

Note:- (why we square things?)  
→ absolute error captures the error value  
but not the variation in the error  
vs. squared error that captures  
variation in the error.

S-I :- D1: chapter 1, D2: chapter 2, D3:

Self-learning  
Based. agents

ANN :- (Artificial Neural Network)

- ① Underfitting ≠ Overfitting
- ② ANN architecture
- ③ Back propagation Algorithm

Training data:- The data containing examples to learn the parameters of a classifier.  
e.g. OR data to learn weight values.

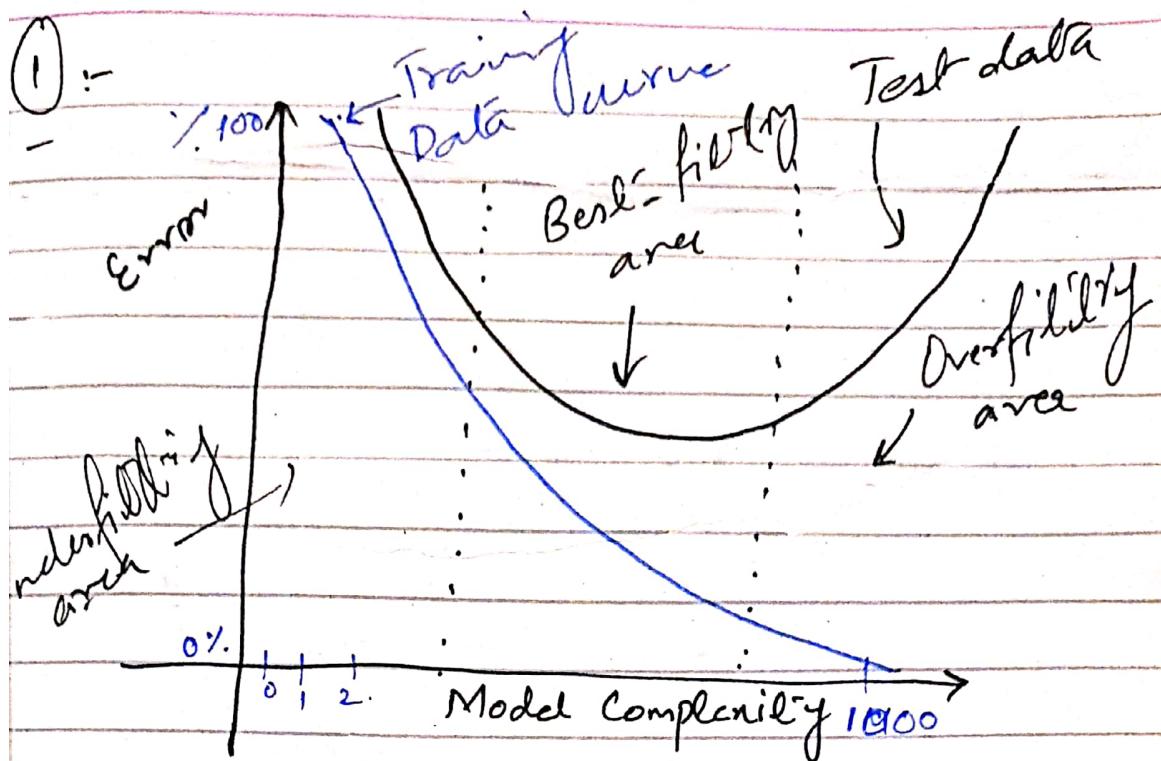
Test data:- the data containing previously unseen examples to test a classifier's performance.

Training

	$x_1$	$x_2$	$f_d$
$E_1$	0	0	0
$E_2$	0	1	1

Test

	$x_1$	$x_2$	$f_d$
$E_1$	1	0	?
$E_2$			

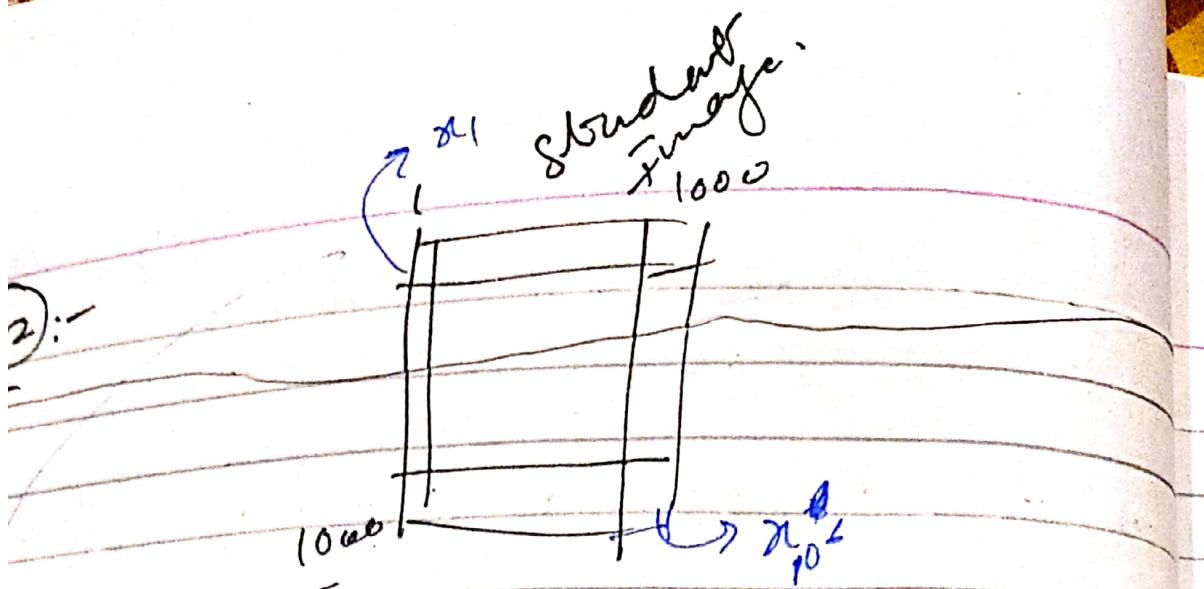


→ No. of perceptions ↑  
 No. of weights (features) ↑

Underfitting :- (Need to increase complexity)  
 High error on the  
 training & the test data.

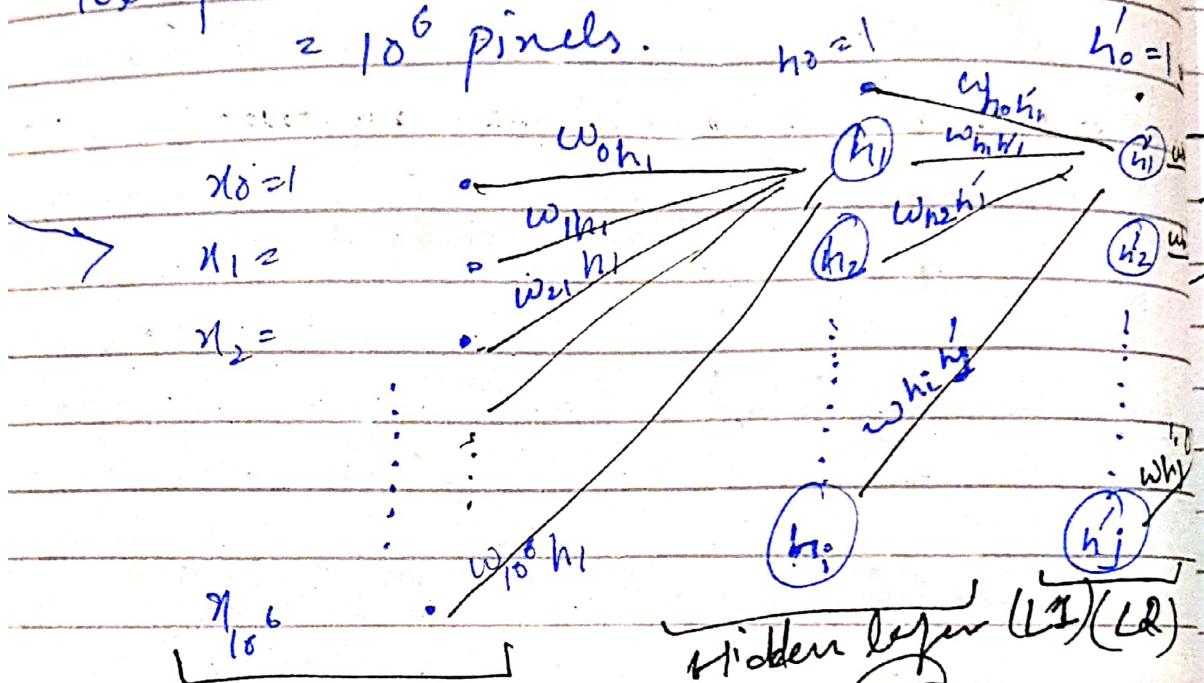
Overfitting :- (Need to reduce complexity)  
 Low error on the  
 training data & high error on  
 the test data.  
 Reduce or increase the training data.

Best-fitting - Acceptable error  
 on both the training & test  
 data.



Total pixels = 1,000,000

$= 10^6$  pixels.



Input layer  $\rightarrow$  data + 1

Bias.

→ Why Bias?

In case of a suppressed signal, we need to elevate the volume.

\* → Almost all the time only 2 to 3 hidden layers are sufficient.

Right

$0_1 \ 0_2 \ 0_3 \ 0_4 \ 0_5 \ 0_6$       chris  
 0 0 0 0 0 1 — 19P-0001  
 0 0 0 0 0 1 0 — 19P-0002

① 0/1

② 0/1

⋮

⑥ 0/1

Output Layer

0's are 6

Because we have 60 different bits to classify 55 among 35 students, so we need

\* → So, when image of chris is given the output will be —

①	— 0
②	— 0
③	— 0
④	— 0
⑤	— 0
⑥	— 1

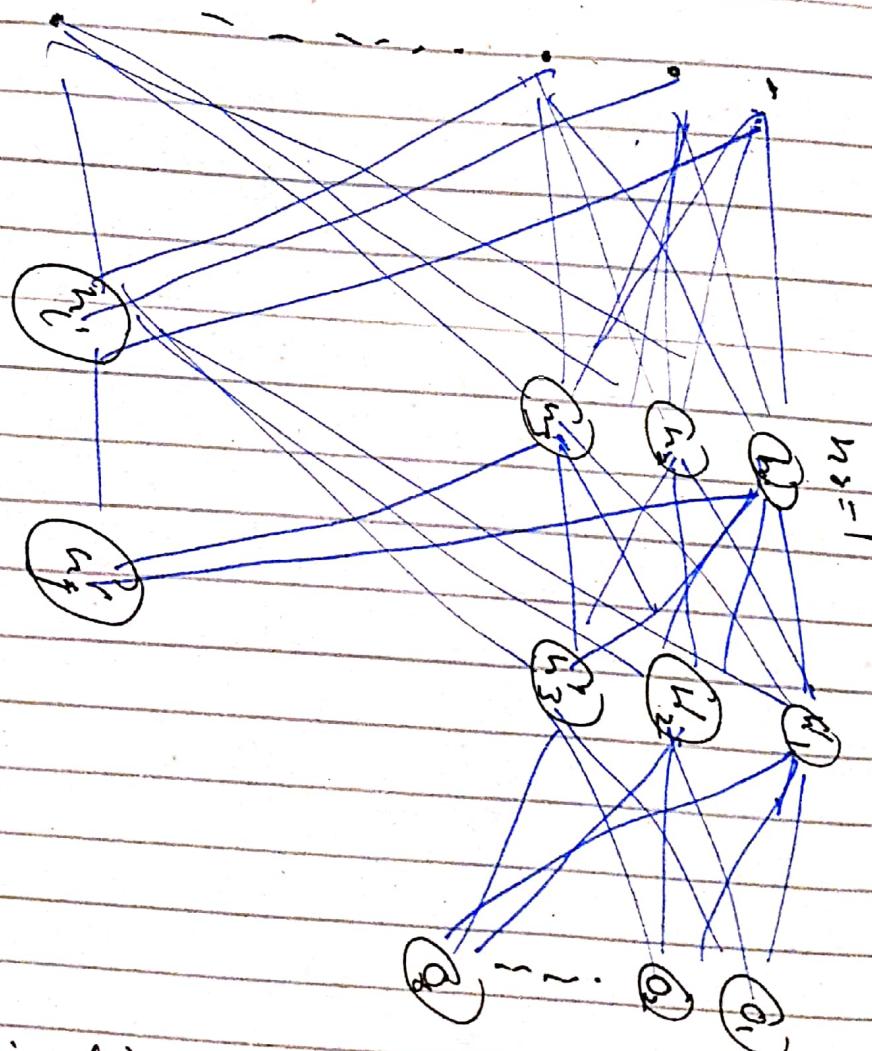
→ 6 bits (64 bits) using more bits will result in overfitting.

→ Perceptron count of Hidden layer & Layer counts are dependent on the Over, underlying graph. So start with low & keep observing.

ANN  
complete architecture.

In same algorithm but we use  
back propagation, by propagating  
from output to see how  
each hidden perception affects/contributes

$$\text{of } w_{oh} = -\eta \cdot \frac{\partial E}{\partial w_{oh}}$$



validation dataset is used to  
fine-tune our hyper parameters.  
All that are representative of your  
purpose.

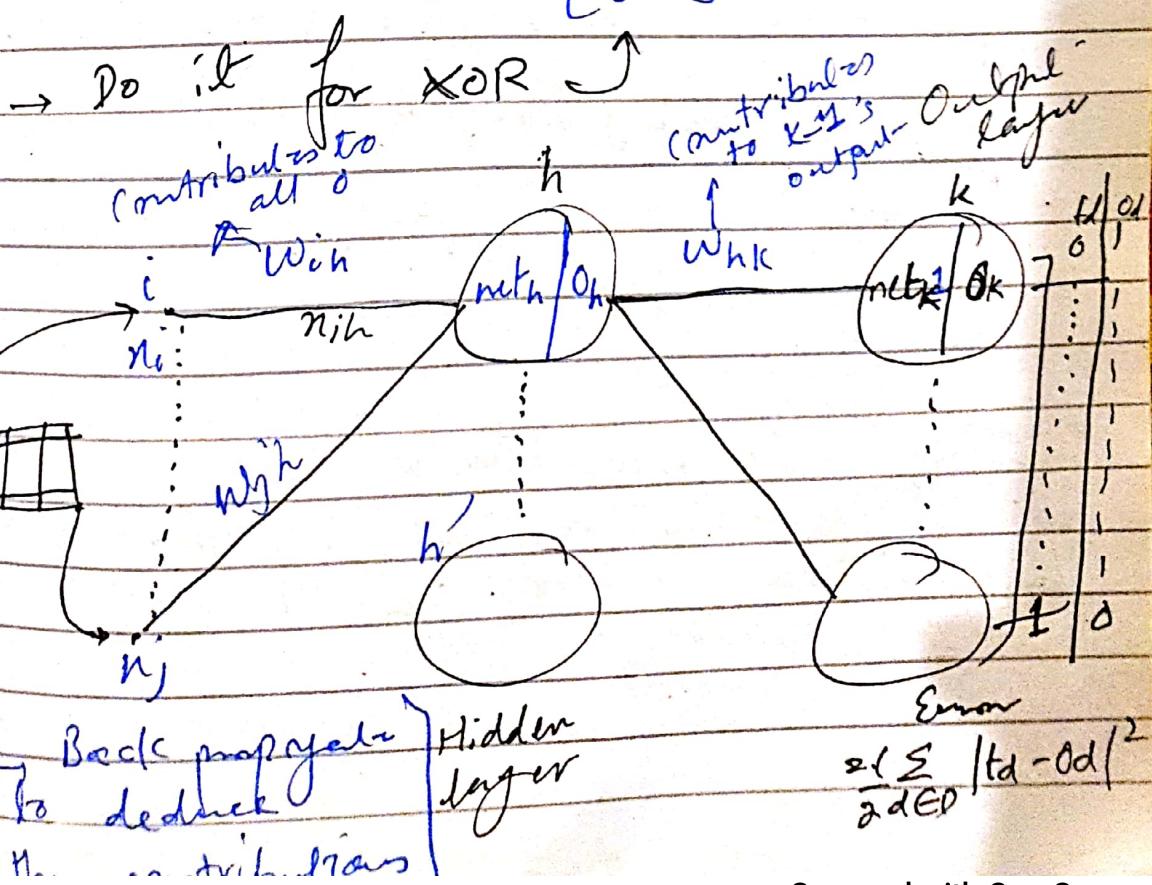
# Artificial Neural Networks

## Back propagation Algorithm

- ① Initialize all weights randomly i.e  $w_{ij}$
- ② Until convergence, Do
- ③ Input it to network & compute network outputs.
- ④ For each output unit  $k$ .  
 $\delta_k \leftarrow o_k (1 - o_k) (t_k - o_k)$
- ⑤ For each hidden unit  $h$ .  
 $\delta_h \leftarrow o_h (1 - o_h) \sum_{k \in \text{outputs}} w_{hk} \delta_k$
- ⑥ Update each network weight  $w_{ij}$

