\* **Fuzzy Logic :-**

$\rightarrow$ description of statement

eg- The person is tall.

$$\left\{\begin{array}{l}\rightarrow \text{Membership function}\\ \hookrightarrow \text{Defuzzification}\end{array}\right.$$

What is the (measurement)/ % of a statement to be T/F.

$0-1 \longrightarrow$ fuzzy node $\longrightarrow 0^{*}$

$$\left.\begin{array}{l}0.8\\ 1^{*}\end{array}\right\}$$

A = Diet is high

The truth value of (diet is high) = 0.7

0.7

500    Calories consumed

And $\rightarrow$ weaker condition $\rightarrow T(A \wedge B) = \min\{T(A), T(B)\}$

OR $\rightarrow$ Stronger condition. $\rightarrow T(A \vee B) = \max\{T(A), T(B)\}$

eg :- (slide)

a) $f_{\text{diet high}} = \frac{1}{5000} x = \frac{3000}{5000} = \boxed{0.6} \rightarrow$ membership for diet high for this patient = 0.6

b) $f_{\text{exercise high}} = \frac{1}{2000} x = \frac{1000}{2000} = 0.5$

**Fuzzy rules**

1) $T_{\text{with}}(\text{diet low}) \wedge \text{Truth}(\text{ex high}) = \min(0.4, 0.5) = 0.4$

2) $T(\text{diet high}) \vee T(\text{Ex low}) = \max(0.6, 0.5) = 0.6$

3) ⊘ Balanced $\Rightarrow$ Risk low = 0.4

4) Unbalanced $\Rightarrow$ Risk high = 0.6

$$T(\text{high risk}) = \frac{1}{125} x \qquad (x = \text{likelyhood of heart disease})$$

$$0.6 = \frac{x}{125} \Rightarrow \boxed{x = 75} \qquad \left| \begin{array}{l} f(\text{risk low}) = 0.8 - \frac{x}{125}\\ \boxed{x = 50} \end{array} \right.$$

intersection / And



Patient
Specific

high
risk

0.6

0.4

75

$x$

low
risk

0.4

50

$x$

Aggregated
Risk function

↳ Truth value (Risk high) 0.6

V

Truth value (Risk low) 0.4



20  40  50 75  100

$x$

likelihood of
→ heart disease

⟶ **Defuzzification :-**

likelihood of heart disease of this person ?

$$\int_{0}^{100} f(aggregated\ risk)\,dx = \int_{0}^{50} 0.4\,dx + \int_{50}^{75} \frac{x}{125}\,dx + \int_{75}^{100} 0.6\,dx$$

$$= 20 + \left[\frac{x^2}{2}\right]_{50}^{75} + 15$$

$$= \boxed{47.5\%}$$

→ Evaluate the risk

→ likelihood ⟨ defuzzyfication

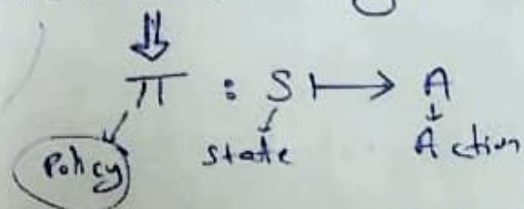\* **State machine.**

→ Bayesian classifier → Confidence/ Belief score

Softmax ⟶ entropy → uncertainity

# *MPP (Markov Pecision Process):-

→ Planning in Uncertain Environment

→ Learning ⤳ interaction

⟹ Renframent learning ● vs supervised learning
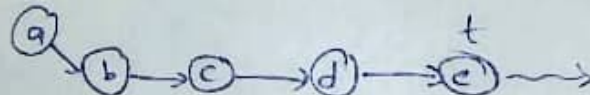
⬇

$$\pi : S \longmapsto A$$

(Policy) state     Action

→ Take action A to affect states.

→ Value function

Predicted value

↳ Can't figured out immediately.

→ State, Actions

Transition Probability :- $Pr\left(S_{t+1} \mid S'_t, A_t\right)$

* Markov Property ⓐ ⓑ→ⓒ→ⓓ—ⓔ ⤳ $^t$

⬇

Given the current S. & A, the next state is independent of all privious state & actions.

⬇

So, memory-less

Reward : $R(s)$ → real value

Find a policy : $\pi : S \longrightarrow A$ to maximize reward

maximize expected reward $E\left[r'_t \mid \pi, S_t\right]$ for all states

Constraint: The agent has "t" time step to complete the goal.

$$E\left[\sum_{k=0}^{t} r_k^i \mid \pi, S_o\right] \rightsquigarrow \text{maximize}$$

$\hookrightarrow t$ timestep

$\rightarrow$ if $t = \infty \rightarrow$ infinte time horizon.

Sooner $\downarrow$ , more reward $\uparrow$

$$E\left[\sum_{k=0}^{\infty} \gamma r_k^i \mid \pi, S_o\right] \rightarrow \text{discount factor}$$

