

Artificial Intelligence

- Goals of AI :-

- (1) Building a machine that can do human intelligence task
- (2) Intelligent connection b/w perception and action

AI → Artificial [Man-made]
AI → Intelligence [Power of thinking]

Advantages of AI

- (1) Accuracy ↑ Error ↓
- (2) Fast Decision Making
- (3) Reliability is more
- (4) usefulness is Risky

Disadvantages of AI

- Cost ↑
Can't perform beyond the limit
No feelings/Emotions
dependency ↑

Reasons of Boost in AI:

- (1) Software or device can be made to solve Real-time Problems
- (2) Creation of Virtual assistant

[SIRI, CORTANA]

Classification of AI :-

- (1) weak AI → Can't perform beyond its field or limitations.

Examples → Flying machine

→ using logic

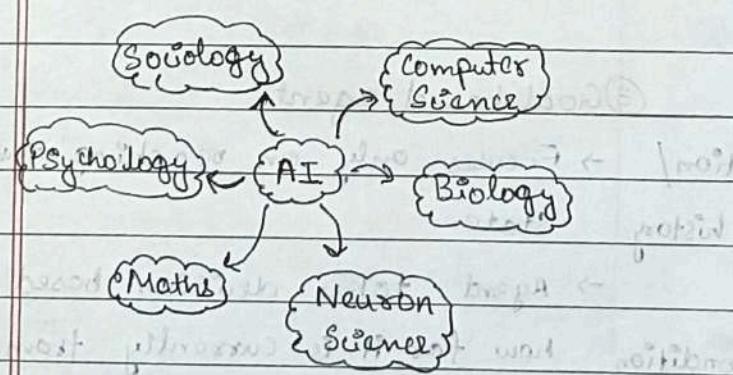
→ Apple Siri

→ Playing chess

- (2) Strong AI → It's the Study & design of machines that Simulate human mind to perform intelligent tasks.

It includes borrowing ideas from Psychology and Neuroscience.

- Forgetting things, Genetics, Languages they keep all this in mind before performing an action.



AI composed of :-

③ Evolutionary AI: It's the study & design of machines that simulate (follow) simple creatures and evolve (learn) for e.g. ants, bees.

④ Super AI → ~~an impossible concept~~

- It's a hypothetical concept for now.
- Designing a machine that can perform better than humans.

★ AI agents:

- AI System is composed of
 - Agent
 - Environment

An agent is anything like:

- Perceiving its environment through **sensors**
- Acting upon the environment through **actuators**

Types of an agent

Limitations of Simple Reflex Agent

- Very limited intelligence
- No knowledge about non-perceptual parts of state
- Can go into infinite loop.

② Model-based agent

It makes decisions based on its current perception of the environment, taking into account its internal model or representation of the world.

Unlike simple reflex agents, that only consider the current percept, model-based reflex agent consider the history of past percepts and actions.

For example → Vacuum cleaning task.

The model agent's internal model would include information about:

- 1) The current location of the agent
- 2) The state of the surrounding area
- 3) Its goal to keep the environment clean

① Simple Reflex Agent

- Works only on current situation / Perception and ignores the history of previous state.
- It will only act on if the condition is true.

③ Goal based agent

- Focuses only on reaching the goal state.
- Agent takes decision based on how far it is currently from the Goal State.
- Every action is taken to minimize distance to Goal State.
- more flexible agent.

④ Utility - Based Agents

- useful when there are multiple possible alternatives and agent has to choose in order to perform best action.
- Act based not only on goals but also the best way to achieve goal.

⑤ Learning Agents

- can learn from its past experiences
- Starts to ACT with basic knowledge and then able to act by adapting learning.
- components:

a) Learning Element → Makes improvement in System by learning from environment

b) Critic → Gives feedback about agent's Performance based on Standards.

c) Performance element → Selects the action to perform

d) Problem Generators → Suggests the action

★ PEAS : Grouping of AI Agents used to group similar type of agents together.

PEAS

- Performance measure → It's the output we get from an agent
- Environment → All Surrounding things and conditions
- Actuators → devices through which an agent performs
- Sensor → Devices from which agent perceives observation from environment

Example, Self driving car

- (P) → comfort, Safety, Time, Legal driving
- (E) → Condition of roads, Crossing, Traffic Signals
- (A) → Steering, Brakes, Horn, Accelerators
- (S) → Camera, GPS, Speedometer

★ Classification of Environment

① Accessible, and Inaccessible

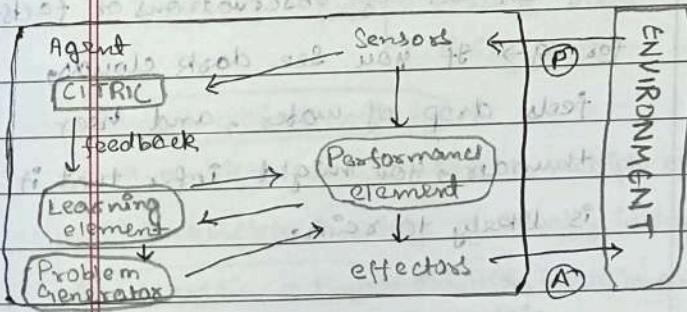
Agent can obtain complete & accurate info about states environment

→ for eg → Room with temp as state (accessible)

→ Event on Earth (Inaccessible)

② Deterministic, or Non Deterministic,
[stochastic]

Agent's Current State and Selected action can completely determine environment's next state. Doesn't have to worry about uncertainty



③ Static or Dynamic

Environment doesn't change its state with passage of time.

for eg → Crossword Puzzles (Static)
→ car driving (Dynamic)

④ Discrete or Continuous :-

finite no. of percepts and actions

for eg → Chess game (Discrete)
→ Self driving (Continuous)

⑤ Observable or Partially Observable

Agent Sensor can access complete state of environment at each point of time

In this some part of environment is not in the reach of agent

⑥ NLP

- understanding
- language generate
- Lang translation

⑦ Reasoning

⑧ Robotics (Locomotive)

⑨ Formal task

Its little complex compare to mundane but not as complex as Expert systems.

- Mathematical Equations
- Geometry
- Logic
- Game theory

⑩ Expert task → Its more complex than mundane and formal task.

- Engineering
- Manufacturing
- Monitoring
- Scientific, finance, Medical.

★ Different tasks in AI

⑪ Mundane (Easiest to learn)

→ Machine / Human learns Mundane (ordinary) tasks since their birth.

→ Learn by → Perception

- Speaking
- Using language
- Locomotives

⑫ Perception

- Computer vision
- Speech, Noise, etc.

NOTES

• Abduction → Its a form of reasoning where one generates the best possible explanation / hypothesis to account for a given set of observations or facts.

for eg → If you see dark clouds, feels drop of water, and hear thunder, you might infer that it is likely to rain.

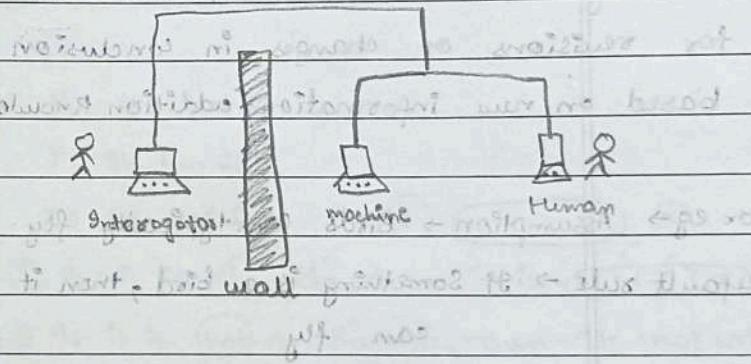
- Induction → It's a mode of reasoning in which general conclusions or patterns are derived from specific observation or examples.
for eg → If you observe multiple instances of apples falling from trees and each time you see it, they fall toward the ground, you might induce that all apples fall toward the ground.

Uncertainty may be caused by problems with data such as:

- Missing / ~~or~~ unavailable data
- Unreliable / ambiguous data
- Imprecise / Inconsistent representation of data

- deduction → It's a mode of reasoning where conclusions logically follow from given premise or statements
for eg → All men are mortal (general premise). Socrates is a man (specific premise). Therefore, Socrates is mortal (deduction conclusion)

★ Turing Test → Introduced by Alan Turing in 1950
→ It's used to determine whether or not machines can think intelligently like humans
→ There will be a (human Interrogator) on one side of wall and other side a machine and human
→ It states that if the human Interrogator cannot distinguish the response given by machine and human then the machine has passed the test & is intelligent.
→ i.e. if the answer is same = machine is intelligent.



Sources of uncertainty

- ① Uncertain inputs → Missing Data
→ Noisy Data

★ AI technique
Also known as ~~as~~ It refers to a specific method, algorithm, or approach used in the field of AI to enable machine to perform intelligent tasks, learn from data, make decisions or solve problems.

- ② Uncertain knowledge →
→ multiple causes leads to multiple effects
→ incomplete knowledge of causality in domain
- ③ Uncertain outputs → Abduction, induction are uncertain
→ Default reasoning, incomplete deduction inference

Single-State Problem	multiple-State Problem	
→ exact prediction is possible	→ semi-exact prediction is possible	that specifies Tweety as a penguin, which is a flightless bird; we have new information that contradicts our initial assumption.
→ state is exactly known after any action sequence	→ not exactly known, but limited to a set of possible states.	→ one disadvantage → can't use for theorem prove.
→ info can be achieved through sensors	→ we can use reasoning to gather information.	• Default Reasoning It's very common form of non-monotonic reasoning. Conclusions are drawn based on what is most likely to be true!
→ consequences of action is known to the agent	→ consequences are not always known to the agent	→ Default logic is a way of reasoning that allows us to make assumptions or "best guess" conclusions when we don't have complete information. These assumptions are called default rules as they represent what we typically expect to be true in certain situations.
→ it's simple but restricted	→ less restricted, more complex	(Miscellaneous notes)
→ for eg → Vacuum world	→ for eg → Vacuum world but no sensors	→ for example, let's say we have a default rule that says, "if it's raining outside, people usually carry umbrellas." In default logic, we would assume that if we see someone walking outside on a rainy day, they are likely carrying an umbrella, even if we didn't actually see the umbrella.

Notes:

- Non-monotonic → it refers to a type of reasoning system or formalism that allow for revisions or changes in conclusion based on new information/addition knowledge

for eg → Assumption → Birds can typically fly.

Default rule → If something is a bird, then it can fly

Let's consider the following statements:

Statement 1: Tweety is a bird

Statement 2: Tweety is a penguin

Based on our default rule, we would initially conclude that Tweety can fly because it's a bird. However, we encounter Statement 2

→ for example, let's say we have a default rule that says, "if it's raining outside, people usually carry umbrellas." In default logic, we would assume that if we see someone walking outside on a rainy day, they are likely carrying an umbrella, even if we didn't actually see the umbrella.

• Monotonic Reasoning

Once the conclusion is taken, then it will remain same even if we add some other info to existing info in our knowledge.

→ Few disadvantages of monotonic:-

- can't represent real world scenarios
- New knowledge from the real world can't be added.

Characteristics :-

1) Monotonic PS → we can't apply a new rule after applying a rule later on.

2) Partially commutative PS → It's a type of

Production System where the order of a rule application can be rearranged without affecting the final result.

★ Production System

Helps in structuring AI programs in a way that facilitates describing and performing the search process.

Production System consist of :-

1) Set of Rules → what action to take if the condition is true.

2) Knowledge Base → collection of rules, facts

and information that the system uses to make decisions or perform tasks.

3) Control Strategy → It determines how the production system selects and applies rules from the rule base to achieve its goals or objectives.

4) Rule Applier → It applies rules on the control strategy.

Steps to solve the problem :-

→ first reduce problem so that it can be shown in a precise statement.

→ Problem can be solved by searching a path through space. [Start → Goal]

→ Solving process can be modelled.

3) Commutative PS → Both monotonic and Partially commutative.

★ Propositional Calculus

	word	Symbol	example
not	not	\neg	$\neg X$ (not X)
and	and	\wedge	$X \wedge Y$
or	or	\vee	$X \vee Y$
implies	implies	\rightarrow	$(X \rightarrow Y)$ if X then Y
if and only if	if and only if	\leftrightarrow	$(X \leftrightarrow Y)$.