$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

where  $\Delta w_{ij} = \eta \delta_j x_{ij}$



Weight update for output layer.

$$\Delta w_{hk} = \eta \cdot \frac{\partial E}{\partial w_{hk}} = \eta \left( \frac{-\partial E}{\partial \text{net}_k} \right) \frac{\partial \text{net}_k}{\partial w_{hk}}$$

$\delta_k^2 \circ_k$

$$\delta_k \circ_k = \frac{\partial E}{\partial \text{net}_k} = \frac{\partial}{\partial \text{net}_k} \sum_d (t_d - o_d)^2$$

$$= \sum_d (t_d - o_d) \cdot \underbrace{(o - \cancel{o_d})}_{\partial \text{net}_k}$$

$$\Rightarrow \quad \quad \quad (o - o_d(1 - o_d))$$

$$\frac{-\partial E}{\partial \text{net}_k} \rightarrow \frac{-}{\cancel{o}} \quad o_d(1 - o_d)$$

$$\circ_k = \underbrace{\left( \sum_d (t_d - o_d) o_d (1 - o_d) \right)}_{\text{for all training examples}}$$

$$\delta_k \leftarrow \circ_k (1 - \circ_k) (t_k - \circ_k)$$

\* This is the derivative of sigmoid,  
so it can be replaced by other  
functions' derivatives such as relu etc.

→ for each  $b_k$  in the input example once you get  $\sigma_i \rightarrow \sigma_b$  then you can back propagate to find out contribution of error - in hidden layer.

→ See back that:-

$$net_k = w_{hk} \cdot o_h + \dots + w_{jk} \cdot o_j$$

$$\frac{\partial net_k}{\partial w_{hk}} = o_h$$

$$\Delta w_{hk} = \eta \left( \frac{-\Delta E}{\partial net_k} \right) \cdot \frac{\partial net_k}{\partial w_{hk}}$$

$\downarrow$                        $\downarrow$   
 $o_k$                        $o_h$   
 $\uparrow$   
 $o_k (1-o_k) (t_k - o_k)$

Now for :- (weight update for hidden layer)

$$\Delta w_{ih} = -\eta \cdot \frac{\partial E}{\partial w_{ih}} = \eta \left( \frac{-\Delta E}{\partial net_h} \right) \cdot \frac{\partial net_h}{\partial w_{ih}}$$

$\downarrow$                        $\downarrow$   
 $o_k$                        $x_{ih}$

$$\frac{\partial E}{\partial net_h} = \frac{1}{n} \sum_{k=1}^n (t_k - o_k)^2$$

$$= \frac{\partial E}{\partial o_k} \cdot \frac{\partial o_k}{\partial net_h} = \frac{\partial E}{\partial net_k} \cdot \sum_{k=1}^n \frac{\partial net_k}{\partial net_h}$$

$$\frac{\partial E}{\partial w_{hk}} = \frac{\partial E}{\partial o_k} \cdot \sum_{\text{out}_k} \frac{\partial \text{net}_k}{\partial o_h} \cdot \frac{\partial \text{net}_k}{\partial w_{hk}}$$

$$\frac{\partial E}{\partial w_{hk}} = \sum_{\text{out}_k} \frac{\partial E}{\partial o_k} \cdot \frac{\partial \text{net}_k}{\partial \text{net}_h}$$

$$= \sum_{\text{out}_k} (-\delta_k) \cdot \frac{\partial \text{net}_k}{\partial o_h} \cdot \frac{\partial \text{net}_k}{\partial w_{hk}}$$

$$= (-\delta_k) \cdot (w_{hk}) \cdot \delta_o (1 - \delta_o)$$

$$\delta_h = \delta_o (1 - \delta_o) \cdot \sum_{\text{out}_k} \delta_k \cdot w_{hk}$$

$$\Delta w_{hk} = \eta \cdot \delta_h \cdot n_{ip}$$

$\uparrow$  Backpropagation

Output:

Hidden, ...,  $\rightarrow$  hidden

end of  
forward pass  
for removal of network

Weight Update for Hidden Layer:-

$$\text{Error} = \frac{1}{2} \sum_{d \in D} \sum_{j=1}^{k'} (t_j^{(d)} - o_j^{(d)})^2$$

$$\sigma_k = -\frac{\partial E}{\partial \text{net}_k}; \quad \sigma_h = -\frac{\partial E}{\partial \text{net}_h}$$

$$\Delta w_{ih} = -\eta \cdot \frac{\partial E}{\partial w_{ih}} = \eta \cdot \left( \frac{-\partial E}{\partial \text{net}_h} \right) \cdot \frac{\partial \text{net}_h}{\partial w_{ih}}$$

$$\frac{\partial E}{\partial \text{net}_h} = \frac{\partial}{\partial \text{net}_h} \left( \frac{1}{2} \sum_{d \in D} \sum_{j=1}^{k'} (t_j^{(d)} - o_j^{(d)})^2 \right)$$

$$\frac{\partial E}{\partial \text{net}_h} = \sum_{K \in \text{outputs}} \frac{\partial E}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial \text{net}_h}$$

$$= \sum_{K \in \text{outputs}} \frac{\partial E}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial o_h} \cdot \frac{\partial o_h}{\partial \text{net}_h}$$

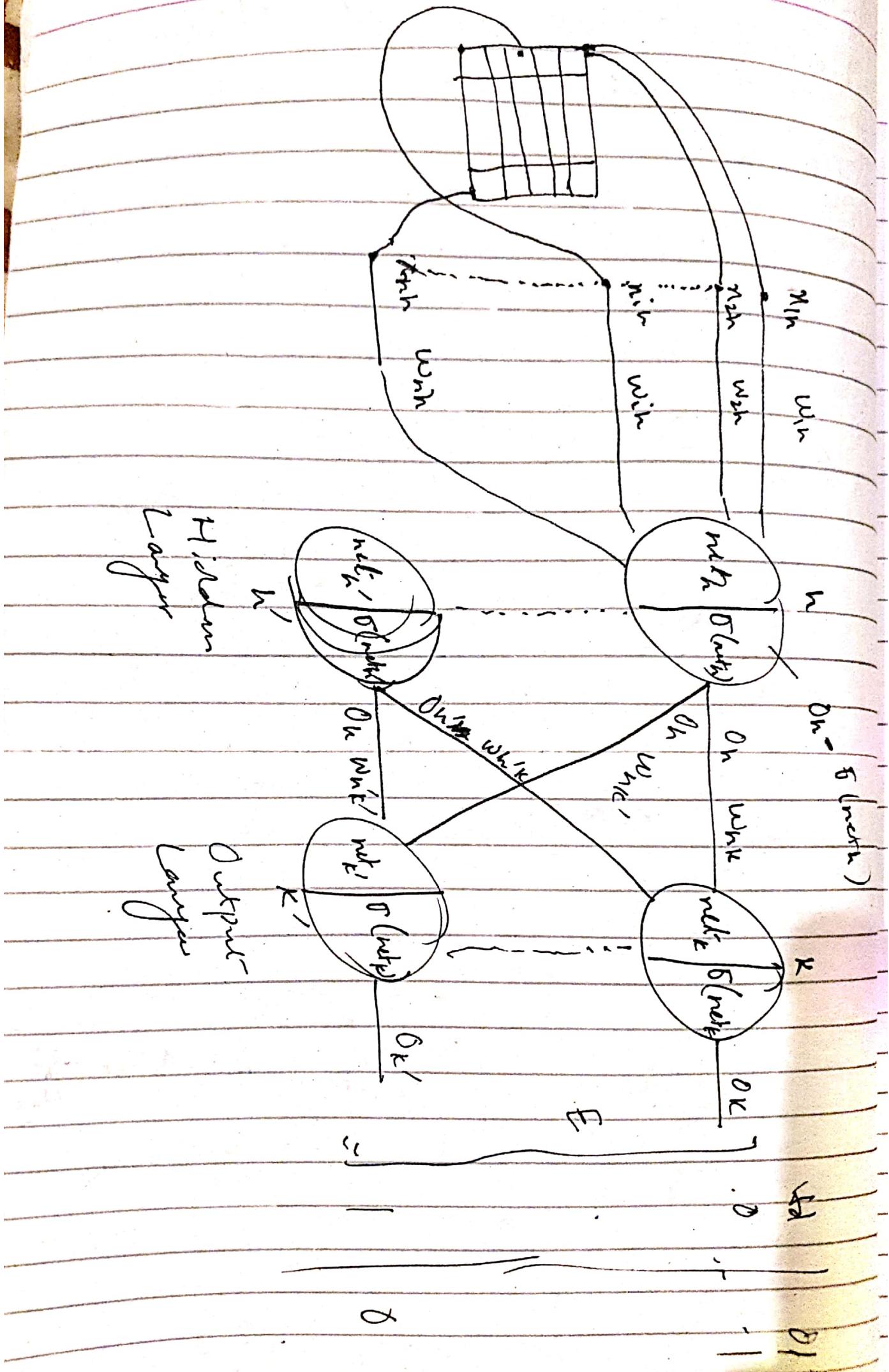
$$= \sum_{K \in \text{outputs}} (-\sigma_k) \cdot w_{kh} \cdot o_h(1-o_h)$$

$$\Delta w_{ih} = +\eta \cdot \frac{\partial E}{\partial \text{net}_h} \cdot \frac{\partial \text{net}_h}{\partial w_{ih}}$$

$$\sigma_h = \sum_{K \in \text{outputs}} \sigma_k \cdot w_{kh} \cdot o_h(1-o_h)$$

fractional  
gradient

$$\Delta w_{ih} = \eta \cdot \sigma_h \cdot \sigma_{ih}$$



# Artificial Intelligence:-

## Hypothesis Spaces:-

\* Hypothesis: A hypothesis is a function that approximates the true function (ground truth). by minimizing error.

→ If hypothesis  $f_h = \text{true } f_{gt}$  for all data, then error = 0%.

e.g OR,

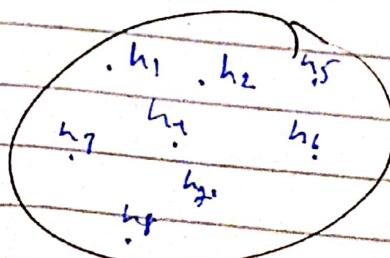
$x_1$	$x_2$	$t$	$h_1$	$h_2$
0	0	0	1	0
0	1	1	0	1
1	0	1	1	1
1	1	1	1	1

50% 0%.

→ Every eg to predict the target data is a hypothesis.

## Hypothesis Space:-

= is a huge space containing different hypothesis. ( $N$  number)



\* → Making a hypothesis is hard from the entire hypothesis space that is optimal & minimizes the error.

e.g. - for a binary one feature,

$n_1$	$h_1$	$h_2$	$h_3$	$h_4$
0	0	0	1	1
1	0	1	0	1

→ If target variable is binary, 4 hypothesis are possible in the hypothesis space -

for 2 features:-

$n_1$	$n_2$	$h_1$	$h_2$	$h_3$	...	$h_{16}$
0	0	0	0	0	...	1
0	1	0	0	0	...	1
1	0	0	0	1	...	1
1	1	0	1	0	...	1

General:-

$n_1$	$n_2$	$n_3$	$h_1$	$h_2$	...	$h_{2^k}$
0	0	0				
0	0	1				
0	1	0				
0	1	1				
1	0	0				
1	0	1				
1	1	0				
1	1	1				
⋮	⋮	⋮	⋮	⋮	⋮	⋮

\* Hypothesis space size for binary features:

$$\text{Hypothesis space size} = 2^{\text{# of features}} \rightarrow \begin{array}{l} \text{size of truth table} \\ 2^2 = 2^1 = 4 \\ 2^3 = 2^2 = 16 \\ 2^8 = 2^3 = 256 \end{array}$$

possible hypothesis ↪

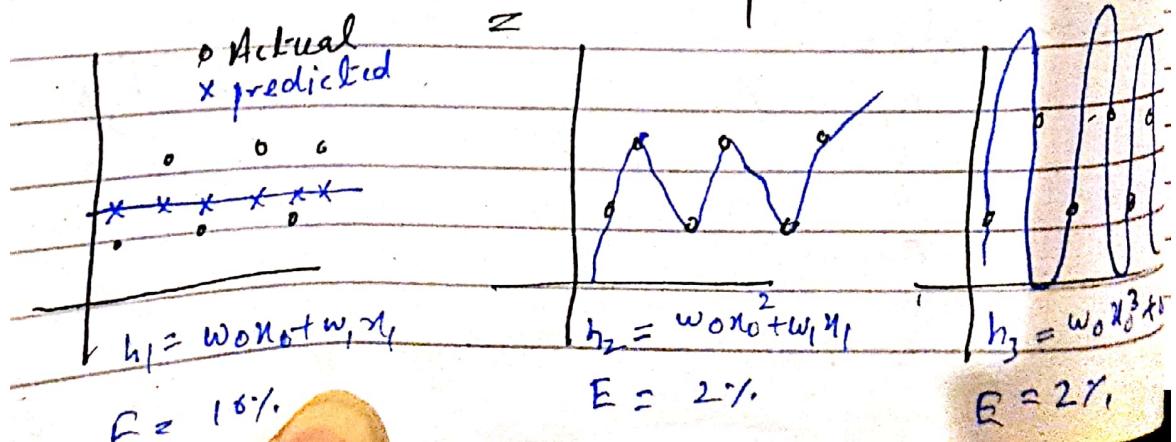
Generic (Non-binary) :-

$$* \quad \begin{array}{l} (\# \text{ of features}) \\ (\text{Types of features}) \\ (\# \text{ of classes / labels}) \end{array}$$

$$\left. \begin{array}{l} 10 \text{ possible values of labels} \\ 7, 1, 6, 12 \text{ features} \\ 5 \text{ number of } n \text{ features} \end{array} \right\} \Rightarrow \binom{10}{10}^{(5)}$$

→ Gradient descent also searches in this space by updating weights.

\* Okhlem's Razor Principle :-



→ If 20% error is acceptable  
 in all the hypothesis above,  
 fit the data-points equally/  
 acceptably then choose the  
 simplest

### Razor

Ockham's principle:-

→ If multiple hypothesis give acceptable error then choose the simplest one. (Wine fitter in our example)



\* Instance Based Classifiers:- (lookup also)

① K-Nearest Neighbor (Non-parametric)

Training:

	$x_1$	$x_2$	$x_3$	label
E <sub>1</sub>	0	0	0	A
E <sub>2</sub>	0	1	0	B
E <sub>3</sub>	0	0	1	B
E <sub>4</sub>	1	1	0	C

Test:

	$x_1$	$x_2$	$x_3$	label
	0	1	0	?