$\gamma = 0.9$

eg :-
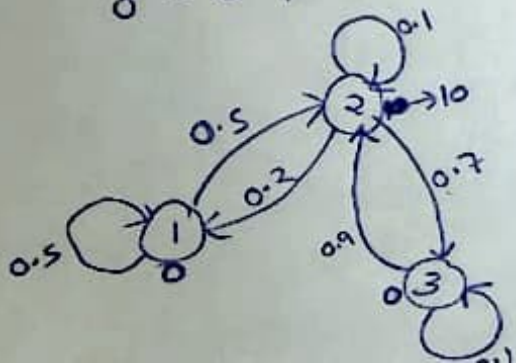
Immediate reward
+
Future reward

Value of state

$$V(s) = R(s) + \cancel{\text{}}$$

$$\gamma \sum_{s'} P(s' \mid s) V(s')$$

+RL   Markov  Chain  MDP
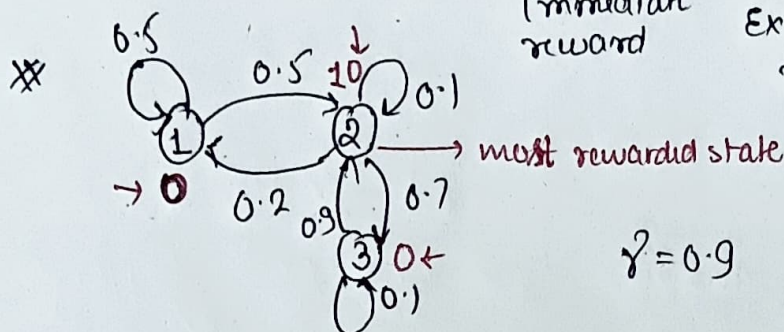
→ S
→ A
→ TP
→ R
→ V
→ γ   discount factor

Reward : $R(s)$

<u>Value of a State</u>

$V(s)$ :  How much total reward do
we expect to get if we start
from state $s$

$$V(s) = \underbrace{R'(s)}_{\substack{\text{Immediate} \\ \text{reward}}} + \underbrace{\gamma * \sum_{s'} P(s'/s) \cdot V(s')}_{\substack{\text{Expected long term} \\ \text{reward}}}$$

\#



→ most rewarded state

$\gamma = 0.9$

∴V.

Assume   $\gamma = 0.9$

* $V(1) = 0 + 0.9 \left( \underbrace{0.5 \times V(1) + 0.5 \times V(2)}_{P(s_i)} \right)$

$V(2) = 10 + 0.9 \left( 0.1 \times V(2) + 0.2 \times V(1) + 0.7 \times V(3) \right)$

$V(3) = 0 + 0.9 \left( 0.1 \times V(3) + 0.9 \times V(2) \right)$

$V(1) = 40.5$        $\boxed{V(2) = 49.5}$   $V(3) = 49.1$

Discount factor $\gamma$ determines:
the weight given to future rewards compared to
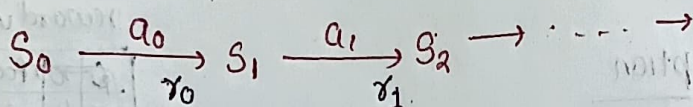current reward
(immediate)

## SUPERVISED LEARNING

Task: $\pi : S \rightarrow A$    learn from experience
Policy

State space
↑
Reinforcement    explore $\rightarrow$ RL algo
learning is ligin
learning from
interaction / experience

★ what the agent tries to most optimize?

$$S_0 \xrightarrow[r_0]{a_0} S_1 \xrightarrow[r_1]{a_1} S_2 \rightarrow \cdots \rightarrow$$

⟹ The total future discounted reward

$$S_0 \xrightarrow[r_0=5]{a_0} S_1 \xrightarrow{a_1} S_a \rightarrow \cdots$$
$a_K \rightarrow r_K = 9$
$a_j \rightarrow r_j = 12$

greedy choice

_____

$$V^\pi(S_t) = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} \cdots$$

greed    $- \rightarrow$ reward

\* what would be the optimal policy

immediate value future reward

$$\pi^*(s) = \underset{a}{argmax} \{ r(s,a) + \gamma . v^*[\gamma(s,a)] \}$$

$$\pi^*(s) = \underset{a}{argmax} \{ \_\_\_ \quad . \quad \}$$



Agent

maximum
reward if action a
is taken

reward is delayed

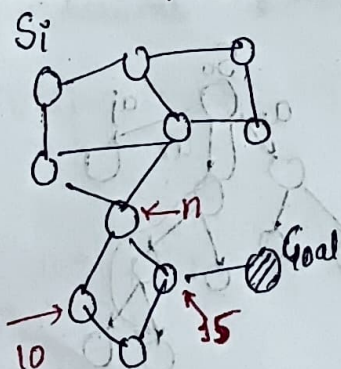| | food |
|---|---|
| | 0 | 0 |
| mouse | 0 | 1 |

## Assumption

① we know $v^*$

② reward is deterministic

③ $s_{t+1} = \gamma(s_b, a)$, transition is deterministic

## Shortest path

$s_i$



Goal

10

$\xrightarrow{}$

$35$

$\leftarrow n$

$$v(s) = \frac{1}{distance}$$

⟹ if you know $v(s)$,
the problem is trivial

immediate
↓
$r + \gamma \searrow$ ← Future

$\gamma = 0.9$

$v^*(s)$ value

**optimal:** Take highest $v^*$ of neighbouring state.

▷ Agent cannot see the complete state space

## Solving a RL problem

→ Dynamic Programming
→ Monte carlo
→ Temporal

Compute $v^\pi$ from $\pi$

↕

Improve $\pi$ based on $v^\pi$

## Bellman Equation

$$v^*(s) \leftarrow \max_a \left[ r(s,a) + \gamma v^*(\delta(s,a)) \right]$$

Immediate ↑    Future ↑

$$\text{Learn} \left( r(s,a), \delta(s,a) \right]$$

**Q-function:** Learns good state-action pair



$$10 \rightarrow Q_A = 10$$
$$\text{(S)} \quad 12 \rightarrow Q_B = 12$$
$$50 \rightarrow Q_C = 50$$

Good state action pair depend on max value of Q-function

**Q-Learning**

$$\pi^*(s) = \arg\max_a Q(s,a)$$
$$v^*(s) = \max_a Q(s,a)$$

$$\boxed{Q(s,a)}$$