

Proof of Gittins Index theorem

First, we prove that there always exist a **Markovian stationary policy** which is optimal. Markovian stationary policy:

- ▶ decision depends only on state: does not change with time
- ▶ mapping from state space to (distribution over) actions

$$a_t^* = \pi(s_{1,t}, \dots, s_{N,t})$$

In fact, Gittins index suggests further that optimal policy is of form

$$a_t^* = \arg \max_{a=1, \dots, N} G_a(s_{a,t})$$

Existence of Markovian stationary optimal policy is in fact true for more general “Markov Decision Processes”.

└ Marovian bandits, Gittins Index

└ Proof of Gittins Index theorem

First, we prove that there always exist a **Markovian stationary policy** which is optimal. Markovian stationary policy:

- ▶ decision depends only on state: does not change with time
- ▶ mapping from state space to (distribution over) actions

$$a_t^* = \pi(s_{1,t}, \dots, s_{N,t})$$

In fact, Gittins index suggests further that optimal policy is of form

$$a_t^* = \arg \max_{a=1, \dots, N} G_a(s_{a,t})$$

Existence of Markovian stationary optimal policy is in fact true for more general "Markov Decision Processes".

We want to maximize

$$\mathbb{E}\left[\sum_{t=1}^{\infty} \gamma^{t-1} r_{a_t}(s_{a_t}, t)\right]$$

where $a_t = \pi_t(\mathcal{H}_t)$, for any history and time dependent policy π .

Let $s_t = (s_{1,t}, \dots, s_{N,t})$. Define $r_a(s_t) := r_a(s_{a,t})$. For bounded rewards, and $\gamma < 1$, this is same as

$$\lim_{T \rightarrow \infty} \mathbb{E}\left[\sum_{t=1}^T \gamma^{t-1} r_{a_t}(s_t)\right]$$

(under above conditions, the limit exists, and by bounded convergence theorem, expectation and limit are interchangeable. In fact, if the reward is bounded, then the limit also exists for average reward criterion.)

Finite Horizon

Let us first consider any **finite horizon** T .

Define value of a state $V_{\ell, T}^{\pi}(s)$ as the expected discounted reward achieved by policy π given $s_{\ell} = s$, and policy π is followed from time ℓ to T .

$$V_{\ell, T}^{\pi}(s) = \mathbb{E}\left[\sum_{t=\ell}^T \gamma^{t-\ell} r_{\pi_t(H_t)}(s_t) \mid s_{\ell} = s\right]$$

Let us first consider any **finite horizon** T . Define value of a state $V_{\ell, T}^{\pi}(s)$ as the expected discounted reward achieved by policy π given $s_{\ell} = s$, and policy π is followed from time ℓ to T .

$$V_{\ell, T}^{\pi}(s) = \mathbb{E}\left[\sum_{t=\ell}^T \gamma^{t-\ell} r_{\pi_t(H_t)}(s_t) \mid s_{\ell} = s\right]$$

Below we derive a recursive formula to compute V_{ℓ}^{π} : Let $a_{\ell} = \pi_{\ell}(\mathcal{H}_{\ell})$.

$$\begin{aligned} V_{\ell, T}^{\pi}(s) &= \mathbb{E}\left[\sum_{t=\ell}^T \gamma^{t-\ell} r_{\pi_t(H_t)}(s_t) \mid s_{\ell} = s\right] \\ &= \mathbb{E}\left[r_{a_{\ell}}(s) + \sum_{t=\ell+1}^T \gamma^{t-\ell} r_{\pi_t(H_t)}(s_t) \mid s_{\ell} = s\right] \\ &= \mathbb{E}\left[r_{a_{\ell}}(s) + \sum_{s'} \sum_{t=\ell+1}^T \gamma^{t-\ell} \mathbb{E}[r_{\pi_t(H_t)}(s_t) \mid s_{\ell+1} = s'] p_{a_{\ell}}(s, s') \mid s_{\ell} = s\right] \\ &= \mathbb{E}\left[r_{a_{\ell}}(s) + \sum_{s'} \gamma \mathbb{E}\left[\sum_{t=\ell+1}^T \gamma^{t-\ell-1} r_{\pi_t(H_t)}(s_t) \mid s_{\ell+1} = s'\right] p_{a_{\ell}}(s, s')\right] \\ &= \mathbb{E}\left[r_{a_{\ell}}(s) + \gamma \sum_{s'} V_{\ell+1, T}^{\pi}(s') p_{a_{\ell}}(s, s')\right] \end{aligned}$$

Dynamic Programming for finite horizon

Claim: Let $V_{t,T}^*(s)$ be the value of state s from time t to T for optimal policy π^* , then

$$V_{\ell,T}^*(s) = \max_a \left\{ r_a(s) + \gamma \sum_{s'} V_{\ell+1,T}^*(s') p_a(s, s') \right\}$$

Note that this claim provides dynamic programming solution for finite horizon.

Also, above claim implies that optimal policy is a **Markovian** policy: only depends on current state and not past history. But for finite horizon, this is not stationary: value of next state depends on number of time steps remaining.

Claim: Let $V_{t,T}^*(s)$ be the value of state s from time t to T for optimal policy π^* , then

$$V_{t,T}^*(s) = \max_a \left\{ r_a(s) + \gamma \sum_{s'} V_{t+1,T}^*(s') p_a(s, s') \right\}$$

Note that this claim provides dynamic programming solution for finite horizon.

Also, above claim implies that optimal policy is a **Markovian** policy: only depends on current state and not past history. But for finite horizon, this is not stationary: value of next state depends on number of time steps remaining.

Here, we prove the dynamic programming equations. We prove that for any π, ℓ, s , $V_{\ell,T}^\pi(s) \leq V_{\ell,T}^*(s)$. Proof by backward induction on ℓ . For $\ell = T$, $V_{T,T}^\pi(s) = r_{a_T}(s) \leq \max_a r_a(s) = V_{T,T}^*(s)$. Now, assume above claim is true for $\ell + 1$. Then,

$$\begin{aligned} V_{\ell}^\pi(s) &= \mathbb{E}[r_{\pi(H_\ell)}(s) + \gamma \sum_{s'} V_{\ell+1}^\pi(s') p_{\pi(H_\ell)}(s, s')] \\ &\leq \max_a \left\{ r_a(s) + \gamma \sum_{s'} V_{\ell+1}^\pi(s') p_a(s, s') \right\} \\ &\leq \max_a \left\{ r_a(s) + \gamma \sum_{s'} V_{\ell+1}^*(s') p_a(s, s') \right\} \\ &= V_{\ell}^*(s) \end{aligned}$$

Infinite Horizon

Now, define

$$V^\pi(s) = \lim_{T \rightarrow \infty} V_{1,T}^\pi(s)$$

. From the dynamic programming equations, it is easy to see that $V_{\ell,T}^*$ depends only on the number of remaining time steps, i.e.,

$$V_{\ell,T}^*(s) = V^*(\ell - 1, T - 1) = \dots = V^*(1, T - \ell + 1)$$

A rigorous proof can be obtained by induction on ℓ . Therefore, for any ℓ ,

$$\lim_{T \rightarrow \infty} V_{\ell,T}^*(s) = \lim_{T \rightarrow \infty} V_{1,T-\ell+1}^*(s) = V^*(s)$$

Therefore, for infinite horizon:

$$V_{\ell,\infty}^*(s) = V^*(s) = \max_a \left\{ r_a(s) + \gamma \sum_{s'} V^*(s') p_a(s, s') \right\}$$

Above are called **Bellman equations**. Note that this implies that optimal policy is **Markovian and stationary**.