

Policy Iteration

$s=0$	$s=1$	$s=2$	$s=3$
$s=4$	$\cancel{s=5}$	$s=6$	$\cancel{s=7}$
$s=8$	$s=9$	$s=10$	$s=11$

Given policy

$$v(6) = \text{rewards following policy } [6, 10, 11, 7]$$

$$= \gamma(6) + \gamma(10) + \gamma(11) + \gamma(7)$$

For any state s at time $t=0$, utility can be written as:-

$$v(s_{t=0}) = \text{rewards following policy } [s_{t=0}, s_{t=1}, s_{t=2}, s_{t=3}, \dots]$$

$$= \gamma(s_{t=0}) + \gamma(s_{t=1}) + \gamma(s_{t=2}) + \dots$$

$$= \sum_t \gamma(s_t)$$

$$v(6) = \gamma(6) + \gamma(10) + \gamma(11) + \gamma(7)$$

$$\text{utility of } \underline{s=10} = \gamma(10) + \gamma(11) + \gamma(7)$$

$$v(6) = \gamma(6) + v(10)$$

$$v(s_{t=i}) = \gamma(s_{t=i}) + v(s_{t=i+1})$$

↑

utility of a state = immediate reward of
that state + utility of its successor
state following the given policy.

$$v(s) = \gamma(s) + v(s')$$

Discounted rewards :

$$s=0$$

By following the policy, the future states are $[0, 1, 2, 1, 2, 1, 2, \dots]$ which is a loop of 1 and 2 forever.

If following the definition of utility

$$\begin{aligned} v(0) &= \gamma(0) + \gamma(1) + \gamma(2) + \gamma(1) + \dots \gamma(2) \\ &= -\infty \end{aligned}$$

Address the concept of Discounted

Rewards.

"Reward received in the future is less valuable than reward received

right now." — Introduce discount factor γ

Additive rewards:-

$$v(s_{t=0}) = \sum_t \gamma^t r(s_t)$$

Using discounted rewards:

$$v(s_{t=0}) = \gamma^t \sum_t r(s_t)$$

finite no. of states = finite
no. of rewards

$0 < \gamma < 1$
utility value
is guaranteed
to be
bounded

$v(s_{t=0})$ = discounted rewards following policy $[s_{t=0}, s_{t=1}, s_{t=2}, s_{t=3}, \dots]$

$$= r(s_{t=0}) + \gamma r(s_{t=1}) + \gamma^2 r(s_{t=2}) + \dots$$

$$v(s) = r(s) + \gamma v(s')$$

$v(s_0)$ = discounted rewards following policy

$$[s_{t=0}, s_{t=1}, s_{t=2}, s_{t=3}, \dots]$$

$$= \gamma r(s_{t=0}) + \gamma r(s_{t=1}) + \gamma^2 r(s_{t=2}) + \gamma^3 r(s_{t=3}) + \dots$$

$$\leq \pi_{\max} + \gamma \pi_{\max} + \gamma^2 \pi_{\max} + \gamma^3 \pi_{\max} + \dots$$

$$= \pi_{\max} \left[\frac{1}{1-\gamma} \right]$$

Add Back Stochastic :-

Deterministic Env

$$s = 6 \quad \text{determined path } [6, 10, 11, 7]$$

$$v(6) = \gamma(6) + \gamma(10) + \gamma(11) + \gamma(7)$$

In Stochastic Env : not guaranteed to follow

path $[6, 10, 11, 7]$

$$\text{prob: } (0.8)^3 = 0.512$$

Define Utility in a stochastic Environment:-

$$v(s_0) = E \left[\sum_t \gamma^t \pi(s_t) \right]$$

E is the expectation w.r.t probability distribution over possible paths following a certain policy

$$v(s) = \pi(s) + \gamma \sum_{s'} P(s' | s, a=\pi(s)) v(s')$$

$$v(s) = \pi(s) + \gamma \sum_{s'} P(s' | s, a=\pi(s)) v(s')$$

\triangleright $s=0$	\triangleright $s=1$	\triangleleft $s=2$	 $s=3$
\triangle $s=4$	 $s=5$	\triangleright $s=6$	 $s=7$
\triangleright $s=8$	\triangleleft $s=9$	\triangleright $s=10$	\triangle $s=11$

Given policy

$s=6$:- to reach $s'=10$ Probability = 0.8
 to reach $s'=7$ Probability = 0.1
 to stay in 6 Probability = 0.1

