# Cauchy Sequence and Green's Equation

MDP: $S, A, p(s'|s,a), E(r|s,a,s'), \gamma$

Bellman eqn:

① — $v^{\pi}(s) = \sum_{a} \pi(a/s) \sum_{s'} p(s'|s,a) \left[ E(r|s,a,s') + \gamma v^{\pi}(s') \right]$

$v^{*}(s) = \max_{a} \sum_{s'} p(s'|s,a) \left[ E(r|s,a,s') + \gamma v^{*}(s') \right]$

Assumption: optimal policy is deterministic.

∴ Instead of $\pi(a/s)$; we will say $\pi(s)$

denotes a deterministic policy

↳ the action we take in state $s$

In fact, we can also say that if an MDP has an optimal policy, there exists atleast one deterministic optimal policy.

Finite MDP - S & A are both finite. Expectations ie
$$E(r|s,a,s')$$ are bounded.

Then we can think of
$v^\pi$ & as a vector with $|S|$ components.

$$\|v\| = \sup_{s \in S} |v(s)|$$

$$= \max (|v(s)|) \quad \forall s \in S$$

max norm,

$\|x\| = 0$ iff $x = 0$

$\|\alpha x\| = \alpha \|x\|$

$\|x + y\| \le \|x\| + \|y\|$

complete normed vector space:

- cauchy sequence: $x_1, x_2, x_3, \cdots$

For every +ve real $\epsilon > 0$, $\exists N \in Z^+$ s.t. $\forall m, n > N$

$$\|x_m - x_n\| < \epsilon.$$

Basically a sequence in which the successive elements are coming closer and closer to each other.

If every cauchy sequence in a normed vector space converges to a pt. in the vector space, then we call the vector space as a complete normed vector space

If every cauchy sequence is convergent, then the vector space is complete.

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

↳ reward expected - one step reward - starting from s & A following $\pi$.

$$P_\pi(j|s) = \sum_a \pi(a|s) P(j|s,a)$$

↳ prob. that I end in state $j$ in one step, starting from s, using policy $\pi$.

⊛ $\sum_{s'} P(s'|s, \pi(s)) . E(r|s, \pi(s), s')$ ⎤ for deterministic policy.

⊛⊛ $P(j|s, \pi(s))$ ⎦

$r_\pi \longrightarrow$ again $|S|$ dimensional vector.

$P_\pi \rightarrow |S| \times |S|$ dimensional stochastic matrix.

$0 \leq \gamma < 1$

↳ all values $\geq 0$
all rows sum up to 1
the each row is a discrete
prob. dist.

$$r_\pi + \gamma P_\pi \cdot v^\pi = v^\pi$$

↙ immediate reward for playing an action

→ payoff for ending in some state j.

$P_\pi$ determines where we land
$v^\pi$ determines the payoff
$v^\pi$ - like a terminal cost. It's the expected return staring from the state & follow-i policy $\pi$.

→ Say it's like a one-step problem, we get reward $r_\pi$ for the transition along the way but also a terminal cost $v^\pi$ for landing in some state that we did.

Look at it like: my decision making problem ended with taking one decision. we made the choice according to $\pi$. we want $v^\pi$ to be the total cost.

Essentially, we want to solve:

② $-$ $r_\pi + \gamma P_\pi \cdot v^\pi = v^\pi$. → which actually comes from algebraically simplifying ①.

↳ This is a matrix eqn.

$$\Rightarrow v^\pi = (I - \gamma P_\pi)^{-1} r_\pi$$

↓ eigenvalues are all 1

→ $P_\pi$ is stochastic
⇒ largest eigenvalue is 1
→ $\gamma P_\pi$ largest eigenvalue is $< 1$

∴ The entire matrix $I - \gamma P_\pi$ cannot have a 0 eigenvalue.

⇒ The determinant exists & matrix is invertible and non-zero

⇒ There is a unique soln. for $v^\pi$

Sometimes called the Green's equation.

Another way to work at this is, take ②. start with some $v_0^\pi$ - substitute and find $v_1^\pi$. sub this back and keep going, find $v_2^\pi, v_3^\pi \cdots$ and finally - it should converge to $v^\pi$ and that will be unique. we can show this.

## Banach Fixed Point Theorem:

$$L_\pi : \quad V \to V$$

space of all value fns

A complete normed vector space of all value fn.

$$L_\pi V = r_\pi + \gamma P_\pi V$$

some element in $V$. not necessarily a value fn. in $V$. when I say it's a value fn., I am implicitly saying that there is a policy for which that is the expected return which need not be the case.

This is what the Bellman eqn says.

$$L_\pi V^\pi = V^\pi \qquad \text{] This is what the Bellman eqn says.}$$

→ $V^\pi$ is a fixed point of $L_\pi$. We apply $L_\pi$ on it and it doesn't move.

