

# Notes on Reinforcement Learning

ERMO WEI

## Contents

Actor-Critic Method.....	2
Bellman Equation.....	2
Bellman Operator.....	2
Boltzmann Selection .....	2
Contraction .....	2
Data Efficiency for Policy Gradient Method.....	3
$\epsilon$ -greedy .....	4
Ergodic MDP .....	4
Experience Replay .....	5
Function Approximation .....	5
GLIE .....	6
Inverse Reinforcement Learning .....	6
Lipschitz Continuity .....	6
LSTD.....	6
Markov Decision Process.....	6
Metric Space.....	7
Monte Carlo Method.....	8
On-policy and Off-policy.....	8
Policy Gradient Reinforcement Learning .....	8
Policy Iteration .....	10
Stationary Distribution .....	10
Stochastic Game .....	10
Temporal Difference Method .....	10
Value Iteration .....	10
Vector Space .....	10

**Actor-Critic Method** Still working on it ...

**Bellman Equation** Due to the Markov property of MDP, we can define the value function of a policy recursively to be

$$\begin{aligned} V^\pi(s) &= E_\pi\{R_t | s_t = s\} \\ &= E_\pi\{\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | s_t = s\} \\ &= E_\pi\{r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} | s_t = s\} \end{aligned}$$

If policy is stochastic, then we have

$$\begin{aligned} V^\pi(s) &= \sum_a \pi(s, a) \sum_{s'} \Pr(s' | s, a) [R(s, a) + \gamma E_\pi\{\sum_{k=0}^{\infty} \gamma^k r_{t+k+2} | s_{t+1} = s'\}] \\ &= \sum_a \pi(s, a) \sum_{s'} \Pr(s' | s, a) [R(s, a) + \gamma V^\pi(s')] \\ &= \sum_a \pi(s, a) \{R(s, a) + \sum_{s'} \Pr(s' | s, a) \gamma V^\pi(s')\} \end{aligned}$$

If policy is deterministic, then the Value function of policy  $\pi$  at state  $s$  is

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

Thus, the Bellman optimality equation can be defined as

$$V^*(s) = \max_a \{R(s, a) + \gamma \sum_{s'} \Pr(s' | s, a) V^*(s')\}$$

Detail see (Barto, 1998)

**Bellman Operator** Consider a [MDP](#), Bellman Operator for policy  $\pi$  can be defined as a map from a value function to another value function

$$T^\pi : \mathbb{R}^N \rightarrow \mathbb{R}^N$$

where  $N$  is the cardinality of state space, i.e.  $N = |S|$ .

Still working on it ...

**Boltzmann Selection** Still working on it ...

**Contraction** The operator  $F$  (function  $f$ ) is a  $\alpha$ -contraction ( $0 \leq \alpha < 1$ , called [Lipschitz constant](#) for  $F$ ) with respect to some norm  $\|\cdot\|$  if

$$\forall X, \bar{X} : \|FX - F\bar{X}\| \leq \alpha \|X - \bar{X}\| \quad \text{or} \quad \|f(X) - f(\bar{X})\| \leq \alpha \|X - \bar{X}\|$$

- Theorem 1. The sequence  $X, FX, F^2X, \dots$  converges for every  $X$ . e.g.  $X, f(X), f(f(X)), \dots$  converges for every  $X$

Proof:

Useful fact : Cauchy sequences : If for  $x_0, x_1, x_2, \dots$ , we have that

$$\forall \epsilon, \exists K : \|x_M - x_N\| < \epsilon \quad \text{for} \quad M, N > K$$

then we call  $x_0, x_1, x_2, \dots$ , a Cauchy sequence.

If  $x_0, x_1, x_2, \dots$ , is a Cauchy sequence, and  $x_i \in \mathbb{R}^n$ , then there exists  $x^* \in \mathbb{R}^n$  such that  $\lim_{i \rightarrow \infty} x_i = x^*$ .