Therefore, the utility of $s=6$ is :-

$$v(6) = \gamma(6) + \beta [0.8 v(10) + 0.1 v(6) + 0.1 v(7)]$$

$$v(0) = \gamma(0) + \beta [0.8 v(1) + 0.1 v(4) + 0.1 v(0)]$$

$$v(1) = \gamma(1) + \beta [0.8 v(2) + 0.1 v(1) + 0.1 v(1)]$$

$$v(2) = \gamma(2) + \beta [0.8 v(1) + 0.1 v(6) + 0.1 v(2)]$$

$$v(3) = \gamma(3) + \beta [v(3)]$$

$$v(4) = \gamma(4) + \beta [0.8 v(0) + 0.1 v(4) + 0.1 v(4)]$$

$$v(5) = \gamma(5) + \beta [0.8 v(10) + 0.1 v(5) + 0.1 v(7)]$$

$$v(7) = \gamma(7) + \beta [v(7)]$$

$$v(8) = \gamma(8) + \beta [0.8 v(9) + 0.1 v(4) + 0.1 v(8)]$$

$$v(9) = \gamma(9) + \beta [0.8 v(9) + 0.1 v(8) + 0.1 v(10)]$$

$$v(10) = r(10) + \gamma [0.8v(11) + 0.1v(6) + 0.1v(10)]$$

$$v(11) = r(11) + \gamma [0.8v(7) + 0.1v(10) + 0.1v(11)]$$

11 Equations, 11 unknowns.

 Policy Evaluation
 ↘

Determine the utility of each state

DP approach - iterative policy evaluation
or sweep.

- ① Initialize the utility of each state = 0
- ② Loop through the states
- ③ For each state, update the utility

using the equation

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

Eg:- set $\gamma = 0.5$

loop through $s=0$ to $s=11$, we get

$$\begin{aligned}
 v(0) &= -0.04 \\
 v(1) &= -0.04 \\
 v(2) &= -0.056 \\
 v(3) &= 1 \\
 v(4) &= -0.056 \\
 v(5) &= -0.04 \\
 v(6) &= -0.04 \\
 v(7) &= -1 \\
 v(8) &= -0.0428 \\
 v(9) &= -0.04214 \\
 v(10) &= -0.042 \\
 v(11) &= -0.0421
 \end{aligned}$$

in-place sweep
 use updated utility
 in the same
 sweep

$$\begin{aligned}
 v(1) &= r(1) + \gamma [0.8v(2) + 0.1v(1) \\
 &\quad + 0.1v(1)]
 \end{aligned}$$

$$\begin{aligned}
 &= -0.04 + 0.5 [0.8 \times 0 + 0.1 \times 0 + 0.1 \times 0] \\
 &= -0.04
 \end{aligned}$$

in-place sweep

$$\begin{aligned}
 v(2) &= r(2) + \gamma [0.8v(1) + 0.1v(6) + 0.1v(2)] \\
 &= -0.04 + 0.5 [0.8(-0.04) + 0.1 \times 0 \\
 &\quad + 0.1 \times 0] = -0.056
 \end{aligned}$$

We repeat sweeping for several
 times until the changes of utility

values between consecutive sweeps

are marginal

Utility after 11 Sweeps :-

$s=0$	$s=1$	$s=2$	$s=3$
$v = -0.0831$	$v = -0.0872$	$v = -0.0963$	$v = 1.9990$
$s=4$		$s=6$	$s=7$
$v = -0.0814$		$v = -0.3329$	$v = -1.9990$
$s=8$	$s=9$	$s=10$	$s=11$
$v = -0.0930$	$v = -0.1110$	$v = -0.4413$	$v = -0.9070$

Here :- For $s=6$

$$v(6) = r(6) + \frac{1}{2} [0.8v(10) + 0.1v(6) + 0.1v(7)]$$

$$v(6) = -0.04 + 0.5 [0.8(-0.4413) + 0.1(-0.3329) + 0.1(-1.9990)]$$

$$\approx -0.3329$$

Sweep approach is alternate method
to solve linear equations

Policy Improvement

$S = 0$	$S = 1$	$S = 2$	$S = 3$
$V = -0.0831$	$V = -0.0872$	$V = -0.0963$	$V = 1.9990$
$S = 4$		$S = 6$	$S = 7$
$V = -0.0814$		$V = -0.3329$	$V = -1.9990$
$S = 8$	$S = 9$	$S = 10$	$S = 11$
$V = -0.0930$	$V = -0.1110$	$V = -0.6613$	$V = -0.9070$

① After we get the utility for this policy, it is time to improve the policy.

