

④

Project Task

- ▷ What is optimality?
- ⇒ It means choosing actions on which you will be able to maximise your rewards.
- In a word it is about making the best decisions at every step.

Optimality

Policy optimality

→ A policy π^* will be optimal if it provides the highest expected return.

$$\rightarrow \pi^* > \pi$$

$$\boxed{H^* = \arg \max_{\pi} V^{\pi}(s)}$$

○ π^* → Policy

○ V^{π^*} → Value func under Policy π^*

Value Func optimality

→ $v^*(s)$ - The max possible expected return
 → $\phi^*(s, a)$ - The max expected return from taking action a in state s & acting optimally.

$$H^* = \begin{cases} 1, & \text{if } a = \arg \max v^{\pi}(s) \\ 0, & \text{otherwise} \end{cases}$$

2) Expectation

→ In a word we can say, it is a way to handle uncertainty

It gives us the average outcomes we can expect from a random event.

So,

$$E[x] = \sum x \cdot P(x)$$

$P(x) \rightarrow$ Probability distribution
 $x \rightarrow$ ~~Random~~ Variable

Bellman Expectation Eqn

$$V^*(s) = E_x [R_{t+1} + \gamma V^*(s_{t+1}) \mid s_t = s]$$
 or → Distribution Factor
○ $E_x[\cdot] \rightarrow$ Expected value



Bellman Optimal Eqn

$$V^*(s) = \max_a E [R_{t+1} + \gamma V^*(s_{t+1}) \mid s_t = s, A_t = a]$$

expected return for each possible action a . Then the action which will give maximum expected return is chosen.

Now,

$$q_{\pi}(s, a) = r(s, a) + \gamma \sum_{s'} P(s'|a, s) \sum_{a' \in A} \pi(a'|s') q_{\pi}(s', a')$$

$$\begin{aligned} q_{\pi}^*(s, a) &= r(s, a) + \gamma \sum_{s'} P(s'|a, s) \underbrace{\sum_{a' \in A} \pi^*(a'|s') q_{\pi}^*(s', a')}_{\max_{a'} q^*(s', a')} \\ &= r(s, a) + \gamma \sum_{s'} P(s'|a, s) \max_{a' \in A} q^*(s', a') \end{aligned}$$

Exercise | Page no → 67
3.25

* Equation of $v^*(s, a)$ in terms of $q^*(s, a)$ is

$$v^*(s, a) = \max q^*(s, a)$$

$v^*(s, a) \rightarrow \text{Max}^m \text{return}$

$q^*(s, a) \rightarrow \text{Expected return}$

So, from here we can say that

$$\begin{aligned} v^*(s) &= \max \text{return form } (s, a) \\ &= \max q^*(s, a) \end{aligned}$$

3.26

When starting in state s , taking action a & following optimal policy π^* ,

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

The four-argument transition func

$$q^*(s, a) = \sum P(s', r \mid s, a) \cdot [r + \gamma V^*(s')]$$

3.27

① We know,

$$V^*(s) = \max q^*(s, a)$$

A/c to the optimal policy π^*

$$\pi^*(s) = \text{action } a \text{ that gives } q^*(s, a) = V^*(s)$$

$$\pi^*(s) = \arg \max q^*(s, a)$$

3.28

We know,

$$q^*(s, a) = \sum_{s'} p(s' \mid s, a) [r(s, a, s') + \gamma V^*(s')]$$

&

$$\pi^*(s) = \arg \max q^*(s, a)$$

So, we can say

$$v^*(s) = \arg \max \sum_{s'} p(s'|s,a) [r(s,a,s') + \gamma v^*(s')]$$

3.28 3.29

Bellman Expectⁿ for $v^*(s)$:-

$$v^*(s) = \sum a \pi(a|s) \sum s' p(s', r|s, a) [r + \gamma v^*(s')]$$

Bellman Optimality

3.29

3-argument transition probn

$p(s'|s,a)$: the probⁿ of transitioning $s' \rightarrow s$ after action a

so,

i) Bellman expectⁿ for v^*

$$v^*(s) = \sum a \pi(a|s) [r(s,a) + \gamma \sum s' p(s'|s,a) v^*(s')]$$

ii) Bellman Optimality eqn for v^*

$$v^*(s) = \max [r(s,a) + \gamma \sum s' p(s'|s,a) v^*(s')]$$

(iii)

