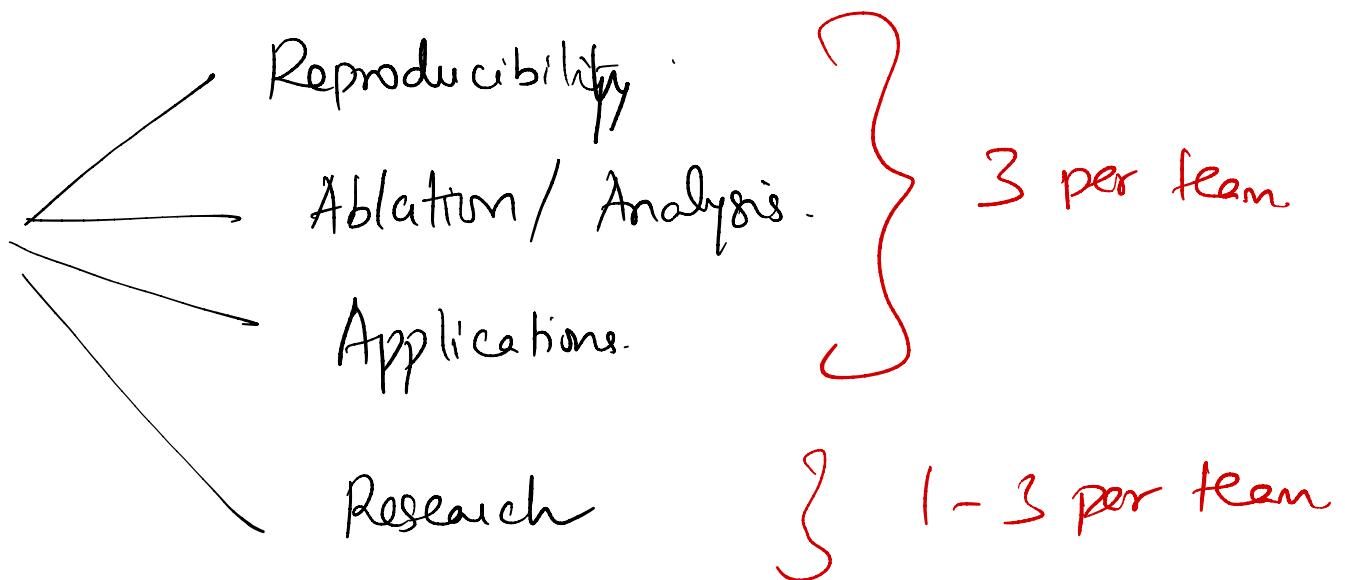
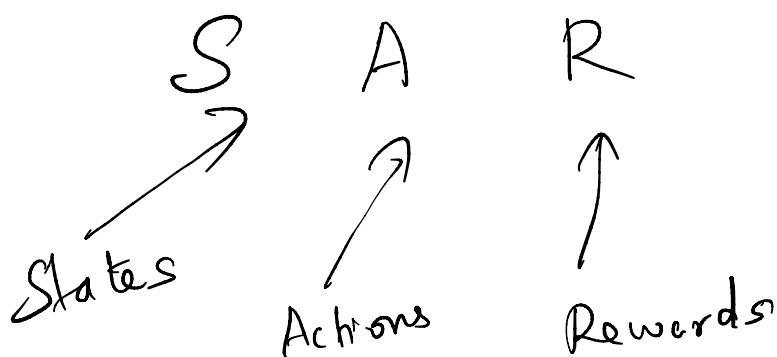


Lecture - 03



Finite Markov Decision Process



$$P(s', r | s, a)$$

dynamics
fn.

Markov Property)

$$P(S_t, R_t | S_{t-1}, A_{t-1})$$

Return

episodic problems.

