

# Artificial Intelligence

03/01/2025

- Imparting normal intelligence into a machine is called Artificial Intelligence.
- Constraint of time & space → which intelligence is to be applied when and where.
- Thinking power → not possessed by machines.

Eg:  $2001 \times 1999 \rightarrow$  diff. to calculate

$$(2000+1)(2000-1) \leftarrow \text{intelligence, only applicable here; not everywhere.}$$

Brain is divided into 2 hemispheres: —

Left

classical AI.

Right

statistical power  
(primitive; existed before AI)

Book Stuart Russell: Artificial intelligence: A modern approach

Motivations

- AI is related to philosophy, psychology and intuitive thinking.
- AI problems are mostly non-deterministic → it contains uncertainty with approximate solution.
- AI is breadthwise not depth: AI is applicable more on a general level rather than specific cases.
- The specific part is not entirely discarded, and a part of it is taken along with the general part.

Heuristic  $\rightarrow$  guessing power of human beings.  
(not deterministic)

- ↳ Random search.
- ↳ contains both risk and gain
- ↳ more thinking power to obtain better solution.

Eg) Block wall problem

B on top of C on top of A.

Planning to remove B (first as its top is empty).

Reinforcement learning  $\rightarrow$  is not data driven.  
based on function generating reward/punishment  
and transitions states accordingly.

Markov model  $\rightarrow$  Next state depends on previous state.

Reinforcement learning involves one or more actions to move from one state to another. IA and choosing the best one. (it is more than Markov learning).

Agent  $\rightarrow$  Human/Robot

Turing Test

$\rightarrow$  To determine whether a computer can think like a human.

$\rightarrow$  Natural language, knowledge, reasoning power and learning power :- criterion to pass the test.

Expert system :- Human reasoning system through rules which are provided to machine.

Eg: If it is sunny, do not take the umbrella.

• Mycin  $\rightarrow$  expert system used to prescribe antibiotics for diseases.

• Apple had developed AI that defeated human in chess (1997)

### Phases of AI

- First  $\rightarrow$  logic based  $\Rightarrow$  deterministic
- Second  $\rightarrow$  probabilistic based
  - $\hookrightarrow$  involves uncertainty of real world.
- Third  $\rightarrow$  neural nw based
  - $\hookrightarrow$  requires huge amount of data.

### Intro

AI  $\rightarrow$  Simulation of human intelligence processes by developing computing systems that have the ability to ~~not~~ perform tasks like humans

### Humans

- $\rightarrow$  Not always rational - biased / emotion driven
- $\rightarrow$  Doesn't take max. goal always, takes risk sometimes.

### AI

- $\rightarrow$  Always rational (doing the right thing)
- $\rightarrow$  Takes maximum goal (benefit)

LOGIC + Domain knowledge is deterministic.

### Knowledge representation

- Rule
- Graph
- Vector

Also involves relation between the components of knowledge

Robots are sent where there is crisis to human life (mining/war/etc.)

AI → computation, logic, probability

### Search, searching

#### Simple searching

no opponent; just follow a path

#### game searching (used in AI)

multiple no. of opponents (agents) whose actions/movements are known to each and every agent.

### Search

Heuristic search is the fundamental technique of AI.

#### Blind (uninformed)

- All paths need to be explored.
- Info not available e.g. BFS

#### Informed (Heuristic search)

→ Guess what is ahead & that info is used.  
e.g. expand a particular node which is best at that level.

### Knowledge

#### Declarative

→ Based on facts & questions on it

#### Procedural

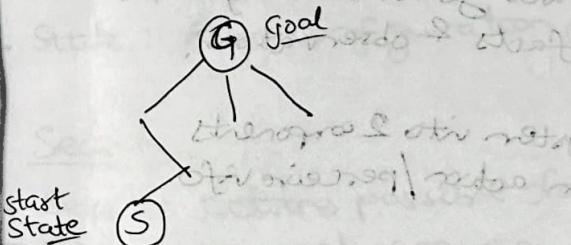
→ Based on 'How'  
→ Needs a fact to reach to sol  $\hookrightarrow$  embedded with declarative.

### Planning

Constructing a sequence of actions that achieves a given set of goals.

e.g. LISP, PROLOG → programming lang;  
goal is given, have to reach to start state.

(Backward search)



### AI planning

Tackles uncertainty in planning that could lead to non achievement of goals.

Predicate → Quantifiable; Relational

Proposition → Yes/ No → not simultaneously

Modifying the environment involves removing obstacles.

### Forms of knowledge

- Empirical
- Scientific
- Intuition
- Authoritative

### Logical Reasoning

- Deduction: From given facts infer another fact
- Abduction: The another fact deduced is possibly true from a given set of facts.
- Induction: Involves generation of rule from facts & observations.

AI divides the system into 2 components

- Agent : perform action / perceive info.
- Environment.

Artificial Narrow Intelligence (weak AI)

Artificial Generative Intelligence (Strong AI)

Artificial Super Intelligence (hypothetical)

Goal-oriented or onlooker intelligent

## Problem Solving

Lec 2

I/P: Set of states  
operator  
start state

( $S_1, S_2, S_3, \dots$ )

START

(SS)

operator

NEXT

(NS<sub>1</sub>) (NS<sub>2</sub>)

each has a cost involved.  
- best next state is chosen accordingly.

### Problem Representation

Broad, generalisation approach to solve a large class of problems.

State space search → to model the problem

State: Any configuration (valid).

### Search

- Several actions possible from SS to NS (many)
- Simulate each sequence of actions to reach the goal
- Calculate the cost of each of these paths from goal to start state.
- Finding such a best sequence is called search. (problem-solver does this)

→ If action is not driven by objective fn, it may fail to reach the goal, like infinite loop.

Eg: Water jugs problem

(x, y)

- Two jugs contain x & y units of water.
- let J<sub>1</sub> allow 3 gallons & J<sub>2</sub> allows 4 gallons at full
- There is no intermediate marking on the jug.

Start state : (0, 0) OR (3, 4)

Suppose goal : (2, y)

Not measurable due to lack of markings.

Let state = (0, 0)

Intermediate states: (3, 0), (0, 4)

Goal state: (2, y)

(0, 0)

actions

(0, 0)

(0, 4)

1) fill

2) empty

3) pour from one to another

(3, 1)

Lab: one such standard problem.

out of which first node is not explored

### Problem Solver

1) Initial State

2) Set of Actions

3) Successor function → explicit / Implicit

4) Goal Test

5) Path Cost

Applying some cost simplifies the process of solving the problem.

### Search

Features:-

Completeness

Time complexity

Space complexity

Optimality

In BFS, time complexity is more

In DFS, space complexity is less, but it is not guaranteed to find the solution.

## Uninformed Search

### problem formulation

- Configuration/state: Representation like  $(x, y)$  for water jug problem.
- Constraints
- Rules: Actions to be performed like fill or empty
- Initial state
- Goal state.
- Implicit state: Intermediate.
- Valid sol $\sqsupseteq$
- General Algorithms  $\rightarrow$  all algos of trees are applicable to search

- Date: 1/1/2025
- Problem of implementation in memory and time.
  - Incompleteness: Sol $\sqsupseteq$  possible so whether it can be achieved  
Eg: If path cost is given as a constraint, the problem gets more difficult.
  - $\Rightarrow$  Change search problem
  - Uncertainty: Explicitly not mentioned, vague ~~constraints~~ representations.  
Eg: Self-driving car (in adverse weather conditions)
  - Try to remove most uncertainty  $\rightarrow$  hence ~~most~~ problems involve approximation.

## Searching

- Search tree: - Possible action sequences starting from initial sequence
- Unexpanded node  $\rightarrow$  node that has not been verified  $\rightarrow$  needs to be expanded.
- Set of actions are applied to current state to generate a new set of states.

### Strategies

- Uninformed
- BFS
- Uniform cost
- DFS
- Depth limited.

game problem  
is irreversible

### Disadv of BFS (Application)

- Redundant paths can't be avoided & actions are reversible (in most cases)
- frontier  $\rightarrow$  ~~the~~ set of leaf nodes available for expansion at any given point.

disadv can be combatted through maintaining a data structure which remembers the explored paths.

### process

Generated  $\rightarrow$  Expanded  $\rightarrow$  Examined  $\rightarrow$  Discarded

### Graph search

Issue: child may have multiple parents  
State-space graph is divided into explored and not-explored.

### In graph

- Black  $\rightarrow$  explored
- White  $\rightarrow$  frontier
- Grey  $\rightarrow$  unexplored

Systematical examination is necessary to avoid exploration issues

Combinatorial explosion :- By checking all the paths it is difficult to find a solution.

Eg:- Travelling Salesman problem with large no. of cities.

### foreach node (attributes)

- state
- parent
- action
- path cost

In uninformed search there is non-variable path cost, in general.

$\rightarrow$  Inclusion of time introduces constraints to the problem.

### Node & state (diff)

- o Eg: Airport  $\rightarrow$  can be a node
- o Physical location  $\rightarrow$  state

Note is a representation

### features of AI (slide 19)

- $\rightarrow$  graph is represented through initial state, actions, transition model (it is infinite)
- $\rightarrow$  complexity is expressed through:
  - b  $\rightarrow$  branching factor: 2 for BFS.
  - d  $\rightarrow$  depth of state space tree.
  - m  $\rightarrow$  max. length of any path.

Shallow  $\rightarrow$  General

Depth  $\rightarrow$  Specific

Explored:  $O(b^{d-1})$

~~unexplored:  $O(b^d)$~~

frontier:  $O(b^d)$

Space complexity:  $O(b^d)$ .

