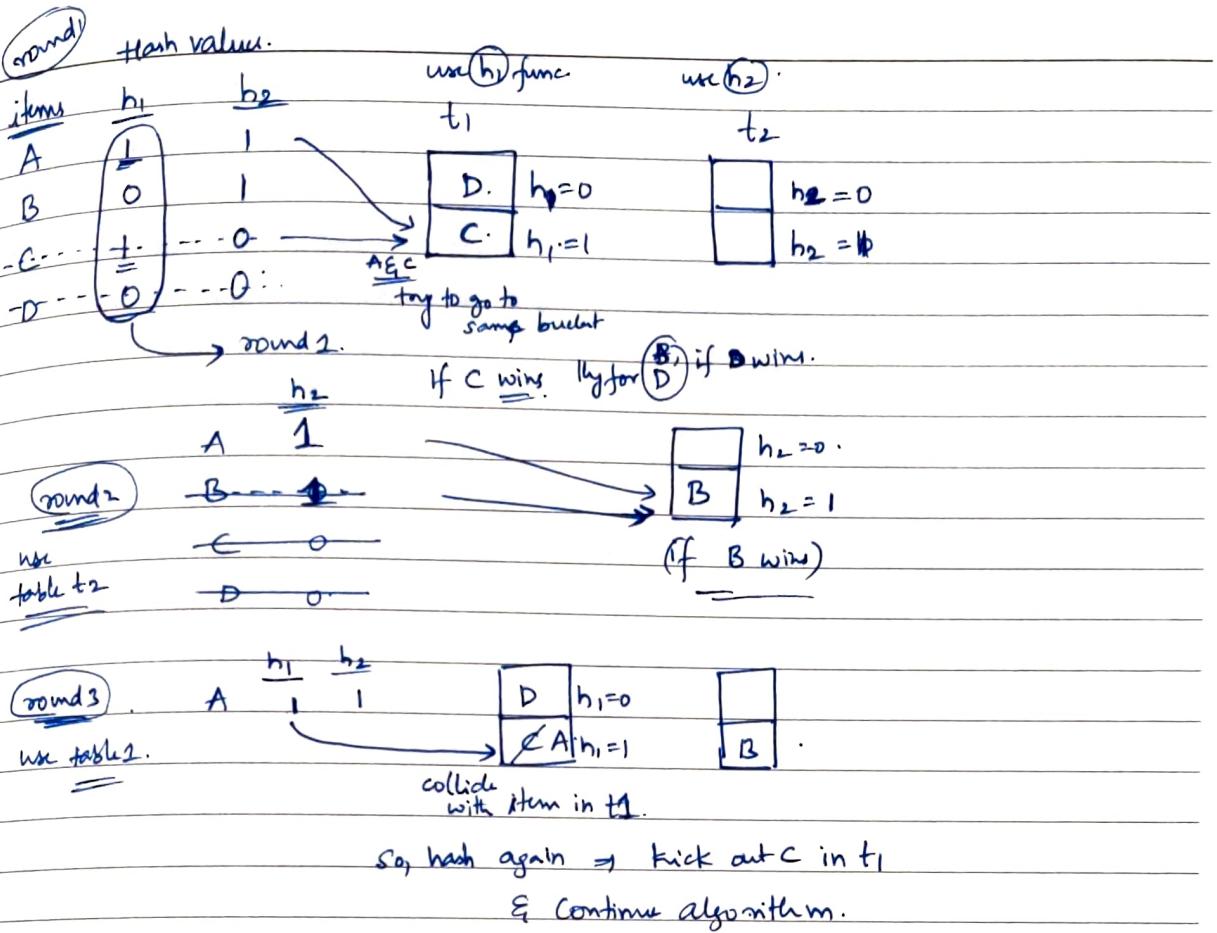
- multiple hash tables (not 1).
- multiple hash functions (1 per hashtable).

procedure :-

1<sup>st</sup> iteration - all items use  $h_1$  & write to  $t_1$

2<sup>nd</sup> (some fail)

- 2<sup>nd</sup> iter. - Those that fail use  $h_2, t_2$ .  
 repeat.



