

## Blind vs Heuristic Search

Date: \_\_\_\_\_

uninformed

strategies have no additional

info except available in prob definition.

BFS, DFS [no guidance]

informed - know whether one non goal state is more promising than others

A\*, Heuristic Search.

BFS: Expand shallowest unexpanded

nodes + new nodes are

inserted at the end of fringe.

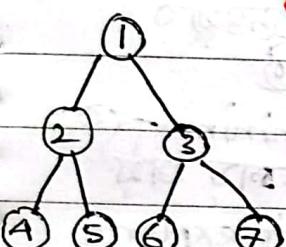
serially check

1. Fringe = (1)

2. Fringe = (2, 3)

3. Fringe = (4, 5)

4. Fringe = (4, 5, 6, 7)



\* explores the nodes at same level before exploring next level.

\* unweighted graph → shortest

Advantage path normally (optimal)

\* Simple → Solution

Disadvantage

Memory, time, large graph → large part search

When we use:

① no info. abt goal / search space.

② small search space

③ guaranteed completeness is prioritized not performance.

A\* Search: source node (actual cost)  $\rightarrow g(n)$ ; goal node (estimated cost)  $\rightarrow h(n)$ ;  $f(n) = g(n) + h(n)$ .

\* Finds shortest paths in a weighted graph

from a node → cost

Example - shortest path finder

Problems: in maps.

\* Poor heuristic value = inefficient

\* consumes memory

Solution: ① Admissible Heuristic; truecost →

② use iterative deepening A\*  
Simplified memory bound  
→ to use less memory.

Memory

Use when:

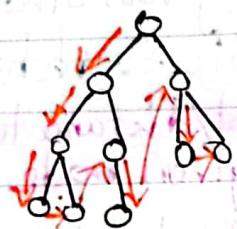
① good heuristic

② complex + large search space.



DFS: \*Explores as far (in depth) as possible before backtrack.

example: puzzles + sudoku:



Problems: ① doesn't guarantee shortest path!

② Gets stuck in depth (if no goal)  
never finding the goal.  
never terminates

Solutions:

Depth limited search: k-limit

set cost function add  $f(n)$  node

मूले searched एवं ना

3 possible outcome:

- ① Solution
- ② Failure (solution not true)
- ③ cut off (no solution in depth limit - but exist वाले)

### Greedy Best First Search

- \* Uses only Heuristic value to choose next node.
- \* selects एवं तरीका node to goal.
- \* 8 puzzles.

Problems: ① fast but not shortest path नहीं गारंटी

② Fails if  $\rightarrow$  local minima - यदि अंतर्गत घास  
or fails to explore alternative path

Solutions:

① cost function add  $f(n)$   
A\* search - यह विकल्प

② if get stuck = backtrack  
and explore others.



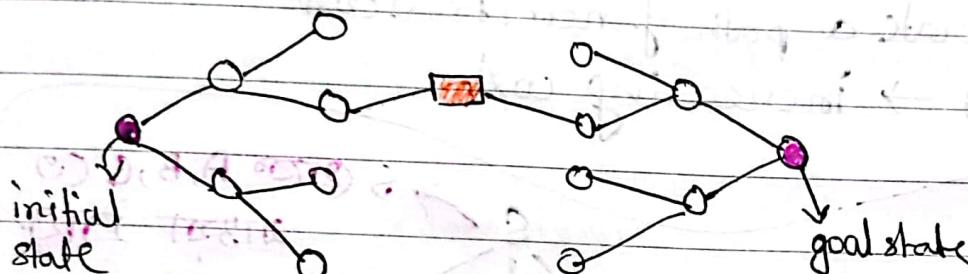
\* Bidirectional Search : To fix BFS (Reduces search space)

■ RUN 2 SIMULTANEOUS SEARCHES .

↳ ① forward from initial state

② backward " goal state .

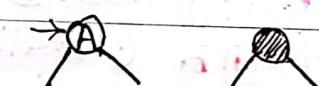
■ Check if the node is in other search tree  
■ if found stop .



\* Iterative deepening Search :

Depth limited version of DFS is run repeatedly with increasing depth limits until goal is found .

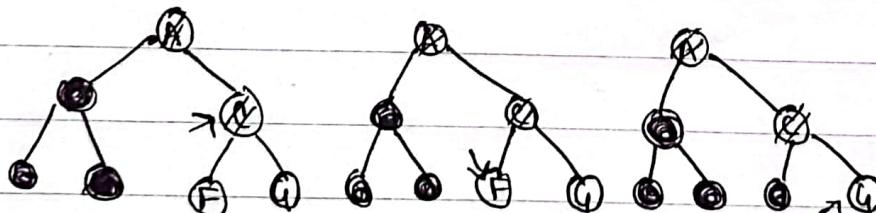
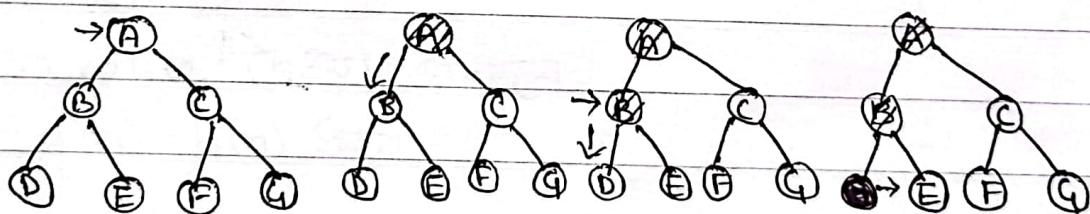
limit = 0 .



limit = 1

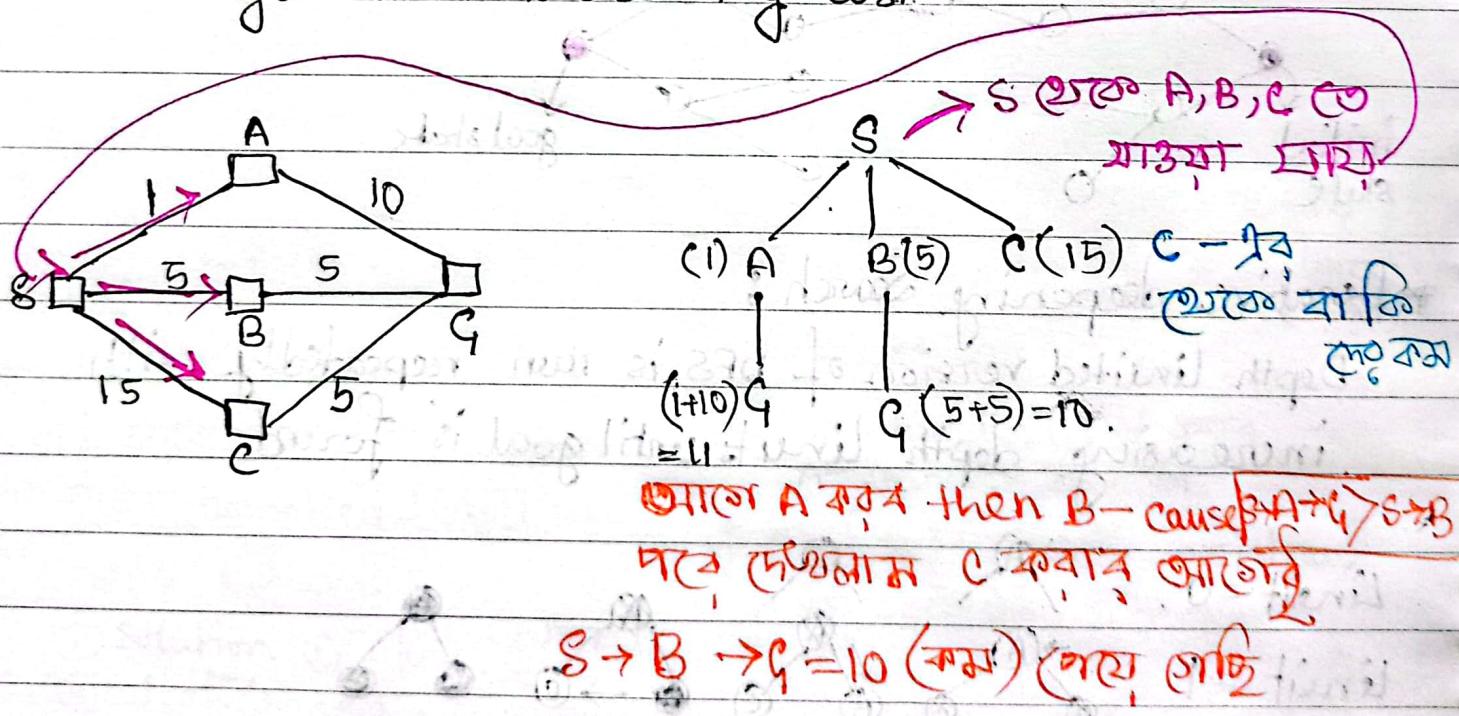


limit = 2 .