Issue at depth d, takes a long time for debugging / finding out.

> It is exponential search problem so hence it is very rarely able to solve unstructured problem.

> Since we are exploring each layer at a time, it is complete (though memory & space complexities would be high)

> Why is the term  $b(b^d - 1)$  added to time complexity?

Ans) Because time complexity takes worst case sol.

If we find the goal at a depth it needs to be expanded and then needs to found,

so this expansion accounts for the extra term after  $b^d$  in time complexity.

BFS is optimal as we systematically explore the nodes layer by layer.

Hence  $SOL$  is always obtained & it is a complete problem.

> BFS does not involve any cost.

### uniform cost search

- Nodes closer to the goal are assigned lower cost.
- Find path where cumulative sum of costs is least.
- No. of steps is not considered, but the total cost is adjudged.

### properties

→ Infinite loop may be created if there exists path(s) with zero cost, i.e., equivalent to  $No\ Op$ . → some min cost must hence be provided. This guarantees completeness.

→ Search is guided by cost rather than depth ( $b^d$ ).

$$O(b^{1 + \lceil C/e \rceil})$$

if all costs are equal,  
 $b^d = 1 + \lceil C/e \rceil$

uniform cost ~~search~~ does more work as it checks all paths from  $S$  to  $G$ .  
So start ;  $G \rightarrow$  goal

### Time Complexity

slide 28

$b=2$

Path : ABG

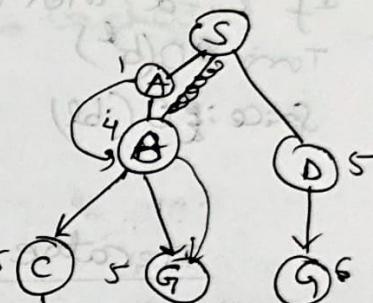
$$\text{Path cost} = 1 + 4 + 5 = 10.$$

$\epsilon = 1$

$$T.C. = O(b^{1 + \lceil C/\epsilon \rceil})$$

$$= O(b^{1 + \lceil 10/1 \rceil})$$

$$= O(2^{11}) \approx O(2048)$$

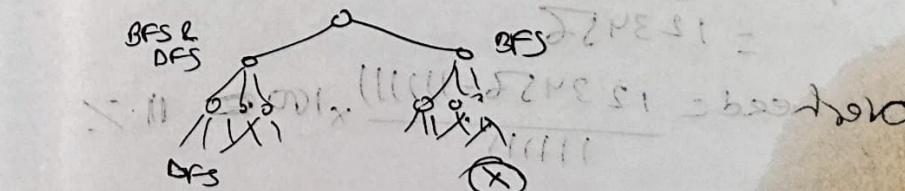


### DFS

- If any rule becomes redundant, backtracking is done → hence it is not optimal.
- Depends on choice of depth

Time:  $O(b^m)$  ← if solutions are deep, it is faster than BFS. ( $m > d$  ?!)

Space:  $O(bm)$  ← linear



## Depth-limited search (DLS)

- DFS with limit  $d$
- Tackles drawback of infinite path.

If  $d \neq l$  then it is not complete.

$$\text{Time} = O(b^d)$$

$$\text{Space} = O(bd)$$

## Iterative deepening (IDS)

- Gradually increasing the limit, until a goal is found.
- Combines benefits of BFS & DFS.  
memory:  $O(bd)$ .

### EXAM file slide 47

No. of nodes generated in DLS :-

$$N_{DLS} = b^0 + b^1 + \dots + b^{d-1} + b^d$$

$$N_{DLS} = (d+1)b^0 + db^1 + (d-1)b^2 + \dots + 3b^{d-2} + 2b^{d-1} + 1b^d$$

$$b = 10, d = 5$$

$$N_{DLS} = 1 + 10 + 100 + 1000 + 10000 + 100000 \\ = 111111$$

$$N_{IDS} = 6 + 50 + 400 + 3000 + 20000 + 100000 \\ = 123456$$

$$\text{Overhead} = \frac{123456 - 111111}{111111} \times 100 \approx 11\%$$

based on  $b$  &  $d$  (2 parameters) / numerical exam

DLS vs DFS

DLS vs IDS

IDS vs BFS

- Q) If  $b$  is more, which algo to choose?  
If  $d$  is less, which algo to choose?

## IDS

$$\text{Time} = O(b^d)$$

$$\text{Space} = O(bd)$$

» In exam, create comparative table & if values of  $b$  &  $d$  are given then compute values.

## Informed search

Heuristic function → Based on knowledge we have to develop heuristic function.

Randomly without information → uninformed search.

Heuristic fn, if not designed properly, may fail to find the soln due to limited knowledge. (mis-utilisation of knowledge)

Provides an estimate of solution cost.

→ Exploring the path by exploiting the knowledge.

This may lead to many wrong goal states.

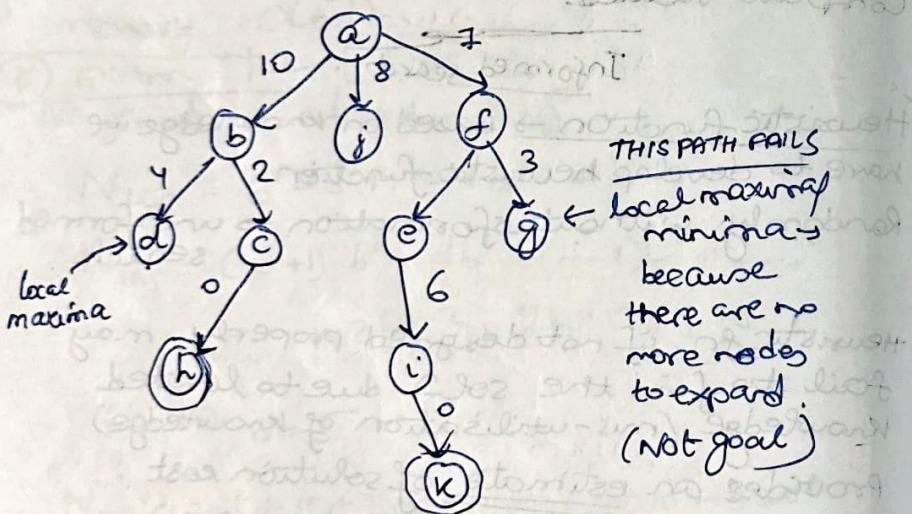
## Hill climbing Algo

- Generate & Test Algo : Generate a fn → if test fails then generate another.
- from initial state & use heuristic fn to generate another state.

$f(\text{start}) \dashrightarrow f(n)$

If  $f(n) < f(\text{start})$  then expand and check if  $f(n)$  is the goal state.

At goal state, value of heuristic function is zero.



## Disadv

- Since we are expanding only one node it may get trapped at local maxima.

## Gradient Descent (Steepest Hill climbing)

As opposed to simple hill climbing, it checks all successor states.

At a local maxima/minima, every move around it produces worse results (trapped).

## Soln

- Give up
- move from the trap through a big jump → in order to find another place where slope is available. Jump means to randomly select another step (which may be trapped again).

In a valley (plateau) → flat line → no slope is available & the algo does not move.

Ridge → multiple local maxima → multiple steps are required to jump from a local maxima/minima.

## Sln

- Backtrack to earlier node in diff-direction.

## Block world problem

If a block is above another desired block as per goal state then +1 (got) dots that may for wrong block it is -1 (missed block)

Since we are only considering block below it, we are using local heuristic

This may lead to many wrong goal states.

due to incorrect heuristic function selection.

- changing the heuristic function ~~to using~~ global info, then it gives better path with unique states.

### Simulated Annealing (SA)

Depends on how a metal is heated & based on its properties it reaches maximum energy level.

Here, heat  $\rightarrow$  knowledge.

After heating it to highest energy level, it is cooled to make it crystallised.

This is the philosophy behind SA.

The cooling down (i.e. change of energy) is done probabilistically to the point where no more change is possible (global minima).

Controlled Randomness

It enables more thorough exploration

From start state (top), the downward move is made to some state, before the randomness is applied.

$\Rightarrow$  Layer wise search in controlled manner.

If change of weight is high, it may overshoot. If it is low, traverser is more, but chance of overshoot is less.

Accordingly energy is modified in such a controlled manner.

$$p = e^{-\Delta E/T} \quad : \text{prob. for transition to higher state.}$$

$\Delta E \rightarrow$  change in energy

$T \rightarrow$  Temperature (Random)  $\rightarrow$  based on annealing schedule

Generate random no. b/w [0,1], if  $r < p$ , then proceed to change.

### Best-first search (OR graph)

$\Rightarrow$  Each path is equally important.

The nodes that were not expanded in a previous iteration remains in contention for expansion in the subsequent iterations. (the ~~to~~ node having lowest value of heuristic cfr is expanded).

$\Rightarrow$  Combination of DFS and BFS.

Make 2 lists: open & closed.

↑ already examined.

87	4
89	4
82	4
0	1 (4)

## Evaluation function

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

$g(n)$  = exact cost of cheapest path so far

$h(n)$  = estimated cost of cheapest path to goal

$f(n)$  = estimated total cost of cheapest path

At goal state  $h(n)=0$

In UCS → path cost is considered,  $h(n) \geq 0$ .