Bellman Expectation

Bellman Expectation Eqn $\hat{q}^n(s, a)$

$$\hat{q}^n(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) \sum_{a'} \hat{q}^n(a'|s')$$

(iv)

Bellman Optimality Eqn $q^*(s, a)$

$$q^*(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) \max q^*(s')$$

Part 2

OK, now I am ~~now~~ considering that, I am ~~current~~ currently in a ~~cave~~ jungle and try to find my ~~out~~ ancestral treasure. which is ~~under a tree~~ in a room.

So in a ^{cave} jungle there are lots of rooms and to reach ~~the~~ there, I can go in different directions. Some paths lead to dead end some to the treasure. Now I have a map of the ^{cave} jungle but I don't know which

direction or action lead to the best outcome overall.

So for me the challenges are to the value of the rooms of the cave that will decide my future rewards.

→ Value Iteration

Here I start with guess. For each room, I consider every possible move. For each move I will calculate that if I take that action, I will get an immediate reward $R(s, a)$ & the discounted value of the room I land in next, $r \sum_{s'} P(s'|s, a) V(s')$.

From here I will pick the action that gives me the highest combined value.

So, V-I eqn is

$$V_{k+1}(s) = \max \left[R(s, a) + r \sum_{s'} P(s'|s, a) V_k(s') \right]$$

Now I repeat this over & over, updating my estimate of how good each room is. Eventually, I will be able to find the true best values v^*

So,

$$\pi^*(s) = \arg \max [R(s,a) + r \sum_{s'} P(s'|s,a) v^*(s')]$$

Policy Iteration

① Evaluation of my current Policy

So I follow my current plan using

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

so, by solving this, I will get a reliable estimate of how good my current behavior is.

② Some Improvement

So, now I am thinking that it will be quite better if look at each room & check that my action will lead to the highest expected return right now or would that be different from my current plan.

So after all that stuff, I have just updated my policy by:

$$\pi_{\text{new}}(s) = \arg \max [R(s,a) + r \sum_{s'} P(s'|s,a) V^{\pi}(s')]$$

Why my ideas will work.

Here each update uses the ~~the~~ best possible action at that step.

Evaluating my current policy gives you precise knowledge & how good it is. since there are only finitely many policies, so I eventually will be able to reach the best one.

→ G_tPI is basically ~~G~~ Generalized Policy Iteration. So it is a framework how policy ~~eval~~ evaluation & Policy improvement work.

~~G~~ So now let's be my policy & ~~v~~ be my value function.

Here in G_tPI I will evaluate current Policy & improve my policy based on ~~my~~ my current value estimates.

Value Iteration

A/c to Bellman optimality operator

$$V_{k+1}(s) = \max \left[R(s,a) + \gamma \sum P(s'|s,a) V_k(s') \right]$$

Now, let

$$V_{k+1} = T^* V_k$$

| \$T^*\$ operator \$\Rightarrow\$ Discounted factor
 | \$R(s,a)\$ \$\Rightarrow\$ Reward func
 | \$P(s'|s,a)\$ \$\Rightarrow\$ Transition func
 | Finite action \$\rightarrow A\$.

so,

$$(T^* V)(s) = \max \left[R(s,a) + \gamma \sum_{s'} P(s'|s,a) V(s') \right]$$

Now, let \$V\$ and \$W\$ be two value func
We ~~show~~ have to show,

$$\|T^* V - T^* W\|_\infty \leq \gamma \|V - W\|_\infty$$

Let's look at states, -

$$|T^* V(s) - T^* W(s)| = \left| \max_a Q_V(s,a) - \max_a Q_W(s,a) \right| \leq \max_a |Q_V(s,a) - Q_W(s,a)|$$

where,

$$Q_V(s,a) = R(s,a) + \gamma \sum_{s'} P(s'|s,a) V(s')$$

$$Q_W(s,a) = R(s,a) + \gamma \sum_{s'} P(s'|s,a) W(s')$$

So :

$$|\phi_v(s,a) - \phi_w(s,a)| = \gamma \left| \sum_{s'} P(s'|s,a) (v(s') - w(s')) \right| \leq \gamma \|v-w\|_\infty$$

