

1 Dynamic Programming

2nd Mid Term.
31th March.
11:30 - 1:00

Policy Evaluation

Solve value function for a given policy π

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

(Bellman Equation for Policy)

$$V^\pi = r^\pi + \gamma P^\pi V^\pi$$

$$r^\pi = (r^\pi(s), s \in S) \quad \text{s.t}$$

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

$$P^\pi = [[P^\pi(s, s')]]_{s, s' \in S}$$

$$P^\pi(s, s') = \sum_{a \in A(s)} \pi(a|s) \sum_r p(s', r|s, a)$$

$$\sum_{s' \in S} p^\pi(s, s') = 1$$

$$P^\pi(s, s') \geq 0 \quad \forall s, s' \in S$$

$$P^\pi(s, s') \geq 0$$

$$(I - \gamma P^\pi) V^\pi = r^\pi$$

$$\text{or} \quad V^\pi = (I - \gamma P^\pi)^{-1} r^\pi$$

Iterative Policy Evaluation

↳ value Iterates for a given policy π

$$V_{n+1}(s) = \sum_{a \in A(s)} \pi(a|s) \sum_{s', r} p(s', r|s, a) (r + V_n(s'))$$

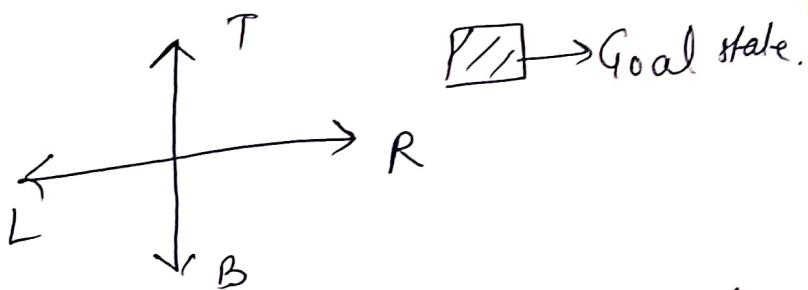
Starting for same V_0

$$V_n \rightarrow V^\pi \text{ as } V \rightarrow V_n$$

Consider a 4x4 grid world.

Example

1	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15



Transitions happens when you are in non shaded region

on all transitions.

$$R_t = -1$$

Non terminal states = 1, 2, ..., 14.

$$P(6, -1 | 5, R) = 1$$

↑
Right = action

Probability that u are in state 5 and go to state 6 with $R_t = -1$ and action = Right.

$$P(7, -1 | 7, R) = 1$$

$\neq r \in R$.

$$P(10, 0 | 5, R) = 0$$

We apply iterative policy evaluation for the random policy.

Picks all directions w.p. prob $1/4$.

$K=0$	$\gamma = 1$
0	0
0	0
0	0
0	0

///	↔	+	+
+	+	+	+
+	+	+	+
+	+	+	///

→ becoz u can take any action and a will be at same state 0 w.r.t U_K .

greedy Policy w.r.t U_K .

U_K for random Policy.

$K=1$	$\gamma < 1 = 0.9$
0	-1
-1	-1
-1	-1
-1	-1

	←	↔	+
↑	+	+	+
+	+	+	↓
+	+	→	

↑
↓
↔
→

$K=2$.

0	$\stackrel{1=3}{-1.75}$	-2	-2
-1.7	-2	-2	-2
-2	-2	-2	-1.7
-2	-2	-1.7	0

	←	←	+
↑	↑	+	↓
↑	+	↓	↓
↔	→	→	

Reward.

$$\text{scratched} \cdot -1 + \frac{1}{4} \times 0 + \frac{3}{4} \times -1 = -1.75.$$

1 prob. of moving left $\left(\frac{1}{4}\right)$

$K=3$

0	1	2	3
4	-2.4	-2.9	-3
-2.4	-2.85	-3	-2.9
-2.9	-3	-2.85	-2.9
-3	-2.9	-2.4	0

	\leftarrow	\leftarrow	\nwarrow
\uparrow	\uparrow	\leftrightarrow	\downarrow
\uparrow	\uparrow	\downarrow	\downarrow
\uparrow	\rightarrow	\rightarrow	

lowest cost path.

