# MARKOV DECISION PROCESS

Agent

$\eta \in$ Reward

$a \in$ Action

$s \in$ State

Environment
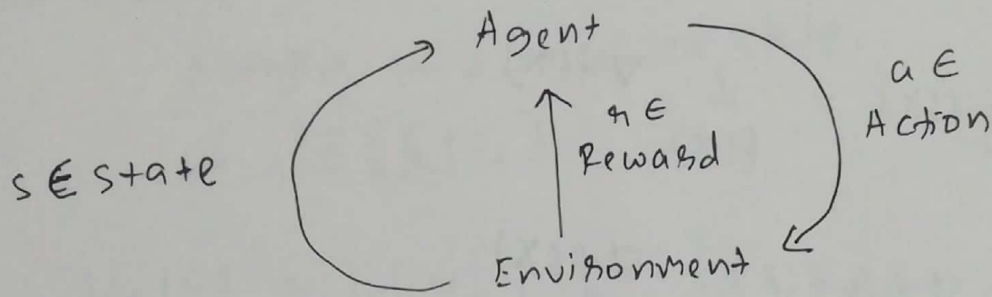
$P(S_{t+1} | S_t, a_t)$  : Formalization

$W_o \leftarrow$ random

loop:
    sample N Sessions
    elite sessions

$$W_{i+1} = W_i + \alpha \nabla \left[ \sum_{S_i, a_i \in Elite} \log \pi W_i (a_i | s_i) \right]$$

RNN:

$$S_t^a = \sigma \left( S_{t-1}^a W_S + 0_t W_o \right)$$

Expected Reward

$$J = \int N(\theta | \mu, \sigma^2) \int P(\tau | \theta) R(\tau) \, d\tau \, d\theta$$

$\theta = $ State      $\tau = $ Action

maximizing $J = \nabla J = \dfrac{\partial J}{\partial ?}$

So, $\boxed{\nabla J = \dfrac{\partial J}{\partial \mu \partial \sigma^2}}$

$$\nabla J = \int \nabla \left[ N(\theta \mid \mu, \sigma^2) \right] \int P(\tau \mid \theta) R(\tau) d\tau d\theta$$

$$\nabla \log f(x) = \frac{1}{f(x)} \nabla f(x)$$

$$f(x) \nabla \log f(x) = \nabla f(x)$$

$$\nabla N(\theta \mid \mu, \sigma^2) = N(\theta \mid \mu, \sigma^2) \nabla \log \left( N(\theta \mid \mu, \sigma^2) \right)$$

$$\nabla J = \int N(\theta \mid \mu, \sigma^2) \nabla \log \left( N(\theta \mid \mu, \sigma^2) \right) \int P(\tau \mid \theta) R(\tau) d\tau d\theta$$

$$\nabla J \cong \frac{1}{N} \sum_{i=0}^{N} \nabla \log N(\theta \mid \mu, \sigma^2) \sum R(s, a \ldots)$$

Strategy for optimization:

① Guess initial $M_0, \sigma_0^2$

② Forever run:
$$\nabla J \approx \frac{1}{N} \sum_{i=0}^{N} \nabla \log \left( N(\theta \mid \mu, \sigma^2) \right) \sum R(s, a, s' \ldots)$$
$$s, a, s' \ldots \in \tau_i$$

updates:
$$\mu = \mu + \alpha \nabla_\mu J \qquad \sigma^2 = \sigma^2 + \alpha \nabla_{\sigma^2} J$$

$$\underset{N(\theta \mid \mu, \sigma^2)}{\text{argmax}} E \qquad R = \underset{N(\theta \mid \mu, \sigma^2)}{\text{argmax}} E \qquad \frac{R - mean}{std}$$

zero reward $\qquad A = \dfrac{R - E(R)}{var(R)}$

# Bellman equations:

state value function $v(s)$

return conditional on state.