## 2) Filters And Sketching:

Bloom Filters :- (Space Efficient "probabilistic" data structure).

Note: In order to compare 2 strings 'u' we just compare their hash values & if they are same  $\rightarrow$  then it's highly likely that 2 strings are same. hash value.

(In general we pass strings to multiple hash func. & get their hash values).  
( $h_1, h_2, h_3, \dots$ )

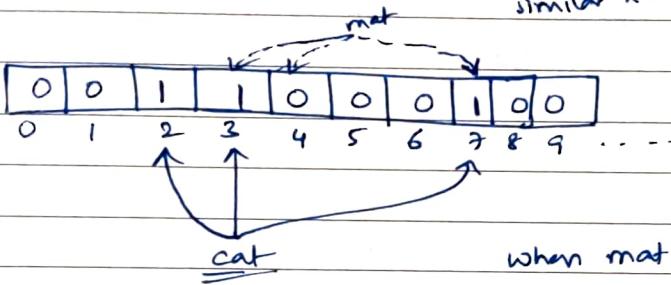
(if "ele" present in the set)  $\rightarrow$  yes (false positive are possible).

Ex: cat  $\rightarrow$  panthrough  $\frac{\text{cat}}{h_1 \rightarrow 3}$   $\frac{\text{mat}}{h_2 \rightarrow 7}$   $\frac{\text{mat}}{h_3 \rightarrow 2}$

$\rightarrow$  no (if ele is not there but bloom filter says it is there).  
[false negative not possible]

( $\because$  2<sup>nd</sup> char is same)

assumption: similar hash func. gets this value  
(cat is same)



when mat is being searched.

we get  $\boxed{1, 0, 1}$   
 $\boxed{3, 4, 7}$  bits

So, set (index of mat) to 1.

$\because$  not all are 1. we can clearly say "mat" is not searched.

0 0 1 1 1 0 0 1 0 0.

now if 'mat' is searched  $\rightarrow$  we get 111. it's there.

Bloom filter  $\rightarrow$  has func. has "rect"  $\rightarrow$  if 70% are 1's. (then bloom filter might have lower probability to say yes/no).  $\rightarrow$  then we can rect().

Now if some word

like 'fat'  $\rightarrow$   $h_1 \rightarrow 3$

$h_2 \rightarrow 4$

$h_3 \rightarrow 2$

0	1	2	3	4	5	6	7	8	9
0	0	1	1	1	0	0	1	0	0

all are 1's  $\rightarrow$  already searched. (false positive)  
(only drawback).

[No false  
-ve's]

$\rightarrow$  i.e. (ele is there but bloom filter says is not present)]

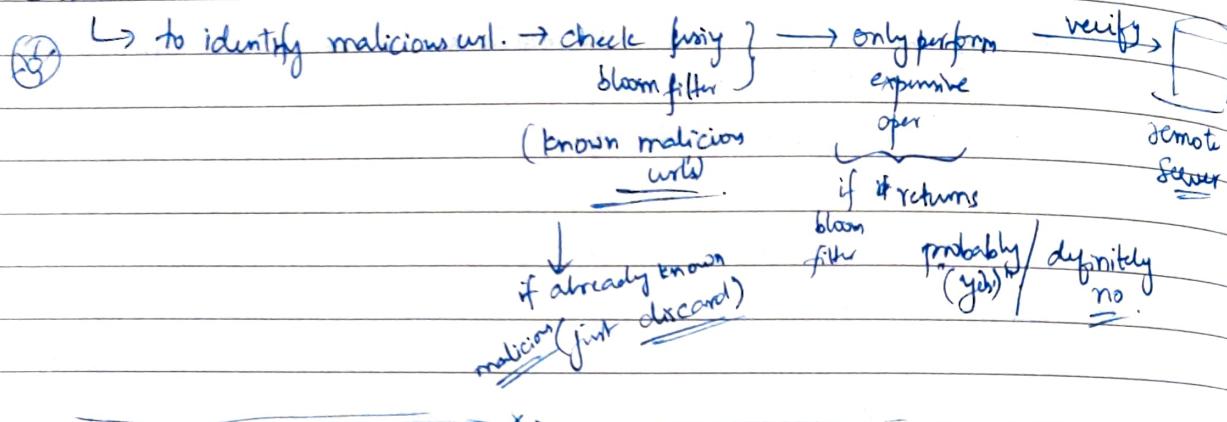
not such case in bloom filter.