$$\begin{aligned}
 & \text{①} + \frac{1}{4}x_0 \rightarrow \\
 & -1 + \frac{1}{4}x_0 + \frac{1}{4}x^{-1.7} + \frac{2}{4}x^{-2} \Rightarrow s=1 \\
 & -1 + \frac{1}{4}x^{-1.7} + \frac{3}{4}x^{-2} \quad - s=2 \\
 & -1 + \frac{1}{4}x^{-2} + \frac{1}{4}(2x^{-2}) \quad s=5.
 \end{aligned}$$

$K=10$

0	-6.1	-8.4	-9
-6.1	-7.7	-8.4	-8.4
-8.4	-8.4	-7.7	-6.1
-9	-8.4	-6.1	0

	\leftarrow	\nwarrow	\downarrow
\uparrow	\uparrow	\leftrightarrow	\downarrow
\uparrow	\uparrow	\rightarrow	\downarrow
\uparrow	\rightarrow	\rightarrow	

$K=\infty$

0	-14	-20	-22
-14	-18	-20	-20
-20	-20	-18	-14
-22	-20	-14	0

	\leftarrow	\leftarrow	\uparrow
\uparrow	\uparrow	\uparrow	\downarrow
\uparrow	\uparrow	\rightarrow	\downarrow
\uparrow	\rightarrow	\rightarrow	

converged.

The greedy policy converge at $K=3$ only, but the optimal policy v_K converge at $K=\infty$.

Exercise

(1) If π is the equiprobable random policy, what is $q_{\pi}(11, \text{down})$? in $K=3$.

state.

q_{π} is optimal policy.

Solⁿ $q_{\pi}(11, \text{down}) = -1$

(2) $q_{\pi}(7, \text{down})$? in $K=\infty$.

Solⁿ $q_{\pi}(7, \text{down}) = -1 - 14 = -15$.

$v_{\pi}(11) = -14$, $q_{\pi}(11, \text{down}) = -1$

As $q_{\pi}(11, \text{down}) > v_{\pi}(11)$

$q_{\pi}(7, \text{down}) = -15$

$v_{\pi}(7) = -20$

$q_{\pi}(7, \text{down}) > v_{\pi}(7)$

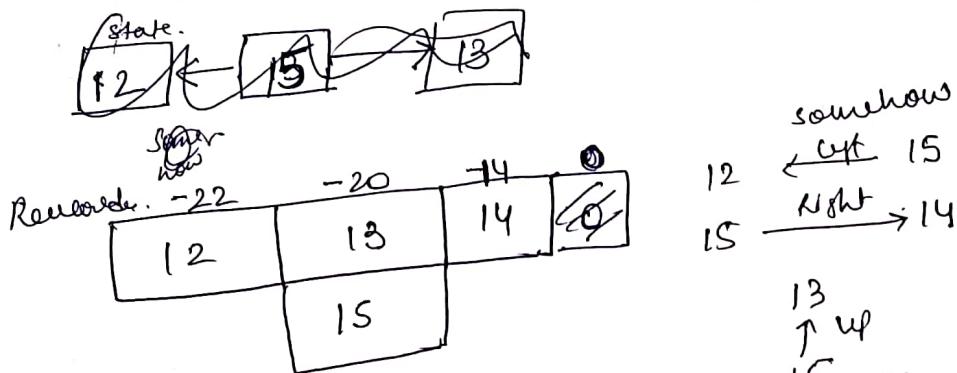
equiprobable
random Policy
can go anywhere

Hence q_{π} Policy is better than v_{π} Policy.

Likewise— $q_{\pi}(6, \text{left}) = -19$
 $v_{\pi}(6) = -20$

$$2\pi(\theta, \text{right}) = -21$$

Exercise we now add a new state 15 to the grid world just below state 13.



somehow
 12 $\xleftarrow{\text{up}} 15$
 15 $\xrightarrow{\text{right}} 14$
 13 $\uparrow \text{up}$
 15 $\downarrow \text{down}$

find $V_\pi(15)$ for the equiprobable policy
 every transition reward = -1
 consider the grid world for $\gamma = 0.9$.

$$V_\pi(15) = -1 - \frac{1}{4} \left[\begin{array}{c} 22 \\ \text{up} \\ \text{left} \end{array} + \begin{array}{c} 20 \\ \text{up} \\ \text{right} \end{array} + \begin{array}{c} 14 \\ \text{right} \\ \text{up} \end{array} \right] + \frac{1}{4} V_\pi(15)$$

$$\frac{3}{4} V_\pi(15) = -15$$

$$V_\pi(15) = -\frac{60}{3} = -20$$

Suppose now that dynamics of 13 is also changed such that action down from 13 takes agent to state 15. find $V_\pi(15)$.