\* Uniform cost Strategy: uninformed search algo.  
USES lowest cumulative cost  
to find path from source to destination.

Suppose  $\text{cost} \geq \epsilon > 0$ .  
 •  $\text{wst}$  is the cost of path to each fringe node  $N$ ;  $g(N) = \sum \text{wst}$  of all step.  
 Goal: Generate a path of minimal cost.  
 Fringe sort  $\rightarrow$  increasing cost.



## Informed Search.

Date:

$A^*$  → is best + quick + easy to code + often work VERY well.  
With  $A^*$  vs Greedy Best-first ask করো \*\*\*

### WHY USE HEURISTICS?

- ① add "smarts" to solution
- ② significant speed-ups in practice (fast solving)
- ③ BUT WORST-CASE EXPONENTIAL TIME COMPLEXITY

Goal node - হলে  $h(n) = 0$   
Source node - হলে  $g(n) = 0$   
কারণ



Evaluation function  $f(n)$ ,  
→ provides estimate for total cost  
 $g(n) = \text{source cost} + n \text{ node}$   
Actual cost ACTUAL

$h(n) = \text{goal node cost} + n \text{ node}$   
- cost ESTIMATED

ESTIMATED TOTAL cost

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

SEARCH Efficiency depends on heuristic quality

যত ভালো তত fast search করবে।

Heuristic Definition: A commonsense RULE (on set of rules) intended to increase the probability of SOLVING some PROBLEMS.

GREEDY

BEST-FIRST SEARCH

↳  $h(n)$  → value যাতে একবে

আমরা  $f(n) = h(n)$  করব

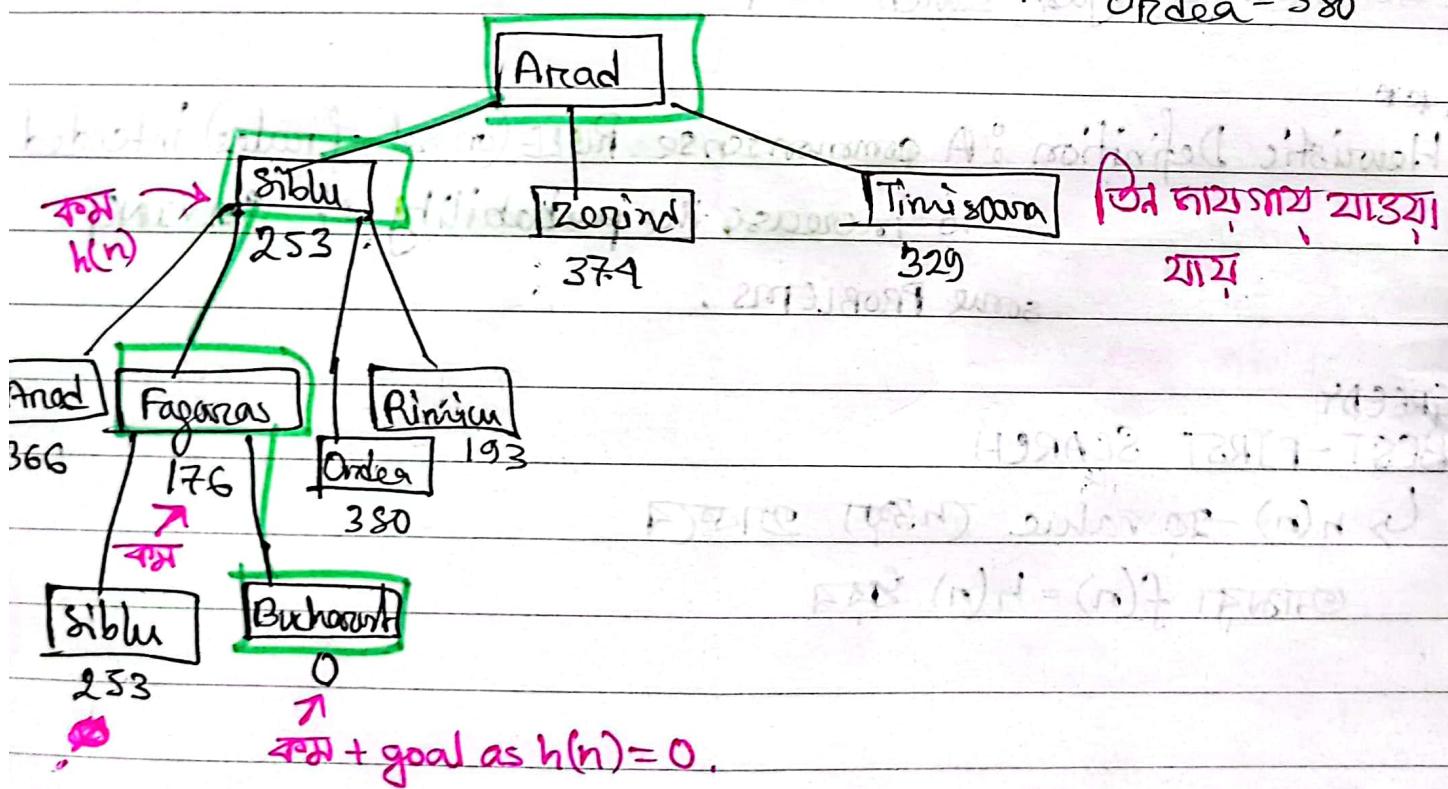
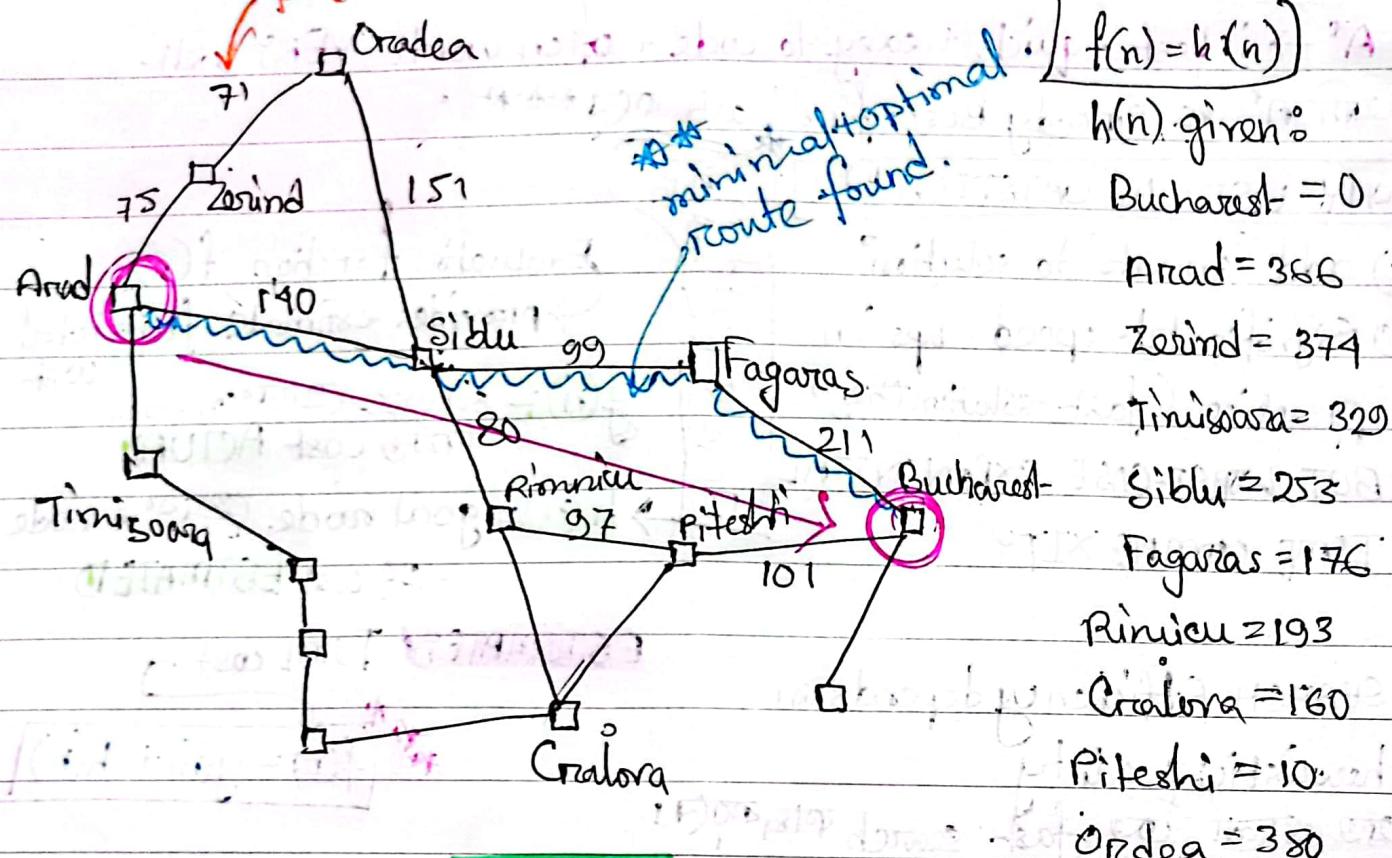
O - O(N log N)