Trade off: less mem.  $\rightarrow$  (so, some false pos.)  $\rightarrow$  (but no false -ve's).

m  $\rightarrow$  bit array  
k  $\rightarrow$  hash func.  
n  $\rightarrow$  strings.

Using bloom filter reduces caching workload.  $\rightarrow$  ↑ caching hit rate.  
 ↳ (check url)  $\rightarrow$  (only caching seen pg on second request).

chromosome:



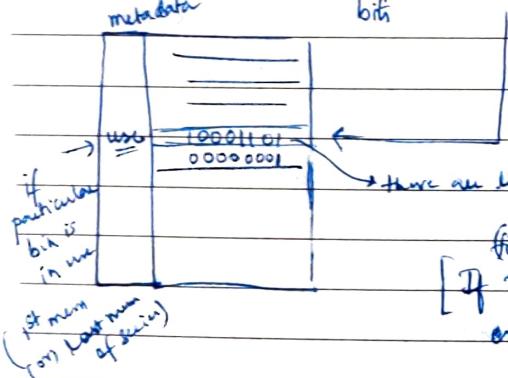
② Quotient Filter: approximate member query data structure.

Squeakr: an exact & approx. k-r counting system

Counting quotient filter

hash (Kmer)  $\rightarrow$  [1011010110001101]

metadata      upper bit      location      remainder (lower bits)



$\rightarrow$  you can use for querying

also counting (based on val in not location)

[If new hash produces same location  $\rightarrow$  then go to next empty slot.]

## ④ Locality Sensitive Hashing (LSH):

ex: similar docs  $\rightarrow$  similar hash-code.

For each doc 'd'  $\rightarrow$  generate k-bit hash code.  
insert doc in hash table.

\* If it's collision  $\rightarrow$  possible duplicate (compare docs in same bucket)

(sometimes it might be accident not same doc)

Can miss near-duplicates

- similar hash-code  $\neq$  same bucket

repeat L times w. diff hash tables (randomized)

cut 3D space by a plane making some angle w.r.t. axis.

$\begin{array}{c} \overrightarrow{H_1} \\ \overrightarrow{H_2} \\ \overrightarrow{H_3} \end{array}$

$d_1 \quad 0100 \quad 1011 \quad 0110$

$d_2 \quad 1100 \quad 1011 \quad 0010$

$\downarrow$

not (collision).

$\rightarrow$  Can generate random hyper planes (for generating hashcodes)

$\therefore$  mainly want: (similar hashcode for nearby points).

$\rightarrow$  Then do comparisons within same bucket  $\rightarrow$  to eliminate false positives.

$\rightarrow$  Repeat with diff hyper planes. (some x no. of times). (if a, d are far apart in actual).

cost: npts, d-dimensions, k-hyperplanes.

$\Delta K$   $\rightarrow$  find bucket where pt lands.

$2^k$   $\rightarrow$  hashcodes /  $\Delta K$  diff regions (due to intersection of hyperplanes)

$(N/2^k)$  - pts in bucket

$(DN/2^k)$   $\rightarrow$  cost of comparison (inside a bucket)

(repeat L times) - LSH:  $L\Delta K + LDN/2^k \rightarrow \log N$  if  $k = \log N$ .

$O(1)$

Jaccard Similarity:-

$$\text{Sim}(C_1, C_2) = \frac{\text{size of intersection}}{|C_1 \cup C_2|}$$

$\hookrightarrow$  union.

min hashing:-

	permuted	c1	c2	c3	c4	Input matrix
1	$\xrightarrow{4}$	1	0	1	0	
3	$\xrightarrow{2}$	4	1	0	0	1
7	$\xrightarrow{1}$	7	0	1	0	1
6	$\xrightarrow{3}$	6	0	1	0	1
2	$\xrightarrow{5}$	6	0	1	0	1
4	$\xrightarrow{7}$	1	0	1	0	
5	$\xrightarrow{6}$	1	0	1	0	
4	$\xrightarrow{5}$	5	1	0	1	0

hash function  
signature matrix

$\boxed{2 \ 1 \ 2 \ 1}$   $\rightarrow$  formed based on 1st permuted rows.