$$E[G_t | S_t]$$

$$V_\pi(s) = E[G_t | S_t = s]$$
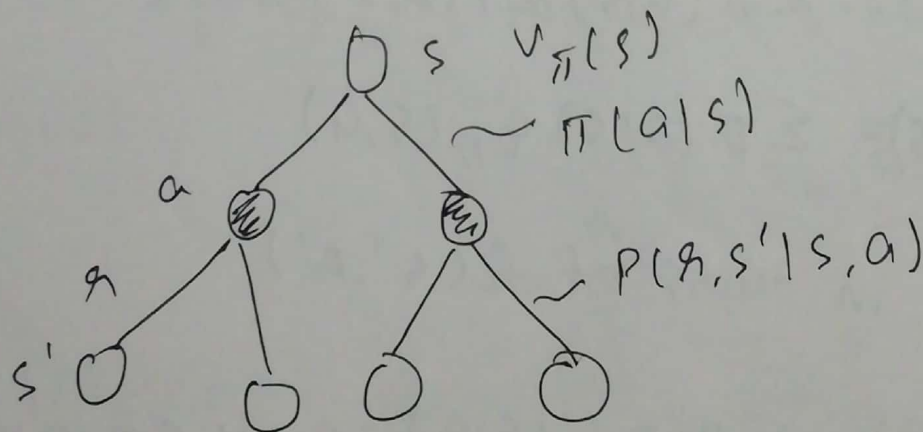
$$= E_\pi [R_t + \gamma G_{t+1} | S_t = s]$$

Stochasticity in Policy and environment

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{\mathcal{R},s'} P(\mathcal{R}, s' | s, a) [\mathcal{R} + \gamma E_\pi[G_{t+1} | S_{t+1} = s']]$$

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{\mathcal{R},s'} P(\mathcal{R}, s' | s, a) [\mathcal{R} + \gamma V_\pi(s')]$$

Bellman expectation equation for $v(s)$

$$V_\pi(s) = E_\pi [R_t + \gamma V_\pi(S_{t+1}) | S_t = s]$$



$s \quad V_\pi(s)$

$\pi(a|s)$

$P(\mathcal{R}, s' | s, a)$

Backup Diagram

Action value function:

expected return conditional on state & action.

$$q_\pi(s,a) = E_\pi[G_t \mid S_t = s, A_t = a]$$

$$\uparrow$$

no policy stochasticity at first step.

$$q_\pi(s,a) = E_\pi[R_t + \gamma G_{t+1} \mid S_t = s, A_t = a]$$

$$q_\pi(s,a) = \sum P(r, s' \mid s, a)[r + \gamma E[G_{t+1} \mid S_{t+1} = s']]$$

$$q_\pi(s,a) = \sum_{r,s'} P(r, s' \mid s, a)[r + \gamma V_\pi(S_{t+1})]$$

$$q_\pi(s,a) = \sum_{r,s'} P(r, s' \mid s, a)[r + \gamma V_\pi(s')]$$

$V_\pi$ in terms of $q_\pi$

$$V_\pi(s) = \sum \pi(a \mid s) \sum P(r, s' \mid s, a)[r + \gamma V(s')]$$

$$V_\pi(s) = \sum \pi(a \mid s) q_\pi(s,a)$$

$q(s,a)$ in terms of $q(s', a')$

$$q_\pi(s,a) = \sum P(r, s' \mid s, a)[r + \gamma \sum \pi(a' \mid s') q_\pi(s', a')]$$
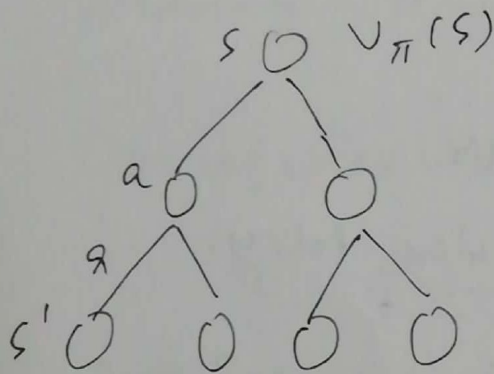
Comparission of Policies: $\pi$, $\pi'$

$$\pi \geq \pi' \quad \text{if} \quad V_\pi(s) \geq V_{\pi'}(s) \; \forall s$$
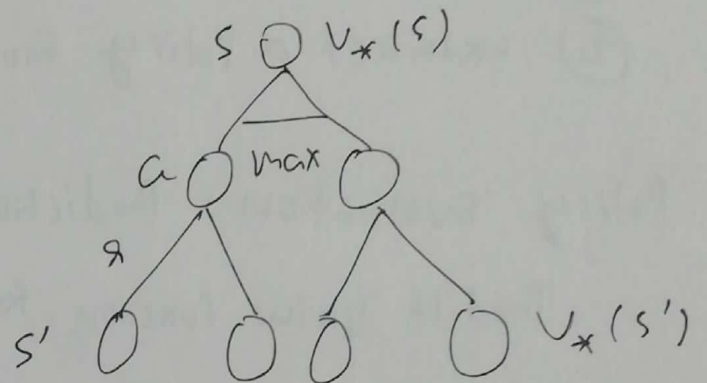
In an MDP atleast 1 $\pi$ exist satisfying above condition. i.e

In any finite MDP There is always atleast one deterministic optimal policy.