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

continuing tasks

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$0 \leq \gamma \leq 1$ is called the discount rate

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots$$

$$= R_{t+1} + \overline{\gamma (R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots)}$$

$$G_t = \boxed{R_{t+1}} + \boxed{\gamma G_{t+1}}$$

Note: In Continous tasks,

Sum of ~~rewards~~ is infinite number
of rewards,

The sum is still finite.

Ex: Constant + 1 reward.

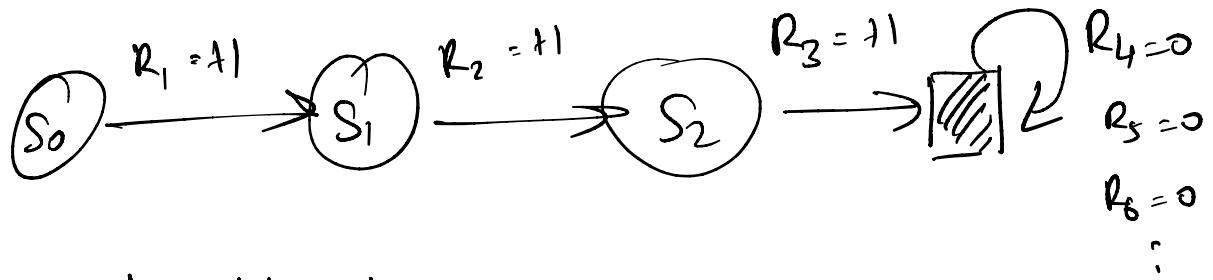
$$G_t = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$$

episodic

Continuing.

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T$$

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$



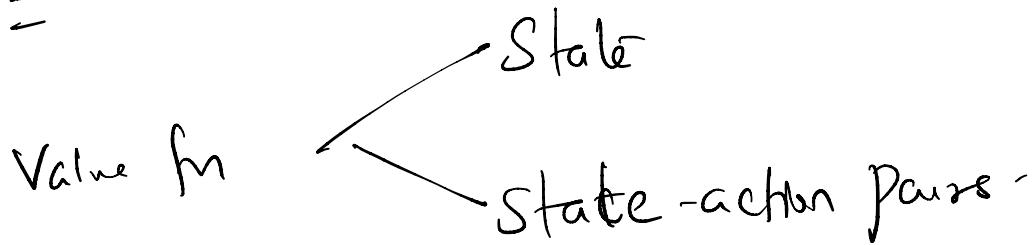
$$G_t = \sum_{k=t+1}^T \gamma^{k-t-1} R_k$$

$$\underline{T = \infty} \quad \text{or} \quad \underline{\gamma = 1}$$

MDP, dynamics, returns

Policies and Value fn's:-

Goal: estimate value for a_1



Policy: mapping from states to
prob. of selecting each possible action.

π (als)

State value fn:

$V_{\pi}(s) =$ Value fn of a state s
under a policy π .

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t \mid s_t = s]$$

$$= \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid s_t = s \right]$$

$\forall s \in S$

Note:

$$V_{\pi} (\text{Terminal state}) = 0.$$

Action-value fn:

$q_{\pi}(s, a)$ = value of taking action a
in state s .

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t \mid S_t = s]$$

$$q_{\pi}(s, a) = \mathbb{E}_{\pi} [G_t \mid S_t = s, A_t = a]$$

$$= \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$$

Bellman equation for value fns:-

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t \mid S_t = s]$$

$$= \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} \mid S_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a)$$