\* → lookup table will search the table for the test input so it will assign B label. But it will fail in case no match is found. So, we go for KNN.

## KNN :-

- 
- ① Training data
- ② The value of K (usually odd, to avoid ties)
- ③ Distance metric (We will use Hamming distance).

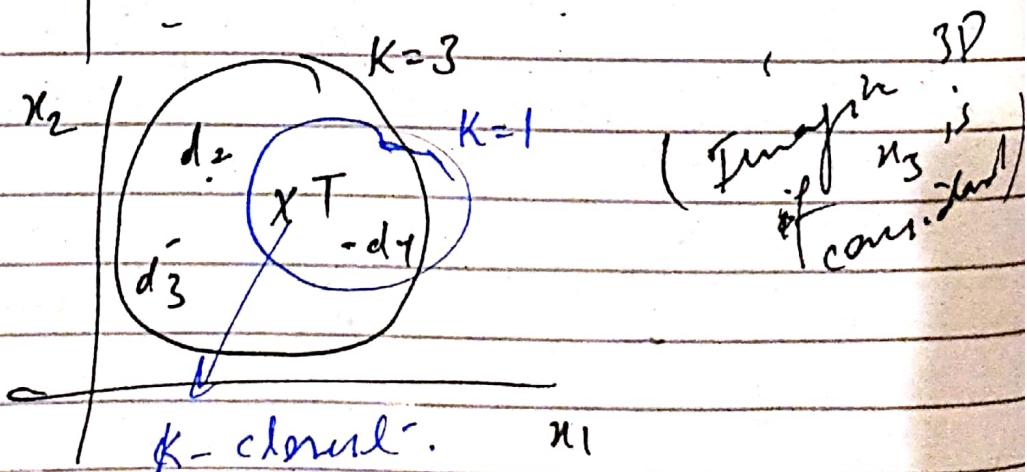
e.g.  $R =$

$$\begin{array}{c} E_i \\ T \end{array} \begin{array}{cccc} & 0 & 0 & 0 \\ & 1 & 1 & 1 \end{array}$$

$$d_1 = 3 = 1 + 1 + 1$$

distance	$n_1$	$n_2$	$n_3$	label
$d_1$	3	:	:	A
$d_2$	2	:	:	B
$d_3$	2	:	:	B
$d_4$	1	:	:	C

Previous table



	$n_1$	$n_2$	$n_3$	label
<del>T</del>	1	1	1	?

For  $K=1$ ; label = C

For  $K=3$ ; label = B

think to group  
with

Techniques to break tie:-

- ① Decrease k.
- ② Take avg distance of the tied points within & then choose the smallest- est avg label.



Full data → Training data

↓ Validation data

→ Slight from best fit point → Best fit point

- ③ Select a class randomly.  
(if no other choice).

\* Distance Metric:-

① Hamming Distance: Binary attributes

② Euclidean Distance: if attribute units are the same -

Ex if ~~attribute~~ Rainy → 0  
Wet → E

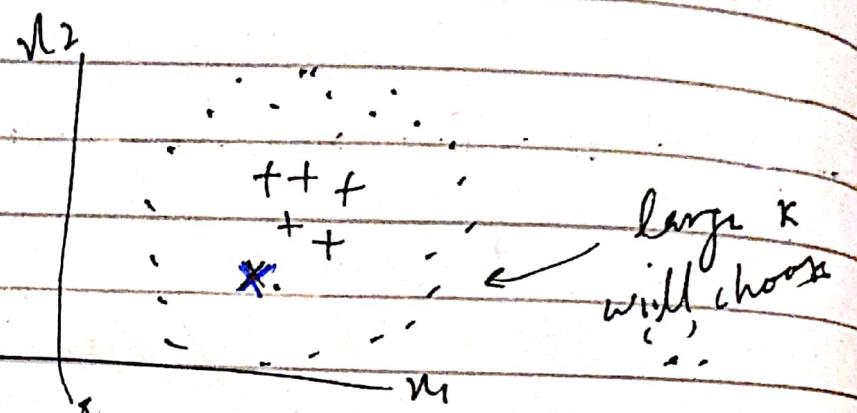
$$d(Q, E) = \sqrt{\sum_{i=1}^{\text{attribute}} (Q_i - E_i)^2}$$

③ Manhattan Distance: attribute unit are real numbers & not the same'

$$d(Q, E) = \sum_{i=1}^{\text{attribute}} |Q_i - E_i|$$

Choosing value of  $K$ :-

- ① If  $K$  is too small, the classifier is sensitive to noise.
- ② If  $K$  is too large the classifier will classify into the majority class.



small  $K$  might also  
choose wrongly due  
to a noise !.

height	distance
$n_1$	$n$
1m	500m
1.5m	800m
1.4m	700m
1.6m	900m

In this case  $n_2$  has a much higher value. & so  $n_2$  will have a higher effect. Hence, we need to normalize it.

## Scaling / Normalization Issues:-

$$P(2, 3) \longrightarrow ?$$

$n_1$	$n_2$	label
-4	3	A $\leftarrow \sqrt{35}$
8	3	B $\leftarrow \sqrt{16}$

- Euclidean Distance from both
- Classified into the label B.

→ ~~also~~ also change dimensions in some cases e.g. convert meters into ~~meters~~ or meters into km.

$n_1$	$n_2$	label
8	3	A $\leftarrow \sqrt{86^2}$
-12	3	B $\leftarrow \sqrt{14^2}$

- Classified into the label A.

(Reasons: So that one attribute doesn't dominate some other  
Ways: attr.)

- ① Division by the maximum.

$$x_i = [x_i - \min(n_i)] / \max(n_i)$$

② Z-Score method - new.

$$z_i = \frac{|x_i - \bar{x}_i|}{\sigma}$$

## Problems with Euclidean Measure (KNN)

① → KNN works well in low dimensional spaces but for higher dimensions it suffers from curse of dimensionality.

② Higher dimensions increase the distances.

③ Higher dimensions make datapoints e.g. outliers close to each other. Thus decreases prediction accuracy.

\* KNN also \* are lazy learners unlike eager learners such as Neural networks.

## Chapter # 03:-

### - Problem Solving by Searching:-

#### \* Searching Agents:-

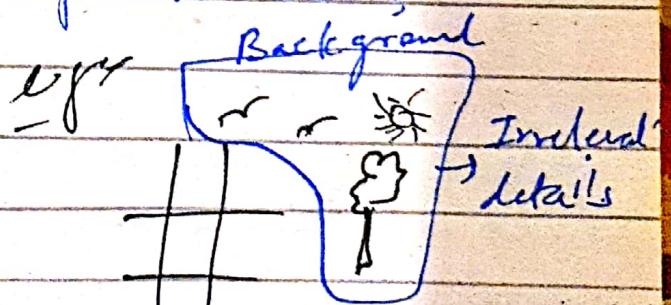
- ① They have a goal to achieve.
- ② The effect of actions is known, i.e., the agent knows how the world evolves.
- ③ The agent chooses a series of actions to achieve its goal.  
e.g.: Tic Tac Toe action sequence.

#### \* Initial Simplifying Assumptions:-

- ① Fully observable.
- ② Deterministic
- ③ Static
- ④ Discrete
- ⑤ Known

#### \* Abstraction:-

Remove irrelevant details to create a 'high level' abstract representation.



#### \* Components of a Search algorithm:-

A search problem should have the following five components:

- ① Start state

## ① Start State:-

First or root state to start a search algorithm.

## ② Actions,-

From any arbitrary state  $s_i$ , all actions from  $s_i$  could be listed.

## ③ Transition Model:-

The next state could be determined by choosing an action.

$$s_j = T(s_i, a_j)$$

## ④ Path Cost:-

- The cost of choosing an action.

## ⑤ Goal State:-

- The goal state is known.

### State Space

State space is the set of all (reachable) connected states. So, the state space can be smaller than the possible number of states (unreachable).

\* To check examples in slides.

Vacuum world state space

graph :-

S-I generation

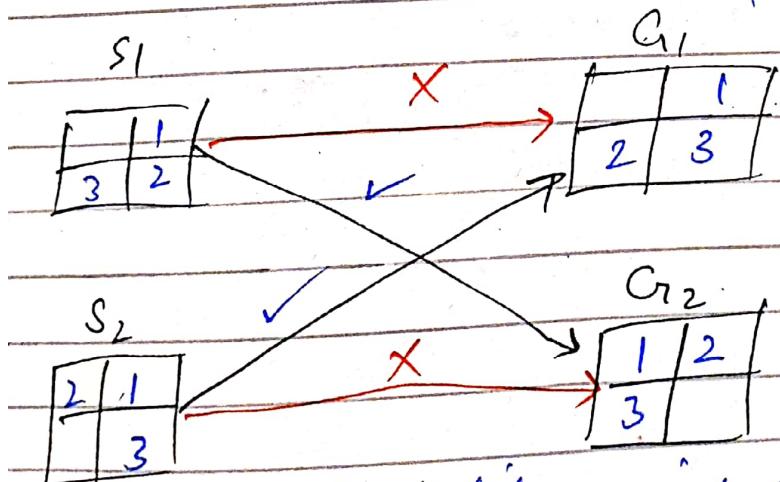
$$n \times 2^n \Rightarrow 2 \times 2^2 = 8$$

→ State changes based on the  
diri-states in vacuum loading

8-puzzle :-  
state space analysis

q! / 2

only even permutations  
are solvable  
given  $m = n = 2$   
 $m \times n$  board

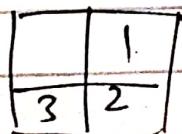


Small tile coming after  
inversions = tile 1 + tile 2 + tile 3

$$\Rightarrow 0 + 0 + 1$$

(Only can be)

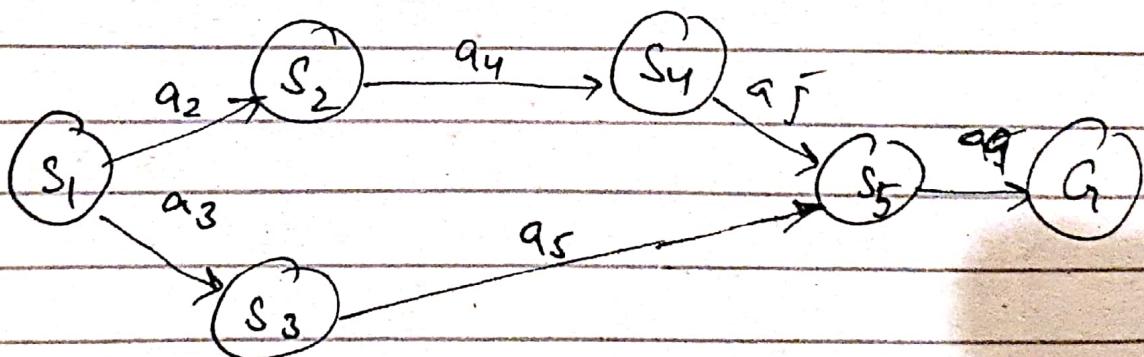
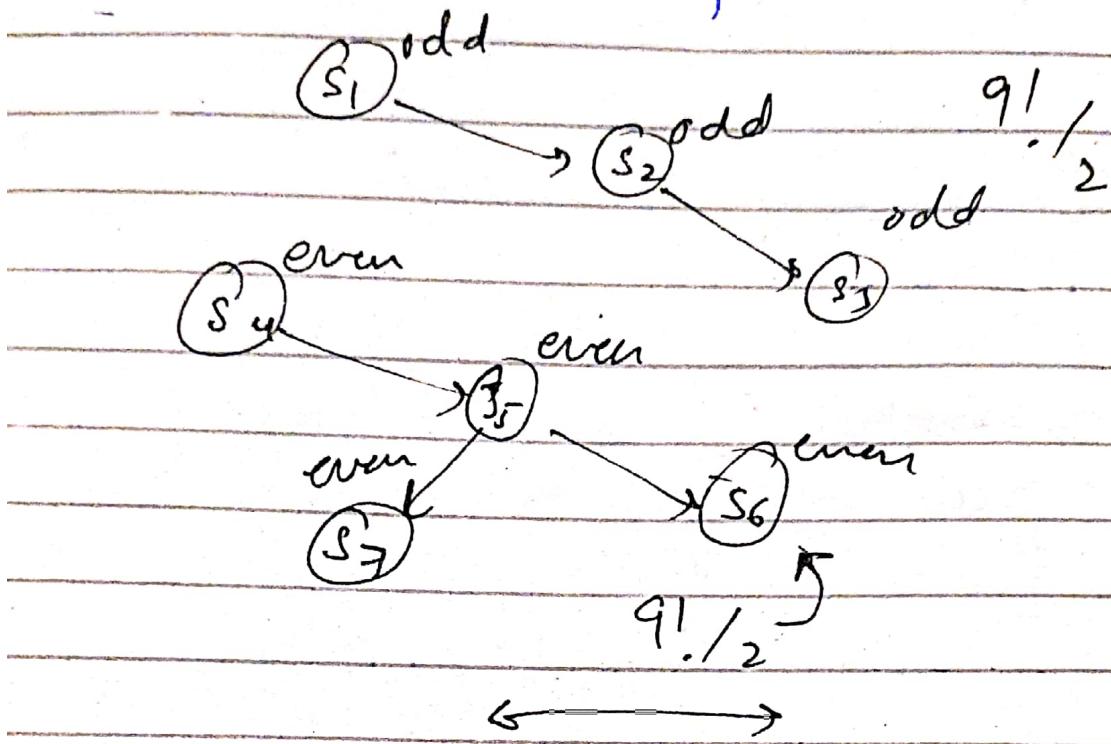
reached by odd parity  
start state).



# Inversion in 8-puzzles

Inversion is

the count of tiles occurring after <sup>ith</sup> tile <sup>a</sup> in 8-puzzle.