In Greedy (best-first),  $g(n)=0$ .

path cost is not important

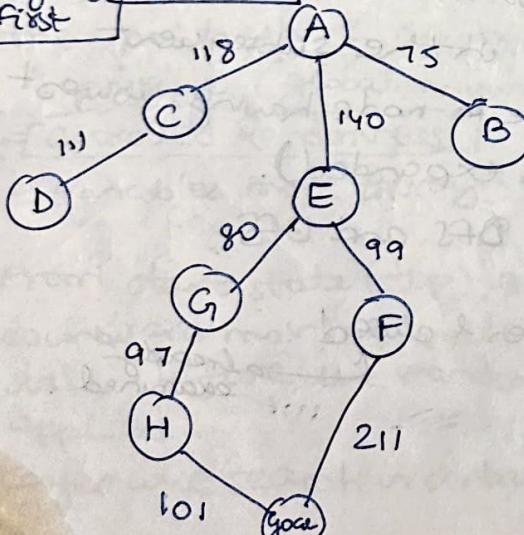
Admissible property → it never overestimates or underestimates the path cost.

## Greedy best-first

Time:  $O(b^m)$

Space:  $O(b^m)$

### Greedy best first Assignment



STATE	$h(n)$
A	366
B	374
C	329
D	244
E	253
F	178
G	193
H	98
goal	0

$$f(n) = h(n) = \text{estimated dist from start}$$

## A\* algorithm

17/01/2025

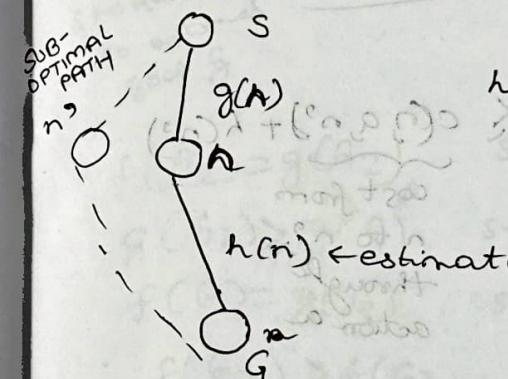
Used for Robot Path planning.

Considers the best of both Dijkstra's and best-first search (greedy).

Measures

→ Distance from start

→ Smart estimate of remaining distance to goal



$h^*(n) \leftarrow$  exact cost

optimal path  
 $h(n) \leq h^*(n)$

Graph search algo

Open List: Nodes to be eval. Sorted by  $f(n)$  in asc.

Closed List: ~~eval~~ Already eval. nodes.

If no path exists then reconstruct the best fit heuristic algorithm.

### PROPERTIES

I. Admissibility property → Never overestimate the cost to reach the goal  
 $h(n)$  is an estimate.

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

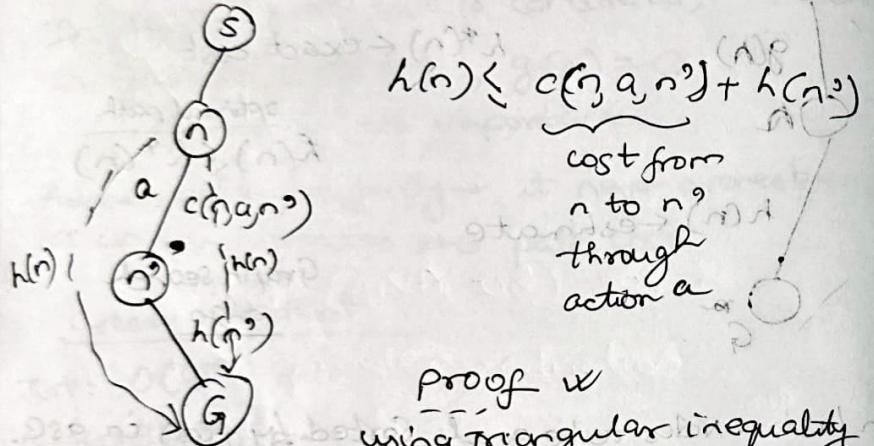
↑  
exact cost to reach n

→ estimated cost to goal from n.

## Admissible

1.  $h(n) \leq h^*(n)$
2.  $h(n) \geq 0$  &  $h(s) = 0$
3.  $h(n) \leftarrow$  straight line distance  
← minimum distance (no overestimation)

## I. consistency/Monotonicity



Start state has zero path cost.

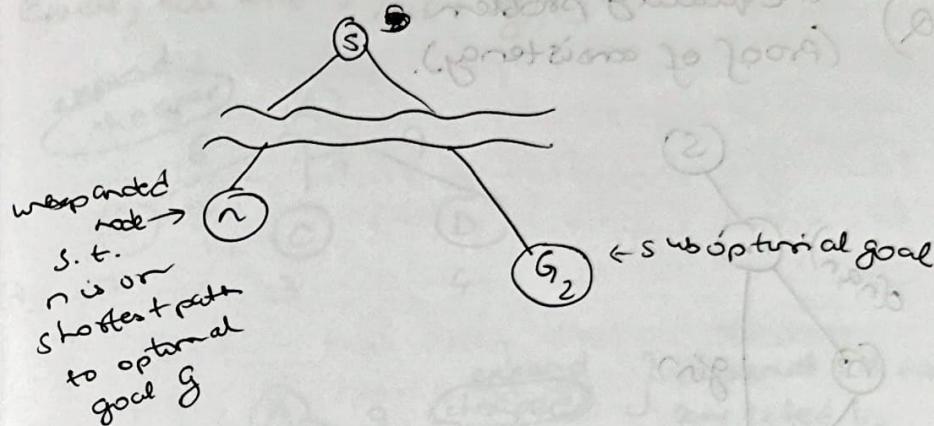
All nodes remain in contention for expansion

In A\* is 418 or 450 the suboptimal solution?

of 450 better since it's not loop

at least two nodes are better

## Optimality of A\* (proof)



$$f(G_2) = g(G_2)$$

$$g(G_2) > g(G)$$

$$f(G) = g(G)$$

$$\therefore f(G_2) > f(G)$$

$$\text{since } h(G_2) = 0$$

since  $G_2$  is suboptimal

$$\text{since } h(G) = Q$$

Hence  $f(G_2) > f(n)$ , and A\* will never select  $G_2$  for expansion.

$$f(G_2) > f(G)$$

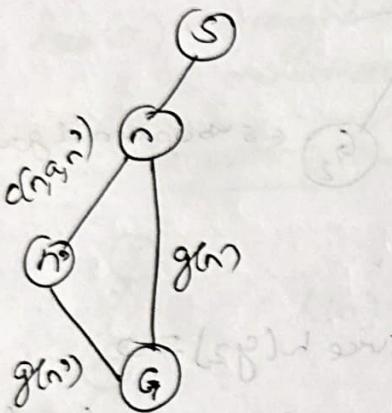
$$h(n) \leq h^*(n) - \text{since } h \text{ is admissible}$$

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

$$f(n) \leq f(G)$$

justify for optimality  
loop due out go

(HW Q) >> Consistency proof using triangular inequality problem.  
(Proof of consistency).



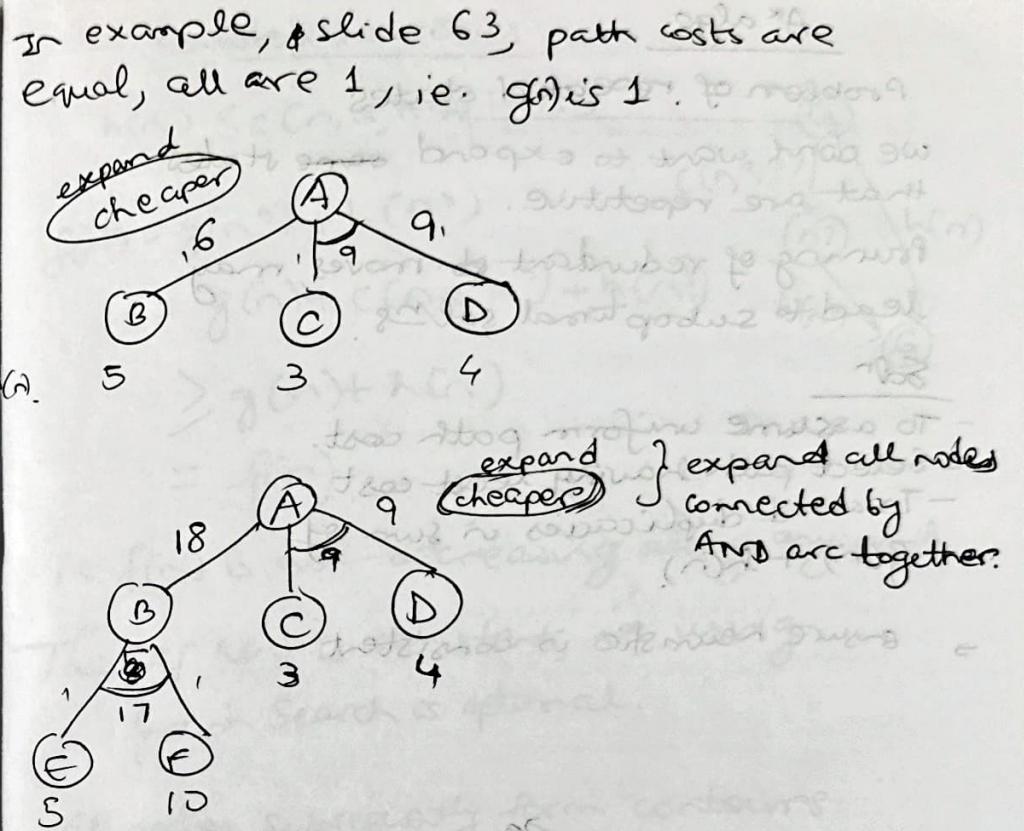
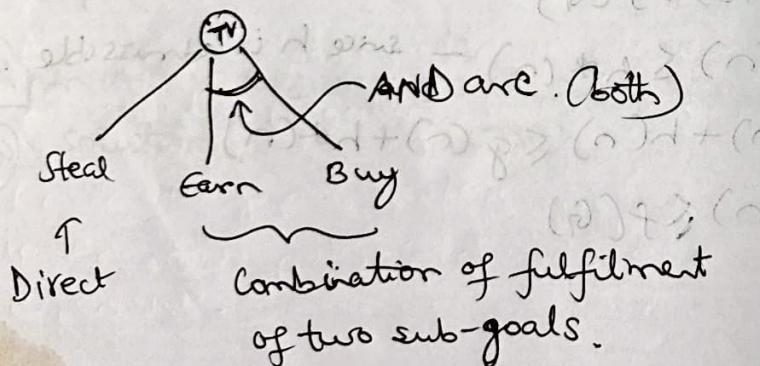
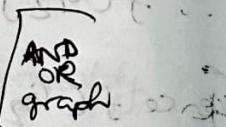
### AO\* (AND-OR) search

- used for problem reduction.