Ahsanullah University of Science and Technology

$g(n)$  - यहाँ लागत है

Date:

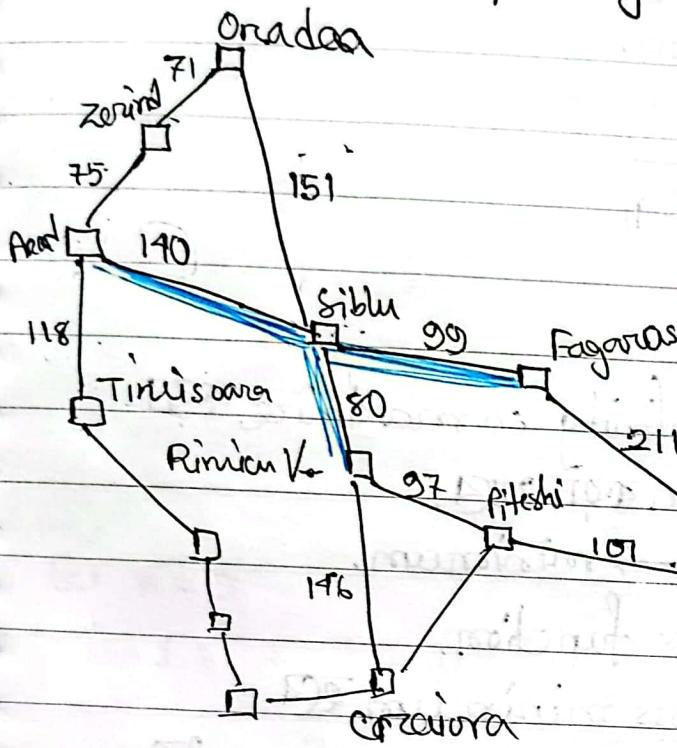


# A\* Search (Best one)

Date:

↳ avoids expensive path [Cheap Thills 😊]

$$\text{evaluation } f(n) = g(n) + h(n)$$



$h(n)$ :

$$\text{Bucharest} = 0$$

$$\text{Arad} = 366$$

$$\text{Zerind} = 374$$

$$\text{Timisoara} = 320$$

$$\text{Sibiu} = 253$$

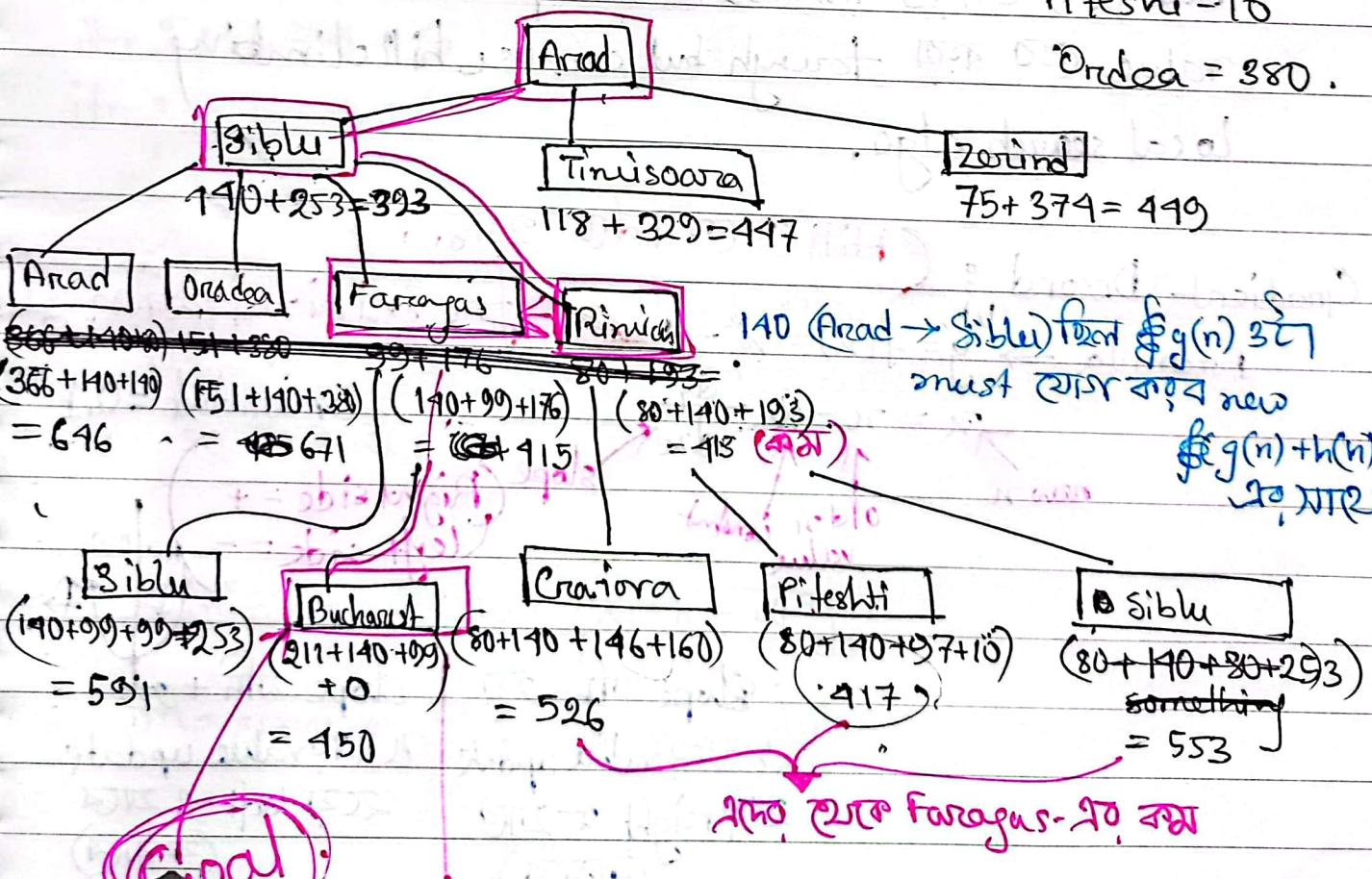
$$\text{Fagaras} = 176$$

$$\text{Rimnicu} = 193$$

$$\text{Craiova} = 160$$

$$\text{Pitesti} = 10$$

$$\text{Oradea} = 380$$



# 8 Queens

## HILL CLIMBING LOCAL SEARCH + OPTIMIZATION STRATEGY

Optimization : maximization

minimization

$$f(x) = x^2 - 4$$

Suppose ~~fito~~ shape graph

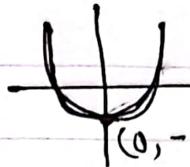
जानि ना

आनि चाहौं Optimize करते

~~द्वितीय~~ function द्वारा maximum value inf possible हो तरीके minimum

$x=2$  value -4 होल  $f(x)=0$

~~द्वितीय~~ निलम  $f(x) = x^2 - 4$  loss function.



