

string matching :

* String matching :

- finding the substring (pattern) in main string (Text)
- patter (P) and Text (T)
- size of P = m
size of T = n
- $x[z]$: string x is suffix of z
 $y[z]$: string y is prefix of z
- Algorithm of string matching
 - ① Naive - string matching
 - ② Rabin - karp
 - ③ KMP

① Naive - String matching :

- Also called Brute force string matching
- simple technique
- Complexity : $(n-m+1)n$
 $= O(nm)$

→ Algorithm :

T = Text

P = Pattern

then

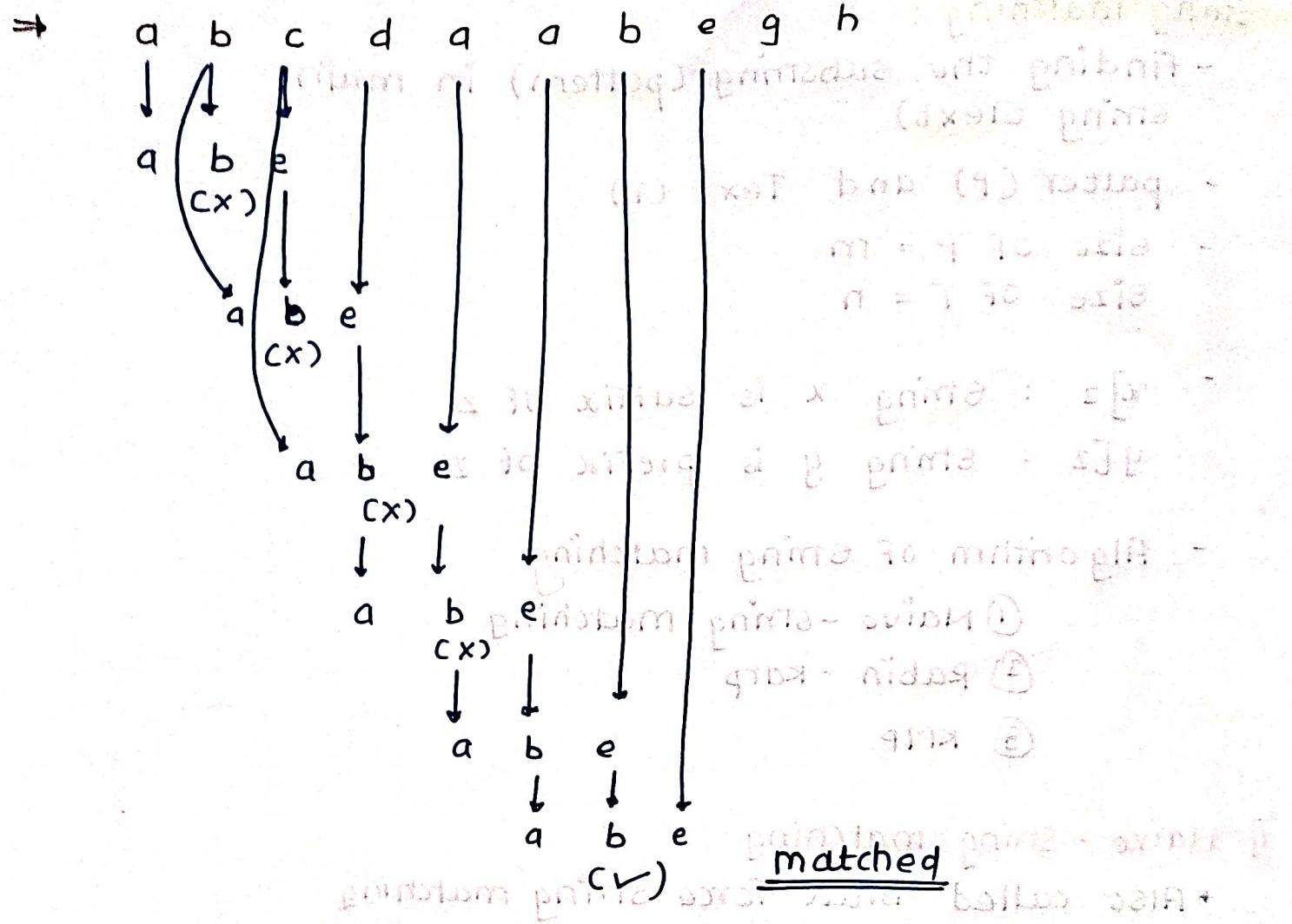
Naive_string_matcher (T, P)