- Informed search algo.

AND  $\rightarrow$  all probs to be solved.

OR  $\rightarrow$  any prob to be solved



Expansion of path takes place of lower weights at root nodes which have been propagated upwards.

### Adv

- Complete
- No inf loops
- less memory

### Disadv

- Not optimal as it doesn't explore paths when it finds a sol?

## A\* algo

### Problem of repeated states

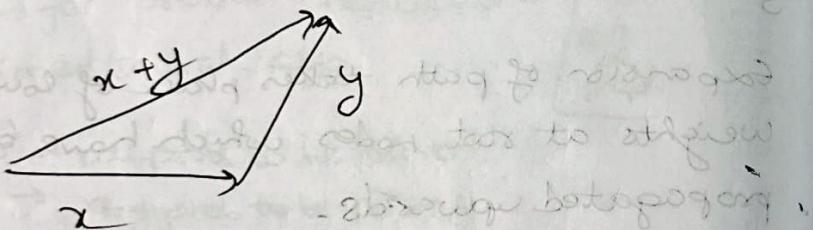
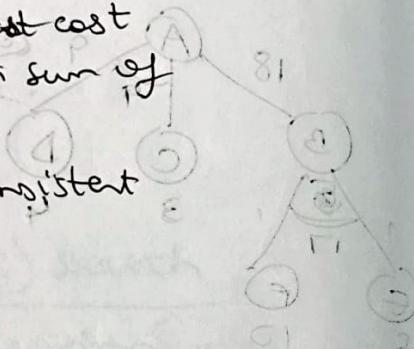
we don't want to expand ~~some~~ states that are repetitive.

Banning of redundant moves, may lead to suboptimal solns.

### Soln

- To assume uniform path cost.
- Select path having least cost.
- To avoid duplicates in sum of  $g(n) + h(n)$ .

→ ensure heuristic is consistent



More no. of parameters increases flexibility. (Features to extract)

Lesser no. of parameters means lesser no. of paths and more expensive.

### Consistent heuristics

$$h(n) \leq c(n, a, n') + h(n')$$

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

$$= g(n) + c(n, a, n') + h(n')$$

$$> g(n) + h(n)$$

$$= f(n)$$

i.e.  $f(n)$  is non-decreasing along any path.

Th: If  $h(n)$  is consistent, A\* using graph search is optimal.

All nodes subsequently form contours increasing f-value → expansion

### Adv

Complete

Time :  $O(b^d)$

Space : All nodes in memory.

Optimal : Yes



Q)  $y$  vs  $x$ , Robot planning problem.

It tries to avoid obstacles

Find the cost using traditional and A\* algorithm.

Divide it piecewise

Draw it in a grid structure.

Each grid is a state.

Check obstacles in each state.

Take variable heuristic for with variable parameter.

### Agent and Environment

Agent  $\rightarrow$  Human / Robot

Environment  $\rightarrow$  Consisting of many states and obstacles from which agents learn.

- There will be a feedback to agent to follow a path
- Do not need huge amt. of data

Calculate 2 signals

- $\rightarrow$  Reward
- $\rightarrow$  Penalty

} Feedback signal

### Intelligent Agents

$\rightarrow$  Rationality

$\rightarrow$  Difficulty nature of environment affects behaviour of agent.

$\rightarrow$  Agent perceives entire sequence & takes actions accordingly

$\rightarrow$  Markov decision process  $\leftarrow$  Next state depends on previous state only,

#### Evaluate

Multiple actions can lead to a state. Evaluation of best from those actions.

Agent  $\rightarrow$  sensor  
 $\rightarrow$  actuator

### Agent $f_r$

$$f: P^* \rightarrow A$$

$\downarrow$  percept histories  
 $\downarrow$  actions

Agent = Arch + Prog.

$\downarrow$  actuator  
 $\uparrow$  sensors  
from sensors & direct action.

State of env changes thru action.

Mapping a state to action  $\rightarrow$  Condition - Action Rule

Lookup table  $\rightarrow$  Sequence of actions

### Vacuum Cleaner

State : [locn, condition]  
eg [A, dirty]

Rationality

Rule  $\downarrow$  percept  
If (Condition) then (Action)

from perception through sensors like camera, etc.

### Rationality

Taking action at the right point of time

→ Reasonable : Reasoning power

→ sensible : Able to explain it

→ Good sense of judgement.

Performance measured through objective fn.

(Eg: amount of dirt cleaned up OR amt of time taken / electricity consumption)  
noise generated for the machine.

performance  $P$  is a fn of various params

$$P = f(a, b, c, d)$$

Need to learn the parameters to maximise performance.

percept sequence → previous, short term, long term.

~~knowledge~~  
Collections of many rules is called knowledge.

performance is expected, not absolute

Enhanced through perfection

Eg: Robo. A J. J.

Knowledge base & sensor to perceive.

### Application of Rationality

PEAS : Performance measure, Env, Actuator, Sensors

Environment → Surrounding conditions

Actuator → takes actions

Sensors → for perceiving

Performance measure → obtained through learning

Q) A problem can be given. Design the PEAS and knowledge base for that case.

### Env types

→ Fully observable : Nothing else needed to take optimal decision.

→ Partially observable : can take internal rep. of state thru knowledge base.

→ Deterministic : Current state affects next state.

→ Stochastic : Current state does not fully influence next state.

→ Strategic : Depends on actions of other agents  
Collaborative & competitive. (games)  
Eg: Robo-soccer      Eg: chess

→ Episodic : like AND-OR graph

↳ Can apply searching algo.

A sequential setting (episode) depends on prev episodes / seq. steps

Eg) Seq - setting → word

episode → sentence

- Multi agent systems → how agents collaborate?
- Static; Env remains unchanged
- Dynamic: keeps track of changes
- Semi-dynamic → Env doesn't change
  - Agent perf-score changes
- Discrete → No. of actions & states is finite
- Continuous → total actions cannot be perceived
- Single-Agent: One agent
- Multi-Agent: perf. of one agent affects the other.
- Taxi-driving is none of the env. types as mentioned in the slides.

Agent for  
Percept seq.  $\xrightarrow{\text{maps}} \text{Action}$

Rule: If (Percept) Then (Action)

Agent for is designed as a Lookup table.  
→ states & actions specified in advance by designer.

Drawbacks of table  
→ Cont. env. makes huge table  
→ long time to build & learn

IID → Identically Independently Distributed  
All samples are independent & need not be learned in sequence.  
Identical features are present for each sample.

Size of lookup table

$P \leftarrow \text{set of possible percepts}$   
 $T \leftarrow \text{lifetime of agent}$  (total no. of percepts received)  

$$\sum_{t=1}^T |P|^t \leftarrow \text{no. of entries}$$

Types of Agent

Simple reflex: Not model based. Reinforcement.  
 Model-based Reflex  
 Goal-based } clutched together  
 Utility-based  
 Learning agents:  
 For sequence agents we use (deep) reinforcement learning.

Types of Agents (based on user):  
 - Human  
 - Robot  
 - Software

→ Simple or partially observable env creates infinite loops. (unavoidable)  
 Make a big jump or backtrack any state if stuck in loop at plateau

Not Data-driven approach

Why do we need huge amount of data?  
what will machine obtain from data?

- SVM, Bayesian

↑ not much stochastic; mathematics based

more data gives better performance  
through non-redundant approach

More data leads to more variability  
(through more right)

Same i/p or less no. of v/fp doesn't  
induce variability.

Reinforcement learning does not require  
huge amt of data. It learns the  
env & how to change state based on  
current position.

Deep Reinforcement learning → ~~any~~

Actions are restricted.

Many states allows learning of  
different actions.

Knowledge base involves condition-action rules

It'sadv of simple

→ Very limited intelligence.  
→ Response to changes in env  
→ Fully observable.

Q) Generate rules for self driving car (PEAS)

Model-Based Reflex agent

→ guided by a model

→ Maintains internal state → representation of  
unobserved aspects of current state depending  
on percept history.

Depends on:

→ change of env in its own

→ Change of env through actions / ~~other~~ agents.

⇒ Rules for "driving back home" (model based)

Goal-based agent

→ goal-driven ; previous ones only based on state

→ objective: Reduce distance from goal.

↳ searching and planning.

→ ~~meas~~ perf. measure thru reward & penalty

→ What happens and is it nearer to goal?

Q) Not only applying ~~break~~ but checking light also.

→ Setting of goal in addition to previous.

→ Reasoning under uncertainty.

Utility Based Agent (Also goal based)

- Mapping to real no. (associated degree of happiness)  
is a utility function.