Banach Fixed Point Theorem:

Suppose $U$ is a Banach space (complete normed vector space) and

$$P : U \to U \text{ is a } \underbrace{\text{contraction mapping}}_{Tu \text{ and } Tv \text{ will be close to each other than } u \text{ & } v \text{ were.}} ; \text{ then.}$$

∃ a unique $V^*$ in $U$, s.t. $Tv^* = V^*$ and for arbitrary $V^0$ in $U$), the sequence $\{V^n\}$ defined by

$$V^{n+1} = P v^n = P^{n+1} V^0 \quad \text{converges to } V^*$$

$$\pi'(1) = a_2$$
$$\pi'(2) = a_2$$

$$\pi(1) = a_1$$
$$\pi(2) = a_1$$

$$V^\pi(1) = 1 + \gamma V^\pi(2)$$
$$V^\pi(2) = 1 + \gamma V^\pi(1)$$
we get $V^\pi(1) = V^\pi(2) = \dfrac{1}{1-\gamma}$





This is for a deterministic policy. You can write the eqn and do the same thing for a stochastic policy as well.

So for any policy we choose the $\pi, \pi'$ and all (here we are only 4 possible), we will be on one of the vertices of the square. These 4 points is our $V$.

〈back to the theorem〉

$T$ is a contraction if

$$\|Tu - Tv\| \le \lambda \|u-v\| \; ; \quad 0 \le \lambda < 1 \quad \forall u, v \in U$$

we need to show that $L_\pi$ is a contraction mapping. Then the rest of the theorem just follows. And we get a lot of information. we will see. Because we know $v^\pi$ is a fixed point, automatically from the theorem, we get that repeatedly applying $L_\pi$ will take us to the fixed point that is also unique as per the theorem

↳ By the way, if $\lambda = 1$, $T$ need not be identity. Just the distance has to be preserved. We can swap $u$ & $v$ also.

Proof:

$$\| v^{n+m} - v^n \| \le \| v^{n+m} - v^{n+m/2} \| + \| v^{n+m/2} - v^n \|$$

(Δ inequality)

↳ This we can do for any three points. we don't need them to be in a sequence. So in fact, we can keep applying the Δ inequality to say: (no technique &

$$\| v^{n+m} - v^n \| \le \sum_{k=0}^{m-1} \| v^{n+k+1} - v^{n+k} \|$$

$$= \sum_{k=0}^{m-1} \| T^{n+k} v' - T^{n+k} v^0 \| \qquad \text{(from the theorem)}$$

$$\le \sum_{k=0}^{m-1} \lambda^{n+k} \| v' - v^0 \|$$

$$= \frac{\lambda^n (1 - \lambda^m)}{1 - \lambda} \| v' - v^0 \|$$

As $n$ & $m$ becomes large, this is going to become smaller and smaller. ∴ we can say the sequence $\{v^n\}$ is cauchy.

Because it's a Banach space, $\{v^n\}$ is convergent.

## Convergence Proof

$$0 \le \|Tv^* - v^+\| \le \|Tv^* - v^n\| + \|v^n - v^+\|$$
$$= \|Tv^* - Tv^{n-1}\| + \|v^n - v^+\|$$
$$\le \lambda\|v^* - v^{n-1}\| + \|v^n - v^+\|$$

We want to prove that $Tv^+ = v^+$.
But we are assuming that $v^+$ is the convergent pt - of $[v^n]$.

$\underbrace{}$ will go to 0 because $[v^n]$ is cauchy. As $n \to \infty$, it converges to $v^+$ and $\|v^* - v^{n-1}\| \to 0$

$\underbrace{}$ will go to 0

$\therefore$ $\|v^* - v^n\| \to 0$ and similarly $\|v^n - v^+\| \to 0$ too.

$\therefore \|v^* - v^n\| \to 0$ as $n \to \infty$

$$0 \le \|Tv^+ - v^*\| \le 0$$

$\Rightarrow Tv^+ = v^*$.

Let $u^*$ & $v^0$ be two fixed pt's

$$\|Tu^* - Tv^+\| \le \lambda\|u^* - v^+\|$$
$$\Rightarrow \|u^* - v^+\| \le \lambda\|u^* - v^*\| \quad \text{and} \quad 0 \le \lambda < 1$$

$\therefore$ For all such $\lambda$, the only possibility is

$$\Rightarrow \quad u^* = v^+.$$

So now we have only shown that the Banach fixed point theorem is true.

We still need to show that $L_\pi$ is a contraction.