$$[\gamma + \gamma \mathbb{E}_{\pi} [G_{t+1} \mid S_{t+1} = s']]$$

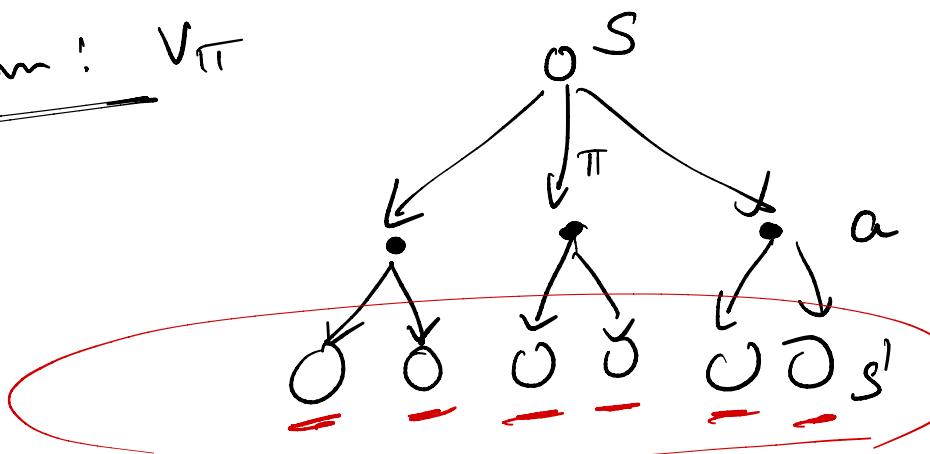
Bellman eqnath.

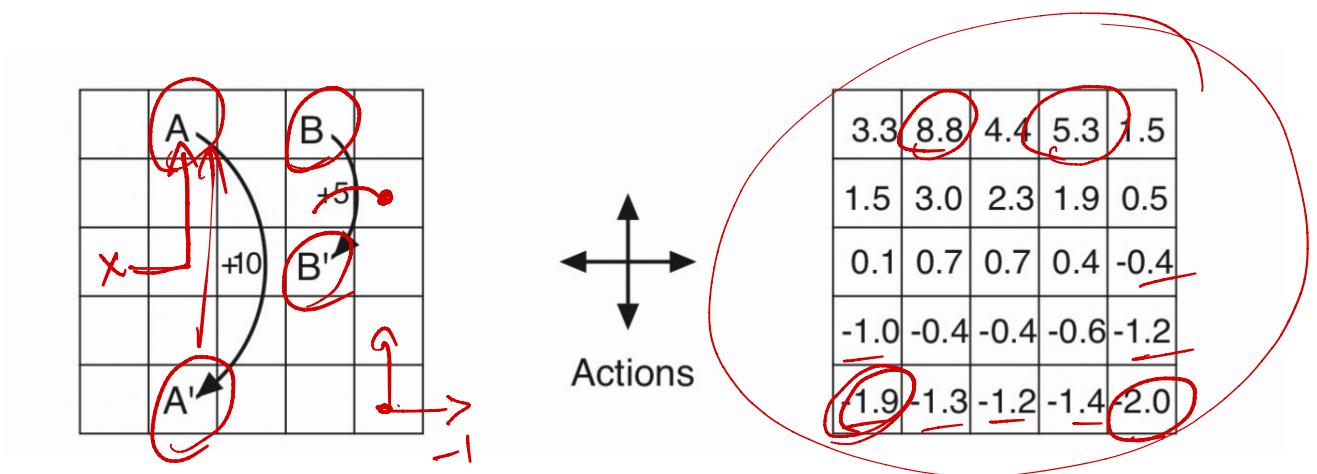
$$V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi}(s')]$$

$\forall s \in S$

unique soln: V_{π^*}

backup diagram: V_{π}

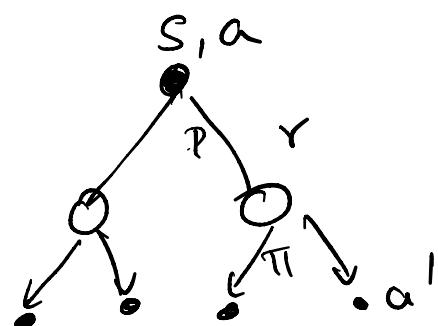




equiprobable random
Policy.

Bellman eqn for $q_{\pi}(s, a)$:

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



Optimal Policies / Optimal value fn.

