## 2.

Value function of a state following a stochastic policy: -

$$V^{\pi}(s) = E_{\tau \sim \pi} \left[ R(\tau) \mid S_0 = s \right]$$

$$\therefore V^{\pi}(A) = E_{\tau \sim \pi} \left[ R(\tau) \mid S_0 = A \right]$$

gives the prob of each action $a_i$ in state A $a_i = \{ down$ right$\}$

$$= \sum_i R(\tau_i) \, \pi(a_i \mid A)$$

Return of each trajectory $\tau_i$  $\tau_i = \{ \tau_1, \tau_2 \}$

$$= R(\tau_1) \, \pi(down \mid A) +$$
$$R(\tau_2) \, \pi(right \mid A)$$

$$= 4(0.8) + 3(0.2) = 3.6.$$

Ex for value function ( Pg : 33)

Assume we have 2 policies $\pi_1$ and $\pi_2$.

Let $V^{\pi}(s)$ value of state s, using policy $\pi_1$ is

$V^{\pi_1}(s) = 13$ ; $V^{\pi_2}(s) = 11$. $\therefore$ Optimal value of

$s = 13$. The policy that gives the optimal value for a state s, is the optimal policy $\pi^*$.

Hence $\pi_1$ is the optimal policy.

If we have the value table of two such states giving us the max (optimal) value in each of them.

| State | Value |
|-------|-------|
| $S_0$ | 7 |
| $S_1$ | 11 |

∴ It's better to be in $S_1$ ∵ $S_0$ has less state value.

∴ $S_1$ is the optimal state.

∴ we can find the optimal state from a value table.

Q function:-
$$Q^{\pi}(S, a) = [R(\tau) | S_0 = S, a_0 = a]$$

In value fn, we find the value of a state, " Q fn, we find the value of a state-action pair.

Ex: Pg: 34:
$$Q^{\pi}(A, down) = \left( R(\tau) | S_0 = A, a_0 = down \right.$$

$$= 1 + 1 + 1 + 1 = 4.$$

$$Q^{\pi}(D, right) = \{ R(\tau) | S_0 = D, a_0 = right \}$$

$$= 1 + 1 + 1 = 3.$$

Expected Q function:-
$$Q^{\pi}(S, a) = \underset{\tau \sim \pi}{E} [R(\tau) | S_0 = S, a_0 = a]$$

Q fn depends on the policy. ∴ there will be different Q-values for a $(S, a)$ pair, depending on the policy. The optimal Q fn

Pg: ②

for a $(S, a)$ pair is the one that gives the max (optimal) $Q$-fn value among all.

$$\therefore \quad Q^*(S, a) = \max_{\pi} Q^{\pi}(S, a)$$

$\therefore$ The optimal policy $\pi$, for a $(S, a)$ is the one that gives the max $Q$-fn.

• $Q$ fns can be visualized using a $Q$-table.

Ex: find the optimal policy $\pi$, for the given $Q$-table.

| State | action | Value |
|-------|--------|-------|
| $S_0$ | 0 | 9 |
| $S_0$ | 1 | 11 |
| $S_1$ | 0 | 17 |
| $S_1$ | 1 | 13 |

$\longrightarrow$

optimal policy

| State | Action |
|-------|--------|
| $S_0$ | 1 |
| $S_1$ | 0 |

$\rightarrow$ det. poly

$\therefore$ Recap:- $\quad V^{\pi}(S) = [R(\tau) \mid S_0 = S]$ : returns

for stochastic policy

$$V^{\pi}(S) = \mathop{E}_{\tau \sim \pi} [R(\tau) \mid S_0 = S]$$

$$= \sum_i {}^{'}R(\tau_i) \, \pi(a_i, s)$$

Pg: ③

$$V^*(s) = \max_{\pi} V^{\pi}(s) \quad \cdots \quad \begin{cases} \text{optimal value of} \\ \text{a state } s \end{cases}$$

given value table of each state, we can find the optimal state.

R – function.

$$Q^{\pi}(s,a) = \left[ R(\tau) \mid s_0 = s, a_0 = a \right]$$

$$Q^{\pi}(s,a) = \mathop{E}_{\tau \sim \pi} \left[ R(\tau) \mid s_0 = s, a_0 = a \right]$$

$$Q^*(s,a) = \max_{\pi} Q^{\pi}(s,a)$$

Given Q–table, we can find the optimal policy.

Bellman Eqn of the value fn for a state for a deterministic envt:-

$$V^{\pi}(s) = R(s, a, s') + \gamma V(s')$$

for a stochastic envt. — Each state has different next states for an action a'

$$V^{\pi}(s) = \sum_{s'} P(s' \mid sa) R(s, a, s') + \gamma V(s')$$

$s'$ trans. prob of next state.     backup     Pg : (4)

$$V(s_1) = \left[ P(s_2|s_1, a_1) R(s_1, a_1, s_2) + \gamma V(s_2) \right] +$$

$$\left[ P(s_3|s_1, a_1) \quad P(s_1, a_1, s_3) + \gamma V(s_3) \right]$$

$$= 0.7 \ R(s_1, a_1, s_2) + \gamma V(s_2) +$$

$$0.3 \ R(s_1, a_1, s_3) + \gamma V(s_3)$$

$$\therefore \quad \boxed{V^{\pi}(s) = \sum_{s'} P(s'|s, a) \left[ R(s, a, s') + \gamma V(s') \right]}$$

⇓ for a stochastic

but deterministic policy cwt

what if the

ⓧ next state is associated with the cwt.

