→ DP is used to solve MDP's and similar problems.
DP aim to find an optimal policy

# Dynamic Programming

- DP is a method of solving complex prob. by breaking them down into sub problems. The solutions to the sub-problems are combined to solve overall problem. The basic assumption is that the model of the env is given.

- The two req. properties of DP are —

1) Optimal substructure —
   optimal sol^n of the sub-problem can be used to solve the overall problem.

2) Overlapping sub-problems —
   sub-problems occur many times. sol^n can be cached & reused.

- Two popular DP algos are —
  a) Policy iteration        } to find optimal
  b) Value iteration.        } policy of a MDP

## POLICY ITERATION:

- Includes iteratively improving a policy until it converges to the optimal policy.
- After initit Initializing a random policy, this method iterates bet^n 2 steps
  - Policy evaluation
  - Policy improvement

(P.T.O.)

$$V_\pi(s) = E[R_\pi(s)]$$

$$q_\pi(s) = r + \gamma E(q_{\pi+1}, a)$$
$$q_\pi(s) = r + \gamma V_\pi(s')$$

## Policy Evaluation :

i) Evaluate the current policy by computing evaluating the expected value of rewards obtained by following the policy.

$$E[R_\pi(s)]$$

ii) Measures how good the policy is by calculating all the state value function $V_\pi(s)$ for all states until the state value function is converged.
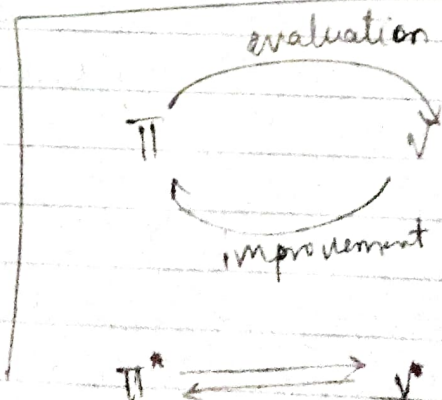
iii) Also known as the prediction problem.

## Policy Improvement :

i) Updates the policy to be greedy w.r.t. current q-value function

ii) This means the policy chooses the action with max that maximises the expected reward from the current state $q_\pi(s,a)$ for each state

$q_\pi(s,a)$

$$\pi'(a|s) = \begin{cases} 1, & a = \text{argmax}_a \left[ \sum_{s'} \sum_r P(s',s|s,a)[r + \gamma V_\pi(s')] \right] \\ 0, & \text{otherwise} \end{cases}$$

iii) Control Problem



evaluation

$\pi$

improvement

$$\pi^* \longleftrightarrow V^*$$

## VALUE ITERATION:

- simpler than policy iteration
- only performs one step of policy evaluation in each iteration.
- the optimal value funch is defined
- It is an iterative algo that starts with an arbitary value funtn and updates it until it converge to optimal value funtn.
- optimal value funtn... max. expected reward of obtained from each state in the given policy.

$$V(s) = \max_a \left[ \sum_s, \sum_\partial P(s', s|s, a)[\partial + \gamma v_\bullet(s)] \right]$$

s.t.,

$V(s)$ - current estimate of the optimal value funtn for state s

$P(s', s|s, a)$ - model of the equam.

$\gamma$ - discount fater

$V(s')$ =

|  | (PI) Policy Iteration | (VI) Value Iteration |
|---|---|---|
| strength | i) converge to optimal policy faster than VI in same cases | i) converges to optimal value funth faster than PI |
| | ii) Can handle larger class of MDP's including those with stochastic policies | ii) Can handle MDP's with infinite horizon |
| | iii) More stable & less sensitive to discount factor | iii) Simpler & easy to implement. |
| Weakness | i) Can be computationally expensive | i) sensitive to the choice of discount factor |
| | ii) can get stuck in local optima. | ii) cannot handle stochastic values directly |
| | iii) May not converge if the MDP has infinte horizon | iii) May require a large no. of iterations to converge |