Proof: Assume  $N > M$ .

$$\begin{aligned}
\|F^M X - F^N X\| &= \|\sum_{i=M}^{N-1} (F^i X - F^{i+1} X)\| \\
&\leq \sum_{i=M}^{N-1} \|F^i X - F^{i+1} X\| && \text{by triangle inequality} \\
&\leq \sum_{i=M}^{N-1} \alpha^i \|X - FX\| && \text{use the condition in Theorem} \\
&= \|X - FX\| \sum_{i=M}^{N-1} \alpha^i \\
&= \|X - FX\| \frac{\alpha^M}{1 - \alpha}
\end{aligned}$$

As  $\|X - FX\| \frac{\alpha^M}{1 - \alpha}$  goes to zero for  $M$  going to infinity, we have that for any  $\epsilon > 0$  for  $\|F^M X - F^N X\| \leq \epsilon$  to hold for all  $M, N > K$ , it suffices to pick  $K$  large enough. Hence,  $X, FX, \dots$  is Cauchy sequence and converges.

- Theorem 2 (Banach fixed-point theorem).  $F$  has a unique fixed point  $X^*$  which satisfies  $FX^* = X^*$  and all sequences  $X, FX, F^2X, \dots$  converge this unique fixed point  $X$ .

Proof:

Suppose  $F$  has two fixed points. Let's say

$$\begin{aligned}
FX_1 &= X_1 \\
FX_2 &= X_2
\end{aligned}$$

this implies,  $\|FX_1 - FX_2\| = \|X_1 - X_2\|$ . At the same time we have from the contractive property of  $F$

$$\|FX_1 - FX_2\| \leq \alpha \|X_1 - X_2\|$$

Combining both gives us

$$\|X_1 - X_2\| \leq \alpha \|X_1 - X_2\|$$

Hence,  $X_1 = X_2$ .

Detail see (Conrad, 2014)

**Data Efficiency for Policy Gradient Method** In [policy gradient](#) method, when we try to evaluate the expected return of a policy  $\pi$ , we have to run the policy several time to be able to have reasonable estimation of how good the policy is. Let  $J(\theta)$  denote the expected return of policy  $\pi_\theta$ , and  $\tau$  denote the trajectory or rollout or history (these terms are interchangeable) of executing the policy  $\pi_\theta$ , then we know

$$\begin{aligned}
J(\theta) &= E_{P(\tau; \theta)} [\sum_{t=0}^H \gamma^t R(s_t, u_t) | \pi_\theta] \\
&= \sum_{\tau} P(\tau; \theta) R(\tau)
\end{aligned}$$

where  $P(\tau; \theta)$  is the probability of having a trajectory  $\tau$  by following policy  $\pi_\theta$  and  $R(\tau)$  is just the accumulated reward of that trajectory. In [policy gradient](#) method, after we evaluate the policy  $\pi_\theta$ , we may want to improve it and use a new policy. Thus, the sample we collected during the process of following policy  $\pi_\theta$  are discarded. However, it will preferable if we can reuse the data gathered of following one policy to estimate the value of following another policy. The method "likelihood ratio" estimation make this data reuse possible.

In practice, if we want to evaluate  $J(\theta)$ , we may want to draw the rollout samples from the distribution induced by policy  $\pi_\theta$ . After taking  $N$  samples  $\{\tau_0, \tau_1, \dots, \tau_N\}$ , we have a unbiased estimator:

$$\hat{J}(\theta) = \frac{1}{N} \sum_i R(\tau_i)$$

Imagine, however, instead of  $\pi_\theta$ , we only have  $\pi_{\theta'}$  available, we can do some trick like this

$$\begin{aligned}
J(\theta) &= \sum_{\tau} P(\tau; \theta) R(\tau) \\
&= \sum_{\tau} P(\tau; \theta) \frac{P(\tau'; \theta')}{P(\tau'; \theta')} R(\tau) \\
&= \sum_{\tau} P(\tau'; \theta') \frac{P(\tau; \theta)}{P(\tau'; \theta')} R(\tau) \\
&= E_{\tau'} \left[ \frac{P(\tau; \theta)}{P(\tau'; \theta')} R(\tau) \right]
\end{aligned}$$

Now, we can estimate the  $J(\theta)$  with respect to the the distribution induced by  $\pi_{\theta'}$ . This method of estimating one expectation with respect to another distribution is called "importance sampling". So, we can rewrite the  $\hat{J}(\theta)$  as

$$\hat{J}(\theta) = \frac{1}{N} \sum_i R(\tau_i) \frac{P(\tau; \theta)}{P(\tau'; \theta')}$$