$\pi \geq \pi'$ if and only if $V_\pi(s) \geq V_{\pi'}(s)$

$\forall s \in S$

Optimal Policy: π^* .
'n' states
'k' actions

V^* — optimal state value fn. n^k .

$$V^*(s) = \max_{\pi} V_{\pi}(s) \quad \forall s \in S$$

$$q^*_*(s, a) = \max_{\pi} q_{\pi}(s, a) \quad \forall s \in S \\ a \in A$$

$$q^*_*(s, a) = \mathbb{E} \left[R_{t+1} + \gamma V^*(s_{t+1}) \mid S_t = s, A_t = a \right]$$

$$V^*(s) = \max_{a \in A} q_{\pi_*}(s, a)$$

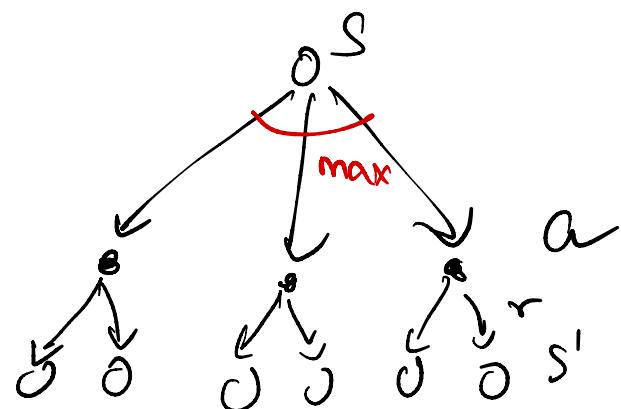
$$= \max_{a \in A} E_{\pi_*} [G_t | S_t = s, A_t = a]$$

Bellman optimality eqn for V_x

$$= \max_a E_{\pi_*} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a]$$

$$V_x(s) = \max_a E [R_{t+1} + \gamma V_x(s_{t+1}) | S_t = s, A_t = a]$$

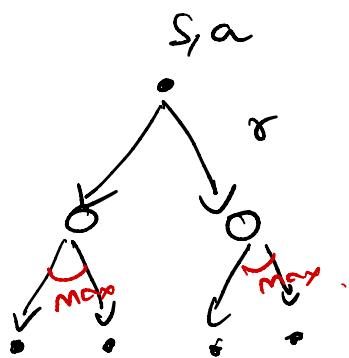
$$V_x(s) = \max_a \sum_{s', r} P(s', r | s, a) [r + \gamma V_x(s')]$$