8-puzzle  $\rightarrow$  Millions

16-puzzle  $\rightarrow$  Hundreds of Millions

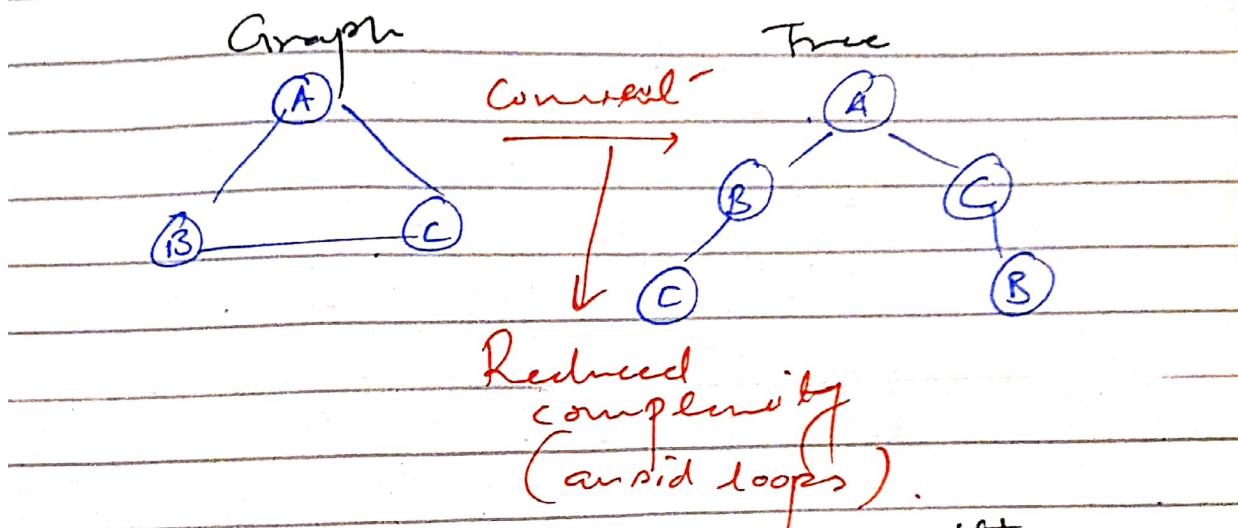
Chess  $\rightarrow 10^{120}$

Vacuum Cleaner  $\rightarrow n \times 2^n$

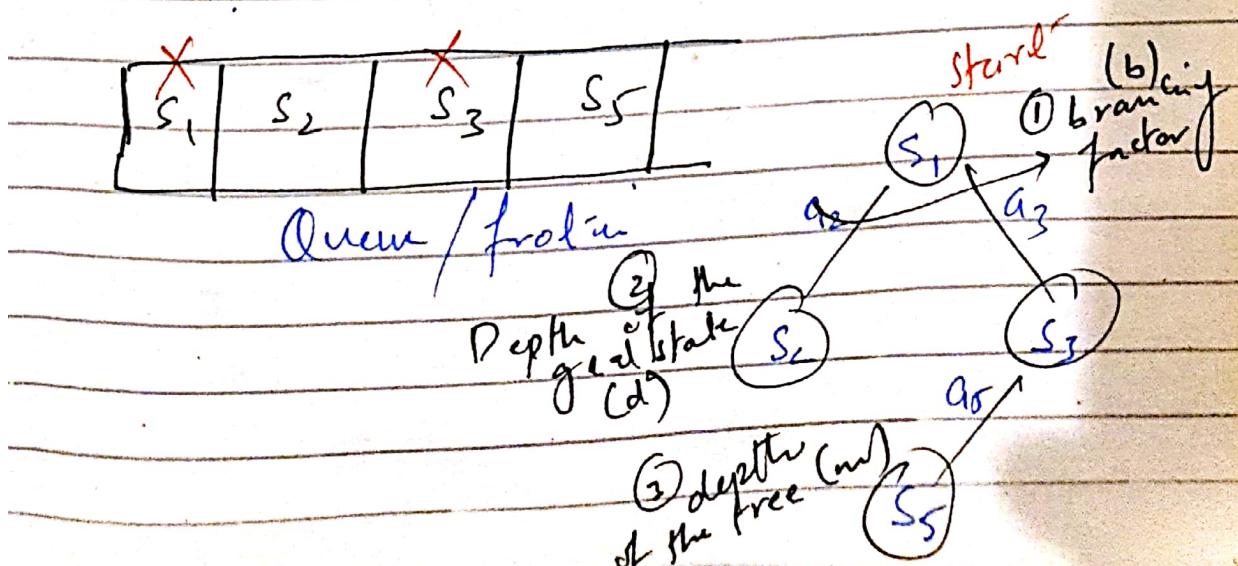
\* Not possible using conventional methods (Dijkstra).

## Practical Solution -

→ Start with one state & keep expanding by maintaining a queue. Goal is to ~~get~~ as few of states as possible expand.



Tree - Search Algorithm -  
 fn (problem, strategy) → Different b/w algorithms.  
 → returns a solution or failing  
 initialize the search tree using  
 the initial state of problem.  
 :  
 : → In slides.



Differentiability b/w algorithms on the basis of :-

① Time Complexity:-

= The amount of time an algorithm takes to find a goal state usually represented in terms of number of nodes.

$$O(b^d); O(b^m)$$

② Space Complexity:-

= The maximum # of nodes an algorithm stores in memory to find a solution.

③ Completeness:-

if Algorithm<sup>2</sup> should return a solution if one exists.

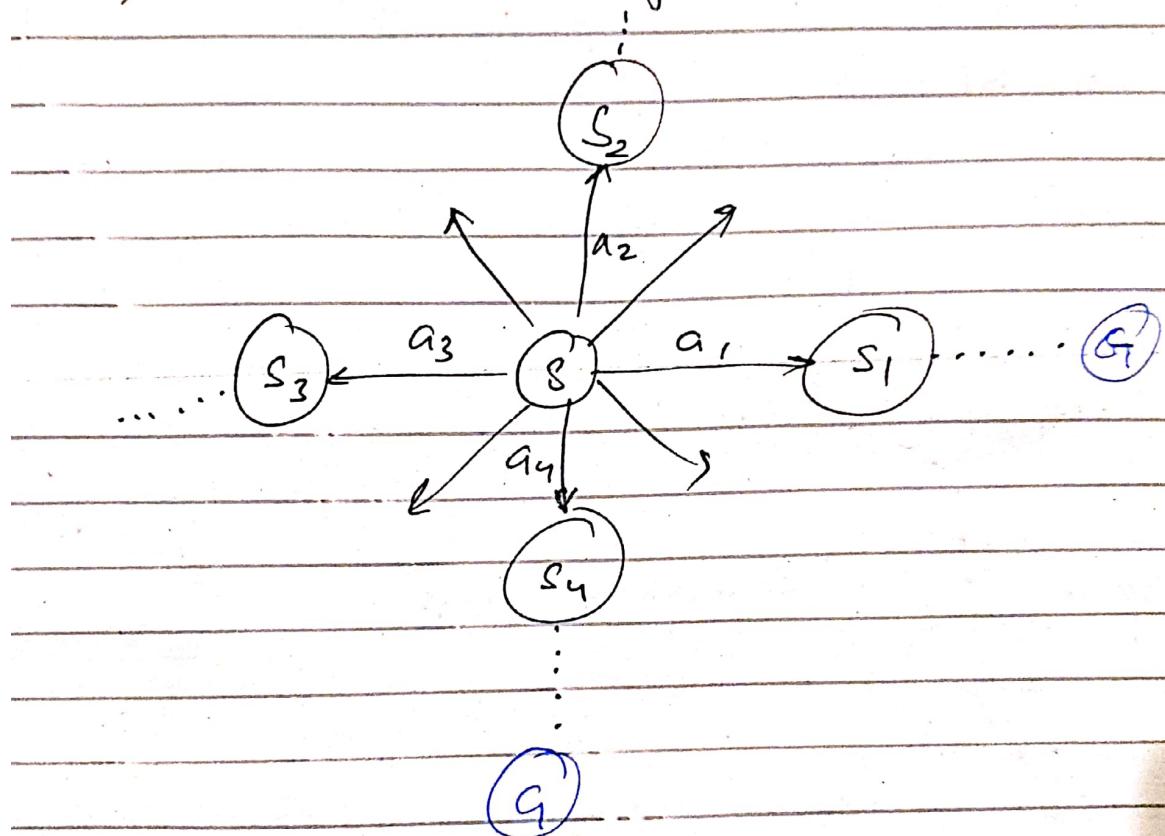
④ Optimality-

= Should return least-cost-solution (in case of multiple solutions).

## \* Uninformed Search Strategies:-

In which the algorithm only has search components & no other information about the state space.

→ 5 search components.



→ Such algorithms expand in all directions, Hence expensive -

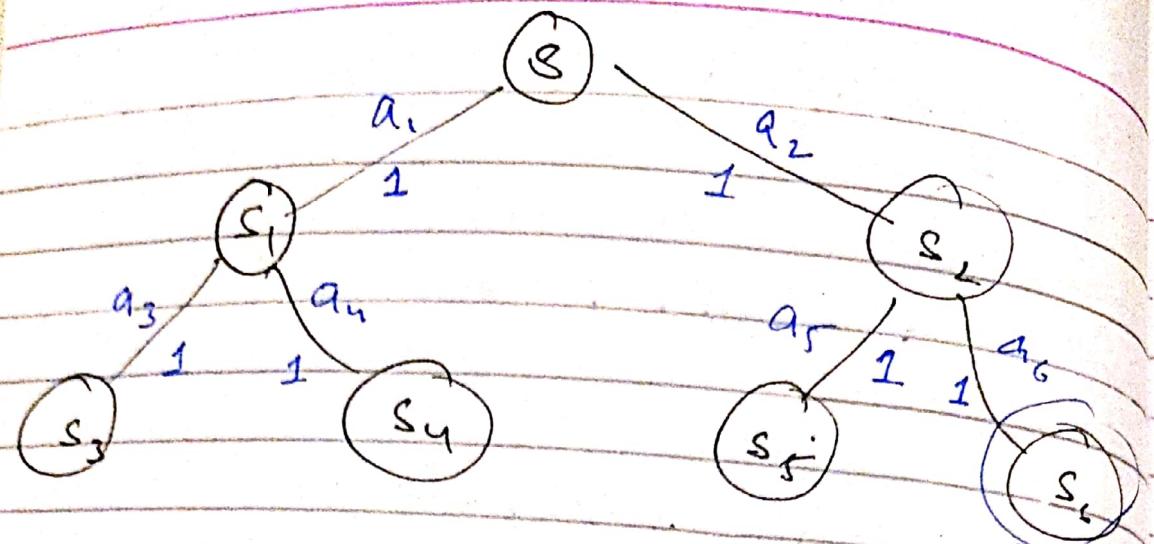
→ Also, called blind searches.

Algorithms include:-

BFS → ① Breadth First Search

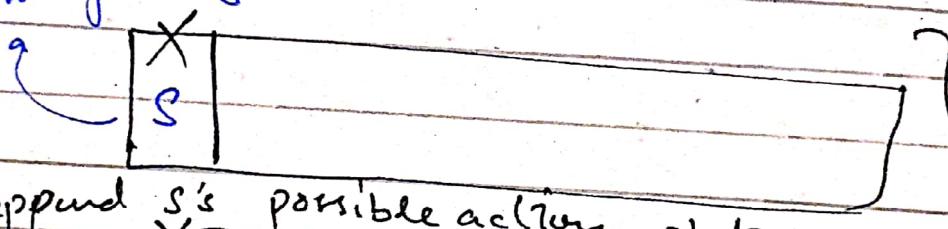
DFS → ② Depth =  $\infty$   
③ Uniform cost - search  
④ Iterative deepening search

⑤ Bi-directional search

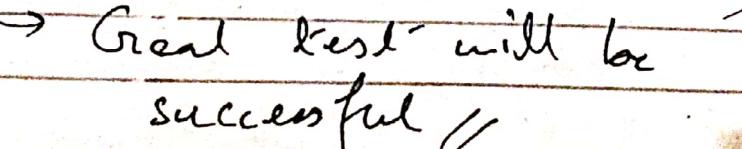
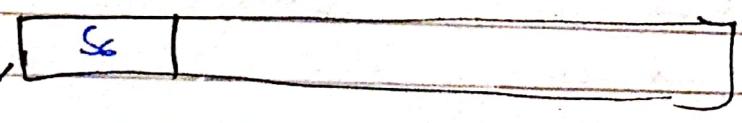
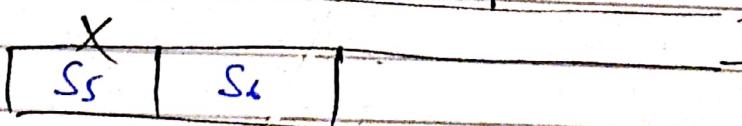
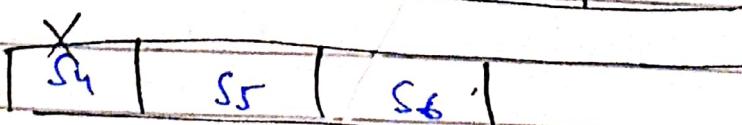
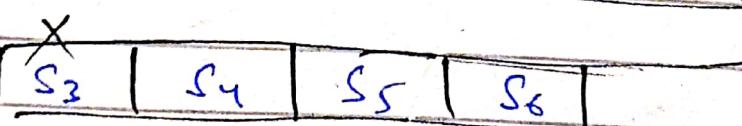
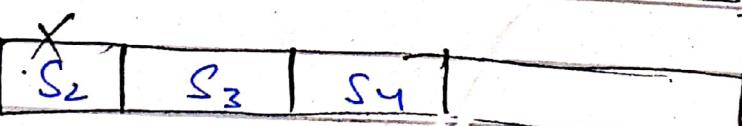


for - BFS :-

This is the worst-case scenario since - BFS uses a queue (FIFO)  
not goal (goal test) dependent



append S's possible actions states



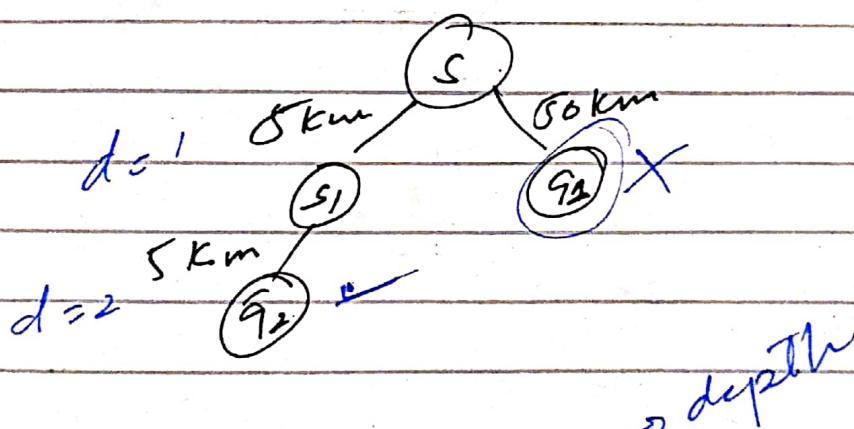
Apply goal test & return expansion degree

Goal test will be successful //

→ Using this question, you can generalize the BFS - spanning tree (search space tree)

① Completeness → yes, if goal state is at finite steps / depth.

② Optimal → yes, if path cost is a non-decreasing fn of depth.