$\underbrace{\quad}_{\text{val. band}}$

on 1st per rows

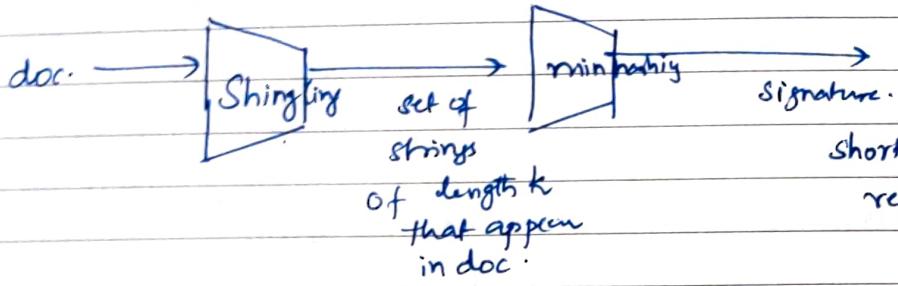
$\boxed{2 \ 1 \ 4 \ 1}$   $\rightarrow$  based on 2nd per rows

$\boxed{1 \ 2 \ 1 \ 2}$   $\rightarrow$  ... 3rd.

dis: If 1 billion rows  $\rightarrow$  no. of per.  $\rightarrow$  very huge  $\rightarrow$  (exceed disk space (4GB))  $\rightarrow$  accessing permuted rows leads to thrashing.

So, use hash fun. for permuting rows.  $h_i(r)$ .  $\rightarrow$   $i^{\text{th}}$  permutation.

### Min hashing



short integer vector that represent the sets  $S_i$  reflect their similarity

<u>elements</u>	$S_1$	$S_2$	$S_3$	$S_4$	$\rightarrow$ docs.
a	1	0	0	1	$S_1 = \{a, d\}$
b	0	0	1	0	$S_2 = \{c\}$
c	0	1	0	1	$S_3 = \{b, d, e\}$
d	1	0	1	1	- Characteristic matrix.
e	0	0	1	0	in doc?

if permuted  
order  
is  $\rightarrow$

b	0	0	1	0	min hash value	$a   c   b   a$
e	0	0	1	0		
a	1	0	0	1		
d	1	0	1	1		
c	0	1	0	1		

based on where we see 1st one.

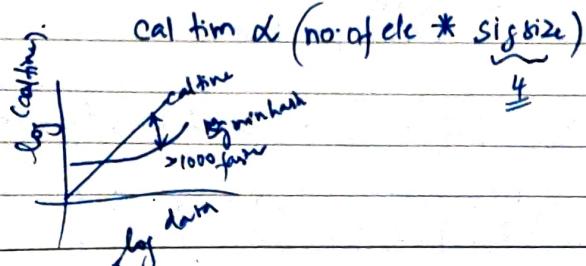
we see Min hash produces similar result as jaccard similarity

row	$S_1$	$S_2$	$S_3$	$S_4$	$\frac{h_1}{2+1 \bmod 5}$	$\frac{h_2}{3x+1 \bmod 5}$	signature matrix
0	1	0	0	1	1	1	$\begin{matrix} s_1 & s_2 & s_3 & s_4 \\ h_1 & 1 & 3 & 0 & 1 \\ h_2 & 0 & 2 & 0 & 0 \end{matrix}$
1	0	0	1	0	2	4	
2	0	1	0	1	3	2	
3	1	0	1	1	4	0	
4	0	0	1	0	0	3	

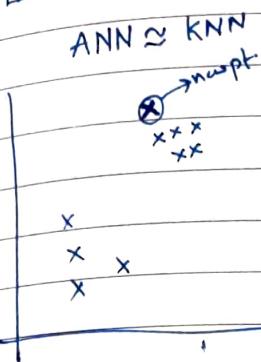
### Weighted min hash.