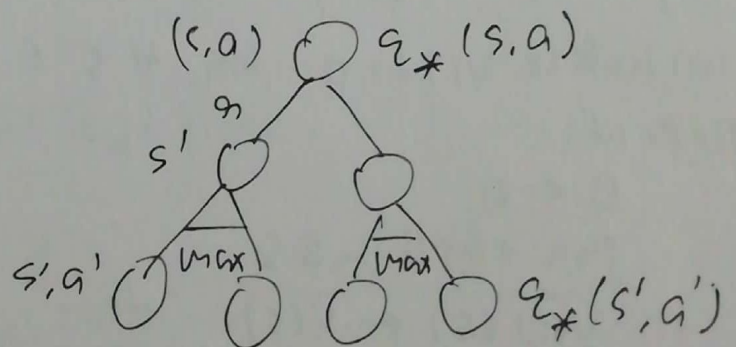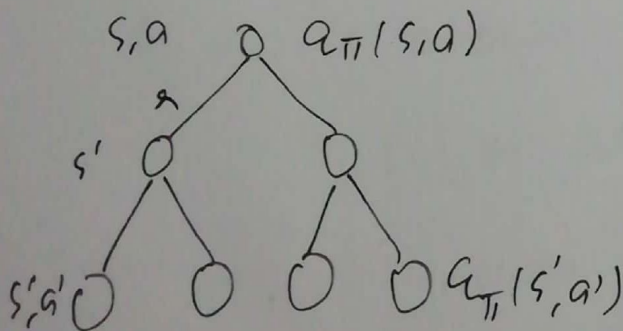$$V_*(s) = \max_\pi V_\pi(s) \qquad q_*(s,a) = \max_\pi q_\pi(s,a)$$



Bellman expectation equation



Bellman optimality equation for $V_*(s)$

$$V_*(s) = \max_a \sum_{r,s'} P(r,s'|s,a)[r + \gamma V_*(s')]$$

$$= \max_a E_\pi[R_t + \gamma V_*(S_{t+1}) | S_t = S, A_t = a]$$





$$q_*(s,a) = E_\pi[R_t + \gamma \max_{a'} q_*(S_{t+1}, a') | S_t = S, A_t = a]$$

$$= \sum_{r,s'} P(r,s'|s,a)[r + \gamma \max_{a'} q_*(s',a')]$$

How to find optimal Policy?

Model & value based.

Model based setup:

Model of world is known ie

$P(r, s' | s, a)$ for all $r, s', s, a$ is known.

Value based setup:

① Build or estimate a value

② extract a Policy from value.

Policy Evaluation : Prediction Problem.

Predict value function for a Particular Policy.

$$V_\pi(s) = \sum \pi(a|s) \sum P(r, s' | s, a) [r + \gamma V_\pi(s')]$$

$$= \mathbb{E}_\pi [R_t + \gamma V_\pi(s_{t+1}) | S_t = s]$$

Algorithm iterative:

Input $\pi$, Policy to be evaluated.
Initialize $V(s) = 0$ for $\forall s \in S$.
Repeat:
 $\Delta \leftarrow 0$
 For each $s \in S$
  $V_{old}(s) \leftarrow V(s)$
  $V(s) = \sum \pi(a|s) \cdot \sum P(r, s' | s, a) [r + \gamma V_\pi(s')]$
  $\Delta \leftarrow max(\Delta, |V_{old}(s) - V(s)|)$
untill $\Delta < \theta$ (a small Positive no')

Policy improvement:

$$\pi'(s) \leftarrow \underset{a}{\text{argmax}} \sum P(\pi, s'|s, a)[\pi + \gamma \underset{\pi}{V(s')}]$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{q_\pi(s, a)}$$

This procedure is guaranteed to produce better policy.