Note, in the equation above, we have a term  $\frac{P(\tau; \theta)}{P(\tau'; \theta')}$ . Since we don't have the model of the world, it's not possible to directly compute  $P(\tau; \theta)$ . But if we expand the fraction term, we can see that

$$\frac{P(\tau; \theta)}{P(\tau'; \theta')} = \frac{\prod_{t=0}^T \pi_\theta(u_t | s_t)}{\prod_{t=0}^T \pi_{\theta'}(u_t | s_t)}$$

Thus, we should be able to calculate the likelihood  $\frac{P(\tau; \theta)}{P(\tau'; \theta')}$  for any two policies  $\theta$  and  $\theta'$ . Actually, we can have a mixture distributions, where  $P(\tau; \theta)$  is replaced by  $\frac{1}{N} \sum_j P(\tau_j | \theta_j)$  where  $\tau_j$  are trajectory of following  $\theta_j$ , then we can make use of multiple source of distributions.

**$\epsilon$ -greedy**  $\epsilon$ -greedy is a common used exploration strategy for Model-free reinforcement learning. This strategy can be seen as a combination of greedy strategy and random strategy. Every time when agent choose an action to perform, it choose greedily with probability  $1-\epsilon$  (exploitation), and randomly with probability  $\epsilon$  (exploration). However, based on the changes of  $\epsilon$  value, this method have two versions. First, the  $\epsilon$  value can decrease as the learning process goes on, this is the *decay exploration* method. Another one which are more commonly used is that we fix the value of  $\epsilon$ .

The difference between those methods is that, the first one can be **GLIE** if the  $\epsilon$ -value goes 0 eventually but the second one cannot. Suppose we have a counter of how many times a state have been visited,  $n_t(s)$  and a constant  $c$ . As long as the  $\epsilon$  value for the a state  $\epsilon_t(s) = \frac{c}{n_t(s)}$  where  $0 < c < 1$ , this method can be consider **GLIE**. However, in practice, we usually use fixed value for  $\epsilon$  and after the convergence of the estimation value of action, we switch to greedy selection.

**Ergodic MDP** An MDP is said to be ergodic if for each policy  $\pi$  the Markov chain induced by  $\pi$  is ergodic. We are giving the definition of Ergodicity below. But first, we will give some auxiliary definitions. Several definition of Markov chain:

- **Reducibility**  
A Markov chain is said to be irreducible if it is possible to get to any state from any state.
- **Aperiodicity**  
A state  $i$  has period  $k$  if any return to state  $i$  must occur in multiples of  $k$  time steps. For example, suppose it is possible to return to a state in  $\{6, 8, 10, 12, \dots\}$  time steps, then  $k$  would be 2, even though 2 does not appear in this list. If  $k = 1$ , then the state is said to be aperiodic: returns to state  $i$  can occur at irregular times.

- Recurrence

A state  $i$  is said to be transient if, given that we start in state  $i$ , there is a non-zero probability that we will never return to  $i$ . State  $i$  is recurrent (or persistent) if it is not transient. Recurrent states are guaranteed to have a finite hitting time.

Here we have the definition of Ergodicity.

A state  $i$  is said to be ergodic if it is aperiodic and positive recurrent. In other words, a state  $i$  is ergodic if it is recurrent, has a period of 1 and it has finite mean recurrence time. If all states in an irreducible Markov chain are ergodic, then the chain is said to be ergodic.

It can be shown that a finite state irreducible Markov chain is ergodic if it has an aperiodic state. **A model has the ergodic property if there's a finite number  $N$  such that any state can be reached from any other state in exactly  $N$  steps.** For example, we will have  $N = 1$  if we have a fully connected transition matrix where all transitions have a non-zero probability. **Additionally, an stationary distribution  $d^\pi$  of states exists and is independent start state  $s_0$ .**

Detail see (Ortner, 2007).

**Experience Replay** (Vanseijen and Sutton, 2015) (Riedmiller, 2005)

**Function Approximation** Function approximation is a way of combining supervised learning algorithms with learning procedures in Reinforcement Learning to be able to generalize value function into large scale MDP. There are three different objective function for function approximation.