③ Time Complexity =  $\frac{1^0 + 2^1 + 3^2}{1+2+4}$

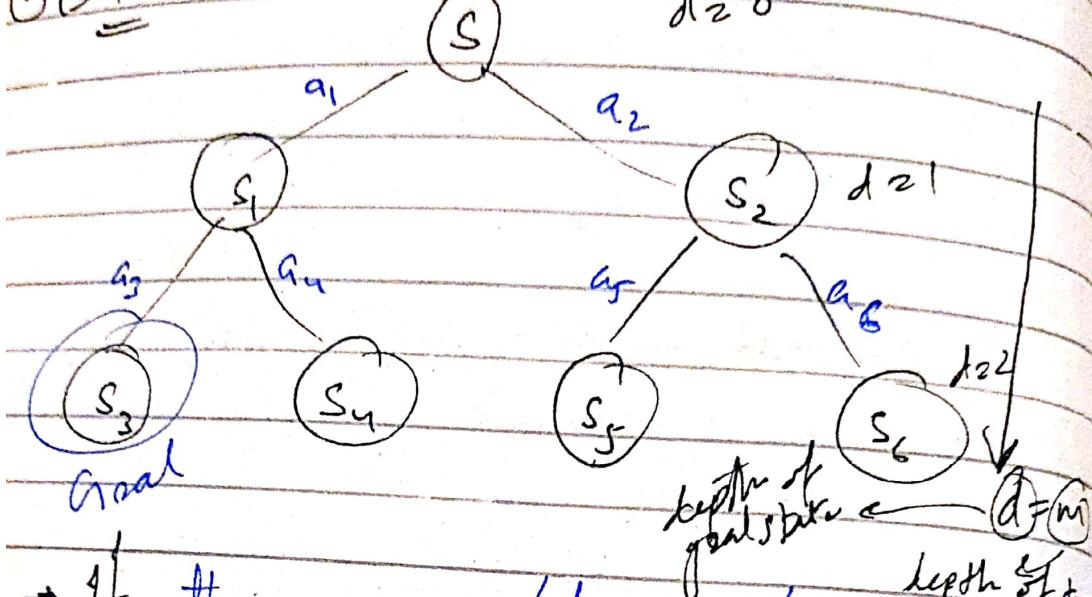
$$\text{branching factor} \leftarrow \frac{b^0 + b^1 + b^2}{b^0 + b^1 + \dots + b^d} = O(b^d)$$

~~In terms of # of nodes.~~

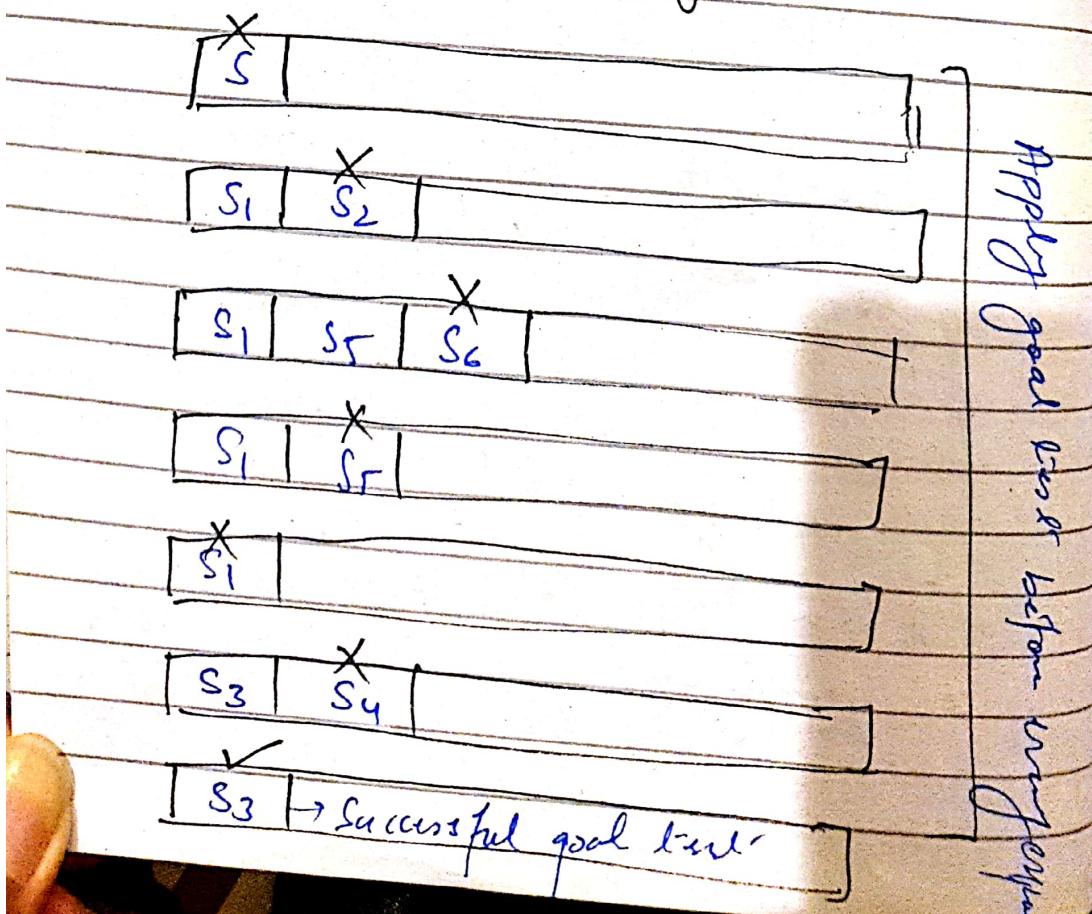
④ Space Complexity =  $O(b^{d+1})$

When to use BFS? when the solution/goal is shallow i.e d is small.

② DFS: (worst case)

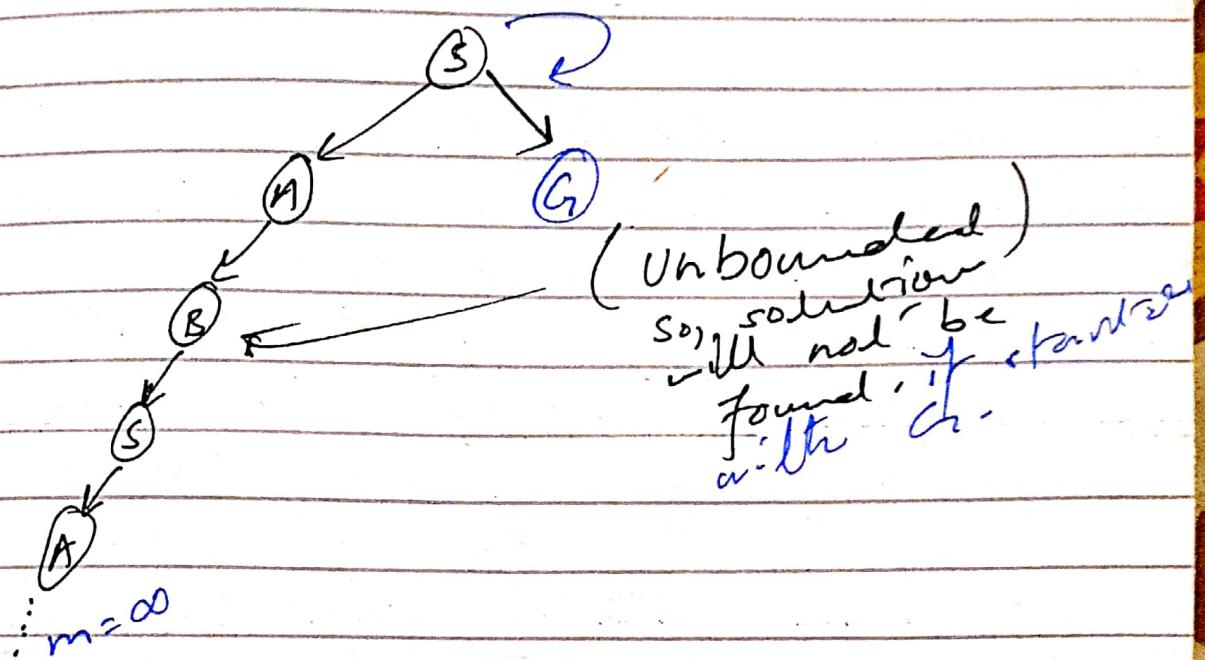


→ If the queue/frontier/fringe is switched from FIFO to LIFO (stack) the algorithm will become DFS.

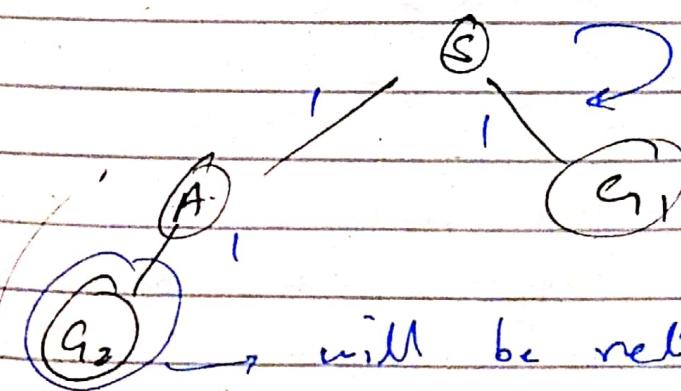


→ BFS explores shallowest node first  
 vs. DFS explores deepest node //.

① Completeness:- For trees with unbounded depth, completeness is false.  $m = \infty$



② Optimality is false, if suboptimal goal node in the first deepest path.



### ③ Time Complexity:-

$$\# \text{ of nodes} = 1 + 2 + 4 + 8 + \dots + 32$$

$\overset{b=2}{\underset{m}{=}}$   $m_1 m_2 m_3 \dots m_{\max}$   
 $= b^0 + b^1 + b^2 + b^3 + \dots + b^m$   
 $= O(b^m)$  //  $\leftarrow$

### ④ Space Complexity

$$\# \text{ of nodes} = 1 + 2 + 2 + 2 + 2 + \dots + 2$$

$\overset{d=2}{\underset{m}{=}}$   $d_1 d_2 d_3 \dots d_m$   
 $= b^0 + b^1 + b^2 + b^3 + b^4 + \dots + b^m$   
 $= 1 + m \times b$   
 $= O(b^m)$  //  $\leftarrow$

$b=10$

$m=12$

$d=5$

Node = 1 Kbytes

$d=1$

Space Complexity

120 Kb

Time Complexity

$10^{12}$  s

$d=2$

120 Kb

$10^{12}$

$d=3$

120 Kb

$10^{12}$

$d=4$

120 Kb

$10^{12}$

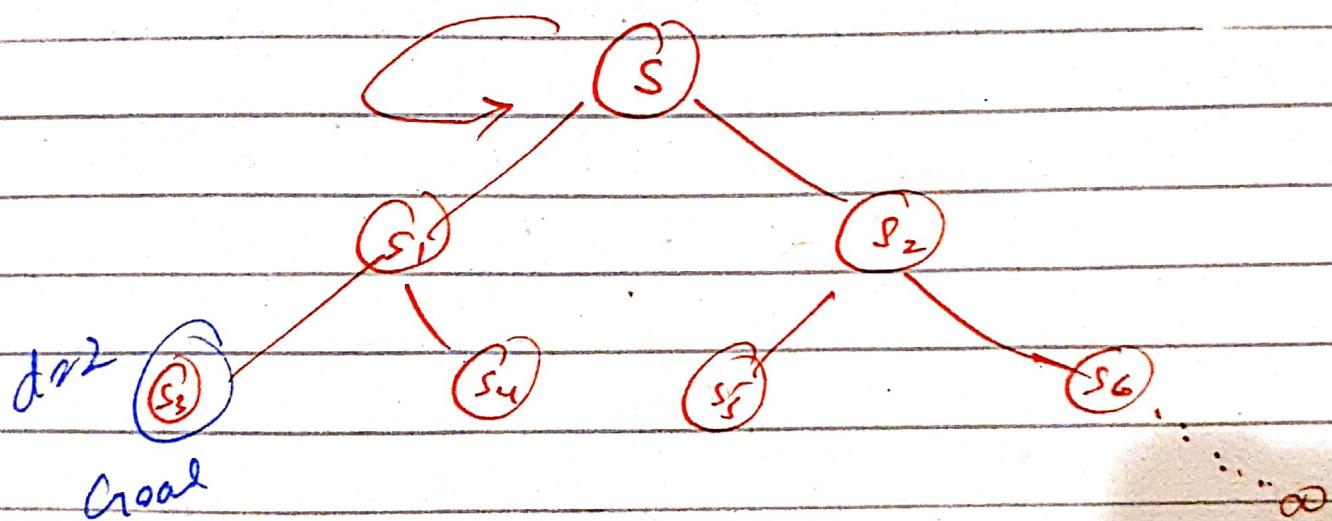
$d=5$

120 Kb

$10^{12}$

## → Iterative Deepening Search

- ① Completeness Guarantees optimality for problems with non-decreasing path cost.
- ② Have the space complexity of DFS i.e., linear.
- ③ slightly higher time complexity.



Solution = null;

D = 0;

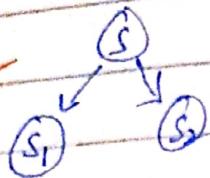
while (solution == null) {  
 solution = DFS(problem, D)

D = D + 1

}

For  $D=0$ :-  $S \rightarrow S$

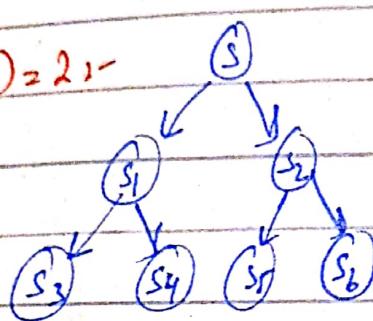
For  $D=1$ :-



$S$	$S_1$	$S_2$
$S_1$	$S_2$	
$S_1$		

$k$

For  $D=2$ :-



$S$	$S_1$	$S_2$
$S_1$	$S_2$	
$S_1$	$S_5$	$S_6$
$S_1$	$S_5$	
$S_1$	$S_3$	$S_4$
$S_3$		

Goal found,

① Completeness:- True, if solution exists at finite steps.

② Optimality:- True, if non-decreasing path cost.

③ Space Complexity

$$\# \text{ of nodes} = 1 + b + b + \dots + b^d$$

$$= 1 + bd$$

$$= O(bd)$$

## (M) Time Complexity:-

$$\text{Time} \rightarrow b^0 + b^1 + b^2 + \dots + b^{d-1} + b^d$$
$$db^0 + (d-1)b^1 + (d-2)b^2 + \dots + 2(b^{d-1}) + b^d$$
$$\Rightarrow db^0 + (d-1)b^1 + (d-2)b^2 + \dots + 2b^{d-1} + b^d$$

Exact  
# of  
nodes  
Complexity  $\rightarrow O(b^d)$

## Comparison of Time Complexity:-

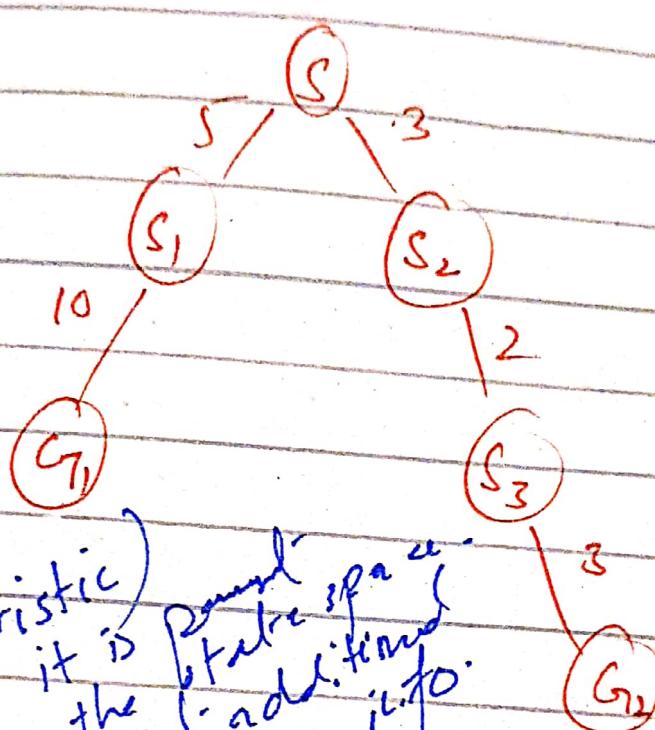
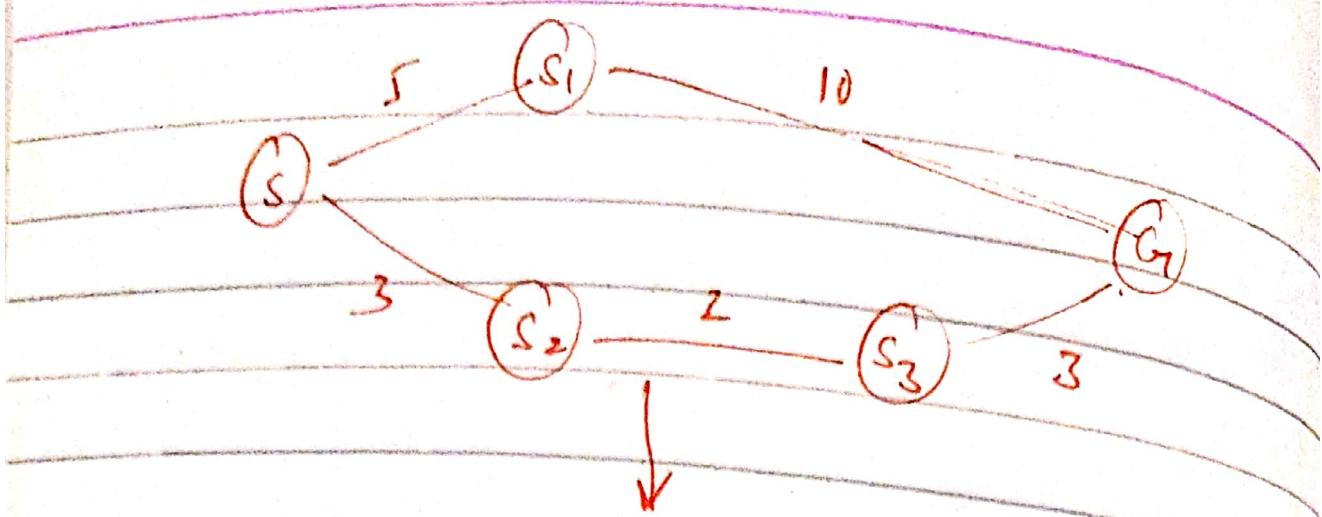
For  $b=10, d=5$

BFS = 111,110 nodes

IDFS = 123,450 nodes (Not big of a difference)

## (u) Uniform Cost Search:-

- ① Suitable for problems with varying costs. (randomly distributed cost).
- ② Gives an optimal solution in such state spaces.
- ③ The queue is prioritize at  $g(n)$  cost, where  $g(n)$  is the cost of node  $n$  from the root node (initial state).