Policy is associated with an action

stochastic cwt?? When we take an action 'a' in a particular state 's', the next state s' is random.

Stochastic policy:? When we are in a state 's', the next action we take is random.

with stochastic Policy: Action in each state has randomness.

$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) \; R(s,a,s') + \gamma V(s')$$

a ⟶ stochas. policy

$s'$ ⟶ stochas. envt

backup

Using expectations:

$$V^\pi(s) = \mathop{E}_{\substack{a \sim \pi \\ s' \sim P}} \left[ R(s,a,s') + \gamma V(s') \right]$$

———×———————×—————×——×——————×—

Bellman Equation of the Q function :-

· It is the sum of the immediate reward and the discounted value of the next state-action pair.

$$\therefore Q^\pi(s,a) = R(s,a,s') + \gamma Q(s',a')$$

for det. envt ↗

for stocastic envt :-

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

for stochastic policy :-

$$Q^\pi(s,a) = \sum_{s'} P(s'|s,a)\left[ R(s,a,s') + \gamma \sum_{a'} \pi(a'|s') Q(s', \right.$$

Pg : ⓕ

∴ Bellman Expectation Equation of the Q function is for a $(s,a)$ pair is:

$$Q^{\pi}(s,a) = \sum_{s'} \overbrace{\left[ R(s,a,s') + \gamma \sum_{a'} \pi(a'|s') Q(s',a') \right]}^{P(s'|s,a)}$$

$$\therefore Q^{\pi}(s,a) = \mathbb{E}_{s' \sim P} \left[ R(s,a,s') + \gamma \mathbb{E}_{a' \sim \pi} Q(s',a') \right]$$

## Bellman Optimality Theorem :—

Recap:

$$V^{\pi}(s) = \sum_{a} \pi(a|s) \sum_{s'} \overbrace{\left[ R(s,a,s') + \gamma V(s') \right]}^{P(s'|s,a)}$$

$$Q^{\pi}(s,a) = \sum_{s'} \overbrace{R(s,a,s')}^{P(s'|s,a)} + \gamma \sum_{a'} \pi(a'|s') Q(s',a')$$

Recap:

1) Bellman Eqn for the Value fn of a State s is :—

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

2) Above in expectation form :—

$$V^{\pi}(s) = \mathbb{E}_{\substack{a \sim \pi \\ s' \sim P}} \left[ R(s,a,s') + \gamma V(s') \right]$$

3)

3) Bellman Eqn for the Q fn of $(s,a)$ pair is:-

$$Q^{\pi}(s,a) = \cancel{R[s]} \cancel{R(s,a,s')} + \cancel{\gamma \sum_{a'} \pi(a'|s') Q(s',a')}$$

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

4) Above eqn in expectation form:-

$$Q^{\pi}(s,a) = \mathop{E}_{s' \sim P} \left[ R(s,a,s') + \gamma \mathop{E}_{a' \sim \pi} Q(s',a') \right]$$

Bellman Optimality Theorem:-

We know, the Bellman eqn of the Value and Q fn. is expectation form is

$$V^{\pi}(s) = \mathop{E}_{\substack{a \sim \pi \\ s' \sim P}} \left[ R(s,a,s') + \gamma V(s') \right]$$

$$Q^{\pi}(s,a) = \mathop{E}_{s \sim P} \left[ R(s,a,s') + \gamma \mathop{E}_{a' \sim \pi} Q(s',a') \right]$$

To find the optimal Bellman Value fn:-

The Bellman optimality eqn, gives the optimal

$\beta$. ©

Bellman value and Q fn...

W.K.T, the value fn of a state depends on the policy.

V(s) for a particular state s, varies with the policy.

∴ the optimal Bellman V(s) for a particular state s, is the one that gives the max value among all values of that state.

How to find this $V^*(s)$.?

We compute the optimal Bellman $V^*(s)$, by selecting the action that gives the max value in that state.

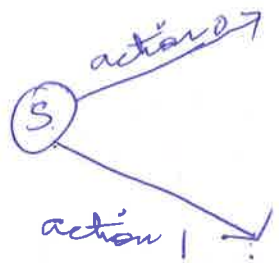We' don't know which action gives the max value

∴ instead of using a policy $\pi$, to select an action, we find the value of the state using all possible actions, and then select the max value of all.

∴ no policy is used, remove the expectation over policy $\pi$, add the max over action.