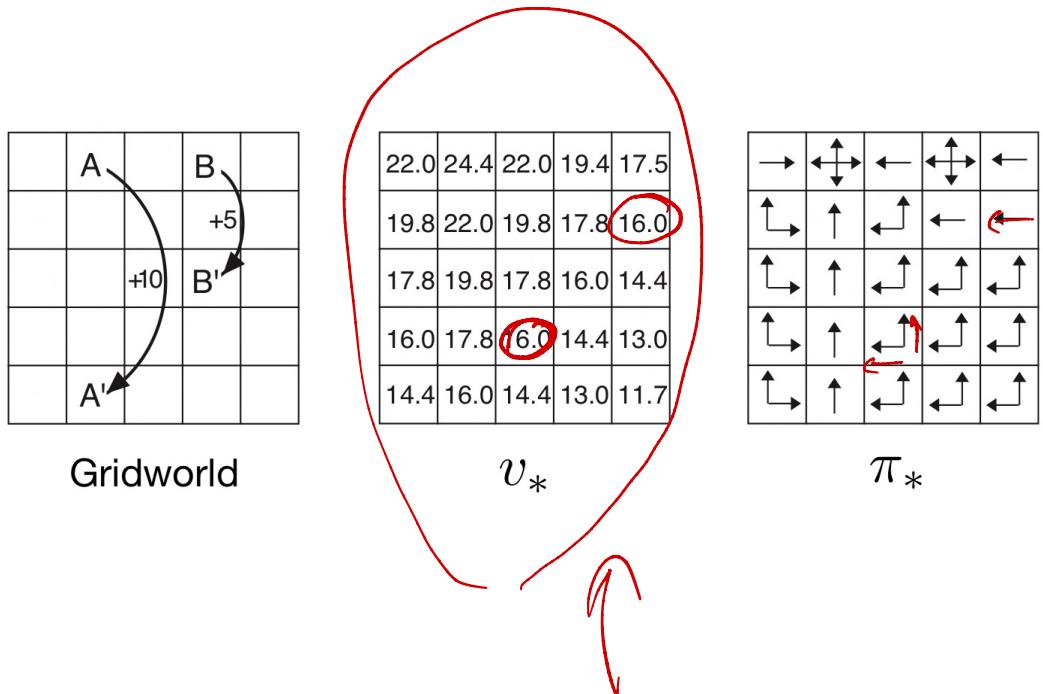


$$q_*(s, a) = \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_*(s_{t+1}, a') \right]$$

$s_t = s, A_t = a$

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





Strong assumption:

- (1) we accurately know the dynamics.
- (2) we have enough compute to complete the s.t.
- (3) Markov property -

Dynamic Programming

Given an MDP, compute the optimal policy.

Finite MDP

$$P(s', r | s, a)$$

~~Key idea:~~ Use the value functions to organize and structure the search for good policies

Bellman Optimality equations - ?

$$\begin{aligned} V_x(s) &= \max_a \mathbb{E} \left[R_{t+1} + \gamma V_x(s_{t+1}) \mid s_t = s, \right. \\ &\quad \left. A_t = a \right] \\ &= \max_a \sum_{s', r'} P(s', r' | s, a) [r + \gamma V_x(s')] \end{aligned}$$

$$\begin{aligned} q_x(s, a) &= \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_x(s_{t+1}, a') \mid s_t = s, \right. \\ &\quad \left. A_t = a \right] \\ &= \sum_{s', r'} P(s', r' | s, a) [\gamma + \gamma \max_{a'} q_x(s', a')] \end{aligned}$$

given π , how to estimate $V_\pi(s)$?

given $V_\pi(s)$, how to find a better policy π' ?

Policy evaluation (Prediction)

Policy improvement

Policy evaluation 1-

Given a arbitrary policy

π , how to compute V_π ?

$$V_\pi(s) = E[G_t | s_t = s]$$

$$= E_\pi[R_{t+1} + \gamma G_{t+1} | s_t = s]$$

$$= E_\pi[R_{t+1} + \gamma V_\pi(s_{t+1}) | s_t = s]$$

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

4.

Iterative sol:

$$v_0, v_1, v_2, v_3, \dots$$

Expected update

$$V_{k+1}(s) = \sum_a \pi(a|s) \sum_{s', r} P(s'|r|s, a) [r + \gamma V_k(s')]$$

$$V_k(s) = V_{k+1}(s) \quad V_0 \quad (\quad \quad \quad)$$

$$V_k = V_{k+1} \quad (\quad \quad \quad)$$

"sweep" \longleftrightarrow iteration.

$$V_k \rightarrow V_{k+1} \quad \text{when } k \rightarrow \infty$$

if $V_{k+1}(s) = V_k(s) \forall s$.

$$V_k = V_{k+1} = V_\pi$$

V_π - is the fixed pt for thus update rule.

$$\{V_k\} \rightarrow V_\pi$$

as $k \rightarrow \infty$.