X: gets hot + polluted

Y: gets humid

Z: gets raining

① If it is humid then it is hot $(Y \rightarrow X)$

② If it is hot and humid, then it is not raining $(X \wedge Y) \rightarrow \neg Z$

$$(X \wedge Y) \rightarrow \neg Z$$

(P, Q, R) are variables

Set of equivalence relations or laws:

Commutative $P \wedge Q \equiv Q \wedge P$, $P \vee Q \equiv Q \vee P$

Associative $(P \wedge Q) \wedge R \equiv P \wedge (Q \wedge R)$, $(P \vee Q) \vee R \equiv P \vee (Q \vee R)$

Double Negation $\sim(\sim P) \equiv P$

De-Morgan $\sim(P \vee Q) = \sim P \wedge \sim Q$, $\sim(P \wedge Q) \equiv \sim P \vee \sim Q$

Absorption $P \wedge (P \vee Q) \equiv P$, $P \vee (P \wedge Q) \equiv P$

Law of contradiction

$P \wedge \sim P \equiv \text{False}$

$$\begin{array}{c} P=1 \\ \sim P=0 \end{array} \quad \begin{array}{c} P=0 \\ \sim P=1 \end{array}$$

Law of excluded middle

$P \vee \sim P = \text{True}$

$$\begin{array}{c} P=1 \\ \sim P=0 \end{array} \quad \begin{array}{c} P=0 \\ \sim P=1 \end{array}$$

Law of Impotency

$P \wedge P \equiv P$

$$\begin{array}{c} P=1 \rightarrow 1 \\ P=0 \rightarrow 0 \end{array}$$

Rules of Inference:

① Modus Ponens \rightarrow If ' P ' and ' $P \rightarrow Q$ ' is given to be true, then we can infer that ' Q ' is true.

P : It's a holiday

Q : The school is closed

$P \rightarrow Q$: If it's a holiday, then school is closed

★ Truth table

It shows how the truth or falsity of a compound statement depends on the truth or falsity of simple statements.

Some of truth table example:

① Negation

P	$\sim P$
T	F
F	T

② AND: $(P \wedge Q) \rightarrow \text{True}$ when P and Q both are true

P	Q	$P \wedge Q$
T	T	T
T	F	F
F	T	F
F	F	F

★ Tautology \rightarrow It's a formula which is always true. $\neg\neg P \equiv P$ $\neg P \vee P \equiv \text{True}$ $\neg P \wedge P \equiv \text{False}$

\rightarrow opposite is contradiction (always false)

② Modus Tollens \rightarrow If $\sim Q$ and $\sim P \rightarrow \sim Q$ (not of P implies not of Q) are given to be true, then we can infer that $\sim P$ is true

Eg \rightarrow Show that $(P \rightarrow Q) \vee (Q \rightarrow P)$ is tautology

P	$\sim Q$	$P \rightarrow Q$	$Q \rightarrow P$	$(P \rightarrow Q) \vee (Q \rightarrow P)$
T	T	F	T	T
T	F	T	F	T
F	T	T	F	T
F	F	T	T	T

$\sim Q = \text{School is not closed}$

if it's not a holiday, then school is not closed

$\sim P$

$\sim P$ its not a holiday (True)

all are true
∴ tautology

★ First order Predicate logic

This helps us describe and analyze the relationship b/w objects and properties using symbols, variables, and logical connectives. It allows us to make statements and draw conclusion based on logical rules.

Key components of FOPL:

1) Objects → It allows us to talk about objects, such as people, animals, or any other entities in given domain. We can represent objects using variables like x, y , or constants like a, b .

2) Predicates → It describes properties or relationships between objects. They are statements that can be either true/false depending on the objects involved.

3) Quantifiers → It helps us express general statements or specify specific instances. for eg → "for all" (\forall), and "there exists" (\exists), to make statements about all objects or atleast one object.

Existential

4) Connectives → and (\wedge), or (\vee) and not (\neg)

Properties of Sentence

Every atomic sentence is a sentence

if If S is sentence, then $\neg S$ is sentence

if If S_1, S_2 are sentences:

(a) then $S_1 \wedge S_2$ is a sentence (conjunction)

(b) $S_1 \vee S_2$ is a sentence (disjunction)

(c) $S_1 \rightarrow S_2$ is a sentence (implication)

(d) $S_1 \equiv S_2$ is a sentence (equivalent)

if If x is variable and S is sentence then

$\forall x S$ is a sentence for all x

if If x is variable and S is sentence then $\exists x S$ is sentence

There exist some x for which x is a sentence

Example of Quantifiers

1) All Boys like football

$\forall x: \text{Boys}(x) \rightarrow \text{Like}(x, \text{football})$

2) Some Boys like football

$\exists x: \text{Boys}(x) \wedge \text{Like}(x, \text{football})$

Representing class membership / class inclusion

① Is a → Represents relationship of class inclusion

② Instance → Represent class membership relationship

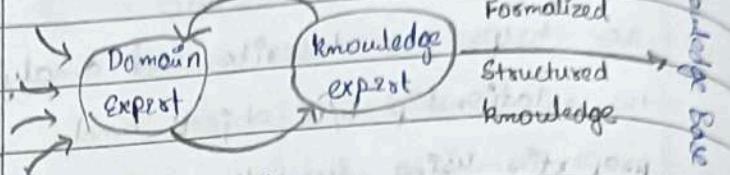
Penguin is a Bird

Eg → Two facts Penguin is Penguin

$\boxed{\text{Bird}}$ all the attributes of Bird are inherited by Penguin class

So, class (Penguin, Bird) \rightarrow ISA

Instance (Pingu, Penguin) \rightarrow Instance



Predicate calculus

Predicates are used that involves constants, variables, relation, function.

can represent
Tall (Aman)

No. of sides ()

Propositional calculus

uses propositions in which complete sentence is denoted by symbol

(Its rainy day) \rightarrow x

can't represent individual entities (Aman is tall)

can't express generalization specialization or pattern

(square has 4 sides)

knowledge from multiple sources

knowledge, concepts, solutions

Tasks in knowledge Base :-

- Collect \rightarrow Acquiring knowledge from expert
- Interpret \rightarrow Review & identify key parts
- Analyze \rightarrow Forming theories & strategies
- Design \rightarrow forming better understanding of problem

Knowledge acquisition techniques :-

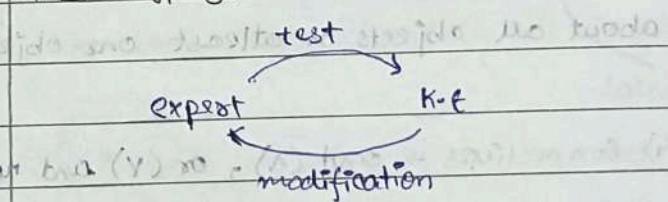
a) introspection \rightarrow Expert acts as both expert and knowledge engineer

b) observation \rightarrow Expert closely observes work

c) induction \rightarrow Converting set of examples into rules.

d) protocol analysis \rightarrow Expert is asked to perform task and verbalized through process

e) prototyping \rightarrow Expert + K-E = System



★ Knowledge acquisition

It's the process of extracting, structuring and organizing knowledge from one or more sources. Knowledge is a collection of specialized facts, procedures and judgments.

rules -

\rightarrow In knowledge acquisition, we acquire knowledge from multiple different sources

\rightarrow In knowledge elicitation, we acquire knowledge from human experts only

f) interviewing \rightarrow Experts verify the knowledge

Knowledge Representation

- Propositional logic
- First order logic
- Rule-based System
- Semantic networks
- Frames

CLASSMATE

Date _____
Page 11

Architecture of knowledge Based System :-

KBS is a System that draws upon knowledge of Human Experts captured in knowledge base to solve problems that normally require human expertise.

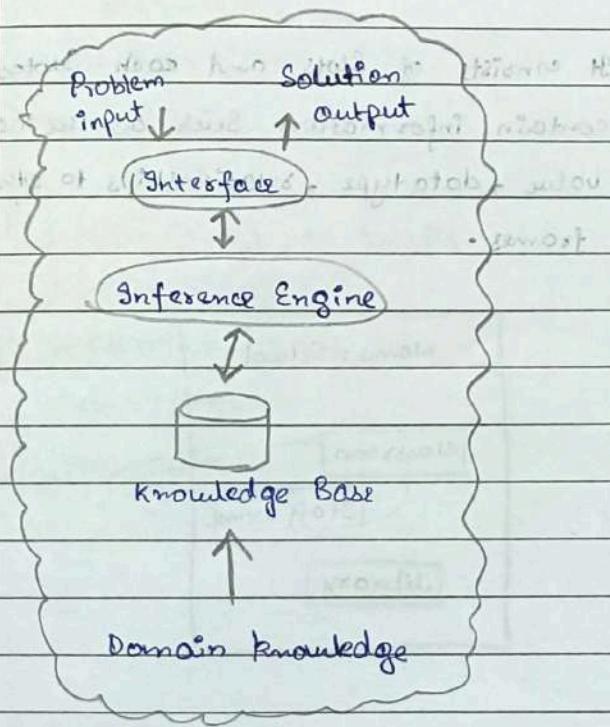
ES is a Specialized type of KBS, not all KBS are designed as Expert Systems.

Components :-

Knowledge Base → Contains organized collection of Facts about a System domain

Inference Engine → Tries to derive answers from Knowledge Base. It's the Brain of any expert System.

User Interface → Enables the user to communicate with KBS.



• Knowledge Engineering → It's the process of digitizing the knowledge and utilizing it to resolve complicated problems. It's an essential part of AI.

Role of Knowledge Engineer is to develop Software or Systems by gathering all relevant information and putting it into systematic format.

Types of knowledge :-

① Procedural knowledge → Describes How to solve Problem

→ Also known as Imperative knowledge

→ Provides direction on how to do something

→ Includes Rules, Strategies and procedures

② Declarative knowledge → Describes what about problem

→ Tells us facts what things are

→ Includes Concepts, Facts, objects.

③ Meta knowledge

→ Describing knowledge about another knowledge

→ used to pick other knowledge that is best suited for solving a problem

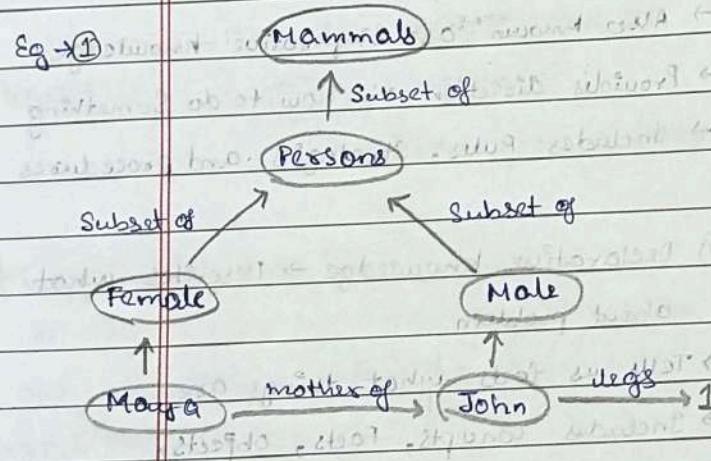
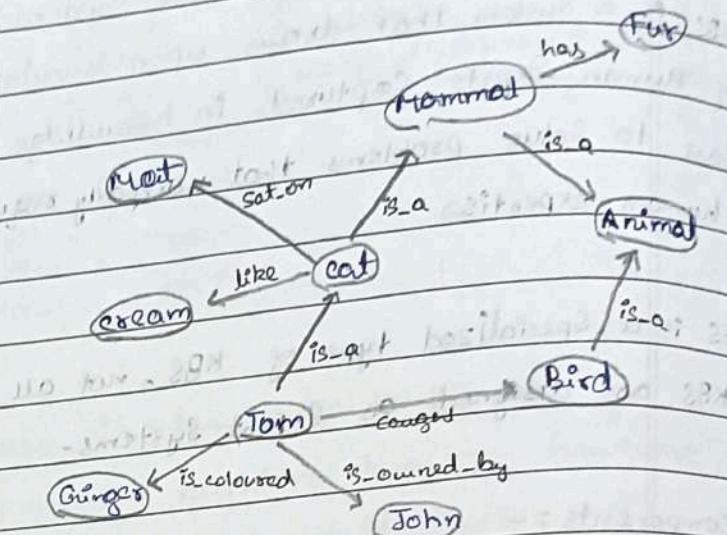
④ Heuristic knowledge → It's often employed in complex & uncertain domains where finding an optimal soln. may be difficult. They don't rely on facts, logical reasoning. It allows to use shortcuts to arrive at a good solution.