1. Direct method

In direct method, we are trying to minimize the square error of our estimation and the “true” value function which we temporarily assume given by an oracle.

$$J(\theta) = \frac{1}{2}(v(s) - V_\theta(s))^2$$

The term  $\frac{1}{2}$  is just for neat math like regular supervised learning, and  $v(s)$  is the “true” value for state  $s$  which we want to approximate. To use the gradient descent method, we take the gradient of the  $J(\theta)$ , which gives us

$$\begin{aligned}\nabla_\theta J(\theta) &= \frac{1}{2} * 2 * (v(s) - V_\theta(s)) \nabla_\theta V_\theta(s) \\ &= (v(s) - V_\theta(s)) \nabla_\theta V_\theta(s)\end{aligned}$$

Note, in the equation above,  $v(s)$  is just a constant number. Since we don't really have an oracle in learning, we need to find a way to estimate the “true” value function. Various methods can be used here. For example, we can use the temporal difference learning rule, thus the gradient becomes  $(r + \gamma V_\theta(s') - V_\theta(s)) \nabla_\theta V_\theta(s)$ , or we can use the monte carlo estimation, then the gradient becomes  $(G_t - V_\theta(s)) \nabla_\theta V_\theta(s)$  or eligibility trace.

However, this direct method has some issues with convergence.

detail here

2. Fixed point method

Another kind of objective function is based on the “fixed point” notion in contraction. With the contraction property, we know the value function will converge to a fixed point

$$V = TV$$

where  $T$  is the Bellman operator. Thus, after convergence, our value function should fulfill the Bellman equation which gives the second objective function

$$\begin{aligned}J(\theta) &= \frac{1}{2}(V_\theta(s) - TV_\theta(s))^2 \\ &= \frac{1}{2}(V_\theta(s) - (r + \gamma V_\theta(s')))^2\end{aligned}$$

when we take the gradient of this cost function  $J(\theta)$ , we have

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \frac{1}{2} * 2 * (r + \gamma V_{\theta}(s') - V_{\theta}(s)) \nabla_{\theta} (r + \gamma V_{\theta}(s') - V_{\theta}(s)) \\ &= (r + \gamma V_{\theta}(s') - V_{\theta}(s)) (\nabla_{\theta} \gamma V_{\theta}(s') - \nabla_{\theta} V_{\theta}(s)) \\ &= (r + \gamma V_{\theta}(s') - V_{\theta}(s)) (\gamma \nabla_{\theta} V_{\theta}(s') - \nabla_{\theta} V_{\theta}(s))\end{aligned}$$

need more stuff here

### 3. Projected fixed point method

**GLIE** GLIE stands for “Greedy in the Limit of Infinite Exploration” (Singh et al., 2000). The learning policies in RL can be divided into two broad categories: a *decay exploration* strategy which become more and more greedy and *persistent exploration* which always maintain a fix exploration rate. The advantage of the first one is that we can eventually converge to the optimal policy. The second one may have the advantage always be adaptive but may not converge to the optimal. (In here, we talk about convergence in the sense that the behavior will become optimal. It is possible that some of the algorithm converge to the correct Q-value but still behave randomly with some probability by using persistent exploration strategy, Q-learning with fix  $\epsilon$ -greedy for example). We may want to consider this in the context of [on-policy&off-policy](#).

If a *decay exploration* strategy has the following two characters:

1. each action is executed infinitely often in every state that is visited infinitely often, and
2. in the limit, the learning policy is greedy with respect to the Q-value function with probability 1.

Then we can consider this decay exploration strategy GLIE. Some example of GLIE include [Boltzmann Selection](#),  $\epsilon$ -greedy.

**Inverse Reinforcement Learning** Reason for using IRL instead of behavior cloning:

- Reward functions can generalize the expert’s behavior to the state space regions not covered by the demonstration. For example, the transition model can change and thus, the optimal policy change as well.
- behavioral cloning cannot be applied whenever the expert demonstrates actions that cannot be executed by the learner (think at a humanoid robot that learns by observing a human expert)

**Lipschitz Continuity** Given two metric space two norm (a function that assigns a strictly positive length or size to each vector in a vector space)