- 1 $n = T.length$
- 2 $m = P.length$
- 3 for $i=0$ to $n-m$
- 4 if $P[i \dots m] = T[i+i \dots i+m]$
- 5 print "pattern match at ith shift"

Ex[±] : Average case

Text : a b c d a a b e g h

pattern : a b e

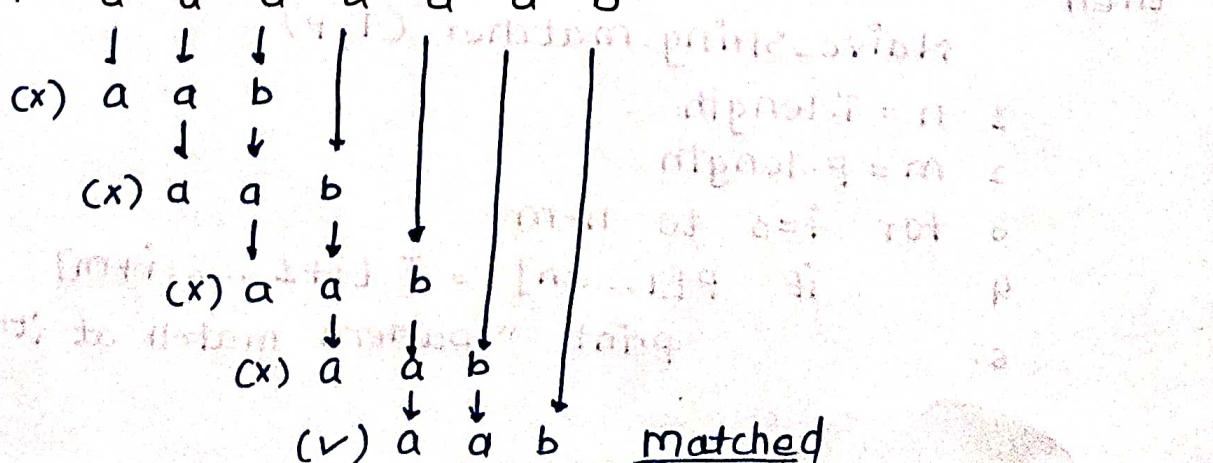


Ex 2 : Best case: pattern = a b c

Text : a b c d a b e a b
 ↓ ↓ ↓
 a b c
 (✓) matched

Ex 3: Worst case: pattern = a a b c d e f g

Text : a a a a a a b



- ② Rabin Karp algorithm:
- String matching algorithm
 - Uses hashing technique to find pattern in the text.
- Terms

Hash function converts a substring of the text into a hash value.

Rolling hash function

updating hash value for each substring by considering the previous hash value.

- There are several way to calculate hash value

[i] way 1:

steps:

- 1) maps each character with no. in text and pattern
(a=1, b=2, c=3 ... z=26)

2) Add all digit assigned to pattern (eg. abe)

Hash function {
(a=1, b=2, e=5)

$$1+2+5 = 8 \quad [\text{Hash value}]$$

3) Now start calculating hash value using same tech. for m-character subsequence of Text.

4) if hash value calculated matches with hash value calculated at step 2
then: pattern match

else

calculate hash value for next m-character sub sequence of Text by shifting 1 to right side.