∴ optimal Bellman Value fn of a state is:_

$$V^*(s) = \max_{a} \underset{s' \sim P}{E} \left[ R(s,a,s') + \gamma V(s') \right]$$

**Example:**

action 0 ⟶ In states $s$, and have 2 possible actions 0 and 1.

$$V^*(s) = \max \left( \begin{array}{c} \underset{s'\sim P}{E} \left[ R(s, 0, s') + \gamma v^*(s') \right] \\ \underset{s'\sim P}{E} \left[ R(s, 1, s') + \gamma v^*(s') \right] \end{array} \right)$$

**Optimal Bellman Q function:-**

w.r.t.

$$Q^\pi(s, a) = \underset{s'\sim P}{E} \left[ R(s, a, s') + \gamma \underset{a'\sim\pi}{E} Q(s', a') \right]$$

Instead of ~~using a policy $\pi$ for~~ choosing the action $a'$ in state $s'$, we choose all possible actions in $s'$, and find the max Q value.

$$\therefore Q^*(s, a) = \underset{s'\sim P}{E} \left[ R(s, a, s') + \gamma \underset{a'}{\max} Q^*(s', a') \right]$$

**Example:**

$s$ ⟶ $a$ ⟶ $s'$ ⟶ 0, 1

$$\therefore Q'(s, a) = \underset{s'\sim P}{E} \left[ R(s, a, s') + \gamma \max \left( \begin{array}{c} Q^*(s', 0) \\ Q^*(s', 1) \end{array} \right) \right]$$

Pg: 10

To summarize:-

i). The Bellman optimality eqn for value fn of a state $s$, is:-

$$V^*(s) = \max_a \; E_{s' \sim P} \left[ R(s,a,s') + \gamma V^*(s') \right]$$

2) The Bellman optimality eqn for Q fn of a $(s,a)$ pair is

$$Q^*(s,a) = E_{s' \sim P} \left[ R(s,a,s') + \gamma \max_{a'} Q^*(s',a') \right]$$

3) we can also expand the expectation & rewrite

$$V^*(s) = \max_a \sum_{s' \sim P} \left[ \overbrace{R}^{P(s'|s,a)} (s,a,s') + \gamma V^*(s) \right]$$

$$Q^*(s,a) = \sum_{s' \sim P} P(s'|s,a) \left[ R(s,a,s') + \gamma \max_{a'} Q^*(s',a') \right]$$

The relationship between the value and Q fns:-

In chapter 1, w.k.t, the value of a state $s$, is the expected return starting from that state $s$, following a policy $\pi$.

$$V^\pi(s) = E_{\tau \sim \pi} \left[ R(\tau) \mid S_0 = s \right]$$

Clearly, the Q value of a $(s,a)$ pair, is the expected return of starting from state $s'$, by performing action, $a$.

A)

$$\therefore \quad Q^{\pi}(s,a) = \underset{\tau \sim \pi}{E} \left[ R(\tau) \mid s_0 = s, \; a_0 = a \right]$$

w.k.t, the optimal value fn of a state, s gives its maximum state-value.

$$V^*(s) = \underset{\pi}{max} \; V^{\pi}(s)$$

the optimal Q fn of a $(s,a)$:-

$$Q^*(s,a) = \underset{\pi}{max.} \; Q^{\pi}(s,a)$$

Is there a relation betn. $V^*(s)$ and $Q^*(s,a)$?

w.k.t. the optimal value fn. of a state is the
max expected returns whoure start from of a state. The optimal
Q fn. is the max. expected value return
when we start from a state and performs an
action `a`.

$\therefore$ the optimal value fn. of a state, is the
max. of all optimal Q value over all possible
actions from a state s, $\therefore$ we can derive
V from Q:

$$\boxed{V^*(s) = \underset{a}{max} \left[ Q^*(s,a) \right]}$$

## Recap :-

Bellman eqn of the value and Q fns. -

$$V^{\pi}(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) \left[ R(s,a,s') + \gamma V^{\pi}(s') \right]$$

$$Q^{\pi}(s,a) = \sum_{s'} P(s'|s,a) \left[ R(s,a,s') + \gamma Q^{\pi}(s',a') \updownarrow \sum_{a'} \pi(a'|s') \right]$$

Bellman optimality eqn of the value & Q fns. -

↳ ( irrespective of any policy )

$$V^{*}(s) = \max_a \sum_{s'} P(s'|s,a) \left[ R(s,a,s') + \gamma V^{*}(s') \right]$$

$$Q^{*}(s,a) = \sum_{s'} P(s'|s,a) \left[ R(s,a,s') + \gamma \max_{a'} Q^{*}(s',a') \right] \to \textcircled{1}$$

If we have optimal value fn $V^{*}(s)$, we can use it to derive optimal Bellman Q fn, ( we can derive Q from V).

$\therefore$ w.k.t: $\quad V^{*}(s) = \max_a Q^{*}(s,a) \quad -- (2)$

Subst (2) in (1), we get

$$Q^{*}(s,a) = \sum_{s'} P(s'|s,a) \left[ R(s,a,s') + \gamma V^{*}(s') \right] \quad --(3)$$

Also we can be

Subs (3) in (2), we get

$$V^*(s) = \max_a \sum_{s'} P(s'|s,a) \left[ R(s,a,s') + \gamma V^*(s') \right]$$

---

**Dynamic Programming.** In RL, we use DP to find the optimal policy, using 2 methods: 1) Value iteration and 2) policy iteration.

DP is a model-based method; ie, to find an optimal policy using DP, we should know the model dynamics. [station transn prob; and reward functions).