Given two metric spaces  $(X, d_X)$  and  $(Y, d_Y)$ , where  $d_X$  denotes the metric on the set  $X$  and  $d_Y$  is the metric on set  $Y$  (for example,  $Y$  might be the set of real numbers  $\mathbb{R}$  with the metric  $d_Y(x, y) = |x - y|$ , and  $X$  might be a subset of  $\mathbb{R}$ ), a function  $f: X \rightarrow Y$  is called Lipschitz continuous if there exists a real constant  $K \geq 0$  such that, for all  $x_1$  and  $x_2$  in  $X$ ,

(Eriksson et al., 2013)

**LSTD** Still working on it ...

Related paper

**Markov Decision Process** blah blah here

Markov Decision Process can be seen as a extension of Markov Chain with additional action set (allowing selection) and reward function (motivation). It can be reduced to Markov chain if we have only one action per state and same reward for all the state.

**Metric Space** Still working on it ...

**Definition (Metric Space, metric).** A *metric space* is a pair  $(X, d)$ , where  $X$  is a set and  $d$  is a metric on  $X$  (distance function on  $X$ ), that is, a function defined<sup>1</sup> on  $X \times X$  such that for all  $x, y, z \in X$  we have:

1.  $d$  is real-valued, finite and nonnegative.
2.  $d(x, y) = 0$  if and only if  $x = y$
3.  $d(x, y) = d(y, x)$  (Symmetry)
4.  $d(x, y) \leq d(x, z) + d(z, y)$  (Triangle inequality)

In this definition,  $X$  is called the underlying set of  $(X, d)$ . Its elements are called points. Here are some examples of metric space:

- Real line  $\mathbb{R}$

The set of all real numbers and  $d(x, y) = |x - y|$

- Euclidean plane  $\mathbb{R}^2$

Suppose we have  $x = (\epsilon_1, \epsilon_2)$  and  $y = (v_1, v_2)$ . The euclidean metric defined by

$$d(x, y) = \sqrt{(\epsilon_1 - v_1)^2 + (\epsilon_2 - v_2)^2}$$

- Euclidean space  $\mathbb{R}^n$

If  $x = (\epsilon_1, \dots, \epsilon_n)$  and  $y = (v_1, \dots, v_n)$ , then euclidean metric defined by

$$d(x, y) = \sqrt{(\epsilon_1 - v_1)^2 + \dots + (\epsilon_n - v_n)^2}$$

- Unitary space  $\mathbb{C}^n$

a  $n$ -dimensional unitary space  $\mathbb{C}^n$  is the space of all ordered  $n$ -tuples of complex numbers with distance function

$$d(x, y) = \sqrt{|\epsilon_1 - v_1|^2 + \dots + |\epsilon_n - v_n|^2}$$

- Sequence space  $l^\infty$

This example and the next one give a first impression of how surprisingly the concept of a metric space is. As a set  $X$  we take the set of all bounded sequences of complex numbers; that is, every element of  $X$  is a complex sequence

$$x = (\epsilon_1, \epsilon_2, \dots) \quad \text{briefly} \quad x = (\epsilon_j)$$

such that for all  $j = 1, 2, \dots$  we have

$$|\epsilon_j| \leq c_x$$

where  $c_x$  is a real number which may depend on  $x$ , but does not depend on  $j$ . We choose the metric defined by

$$d(x, y) = \sup_{j \in \mathbb{N}} |\epsilon_j - v_j|$$

where  $y = (v_j) \in X$  and  $\mathbb{N} = 1, 2, \dots$ , and  $\sup$  denotes the supremum (least upper bound).

- Function space  $C[a, b]$

As a set  $X$  we take the set of all real-valued functions  $x, y, \dots$  which are functions of an independent real variable  $t$  and are defined and continuous on a given closed interval  $J = [a, b]$ . Choosing the metric defined by

$$d(x, y) = \max_{t \in J} |x(t) - y(t)|$$

where  $\max$  denotes the maximum, we obtain a metric space which is denoted by  $C[a, b]$ . (The letter  $C$  suggests "continuous.") This is a function space because every point of  $C[a, b]$  is a function.

(Kreyszig, 1989).

---

<sup>1</sup>The symbol  $\times$  denotes the Cartesian product of sets:  $A \times B$  is the set of all ordered pairs  $(a, b)$ , where  $a \in A$ , and  $b \in B$ . Hence,  $X \times X$  is the set of all ordered pairs of elements of  $X$

**Monte Carlo Method** Monte Carlo method is a way of making the prediction in model-free environment. The question it wants to solve is that suppose we have a policy  $\pi$  known, how good is this policy? In this case, we evaluate the policy by giving the method episodes of experience  $\{s_1, a_1, r_2, \dots, s_T\}$  generated by following policy  $\pi$  and wants the value function  $V^\pi$  as output.