If $V_\pi(13)$ changes from its neighboring value also changes and ultimately all values of the grid world changes. $13 \rightarrow 15$ changes values but $15 \rightarrow 13$ do not.

-22	13	-14	0
	14		

Assuming neighbouring values of state 13 do not change. calculating $v_{\pi}(13)$

$$v_{\pi}(13) = -1 - \frac{1}{4}(-22 + 20 + 14) + \frac{1}{4} v_{\pi}(15)$$

$$= -1 - 14 + \frac{1}{4} v_{\pi}(15)$$

$$= -15 + \frac{1}{4} v_{\pi}(15)$$

$$v_{\pi}(15) = -1 - \frac{1}{4} [22 + 14] + \frac{1}{4} v_{\pi}(13) + \frac{1}{4} v_{\pi}(15)$$

$$= -10 + \frac{1}{4} v_{\pi}(13) + \frac{1}{4} v_{\pi}(15)$$

$$= -10 + \frac{1}{4} (-15 + \frac{1}{4} v_{\pi}(15)) + \frac{1}{4} v_{\pi}(15)$$

$$v_{\pi}(15)$$

$$= -10 - \frac{15}{4} + \frac{5}{16} \cdot v_{\pi}(15)$$

$$v_{\pi}(15) = -20$$

$$v_{\pi}(13) = -20$$

But in general it is

not true. ~~that~~

that $v_{\pi}(13) = v_{\pi}(15)$

Mid Term Solution

$$(1) (a) P(T = +ve) = P(T = +ve | S) P(S) + P(+ve | \bar{S}) P(\bar{S})$$

$$= 0.10425.$$

$$(b) P(S | +ve) = \frac{P(+ve | S) P(S)}{\text{Total Probability}}$$

$$\underline{Q2} \quad V_n(i) = \frac{1}{n} \sum_{m=0}^{n-1} \left[r(i, a, x_m^i) + \gamma V_m(x_m^i) \right]$$

$$V_{n+1}(i) = \frac{1}{n+1} \left[r(i, a, x_n^i) + \gamma V_n(x_n^i) + \sum_{m=0}^{n-1} r(i, a, x_m^i) + \gamma V_m(x_m^i) \right]$$

$$= \frac{1}{n+1} \left[r(i, a, x_n^i) + \gamma V_n(x_n^i) \right] + \frac{n}{n+1} \left[V_n(i) \right]$$

$$4s. \quad X_{n+1} = x_n + a(n) \left[h(x_n) + M_{n+1} \right].$$

$$\rightarrow V_{n+1}(i) = \frac{1}{n+1} \left[r(i, a, x_n^i) + \gamma V_n(x_n^i) - V_n(i) \right] + \frac{n}{n+1} \left[V_n(i) \right]$$

$$+ \frac{1}{n+1} \left[r(i, a, x_n^i) + \gamma V_n(x_n^i) - V_n(i) \right]$$

$$+ \frac{1}{n+1} \left[r(i, a, x_n^i) + \gamma V_n(x_n^i) - V_n(i) \right]$$

$$\textcircled{3} (a) V^*(i) = \max_{a \in A_i} Q(i, a)$$

$$(b) q^*(i, a) = \mathbb{E} \left[r(i, a, j) + \gamma V^*(j) \right]$$

$$(c) \pi^*(i) = \arg \max_{a \in A_i} q^*(i, a)$$

$$(d) \pi^*(i) = \arg \max_{a \in A_i} \mathbb{E} \left[r(i, a, j) + \gamma V^*(j) \right]$$

Q4

$$V^*(\text{high}) = \max_{\substack{\{ \text{wait, search} \}}} \left[\begin{array}{l} 3 + \gamma V^*(\text{high}), \\ 0.4 [5 + \gamma V^*(\text{low})] \end{array} \right]$$

