# Reinforcement Learning

↳ taking suitable action to maximize reward in a particular situation

↳ Supervised Learning :-
- ↳ labelled data
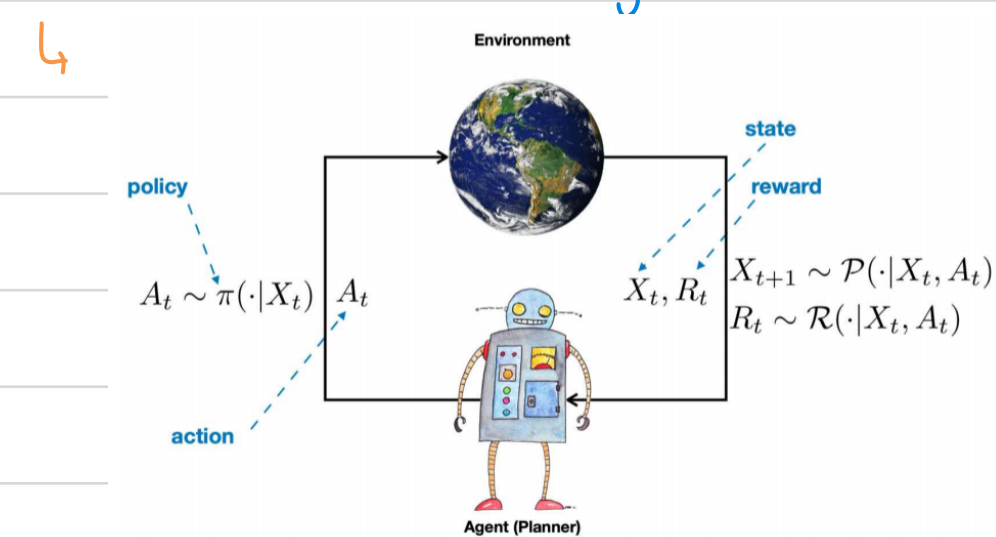  - got target variables

Unsupervised Learning :-
- ↳ labelled data
  - no target variables

Reinforcement Learning :-
- ↳ trial & error
  - reward based learning
  - .

Reinforcement Learning :-

↳



↳ Agent
- decision maker & learner

Environment
- anything outside agent that it interact w/ & attempt to control

State, $X_t$
- variable that summarize whatever has happened to the agent so far

Action, $A_t$
- set of all possible action being on a state, $X_t$

Policy, $\pi$
- indicates action, $A_t$ to be taken on a state, $X_t$
- usually a mapping from state to actions

Transition Probability Kernel
- Describe the dynamic
- Give action's effect in a state

Reward, $R_t$
- reward for simply being in the state, $X$
- a real number specifying the immediate desirability of an action in a state

↳ example    CHESS    ROBOT HAND

| | CHESS | ROBOT HAND |
|---|---|---|
| State, $X_t$ | Position of pieces on board | Position of Ball '' Hand coordinate finger |
| Action, $A_t$ | Movement of a piece | Δ/ past coord. of hand |
| Reward, $R_t$ | +1 ⇒ checkmate +0.5 ⇒ capture a piece +0 ⇒ nothing happen | +1 closer to ball −1 further to ball |
| Policy, $\pi$ | Prob of taking certain move. train NN to output possible move | llow of moving robot hand |

↳ Sequence    an episode

$X_1 , A_1 , R_1 , X_2 , A_2 , R_2 , \dots$
state  action  reward  ...
t=1    t=1    t=1

- terminate after certain number of $t$
- terminate after agent reach certain state
- never terminate

# Markov Decision Process

↳ framework for decision making where outlines are partly random & partly under control of decision maker

↳ 5 state :-
$\{ X , A , P , R , \gamma \}$
state  action  policy  reward  discount

↳ Policy, $\rho$
- state transition kernel
- state action transition kernel

↳ $\bar{\pi} = \{ \pi_1, \pi_2, \pi_3 \}$

policy time t=1

↳ $\bar{\pi}_t ( a_t | X_1, A_1, X_2, A_2, \dots, X_t )$
given
$\underbrace{\qquad\qquad\qquad}_{\text{history } X, A}$

conditional probability distribution given the history states selection

↳ $\pi_t ( A | X_1, A_1, X_2, A_2, \dots, X_t ) = 1$

policy prob. dist. over actions given history sums to 1

↳ Type of Policy
- ↳ 1. Markov Policy (state trans. kernel) depend on $X_t$, $\pi(\cdot | X_t)$
- 2. Deterministic Policy (state trans. kernel) assign mass 1 to action for each state, $X_t$, $\pi_t ( A_t | X_t ) = 1$
- 3. Stochastic Policy → can change (state trans. kernel) assign prob. dist over a given state $X_t$
- 4. Stationary Policy (both) policy dont change over time
  $\pi = \{ \pi, \pi, \pi, \dots, \pi \}$

# Reward

↳ $R_t \sim R ( \cdot | X_t, A_t )$
immediate reward

↳ $r(x, a) = \mathbb{E} ( R | X = x, A = a )$
↳ ave. reward that repeated interaction w/ environment

- receive within one episode as a measure of performance
- a way to maximize long term reward

# Task

↳ ① Finite Horizon Task
- interact for fixed predefined times
- eg. solve a maze w/ fixed steps
- agent goal is to maximize cummulative rewards
- ↳ $G^\pi \triangleq R_1 + R_2 + \dots + R_t$  } sum of rewards  } return function
  $G^\pi \triangleq R_1 + \gamma R_2 + \dots + \gamma^{t-1} R_t$  } discounted sum
  $V^\pi = \mathbb{E} \left[ \sum_{t=1}^{T} \gamma^{t-1} R_t | X_1 = x \right]$  } value function of $\pi$

↳ ② Episodic
- well defined start & end
- terminate when reach a certain state
- eg. chess
- ↳ $G^\pi \triangleq \sum_{t=1}^{T} \gamma^{t-1} R_t$  } return function
  $V^\pi(x) \triangleq \mathbb{E} ( G^\pi | X_1 = t )$  } value function

↳ ③ Continuity
- no endpoint
- eg. control robot for navigation
- ↳ $V^\pi(x) \triangleq \mathbb{E} \left[ \sum_{t=1}^{\infty} \gamma^{t-1} R_t | X_1 = x \right]$  } value function
  $Q^\pi(x) \triangleq \mathbb{E} \left[ \sum_{t=1}^{\infty} \gamma^{t-1} R_t | X_1 = x, A_1 = a \right]$  } action-value function