1) **Value iteration method:** —

a) An optimal policy tells the agent to perform a correct action in each state.

b) To find this policy, first we find the optimal value fn of each state; $\tilde{V}^*(s)$

c) Later use this optimal value fn $V^*(s)$, find the optimal policy [by finding the Q function]

W.K.T the Bellman's optimal value and of the are

$$V^*(s) = \max_a \sum_{s'} P(s'|s,a) \left[ R(s,a,s') + \gamma V^*(s') \right]$$

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

⇓

$$\therefore Q^*(s,a) = \sum_{s'} P(s'|s,a) \left[ R(s,a,s') + \gamma V^*(s') \right]$$

$$\therefore V^*(s) = \max_a Q^*(s,a)$$

& if we know the Q-values of all $(s,a)$ pairs of a particular state, s, Using this, we can find the optimal value fu of that state.

Ex: Q-values of all state-action pairs

State are :-

| State | Action | Value |
|-------|--------|-------|
| $S_0$ | 0. | 2.7 |
| $S_0$ | 1 | 3 |
| $S_1$ | 0 | 4 |
| $S_1$ | 1 | 2 |

Using this, we can find the optimal state values of each state. as

| State | value |
|-------|-------|
| $S_0$ | 3 |
| $S_1$ | 4 |

This is the outline of the value itr algo.

The value iteration Algom:-

<u>Step1</u>: Compute the optimal value fn. of each state iteratively, by taking the max of the Q fns (of all possible actions ~~oper~~ in a state), ie,

$$V^*(s) = \max_a Q^*(s,a)$$

<u>Step2</u>: Extract the optimal policy from the computed $\overset{optimal}{\underset{\wedge}{value}}$ fn.

<u>An informal explan of the value itn algm:-</u>

1. Initialize the value table of all states to zero.

2. ~~Compute & for~~ each state: s, do:

   ~~2.1 find the Q-fn of all possible actions~~

   2.1 for each possible action a in state s do:

      2.1.1 find the Q-value of this (s,a) pairs as

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

      (or)

      $$R(s,a) = \sum_{s'} P_{ss'}^a \left[R_{ss'}^a + \gamma V(s')\right]$$

Pg: (16)

2.1.2    find the max among these Q values
and update that as the Value of
that state    $V^*(s) = \max\limits_{a} Q^*(s,a)$

3. If the value table of 2 consecutive iteration
doesn't change, then, goto the next step to
find the optimal policy.; Else repeat
step 2, using the updated Value table of this recent
iteration.

4. find the Q value of each (s,a) pair using
the optimal value table obtained in step 3
using:

$$Q(s,a) = \sum\limits_{s'} P_{ss'}^{a} \left[ R_{ss'}^{a} + \gamma V(s') \right]$$

5. Extract the optimal policy from the Q-value
table got in step 4.

Ex: Given



We are.

Goal: From state A, reach state C, without visiting B.

2 actions:
0 - left/right
1 - up/down

What is the optimal policy? Ans: perform action 1 in state A.

Given the model dynamics of state A:

| State (s) | Action (a) | Next state (s') | Transn Prob $P_{ss'}^a$ | Reward for a $R_{ss'}^a$ |
|---|---|---|---|---|
| A | 0 | A | 0.1 | 0 |
| A | 0 | B | 0.8 | -1 |
| A | 0 | C | 0.1 | 1 |
| A | 1 | A | 0.1 | 0 |
| A | 1 | B | 0.0 | -1 |
| A | 1 | C | 0.9 | 1 |

Step1: Compute the optimal value function:
 w.k.t. $V^*(s) = \max_a Q^*(s,a)$    1) Initialise value table to
∴ for each state $s$, find $Q(s,a) = \sum_{s'} P_{ss'}^a \left[ R_{ss'}^a + \gamma V(s') \right]$
2)    for all possible actions

3) find the max $Q^*(s,a)$, update the value table.

4) Repeat steps (2) & (3) until convergence

Pg: 18

~~Iteration~~ Iteration 1:

Initialize value table to zero

| State | Value |
|-------|-------|
| A | 0 |
| B | 0 |
| C | 0 |

Soln: In state A: Possible actions are (0, 1)

∴ find $Q(A, 0) = P_{AA}^{0} \{ \cancel{R_{AA}^{0} + \gamma V(A)} \} +$

$Q(A, 1) =$

To using $Q(s, a) = \sum_{s'} P_{ss'}^{a} [ R_{ss'}^{a} + \gamma V(s') ]$

∴ $Q(A, 0) = P_{AA}^{0} [ R_{AA}^{0} + \gamma V(A) ] + P_{AB}^{0} [ R_{AB}^{b} + \gamma V(B) ]$

$\quad + P_{AC}^{0} [ R_{AC}^{0} + \gamma V(C) ]$

$= 0.1 [ 0 + 0 ] + 0.8 [ -1 + 0 ] + \cancel{\xi} 0.1 \begin{bmatrix} 1 + \\ 0 \end{bmatrix}$

$= \cancel{\phantom{0}} 0.8 + 0.1 = -0.7$

$Q(A, 1) = P_{AA}^{1} [ R_{AA}^{1} + \gamma V(A) ] + P_{AB}^{1} [ R_{AB}^{1} + \gamma V(B) ]$