If $q_\pi(s, \pi'(s)) \geq v_\pi(s)$ for all states

Then $V_{\pi'}(s) \geq V_\pi(s)$

Meaning $\pi' \geq \pi$

Convergence:

$$\pi' = \pi \quad \rightarrow \quad V_{\pi'} = V_\pi$$

Then it is optimal.

$$V_{\pi'}(s) = \underset{a}{\text{max}} \sum P(\pi, s'|s, a)[\pi + \gamma V_\pi(s')]$$

Determining optimal policy from $V_*(s)$, $q_*(s, a)$

If $q_*$ is known

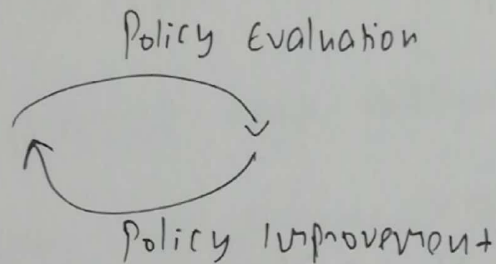$$\pi_*(s) \leftarrow \underset{a}{\text{argmax}} \, q_*(s, a)$$

If $V_*$ is known

$$q_*(s, a)$$

$$\pi_*(s) \leftarrow \underset{a}{\text{argmax}} \sum P(\pi, s|s, a)[\pi + \gamma \underset{*}{V(s')}]$$

In Model free setup we don't know transition probabilities hence we can't convert $V_*(s)$ to $q_*(s)$ so we can't extract $\pi$ from $V_*(s)$ so we can extract it using $q_*(s)$ only.

& Why Precise solution to Bellman equation is not needed and approximation is enough ?

Generalized Policy Iteration: Policy & Value iteration.

Policy Evaluation

Policy Improvement

① Evaluate given Policy
② Improve policy by acting greedily wrt value function.

Policy Iteration:
① Evaluate Policy untill convergence (with some tolerance)
② Improve policy.

Value Iteration:
① Evaluate Policy only with single iteration.
② Improve policy.

Policy iteration : scheme

1. Initialize $V(s)$ and $\pi(s)$ $\forall s \in S$.
2. Perform Policy evaluation (without internal initialization).
3. Policy improvement.

Pseudo code :

Policy_stable $\leftarrow$ TRUE

for each $s \in S$ :

old_action $\leftarrow \pi(s)$

$$\pi(s) \leftarrow \underset{a}{\text{argmax}} \sum P(r, s'|s, a)[r + \gamma V(s')]$$

$q(s, a)$

If old_action $\neq \pi(s)$ then Policy_stable $\leftarrow$ false.

If Policy_stable then stop, return $V \approx V_*$, $\pi \approx \pi_*$;

else go to 2.

Value ituation : Scheme.

Initialize $V$ arbitrarily $V(s) = 0$ for all $s \in S$.

Repeat

$\Delta \leftarrow 0$

for each $s \in S$.

$V_{old} \leftarrow V(s)$

$q(s, a)$

{Bellman optimality equation}
$$V(s) \leftarrow \underset{a}{\text{max}} \sum P(r, s'|s, a)\{r + \gamma V(s')\}$$

$$\Delta \leftarrow \text{max}(\Delta, |V_{old}(s) - V(s)|)$$

until $\Delta < \theta$

output deterministic policy $\pi \approx \pi_*$

$$\pi(s) = \underset{a}{\text{argmax}} \sum_{s', r} P(r, s'|s, a)[r + \gamma V(s')]$$

inturvidiate $V(s)$ in VI might not correspond to any real policy as its just approximation towards evaluation.

Value Ituation (VI) vis Policy Ituation (PI)

VI is fast pur cycle ; PI is slower pur cycle

$$O(|A||S|^2) \qquad ; \qquad O(|A||S|^2 + |S|^3)$$

VI requires many cycle. PI req. few cycle.

Model free Policies:

① Monte - Carlo :

Get all trajectories containing particular $(s, a)$
Estimate $G(s, a)$ for each trajectory.
Average num to get expectation

# Asynchronous Dynamic programming:

in-place iterative DP algo's.

It must compute value of all the states.

Policy evaluation and Policy improvement is interleaved.

like:

loop:

$\qquad D \leftarrow 0$

$\qquad$ loop: ~~for s in S~~ Until every state is evaluated.

$\qquad\qquad$ V.old $\leftarrow$ V(s)

$\qquad\qquad$ $V(s) \leftarrow \max_a \sum_{r,s'} P(r,s'|s,a)[r + \gamma V(s')]$

$\qquad\qquad$ $\Pi(s) = \arg\max_a \sum P(r,s'|s,a)[r + \gamma V(s')]$

$\qquad\qquad$ ~~$\Delta \leftarrow \max(\Delta, |v-V(s)|)$~~

$\qquad$ ~~$\leftarrow$ until $\Delta < \theta$~~

# Generalized Policy Iteration:



$V = V_\Pi$

$V_*, \Pi_*$

$\Pi = Greedy(V)$