- Multi Object optimisation (MOOp)

- Everything is approximated under uncertainty.

Q) Diff b/w types of goal & how is it progressing?

## learning agent

→ learn from past experiences, not in prev.

### Components

- learning element → learning
- critic → feedback
- perf. element → select action
- prob. generator: → suggestive actions

## Multi Agent System (MAS)

- multiple interacting agents
- multiple knowledge bases
- common goal/objective.

Eg: Game theory, my agent-based modelling.  
↓  
same board,  
agents diff.

## Hierarchical agents

Eg: Robo soccer.  
Position of robots fixed;  
just differs in position from goal ⇒  
closer or farther.

## Artificial env.

- To test if w prog. is working or not.
- Program in art. env. is checked against human evaluation.
- Through Turing Test (Generative Advisory Model GAM)
- ? performing t w task & try to explain

## Turing Test

Check b/w human & machine response & if it can be determined separately.

25/01/2025

## Planning

- Regd to reach a particular destination
- finding best route: Tasks at particular time & why.
- Logical side; domain independent.
- Sequence of actions; next action depends on prev. one based on pre-conditions.

### Types

FSSP

Forwards State Space

Planning (Progression)

$S \rightarrow S' \rightarrow \dots \rightarrow G$ .

BSSP

backward

- $G \rightarrow \dots \rightarrow S$  (Regression)
- Check subgoals for consistency.
  - Small branching factor
  - Not sound algorithm

→ Stack planning → (bidirectional search algo) → combines BSSP & FSSP

## Planning system

Eg: Sussman anomaly.  
[Block problem]

### Interleaving

whether to go back a step at some iteration, i.e. undo the operation.

Closer to goal state so that loop / undo operation is minimised.

Precondition for block transfer, both need to be clear above.

- Planning → such that subtasks are in sequence and linked & maintain optimal cost.
- Intermediate moves follow a specific sequence, not randomised.

## Partial-order planning

- Delay during search
- Sequence to be established.

In sock & shoe planning, ~~the~~ sock on either leg is independent but shoe on a particular leg is dependent on sock on that leg beforehand.

The parallel plans are non-committal on which comes first.

## Target stack plan

- Knowledge base reqd for current situation and actions.
- Planning Graph: data structures used in automated planning.

○ State level → Nodes with logical propositions

□ Action → Influences state levels

Edges → State Node to Action

Arcs → Action to State Node.

— Mutex Relationship → cannot coexist simultaneously, reduces complexity.

### In graph

- many levels

- Level → Actions & literals

↳ Representations starting from basic primitives to high-level knowledge representation & planning

### Action

#### Persistence action

[It is carried forward to the next state].

Start: Have(Cake)

Goal: Have(Cake)  $\wedge$  Eaten(Cake)

{PLANNING GRAPH}

#### Point

##### Start

##### To

##### Goal

##### So

##### No Opt

##### Have(Cake)

##### Have(Cake)

##### Eat(Cake)

##### Have(Cake)

##### Eaten(Cake)

##### Eaten(Cake)

##### Mutex among actions

→ Mutex conditions among predicates that have been derived.

#### Mute x Cond

→ Interference

→ Inconsistent effects

→ Competing nodes.

Mutex link defines which literals can't exist together.

Have(Cake) and Eaten(Cake) are mutually exclusive.

→ Generate nodes with further actions till goal literals do not have mutex links (Termination condition)

## Reasoning uncertainty

### Deduction

- formal logic.
- certain
- known premise  $\rightarrow$  Conclusion

### Induction

- informal: does not follow some rules of prob
- uncertain (probabilistic)
- Plausible, believable, probable, reasonable

## Reasoning under uncertainty

$\rightarrow$  creating machines with intelligent behaviour

$\rightarrow$  FOL (First Order Logic): predicate logic  $\rightarrow$  cannot handle uncertainty.

$P \rightarrow$  pair  $\rightarrow$  has different degrees of tolerance and hence induces uncertainty.

### Predicting

- $\rightarrow$  Probabilistic reasoning to determine chance.
- $\rightarrow$  Cause to effect

### Diagnosing

- $\rightarrow$  Effect to cause
- $\rightarrow$  Backward
- $\rightarrow$  forward

Paradigms: Bayesian & Non-Bayesian

### Bayesian Approach

$$P(H|E) = \frac{P(H)}{P(E)} P(E|H)$$

Disadv: Variables are uncorrelated (independent)  
(criticism: not true)

- Determines chance of truth or falsity of event.

Belief  $\rightarrow$  based on ~~personal~~ knowledge & perception

$\rightarrow$  It is a personal value & depends on above 2 parameters. (Subjective value)

$\rightarrow$  New evidence affects the belief.

$\rightarrow$  Prior belief: prior to new evidence

### Bayesian depends on

- Joint probability distribution & Conditional probability distribution  
(Joint occurrences, relatively harder to determine)
- Acyclic directed graph (DAG).

### Belief Networks

- $\rightarrow$  Nodes: Set of random variables (observations)
- $\rightarrow$  Directed links:  $X \rightarrow Y$  influences  $Y$  directly. ( $X \rightarrow Y$ )
- $\rightarrow$  Conditional probability table: effects of parents on each node

Joint probability  $P(X_i | X_1, X_2, \dots)$   
parents immediate parent(s) in graph (DAG)

$\rightarrow$  Only initial variable is given

$\rightarrow$  Revise tree until difference is greater than threshold.

### Knowledge base

If  $X$  then  $Y$  with some probability  $p$ .

Probability of  $E$  is normalized with  $H$  and  $H'$ .

$$P(E) = P(E|H) P(H) + P(E|H') P(H')$$

## Knowledge-based systems

Knowledge is represented in IF-THEN format  
 (uncertainty is induced by assigning probability 'p' to the rules).

Abductive reasoning → giving more emphasis on observations.

If freeze then have cold

IF  $X$  then  $Y$ . with prob  $p$   
 we need to determine  $Y$  without knowing anything about  $X$ . (using Bayes' rule), i.e.  $P(H|E)$

$$P(H|E) = P(E|H) \cdot \frac{P(H)}{P(E)}$$

$$P(H|E^c) = \frac{P(E|H) P(H)}{P(E|H) P(H) + P(E|H^c) P(H^c)} = P(E)$$

### Propagation of Belief

- Each piece of evidence affects only one hypothesis
- In real life, multiple evidences ( $E_1, \dots, E_n$ ) affect multiple hypothesis ( $H_1, \dots, H_m$ )

$$P(H_i | E_{j_1} \cap E_{j_2} \cap \dots \cap E_{j_k}) = \frac{P(E_{j_1} \cap E_{j_2} \cap \dots \cap E_{j_k}) P(H_i)}{P(E_{j_1} \cap E_{j_2} \cap \dots \cap E_{j_k})}$$

Posterior probabilities

### Assumptions

- All hypothesis are mutually exclusive.
- There is at least one hypothesis
- Hypothesis are collectively exhaustive
- Evidences are conditionally independent.

Hence we use Bayesian method

$$P(E_1 \cap E_2 \cap \dots \cap E_n) = P(E_1) \cap P(E_2) \cap \dots \cap P(E_n)$$

- Prior prob of hypothesis given.  $P(H_i)$
- Conditional prob of matrix given for each evidence.  $P(E_j | H_i)$
- Single and combined evidence
- $P(H_i | E_1 \cap E_2)$ : find hypothesis for combination of evidences

### Burglar alarm at home problem

Hypothesis → Alarm is sounded @ Harry's house

Evidences → No Burglary

No Earthquake

John & Mary both called Harry

Joint probability distribution.

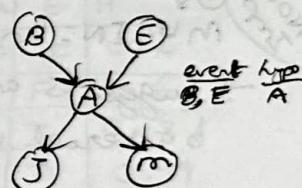
J → John calls

m → Mary calls

A → Alarm sounded

B → Burglary

E → Earthquake



Draw the network of in question.

$$\rightarrow [J, m, A, B, E]$$

$$= p[J|m, A, B, E] \cdot p[m|A, B, E]$$

$$= p[J|m, A, B, E] \cdot p[m|A, B, E] p[A|B, E]$$

$$\vdots$$

$$\rightarrow p[J|m, A, B, E] p[m|A, B, E] p[A|B, E] p[B|E] p[E]$$

$$p(J \wedge m \wedge A \wedge \neg B \wedge \neg E)$$

$$= p(J|A) p(m|A) p(A|\neg B \wedge \neg E) p(\neg B) p(\neg E)$$

parent

→ all values given in table

Combinations of factors for alarm

$$P(A|BE) P(BE) \cap P(A|B'E) P(B'E) \cap P(A|BE') P(BE') \\ \cap P(A|B'E') P(B'E')$$

$$P(A'|B) = P(A'|BE) + P(A'|BE')$$

Inference from belief & w

→ cause to effect → Effect to cause.

### Adv & Disadv of Bayesian

- Sound theoretical foundation
- Requires large size of samples to obtain acc. prob.
- Consistency & comprehensiveness of values

Y/N  
No  
Yes

MYCIN: Uncertainty Handling. (Stanford 1972)

→ Suggests antibiotics and dosages of  
bacterial infections

LISP  
program

- Knowledge Base
- Patient Database (large sample size)
- Product of probabilities not great option

### Utility Theory

- in alternative framework
- Determines how AI systems make choices among options based on perceived utility

Objective → Certain

Subjective → depends on perception

$x \rightarrow$  outcome

$U(x) \rightarrow$  fn.  $x \rightarrow$  utility value

$P(x) \rightarrow$  prob.