Applications :

① IC design

② Telecom networks optimization

③ Automatic programming

④ Factory floor layout design

⑤ Job-shop scheduling

⑥ Vehicle routing

⑦ Portfolio management

- NO Path cost considered

- optimum evaluation/objective/fitness  $\rightarrow$  function.

$x=2$  value होल loss minimum होल

but  $f(x)$  दोबार intense पर्सी expression होल  $x=2$ ,

value होल कर्ता tough but can use hill climbing

local search algo.

Gradient Descent : (Hill Climbing किम ग्राफ़)

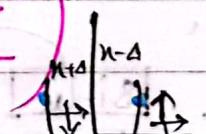
Formula  $\rightarrow$   $y = f(x) = x^2 - 4$  ~~for minimization~~

$$x = x - \alpha \frac{dy}{dx}$$

new  $x$       old  $x$       constant  
 value      value      value

$\alpha$  constant = 0.1

slope (Right side = +)  
 (Left side = -)



Slope  $\frac{dy}{dx} = 2x$

$x=2$  value update

(पास्यु)

Slope  $\frac{dy}{dx} = 2x$

$x=3$  value update

(कर्ता)



$$-\alpha \frac{dy}{dx} \rightarrow \Delta$$

~~20-Slope of function~~

$\alpha \frac{dy}{dx}$  if positive  $\Rightarrow$  then  $x - (\alpha \frac{dy}{dx}) \rightarrow x - \alpha \frac{dy}{dx}$  left side  
 $\alpha \frac{dy}{dx}$  " neg " " "  $x + (-\alpha \frac{dy}{dx}) \rightarrow x + \alpha \frac{dy}{dx}$  right side

Now maths:

$$y = f(x) = x^2 - 4$$

$$x = x - \alpha \frac{dy}{dx} \quad [\because \alpha = 0.1]$$

let  $x = 5$

$$\text{itro 1: } x = 5 - (0.1(2x)) = 5 - (0.1 \times 10) \\ = 5 - 1 = 4$$

$$\text{itro 2: } x = 4 - (0.1 \times 8) = 3.2$$

$$\text{itro 3: } x = 3 - (0.1 \times 6) = 2.4$$

বেসিন VS হাইল্ট টেক্সট  
প্রক্রিয়া কর্মসূচি - ১ ০ এবং ০ এর মধ্যে Algo আর কর্মসূচি না

as merge এর জন্যে also ' $\alpha \frac{dy}{dx}$ ' এর জন্যে ০ আহরণ

so ০-০ হচ্ছে।

Minimize কর্মসূচি  $x = x - \alpha \frac{dy}{dx}$

Maximize কর্মসূচি  $x = x + \alpha \frac{dy}{dx}$



## Features of Hill Climbing Local Search:

- ① Uses SINGLE CURRENT state
- ② INVESTIGATE NEIGHBOURS — generated by successor f<sup>n</sup>.
- ③ If optimal/sufficient value found → Terminates
- ④ Doesn't involve Goal State
- ⑤ Greedy best first search — ~~it's~~ — SELECTS FROM NEIGHBOURS — ~~its~~ heuristic function ~~it's~~ that leads to the "top of the Hill" quickly.
- ⑥ Little MEMORY required.
- ⑦ Better for large + infinite type search space

## Features of ENVIRONMENT:

- ① Hill top : Optimum (Global maximum)
  - ② Foothill top : Not optimum (local maximum); less promising than current state.  
Might be a lot of foothills.
  - ③ Plateau top : Not optimum (local maximum), all successors have same value.
  - ④ Shoulder : like plateau, but at some point goes up.
- 



## Measures to face odds:

- Sideways move : Shoulder एठा, खें, खेठे particular direction -> चातड्याठी try करा
- Random restart hill climbing : local maximum -> आ॒काय जाने - begin with new randomly generated state. Good - if local maxima कम बऱ्हजा उल्लंघन,
- Stochastic Hill Climbing : Choose randomly from uphill moves.
- First choice Hill Climbing : Choose first randomly generated successor that is better than the current state.

### Why named HILL CLIMBING?

- Moves in the direction of increasing value to find the peak of the mountain or the best solution and terminates once reached.

### Why Greedy local Search?

- Looks to its good immediate neighbour state and not beyond - is iterative + starts with arbitrary solution - finds better solution with minimal change.



# Genetic Algorithm.

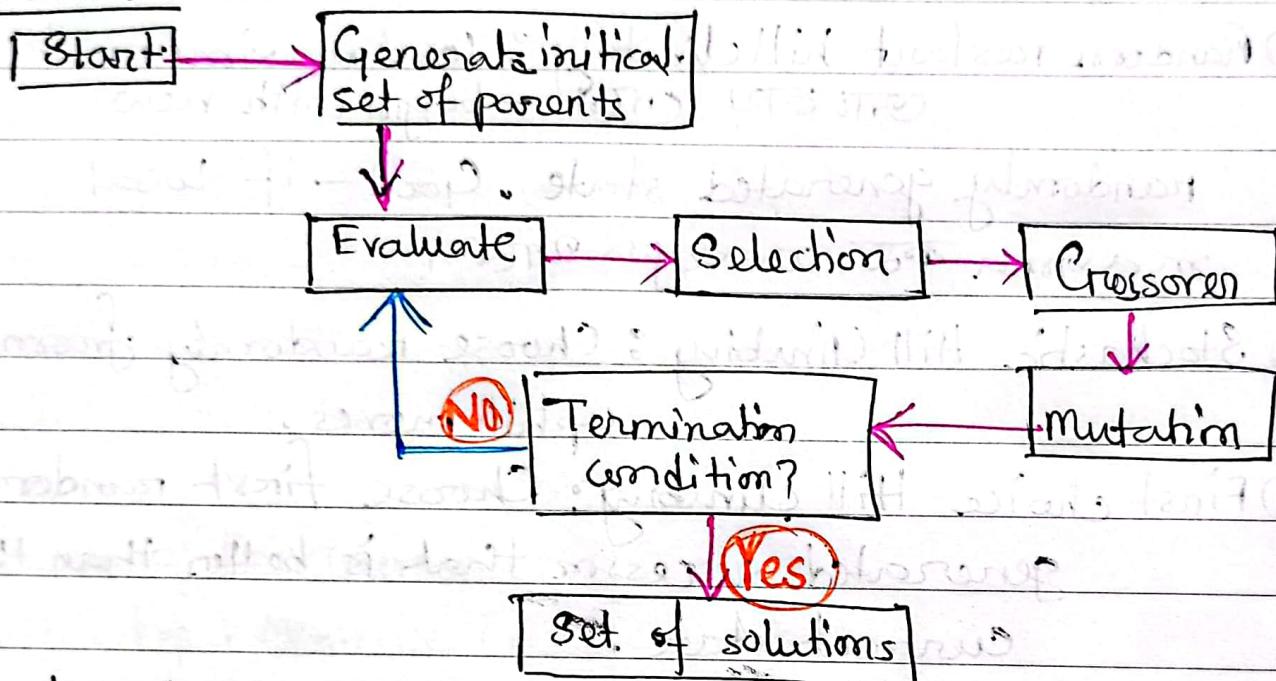
optimization problem solver

Date:

↳ solves optimization problems both unconstrained +  
unconstrained natural selection.

Derived from biological evolution.