⑤ Structural knowledge →

- gts the basic knowledge to problem Solving
- describes relationships b/w Various concepts
- describes an expert overall mental model of a problem

★ Semantic Network

- Graphical notation for representing knowledge in interconnected nodes/pattern.
- Popular in AI and NLP: It represents knowledge or Semantics.



★ Frames

It's a knowledge representation technique used to organize and structure information about objects, concepts, or entities in a hierarchical manner. Frames provide a way to represent and store knowledge about a particular domain or subject area.

It consists of slots and each slot can contain information such as the name, value, datatype, relationships to other frames.

Eg → ②

(Tom is a cat)

Tom caught a bird

Tom is owned by John

Tom is ginger in colour

Cats like cream

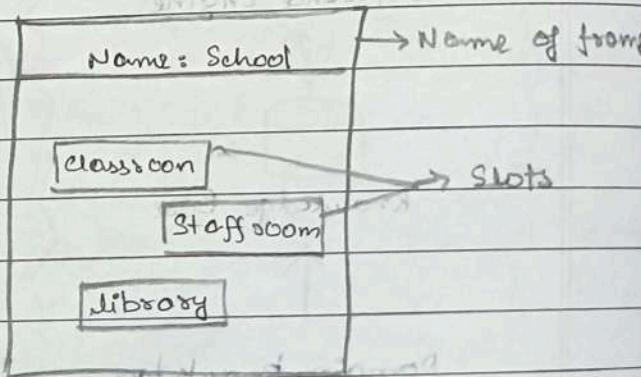
The cat sat on the mat

A cat is a mammal

A bird is an animal

All mammals are animals

Mammals have fur



Operators	Precedence
()	3
* /	2
+ -	1

∴ It's easy to evaluate postfix for the machine hence we convert infix to postfix / prefix

classmate

Date _____
Page 13

Infix, Prefix and Postfix

Notations to write an Expression

Eg → Infix	Prefix	Postfix
$a \times b$	* ab	ab*
$a - b$	- ab	ab -
$A * (B + C) * D$	ABC + * D *	
$x - y * z$	* tab - cd	ab + cd - *
$P - q - r / a$	- x * y z	xyz*-
$(m - n) * (p + q)$	--pq / ra	
		mn-pq+*

① $x - y * z$ to prefix and postfix

Prefix	Postfix
$(x - (y * z))$	$(x - (yz*))$
$(x - (*yz))$	$(x - (yz*))$
$- x * y z$	$xyz*-$

② $P - q - r / a$ to prefix

$\rightarrow ((P - q) - (r / a))$
$\rightarrow ((- Pq) - (r / a))$

$\rightarrow -- Pq / ra$

③ $(m - n) * (p + q)$ to postfix

$((m - n) * (p + q))$

$((m - n) * (pq +))$

$((m - n) (pq +) *)$

$mn - pq + *$

Different Sorting techniques:

① Radix Sort

Here, we don't do any comparison between the numbers. It sorts elements based on their digits or characters.

Eg → 904	001	001	001
046	062	904	005
005	904	005	046
074	074	046	062
062	005	062	074
001	046	074	904

② Stable vs Unstable Sort

3	5	2	1	5'	10	unsorted array
5	1	2	3	5	5'	Stable Sort
6	1	2	3	5	5'	Unstable Sort

In case of Stable Sort, the order of the element in which they occur in a list or array is maintained even after the sort.

③ Internal Sorting

In this style of sorting, the complete data will be there in memory at a time. Because of this the sorting can be completed in memory itself.

Different types of Internal Sort:

- bubble Sort
 - Selection Sort
 - Insertion Sort
 - Shell Sort
 - Quick Sort
 - Merge Sort
 - Heap Sort
 - Radix Sort
- $O(n^2)$
- most $\rightarrow O(n^2)$
- $O(n \log n)$
- $O(n)$

⑥ Bucket Sort

→ work on floating point numbers between range 0.0 to 1.0
 → inputs should be uniformly and independently distributed across (0,1) to get a running time of $O(n)$

Eg $\rightarrow 0.79, 0.13, 0.64, 0.39, 0.20, 0.89, 0.53, 0.42, 0.06, 0.94$

⑦ External Sorting

If the Data collection is too big to fit into the memory then sorting can not be completed in memory alone. Such techniques are called External sorting.

⑧ Counting Sort

→ If provide the size of the input

and the range of the input

(that will be given in the Question)

→ Suppose, inputs $\rightarrow 2, 1, 2, 3, 1, 2, 4$
 Range $\rightarrow (1 \leftrightarrow 5)$

0	(0.08)	$(p+q) \times (n-m)$
1	(0.13)	∴ This is the best case where elements are uniformly distributed.
2	(0.20)	$\rightarrow O(1) \times m = O(m)$
3	(0.39)	$\rightarrow O(n) \times x = O(nx)$
4	(0.42)	
5	(0.53)	
6	(0.64)	
7	(0.79)	
8	(0.89)	
9	(0.94)	

0	x fitting of $(p+q) \times (n-m)$
1	$((p+q) \times (n-m))$
2	$((+ \text{most case}) \rightarrow O(n) \times (n-m))$
3	$\times (+pq) (-m \rightarrow O(n^2))$
4	$x + pq - m$
5	
6	
7	0.74 — 0.78 — 0.79 — 0.795
8	
9	
10	

1	11	It works in linear time but the range is predefined, we can't go outside the range. we might have to take a lot of space unnecessarily
2	111	
3	1	
4	1	
5	0	

No. of Occurrence

∴ Counting Sort $\rightarrow 1, 1, 2, 2, 2, 3, 4$

$\therefore TC \rightarrow O(n+k)$ SC $\rightarrow O(k)$

★ Scripts → used for knowledge representation.	Advantages
→ It's a structure that provides a set of circumstances which could be expected to follow on from one another.	→ Event Prediction is Possible
→ <u>Condition</u> considered to consists of a no. of slots or frames but with more specialized roles.	Disadvantages
	→ Less general than frames
	→ May not be suitable to represent all kind of knowledge.

★ Bayes's Theorem

Describes the probability of an event, based on prior knowledge of conditions.

Components of a Script:

① Roles → Person involved in Event.

(Student)

$$P(B/A) = \frac{P(A/B) \cdot P(B)}{P(A)}$$

② Props → objects involved in Event

(Pen, Answer Sheet)

③ Entry Conditions → Conditions that needs to be satisfied before event occur in Script (ID card should be there)

P(B/A) → likelihood (Prob. of evidence)

P(A/B) → Posterior (Prob. of A when B is true)

P(B) → Marginal Prob. (Prob. of evidence)

P(A) → Prior Prob. (Prob. of hypothesis)

④ Results → Condition that will be true after event in script.

⑤ Traces → Variations on the Script.

(Exam centre)

Q: What's the probability that person has disease dengue with neck pain?

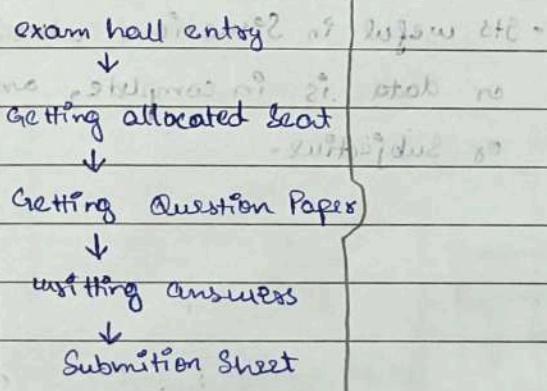
Given → 80% of time dengue causes neck pain

$$P(a/b) = 0.8$$

$$\rightarrow P(\text{dengue}) = 1/30000 \rightarrow P(b)$$

$$\rightarrow P(\text{neck pain}) = 0.02 \rightarrow P(a)$$

⑥ Scenes → Sequence of events that occurs



Sol: a = Proposition that person has neck pain

b = Person has dengue

$$P(b/a) = \frac{P(a/b) \cdot P(b)}{P(a)}$$

$$\Rightarrow 0.8 \cdot \frac{1/30000}{0.02} = 0.00133$$

Applications of Baye's theorem in AI:

- ① Robot/Automatic machine's next step is calculated based on previous step.
- ② Forecasting → weather
- ③ helps in solving Monty Hall problem

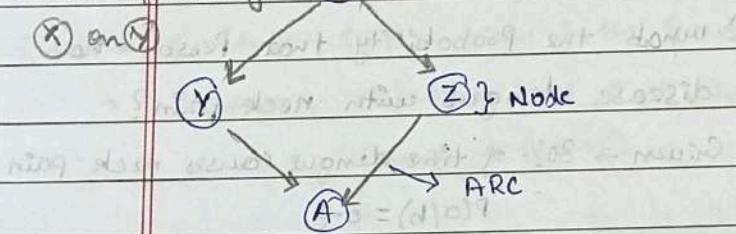
★ Bayesian Belief N/W in AI

• defines probabilistic independence

It's a probabilistic graphical model which represents a set of variables and their conditional dependencies using a directed acyclic graph.

- It consists of Directed Acyclic graph and
- Table of conditional Probability
- It also consists of Nodes and arc

Direct influence of (X) } Parent of Y, Z



(Independent of (X))

Problem →

B	E	$P(A B, E)$
T	T	.95
T	F	.94
F	T	.29
F	F	.001

B	E	$P(B)$
		.001

E	$P(E)$
	.002

A	$P(M A)$
T	.70
F	.01

Q) what is the probability that the alarm has sounded but neither a burglary nor an earthquake has occurred, & both John and Mary call?

$$P(\text{alarm} \wedge \neg b \wedge \neg e) \Rightarrow P(j/a) P(m/a) P(a \wedge \neg b, \neg e) P(\neg b) P(\neg e)$$

↓
j is called
m is called
a is there
no b, no e

$$\Rightarrow 0.90 \times 0.70 \times 0.001 \times 0.999 \times 0.998$$

$$\Rightarrow 0.00062$$

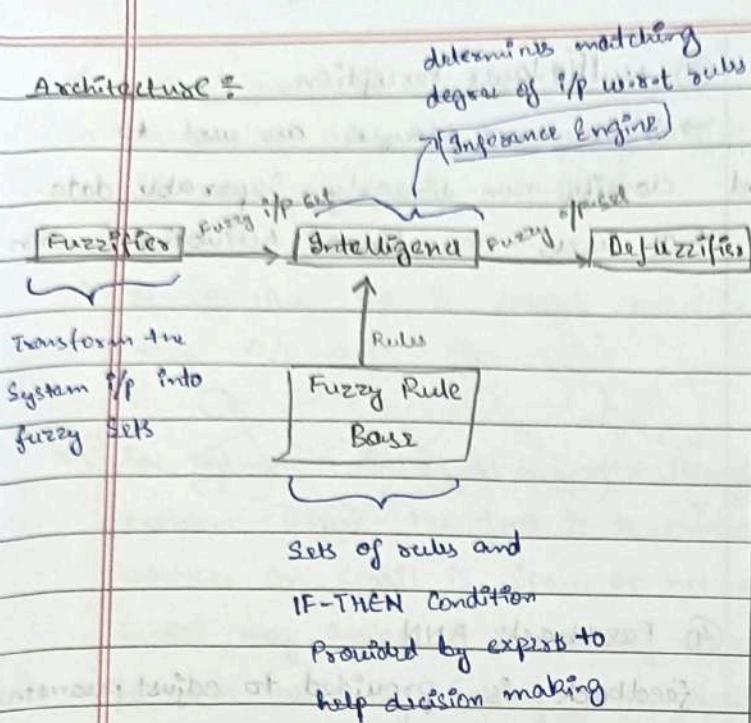
• Probability based Belief

★ Fuzzy logic

• It's a type of reasoning that deals with uncertainty and imprecision. It allows for the representation of vague concepts by assigning degrees of truth to statements or variables, rather than strict true or false values.

• In simple terms, fuzzy logic mimics human decision-making by considering shades of grey instead of relying solely on black and white.

• It's useful in situations where information or data is incomplete, ambiguous, or subjective.



★ Learning

Learning refers to ability of a machine or computer system to acquire knowledge or skills on its own, without explicit programming for every specific task.