- Complexity for avg. cases = $O(n-m+1)$

complexity due to drawback = $O(nm)$

Drawback: Supurious hit occurs where hash value of pattern matches with hash value of m-character subsequence of text but actual pattern is different.

i) Ex : Text : a b c d a b c e
 pattern : a b e

→ Text : a b c d a b c e
 (1) (2) (3) (4) (1) (2) (3) (5)

pattern : a b e
 (1) (2) (5)

∴ Hash value of pattern = $1+2+5=8$

→ Text : a b c d a b e e
 (1) (2) (3) (4) (5) (2) (3) (5)

$\times 6$ (1 + 2 + 3) and 7 times 10⁰
 $\times 9$ (2 + 3 + 4) and 7 times 10¹
 \checkmark (8) 3 + 4 + 1 and 7 times 10²

∴ Here hash value matches

but if we check actual pattern
 (a b e) existing at c d i a \neq a b e
 (m-char) — spurious hit

$\times 7$ 4 + 1 + 2 and 7 times 10⁰
 $\times 8$ 1 + 2 + 5 and 7 times 10¹

a b e = a b e (matched)

[ii] way 2

i) map each char with no. (present in Text & pattern)
 $(a=1, b=2 \dots z=26)$

ii) Hash function:

$$\Rightarrow P[1] \times 10^{(m-1)} + P[2] \times 10^{m-2} + \dots + P[m] \times 10^0$$

$$\text{ex : } a \ b \ c = 1 \ 2 \ 3$$

$$\Rightarrow 1 \times 10^2 + 2 \times 10^1 + 3 \times 10^0 = 123 \quad [\text{Hash value}]$$

iii] Now in same fashion generate hash value for m-character subsequence of Text
 if value match
 - pattern match

else
 - check for next m-char subsequence
 by shifting 1 to the left.

ex Text : c c a c c a d b a c
 pattern : d b a

→ Hash value for pattern: $\begin{matrix} d & b & a \\ (c_4) & (c_2) & (c_1) \end{matrix}$
 $4 \times 10^2 + 2 \times 10^1 + 1 \times 10^0 = 421$

- i) cca = $3 \times 10^2 + 3 \times 10^1 + 1 \times 10^0 = 331$
- ii) cac = $3 \times 10^2 + 1 \times 10^1 + 3 \times 10^0 = 313$
- iii) acc = 133
- iv) cca = 331
- v) cad = 314
- vi) adb = 142
- vii) dba = 421 ✓ match dba = dba
- viii) bac =

∴ complexity = $O(ntm)$

iii] way 3:
 → Same as way 2 but instead of hash function
 use **rolling hash function** to generate hash
 values. (i.e. use previous hash values)

ex: Text : c c a c c a d b a c
 pattern : d b a

→ Hash value for pattern

$$\begin{matrix} d & b & a \\ 4 & 2 & 1 \end{matrix}$$

$$= 4 \times 10^2 + 2 \times 10^1 + 1 \times 10^0 = 421$$

- (3) ccq = $3 \times 10^2 + 3 \times 10^1 + 1 \times 10^0 = 331$
 (3) cac = $(331 - 3 \times 10^2) \times 10 + 3 \times 10^0 = 313$
 (1) acc = $(313 - 3 \times 10^2) \times 10 + 3 \times 10^0 = 133$
 (3) cad = $(133 - 1 \times 10^2) \times 10 + 1 \times 10^0 = 331$
 (3) ccq = $(331 - 3 \times 10^2) \times 10 + 4 \times 10^0 = 314$
 (2) dbq = $(314 - 3 \times 10^2) \times 10 + 2 \times 10^0 = 142$
 (1) bac = $(142 - 1 \times 10^2) \times 10 + 1 \times 10^0 = 42$ ← matched

Here rolling hash function is used.

- iv] As we are taking power, so we may get very large values in some cases, so to avoid such cases we can use mod operation while hash value calculation.

- ③ KMP algorithm:
- knuth-Morris-Pratt algorithm for string matching
 - It avoids unnecessary character comparisons by utilizing information from previous matches.
 - We need to calculate prefix table using prefix function.
 - prefix function(n) calculates the length of the longest proper prefix that is also proper suffix of given substring.