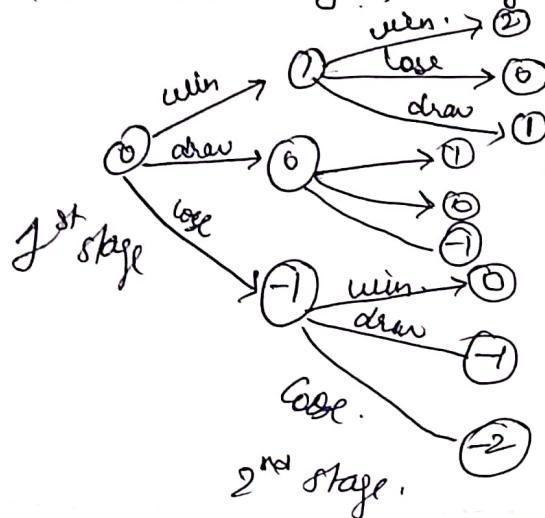
$$V^*(\text{low}) = \max_{\substack{\{ \text{search,} \\ \text{recharge,} \\ \text{wait} \}}} \left[\begin{array}{l} 0.2 [5 + \gamma V^*(\text{low}), \\ \text{search.} \\ 0.5 [-3 + \gamma V^*(\text{high}), \\ 0 + \gamma V^*(\text{high}), \\ 3 + \gamma V^*(\text{low})] \end{array} \right]$$

$$0.8 \cancel{0.3} \div$$

Q5 State space = ? difference of in points between
btw you & opponent.

Action state = $\{ \text{Timed, Bold} \}$.

Terminal single stage = +1, 0, -1



Rewards are given at the end of 2 games.

$$J^*(1) = \max_{\{Bold, Timid\}} \underbrace{[0.3[0+1] + 0.7[0+0],}_{\text{Bold!}} 0.6[0+1] + 0.4[0+0]}_{\text{Timid}}$$

Optimal action is Timid.

$$J^*(0) = \max_{\{Bold, Timid\}} [0.3[0+1] + 0.7[0-1], 0.6[0+0] + 0.4[0-1]]$$

Bold. ~~0.4~~

$$J^*(-1) = \max_{\{B, T\}} [0.3[0+0] + 0.7[0-1], 0.6[0-1] + 0.4[0-1]]$$

~~1~~ -10.

$$\begin{aligned}
 J_K^*(0) &= \max_{\{B, T\}} \left[0.3 [0 + J_2^*(1)] + 0.7 \right. \\
 &\quad \left. [0 + J_2^*(-1)] \right], \\
 &0.6 [0 + J_2^*(0)] + 0.4 [0 + J_2^*(-1)] \\
 &0.6 (-0.1) + 0.4 (-0.7) \\
 &= -2 + -2.8.
 \end{aligned}$$

Given any Policy π , if we find an action a in state s , such that $q_\pi(s, a) \geq v_\pi(s)$, then it is advantageous to pick action a in state s .

suppose we have a new policy π' such that $\pi'(i) = \pi(i) \quad \forall i \neq s$

$$\pi'(s) = a$$

while $\pi(s) \neq a$

Q) Is $v_{\pi'}(s) \geq v_\pi(s)$? $\forall s \in S$.

$$\begin{aligned}
 v_{\pi'}(s) &\leq q_{\pi'}(s, a) \\
 &= E \left[R_{t+1} + \gamma v_{\pi'}(s_{t+1}) \mid s_t = s, A_t = \pi'(s) \right] \\
 &= E_{\pi'} \left[R_{t+1} + \gamma v_{\pi'}(s_{t+1}) \mid s_t = s \right] \\
 &\leq E_{\pi'} \left[R_{t+1} + \gamma q_{\pi'}(s_{t+1}, \pi'(s_{t+1})) \right] \\
 &= E_{\pi'} \left[R_{t+1} + \gamma E_{\pi'} \left[R_{t+2} + \gamma v_{\pi'}(s_{t+2}) \right] \right] \\
 &\quad \left. \begin{aligned} s_{t+1}, A_{t+1} &= \pi'(s_{t+1}) \\ s_t &= s \end{aligned} \right]
 \end{aligned}$$

$$\begin{aligned}
&= E_{\pi^1} \left[R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{\pi^1}(s_{t+2}) \mid s_t = s \right] \\
&\leq E_{\pi^1} \left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 V_{\pi^1}(s_{t+3}) \mid s_t = s \right] \\
&\quad \vdots \\
&\leq E_{\pi^1} \left[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid s_t = s \right] \\
&= V_{\pi^1}(s)
\end{aligned}$$

$$\pi'(s) = \operatorname{argmax}_a q_{\pi}(s, a)$$

$$\begin{aligned}
\pi'(s) &= \operatorname{argmax}_a E \left[R_{t+1} + \gamma V_{\pi^1}(s_{t+1}) \mid s_t = s, A_t = a \right] \\
&= \operatorname{argmax}_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma V_{\pi^1}(s')]
\end{aligned}$$

[Policy Improvement.]