(Not heuristic)  
as it is prime space  
in the total addition  
in the path add it to.

$g(n)$  = Cost of node  $n$  from root.

$$g(S) = 0$$

$$g(S_1) = 5$$

$$g(S_2) = 3$$

$$g(S_3) = 5$$

$$g(G_1) = 15 //$$

$$g(G_2) = 8 //$$

Priority  $g(u)$  :-

	$\Gamma(s)$	
$g(u)$	0	
	$\Gamma(s_1 \mid s_2)$	expand $\min g(u)$
$g(u)$	5 3	cost.

	$\Gamma(s_1 \mid s_3)$	same $g(u)$ , just choose one randomly (no constraint)
$g(u)$	5 5	

→ However, usually, we expand newer ones because we started from this path (optimal) so keep exploring it.

Let's explore  $s_1$  (contradictory).-

	$\Gamma(s_3 \mid G_1)$	
$g(u)$	5 15	

→ even though we have the goal in the queue, we still expand the  $\min g(u)$  cost.

	$\Gamma(G_1 \mid (G_2))$	
$g(u)$	15 8	

→ Before every expansion, apply the goal test. So optimal soln.

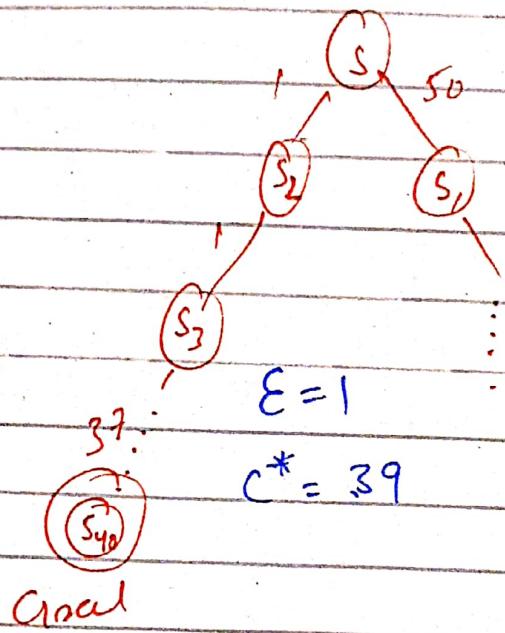
$$G_2 \rightarrow S_3 \rightarrow S_2 \rightarrow S$$

① Completeness:- Yes, if goal is at finite steps.

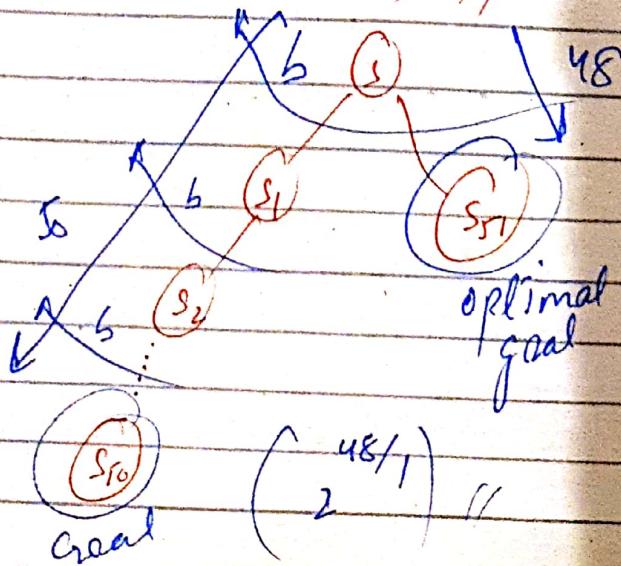
② Optimality:- Yes

③ Time Complexity:-  $C^*$  is the optimal goal cost.

$\epsilon$  is the minimum step cost.



$$O(b^{C^*/\epsilon}) //$$



create

$$(2^{48/1}) //$$

## ⑤ Bi-directional :-

Actions are reversible in goal is a concrete state, not abstract.

→ Read yourself //



CH#3 //

## Informed Searches:-

\* Heuristic ftns

\* Best first Searches

i) Uniform Cost Search (Uninformed)

ii) Greedy Best first Search

iii) A\* Search

iv) Memory Bounded A\*

v) Recursive Best First Search.

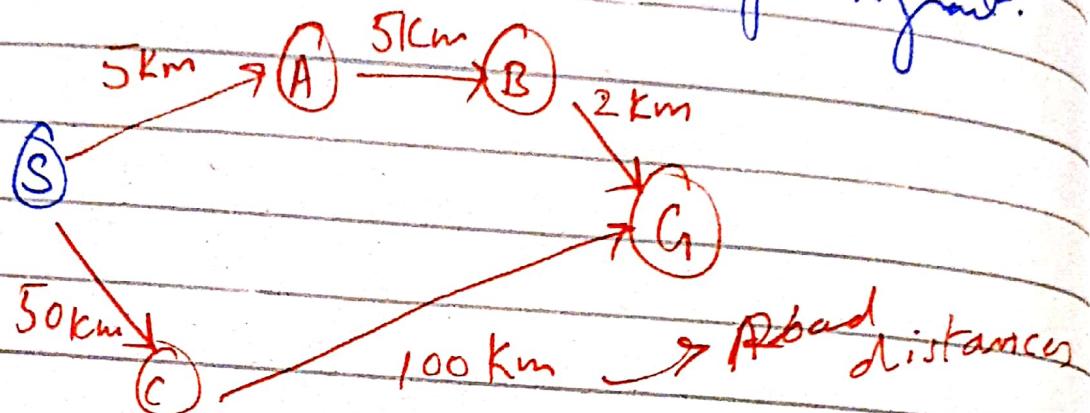
→ Uninformed searches expand in all directions. If goal takes 30 steps with a branching factor of 4. So,  $4^{30}$  nodes will be searched, which is expensive.

→ Additional information in informed searches will allow you to expand more in the goal's direction.

## \* Heuristic:-

- ① Rule of thumb to find a solution.
- ② Any additional information that speeds up the algorithm to find a solution.

Example:- Path finding Agent.



$h_1(n) \approx$  Straight line distance to goal state.

$$h_1(s) \approx 6 \text{ km}$$

$$h_1(A) = 4 \text{ km}$$

$$h_1(B) = 1 \text{ km}$$

$$h_1(C) = 60 \text{ km}$$

$$h_1(G) = 0 \text{ km}$$

$$| s \cdot L \cdot D < R \cdot D$$

$$\forall n, h(n) \leq h^*(n)$$

Admissible

Heuristic,,

→ This is one heuristic, there can be multiple heuristics.

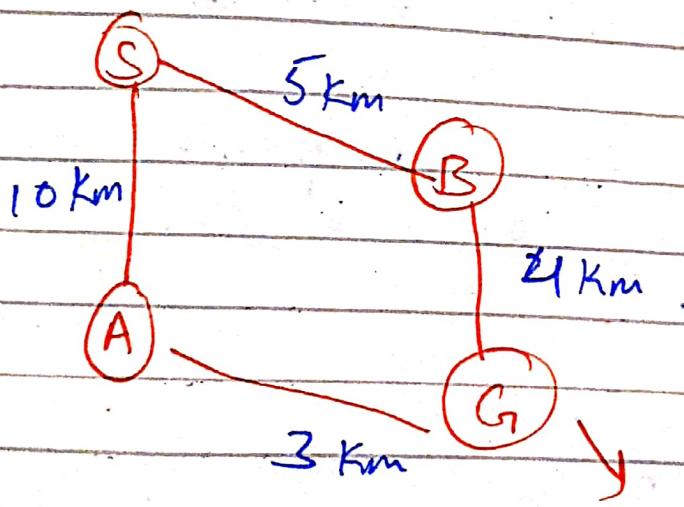
Another heuristic,

$h_2(n)$  = Railroad distances to goal state.

$h_3(n)$  = Aeroplane " "

$h_4(n)$  = Distance by foot.

① Greedy Best First Search



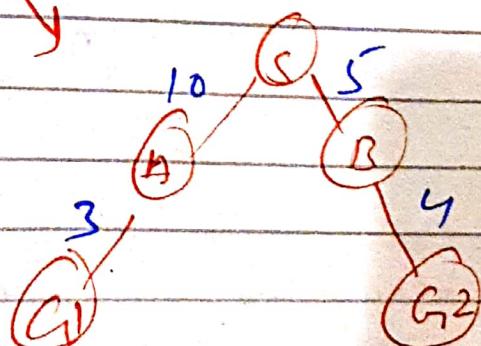
$h(n)$  = straight line distance.

$$h(S) = 7 \text{ km}$$

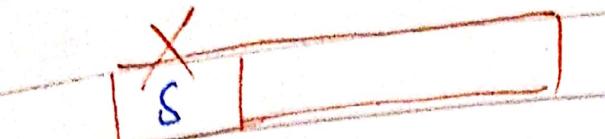
$$h(A) = 2 \text{ km}$$

$$h(B) = 3 \text{ km}$$

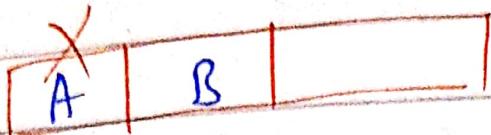
$$h(G) = 0 \text{ km}$$



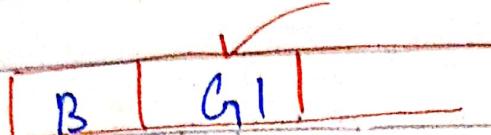
→ Queue, priority on  $h(n)$ .



$h(n) \neq$



$h(n) = 2 \quad 3$

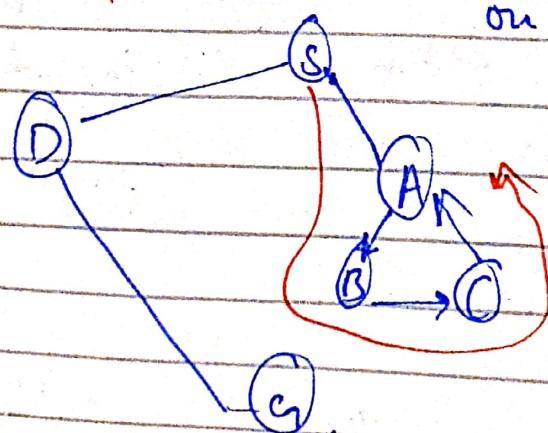


$h(n) = 3 \quad 0$

→ Goal is found but  $C_1$  is not the optimal solution.

(1) Optimality:- Not optimal.

(2) Completeness:- Not complete based on heuristic.



→ However, if loops are entered for using trees then complete

Doubt check this

③ Time Complexity:-  $b^0 + b^1 + b^2 + \dots + b^m$

$O(b^m)$

④ Space Complexity:-

$O(b^m)$

Designing a heuristic function:-

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

S<sub>1</sub>                    S<sub>2</sub>                    G

→ Assuming both are reachable. Which seems closer? -

→ They seem closer, more # of tiles off the actual location in S<sub>1</sub>.

$h_1(n) = \# \text{ of misplaced tiles}$

→ Lower heuristic is better.

① We relax the rules/constraints of the game/environment.

e.g.: If move empty block in the misplaced tiles directly, then how many steps. i.e. 3.

$h_2(n) \rightarrow$

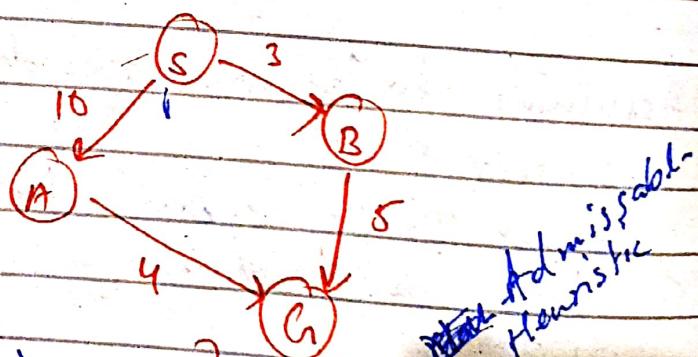
$h_2(n)$  = manhattan distance of misplaced tiles from actual position in goal state.

$$\begin{array}{l|l} h_1(s_1) = 3 & h_2(s_1) = 4 \\ h_1(s_2) = 7 & h_2(s_2) = 18 \end{array} \quad \begin{array}{l} 2+1+1+0+0+0... \\ 3+1+4+2+0+2+1 \\ +3 \end{array}$$

## (ii) A\* Search:-

→ Queue prioritizes on  $f(n)$ , where

$$f(n) = g(n) + h(n)$$



$$g(S) = 0 \quad h(S) = 7$$

$$g(A) = 10 \quad h(A) = 3$$

$$g(B) = 3 \quad h(B) = 4$$

$$g(C) = 8 \quad h(C) = 0$$

$h(n)$  = straight line distance

$f(n)$  = Distance of  $n$  from root 'S'.

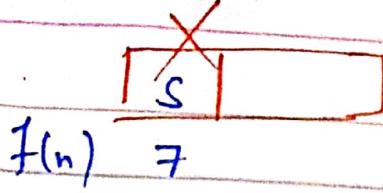
✓ or convert  
into breath,  
so get both,

$$f(S) = 7$$

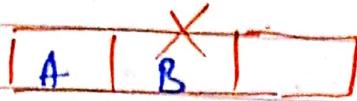
$$f(A) = 13$$

$$f(B) = 7$$

$$f(C) = 8$$



$f(n) = 7$



$f(n) = 13 \quad 7$



$f(n) = 13 \quad 8$

Note:-

In case of ties, choose random or follow the same path.

Optimal path:-  $S \rightarrow B \rightarrow G_1$

→ Now, if you choose a non-admissible heuristic, it will not give optimal path.

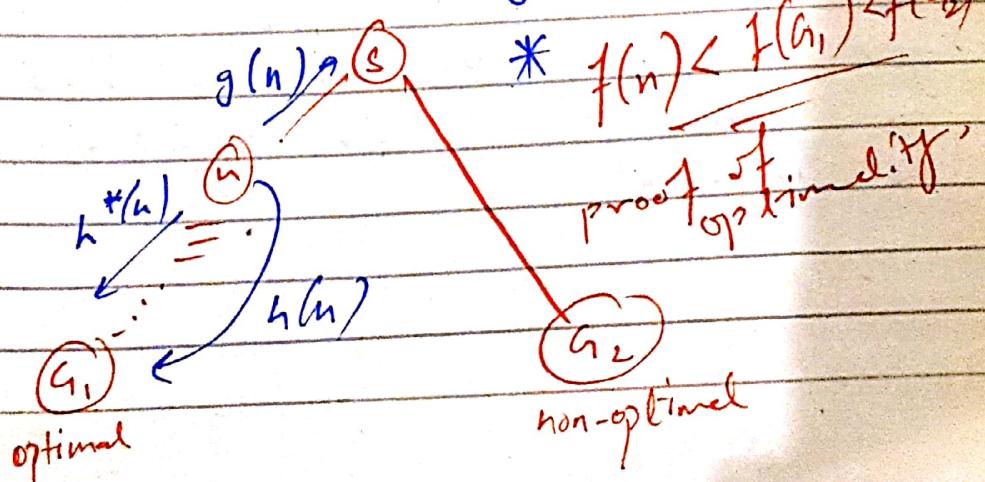
① Optimality:-  $A^*$  is optimal, provided

For trees - i)  $h(n)$  is admissible i.e  $h(n)$  never overestimates the actual cost.

$$\forall n; h(n) \leq h^*(n)$$

heuristic      actual cost

→ Admissibility is for trees.