Prefix table example : using Prefix function (π)

①	string :	a	b	a	b	d
		0	0	1	2	0

- i) a x
 ii) ab = a b x
 iii) ab a = ab a ba x

- iv) abab = a d ✓
 v) ababd = a d x
 vi) ab ab ✓ ab ab ✓ ②
 vii) aba bab x

(2)

a	b	q	b	a	cd
0	0	1	2	3	0 1

- 1) $a \rightarrow x$
- 2) $ab \rightarrow a b x$
- 3) $aba \rightarrow a a \checkmark \textcircled{1}$
 $ab ba x$
- 4) $abab \rightarrow a b x$
 $ab ab \checkmark \textcircled{2}$
 $ab ab ab x$

- 5) $ababa \rightarrow a a \checkmark$
 $ab ba x$
 $aba aba \checkmark \textcircled{3}$
 $abab baba x$

6) $abab ac \Rightarrow a c x$
 $ab ac x$
 $ab a bac x$
 $abab abac x$
 $abab a babac x$

7) $ababaca \Rightarrow a a \checkmark \textcircled{1}$
 $ab ca x$
 $ab a aca x$
 $ab ab baca x$
 $ababa abaca x$
 $ababac baba cax$

Steps for KMP algo :

- 1] For given pattern calculate prefix table (preprocessing)
- 2] use two pointers (initialization)
- for Text for pattern
 $(i) = 1$ $(j) = 0$
- 3] Matching : compare characters : while iterating through Text (T) and Pattern (P)
- If characters match increment both i & j
 - If character mismatch update j using info from table. (it avoids unnecessary char comparison)

Example : a b a b c a b c a b a b d (Text)
 ab a b d (pattern)

\Rightarrow i) Prefix table:

a	b	a	b	d
0	0	1	2	0

Text : $\begin{array}{ccccccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 & 15 \\ \text{T} & \text{a} & \text{b} & \text{a} & \text{b} & \text{c} & \text{a} & \text{b} & \text{a} & \text{b} & \text{a} & \text{b} & \text{a} & \text{b} & \text{c} \end{array}$

(T) $P=1$

Pattern : $j=0$

(P)

a	b	a	b	c	a	b	a	b	d	a	b	a	b	c
$\pi \rightarrow 0$	0	1	2	0										

if $T[i] == P[j+1]$
then $i++$
 $j++$

$\Rightarrow T[0] = p[0+1]$

$a = a \checkmark$

$T[1] = p[1+1]$

$b = b \checkmark$

$T[2] = p[2+1]$

$a = a \checkmark$

$T[3] = p[3+1]$

$b = b \checkmark$

$T[4] = p[4+1]$

$a = a \checkmark$

$T[5] = p[5+1]$

$c = a \times \text{if } j=0$

$T[6] = p[6+1]$

$b = b \checkmark$

$T[7] = p[7+1]$

$a = a \checkmark$

$T[8] = p[8+1]$

$a = a \checkmark$

$T[9] = p[9+1]$

$b = b \checkmark$

$T[10] = p[10+1]$

$a = a \checkmark$

$T[11] = p[11+1]$

$b = b \checkmark$

$T[12] = p[12+1]$

$a = a \checkmark$

$T[13] = p[13+1]$

$d = d \checkmark$

$T[14] = p[14+1]$

$d = d \checkmark$

$T[15] = p[15+1]$

$c = c \checkmark$

else if $C_{\pi(j)} == 0$

$j = \pi[j]$

$i++$

$j = \pi[j]</math$

Convex polygon :

- If you connect any two points within a polygon, all the points on that line are also within the polygon.

Problem statement :

Given set of points (n - with x and y co-ordinates), find whether the polygon defined is convex?

- order of points is imp
- finding if the given point is within the polygon ?

Solution

- 1] i] Forming triangles :
 - forming triangles from given set of n points
 - there can be several triangles formed by connecting combinations of three points.

2] Convexity check :

- check triangles are convex or not?

3] Assumption :

- soln assume that checking whether a point lies within a triangle takes 1 unit of time

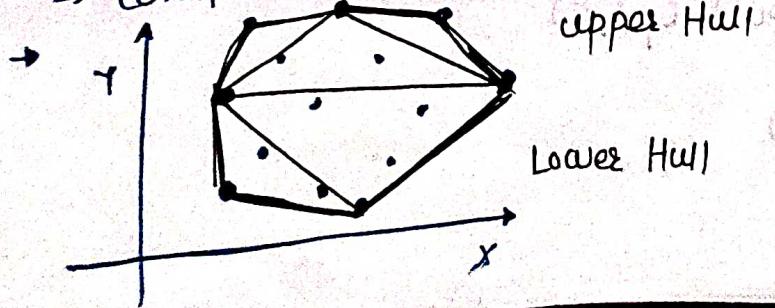
complexity : for forming all possible triangles = $O(n^3)$