Suppose π' [the greedy policy] is as good as (but not better than) the old policy π , then

$$V_{\pi^1}(s) = V_{\pi^1}(s) \quad \forall s \in S$$

Then

$$\begin{aligned}
V_{\pi^1}(s) &= \max_{a \in A(s)} E \left[R_{t+1} + \gamma V_{\pi^1}(s_{t+1}) \mid s_t = s, A_t = a \right] \\
&= \max_{a \in A(s)} \sum_{s', r} p(s', r \mid s, a) [r + \gamma V_{\pi^1}(s')] \\
&= \boxed{V_{\pi^1} = V^* = V_{\pi}}
\end{aligned}$$

Consider a state s where there are a_1, a_2, \dots, a_n as the feasible actions.

Let a_1, a_2, a_3 be the maximizing actions.

Consider a stochastic policy.

$$\pi'(s) = \begin{cases} a_1 & \text{w.p. } p_1 \\ a_2 & p_2 \\ a_3 & p_3 \end{cases}$$

$$\text{s.t. } p_1 + p_2 + p_3 = 1 \quad [\& p_1, p_2, p_3 \geq 0]$$

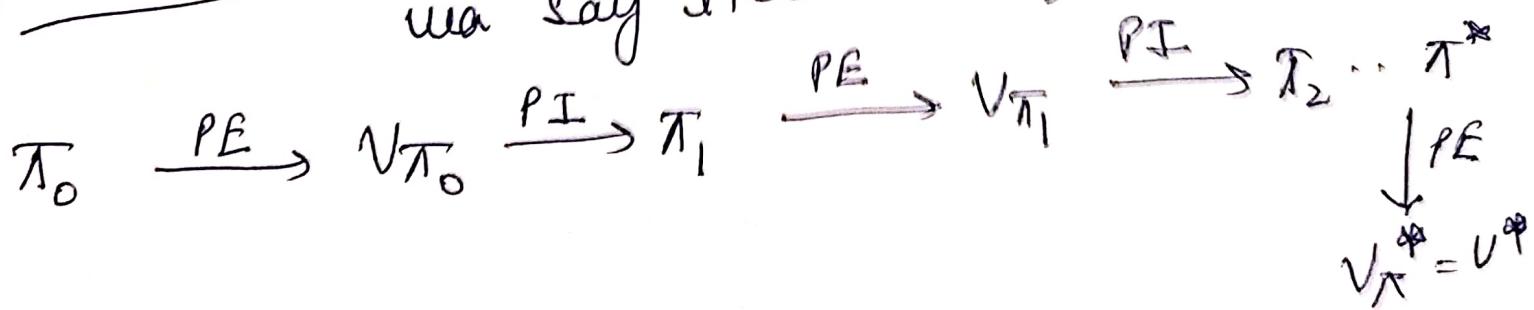
Then $\pi'(s)$ gives the optimal actions in state s regardless of p_1, p_2, p_3 .

Policy Iteration

Involves two nested loops

Outer loop policy improvement

Inner loop policy evaluation. (can be done via say iterative PE)



V^* is the optimal value policy.

π^* is the optimal policy.

Consider a state s where there are a_1, a_2, \dots, a_n as the feasible actions.

Let a_i^*, a_j^*, a_k^* be the maximising actions.