Maximum Expected Utility (MEU): optimises the decision-making principle.

- orderability: if  $U(x) > U(y)$ , then  $x$  is preferred over  $y$

exam Q) write short note on "utility".

= transitivity

- Continuity

- Substitutability

- Monotonicity

- Decomposability

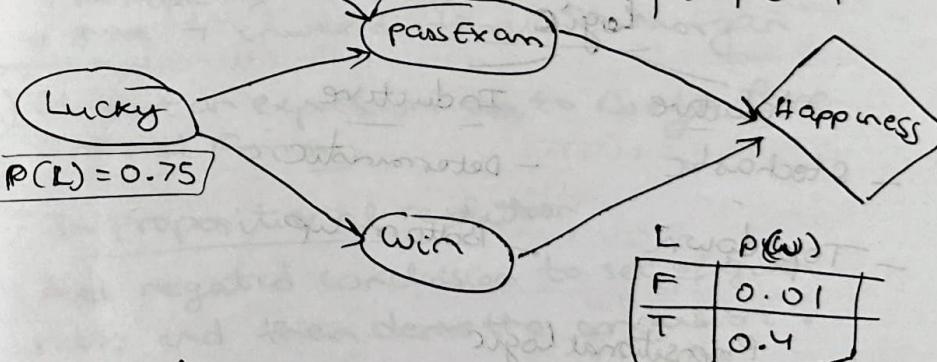
Q)

### Assignment

L	S	P(E)
F	F	0.01
T	F	0.5
F	T	0.9
T	T	0.99

$$P(\text{study}) = 0.8$$

W	E	H
F	F	0.2
T	F	0.6
F	T	0.8
T	T	0.99



Find  $P(H)$

Joint.

Probability of happiness for all conditions

Get joint & product  
select specific - distributed

Propositional Logic

Proof or validation behind any reason provided

Inference engine  $\leftarrow$  domain-independent

Knowledge base  $\leftarrow$  domain specific

Logic

Formal languages for representing info so that conclusions can be drawn.

entailment  
 $\text{KB} \models a$

} one thing follows  
another

- relations based on sentences based on semantics.

models

formally structured words from which we can evaluate the truth.

Logic

- Deductive
- Stochastic
- Top down

- Inductive
- Deterministic
- Bottom up

Propositional Logic

proposition: declarative statement which is either true or false.

Symbolic logic  $\rightarrow$  because it uses symbols

Tautology  $\rightarrow$  always True

Contradiction  $\rightarrow$  always False

propositionsAtomic

true/false sentences

Compound

- Combination of many atomic propositions.

$$P \rightarrow Q \equiv \neg P \vee Q$$

Validity  $\rightarrow$  Sentence is true in all models.

Satisfiable  $\rightarrow$  Sentence is true in some model, at least one

Unsatisfiable  $\rightarrow$  Sentence is true in no model

CNF) to convert to (disjunctive) normal form

~~Sum~~ Product of sum:

clauses individually true makes whole true.

Convert to CNF

$\rightarrow$  Remove  $\leftrightarrow$  to  $(P \rightarrow Q) \wedge (Q \rightarrow P)$

$\rightarrow$  Remove  $\Rightarrow$  to  $\neg P \vee Q$

$\rightarrow$  move  $\neg$  inwards thru De-Morgan

②) Convert an expression to Conjunctive Normal Form.

In propositional resolution

Add negated conclusion to set of given rules and then derive the conclusion.

## Game Search

31/01/2025

Prev search: - Agent completely controls environment.

- Deterministic
- State doesn't change (unless changed by robot)
- Single path from start to goal.

Game Search - Multi-agent

- Stochastic
- Non-deterministic path.

Every agent considers actions of other agents (Cooperative/Competitive) & optimises its own benefit.

> Adversarial search problem  $\rightarrow$  Game Search

- Impact of one agent significantly impacts other agents; in a Multi-Agent Environment (MAE)
- At the end, utility values are equal and opposite, for 2 agents with alternating actions.

Eg: DeepBlue by Apple.

- Easy game state,
- limited action rules for agents
- Precise rules define outcomes.

Eg: Chess

- Avg branching factor = 35.
- 50 moves each player, search tree has  $35^{100}$  nodes.

Brute force becomes inefficient here.

Typical Search  
→ NO opponent

game search  
→ unpredictable opponent.

Zero sum  $\rightarrow$  sum of player & opponent's move is zero.  
Win  $\rightarrow +1$ , lose  $\rightarrow -1$ , draw  $\rightarrow 0$ .

Eg: Rock, Paper, Scissors

Sum at each cell is zero, so zero-sum.  
If draw, both zero; else  $+1/-1$  or vice versa.

Eg: Battlebots

State-space representation

0,0		10,8
8,10		0,0

Tic-tac-toe

$\rightarrow$  2 players: MAX & MIN

$\rightarrow$  MAX wants to maximise utility & vice-versa  
 $\rightarrow$  MAX plays first

for terminal node: Max  $\rightarrow u(t)$   
Min  $\rightarrow -u(t)$  } zero-sum.

a  $\rightarrow$  b, c, d  $\rightarrow$  ...  
Recursive strategy.

Ply: Half-move.

From first max node expand 3 possible min nodes, then moves to terminal. Backtrack using smallest value (min) and evaluate.

MAX will choose largest of above obtained smallest values.  
using DFS  $\rightarrow$  This computation takes place.

(Minimax Algorithm)

## Minimax Algo -

uses DFS

Depth  $\leftarrow m$

legal moves at each point  $\in b$ .

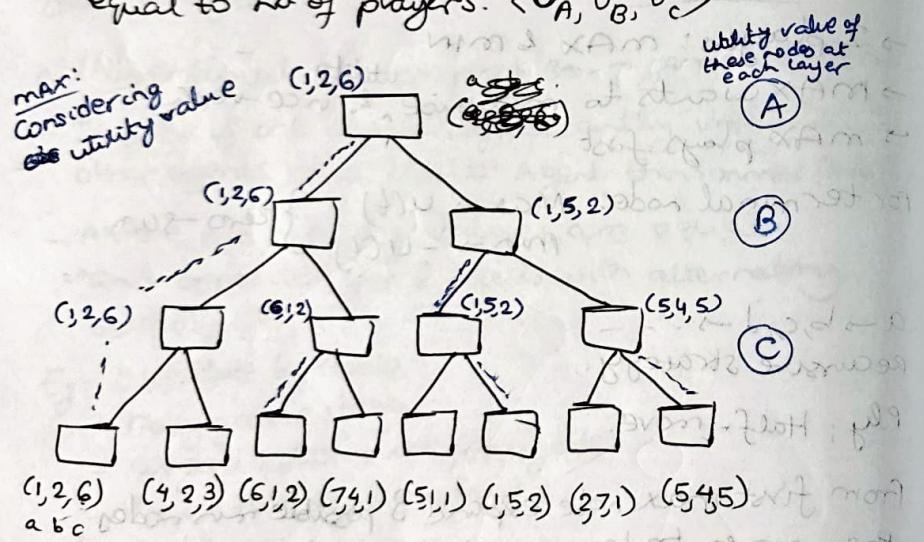
Time complexity:  $O(b^m)$ .

Space:  $O(bm)$   $\leftarrow$  all actions together.

:  $O(m)$   $\leftarrow$  actions one at a time

## In multiplayer games

$\rightarrow$  Single value of each node is replaced with a vector of values having dimensions equal to no of players.  $\langle v_A, v_B, v_C \rangle$



## minimax

Complete: Yes (if tree is finite)

Optimal: Yes (against optimal opponent)

Time =  $O(b^m)$

Space =  $O(bm)$

## Alpha-Beta Pruning Algo

- $\rightarrow$  prune branches that need not be expanded again, i.e. need no more examination.
- $\rightarrow$  For MIN, pruning trees having higher than that value (atmost)  $\leftarrow \beta$ -cuts
- $\rightarrow$  For MAX: (atleast)  $\leftarrow \alpha$ -cuts (prune lower values)

After generating value for only some of  $n$ 's children, we can never reach ' $n$ ' in a minimax strategy.

Alpha: highest value found till any point.

Starts with  $-\infty$

Beta: lowest value choice found at any point.

Starts with  $+\infty$ .