- for checking if a point lies within it = $O(n)$

$$\therefore \text{Total} = O(n^4)$$

- 2] Divide and Conquer (D and C) Method :

- Given a set of points we need to find convex Hull
- complexity : $O(n \log n)$



[smallest convex polygon that contains all the points in the plane]

Solution approach :

- 1) select extreme points (left, right)
- 2) Join left and right point (so that line joining those points divides the space in lower/Upper hull.)
- 3) Upper Hull :
 - select triangle formed by these two points & third point (so that area will max .)
 - Now consider two small hulls above the sides of triangle
 - Repeat till no point will be outside.

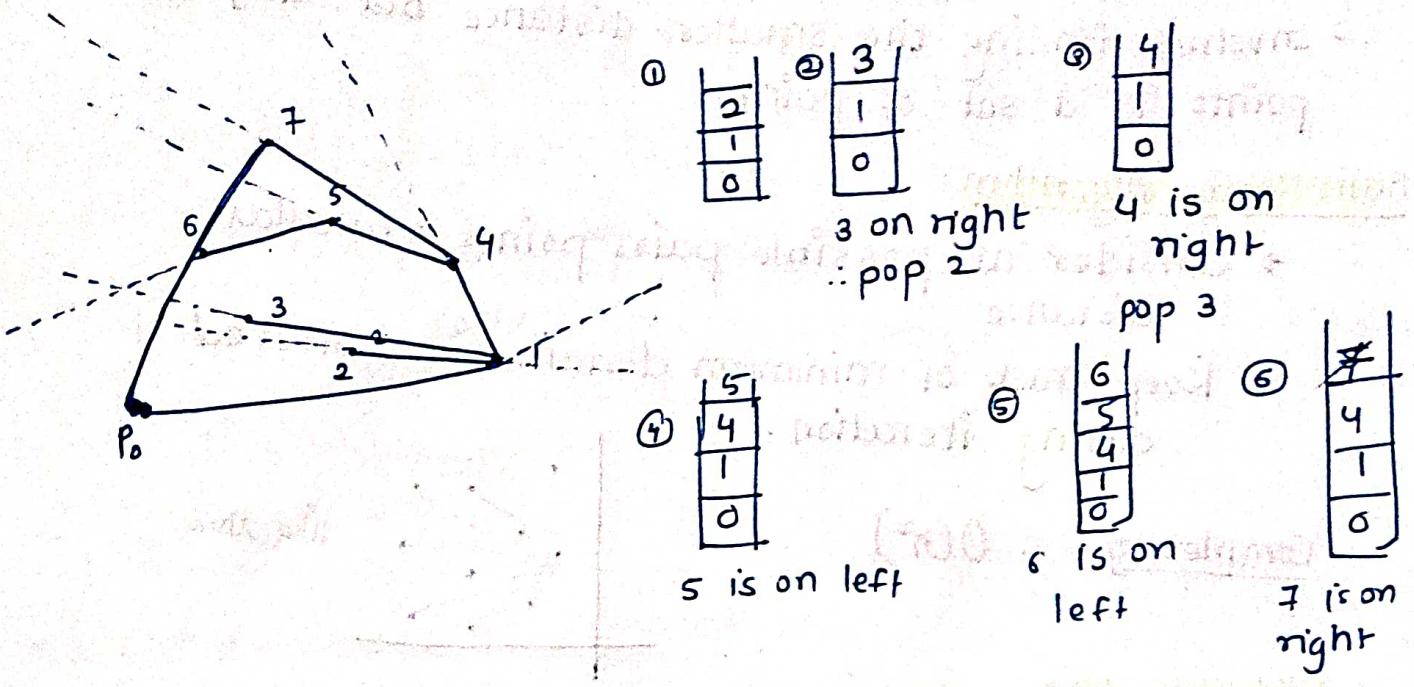
- 4) Repeat same for lower hull.