$\quad + P_{AC}^{1} [ R_{AC}^{1} + \gamma V(C) ]$

$= 0.1 [ 0 + \cancel{\phantom{0}} 0 ] + 0 [ -1 + 0 ] + 0.9 [ \cancel{0} + 0 ]$

$= 0 + 0 + 0.9 = 0.9$

∴ updated value table:

| State | Value |
|-------|-------|
| A | 0.9 |
| B | 0 |
| C | 0 |

for State B:

$R(B,0) = P$

$R(B,1) =$

table at end of itn1 (assume) is:-

Similarly, if we know the model dynamics of state B & C we can find $Q[B,0) \; Q(B,1) \; Q(C,0)$ and $R(C,1)$. finally updated state value

| State | Value |
|-------|-------|
| A | 0.9 |
| B | -0.2 |
| C | 0.5 |

Stn2: Using the value table of itn1, find R, values of all states w.r.to all possible actions in each state:

Ex: On state A:-

$$R(A,0) = P_{AA}^{0}[R_{AA}^{0} + \gamma V(A)] +$$
$$P_{AB}^{0}[R_{AB}^{0} + \gamma V(B)] +$$
$$P_{AC}^{0}[R_{AC}^{0} + \gamma V(C)]$$

$$= 0.1(0+0.9) + 0.8(-1\overset{-0.2}{\cancel{+0}}) + 0.1(1+0.5)$$

$$= -0.72$$

$$R(A,1) = P_{AA}^{1}[R_{AA}^{1} + \gamma V(A)] + P_{AB}^{1}[R_{AB}^{1} + \gamma V(B)] +$$
$$P_{AC}^{1}[R_{AC}^{1} + \gamma V(C)]$$

$$= 0.1(0+0.9) + 0.0(-1-0.2) + 0.9(1+0.5)$$

$$= 1.44$$

Similarly we find the Q-value of all states & update the state table as [Assumption]

Value table from itn2.

| State | Value |
|-------|-------|
| A | 1.44 |
| B | -0.50 |
| C | 1.0 |

**Sn 3:** Repeat the same steps. While computing Q values, use the updated value table got from the previous itn.

Assume value table from itn3

| State | Value |
|-------|-------|
| A | 1.94 |
| B | -0.70 |
| C | 1.3 |

Convergence?

Keep repeating, until the value table b/t 2 consecutive itns doesn't change or changes by a very small (fraction based on a threshold)

Assume value table from itn4 is

| State | Value |
|-------|-------|
| A | 1.95 |
| B | -0.72 |
| C | 1.3 |

i.e. the diff is small, we take this as the optimal value table.

Next, step 2: to extract the optimal policy from this optimal value table.

Step2- Extract the optimal policy from the
optimal value table (function) from step1.

Assume Optimal value table is

| State | Value |
|-------|-------|
| A | 1.95 |
| B | -0.72 |
| C | 1.3 |

Use the $\hat{R}$ fn to compute the
for a state s,                              Policy.
2.1) find the $Q$-value of all $(s,a)$ pairs as:

$$\hat{Q}(s,a) = \sum_{s'} P_{ss'}^a [R_{ss'}^a + \gamma V(s')]$$

[Use the optimal value fn to find $V(s')$]
                                    of step1

2.2) Extract the optimal policy, by selecting the action
that has the max $Q$ value in a state.

$$\Pi^* = \underset{a}{argmax}\ Q(s,a)$$

for state A :-

$$Q(A,0) = P_{AA}^0 [R_{AA}^0 + \gamma V(A)] + P_{AB}^0 [R_{AB}^0 + \gamma V(B)] +$$
$$P_{AC}^0 [R_{AC}^0 + \gamma V(C)]$$

$$= 0.1 [0 + 1.95] + 0.8(-1 -0.72) +$$
$$0.1 [1 + 1.3]$$

$$= -0.951$$

$$Q(A,1) = 2.26.$$

w.k.t. the optimal value fn is and Q fn is

$$V^*(s) = \max_a \sum P(s'|s,a) \left[ R(s,a,s') + \gamma V^*(s') \right]$$

w.k.t, given value fn, we can derive the Q fn..

$$Q^*(s,a) = \sum_{s'} P(s'|s,a) \left[ R(s,a,s') + \gamma V^*(s') \right]$$

∴ Q table is

| State | Action | Value |
|-------|--------|-------|
| A | 0 | −0.95 |
| A | 1 ✓ | 2.26 ✓ |
| B | 0 | −0.5 |
| B | 1 ✓ | 0.5 ✓ |
| C | 0 | −1.1 |
| C | 1 ✓ | 1.4 ✓ |

from this Q table, pick the action in each state that has the max. value as an optimal policy.

∴ In state A, the optimal policy is action 1, ie, moving down..

# Solving the frozen lake problem with value itn.

| | | | |
|---|---|---|---|
| S | F | F | F |
| F | H | F | H |
| F | F | F | H |
| H | F | F | G |

← Start (S)
← Goal

S .. starting state

F .. frozen state

H ... hole state

G .. goal state

states are encoded from 0 to 15.