Procedure :



Pseudo code :

1. START
2. GENERATE the initial population
3. Compute fitness
4. REPEAT
5. SELECTION
6. CROSSOVER
7. MUTATION
8. COMPUTE FITNESS
9. UNTIL population has CONVERGED
10. STOP



Ahsanullah University of Science and Technology

\*\* Importance of crossover: Without crossover there's just local mutations — won't change at all or make things slow; It will be hard to get population out of local optimum.

\*\* Importance of mutation: Adds diversity of a population. Increases the likelihood of generating individuals with better fitness values.

Suppose a parent has fitness 5 and 6, they will give offsprings near to their fitness but two parents with 20, 6 fitness will either mutate give better fitness or bad but with diverge faster.

Complexity: NQueens  $O(N!)$ . use backtracking.

$$\text{GA} \rightarrow O(gnm)$$

$\nwarrow$   
no. of generation

$\nearrow$  size of individuals.  
 $\searrow$  population size

Example:

$$\text{♂} \rightarrow [3 \ 8 \ 4 \ 7 \ 2 \ 3 \ 2 \ 5]$$

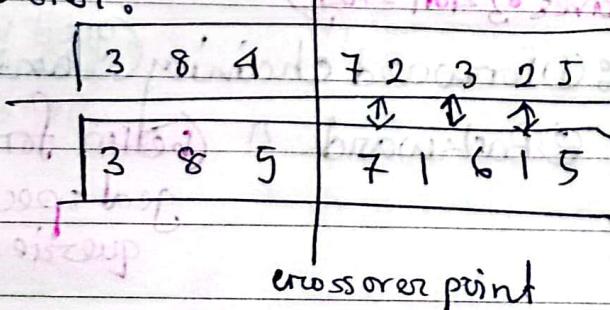
$$\text{♀} \rightarrow [3 \ 8 \ 5 \ 7 \ 1 \ 6 \ 1 \ 5]$$

Before mutation:

$$[3 \ 8 \ 4 \ 7 \ 1 \ 6 \ 1 \ 5]$$

$$[3 \ 8 \ 5 \ 7 \ 2 \ 3 \ 2 \ 5]$$

Crossover:



After mutation:

$$[3 \ 8 \ 4 \ 7 \ 1 \ 6 \ 2 \ 5]$$

$$[3 \ 8 \ 6 \ 7 \ 2 \ 3 \ 2 \ 5]$$



# MACHINE LEARNING

Date: 20/01/2024

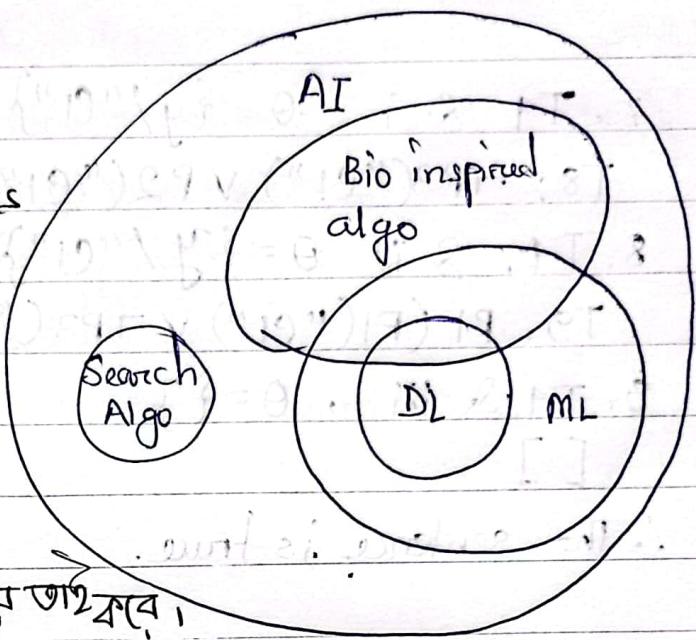
Artificial Intelligence: mimics human intelligence + improves iteratively based on info they collect; faces new situation

\* Why sw not AI?

cause code -এই যা বলা হয়

কিন্তু - does not learn + do not

mimics human. + Dev যা চাষ করে, করে,



Machine Learning: type of AI → allows sw apps to become more accurate at predicting outcomes without being explicitly programmed to do so.  
Not hardcoded — features collect করে ফরেস

\* How ML works?

Dataset is split → Test.

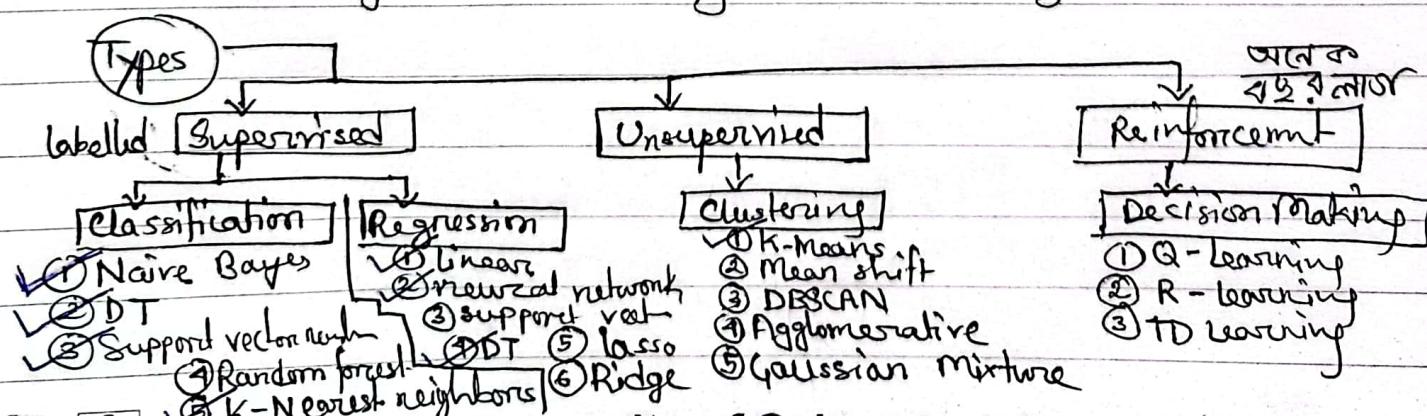
Train

1st → training set input - ফিল্টে  
নির্বাচন করা → 1st epoch.

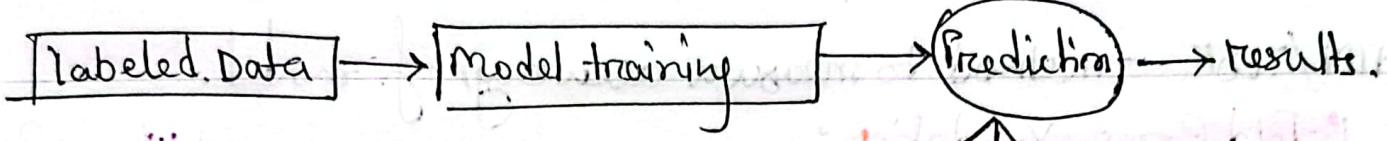
প্রতিপথ prediction evaluate করুন। তোলো (loss এর)  
স্থায়ীভাব (n বৰ্ষী) .

if accuracy is acceptable ML algo is deployed.

বাইলে ML algo is trained again with augmented dataset.



Ahsanullah University of Science and Technology



Supervised Learning: Operator provides dataset - with input outputs.  
Algo finds a method - to arrive at labeled.

### Classification

discrete values [funk/mol, yes/no, true/false etc]

### Regression

continuous value [price, salary, age etc]

Unsupervised Learning: not not labeled but analyzes data for the hidden structures within it. — determines correlations.

Reinforcement Learning: agent is able to perceive + interpret environment = take actions + learn through trial + error.

Can be applied to trajectory optimization, motion planning, dynamic pathing controller optimization, scenario-based plan learning.