Other methods :

I] Graham scan :

Algo :

- 1) Find (leftmost & rightmost point) P_0 .
- 2) for remaining points, sort by polar angle in counter clockwise around P_0 .
- 3) Let s is empty stack
- 4) push (P_0, s)
- 5) push (P_1, s)
- 6) push (P_2, s)
- 7) for $i=3$ to n — as 3 points are already
done
if (p_i is left of top of stack)
then push (p_i, s)
else
pop (s)
- 8) push (p_i, s)
- 9) return s .



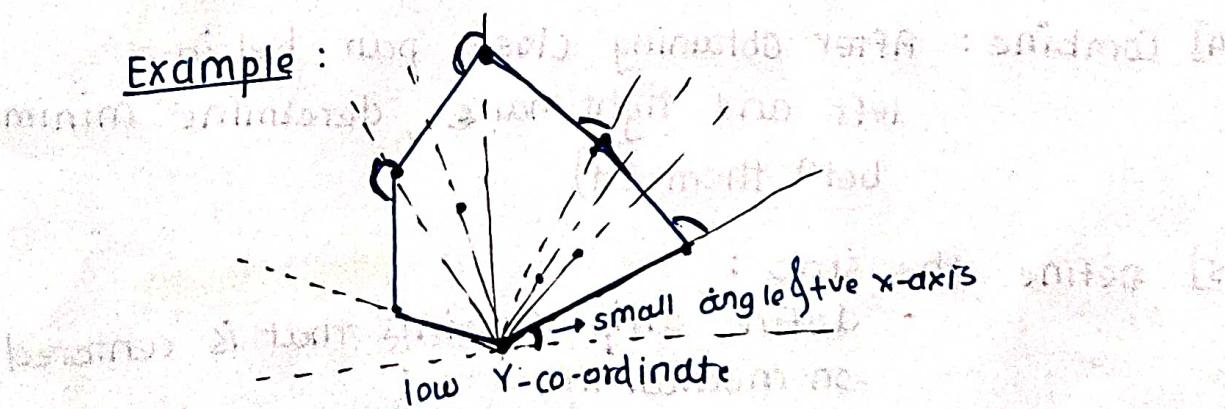
Complexity : $n \log n$

2] Jarvis march :

Algo : Note lowest & highest point

- start from lowest (point with low Y-co-ordinate)
- At every point search for point which makes lowest polar angle with +ve x-axis

Example :



complexity $n * h$ - points defining convex polygon

no. of steps = no. of vertices

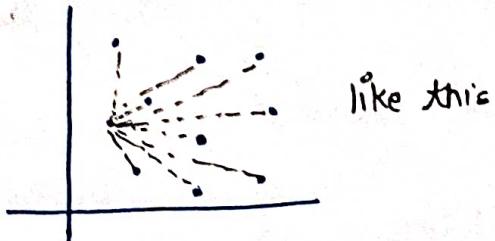
Find closet pair of points :

- involves finding the smallest distance bet'n any two points in a set of points.

Brute force Algorithm:

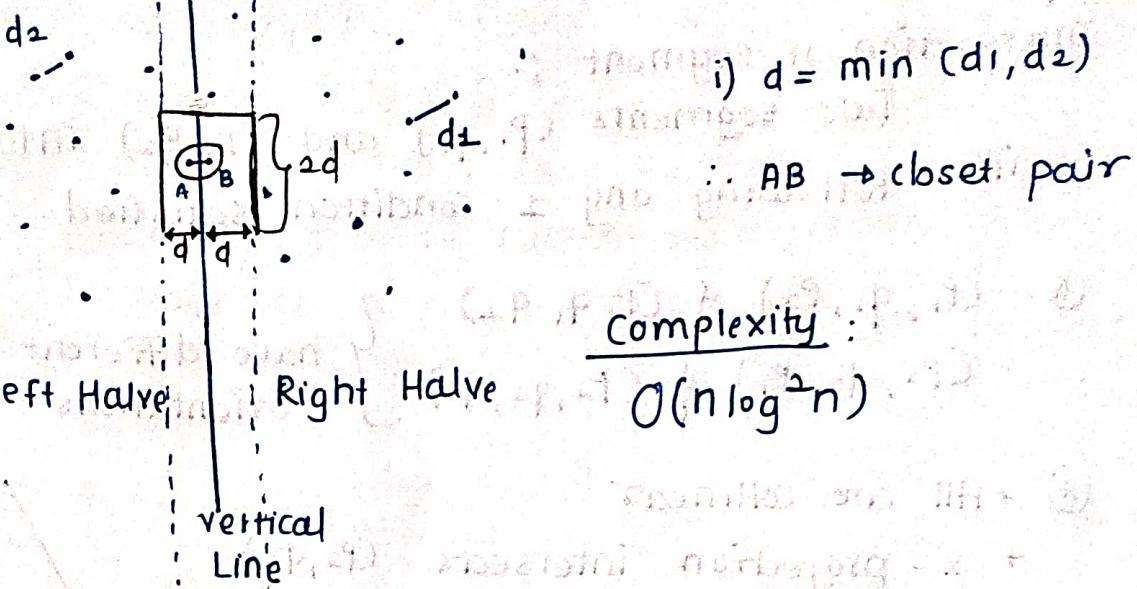
- Consider all possible pairs points : calculate distance
- Keep track of minimum distance encountered during iteration.