* Utility represents how good a state is.

② When choosing ACTIONS for a state, prefer successor states with higher utility

③ when $S = 2$

potential successor state include 1, 3, 6

Clearly $S = 3$ is best choice since

$$V = 1.9990$$

$S=0$	$S=1$	$S=2$	$S=3$
\rightarrow	\rightarrow	\leftarrow	
$V = -0.0831$	$V = -0.0872$	$V = -0.0963$	$V = 1.9990$
$S=4$		$S=6$	$S=7$
$V = -0.0814$	\uparrow	$V = -0.3329$	$V = -1.9990$
$S=8$	$S=9$	$S=10$	$S=11$
\rightarrow	\downarrow	\rightarrow	\uparrow
$V = -0.0930$	$V = -0.1110$	$V = -0.4413$	$V = -0.9070$

Since MDP is stochastic, selecting a preferable successor state does not guarantee we will reach it.

Rather than successor states, compare actions.

For $S=6$ [actions: UP, DOWN, RIGHT, LEFT]

choosing UP \rightarrow the reward is

$$0.8 V(2) + 0.1 V(6) + 0.1 V(7)$$

$$= 0.8(-0.0963) + 0.1(-0.3329) + 0.1(-1.990) = \boxed{-0.3093}$$

choosing LEFT \rightarrow the reward is best choice

$$0.8 V(6) + 0.1 V(10) + 0.1 V(2) = -0.3201$$

$$0.8(-0.3329) + 0.1(-0.4413) + 0.1(-0.0963)$$

DOWN :-

$$0.8 v(10) + 0.1 v(6) + 0.1 v(7) = -0.5862$$

$$= 0.8 [-0.4413] + 0.1 (-0.3329) + 0.1 (-1.9990)$$

RIGHT:

$$0.8 v(7) + 0.1 v(2) + 0.1 v(10)$$

$$= 0.8 (-1.9990) + 0.1 (-0.0963) + 0.1 (-0.4413)$$

$$= -1.6530$$

Hence update policy of state $s=6$ to UP

Perform this process for each state

$$\pi(s) \leftarrow \operatorname{argmax}_a \left[\sum_{s'} p(s'|s, a) v(s') \right]$$

After policy improvement, the policy looks like :-

$s=0$	$s=1$	$s=2$	$s=3$
\downarrow	\leftarrow	\rightarrow	
$v = -0.0831$	$v = -0.0872$	$v = -0.0963$	$v = 1.9990$
$s=4$		$s=6$	$s=7$
\rightarrow		\uparrow	
$v = -0.0814$		$v = -0.3329$	$v = -1.9990$
$s=8$	$s=9$	$s=10$	$s=11$
\uparrow	\leftarrow	\leftarrow	\leftarrow
$v = -0.0930$	$v = -0.1110$	$v = -0.4413$	$v = -0.9070$

This is a better policy than the previous one. But, this is not the best one.

↳ Repeat process of policy evaluation and policy improvement again and again.

↳ Do not need to initialize utility to zero continue with the utility values, we got. This helps to converge faster

policy evaluation

* update utility values while keeping policy constant

from $\pi[i], v[i-r]$

↓

$\pi[i], v[i]$

policy improvement

* update policy while keeping utility values constant

from $\pi[i], v[i]$

↓

$\pi[i+1], v[i]$

When $\pi[i] = \pi[i+1]$
break; # convergence occurs

In Example from:

R randomly generated policy $\pi[1]$

optimal policy π^*

$\pi[6]$

Effects of discounted factor

Discounted factor γ is used to discount the rewards of future states when adding them together

$$V(S_{t=0}) = \pi(S_{t=0}) + \gamma \pi(S_{t=1}) + \gamma^2 \pi(S_{t=2}) \dots$$

$$= \gamma^t \sum_t \pi(S_t)$$

If we choose 0.1 as γ , we have :-

$$V(S_{t=0}) = \pi(S_{t=0}) + 0.1 \pi(S_{t=1}) + 0.1^2 \pi(S_{t=2}) \dots$$

If we choose 0.9 as γ , we have:-

$$V(S_{t=0}) = \pi(S_{t=0}) + 0.9 \pi(S_{t=1}) + 0.9^2 \pi(S_{t=2}) \dots$$

smaller γ (0.1) emphasizes immediate

reward of current state

larger γ (0.9) emphasizes long term reward
in the future