Actions:
```
left .. 0
down   1
right  2
up .. 3
```

Step1. Compute the optimal value fn.

```
def value_itn (env):
    num_itns = 1000
    threshold = 1e-20
    gamma = 1.0
    value_table = np.zeros (env.observation_space.n)
    for i in range(num_itns):
        updated_value_table = np.copy (value_table)
        for s in range(env.observation_space.n):
```

24(38)

```python
Q_values = [
              sum([Prob * (r + gamma * updated_value_table[s]
              for prob, s_, r, _ in      env.P[s][a])
                   for a in range(env.action_space.n)]
value_table[s] = max(Q_values)

if (np.sum(np.fabs(updated_value_table -
                   value_table)) <= threshold):
    break

return value_table.

```

step2: finding optimal policy. -

```python
def extract_policy(value_table):
    gamma = 1.0
    policy = np.zeros(env.observation_space.n)
    for s in range(env.observation_space.n):
        Q_values = [sum([Prob * (gamma * value_table[s]
                    for prob, s_, r, _ in env.P[s][a]])
                    for a in range(env.action_space.n)]
        policy[s] = np.argmax(np.array(Q_values))   Pg:25
    return policy
```
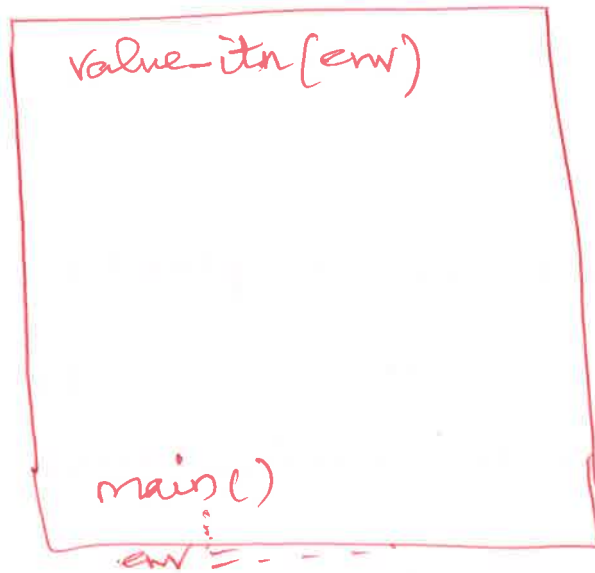
Flow diagram...

value_itn (env)

main ( )

env = - - - -

optimal_ = value_itn (env)

def value_iteration (env) :
    :
    return optimal_value f ~~value_table~~

def extract_policy ( value_table

    return policy

main ( )

    // initialize ~~the~~ env

    optimal_value_ = value_iteration (env)
        fn

    optimal_policy = extract_policy (optimal
                                value_fn )
    print (optimal_policy)

Policy iteration Method :-

~~Compute the optimal value fn, using~~

In the value itn method :

   1. Compute the optimal value fn, by taking the max over the $Q$ fn (Q values)

   2. Extract the optimal policy from the optimal value fn, got in step 1.

In the policy itn method :

   1. Compute the optimal value fn iteratively, by using the policy.

   2. Extract the optimal policy from the optimal value fn, got in step 1. [ this is the same policy that generated the optimal value fn ]

How to find the value fn of a state for a given policy $\pi$ ?

If $\pi$ is stochastic :

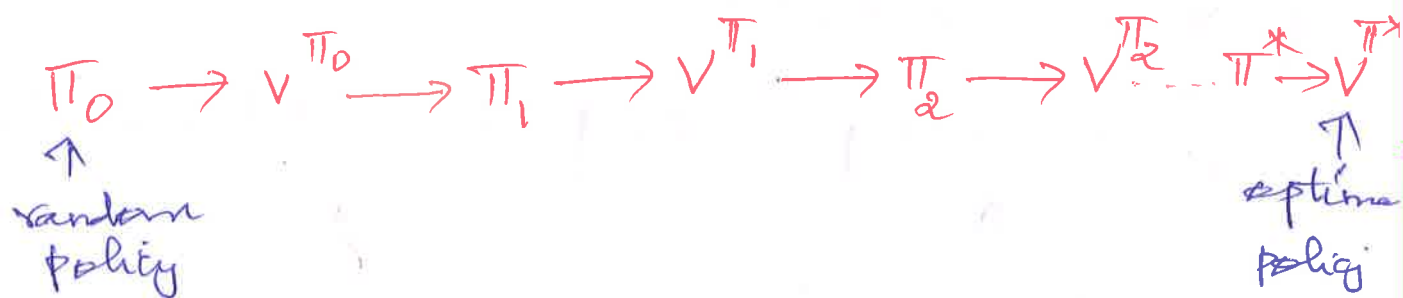$$V^{\pi}(s) = \sum_{a} \pi(a|s) \sum_{s'} P(s'|s,a) \left[ R(s,a,s') + \gamma V^{\pi}(s') \right]$$

If $\pi$ is deterministic :

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

(or) $$V^{\pi}(s) = \sum_{s'} P_{ss'}^{a} \left[ R_{ss'}^{a} + \gamma V^{\pi}(s') \right]$$

Assume the policy is deterministic!