Iterative Policy evaluation

$$V_k(s) :$$

$$V_{k+1}(s) :$$

Iterative Policy Evaluation, for estimating $V \approx v_\pi$

Input π , the policy to be evaluated

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in \mathcal{S}$:

(Sweep)

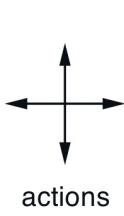
$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

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

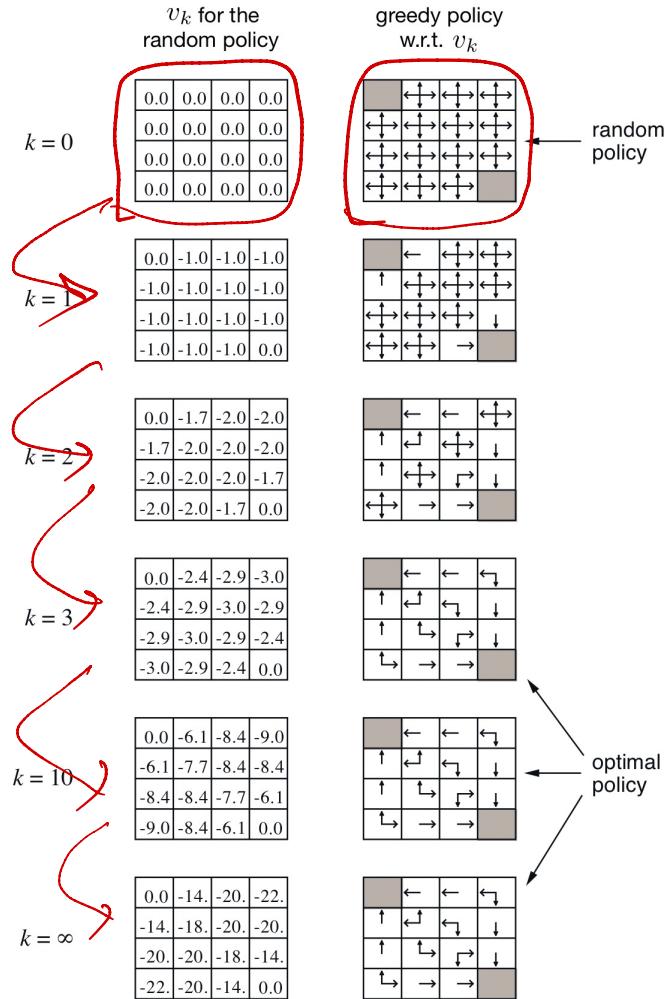
until $\Delta < \theta$

$$\theta = 0.1$$



	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

$R_t = -1$
on all transitions



Policy improvement :-

Compute Val. fns for a policy to help find better policies

π, V_π

's' $a = \pi(s)$ or $a \neq \pi(s)$.

Consider selecting 'a' in 's' and thereafter following the existing poly π .

$$\begin{aligned} q_{\pi}(s, a) &= \mathbb{E} \left[R_{t+1} + \gamma V_{\pi}(s_{t+1}) \mid S_t = s, A_t = a \right] \\ &= \sum_{s', r} P(s', r | s, a) [r + \gamma V_{\pi}(s')] \end{aligned}$$

if $q_{\pi}(s, a) > V_{\pi}(s)$

Policy improvement theorem :- (PIT)

Let π and π' be a pair of deterministic policies such that

If $s \in S$

$$q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s)$$

then the policy π' must be as good as or better than π .

$$\text{i.e. } v_{\pi'}(s) \geq v_{\pi}(s) \forall s \in S$$

Proof:

$$V_{\pi}(s) \leq q_{\pi}(s, \pi'(s))$$

$$q = \mathbb{E}[R_{t+1} + \gamma V_{\pi}(s_{t+1}) \mid s_t = s, A_t = \pi(s)]$$

$$= E_{\pi}, [R_{t+1} + \gamma V_{\pi}(s_{t+1}) \mid s_t = s]$$

$$\leq E_{\pi}, [R_{t+1} + \gamma \underbrace{q_{\pi}(s_{t+1}, \pi'(s_{t+1}))}_{\text{red arrow}} \mid s_t = s]$$

$$= E_{\pi}, [R_{t+1} + \gamma E_{\pi}, [R_{t+2} + \gamma V_{\pi}(s_{t+2}) \mid s_{t+1} \\ A_{t+1} = \pi(s_{t+1})]]$$