Components of learning Systems:

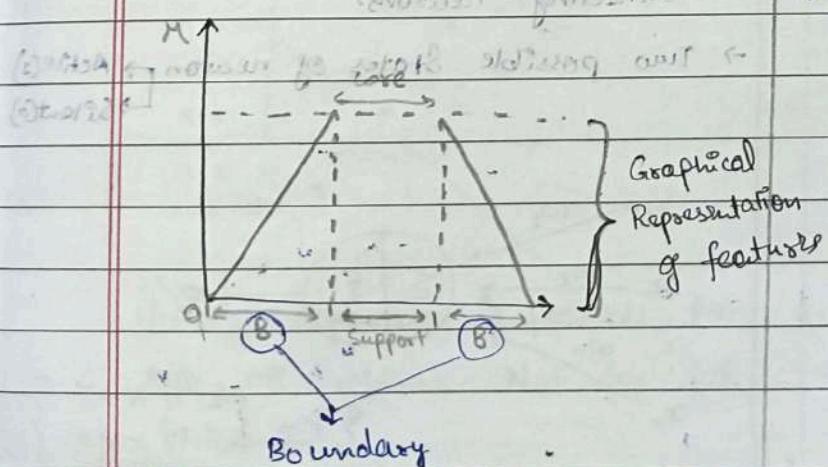
- Learning Components → (main) acquire knowledge, make changes / improvements.
- Performance Components → Task is performed by choosing an action that needs to be taken.
- Problem generator → Suggests Problems / actions that would lead to generation of new examples to improve learning.
- Critic → Gives feedback about performance.

★ Membership function in Fuzzy Systems

- Membership function is used to represent degree of truth in fuzzy logic.
- Represented by Graphical Forms.
- output always lies b/w 0 to 1

Features of membership function:

- i) CORE → If $MF=1$ for maximum truth
- ii) Support → $MF > 0$
- iii) Cross-over → $MF=0.5$ at mid point
- iv) Boundary → $1 > MF > 0$ at outer boundaries

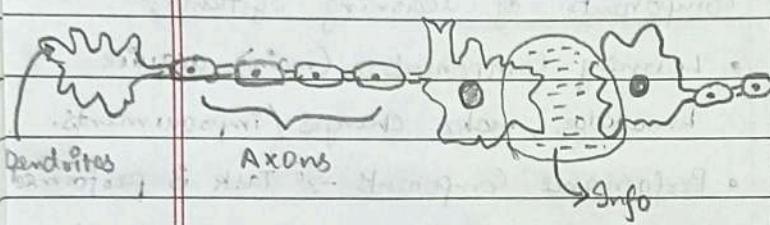


★ Bio-logical model of Neuron (Spiking Neural model)

- Mathematical description of the properties of certain cells in nervous system that generates sharp electrical potentials across their cell membrane.
- Human Brain is composed of Neurone cells called Neurons and they are connected to other thousand cell by Axons.

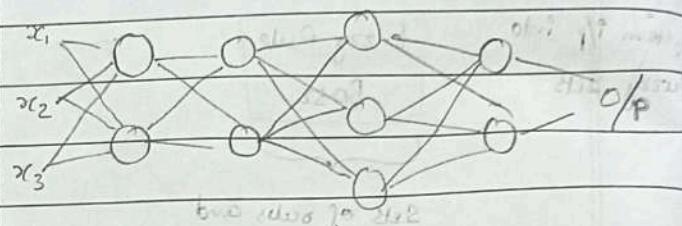
* Bio-logical model of Neuron

→ Stimulus from external environment or input from sensory organ are accepted by Dendrites.



(3) Multi-layer Perception

→ 3 or more layers are used to classify non-linearly separable data.
→ It uses non-linear Activation function



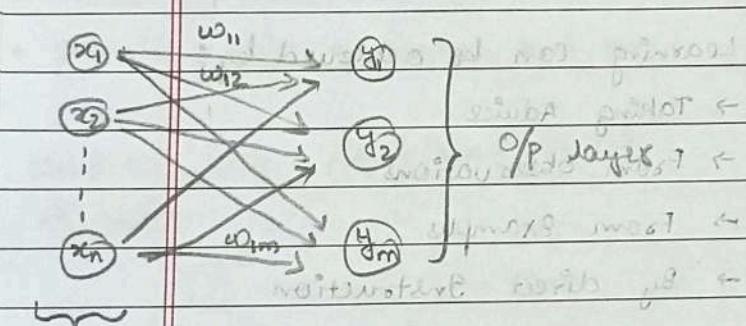
(4) Feedback ANN

feedback is provided to adjust parameters

* Different types of Artificial Neural Network

(1) Single Layer Feed forward network

only two layers → input and output



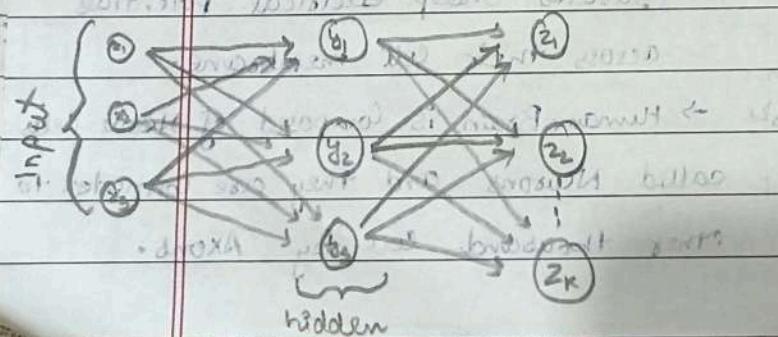
Input layer

→ first layer of the network

(2) Multi-layer feed forward network

→ has a hidden layer in between them

→ computationally more stronger



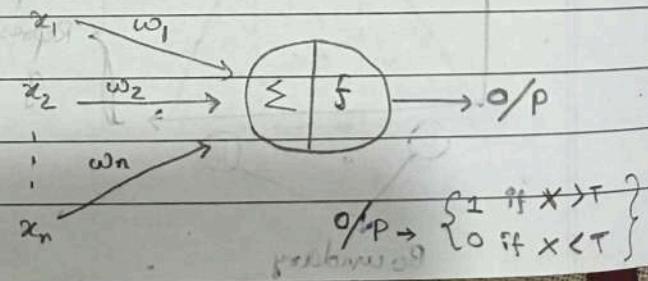
* MCP Neuron Model - (McCulloch Pitts)

→ first mathematical model of biological neuron.

→ Basic building block of neural network

→ Directed weighted Graph is used for connecting neurons.

→ Two possible States of neuron → Active(1) or Silent(0)

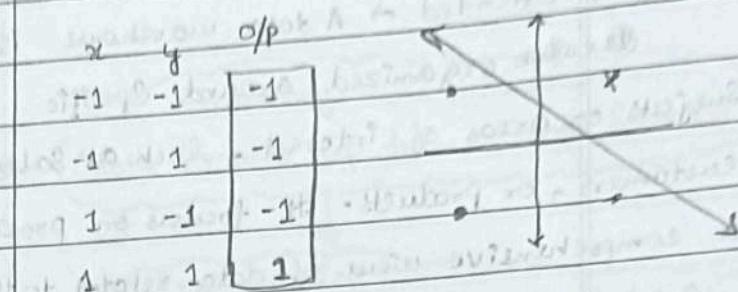


* Bias/ Threshold

- Minimum value of weighted active i/p for a neuron to fire
- If effective i/p is larger than T_g , then O/p $\rightarrow 1$ else 0

AND Problem:

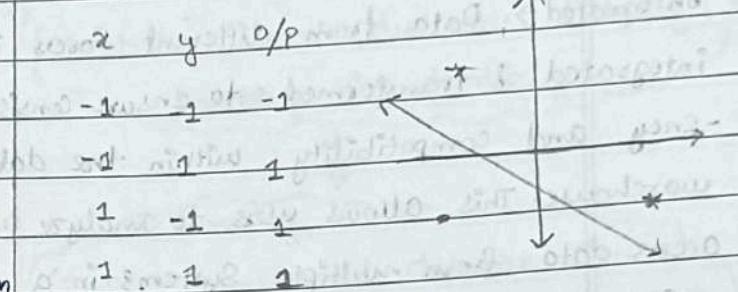
x	y	O/p
-1	-1	-1
-1	1	-1
1	-1	-1
1	1	1



For example → In a binary classification problem where the goal is to predict whether an email is Spam or not, a model may assign a probability score to each mail. By setting a threshold such as 0.5, if the predicted probability of an email being Spam exceeds the threshold, it's classified as Spam; otherwise, it's classified as not Spam.

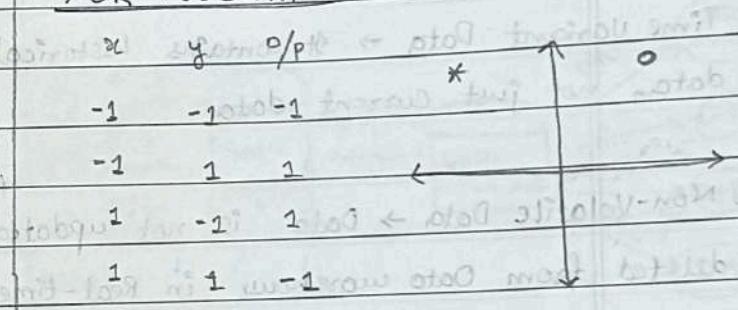
OR Problem

x	y	O/p
-1	-1	-1
-1	1	1
1	-1	1
1	1	1



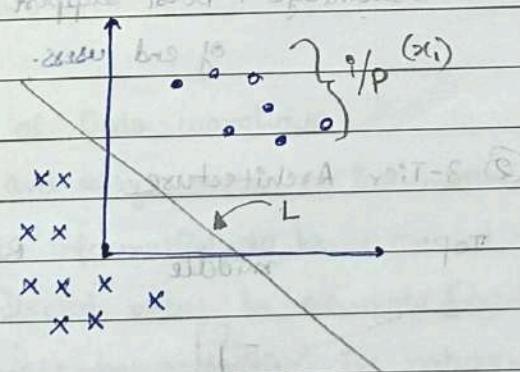
XOR Problem

x	y	O/p
-1	-1	-1
-1	1	1
1	-1	1
1	1	1



* Linearly Separable Patterns

When two classes of patterns that can be separated by decision boundary represented by linear, then they are said to be linearly separable.



Linearly Separable is possible for AND, OR Problems. Not for XOR Problems.

* Data warehouse

In Data warehouse, data is extracted, transformed, and loaded from multiple operational systems, such as database, spreadsheets, or transactional system. The data is then organized, cleaned and standardized to ensure consistency and quality. This process is commonly known as data integration.

Key characteristics:

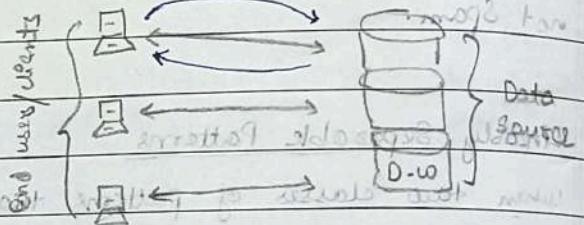
- ① Subject-Oriented → A data warehouse is organized around specific subjects or areas of interest, such as sales, customers, or products. It focuses on providing a comprehensive view of data related to those subjects.
- ② Integrated → Data from different sources is integrated & transformed to ensure consistency and compatibility within the data warehouse. This allows users to analyze and access data from multiple systems in a unified and standardized format.

③ Time Variant Data → It contains historical data, not just current data.