Complexity : $O(n^2)$



Divide and Conquer:

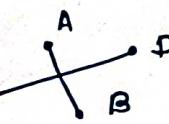
- 1] Sort points : sort points based on their x-coordinates
- 2] Divide : Divide set of points into two halves along the vertical line passing through x-coordinate.
- 3] Conquer : find Recursively, the closet pair of points in each halve. (By apply same D&C tech)
- 4] Combine : After obtaining closet pair in left and right halve, determine minimum bet'n them. (d)
- 5] Define the strip :
 - define strip of points that is centered on median line.
 - It is vertical band that includes points with min distance found.
 - check points that might form a closer pair (which one point lies in left & other in right halve)
 - update min distance accordingly. (if new distance $< d$)



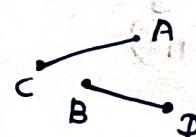
• Pair of Intersecting Segments:

Terminologies :

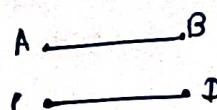
1] Intersecting :



2] Non intersecting :



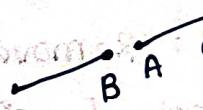
3] Parallel :



4] Colinear / Overlapping :

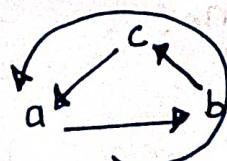


5] Colinear :

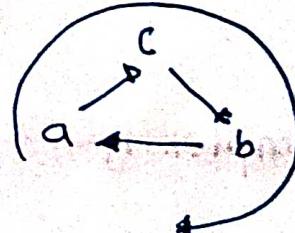


Orientation

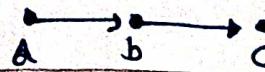
1] Counterclockwise



2] Clockwise



3] Colinear



* So,

Intersection of Segment :-

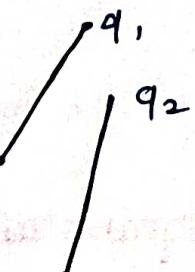
Two segments (P_1, q_1) and (P_2, q_2) intersect iff following any 1 condition satisfied.

A $(P_1, q_1, P_2) \nparallel (P_2, q_1, q_2)$ } have different
 $(P_2, q_2, P_1) \nparallel (P_2, q_2, q_1)$ } orientations

B → All are collinear

→ x-projection intersects (P_2, q_2)

→ Y-projection intersects (P_2, q_2)



Approach

① Brute force Approach:

→ for every line segment check if they intersect.

→ Complexity : $O(n^2)$

using above conditions A & B

② Sweepline Algo:

→ Process events as a sweep line moves from L to R.

Events :

i) Insert : Add to active set, if left endpoint

ii) Remove / Delete : Remove from set, if Right end

iii) Intersect : check for intersection b/w current segment & those are in active set.

Active set :

maintains segments intersecting with sweep line.

→ Complexity : $O(n \log n)$

Steps :

- ① sort all endpoints of seg. L-R
- ② - Add seg to active set \Rightarrow left end
- Delete seg from active set \Rightarrow Right end
- Check intersection of current seg with seg present in active set