$$= E_{\pi}, [R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{\pi}(s_{t+1})] \begin{matrix} s_{t+2} \\ s_{t+3} \end{matrix}$$

$$\leq E_{\pi}, [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} \\ + \gamma^3 V_{\pi}(s_{t+3})] \begin{matrix} \dots \\ s_{t+3} \end{matrix}$$

$$\leq E_{\pi}, [R_{t+1} + \gamma R_{t+2} + \dots]$$

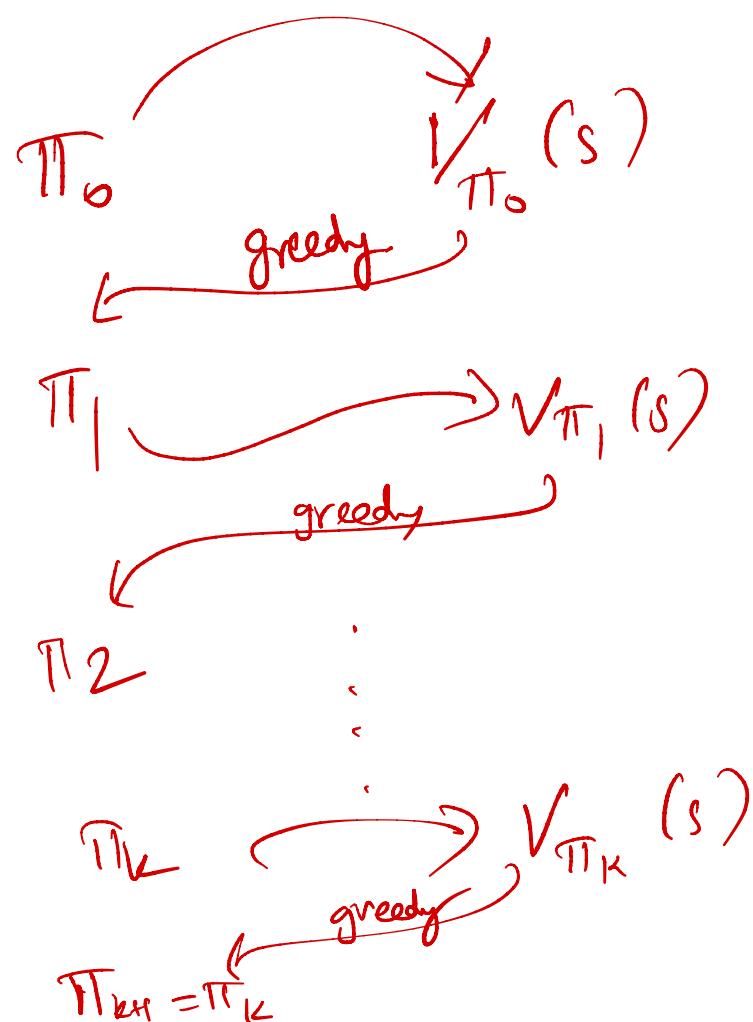
$$= V_{\pi'}(s)$$

Given a value fn. for some Policy π ,

new policy π' / wrt. val. f.

Greedy policy.