1. Start with a random policy $\pi_0$.

2. Initialize the value fn (table) of all states to zeros.

3. Use $\pi_0$, to find an optimal value table iteratively $V^{\pi_0}$. [This will not be optimal ~~policy~~ value table, since it is generated from a random policy]

4. Use $V^{\pi_0}$ to generate a policy $\pi_1$.

5. Compute $V^{\pi_1}$ using $\pi_1$. Check if $V^{\pi_1}$ is optimal. If so, stop. Else generate $\pi_2$ from $V^{\pi_1}$.

6. Use $\pi_2$ to generate a $V^{\pi_2}$. If $V^{\pi_2}$ is optimal stop. Else generate a new policy $\pi_3$... and so on.

$$\pi_0 \rightarrow V^{\pi_0} \rightarrow \pi_1 \rightarrow V^{\pi_1} \rightarrow \pi_2 \rightarrow V^{\pi_2} \ldots \pi^* \rightarrow V^{\pi^*}$$

$\uparrow$
random
policy

$\uparrow$
optimal
policy

How to decide that the value fn is optimal?

If it doesn't change over iterations?

Since the value fn is generated from a policy,
over a series of iterns,
if the policy doesn't change bet^n 2 consecutive
iterations, then it is an optimal policy $\pi^*$.

The value fn generated from $\pi^*$ is the optimal
value fn $V^{\pi^*}$.

Pseudo-code :-

```
Policy = random-policy
for i is range ( num_itns):

    value_fn = compute-value-fn ( policy)

    new_policy = extract-policy ( value_fn)
    if policy == new_policy :
        break
    else
        Policy = new-policy.
```
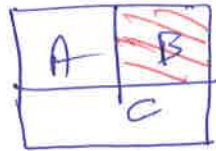
Algorithm for policy itrn.-

1. Initialize a random policy
2. Compute the value fn iteratively, from this policy
3. Extract a new-policy from the value fn of step(2)

4. If the extracted policy is the same as the policy generated in step 2, then stop, else goto step(2) with this new-policy & repeat steps 2 to 4.

EX: [Pg: 118]



States: A ... 0
        B ... 1
        C ... 2

Actions: 0 - left/right; 1 - up/down.

Goal: from A, reach C, without visiting B.

Given model dynamics of State A.

| State s | Action a | Next state s' | $P^a_{ss'}$ | $R^a_{ss'}$ |
|---|---|---|---|---|
| A | 0 | A | 0.1 | 0 |
| A | 0 | B | 0.8 | -1 |
| A | 0 | C | 0.1 | 1 |
| A | 1 | A | 0.1 | 0 |
| A | 1 | B | 0.0 | -1 |
| A | 1 | C | 0.9 | 1 |

Step1: - Initialize a random policy: -

A → 1 ; B → 0 ; C → 1.

Step2: Compute the value fn using this random policy.

$$V^{\pi}(s) = \sum_{s'} P^a_{ss'} [ R^a_{ss'} + \gamma V^{\pi}(s') ]$$

Step ③: Initialize the value table of all states to zero.

| State | Value |
|-------|-------|
| A | 0 |
| B | 0 |
| C | 0 |

Initial value table

$$\gamma = 1$$

Step ④: Itr 1: find the $V(s)$ for a state, only for those actions mentioned in the policy, $\pi$.

$\therefore$   $V(A)$ ·· only for action 1.

$\therefore V(A) = P^1_{AA}\left[R^1_{AA} + \gamma V(A)\right] + P^1_{AB}\left[R^1_{AB} + \gamma V(B)\right]$

$$+ P^1_{AC}\left[R^1_{AC} + \gamma V(C)\right]$$

$$= 0.1\left[0 + 0\right] + 0\left[-1 + 0\right] + 0.9\left[1 + 0\right]$$

$$= 0.9$$

for state B, find only action 0, if we know the model dynamics

$\therefore V(B) = P_B$

Clearly from the model dynamics of states B and C, find $V(B)$ for action 0 and $V(C)$ for action 1.

Assume the value table from itr1 is

| State | Value |
|-------|-------|
| A | 0.9 |
| B | −0.2 |
| C | 0.1 |

This will not be optimal.

Itn2: find $V(S)$ using the value table of from

Itn1. the P

$$V(A) = P'_{AA^*} [R'_{AA^*} + \gamma V(A)] + P'_{AB} [R'_{AB} + \gamma V(B)]$$

$$+ P'_{AC} [R'_{AC} + \gamma V(C)]$$

$$= 0.1 [0 + 0.9] + 0 [-1 - 0.2] + 0.9 [1 + 0.1]$$

$$= 1.08$$

Clearly keep & find $V(B)$ and $V(C)$. Assume value table from itn2 is

| State | Value |
|-------|-------|
| A | 1.08 |
| B | ~0.5 |
| C | 0.5 |

value table from itn3 is

| | |
|---|---|
| A - | 1.45 |
| B | -0.9 |
| C - | 0.6 |

value table from itn4 is

| | |
|---|---|
| A | 1.45 |
| B | -0.9 |
| C - | 0.6 |

← is

∴ No change the optimal final value table from the

Step3: Extract a new policy using this value &
from step2.