①  $f(n) < f(a_1) < f(a_2) \therefore$   
proof of optimality :-

$$f(n) = g(n) + h(n) \rightarrow (a)$$

$$f(a_1) = g(a_1) + h(a_1) //$$

$$f(n) \leq g(n) + h^*(n)$$

since

$$f(n) \leq g(a_1) \quad h(n) \leq h^*(n)$$

$$f(n) \leq f(a_1) // \text{if nodes are the same (ideal case)}$$

$$f(a_1) = g(a_1) + 0 \quad (\because \text{condition})$$

$$f(a_2) = g(a_2) + 0 \quad (\text{if optimality})$$

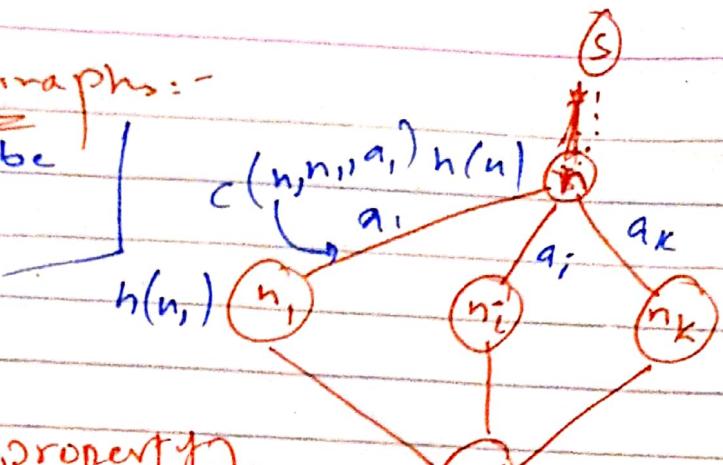
$$g(a_1) < g(a_2)$$

Proved  $f(n) < f(a_1) < f(a_2)$

→ If the above statement holds for trees then, the solution will be optimal;

For Cographs:-

$h(n)$  should be  
consistent.



consistency property :-

$$\rightarrow h(n) \leq c(n, n_i, a_i) + h(n_i)$$

II  $\Rightarrow f(n) < f(n_i)$  :-

proof of consistency

$$f(n) = g(n) + h(n)$$

$$f(n) < g(n) + c(n, n_i, a_i) + h(n_i) \quad \text{(By consistency property)}$$

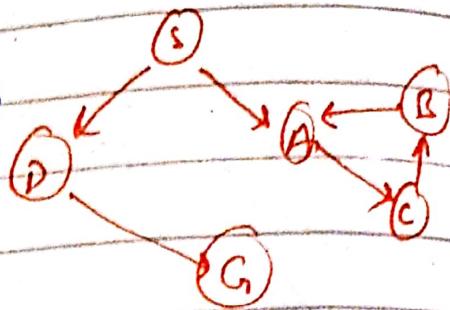
$$f(n) < g(n_i) + h(n_i)$$

$f(n) < f(n_i)$  // proved,

$\rightarrow$  If the above statement holds for graphs for every node, then the solution will be optimal.

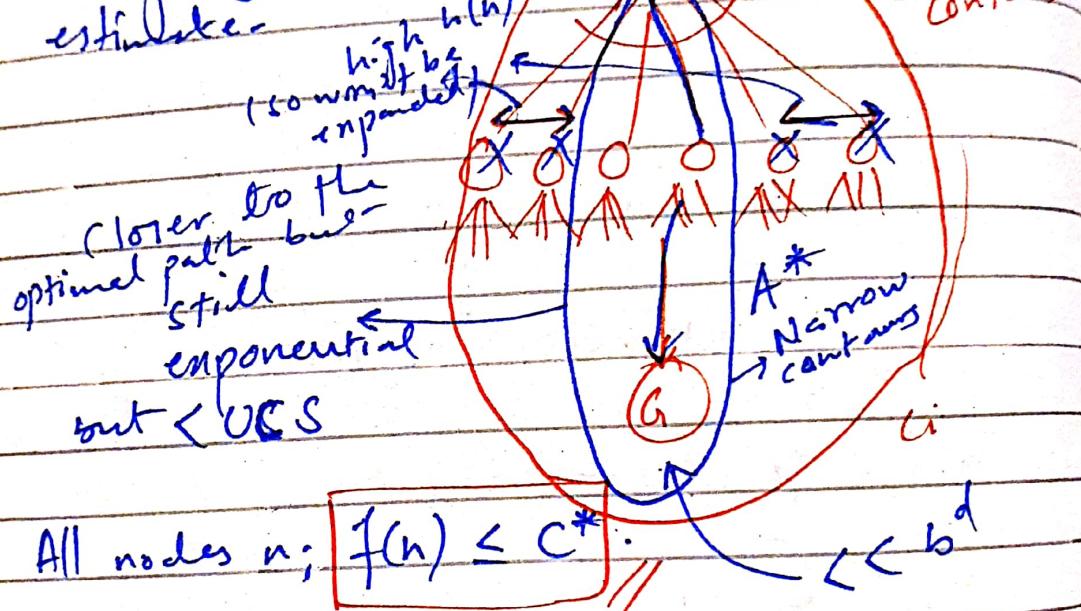
$\rightarrow$  Given these 2 conditions, the solution will be optimal.

② Completeness:- Yes, even in such situations



③ Time Complexity:- ~~Cannal~~ specifies

categorize, depends on estimator



④ Space Complexity = much less than UCS.

→ man of all admissible heuristics are the best. (maybe)

→ if  $h(n)=0$ , then A\* will become UCS.

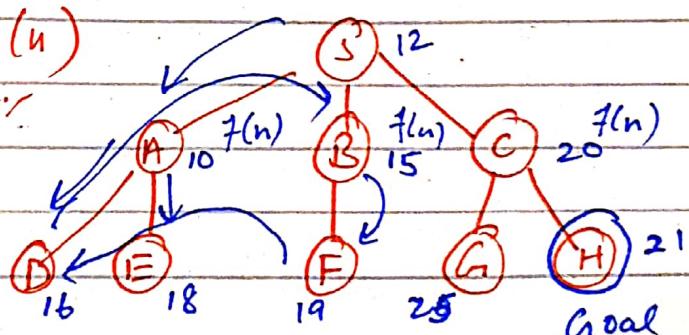
### (iii) Recursive Best First Search,-

- i) A\* search runs out of memory before running out of time.
- ii) RBFS is a memory efficient implementation of A\* search.
- iii) In RBFS Queue is priority on  $f(n) = g(n) + h(n)$
- iv) RBFS searches like DFS i.e. only one depth at a time in a recursive fashion, using a variable called f-value.

$$f(n) = g(n) + h(n)$$

Space Complexity

$$\Theta(bd)$$



$$f(n)$$

$$12 \times$$

S

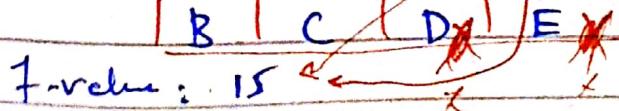
f-value:  $\infty$

$$f(n) := 16 \quad 15 \quad 20 \quad \text{max depth nodes}$$



f-value: 15

$$15 \quad 20 \quad 16 \quad 18$$

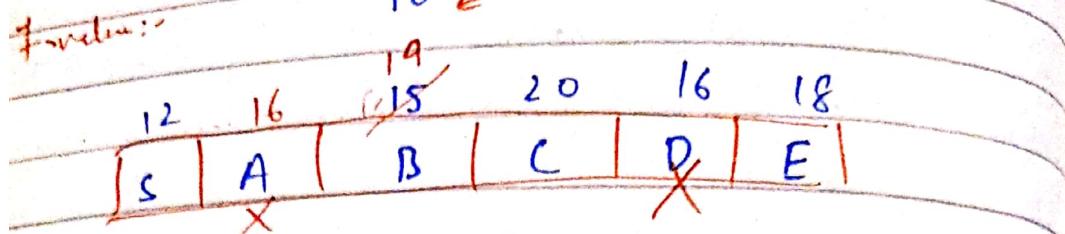
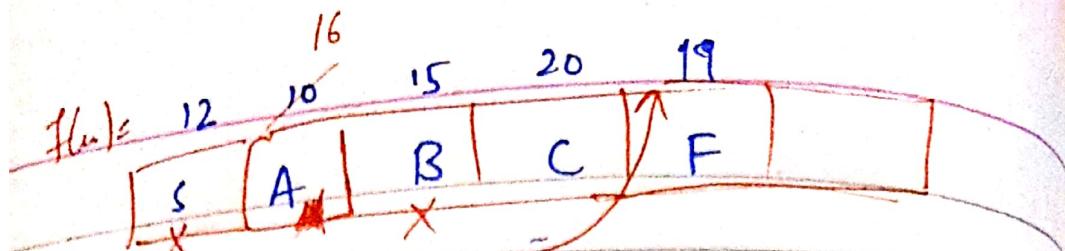


f-value: 15

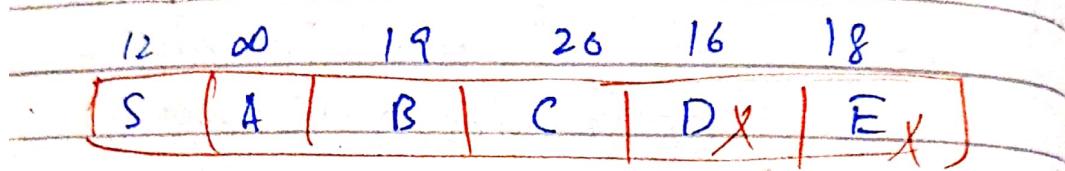
\*

f-value is the alternative path cost.

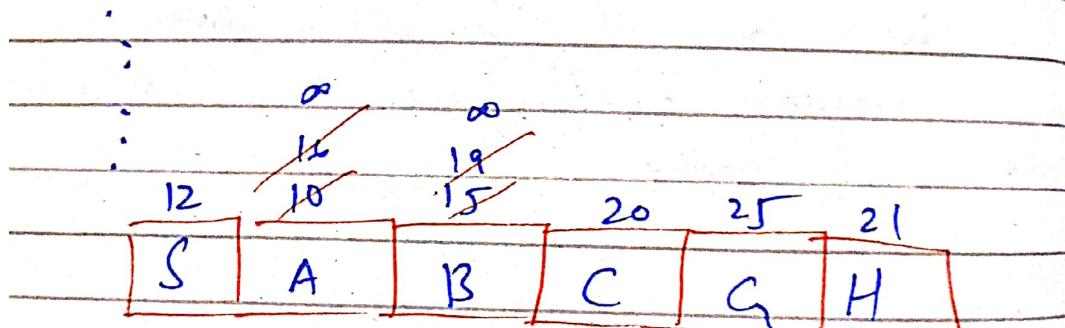
→ The max f(n) exceeds f-value, algo backtracks.



f-value: 19      18



f-value 19



Final values: 19      00      25

→ See Rummie example for better understanding.

Optimality: Yes given the 2 conditions of A\*.

Time: Difficult to characterize  
Depends on  $h(u)$ .

## Notion of dominance:-

If  $h_2(n) \geq h_1(n)$  for all  $n$  (both admissible) then  $h_2$  dominates  $h_1$ . aka.  $h_2$  is better for search.

## Effective branching factor, ( $b^*$ )

$$N+1 = 1 + b^* + (b^*)^2 + \dots + (b^*)^d$$

→ Measure provides a good guide to the heuristic's overall usefulness.

## Chapter # 04:-

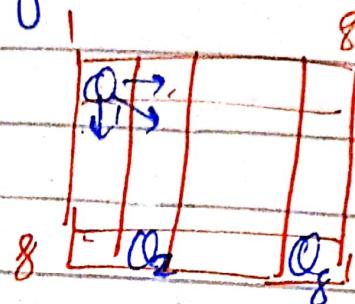
## Local Search & Optimization

- ① Problem Formulation.
- ② Hill climbing & variants.
- ③ Simulated Annealing Algorithm
- ④ Genetic Algorithms

# ① Problem Formulation & Related Co.

i) In local searches the path to goal state is not important & we are interested in goal state.

8-Queen problem

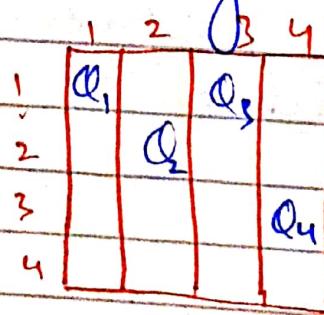


(12 goals)

ii) A state in local searches has a fittness value (maximization problems) or objective value (minimization problems).

~~fittness expressed function~~ → heuristic function.

iii) In local searches we search in the local neighborhood & ignore the previous state after taking an action.



fitness  $f(u)$  - (for any state)  
 $f(u) = \#$  of non-attacking queen pairs

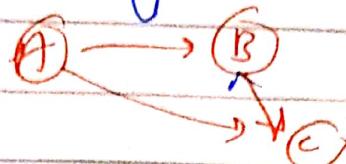
s

→ state aka solution / individual / ...

Mater-

\* for travelling sales man, objective fn can be.

$f(n) = O(n)$  = Distance by visiting each city at most once.



