# RL - Module 1 & 2

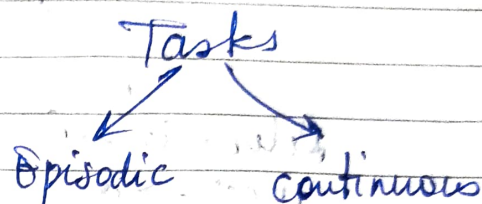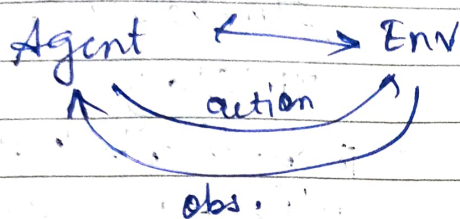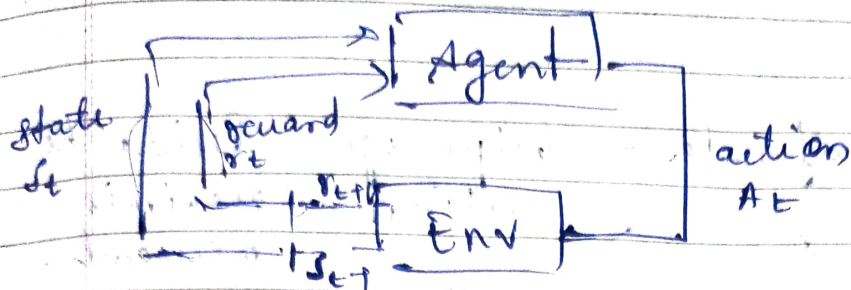**1) RL :**
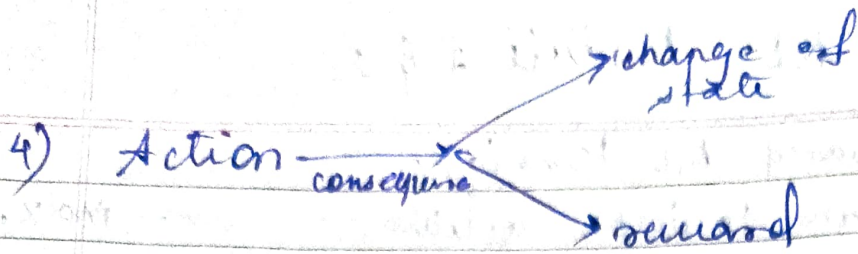- feedback based ML technique
- agent learns to take actions - env - max.
  of rewards
- agent - accomplishing - tasks - env - true/use
  rewards



Tasks → Episodic, Continuous

**2)** - State vector - list of features - helps agent take action.
   ↳ collection of relevant observations.

**3)** Objectives of RL agents -

- episodic tasks - find seq. of actions - make - majority of episodes successful

- continuous task - break into multiple episodes - find actions - maximise avg. rewards - from those episodes.

4) Action ⟶ change of state
consequence ⟶ reward

Action + consequence ⟶ knowledge Base

5) Policy

Policy — set of rules — helps agent decide the action — maximize rewards.

| Deterministic | Stochastic |
|---|---|
| $\pi(s) \to a$ | $\pi(s/a)$ |
| — given state — what action to maximize rewards? | — prob. distribute over actions for each state |
| — always produce same action | — o/p a set prob. of each possible action. |
| — outcome is fully determined | — outcome uncertain |
| — optimal coaition known with certainity | — multiple optimal outcomes. |

6) Exploitation — greedy approach — agent takes the same actions it has taken before which have proven to get some rewards

- best action in the face of current knowledge
- doesn't improve the knowledge of the agent

Exploration — agent focuse on improving the knowledge of the env.
- takes steps it has never taken before to get to know the env.
- reaps long term benefits

7) Markov State:
Markovian Assumption — current state contains all the necessary info about the past states & past actions
- current state captures info abt past state

$$P[S_{t+1}|S_t] = P[S_{t+1}|S_1,\ldots S_t]$$

8) Markov Decision Process
- stochastic decision making process
- uses mathematical framework — model decision making
- MDP evaluate which action the decision maker should take.

9) Value Function:
- Value of a state - total reward an agent
can expect to accumalate - starting from
that state.

State value functn
- $V_\pi(S)$
- Exp. returns starting
from state 's' following 'π'
till we reach terminal
state.
- $V_\pi(S) = E[R_\pi(S)]$
- How good it is to
be in a particular
state?
- evaluation done to
maximize total
rewards

Action value functn.
- $q_\pi(S,a)$
- The value of taking
an action
- Exp. returns start
from state 's'
following π & 
take action 'a'
- $q_\pi(S,a) = E[R_\pi(S,a)]$
- calculates the value
of performing an
action.

10) Optimal Policy:
$$\forall \pi: \pi \geq \pi'$$
$$\pi \geq \pi' \ \ \forall \ s: \ \ V_\pi(S) \geq V_{\pi'}(S)$$

11) Model of the env:
$P(s'\gamma | s,a)$
- describe the behaviour of env in response to
the actions

Model Based
- implicit
- model map consequnts
to altion

Model Free
- explict
- it cannot.

# 12) RL Equations

## 1) State Value function:

$$V_\pi(s) = \mathbb{E}[R_\pi(s)]$$
$$q_\pi(s,a) = \mathbb{E}[R_\pi(s,a)]$$

Value of state 's' $= R_\pi(s)$

Value ~~for each~~ of each action 'a' in state 's' following $\pi(a/s)$ $= R_\pi(s,a)$

$$\therefore R_\pi(s) = \sum_{a \sim \pi(a/s)} \pi(a/s) \, R_\pi(s,a)$$

$$\therefore V_\pi(s) = \sum_{a \sim \pi(a/s)} \pi(a/s) \, q_\pi(s,a)$$

## 2) Action Value function:

$$V_\pi(s) = \sum_a \pi(s/a) \, q_\pi(s,a)$$
$$q_\pi(s,a) = \gamma + V_\pi(s')$$

\# Discount factor $(0 < \gamma < 1)$

$$q_\pi(s,a) = \gamma + \gamma \, V_\pi(s')$$
$$= \gamma + \gamma \sum_a \pi(a/s) \, q_\pi(s',a)$$

$$\therefore q_\pi(s,a) = \sum_{s'} \sum_\gamma p(s'\gamma \mid sa)(\gamma + \gamma V_\pi(s))$$

## 13) Bellman's Eqⁿ of optimality

$$\pi^*(a^*/s) = \begin{cases} 1, & a^* = \text{argmax} \; (q^*(a,s)) \\ 0, & \text{otherwise} \end{cases}$$

$$V_\pi^*(s) = \sum_a \pi^*(a/s) \, q^*(a,s)$$

$$q^*(s,a) = \sum_{s'} \sum_\gamma p(s'\gamma \mid sa)(\gamma + \gamma V_\pi^*(s)).$$