$$\pi'(s) = \arg \max_a q_{\pi}(s, a)$$

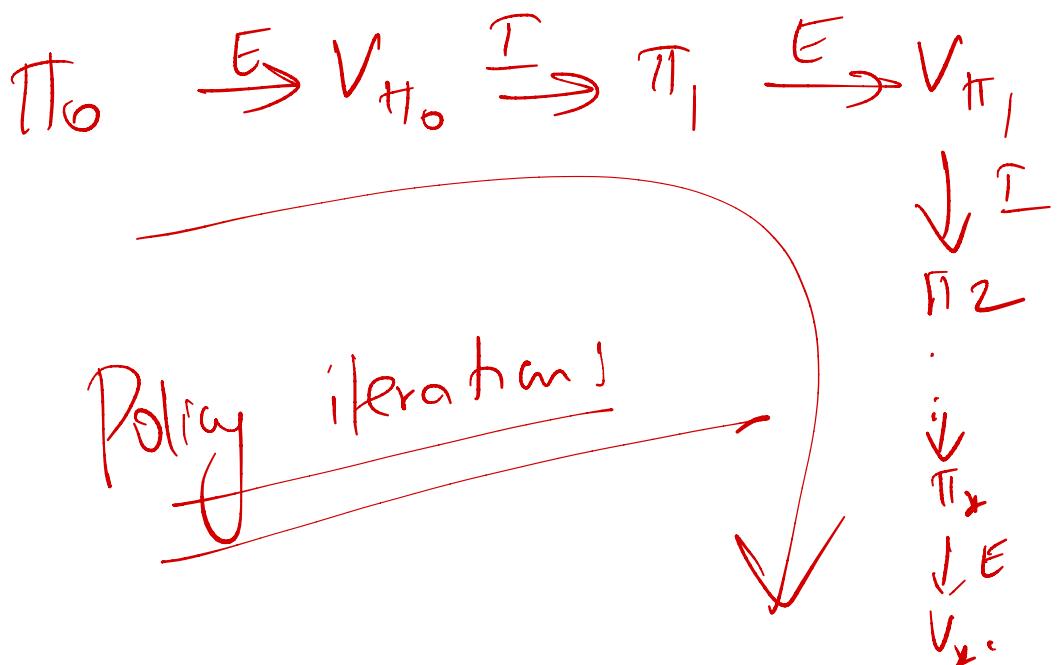


$$V_{\pi} = V_{\pi'}$$

$$V_{\pi'}(s) = \max_a \mathbb{E} \left[R_{t+1} + \gamma V_{\pi'}(s_{t+1}) \right]$$

$s = s'$
 $A_t = a$

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



Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$$

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

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement

$policy\text{-stable} \leftarrow true$

For each $s \in \mathcal{S}$:

$$old\text{-action} \leftarrow \pi(s)$$

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

If $old\text{-action} \neq \pi(s)$, then $policy\text{-stable} \leftarrow false$

If $policy\text{-stable}$, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Value iteration:

$$V_{k+1}(s) = \max_a E[R_{t+1} + \gamma V_k(s_{t+1})] \quad \begin{array}{l} s_t=s \\ A_t=a \end{array}$$

$$= \max_a \sum_{s',r} \sum_{s',r} p(s',r|s,a) [r + \gamma V_k(s')] \quad \forall s \in \mathcal{S}$$

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

```
| Δ ← 0
| Loop for each  $s \in \mathcal{S}$ :
|    $v \leftarrow V(s)$ 
|    $\cancel{V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]}$   
|    $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
until  $\Delta < \theta$ 
```

Output a deterministic policy, $\pi \approx \pi_*$, such that

$$\pi(s) = \arg \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$$

First-visit MC prediction, for estimating $V \approx v_\pi$

Input: a policy π to be evaluated

Initialize:

$V(s) \in \mathbb{R}$, arbitrarily, for all $s \in \mathcal{S}$

$Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Loop forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless S_t appears in S_0, S_1, \dots, S_{t-1} :

Append G to $Returns(S_t)$

$V(S_t) \leftarrow \text{average}(Returns(S_t))$

Monte Carlo ES (Exploring Starts), for estimating $\pi \approx \pi_*$

Initialize:

$\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s \in \mathcal{S}$

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose $S_0 \in \mathcal{S}, A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability > 0

Generate an episode from S_0, A_0 , following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a)$

On-policy first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$G \leftarrow \gamma G + R_{t+1}$

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow$ average($Returns(S_t, A_t)$)

$A^* \leftarrow \arg \max_a Q(S_t, a)$ (with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$