	$h_1$	$h_2$
1	1	1
2	2	2
0	0	0

smallest hash value of each col. reduce no. of direct



## ⑤ Approximate Nearest Neighbour Search:-



In KNN  $\rightarrow$  we find  $k$ -closest <sup>nearest</sup> data pt for the given pt. (physically in space)

Ex: 3 closest word for a given word.

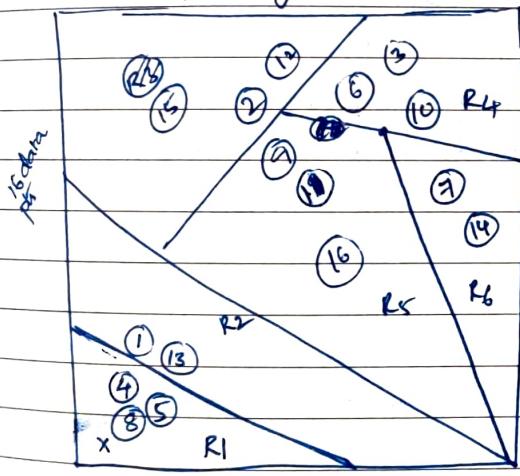
Task: Given a point  $(8)$   $\rightarrow$  find  $K$  closest Neighbours.

But...  $N$  points means  $O(n)$ .

i	0	1	2	$\rightarrow k$ closest if $k=3$
do	$d_1$	$d_2$		keep track of distance of pt's (using $sqrt$ )

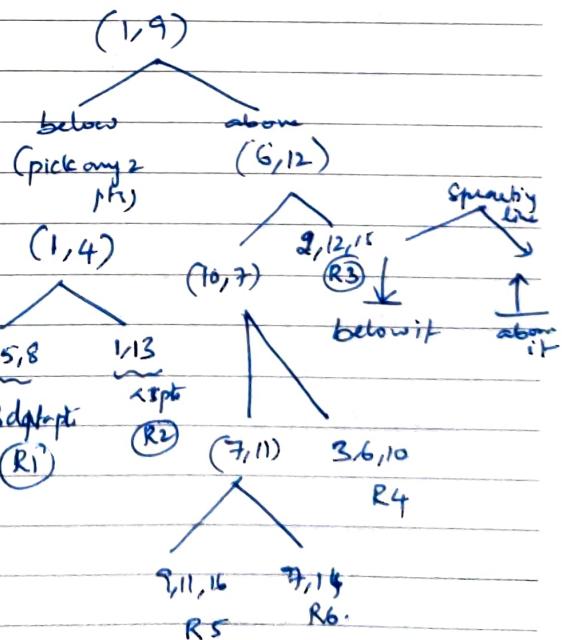
$\Rightarrow$  (distance of given pt & pts in space).

What if we got almost Closest Neighbours:-



(in fraction time).

1. pick any 2 data pts 1, 9 in random.
2. draw a line equi-distant from 2 points.



After finding where the new pt lies?  
 $\rightarrow$  (if 'x' data pt is given <sup>(in which region)</sup>)  
we can just find distance with  $R_i$  regions & find the nearest data point)

$O(\log n)$ .

Note: [Some errors if we take pt near pt(9).]

but we get approx. nearest neighbor.]

Instead of 1 tree  $\rightarrow$  construct several trees (Construct Forest)  
 So, result is based on several constructed trees.  
 & then return the nearest neighbour.

(So, if many trees are constructed  $\rightarrow$  more data structure storage)  
 but still worth it.