### ML life cycle

① Gathering Data

GP TT DM

② Data processing + EDA

③ Train Model

④ Test Model

⑤ Model Deployment

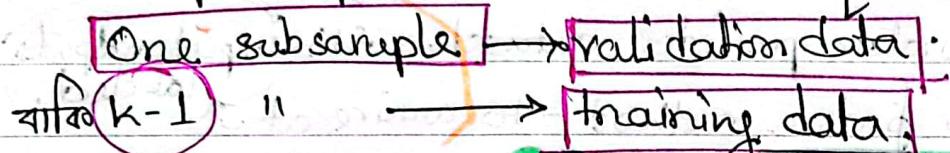
⑥ Model monitoring



**TRAIN / TEST** - method to measure accuracy of model.

### K-Fold Cross Validation :

sample is partitioned into **K** equal subsamples



Cross Validation process is **repeated k times**.

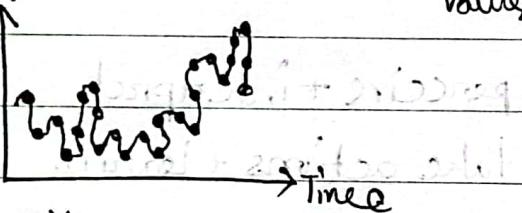
Each **K** subsamples are used only once as validation data

- why use ?
- reduce overfitting
- improves estimation
- handle imbalance
- increase robustness
- avoid data leakage

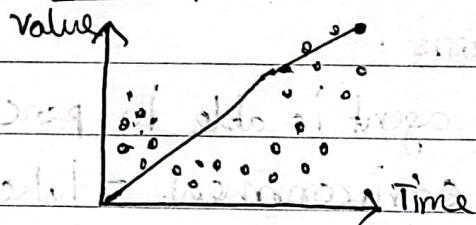

Test data

Training data

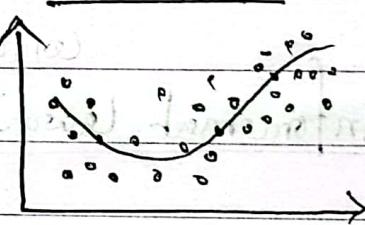
**Overfitted**



**Underfitted**



**Robust / Fit**



\* Fits training data      \* doesn't fit any sets

\* Doesn't fit test      \* remedy : try alternative

\* Model complex      \* ML algo.

\* too many parameters

\* not enough data

\* Ratio of model

complexity to train  
set is high.

→ model memorizes

does good with seen data

but unknown data কেবল সুন্দর  
কি করবে।

### Common reasons:

\* Hidden pattern in collected data

\* Model অনেক type

\* Dataset complexity unknown

\* Overfit - Test → 70%

Train → 99.99%

Overfit → PhD holders

Underfit → All

Good fit → we

for verity.



## Data Pre-Processing

### ① Null value handling

- value missing → delete rows.
- impute mean / median for missing values.
- Categorical columns imputation method.

### ② Feature Selection

- Prob. statement - 20 Domain knowledge + features का नियमितीया का अनुसार चुना जाएगा to select best features. Example : car price prediction problem - (features: manufacturer; license no. etc).
- missing values (due to data corruption + failing to record).
- - imputing the missing value may not match real data, so it's better to drop the column/features

#### \*\*\* FORWARD FEATURE SELECTION technique

↓ (n features selected on prev result)

■ Train those n features + evaluate performance.

■ Best performed feature is finalized.

■ Repeats until desired no. of features are achieved.

finds subset of best performing features for the next

#### \*\*\* How to understand if data is balanced?

→ Suppose dataset of 1000 samples → ~~split into~~ (500) (500)

60 - 40 → balanced

but if (900) (100) → imbalanced

70 - 30 → handleable imbalance

80 - 20 → handleable imbalance

90 - 10 → imbalanced



To combat imbalanced classes in dataset, we use

- ① **Under sampling**: decrease sample size of minority class (remove representative samples)

 useful কুল বাদ দেওয়া যাবে বা।

- ② **Over sampling**: increase sample size of minority size (could introduce noise & redundancy)

 noisy data might get overfilled as they'll replicate  
cleaning noisy data is best to avoid it.

950 - yes

50 - No

500  
50

950  
500

500  
500

under  
(maj -)  
(min +)

over  
(min +)

both applied  
better

\* **AUC-ROC**: curve helps to visualize how well ML classifier performs.

It works → binary classification problems.

Area Under the Curve of (Summary of ROC Curve): measurement of ability of a binary classifier to distinguish btw classes.

Receiver Operating Characteristic

→ plots TPR against FPR

		actual	
		Pos	Neg
pred	Pos	TP	FP
	Neg	FN	TN

$$\text{True Positive Rate / Recall / Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{False Negative Rate, FNR} = \frac{\text{FN}}{\text{TP} + \text{FN}}$$

High TNR  
Low FPR

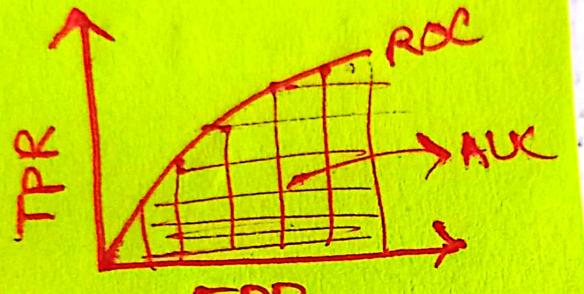
Desired

$$\text{TNR / Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}}$$

$$\text{FPR} = \frac{\text{FP}}{\text{TN} + \text{FP}} = 1 - \text{Specificity}$$



\* AUC-ROC curve helps to visualize how well a classifier performs.



Higher the AUC the better perfm. btw pos & neg classes



② Over sampling

could introduce noise in the sample

noisy data might  
cleaning noisy data

500

500

under

500

(maj -)

## \* Probability of PREDICTIONS

- 2<sup>n</sup> data point - 2<sup>n</sup> actual class predict करना
- or predict the probability of it belonging to diff class.  
(more control over the result)
- positive class classify करने तक data points 2<sup>n</sup> diff thresholds set करने से Sensitivity(TPR) and Specificity(TNR)  
Change 2<sup>n</sup> याएँ।
- यद्यों ना पड़ी threshold better result - फिर बढ़िया होगा — depends on whether to lower FN or FP.

## ML Application

- ① Google translate
- ② Faster route selection (GMNP)
- ③ Automated self driving
- ④ Face recognition smartphone.
- ⑤ Speech recognition
- ⑥ Ads recommendation
- ⑦ Netflix "
- ⑧ Auto friend tagging suggestion
- ⑨ Stock market trading
- ⑩ Fraud detection
- ⑪ Weather Prediction
- ⑫ Medical Diagnosis
- ⑬ Chatbox
- ⑭ ML in agriculture

## ML Benefits

- ① Works **Automation**
- ② Powerful **predictive ability**
- ③ increased in **sales** in the ecommerce market
- ④ **Medical diagnosis - Drug development**
- ⑤ **Robotic medical surgery**
- ⑥ **Finance** : increases productivity enhances revenue, secure transactions.
- ⑦ Modeling data to make **useful decisions**.



## Linear Regression

 Simplex

- \* Data are modeled using straight line.

- \* fits a straight line

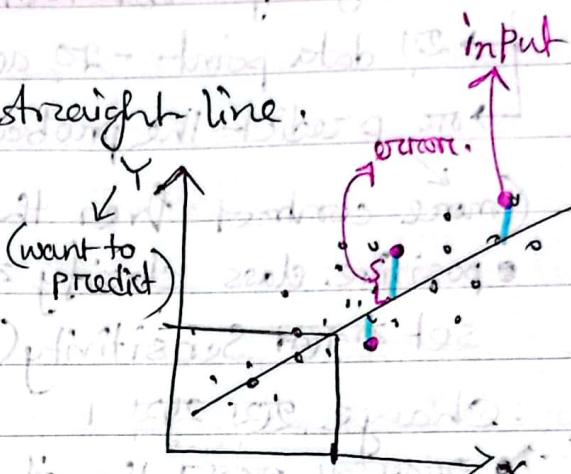
- \* Function,  $y = f(x) = ax + b$

- \*  $x$  - input variable

- \*  $y$  - output variable

- \* looks for how close the data

: to the line are. (or close to the line are) 



4 metrics used to evaluate prediction error rates and model performance:

① MAE : (Mean Absolute Error) - shows difference btwn actual and predicted value extracting avg absolute difference over data set

② MSE : (Mean Squared Error) - squared  $MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|$  with the avg. difference

③ RMSE : (Root Mean Squared Error)  $MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$  square root of MSE

④ R-squared : (Coefficient of determination)

0-1 percentage এর মাত্র

The higher the better

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$$

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

$y \rightarrow$  actual

$\hat{y} \rightarrow$  predicted

$\bar{y} \rightarrow$  mean value of  $y$



## K means Clustering (Unsupervised)

- input variable ~ mean value -এই respect-এ distant measure করবি.