As we know, the value of being in a state  $s$  is the expectation of the discounted rewards received afterwards.

$$V^\pi(s) = E_\pi[r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{T-1} r_T]$$

two methods can be used here : first-visit and every-visit method

**On-policy and Off-policy** An RL algorithm can be essentially divided into two parts, the *learning policy* and *update rule*. The first one is a non-stationary policy that maps experience (state visited, action chosen, reward received) to into a currently choice of action. The second part is how the algorithm uses experience to change its estimation of the optimal value function. In off-policy algorithm, the *update rules* doesn't have relationship with *learning policy*, that is the *update rules* doesn't care the what action agent take. Q-learning can be consider as the off-policy algorithm.

$$Q_{t+1}(s, a) = (1 - \alpha)Q_t(s, a) + \alpha(r_t + \gamma \max_{a'} Q(s', a'))$$

We can see that the Q-value is update based on the  $\max_{a'} Q(s', a')$ , which doesn't depend on the action the agent was taking.

However, if we take a look of SARSA(0), which is very similar to Q-learning.

$$Q_{t+1}(s, a) = (1 - \alpha)Q_t(s, a) + \alpha(r_t + \gamma Q(s', a'))$$

We can see the update is based on the Q-value of the next action of the agent. Thus it is an on-policy algorithm. The convergence condition are heavily depend on the *learning policy*, The Q-value of SARSA(0) can only converge to optimality in the limit only if the learning policy behavior optimal in the limit. The SARSA(0) and Q-learning will be same if we use greedy action selection strategy.

Detail see (Singh et al., 2000).

**Policy Gradient Reinforcement Learning** Instead of determining the action based on value function, another method of determine the action to take is the direct policy search method. It is well known that value-function combined with function approximation are unable to converge to any policy even for simple MDPs (in mean-squared-error, residual-gradient, temporal-difference, and dynamic-programming sense)

Let  $\theta$  denote the parameters of the function approximator, neural network for example, and  $J(\theta)$  the performance of the corresponding policy (e.g. the average reward per step). Then in the policy gradient approach, we can update the policy parameters proportional to the gradient:

$$\Delta\theta = \alpha \frac{\partial J(\theta)}{\partial \theta}$$

Let's begin with the definition of the two formulations:

- Average reward formulation

$$J(\theta) = \lim_{n \rightarrow \infty} \frac{1}{n} E\{r_1 + r_2 + \dots + r_n | \pi\} = \sum_s d^\pi(s) \sum_a \pi(s, a) R(s, a)$$

where  $d^\pi(s)$  is the stationary distribution of the states induced by  $\pi$

where  $\alpha$  is a positive-definite step size.

REINFORCE algorithm also finds an unbiased estimate of the gradient, but without the assistance of a learned value function. REINFORCE learns much more slowly than RL methods using value functions and



has received relatively little attention. Learning a value function and using it to reduce the variance of the gradient estimate appears to be essential for rapid learning.

Likelihood Ratio Methods

$$\begin{aligned}
J(\theta) &= E[r(\tau)] \\
&= \int_{\tau} p_{\theta}(\tau) r(\tau) \\
\nabla_{\theta} J(\theta) &= \nabla_{\theta} \int_{\tau} p_{\theta}(\tau) r(\tau) \\
&= \int_{\tau} \nabla_{\theta} p_{\theta}(\tau) r(\tau) \\
&= \int_{\tau} \nabla_{\theta} p_{\theta}(\tau) \frac{p_{\theta}(\tau)}{p_{\theta}(\tau)} r(\tau) \\
&= \int_{\tau} p_{\theta}(\tau) \frac{\nabla_{\theta} p_{\theta}(\tau)}{p_{\theta}(\tau)} r(\tau) \\
&= \int_{\tau} p_{\theta}(\tau) \nabla_{\theta} \ln p_{\theta}(\tau) r(\tau) \\
&= E[\nabla_{\theta} \ln p_{\theta}(\tau) r(\tau)]
\end{aligned}$$