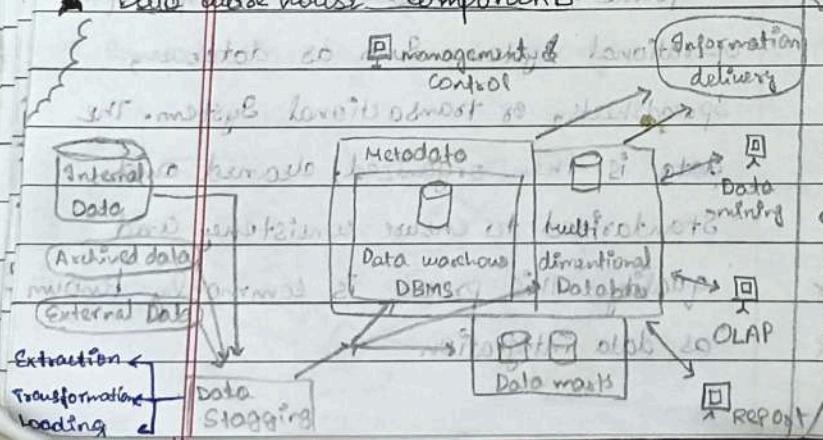
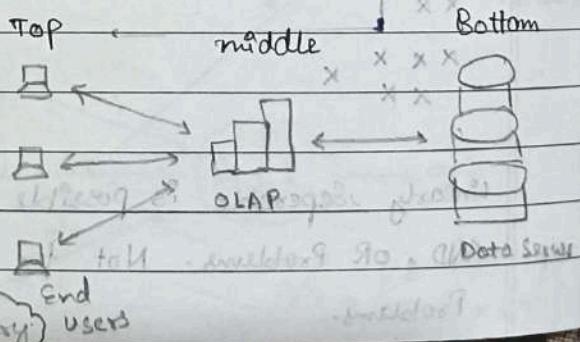
④ Non-Volatile Data → Data is not updated/deleted from Data warehouse in Real-time

⑤ Data Granularity → In Data warehouse, it is efficient to keep data summarized at different levels.

- Management and Control → Coordinates services & activities within Data warehouse
- OLAP → It's a technology and approach used for analyzing and querying multidimensional data from a data warehouse or other data sources. It allows users to perform complex and interactive analysis of data to gain insights and make informed decisions.

★ Different type of Data warehouses**① 2-Tier Architecture**

- Easy to maintain
- Fast communication
- Disadvantage → Doesn't support large no. of end users.

★ Data warehouse Components**② 3-Tier Architecture**

→ It's most widely used

→ Bottom tier → After cleaning, transformation data is loaded

→ Middle tier → OLAP Server

→ Top tier → Front end client layers

③ 4-TIER Architecture

→ end user

→ OLAP

→ Presentation (New layer)

also known as the user interface or client tier, is responsible for interacting with users and presenting information to them. It includes components such as web browsers, mobile apps, or desktop application that users directly interact with.

★ ETL Process → (Extraction, Transform & Load)

Extracting

(Data brought up from external sources)

Transforming

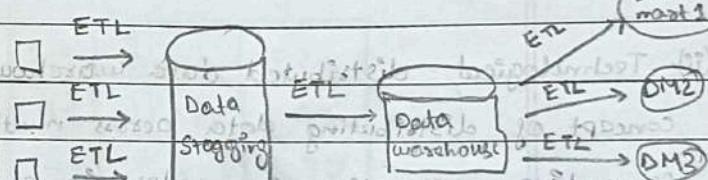
(Converting the extracted data into Standard format)

(Conversion, cleaning duplicates, filtering, binning, standardizing)

Loading

(Putting transformed database into a

large database or data warehouse.)



★ Need for Data warehouse

→ For taking quick & effective decisions

→ Helps users in providing necessary inform.

→ Creating and managing efficient data repository.

Data

Sources

Data warehouse

→ It supports pre-defined operations. (Retrieving of data, insert, delete & update)

OLTP

→ It only supports predefined operations.

★ Goals of Data warehouse

→ Secured and easy access of information to user

→ Consistent information to be provided

→ Data collected must be accurate, verified.

→ Data must be adaptive in nature.

→ Query accesses thousand → few records at a time

or million of records

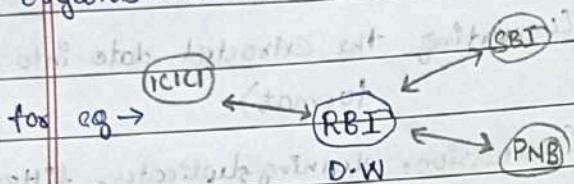
→ Stores historical data → No historical data.

★ Disadvantage → Construction of D-W for large organization is complex task & can take years to complete.

- slow processing speed

★ Distributed data warehouse

In this the data are shared across multiple data repositories, for the purpose of OLAP (Online Analytics) and where each data warehouse may belong to one or more organisation.



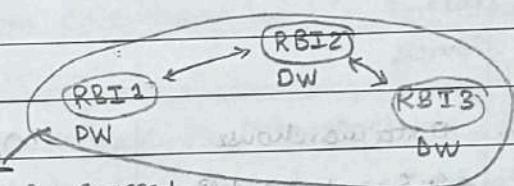
Categories of distributed data warehouse:

(i) Local and Global data warehouse

↓
data of a particular site (SBI, PNB data)
↓
Integrate data (RBI data)

(ii) Technological distributed data warehouse:

concept of distributing data across multiple computing resources or nodes in a network.



If we summarize a report, then all the 3 gonna be treated as 1 entity

(iii) Independently evolving distributed

data warehouse

here multiple data warehouse or data marts evolve independently over time and they have its own schema, data model, data sources and business rules.

★ Data mart

Data mart is a subset of the DW that is usually oriented to a specific purpose.

Reasons for creating Data mart

- Easy access of frequent data
- Improved end-user response time
- Less cost

Different types of Data mart

• Dependent Data mart → In this the Data mart is built by drawing Data from Central Data warehouse that already exists.

• Independent Data mart → Data mart is created without help of Data warehouse

Disadvantage:

- unorganized development
- increase in Data mart size leads to problems such as performance degradation, data inconsistency.

Difference b/w data warehouse & data mart

- Data warehouse stores detailed data, allowing complex analysis across multiple dimensions, while data marts often store summarized and aggregated data optimized for specific business areas.

→ Building and maintaining our data warehouse is complex and time-consuming, involving data while data marts are quicker and easier to implement due to their narrow scope and simple requirements.

★ 5 Steps in data mart:

- ① Designing →
 - understanding the department's requirement
 - Integrating and Summarizing data
 - ensure data quality
 - providing user-friendly access
 - ...

★ Dimensional modelling

It uses a multidimensional schema having measure (1,2,3) and dimension (a,b,c) attributes. Facts are the numerical measures or quantities by which one can analyze relationship b/w dimensions. The relations containing such multidimensional data are called Fact tables.

Bid	tid	number
B1	1	25
B2	2	26

- ② Dimensions → A Dimension table is a table associated with logically related attributes

Bid	Author	Price
B1	ABC	40
B2	XYZ	70

③ Population → (Loading the data)

- getting data from the Source (ETL)
- cleaning and transforming the data.

★ Data warehouse Schemas

① STAR SCHEMA

→ It consists of Fact table with a single table for each dimension

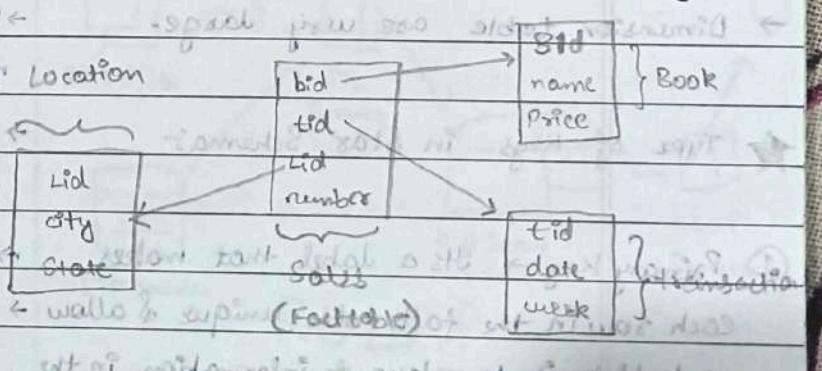
④ Accessing →

- involves putting the data to use.

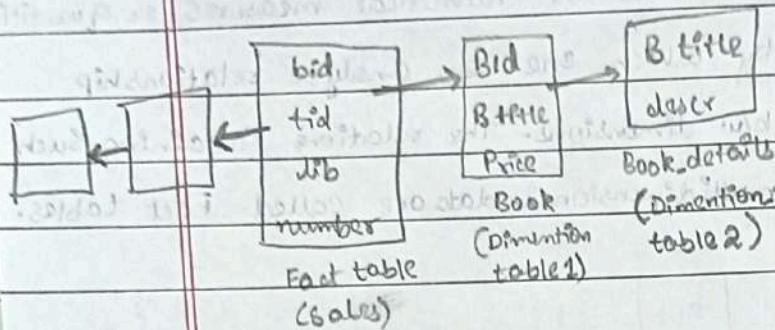
→ simplest & easy algorithm for sub queries
→ most suitable for query processing

⑤ Managing →

- involves managing the data over its lifetime
- Secure the data
- Managing the growth of data

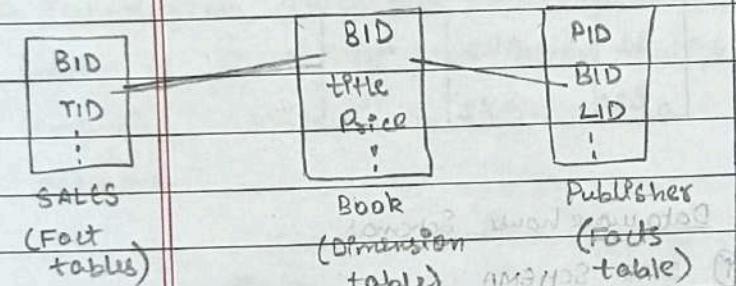


ii) Snowflake Schema → It's an extension of the Star Schema, which organizes data into a central fact table connected to multiple dimension tables.



- Dimension tables are easier to update.
- Disadvantage → Complex Schema.

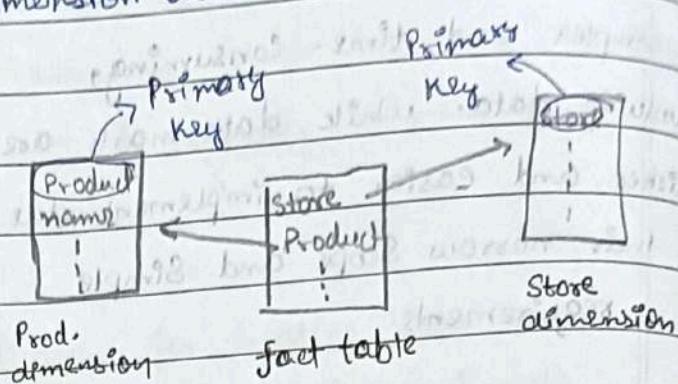
iii) Fact constellation Schema → Also known as Galaxy Schema. In this, multiple Fact tables Share the Dimension tables.



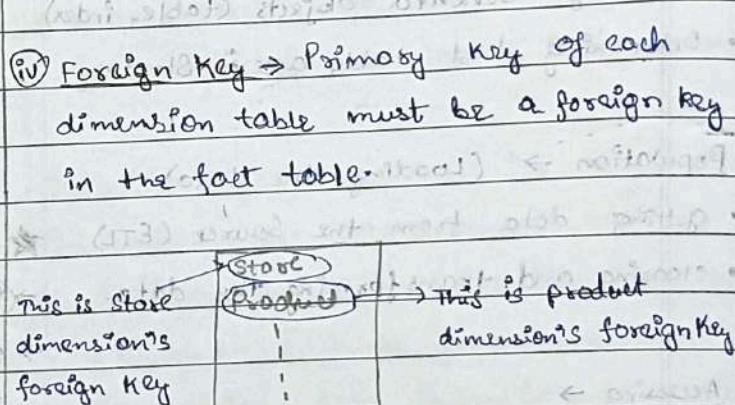
Disadvantage

- complex due to multiple fact tables.
- difficult to manage
- Dimension table are very large.

dimension table.



ii) Surrogate Keys → They are created by the system and have no inherent meaning or relationship to the data they represent. Surrogate Keys help ensure that each record in a database table has a distinct identifier, making it easier to handle and organize the data. For Eg → MUMBAI1123 DELHI123



Multidimensional Data modelling → Data is stored in multidimensional database → RDBMS (Relational database)

Relational Data modelling → Storage of data in multidimensional cubes → 2D Tables

→ Data's in denormalized form → normalized form
→ Mainly used for OLAP (Analytical Purposes) → OLTP (Transactional Purposes)
→ MDX language is used → SQL lang is used

Type of Keys in Star Schema:

① Primary Key → It's a label that makes each row in the fact table unique & allow us to link it to relevant information in the

• OLAP (Online Analytical Processing)

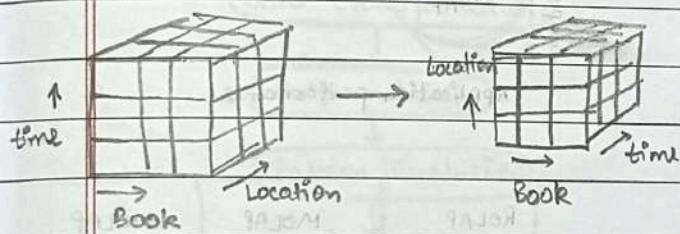
It's a category of software technology that enables analysts managers to analyze the complex data derived from the Data warehouse.