- \* take k as input

- \* partitions the set into  
k subset

- \* Distance is measured

- ↳ mean value of positions of samples in a cluster.

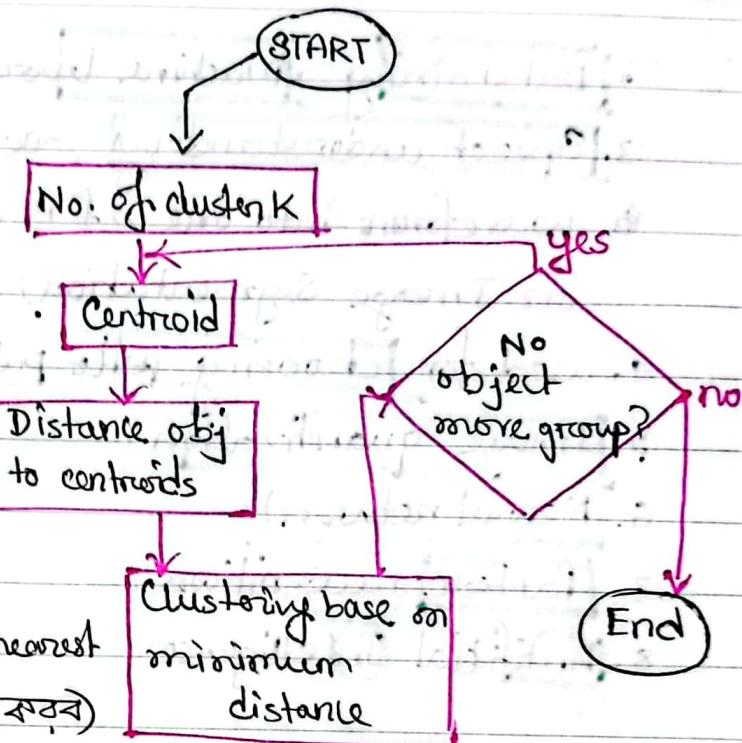
- ↳ center of gravity

Centroid

### 2 major steps:

1. Data assignment: (point যুক্তাকে nearest centroid-এ assign করব)

2. Centroid update step



\*\*\* কার্যন terminal এবং : ① No data points changes cluster

② sum of distance is minimized.

③ iterations reached maximum.

### Weakness:

- ① Small dataset -এ ফোর্মেট সুবিধা নাই

- ② K select করা difficult, cannot visualize at first.

- ③ We never know the real cluster - same data differently input file clustering might differ.

- ④ Sensitive to initial condition. Different initial condition dif same data - তবে different cluster থাকে। local optimum-এ আসেগায়



Application:

- no. of iterations  
no. of obj or points  
K-no. of clusters
- \* 1. Efficient & fast  $\rightarrow O(tkn)$  complexity.
  - 2. Data mining, Machine learning - to use 22J. other fields
  - 3. Speech understanding - acoustic data ; to convert waveforms into one of 15 categories (- Vector Quantization or Image Segmentation)
  - 4. used for choosing palette on old fashioned display devices
  - 5. Image quantization
  - 6. Neural network
  - 7. Pattern recognition
  - 8. Artificial Intelligence
  - 9. Classification analysis
  - 10. Image processing
  - 11. Machine Vision

maths

Suppose we have four samples of medicine (weight, pH value).

Sample : (1, 1), (2, 1), (4, 3), (5, 4) and  $K=2$  given.

Solution:

$K=2$ , so 2 clusters, 2 centroids.

Let initial centroids be (1, 1) (2, 1)

[Step 1] Formula for distance (Euclidean distance) =  $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$   
 $\text{Dist का रूप = } \sqrt{\sum_{j=1}^n (x_j - x_{j1})^2 + (y_j - y_{j1})^2}$  (Cluster 1)

Sample	Dist from Centroid 1	Dist from Centroid 2	Cluster Determined
(1, 1)	$\sqrt{(1-1)^2 + (1-1)^2} = 0$	$\sqrt{(2-1)^2 + (1-1)^2} = 1$	Centroid 1
(2, 1)	$\sqrt{(1-2)^2 + (1-1)^2} = 1$	$\sqrt{(2-2)^2 + (1-1)^2} = 0$	Centroid 2
(4, 3)	$\sqrt{(1-4)^2 + (1-3)^2} = 3.61$	$\sqrt{(2-4)^2 + (1-3)^2} = 2.828$	Centroid 2
(5, 4)	$\sqrt{(1-5)^2 + (1-4)^2} = 5$	$\sqrt{(2-5)^2 + (1-4)^2} = 4.24$	Centroid 2



Step 2

Centroid update:

$$\text{Centroid } 1 = \left( \frac{\sum x}{\text{no. of pt}}, \frac{\sum y}{\text{no. of pt}} \right) = (1, 1)$$

Formula: ~~( $\sum x$ ) / no. of pt, ( $\sum y$ ) / no. of pt~~

Centroid Cluster 1 Cluster 2

$$\text{Centroid } 2 = \frac{2+4+5}{3}, \frac{1+3+4}{3} = \frac{11}{3}, \frac{8}{3}$$

Distance update:

Sample	Dist from Centroid 1	Dist from Centroid 2	Cluster Determined
(1, 1)	$\sqrt{(1-1)^2 + (1-1)^2} = 0$	$\sqrt{\left(\frac{11}{3}-1\right)^2 + \left(\frac{8}{3}-1\right)^2} = 3.14$	Cluster Centroid 1
(2, 1)	$\sqrt{(1-2)^2 + (1-1)^2} = 1$	$\sqrt{\left(\frac{11}{3}-2\right)^2 + \left(\frac{8}{3}-1\right)^2} = 2.36$	Cluster Centroid 1
(4, 3)	$\sqrt{(1-4)^2 + (1-3)^2} = 3.61$	0.47	Cluster Centroid 2
(5, 4)	5	1.89	Cluster Centroid 2

Step-3

Centroid Update : Centroid<sub>1</sub> =  $\left(\frac{1+2}{2}, \frac{1+1}{2}\right) = \left(\frac{3}{2}, 1\right)$

Dist Update:

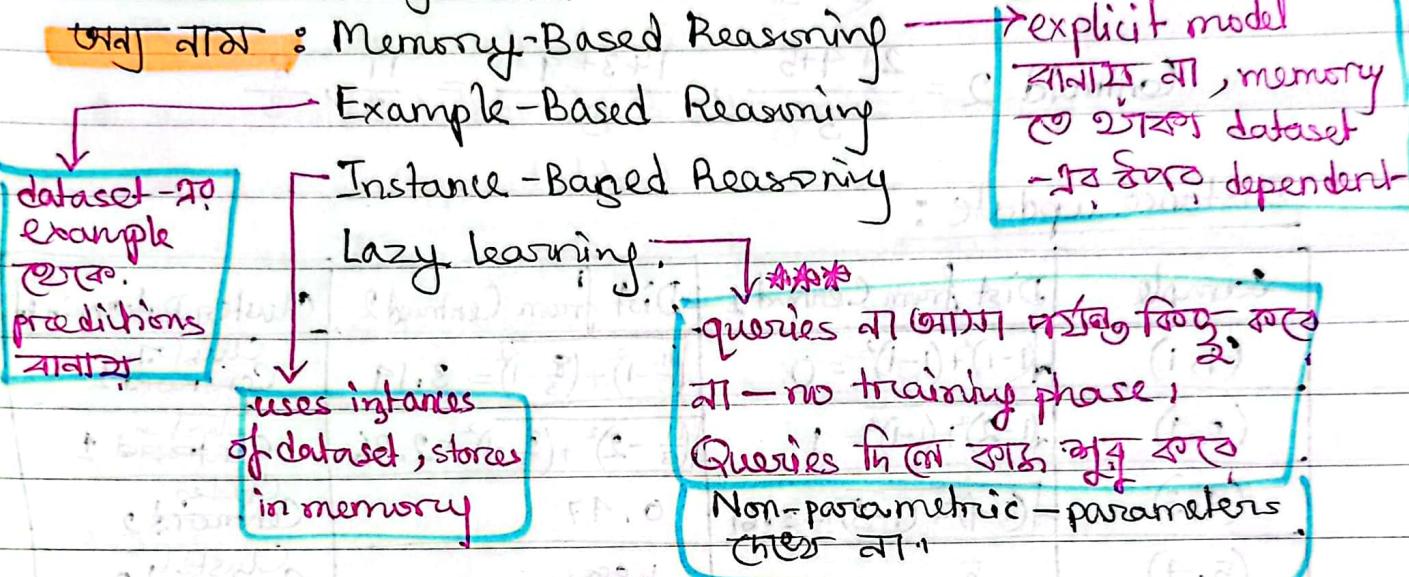
Centroid<sub>2</sub> =  $\left(\frac{4+5}{2}, \frac{3+4}{2}\right) = \left(\frac{9}{2}, \frac{7}{2}\right)$