Policy Iteration

Hence

$$\|T^*v - T^*w\|_\infty \leq \gamma \|v-w\|_\infty$$

Policy

Now

So, Here $\gamma < 1$

Hence contraction is proved

Since here T^* is contracting on the complete metric space ~~(not)~~.

So,

Now, here a unique fixed pt v^* such that $T^*v^* = v^*$ & for any initial value v_0 the iteration $v_{k+1} = T^*v_k \rightarrow v^*$

Then,

So,

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

And

•

Policy Iteration

So let initial policy is π_0 . Policy Evaluation is V^{π_k} & Policy improvement is π_{k+1} .

Now let,

- V^{π_k} : the value func under policy π_k
- π_{k+1} : the greedy improvement of π_k

So,

$$Q^{\pi}(s, \pi'(s)) \geq V^{\pi}(s)$$

Then,

$$V^{\pi'}(s) \geq V^{\pi}(s), \forall s$$

And if inequality is strict in any state, then:

$$V^{\pi'}(s) > V^{\pi}(s) \text{ for at least one } s$$

This Shows :- Improving the policy makes it better or equal in value.

- The seqn V^{π_k} is Monotonic & Bounded

$$V^{\pi_{k+1}} \geq V^{\pi_k}$$

Because policy improvement is

greedy & value func are well-defined

→

Also,

$$V^{\pi_k} \leq V^* \text{ for all } k,$$

⇒ Total

∴ The seqn is monotonic increasing and bounded above, so it must converge

so

#

There are only finitely many policies:

$$|\Pi| = |A|^{|S|}$$

$$\Pi_{k+1} = \Pi_k$$

∴ Policy iteration $\rightarrow \pi^*$ in finite steps

$$V_{k+1}(s) = \max_a \left[R(s, a) + \gamma \sum_s P(s'|s, a) V_k(s') \right]$$

① $n = \text{number of states}$

② $|A| = m = \text{number of actions per state}$

③ $\gamma \in (0, 1) = \text{discount factor}$

Time CX

For each state s , loop over:

① Each action a : m times

② Each next s' : n transitions

$$\rightarrow O(n \cdot m \cdot n) = O(n^2m)$$

\Rightarrow Total iterations needed for ϵ -optimal value:

$$O\left(\frac{\log(1/\epsilon)}{1-\gamma}\right)$$

So overall time \propto

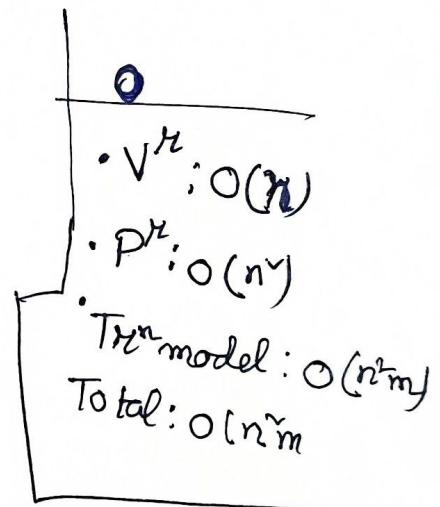
$$O\left(\frac{n^2m}{1-\gamma} \log \frac{1}{\epsilon}\right)$$

Policy Evaluation

$$V^R = R^R + \gamma P^R V^R$$

$$\Rightarrow V^R = (I - \gamma P^R)^{-1} R^R$$

Total Iteration



The number of possible deterministic policies is $|I|^I = m^n$

But in practice, P_I converges very fast
Overall time

$$O(K \cdot (n^3 + nm)) \quad [\text{usually small } K]$$