A B C  
A B C B A

Back to Traveling Salesman problem

$$f(a) = 4C_2 \Rightarrow \\ (\text{Goal}) = 4! \\ \frac{n!}{(n-r)! \times r!} = 6$$

is non-existing pairs

Q <sub>1</sub>	Q <sub>2</sub>
Q <sub>1</sub>	Q <sub>3</sub>
Q <sub>1</sub>	Q <sub>4</sub>
Q <sub>2</sub>	Q <sub>3</sub>
Q <sub>2</sub>	Q <sub>4</sub>
Q <sub>3</sub>	Q <sub>4</sub>

$$f(s) = 3$$

(Our state)

Q <sub>1</sub>	Q <sub>2</sub>	✓
Q <sub>1</sub>	Q <sub>3</sub>	✓
Q <sub>1</sub>	Q <sub>4</sub>	X → 3
Q <sub>2</sub>	Q <sub>3</sub>	✓
Q <sub>2</sub>	Q <sub>4</sub>	X ↗
Q <sub>3</sub>	Q <sub>4</sub>	X

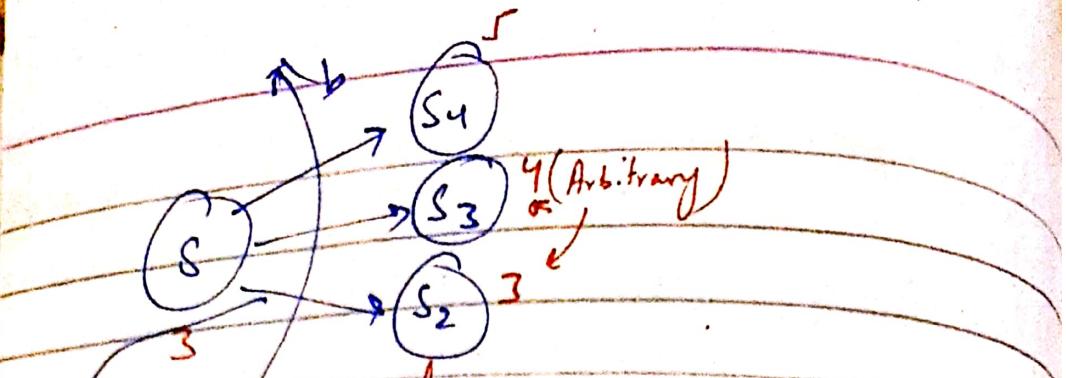
action:- move queen Q<sub>1</sub> to any location (causes state change)

if:-

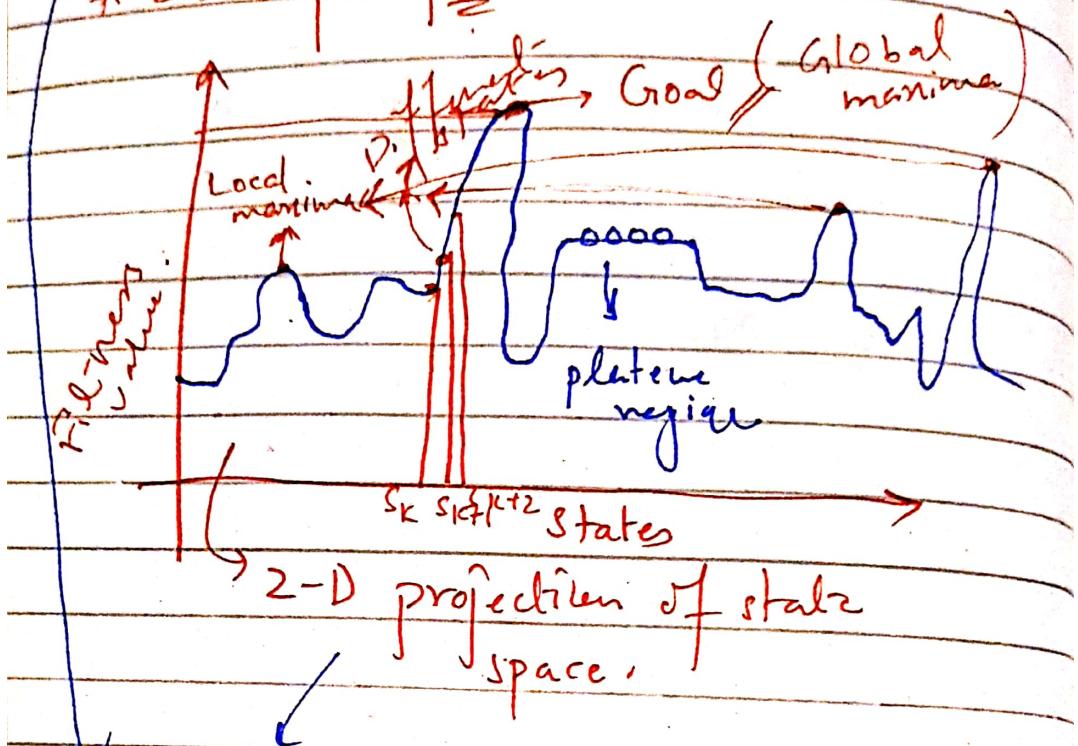
		Q <sub>2</sub>	Q <sub>3</sub>	Q <sub>4</sub>
↓	Q <sub>1</sub>			

S<sub>4</sub>

$$f(s_4) = 5$$



\* Landscape of Local Searches:-



→ In reality, it depends on the branching factor. Using 2-D graph for simplicity.

→ need b+1 dimensional space.

Three types of problems:-

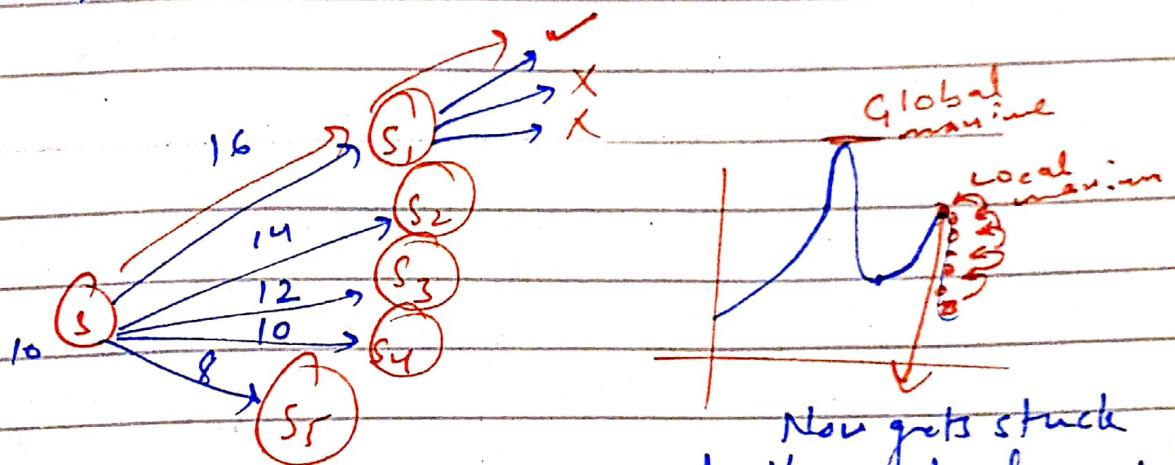
① Local maxima:- Algorithm is stuck in local maxima. (Local region of search landscape having the maximum value but less than global). <for  $10^8$  states>

② Ridges: Series of local maxima connected together.

③ Plateau Region: A region where fitness fun. is flat.

① Hill climbing Algorithm

→ Choose the highest fitness state from the current state.



problem with hill climbing: → ridges & plateaus.

→ So, solution to this problem is resetting using randomness when gets stuck.

Simulated Annealing → Physical randomness  
Cuckoo → Biological randomness  
Ant colony →

Swarm → : each have their own way of activity  
: randomness

## Simulated Annealing

$\Delta E$  is fixed //

High  $T \rightarrow$  high prob  
 $T = 1000 ; P =$

$T = 1 ; P =$

Low  $T \rightarrow$  Low prob

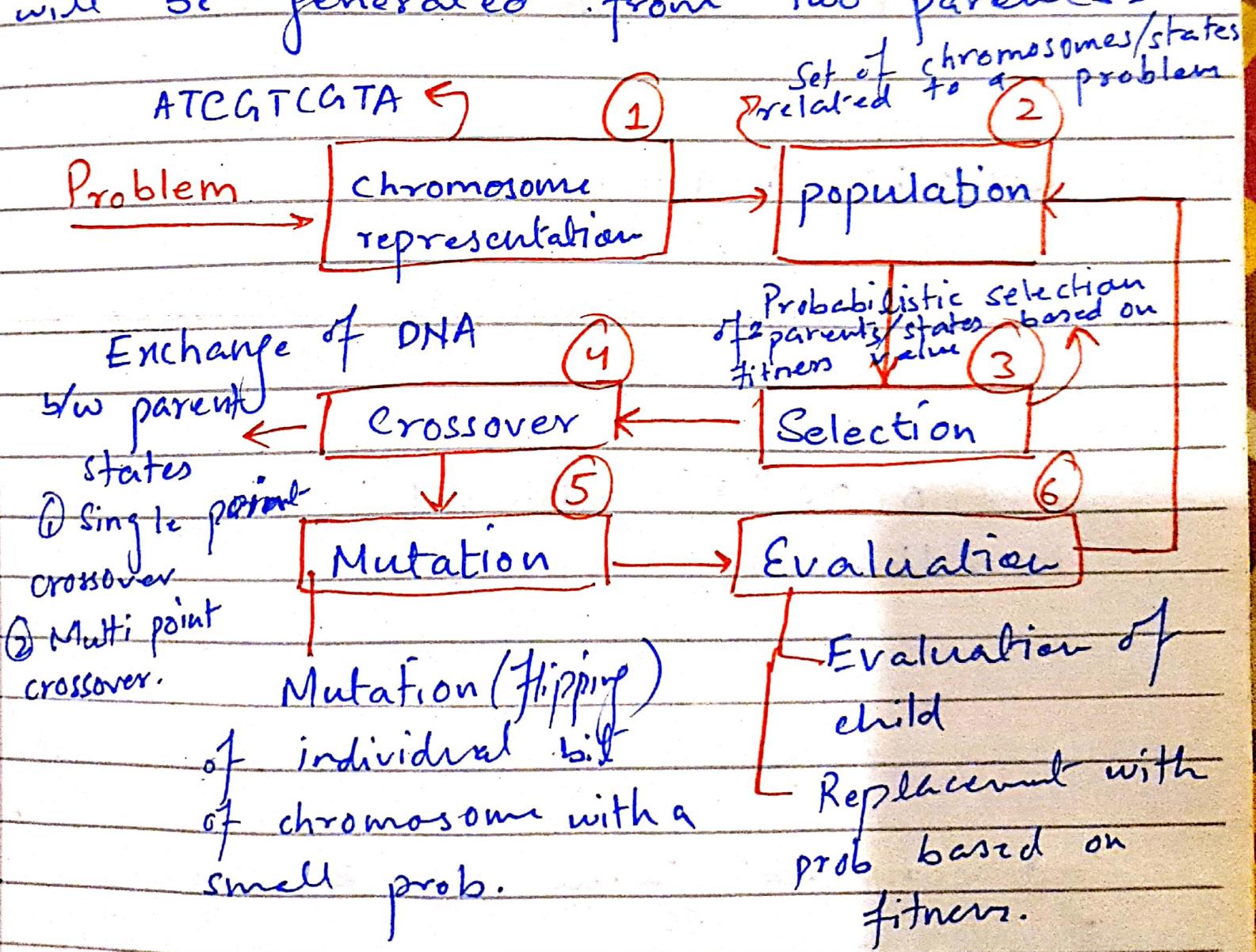
$$P_i = e^{\frac{\Delta E}{T}} \rightarrow \text{Increase in } T$$

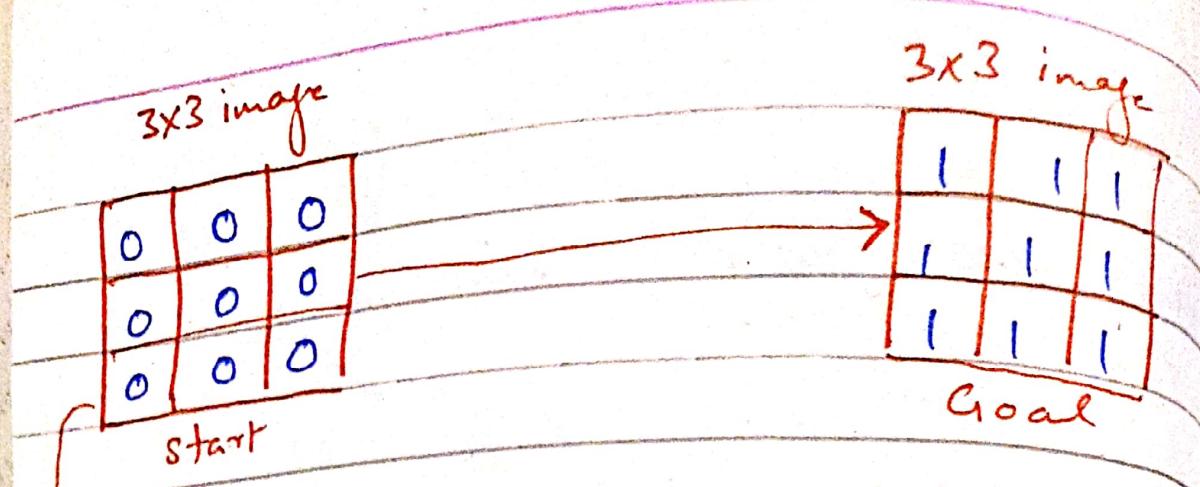
will not decrease  $P_i$  because  $\Delta E$  is negative.

# ★ Genetic Algorithm

→ Randomness inspired from biological evolution.

→ Parallel search algorithm based on sexual reproduction, i.e next state will be generated from two parents.





$f(n)$  = The # of 1s in a state n.

$f(n)$  = The # of 0s in a state n.

Chromosome Representation -

ATCGTCTGA

→ This can be represented in multiple ways.

① 101010010 (Binary)

② 123412341      A=1

T=2

C=3

G=4

→ Convert from 2D to 1D.

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

1	2	3
4	5	6
7	8	9

→ Every chromosome is a solution & has a fitness value

$$f(\text{start}) = 0$$

→ Also given a chromosome, can be converted into a state.

$$S_1 \quad 010 \dots 000 \quad 100 \\ f(S_1) = 2$$

## 2. Population:-

Generative individuals randomly.

Randomly generated population chromosomes

P <sub>1</sub>	0 00	100	. 000	f(n) = # of 1's.
P <sub>2</sub>	100	001	010	
P <sub>3</sub>	001	000	100	
P <sub>4</sub>	000	110	000	

Population ↑ = Diversity

↓ = No diversity

→ Get their fitness value.

$$f(P_1) = 1, \quad f(P_2) = 3, \quad f(P_3) = 2, \\ f(P_4) = 2$$

## 3. Selection:-

$$P_i = \frac{f(P_i)}{\sum_j f(P_j)}$$

$$P_1 = 1/8 = 0.125 = \frac{12.5}{100} \approx 13/100$$

$$P_2 = 3/8 = 0.375 = \frac{37.5}{100} \approx 38/100$$

$$P_3 = 2/8 = 0.25 = \frac{25}{100} \approx$$

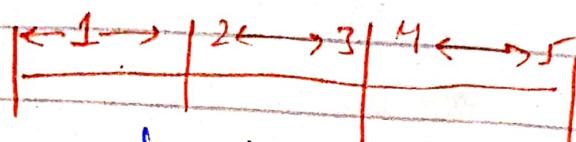
$$P_4 = 2/8 = 0.25 = \frac{25}{100} \approx$$

	$P_1$	$P_2$	$P_3$	$P_4$	
1	13	50	75	100	

→ Randomly generate number e.g. 40, so we select  $P_2$ .

Leftover:  $P_1, P_3, P_4$

$\frac{1}{3}$ ,  $\frac{2}{3}$ ,  $\frac{2}{3}$



→ So randomly select  $P_2$ 's pair, which can be either one of them.

①  $P_2 \rightarrow P_3$

②  $P_1 \rightarrow P_2$

!

④

Crossover

Single point crossover

$P_2$	100 001	010
$P_3$	001 000	100

e.g. ⑥

5/0 1-1

→ Generate random numbers to separate the P.

$c_1$  100 001 100

$c_2$  001 000 010

for  $P_1$  &  $P_2$  e.g. ④

$P_1$	0 0 0	X	1 0 0	0 0 0
$P_2$	1 0 0	X	0 0 1	0 1 0

$C_3$  0 0 0 | 0 1 0 1 0

$C_4$  1 0 0 0 0 0 0 0 0

### ⑤ Mutation--

$$\stackrel{=}{\curvearrowleft} P_m = 0.01$$

Out of 100 bits  $\rightarrow$  1 will be mutated.

e.g.  $c_1$  becomes

1 0 1 1 0 1 1 0 0

### ⑥ Evaluation--

$$\stackrel{=}{\curvearrowleft} f(c_i) \quad f(a) = 9$$

$c_1$	5	→ check for goal state, here
$c_2$	2	
$c_3$	3	
$c_4$	1	

$\rightarrow$  if the goal is found, stop.

Else replace the children in the population based on the fitness of the population (loop).

→ If  $f(a)$  is not known, keep doing it until  $f()$  stops increasing.

2 imp steps → Crossover & mutation

→ because crossover is convergence operator

→ mutation is exploration operator.