Consider a stochastic policy.

$$\pi'(s) = \begin{cases} a_i & \text{w.p. } p_1 \\ a_j & p_2 \\ a_k & p_3 \end{cases}$$

$$\text{s.t. } p_1 + p_2 + p_3 = 1 \quad [\& p_1, p_2, p_3 \geq 0]$$

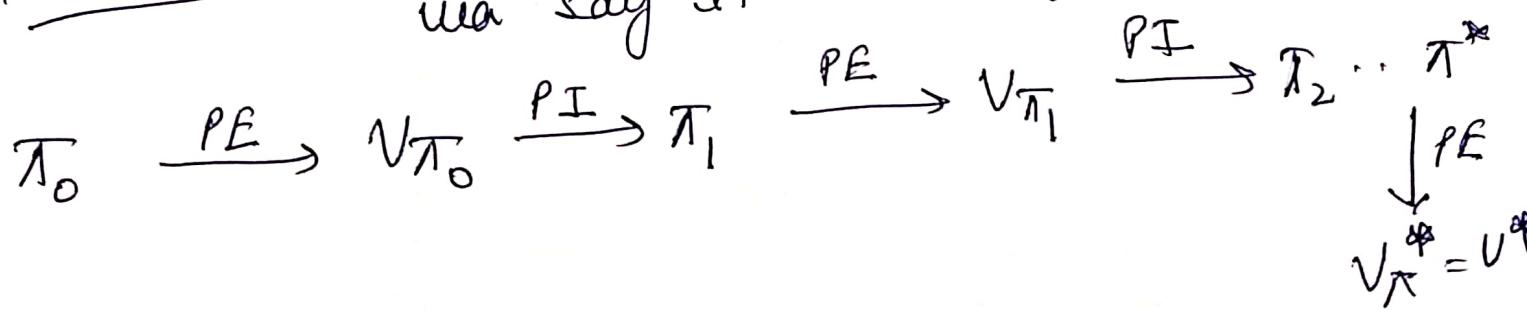
Then $\pi'(s)$ gives the optimal actions in state s regardless of p_1, p_2, p_3 .

Policy Iteration

Involves two nested loops

Outer loop policy Improvement

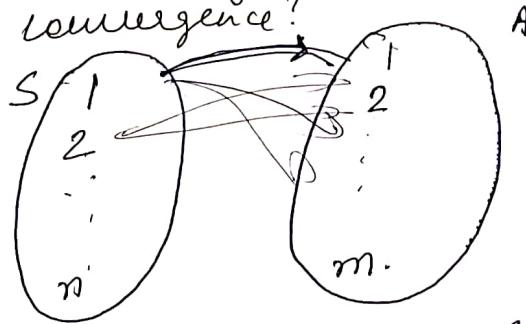
Inner loop policy evaluation. (can be done via say iterative PE)



V^* is the optimal value policy.

π^* is the optimal policy.

Q. how many steps it will take for convergence?



$$\text{no. of functions b/w } S \text{ & } A = m \underbrace{ \times m \times \dots m }_{n \text{ times.}} = \underline{\underline{m^n}}$$

Hence No. of deterministic policies is also finite.
if S & A is finite

Policy Iteration algorithm converges in a finite no. of state [measured in terms of policy improvement]

policy Iteration (using iterative PE) for estimating

$$\pi \rightsquigarrow \pi^*$$

1. Initialize :- $V(s) \in \mathbb{R}$ and $\pi(s) \in A(s)$

arbitrarily for $s \in S$.

2. Policy Evaluation :-

loop $\Delta \leftarrow 0$
loop for each $s \in S$
 $V \leftarrow V(s)$

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

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

until $\Delta < \theta$ [some small positive No.]

3. Policy Improvement

for each $s \in S$,

$$\text{old actions} \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s', r} P(s', r | s, a) [r + \gamma V(s')]$$

If old actions $\neq \pi(s)$

Go back to step 2.

If $V_\pi(s) = V_{\text{old}}(s) \forall s$, then $V_\pi = V^*$ &
 π & old are both optimal.

Policy Iteration

(1) Policy evaluation

Given Policy update π_k , evaluate V_{π_k} as

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

(2) Policy Improvement

$$\pi_{k+1}(s) \operatorname{argmax}_{a \in A(s)} \sum_{s', r} P(s', r | s, a) [r + \gamma V_{\pi_k}(s')]$$

$$\text{Or } \boxed{T V_{\pi_k} = T_{\pi_{k+1}} V_{\pi_k}}$$

$$\max_{a \in A(s)} \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi_K}(s')]$$
$$= \sum_{s', r} p(s', r | s, \pi_{k+1}(s)) [r + \gamma v_{\pi_K}(s')]$$