However, we know that

$$p_{\theta}(\tau) = p(x_0) \prod_{k=0}^H p(x_{k+1}|x_k, u_k) \pi_{\theta}(u_k|x_k)$$

Thus

$$\begin{aligned}
\nabla_{\theta} \ln p_{\theta}(\tau) &= \nabla_{\theta} [\ln p(x_0) + \sum_{k=0}^H (\ln p(x_{k+1}|x_k, u_k) + \ln \pi_{\theta}(u_k|x_k))] \\
&= \nabla_{\theta} [\ln p(x_0) + \sum_{k=0}^H \ln p(x_{k+1}|x_k, u_k) + \sum_{k=0}^H \ln \pi_{\theta}(u_k|x_k)] \\
&= \nabla_{\theta} \ln p(x_0) + \sum_{k=0}^H \nabla_{\theta} \ln p(x_{k+1}|x_k, u_k) + \sum_{k=0}^H \nabla_{\theta} \ln \pi_{\theta}(u_k|x_k) \\
&= 0 + 0 + \sum_{k=0}^H \nabla_{\theta} \ln \pi_{\theta}(u_k|x_k) \\
&= \sum_{k=0}^H \nabla_{\theta} \ln \pi_{\theta}(u_k|x_k)
\end{aligned}$$

Note, if instead of stochastic policy, we are using a deterministic policy, then we have

$$p_{\theta}(\tau) = p(x_0) \prod_{k=0}^H p(x_{k+1}|x_k, \pi_{\theta}(x_k))$$

In this case,

$$\begin{aligned}
\nabla_{\theta} \ln p_{\theta}(\tau) &= \nabla_{\theta} \left[ \ln p(x_0) + \sum_{k=0}^H \ln p(x_{k+1}|x_k, \pi_{\theta}(x_k)) \right] \\
&= \nabla_{\theta} \ln p(x_0) + \sum_{k=0}^H \nabla_{\theta} \ln p(x_{k+1}|x_k, \pi(x_k)) \\
&= \nabla_{\theta} \ln p(x_0) + \sum_{k=0}^H \nabla_{u_k} \ln p(x_{k+1}|x_k, u_k) \nabla_{\theta} \pi(x_k) \\
&= 0 + \sum_{k=0}^H \nabla_{u_k} \ln p(x_{k+1}|x_k, u_k) \nabla_{\theta} \pi(x_k) \\
&= \sum_{k=0}^H \nabla_{u_k} \ln p(x_{k+1}|x_k, u_k) \nabla_{\theta} \pi(x_k)
\end{aligned}$$

Since we need to compute  $\nabla_{u_k} \ln p(x_{k+1}|x_k, u_k)$ , thus we need to know the model of the system

**Policy Iteration** Still working on it ...

**Stationary Distribution** Still working on it ...

**Stochastic Game** Stochastic game can be seen as an extension of [MDP](#).

**Temporal Difference Method**

**Value Iteration**

**Vector Space**

**References**

- A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 1998.
- K. Conrad. The contraction mapping theorem. *Expository paper. University of Connecticut, College of Liberal Arts and Sciences, Department of Mathematics*, 2014.
- K. Eriksson, D. Estep, and C. Johnson. *Applied Mathematics: Body and Soul: Volume 1: Derivatives and Geometry in IR3*. Springer Science & Business Media, 2013.
- E. Kreyszig. *Introductory functional analysis with applications*, volume 81. 1989.
- R. Ortner. Linear dependence of stationary distributions in ergodic markov decision processes. *Operations research letters*, 35(5):619–626, 2007.
- M. Riedmiller. Neural fitted q iteration – first experiences with a data efficient neural reinforcement learning method. In *Machine Learning: ECML 2005*, volume 3720 of *Lecture Notes in Computer Science*, pages 317–328. 2005.
- S. Singh, T. Jaakkola, M. L. Littman, and C. Szepesvári. Convergence results for single-step on-policy reinforcement-learning algorithms. *Machine Learning*, 38(3):287–308, 2000.
- H. Vanseijen and R. Sutton. A deeper look at planning as learning from replay. In *Proceedings of the 32nd International Conference on Machine Learning (ICML-15)*, pages 2314–2322, 2015.