Let $u$ & $v$ be in $V$

$$L_\pi u(s) = r_\pi(s) + \sum_{j \in S} \gamma P_\pi(j/s) u(j)$$
$$L_\pi v(s) = r_\pi(s) + \sum_{j \in S} \gamma P_\pi(j/s) v(j)$$

Let $L_\pi v(s) > L_\pi u(s)$

$$0 \le L_\pi v(s) - L_\pi u(s) \le r_\pi(s) + \gamma \sum P_\pi(j/s) v(j) - r_\pi(s) - \gamma \sum P_\pi(j/s) \cdot u(j)$$

$$= \gamma \sum_j p(j/s) (v(j) - u(j))$$

$$\le \gamma \|v - u\| \left(\boxed{\sum_j P(j/s)} \longrightarrow 1\right)$$

because of max norm

$$= \gamma \|v - u\|$$

$\longrightarrow$ PTO

Similarly, when $L_\pi v(s) < L_\pi u(s)$ we can do.

so, we get:
$$|L_\pi v(s) - L_\pi u(s)| \leq \gamma \|v - u\| \quad \#1$$

⟹ $L_\pi$ is a contraction. ⟶ Pointwise, it is being drawn close by $\gamma$, so even if the max difference would have some down.

Similarly, we have to do this proof for the optimally eqn.

---

## $L_\pi$ ⊕ convergence

$L_\pi$ is a contraction.
$V$ — space of all fn.s that are component-wise bounded
      space of bounded fns.                        ⟹ vector notation

$$\begin{cases} v^* = \max_\pi [r_\pi + \gamma p^\pi v^*] \\ v^*(s) = \max_a \{ E(r|s,a) + \gamma \sum_j (j|s,a) \cdot v^*(j)) \} \end{cases}$$

call this as some operator $L$ ⊕

$\therefore$ $Lv^\oplus \equiv \max_\pi [r_\pi + \gamma p^\pi v^\oplus]$

not a                claim is that $v^*$ is the fixed point of $L$:
linear tr.           claim is $L$ is a contraction:
of $v$.             ⟹ $v^*$ is a unique fixed pt. and if we keep
max is there.        applying $L$, we converge at $v^*$.
                     All this is true because $V$ is a Banach space.

Let $a_s^* \in \text{argmax}_a \{ E(r|s,a) + \gamma \sum p(j|s,a) v(j)) \}$ ⟶ $Lv(s)$
⟶ one of the best action to
    take at $s$.     $\leq (E(r|s,a_s^*) + \gamma \sum p(j|s,a_s^*) \cdot v(j)))$ — $\overline{\text{Pro.}}$ contd.

$0 \leq Lv(s) - Lu(s)$                         We will proceed in
$\underbrace{\qquad}$                          the same way we
to clarify: (for $L$ & $L_\pi$)                did 1st time.
$Lv(s)$ means:

$Lv$ is a fn.
$Lv(s)$ is the output of
the fn. for the argument $s$.
$L$ is not acting on $v(s)$
$L$ takes a fn. and outputs
a fn. $L$ takes $v$ and
outputs $Lv$.

⟶ Pto

$$0 \leq L_v(s) - L_u(s) \leq E(r/s, a_s^{*1}) + \gamma \sum p(j/s, a_s^{*1}) \cdot v(j)$$
$$- \left[ E(r/s, a_s^{*1}) + \gamma \sum p(j/s, a_s^{*1}) \cdot u(j) \right]$$

here in both uses, we are getting rid of the max by using $a_s^*$ as that if no best action. But in the second term, we are using the max corresponding to $v$ and not $u$. We are using $a_s^*$ in both cases. So the second term, is in a sense, not necessarily using the optimal action. The value will be $\leq$ the optimal value and hence when we subtract that term, we get the $\leq$ sign and not $a = $ sign.

$$= \gamma \sum p(j/s, a_s^*) [v(j) - u(j)]$$
$$\leq \gamma \|v - u\| \sum_{j} p(j/s, a_s^*) \qquad (\because \text{max norm})$$

$$0 \leq L_v(s) - L_u(s) \leq \gamma \|v - u\|$$

Illy we can do $L_u(s) - L_v(s)$

we finally get:
$$|L_v(s) - L_u(s)| \leq \gamma \|v - u\| \quad \forall s.$$

$\Rightarrow L$ is a contraction.

Now we have a method to solve an MDP!. Start with an arbitrary value fn, keep solving a for $v$ repeatedly. converge at $v^*$. From an optimal value fn, how to recover an optimal policy?

Run

$$a_s^* \in \underset{a}{\arg\max} \left[ E(r/s, a) + \gamma \sum_{j} p(j/s, a) v^*(j) \right]$$

↳ pick some action from the argmax set.   ↳ optimal value

$\longrightarrow$ PTO