Sample	Dist from Centroid <sub>1</sub>	Dist from Centroid <sub>2</sub>	Cluster Determined
(1, 1)	0.5	4.30	Cluster 1
(2, 1)	0.5	3.54	Cluster 1
(4, 3)	3.20	0.71	Cluster 2
(5, 4)	4.61	0.71	Cluster 2

Step-3 তে Centroid Cluster update হ্যানি So terminate

করবো।

### K-Nearest Neighbors (KNN)



- \* Used in pattern recognition
- \* One of the top data mining algorithm
- \* Non-parametric
- \* Used in statistics and machine learning.

### K-NN Classification:

- \* Output → class membership
- \* Object কে classify করি যাবু কাছের majority neighbours  
- কেবল behavior দেখো, বিলায়
- \* K selection → positive integer + small

### KNN regression:

- Output → property value of object  
↳ value is the avg of k nearest neighbour's value



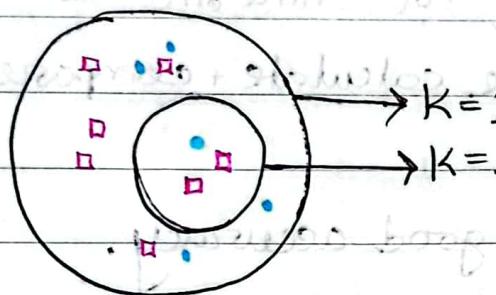
## Classification Approach :-

1. K set করুৰ

2. K-এর value স্বামী neighbours বিটে check কৰুৰ

3. Neighbours - The distance কম = similarities কৰি

" দৈৰ্ঘ্য = " কম



4. Distance measure -  $\exists$  functions (continuous variables) :

$$\text{Euclidean} = \sqrt{\sum_{i=1}^k (x_i - y_i)^2}$$

$$\text{Manhattan} = \sum_{i=1}^k |x_i - y_i|$$

$$\text{Minkowski} = \left( \sum_{i=1}^k (|x_i - y_i|)^q \right)^{1/q}$$

5. ~~Jar~~ যদি কোনো ক্লাসে তাৰা একটা class-এ assigned হৈ  
and soon.

### K-Selection:

\* Too small = sensitive to noise points.

\* " large = will include large K points from other classes too.

\* best way  $\rightarrow K < \sqrt{n}$ ;  $n = \text{no. of examples}$ .



## Strengths of KNN:

1. Simple + intuitive.
2. Any distribution will work.
3. Sample যেমনি এলো করবে তা প্রদর্শিত করবে।

## Weakness:

1. New example classify একটি time নাগে as বর্তনী থেকে  
বাকি দূর distance calculate + compare করা লাভ
2. Choosing k is tricky.
3. Large samples give good accuracy.

## Confusion Matrix (Evaluation)

Given table →

		Predicted Class	
		Positive	Negative
Actual Class	Positive	TP	FN
	Negative	FP	TN

### Formulas:

$$\textcircled{1} \text{ Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\textcircled{2} \text{ Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\textcircled{3} \text{ Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

(যথ যাই করে, divide)

$$\textcircled{4} \text{ FScore} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$



## Naïve Bayes

$P(c|x)$    
 ↓   
 यहाँ तक already. (एव्हाला predictor  $\rightarrow$  Yes/No)   
 यहाँ probability (एव्हाला class)

Posterior probability =  $P(c|x)$

Likelihood " of prediction =  $P(x|c)$

Formula =  $P(c|x) = \frac{P(x|c) * P(c)}{P(x)}$

Class differ फले, while  $P(x)$  <sup>number of prediction</sup> will stay same.

$$\text{So, } P(c_1|x) = P(x_1|c_1) * P(x_2|c_1) * \dots * P(x_n|c_1) * P(c_1)$$

$$P(c_2|x) = P(x_1|c_2) * P(x_2|c_2) * \dots * P(x_n|c_2) * P(c_2)$$

→ प्रमाण बदले करते हुए ना  
 -  $P(x)$  बदल दिये,

$P(c|x)$  - यह value might  
 change but comparison  
 same - ही उपर्युक्त !

## Example: Playing Tennis

### *PlayTennis: training examples*

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

## Math:

Total predictions = 14 (Yes, no feature)

$$P(\text{Play} = \text{yes}) = 9/14$$

$$P(\text{Play} = \text{no}) = 5/14$$

For class ~~to~~ feature - to form PC table:

Learning Phase:

Outlook	Play = yes	Play = No
Sunny	$P(\text{Sunny}   \text{yes}) = 2/9$	$P(\text{Sunny}   \text{No}) = 3/5$
Overcast	$P(\text{Overcast}   \text{yes}) = 4/9$	$P(\text{Overcast}   \text{No}) = 0/5$
Rain	$P(\text{Rain}   \text{yes}) = 3/9$	$P(\text{Rain}   \text{No}) = 2/5$

② Temperature	yes	No
Hot	$P(\text{Hot}/\text{yes}) = 2/9$	$P(\text{Hot}/\text{No}) = 2/5$
Mild	$P(\text{Mild}/\text{yes}) = 4/9$	$P(\text{Mild}/\text{No}) = 2/5$
Cool	$P(\text{Cool}/\text{yes}) = 3/9$	$P(\text{Cool}/\text{No}) = 1/5$

③ Humidity	Yes	No
High	$P(\text{High}/\text{Yes}) = 3/9$	$P(\text{High}/\text{No}) = 4/5$
Normal	$P(\text{Normal}/\text{Yes}) = 6/9$	$P(\text{Normal}/\text{No}) = 1/5$

④ Wind	Yes	No
Weak	$P(\text{Weak}/\text{Yes}) = 6/9$	$P(\text{Weak}/\text{No}) = 2/5$
Strong	$P(\text{Strong}/\text{Yes}) = 3/9$	$P(\text{Strong}/\text{No}) = 3/5$

Test Phase:

let a new instance,

$x' = (\text{Outlook} = \text{Sunny}, \text{Temp} = \text{Cool}, \text{Humidity} = \text{High}, \text{Wind} = \text{Strong})$

( $x'$  situation - 2 yes - 20 probability :)

$$\begin{aligned} P(\text{Yes} | x') &= [P(\text{Sunny}/\text{Yes}) P(\text{Cool}/\text{Yes}) P(\text{High}/\text{Yes}) P(\text{Strong}/\text{Yes}) P(\text{Yes})] \\ &= 2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/11 = 0.0053 \end{aligned}$$

( $x'$  situation 23rd  $\rightarrow$  No - 20 probability :)

$P(\text{No} | x') = [\text{same but with No}]$

$$= 3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/11$$

$$= 0.0206$$

∴  $P(\text{Yes}) > P(\text{No})$  ans Here, "No"