∵ policy ~~is general~~ tells us which is the correct action in each state. This is given by

$$Q(s,a) = \sum_{s'} P_{ss'}^a \left[ R_{ss'}^a + \gamma V[s'] \right]$$

for all $(s,a)$ pairs of a state $s$.

for state A :-

∴ find $Q(A,0)$ : $P_{AA}^0 \left[ R_{AA}^0 + \gamma V(A) \right] + P_{AB}^0 \left[ R_{AB}^0 + \gamma V(B) \right]$

$+ P_{AC}^0 \left[ R_{AC}^0 + \gamma V(C) \right]$

$= 0.1 \left[ 0 + 1.46 \right] + 0.8 \left[ -1 - 0.9 \right) +$

~~$Q(\rightarrow)$~~  $0.1 \left[ 1 + 0.61 \right]$

$= -1.21.$

$Q(A,1) = P_{AA}^1 \left[ R_{AA}^1 + \gamma V(A) \right] + P_{AB}^1 \left[ R_{AB}^1 + \gamma V(B) \right]$

$+ P_{AC}^1 \left[ R_{AC}^1 + \gamma V(C) \right]$

$= 0.1 \left[ 0 + 1.46 \right] + 0.0 \left[ -1 - 0.9 \right] + 0.9 \left[ 1 + 0.61 \right]$

$= 1.59$

Clearly do for states B and C for all actions.

∴ Q table is

| State | Action | Value |
|---|---|---|
| A | 0 | -1.21 |
| A | 1 | 1.59 |
| B | 0 | 0.1 |
| B | 1 | 0.0 |
| C | 0 | 0.5 |
| C | 1 | 0.0 |

P.T.(33)

from this Q table, pick the action in each state, that gives the max. value. This gives a new policy.

$$A \rightarrow L; \quad B \rightarrow D; \quad C \rightarrow D.$$

Step4: check the new policy

If The extracted new policy from step (3) is same as the policy used in step(2), then stop: else goto step2 with this new policy & repeat steps 2 to 4.

---

Solving the Frozen Lake Problem with policy iteration

main()

```
import gym
import numpy as np
env = gym.make('FrozenLake-v0')
optimal_policy = policy_itn(env)
print(optimal_policy)
```

```python
def policy-itn (env):
    num-itns = 1000
    // initialize policy to zero in all states
    policy = np.zeros ( env.observation_space.n)

    for i in range ( num-itns):

        val-fn = compute-val-fn( policy)

        new-policy = extract-policy ( val-fn)

        if (np.all ( policy == new-policy)):

            break

        policy = new-policy

    return (policy)

def compute-val-fn (policy):
    num-itns = 1000
    threshold = 1e-20
    gamma = 1.0
    // init. val. table of all states to zero
    val-tab = np.zeros ( env.observation_space.n

    for i in range ( num-itns):

        updated-val-tab = np.copy (val-tab)

        for s in range (env.observation_
                                    space.n):
            // find value of a state only for action a as in
            a = policy [s]                       the policy
```

$$// \quad V^{\pi}(s) = \sum_{s'} P^a_{ss'} \left[ R^a_{ss'} + \gamma V^{\pi}(s') \right]$$

val-tab[s] = sum(

[prob * (r + gamma * updated-val-tab[s_])

for prob, s_, r, _ in env.P[s][a]])

if (np.sum(np.fabs(updated-val-tab - val-tab)) <=

threshold:

break

return (val-tab)


def extract-policy( val-tab):

// policy is generated using Q(s,a), for all
// actions 'a' in a state 's'; not with resp to any π.

$$// \quad Q(s,a) = \sum_{s'} P^a_{ss'} \left[ R^a_{ss'} + \gamma V(s') \right]$$

gamma = 1.0
// init. policy of all states to zero.
policy = np.zeros(env.observation_space.n)

~~for i in range( num~~

for s in range( env.observation_space.n):

Q-values = [sum([prob * (r + gamma *

val-tab[s_])

for prob, s_, r, _ in env.P[s][a]])

for a in range( env.action_space.n)]

pg. 36

`policy [S] = np.argmax ( np.array (Q-Values))`
// this returns the index of max Q-value ~~that~~
return policy.

**❈ Limitations of DP :-**

DP is a model-based method to find the optimal policy. In both value & policy itn methods, we need to know the model dynamics [transn probabilities] and reward functions.

In value iteration:-

$$V^*(s) = \max_a Q^*(s,a)$$

1. find the optimal val fn of a state 's', by finding the max over all Q functions (s,a) of all actions in state 's'.

$$Q(s,a) = \sum_{s'} P_{ss'}^a \left[ R_{ss'}^a + \gamma V(s') \right]$$

2. Extract the optimal policy from this optimal val fn

In Policy iteration:-

1. find the optimal val fn of a state 's', by ~~finding the~~ using the policy iteratively

$$V^\pi(s) = \sum_{s'} P_{ss'}^a \left[ R_{ss'}^a + \gamma V^\pi(s') \right]$$

Start with a random policy and find the value fn.

2. find The policy that generated this optimal value fn, is the optimal policy.

∴ in both methods, we should know $P_{ss'}^a$.

To find the optimal policy ~~is a stochastic event~~, without ^knowing model dynamics, we use model–free methods ^like the Monte Carlo methods.