• OLTP (Online Transaction Processing)

It refers to a class of Systems and applications designed to manage and process real-time transactions. (Insert, update, delete)

★ Different OLAP operations

- Pivoting → It's the technique of changing from one-dimensional orientation to another. Also called rotation.



- Slice and Dice → They help uncover patterns, trends, and relationship within the data and enable users to make informed decisions based on the insights gained.

→ Slicing is like taking a slice out of the larger dataset to examine and analyze a particular portion of it.

→ Dicing is like cutting the data into different slices and examining each slice separately.

• Rollup and Drill Down

If we come to confined data from large data then its called Rollup.

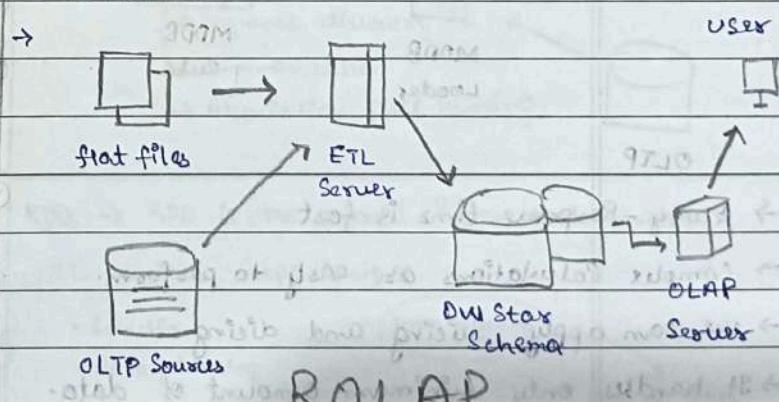
If we go from confine data to large data then its called Drill Down.

★ OLAP Rules By CODD Person

- Multidimensional conceptual view
- Transparency
- Accessibility
- Consistency
- client/server
- Generic dimensionality
- Dynamic Sparse matrix
- Multi user
- unrestricted
- Flexible
- Unlimited

★ ROLAP (Relational OLAP)

- It work with relational database whereas OLAP used to work with multidimensional database.
- Fact and Dimension Table are stored as Relations.



ROLAP

Advantages of ROLAP

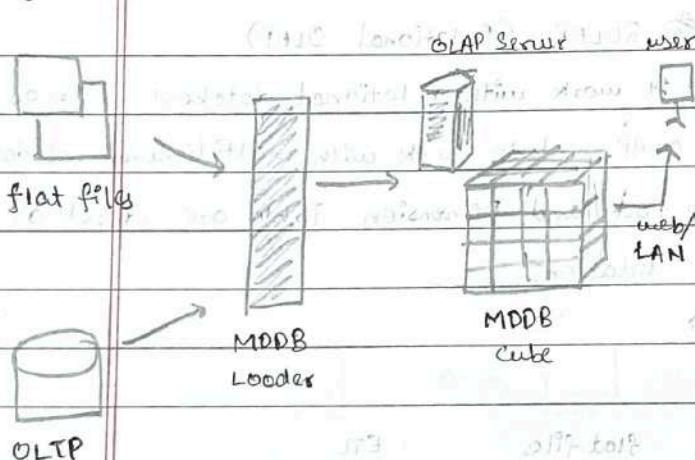
- handle large amount of data
- 2D relational tables can be viewed in multiple multidimensional forms.
- Any SQL reporting tool can access data.
- time needed to load data is less.

Disadvantages of ROLAP

- difficult to perform complex calculation.
- long query time for large data size.

MOLAP (Multidimensional OLAP)

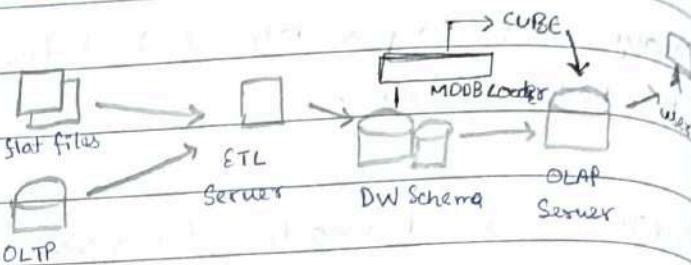
- MOLAP has a 3-tier architecture (OLTP → OLAP → user)
- It's a multidimensional database and storing optimization that helps in better performance
- ! → Array are used. (we can find cell's location using array)
- Compression is also used (so we can store more data in less space)
- Eg → Oracle's Express Server



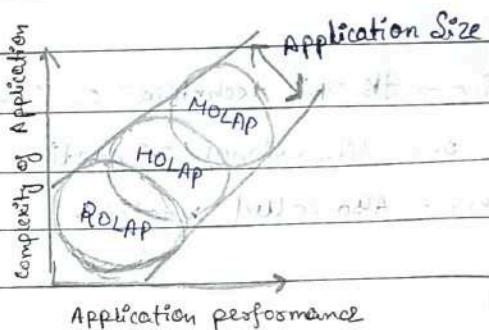
- Query-response time is fast
- Complex calculations are easy to perform.
- we can apply slicing and dicing
- It handles only minimum amount of data.
- Loading the dataset is slow & high complexity.

HOLAP (Hybrid OLAP)

This system includes the best of ROLAP and MOLAP.



→ The complex queries problem that we had in ROLAP can be solved using MDDB Loaders

Comparison of OLAP Servers :-

	ROLAP	MOLAP	HOLAP
i) Detail data storage location	Relational database	MDDB	Relational Database
ii) Aggregate data storage location	Relational database	MDDB Loader	MDDB
iii) Space req.	Large	Medium	Small
iv) Query Response Performance time	Slow	fast	medium
v) Processing time	Slow	fast	fast

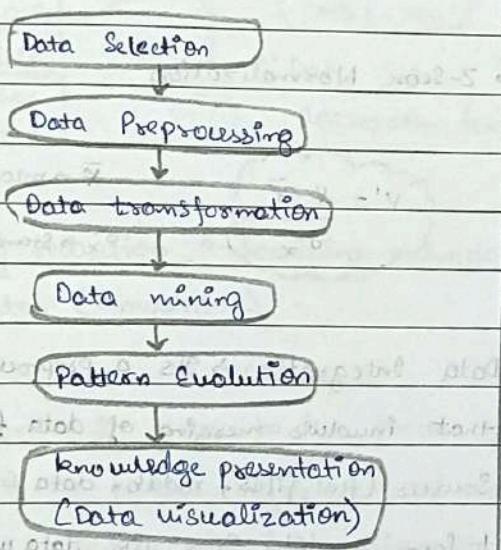
★ Data mining → Process of discovering or mining knowledge from a large amount of data.

→ Another term for Data mining → KDD
(Knowledge Discovering from data)

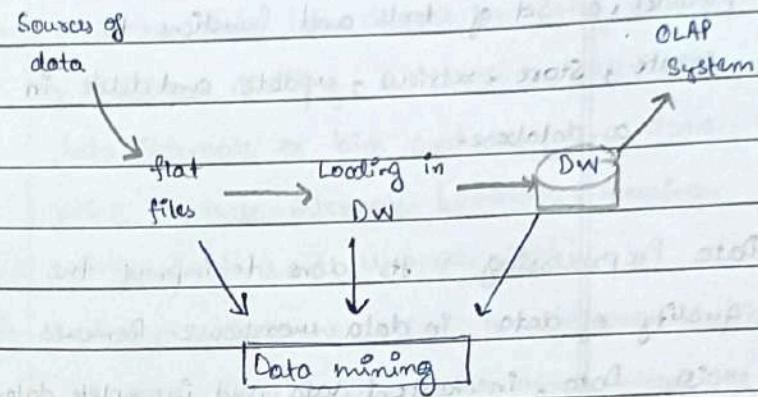
→ Various Data mining technologies are the following:

- ① Statistics
- ② Artificial Intelligence
- ③ Machine Learning

★ Steps in Data mining Process / Phases of KDD in Database:



★ Data mining in the Datawarehouse Environment



★ Data mining applications:

- ① Customer Segmentation → (to understand customers)
- ② Fraud Detection
- ③ Demand Prediction

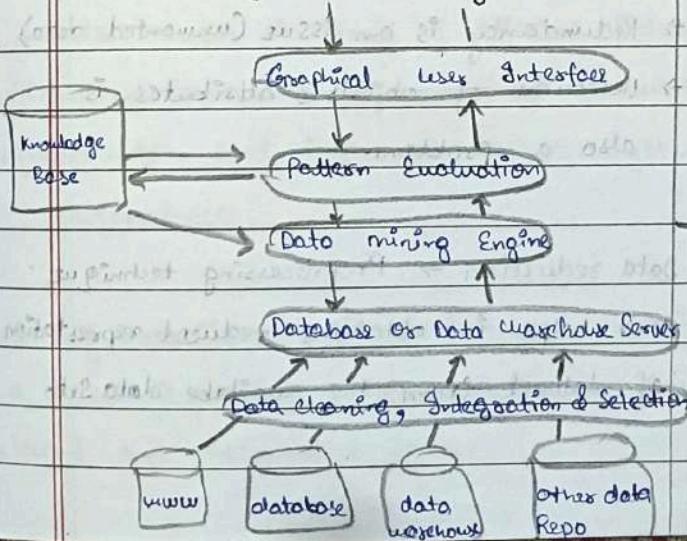
Benefits of Data mining Tasks:

- ① PREDICTIVE → Predict values of Data by making use of known results from a different set of sample data.
 - ↳ Classification
 - ↳ Regression
 - ↳ Prediction
 - ↳ Time Series Analysis

- ② DESCRIPTIVE → Enables you to determine patterns & relationship in a sample data.

- ↳ Clustering
- ↳ Sequence discovery
- ↳ Summarization
- ↳ Association rule mining

★ Architecture of Data mining



KDD → KDD is the overall process that involves

various steps such as

- Selection
- data cleaning
- Preprocessing
- data transformation
- data mining

one of the step in KDD

DBMS → It's a software that helps manage and organize data in a structured way. It provides a set of tools and functions to create, store, retrieve, update and delete in a database.

Data Preprocessing → It's done to improve the quality of data in data warehouse. Removes noisy data, inconsistent data and incomplete data. (Data with missing values)

Data Cleaning → It ~~does~~ cleans the data by filling in the missing values, smoothing noisy data, resolving the inconsistency and removing the outliers.

ways to handle missing data during cleaning

- Manual entry of missing data
- Using attribute mean
- Using most probable value (Predicting the value)
- Using Global constant (NA)
- Ignore the tuple or row

Data transformation → It's a data preprocessing technique that transforms or coordinate the data into alternate forms appropriate for mining. It involves:

- i) Smoothing → removing noisy data
- ii) Aggregation → constructing a data cube (OLAP)
- iii) Generalization → low-level concepts are replaced with high level

⑩ Normalization → Attributes values are normalized by scaling their value so that they fall in specified range. For eg →

40
36
47

← we have this table of numbers. We have converted all the numbers into range (0-1) that's normalization.

• Min-Max Normalization

$$v' = \frac{v - \min_x}{\max_x - \min_x}$$

v' → new value
 v → original value
 \min_x → minimum value of the attribute
 \max_x → maximum value of the attribute

• Z-Score Normalization

$$v' = \frac{v - \bar{x}}{\sigma_x}$$

\bar{x} → mean of attribute(x)
 σ_x → standard deviation

Data Integration → It's a preprocessing method that involves merging of data from different sources (flat files, mddbs, data cubes) in order to form a data store like data warehouse.

- Redundancy is an issue (unwanted data)
- Detection of objects & attributes is also a problem.

Data reduction → Preprocessing technique that helps in obtaining reduced representation of dataset from the available data set.

- Splitting → dividing the attributes in topdown approach
- Merging → dividing the attributes in bottom up approach
- Supervised →
- Unsupervised

classmate

Date _____
Page 29

- Integrity of the original data should even after reduction in data volume
- It should produce same analytics result as on original data.