## ⑥ String Algorithms (Edit distance):-

### Levenshtein Distance [edit distance]

Our key

replace	insert
delete	you're here

	" "	b	e	n	y	a	m	
" "	0	1	2	3	4	5	6	be to form ... take 2 deletions
e	1	1	2	3	4	5	6	① b,e are diff.
p	2	2	2	2	3	4	5	② min (ins, del, 2, 2, 1)
h	3	3	3	3	3	4	5	(ins & del)
r	4	4	4	4	4	4	5	2, 2, 1
e	5	5	5	5	5	5	5	1 ins
m	6	6	5	5	6	6	5	0 del

↓  
 from " "  $\rightarrow$  need to form subproblem cph  
 $\Rightarrow$  3 invocations.

Time complexity :  $O(ab)$   
 Space " :  $O(ab)$ .

if both e, e are same  
 its same as  $(i-1, j-1)$  subproblem

$$\left[ \begin{array}{l} \text{if diff. letters} \\ \hline \end{array} \right] \min \left[ \begin{array}{l} (i-1, j), (i, j-1) \\ (i-1, j-1) \end{array} \right] + 1$$

## ⑦ Burrows Wheeler Transform (BWT):-

BANANA \$  
 $\underbrace{\quad}_{n}$

proc process  
 our database.

[So, given gquery  
 you can find where  
 the word occurs]

Improve file compression

AN  
 NANA  
 BAN  
 NA  
 $\underbrace{\quad}_{K \text{ size}}$

all rotations of sequence :-

BANANA\$

\$BANANA  
A\$BANAN  
NA\$BANA  
ANA\$BAN.  
NANA\$BA  
ANANAS\$B.

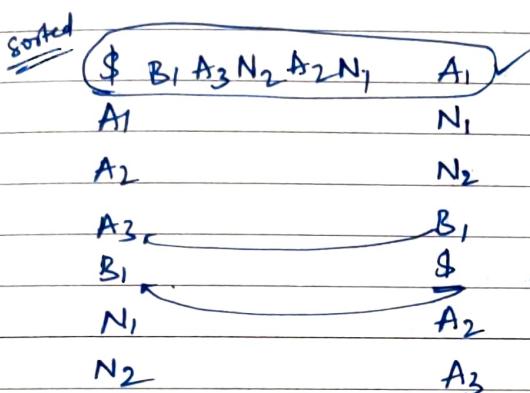
} Sort All rotations  
alphabetically

'\$' is smallest character  
assuming ↗

→ BWT transform  
is last col. of  
matrix

ANNB\$AA

↙ u can  
reconstruct  
the actual  
content of  
matrix



dollar followed by B<sub>1</sub>

pattern matching using the BWT :-

i	First	Last	L2F(i)	for AN
top → 0	\$ <sub>1</sub>	A <sub>1</sub>	1	↓ start here (base word)
1	A <sub>1</sub>	N <sub>1</sub>	5 ✓ → top	→ top A
2	A <sub>2</sub>	N <sub>2</sub>	6 ✓ → bottom	→ bottom A
3	<u>A<sub>3</sub></u>	B <sub>1</sub>	4	→ bottom A
4	<u>B<sub>1</sub></u>	\$	0	↓ 2,3 → bottom A
5	N <sub>1</sub>	<u>A<sub>2</sub></u>	2	try A <sub>2</sub> , A <sub>3</sub> where AN occurred
bott → 6	N <sub>2</sub>	A <sub>3</sub>	3	↓ 2,3 position occurred
Searched → NAN / AN				5 → top 6 → bottom
				N <sub>1</sub> N <sub>2</sub> position
↑ top, bottom would be 6 ⇒ 3			For NA	

## \* Suffix Tree / Arrays :-

Suffix → is a substring at the end of a string of characters.

ex: for horses → e

sc  
ore  
orse  
horse.

So, Suffix array is an array which contains all the "sorted suffixes" of a string.

- 0 camel
- 1 amel
- 2 mel
- 3 el
- 4 e.

1	amel
0	camel
3	el
4	e
2	mel.

↓

} sorted suffixes

→ So, actual 'suffix array' is array of sorted indices.

→ This provides a compressed representation of sorted suffixes without actually needing to store the suffixes.

→ Suffix array provides space efficient alternative to 'suffix tree' which itself is compressed version of tree.

Note: Suffix array can do everything suffix array can do., with some additional information such as (LCP) array

— longest common prefix. —