In backtracking node values pass to upper nodes

$\alpha \geq \beta$  for pruning

At starting node  $[\alpha, \beta] = [-\infty, \infty]$

This is the range of possible values

$\rightarrow$  Root layer is MAX

layer 1  $\rightarrow$  MIN

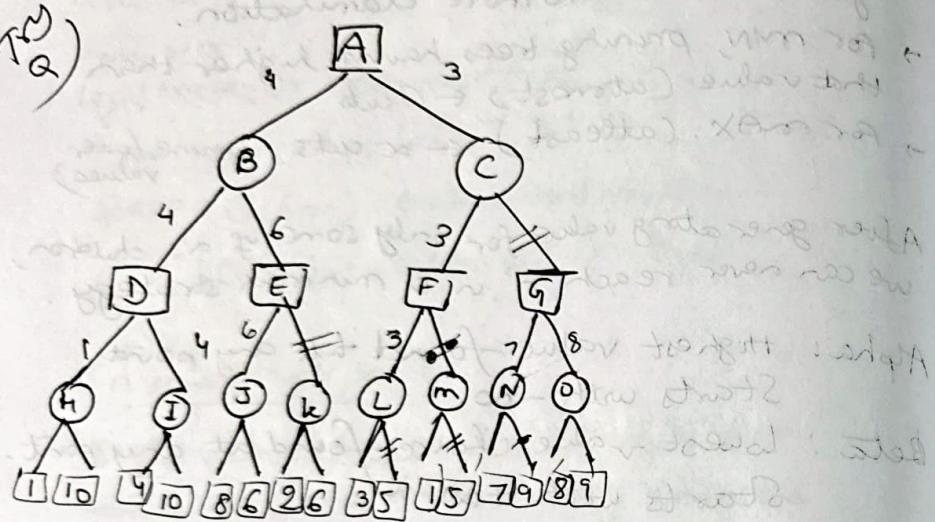
and alternatingly so on

The branch having lower value outside range is pruned and not examined

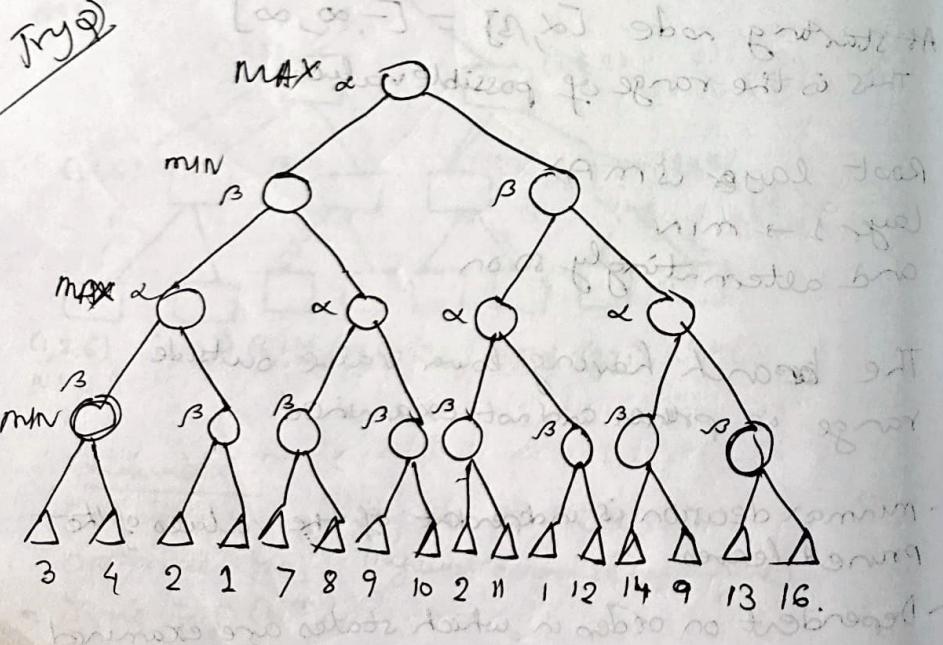
- Minmax decision is independent of the values of the pruned leaves

- Dependent on order in which states are examined.

## Alpha-Beta Pruning Procedure



Pruning only affects time complexity, not final result.



## Reinforcement Learning

03/02/2025

Unsupervised learning → No output values, hence no change of weights, no error calculation

Reinforcement learning is a Trial and Error method. Error is modified through experience, i.e. we learn from errors.

(Cycle riding):  
Agent: Person ; Penalty: Pain  
Environment: Cycle, Road ; Reward: Compliments from friends & parents  
Error: Falling down.

We obtain reward/penalty through error.

This is called evaluation. It can be mapped to the cost function of supervised learning.

State: Start state

Action: Paddling

State (next): Forward / Backward / Falling.

→ State, Action, State<sub>2</sub>

Evaluation / Payoff → method of obtaining reward / penalty.

- Learning procedure: gathering experience.
- Not much data.
- Computational resource very less (GPU not reqd) → Hence less energy required.

### YT Lecture:

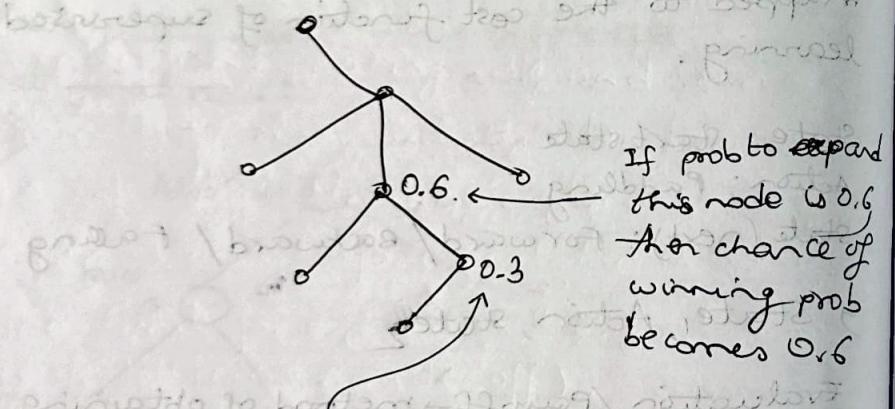
Prof. Swagatam Das (ISI)

→ Reward obtained through Reinforcement

learning is not immediate reward, rather intermediate.

After taking an action, if we determine that a previous action is not good, then we go back to modify the action.

Reward is always expected reward, not exact reward, due to dynamic environment



If prob of winning decreases after a layer, i.e. towards the goal (i.e. there is less scope of adjustment), then the prev value is updated through experience.

Rather than immediate reward, it is cumulative reward (or penalty).

upon reaching a goal state, the path is traced back and learnt.

Evaluation takes place after a certain outcome.

Close to the goal we are more confident to win.

### Temporal Difference Learning

Temporal → Time-variant

~~Action  $A(t+1)$  depends on  $A_t$~~

→  $A(t+1)$  (Action at  $t+1$ ) depends on  $A_t$  (action at  $t$ )

→ Prediction at time  $t+1$  is more accurate than prediction at time  $t$ .

→ From feedback based on error, the previous state has to be modified.

### Expected Reward

- Non-deterministic/Stochastic/Random

[K-means clustering has initialisation problem due to randomness in initialising]

### How to solve randomness

Run the algorithm many times and then take the average.

Now how prev a value is going to change?

For expected reward,  
play the game from start to finish and  
multiple times and take average  
reward over each such path.  
(since winning is not always guaranteed).

> Exploitation: Using the available (repeated) information only for learning. But there is no scope for learning outside that knowledge. It cannot give incremental knowledge. It may lead to missing of some better optimal result/opportunity.

> Exploration: Taking risk and searching (new paths) something else. Initially when knowledge is not available, then exploration is done. After obtaining info, exploitation is done. The limit for exploration & exploitation is a dilemma (trade-off).

→ If opponent is random, then that creates noise (which creates problem).

→ If opponent is not random, and also plays according to fixed rules, and if we are playing Greedy<sup>to win</sup>, then we will be just playing a very very small

-tree/graph  
factor of the game, without exploration → it is the best choice.

So, exploration may not give us the immediate benefit.

Eg): Bandit Problem or chess.

Book: Richardson

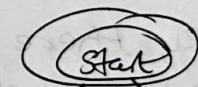
Greedy agent can get stuck in sub-optimal state.

Epilon Greedy Algo

choose a random option with probability of  $\epsilon$ . Set less probability to start exploration, initially, ~~initially~~. After some information is available, exploitation is done.

Exploration:  $\epsilon$

Exploitation:  $1 - \epsilon$



$\epsilon$

$1 - \epsilon$

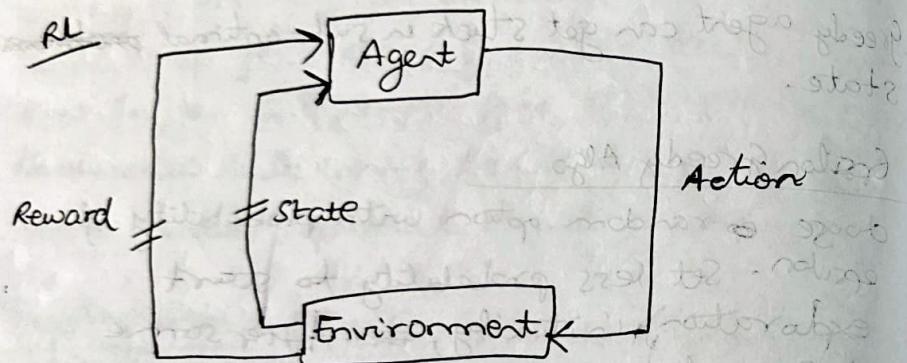
A Reward Signal

## (IRL) Immediate Reinforcement Learning

11/02/2022

In normal reinforcement learning, there is a delay in obtaining reward/payoff. It takes some time to take decision.

→ In IRL, evaluation happens after performing each action, i.e. there is no delay.



→ In IRL, there is no 'state' input, only 'reward' is present as input. Hence, we try to repeat the same action multiple times, and it may not always be successful → hence, the nature of the problem is stochastic.

Since the problem is stochastic, we determine expected reward, not exact reward.

→ At every time ( $t$ ) instance, pick up an action ( $a_t$ ) to get reward ( $R_t$ )

IRL is a simpler version of RL.

- It does not have temporal nature
- Reward is not a scalar evaluation, because evaluation is stochastic
- Keep trying any action multiple times before figuring out which is the right action.
- IRL exploits the same action; it does not involve exploration.

Eg:- IRL is useful in Drug Testing

Apply drug-1 to patient-1 and patient-2

Observation: P-1 is cured, P-2 is not cured.

Can it be concluded that D-1 is effective for P-2 and ineffective for P-2?

It cannot, because P-2 may have other symptoms.

It is repeated. D-1 is applied to P-1 again. However, P-1 may develop newer symptoms and the D-1 may render ineffective.