Data Cube Aggregation

It's a process in which information is gathered and expressed in a summary form. (Statistical)

for eg →

Year 2017	
Hy	Sales
H1	500
H2	300

Years	Sales
2017	800
2018	700

Year 2018	
Hy	Sales
H1	600
H2	100

∴ Data cube
aggregation form

Dimensionality reduction → Removing redundant attributes (unwanted)

Data compression → It refers to the process of reducing the size or storage requirements of data without losing significant information.

It provides benefits in terms of efficient storage utilization, faster transfer, reduced bandwidth requirements, and improved overall system performance.

Numerosity reduction → No. of features or attributes of a dataset is reduced while preserving the most relevant and informative ones.

(Discretisation) and Concept hierarchy

Discretization is the process of transforming continuous data into discrete or categorical form. It involves dividing a continuous attribute into intervals or bins and assigning data points to these intervals. [Splitting, Merging, Supervised and Unsupervised]

→ Concept hierarchy refers to the organization of data in a hierarchical manner, where the data is structured into levels or categories of increasing detail.

→ Think of a concept hierarchy like a tree a tree structure. At the top level, you have a broad category or concept, and as you move down the hierarchy, you encounter more specific subcategories or concepts that provide more detailed information.

[Binning, Histogram analysis, cluster analysis]

→ for eg → consider a concept hierarchy for the attribute 'Time'. At the top level, you may have the concept of 'Year', followed by 'Quarter', 'Month', 'Week' and finally 'Day' at the lowest level. Each level represents a different level of granularity or detail.

Dimensionality reduction → It represents the original data in compressed or reduced form by applying data encoding or transformation.

→ Lossless → If original data can be reconstructed from compressed data without losing any information.

→ Lossy → If original data can be reconstructed from compressed data without losing of information.

Clustering algorithm group similar data points together based on their inherent characteristics or proximity in the feature space.

Numerosity reduction → It reduces the data volume by choosing alternative smaller forms of data representation.

② Dimensionality reduction → (done)

↳ parametric → dimensionality reduction that involves the use of statistical models or determine the most relevant feature in a dataset. (Regression, log-linear models)

Semi-Supervised learning → A blend of supervised and unsupervised learning.

→ Non-parametric → data is stored in the form of histogram, clustering, Sampling.

A small portion of the data is labelled,

Supervised → Input data is provided to the model along with the output target variables is given. The goal is to train the model so that it can predict the output when its given new data.

Similar to supervised learning, where each labelled example consists of input data and corresponding output label.

However, the majority of the data remains unlabeled. The goal is to utilize the labeled data to guide the learning process & leverage the unlabeled data to extract additional information and improve generalization for e.g. medical prediction.

1) Classification → It predict / classifies the discrete values such as, Male / Female True / False etc.

Reinforcement learning → Focus on making decisions based on previous experience.

2) Regression → It predict continuous values such as Price, Salary, etc.

Noise in Data → It's a random error or

Variance in a measured variable.

Unsupervised learning → Here only the input data is provided to the model and it does not rely on predefined labels or target outputs for training.

There are three techniques to remove noise:

① Smoothing by BIN means:

Value of the bin is replaced by mean value (average).

(ii) Bin medians \rightarrow value of bin is replaced by medians value

(iii) Bin boundaries \rightarrow value of bin is replaced by minimum and maximum value

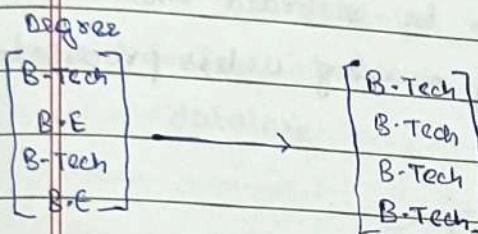
Steps of Data Cleaning Process

B Detection of discrepancy

\hookrightarrow manual error

\hookrightarrow inconsistent data.

(ii) Transformation of data. for eg \rightarrow



Q) Divide the given Sample Data in two (2) clusters using K-Means Algorithm

	Height (H)	Weight (W)	Centroid
1	185	72	(185, 72)
2	170	56	(170, 56)
3	168	60	
4	179	68	
5	182	72	
6	188	77	
7	180	71	
8	180	70	
9	183	84	
10	180	88	
11	180	87	
12	177	76	

E-D for row 3 \rightarrow

$$\sqrt{(168-185)^2 + (60-72)^2} = 20.80$$

$$\sqrt{(168-170)^2 + (60-56)^2} = 4.48$$

Since this is smaller, 3rd row will be in 2nd cluster

as cluster 1 and cluster 2

\rightarrow we will consider first two rows and then solve for the third row to check whether its 1st or 2nd row.

* Memory-Based Reasoning \rightarrow It uses known instances of model to predict unknown instances. When a new record arrives for evaluation, the algo finds neighbours similar to the new record, then uses characteristics of the neighbours for prediction and classification

Link Analysis \rightarrow useful for finding patterns from relationships. Applications of Link analysis are:

(i) Association Discovery \rightarrow Algo finds combination where presence of one item suggests the presence of another.

for eg \rightarrow There are chances that a customer might buy Butter when buying bread.

So there's a link b/w Butter & Bread.

(ii) Sequential Pattern Discovery \rightarrow Discovers patterns where one set of items follows another specific set.

(iii) Similar Time Sequence Discovery \rightarrow It refers to the process of identifying

patterns, trends, or relationships in a sequence of events, or data points that occur over time.

★ Association rule mining → By this, analysts can discover interesting relationships b/w items or attributes within transactional data. It has applications like market basket analysis, recommendation systems.

Parameters →

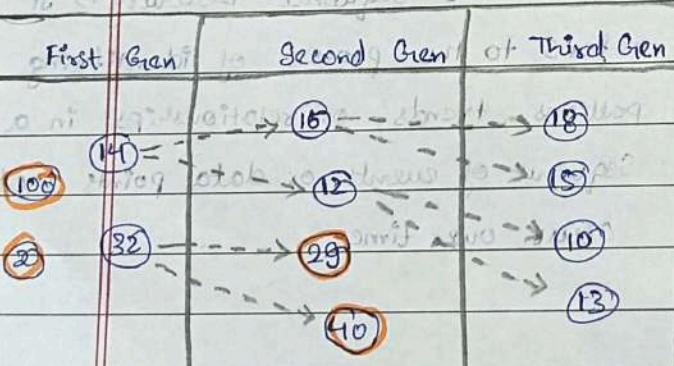
- i) Finding all items that appears frequently in transaction → min. Support Count
- ii) Finding strong associations among frequent items → confidence

But finding strong associations among frequent items.

Genetic Algorithm

It selects the best attributes in every iteration, then it will crossover those attributes and will make few changes with the help of mutation in the order to improve next successive generations.

Suppose we have 4 different attributes & we are finding the most optimised solution



0 → not Selected

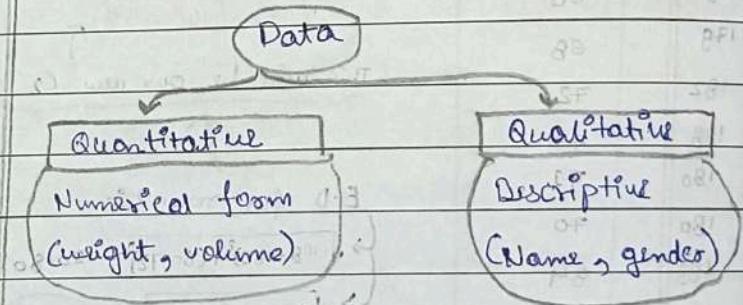
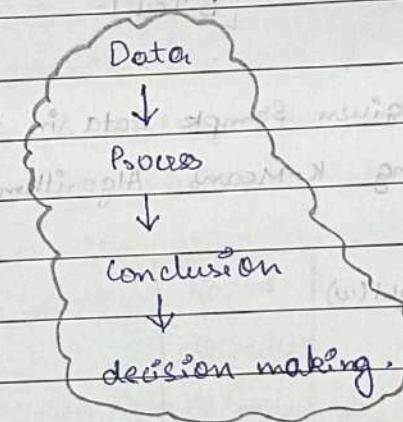
★ DBMS

what is data?

→ Data is a collection of raw, unorganized facts and details like text, observations, figures, symbols and descriptions of things etc.

→ Data is measured in terms of bits and bytes which are basic unit of information in the context of computer storage and processing.

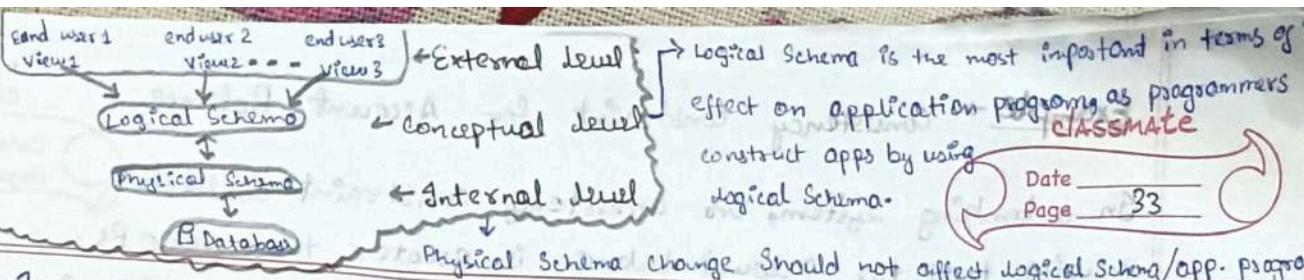
→ Data can be recorded and does not have any meaning unless processed.



what is information?

→ It's processed, organized and structured data

→ It provides context of the data and enable decision making.



Logical Schema is the most important in terms of effect on application programs as programmers construct apps by using logical Schema.

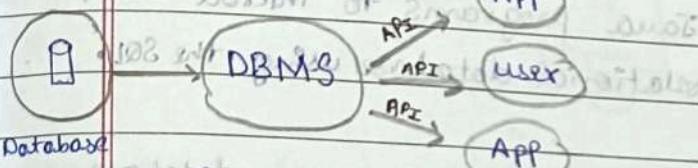
Date _____
Page 33

Database → It's a system where data is stored in a way that it can be easily accessed, managed and updated.

DBMS → It's a collection of interrelated data and a set of programs to access those data.

→ DBMS is the database itself, along with all the software and functionality. It's used to perform different operations, like addition, access, updating and deletion of data.

→ It's an interface between users and the database.



DBMS vs file systems

DBMS and file Systems are both used for managing data, but they have significant differences in terms of structure, functionality and usage.

Here are disadvantages of file Systems:

- Data redundancy and inconsistency
- Difficulty in accessing data
- Data Isolation
- Integrity problems
- * (These are advantages of DBMS)

Schema and different levels in Schema

A Schema refers to the structure or blueprint of a database. It defines how the data is organized, the relationships between different data elements, and the rules or constraints that govern the data.

Different levels in Schema: (abstraction)

i) Physical / Internal Schema →

It describes how the data is stored.

It deals with low-level details, including data structures, file organization.

ii) Conceptual / Logical Schema →

It serves as a high-level abstraction that focuses on the essential aspect of the database design, providing a global perspective of the entire database system.

It represents the overall logical structure and organization of the entire database.

iii) External / View Schema

→ The external Schema represents the view of the database from the perspective of individual users or groups of users.

→ These Schemas provide a simplified and customized view of the database by selecting specific tables, attributes, or views that are relevant to a particular user's needs.

Example → Consistency constraint for Account Balance

classmate
Date _____
Page _____

In a banking system, one consistency constraint could be ensuring that the account balance is greater than 1000 Rs.

Abstraction in DBMS

It refers to the process of hiding complex details and presenting users with a simplified and intuitive view of the database. It allows users to interact with the database without needing to understand the underlying complexities of how the data is stored and managed.

Query language → SQL (Structured

Query language) is the most widely used query language, providing a standardized way to interact with relational database.

SQL combines elements of DDL and DML to define database structure and manipulate data.

Difference b/w Instances and Schemas

→ Collection of information stored in the DB at a particular moment is called an instance of DB.

Interface b/w Java, C and SQL

The interface between Java and SQL is typically facilitated through the use of a Java Database Connectivity.

→ Schema doesn't change frequently. Data may change frequently.

(JDBC) • It provides a standard set of classes and interfaces that allow Java programs to interact with relational database using the SQL.

Data models

It describes the Design of a database at logical level. It describes data, data relationships, data semantics and consistency constraints.

Similarly for C → open database connectivity (ODBC)

for eg → ER model

→ DBA → Database administrator.

- Relationships model
- object-oriented model
- object-relational data model etc.

It's a professional responsible for the design, implementation, maintenance, and management of a database system.

Database language

→ DDL (Data definition language)

→ DBA works at logical level whereas users work at view level. DBA works at internal level as well but mostly at logical level.

include CREATE, ALTER, and DROP

→ DBAs have responsibilities that

→ DML (Data manipulation Language)

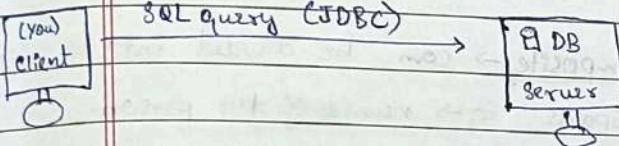
Example of DML statements include

SELECT (to retrieve data), INSERT (to insert new records), UPDATE (to modify existing records), and

include designing and setting up databases, ensuring data security, monitoring performance, resolving issues, maintaining backups, upgrading software, and providing support to users. They play a vital role in managing and maintaining the database environment, ensuring its reliability and efficient operation.

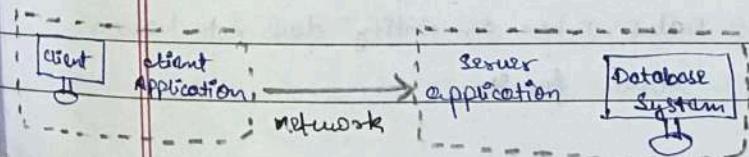
Data Application Architecture

① T1 Architecture → The client, server and DB all present on the same machine.



→ Here, App is partitioned into 2-components.
→ API Standards like ODBC and JDBC are used to interact between client & server.

Here, Client application will communicate to Server application, then Server application will call Database System and retrieve the data. Then Database System will return the data to Server application. The Server application modifies the data and sends it back to Client application.



- T3 architecture are best for WWW applications
- Scalability due to distributed application servers.
- Data integrity, APP server acts as a middle layer between client and DB, which minimize the chances of data corruption.
- Security, client can't directly access DB, hence it's more secure.

T2 Architecture

Difference b/w strong and weak entity-
→ we have this topic in ER model (entity-relationship model)

→ An Entity is a "thing" or "object" in the real world that is distinguishable from all objects.

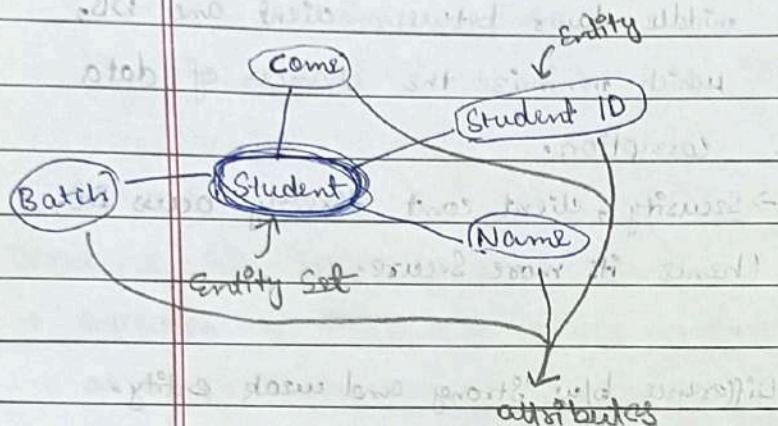
→ The Student entity is a strong entity because it can exist independently and has its own attributes. It has unique identifier (Student ID) that distinguishes it from other Students. On the other hand, the course enrollment entity is a weak entity because it relies on the existence of both the Student and Course entities. It cannot exist on its own and requires foreign keys to link it to the associated Student and Course.

→ To represent the relationship between the strong and weak entities, the database schema would typically include foreign keys in the weak entity table to establish the connection with the strong entities. This ensures data integrity and

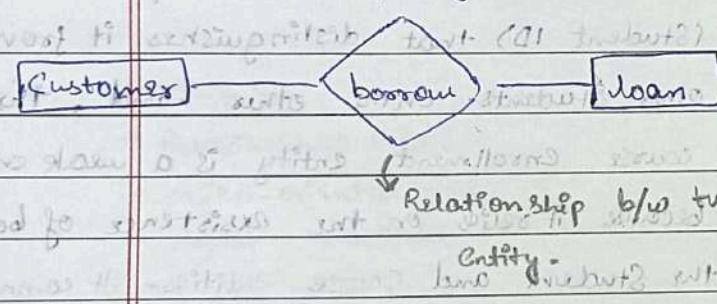
Maintaining the relationship between the entities in the database.

ensures that only integer values can be stored in the attribute.

Entity Set



ER model → using er model, developers can create a clear and structured representation of their data and its relationships, making it easier to design and implement a database system that accurately represents the real-world entities and their connections. for eg →



Attribute consistency constraints → It ensures that values stored in an attribute are there to maintain specific rules or conditions. These constraints help maintain data integrity and ensure that the attribute values meet certain criteria. for eg →

→ Data type constraint → If an attribute is defined as an integer type, the constraint

→ Range constraint → It defines the minimum and maximum values that an attribute can take.

→ Not Null Constraint → It ensures that an attribute cannot have a null or empty value.

Types of Attribute

① Simple → Attributes which can't be divided further. Eg → customer's bank number

② Composite → can be divided into subparts. Eg → name of the person

③ Single-valued → only one value attribute. for eg → Student ID

④ Multi-valued → Attributes having more than one value. for eg → Phone no.

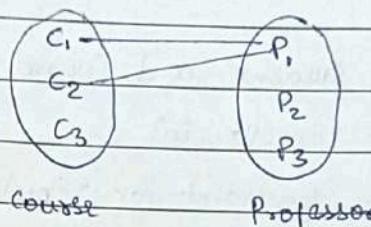
⑤ Defined → value of this type of attribute can be derived from the value of other related attributes. for eg → age of a person can be derived from his D.O.B

⑥ Null value → An attribute takes a null value when an entity does not have a value for it.

Degree of Relationship

- 1) Number of entities participating in a relationship
- 2) Unary → Only one entity participates eg → Employee manages employee.
- 3) Binary → two entities participates eg → Student takes course. (Binary are common)
- 4) Ternary relationship, three entities participates eg → Employee works-on branch, employee works on job.

- Many to One → Entity in A associated with at most one entity in B while entity in B can be associated with N entity in A - foreg → Course taken by professor.



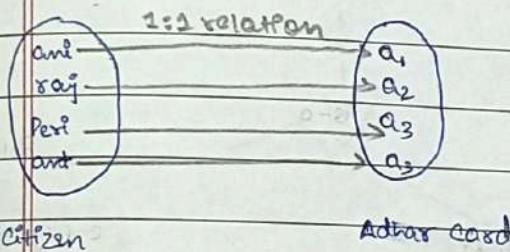
- Many to many → for Eg → Customers buys Product

Mapping Cardinality → Number of entities to which another entity can be associated via a relationship.

Participation Constraints

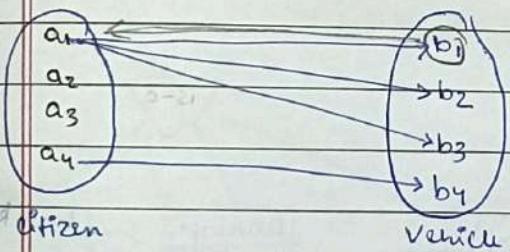
- ① Total participation → Total participation means that every instance of the entity on one side of the relationship must participate in the relationship. It indicates that the participation is mandatory, and without it, the relationship cannot exist.

One to one



- ② Partial participation → It means that not every instance of the entity on one side of the relationship needs to participate. It allows for optional participation, indicating that some instances may choose to participate while others may not.

One to many



→ A person can have many vehicles, but a vehicle can only be owned by single person

→ definition =

Entity in A associated with N entity in B. While entity in B is associated with at most one entity in A.

Eg 1. entities → Employee & Department

- relationship → works in
- If 'works in' is total = It means every employee must work in a department. An employee cannot exist without being associated with a department

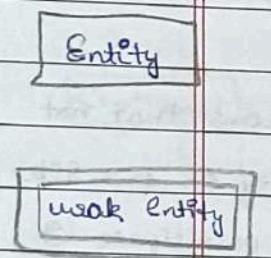
- If participation constraint is partial → It allows for some employees to not be assigned to any department, giving them the option of not participating in the relationship.

(Ex) → entities → Student and Course
relationship → Enrolls in

- If participation constraint for 'Enrolls in' relationship is total → It means every student must be enrolled in at least one course. without enrolling in a course, a student cannot exist in the system.

- If participation constraint is partial → It allows for some students to not be enrolled in any course, giving them the option of not participating in the relationship.

Symbols used in ER Diagram



Attribute

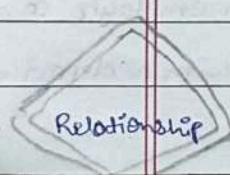
Relationship

Multi-valued Attribute

Primary Key Attribute

Weak Key Attribute

Derived Attribute



Total Participation

Entity

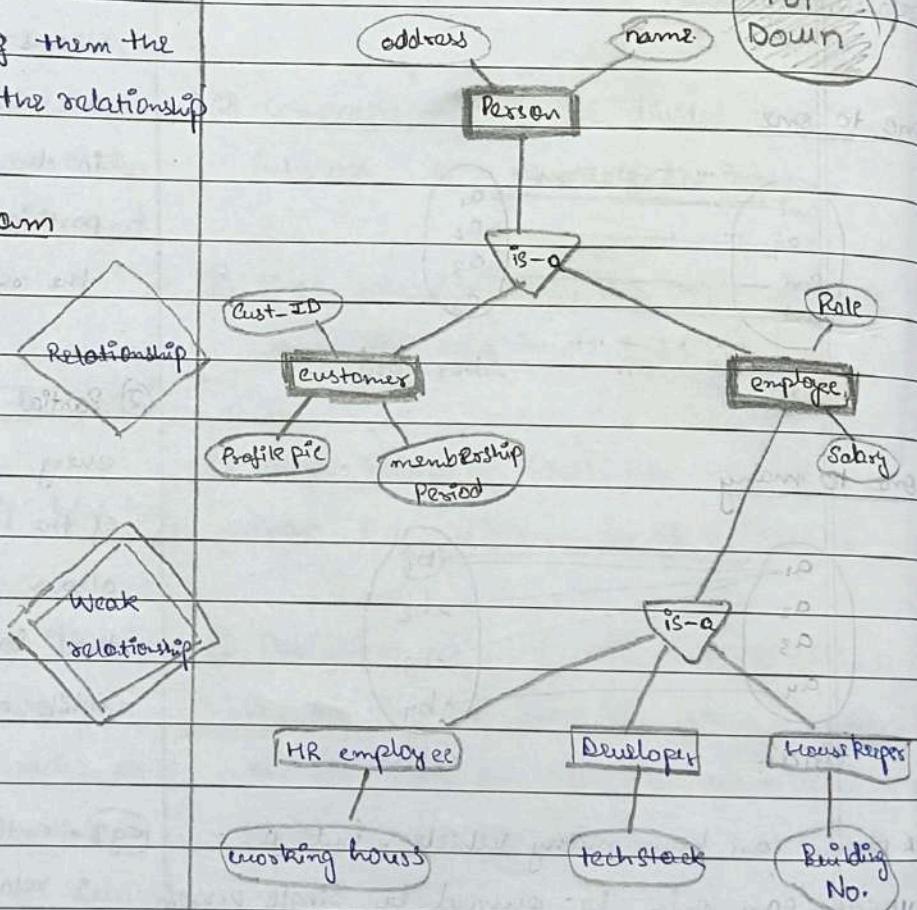
★ Specialisation → Its splitting up the entity set into further sub entity sets on the basis of their functionalities, Specialities and features.

- Its a top-down approach
- we have 'is-a' relationship b/w Superclass and Subclass. Depicted by triangle component.



→ DB designer can show the distinctive features of the sub entities through specialization.

TOP
Down



DBMS is continued after Pg.