Here, the environment is stochastic in nature. Always applying drug-1 is exploitation. Using drug-2 would mean exploration. But there is a dilemma in the no. of times D-1 and D-2 should be applied.

Patient being cured is a reward.

Stopping exploration may lead to missed info, too much exploitation leads to no new testing.

## Formalism

set of actions

$$A = \{1, 2, \dots, n\} \quad \text{e.g.: drug-1, drug-2, \dots}$$

calculating reward which is stochastic in nature.

Say, Associated with each action there is a coin.

~~Let~~

i.e. If there are  $n$  actions, there are ' $m$ ' coins.

Upon taking an action, we toss a coin (not the same for two actions). The coin is fixed for a particular action.

After tossing some probability to coming up for Heads, and some probability to coming up for Tails.

Let, for action 1, tossing a coin, if it is Heads, then evaluation (payoff/reward) is 1, and if it is Tails then reward is 0.

## Probability Distribution

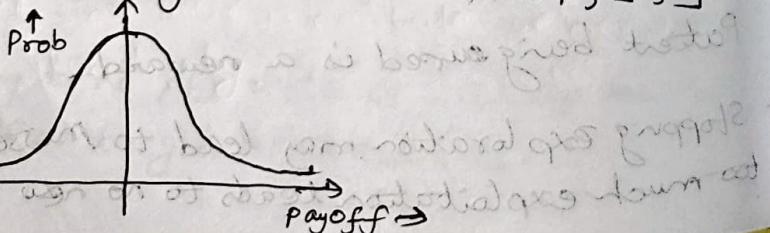
Tossing of a coin is Bernoulli distribution.

According to central limit theory, it is Gaussian distribution.

The probability distribution gives the expected reward.

Process

→ Generate a uniform random number,  $[0, 1]$



→ Keep integrating the area under the curve, called CDF (Cumulative Distribution Function), gives the probability that the random variable  $X$  is less than or equal to the probability.

→ When CDF reaches the random no., then the sample is considered.

For instance, let us have four possible outcomes (payoff)

We want to generate one of those 4 outcomes at random. Each of these outcomes have probability  $P_1, P_2, P_3$  and  $P_4$ .

Step 1: We generate uniform random no.:  $[0, 1]$

Step 2: Check if  $r < P_1$ , then select Outcome-1

Step 3: Check if  $P_1 < r < P_1 + P_2$

Since it is cumulative reward.

Then generate Outcome-2

Step 4: Check if  $P_1 + P_2 < r < P_1 + P_2 + P_3$

Then generate Outcome-3

Step 5: Check if  $P_1 + P_2 + P_3 < r < P_1 + P_2 + P_3 + P_4$

Then generate Outcome-4

## Challenge of IRL

- We do not know the distribution (exactly). To tackle that, we go for pure learning. i.e. where the mathematical modelling fails, we go for algorithms in iterative forms. (this is same as Machine Learning).
- However, we know that there is some distribution, which we are assuming according to which the payoffs are generated. (i.e. we do not know the parameters of the distribution)
- For this, we assume that there is an expected payoff, denoted by  $q^*(a)$

$q^*(1) \rightarrow P_1$ ; for the previous example.

$q^*(2) \rightarrow P_2$ , and so on.

For Gaussian distribution,

$$q^*(1) = \mu_1 - \text{error term}$$

$$q^*(2) = \mu_2 - \text{error term}$$

For IRL, since there is no state, we need to repeat iteratively.

## Multi-arm Bandits

14/02/2025

- Bandits are used in slot machines used to understand tradeoff b/w exploration and exploitation.
- Distribution of reward is unknown.
- Slot machine: pull lever, if symbols match then earn money else lose.
- The choice of lever is random and can be picked in any sequence.

Initial findings: - No. of winnings per 100 attempts per machine.

using this knowledge for choice of lever is exploitation, i.e. using same known knowledge, i.e. we are not learning new knowledge.

Hence, we need to go for exploration on the machine having lower win rate.

Feedback  $\rightarrow$  Reward / Penalty

- Stochastic: due to uncertainty in distribution of reward. Hence, we seek for expected reward.

- Selecting action value is unknown; hence estimate reward is computed through sample-average.

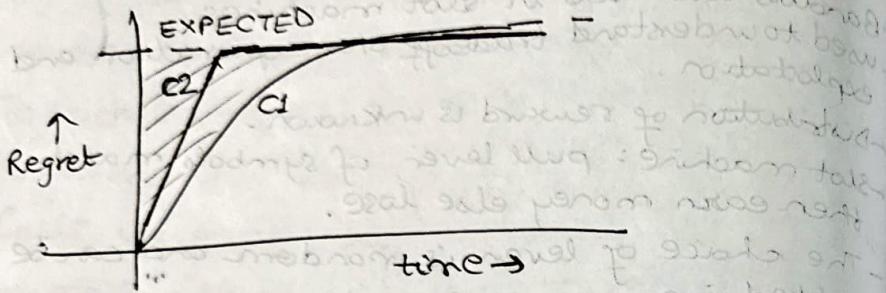
Value for at time t for action a

$$Q_t(a) = \frac{R_1 + R_2 + \dots + R_{N_t}(a)}{N_t(a)}$$

Objective: Over repeated action selections for 1000 times, we need to maximise winnings by concentrating actions on the best lever.

Regret  $\rightarrow$  Diff b/w optimal selection of lever.

2020/2021



Regret for  $C_2 <$  Regret for  $C_1$

→ Choose b/w multiple options, each having unknown payoffs. maximise reward over time.

**Trade-off**

- ↳ Balancing b/w探索(exploration) to learn about their reward distribution.
- ↳ Exploitation of known arms already having high rewards.

This is an iid problem (Independently Identically Distributed)

- $K$  independent arms
- Each arm ' $i$ ' has reward ' $R_i$ ' from unknown distribution with expected value  $\mu_i$ .
- Maximize cumulative reward over  $T$  trials.

Naive method.

→ Probability distribution of rewards is diff for each lever.

### Conflict

- Exploration → Improves current knowledge for long term benefit
- Exploitation → greedy action → exploiting current action-value estimates.

### Epsilon-greedy method

- Balances exploration & exploitation
- ↳  $\epsilon \leftarrow$  probability of choosing to explore
- $1 - \epsilon \leftarrow$  prob. of choosing to exploit

Action at time  $t$ :

$$\begin{cases} \max Q_t(a) & @ 1 - \epsilon \\ \text{any action } a & @ \epsilon \end{cases}$$

**Algorithm**

```

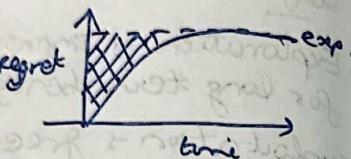
P ← random value
if P < ε :
    random action
else:
    current best action
  
```

If  $\epsilon = 0$ , i.e. greedy approach, it leads to saturation due to lack of exploration  
→ gets stuck for suboptimal actions like local optima.

- for value of  $\epsilon$  \*, 0.01 performs better than 0.1.
- 0.1 reaches optimality earlier but winning chance does not exceed 91%.
  - 0.01 increases slowly, but eventually performs better. (i.e. improvement is continuous)

## Regret minimisation

Trying to minimise area in graph



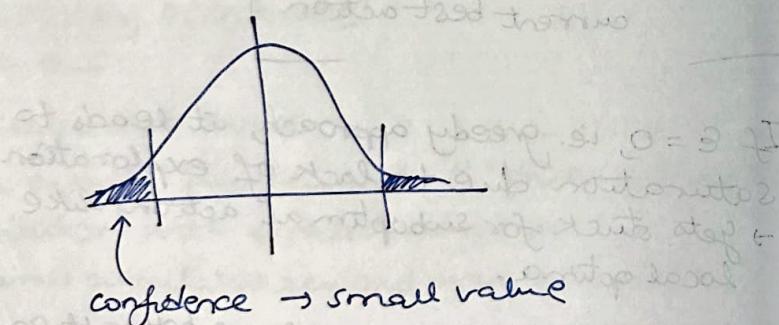
$$\text{Regret} = \frac{\text{payoff of possible action} - \text{payoff of action taken}}{\text{payoff of action taken}}$$

→ Asymptotic correctness:- Playing infinite times would lead to expected reward.

→ Regret optimality :- knowing best arm a priori and selecting that always.

→ probably Approximately Correct (PAC) gives lower bound of regret =  $\log T$ .

## Upper Confidence Bound (UCB)



Deviation from mean by this amount  
is acceptable.

Select actions based on:-

- Estimated rewards
  - Uncertainty / confidence in estimates  
(to be reduced)

+  $\epsilon$  time

$a_t \rightarrow$  action at time t

$a_t$  → reward at time t

$$u(t) = \bar{u}_i t + \sqrt{\frac{2 \log t}{n_i(t)}}$$

uncertainty through entropy measurement.

Algo

Initialisation:- Pull each arm once to gain preliminary info.

- select arm that maximises UCB value

- update pull count & estimated reward for arm

Estimated reward<sub>t</sub>: average of rewards

$$\bar{\mu}_i(t) = \frac{1}{n_i(t)} \sum_{s=1}^t r_{i,s}$$

$r_{i,s} \rightarrow$  reward on  
ith arm  
at s<sup>th</sup> pull

(to be minimised)

b) confidence level:  $\sqrt{\frac{2 \log t}{n_{\text{eff}}}}$

↳ favours arms that have been pulled less no. of times as they have larger confidence interval.

## Adv

- Theoretical guarantee: provides regret bound
- Simple & Efficient
- stationary environment:-

Each slot machine has a program.

The program remains the same during the evaluation process

If it changes, then it becomes non-stationary.