# RL Notes

Ben Barber

May 21, 2024

## 1  Markov Decision Process

A Markov Decision Process (MDP) is a Markov chain with the addition of actions and rewards. It is a discrete-time stochastic control process. MDPs provide a mathematical framework for modeling decision-making situations where outcomes are partly random and partly under the control of a decision maker. MDPs are characterized by:

- States $s \in \mathcal{S}$: The different conditions or situations the system or agent can be in.

- Actions $a \in \mathcal{A}(s)$: Choices or decisions that can be taken in each state.

- Rewards $r \in \mathcal{R} \in \mathbb{R}$: The feedback received after taking actions in certain states.

- Transition Probability $p(s', r|s, a)$: The probability of moving from one state to another, given a specific action.

### 1.1  Policy

The policy $\pi(a|s)$ determines the probability of taking action $a$ in state $s$.

### 1.2  Return

The return $G_t$ is defined as the discounted sum over future rewards.

$$G_t = \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k \tag{1}$$

Where $\gamma$ is the discount factor and $T$ is the terminal step. Expanded form:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{T-t-1} R_T \tag{2}$$

Note: $\mathcal{S}^+$ denotes the state space along with a terminal state.

### 1.3  Trajectory

$$\tau = (S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, \ldots, R_T, S_T) \tag{3}$$

A trajectory $\tau$ is defined as a sequence of states, actions, and rewards associated with one episode of a MDP.

### 1.4  Expectation

$$\mathbb{E}_p[f(\mathbf{x})] = \int p(\mathbf{x})f(\mathbf{x})d\mathbf{x} \quad \text{or} \quad \sum_{\mathbf{x}} p(\mathbf{x})f(\mathbf{x}) \tag{4}$$

Where $p(\mathbf{x})$ is a probability density of $\mathbf{x}$ and $f(\mathbf{x})$ is some scalar function. This should be thought of as the probability weighted average of $f(\mathbf{x})$ over the domain of $\mathbf{x}$.

## 1.5 State Value Function

$$v_\pi(s) = \mathbb{E}_\pi[G_t|S_t = s] \tag{5}$$

The state value function $v_\pi(s)$ gives the expected return when the agent is at state $s$ and follows policy $\pi$.

## 1.6 Action Value Function

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t|S_t = s, A_t = a] \tag{6}$$

The state value function $v_\pi(s)$ gives the expected return when the agent is at state $s$ and takes action $a$, then follows policy $\pi$.

## 1.7 Goal

The goal is to find the optimal policy $\pi_*$ that maximizes the expected return.

$$\pi_* = \max_\pi \mathbb{E}_\pi[G_t] \tag{7}$$

# 2 The Bellman Equation and Generalized Policy Iteration

For any policy $\pi$, all $s \in \mathcal{S}$, and all $a \in \mathcal{A}$, these equations hold:

$$v_\pi(s) = \sum_{a \in \mathcal{A}(s)} \pi(a|s)q_\pi(s, a) \tag{8}$$

This means the state value function is equal to the policy weighted average of the action values.

$$q_\pi(s, a) = \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r|s, a)[r + \gamma v_\pi(s')] \tag{9}$$

This means the action value function is equal to the probability weighted average of the reward obtained in the next step plus the once-discounted value of the next state. Note that this definition assumes *complete knowledge* of the environment, or full access to the probability distribution of the entire MDP ($p(s', r|s, a)$).

The above equations are not the actual Bellman equations, but they can be composed to make them. Substituting (9) into (8) yields the Bellman equation for the state value function:

$$v_\pi(s) = \sum_{a \in \mathcal{A}(s)} \pi(a|s) \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r|s, a)[r + \gamma v_\pi(s')] \tag{10}$$

This relates any state value $s$ to all state values one step away.

Substituting (8) into (9) yields the Bellman equation for the action value function:

$$q_\pi(s, a) = \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r|s, a)[r + \gamma \sum_{a' \in \mathcal{A}(s')} \pi(a'|s')q_\pi(s', a')] \tag{11}$$

This relates any action value $q$ to all action values one step away.

## 2.1 Bellman Optimality

The Bellman optimality equations can be derived from the general Bellman equations with a couple small modifications:

For any *optimal* policy $\pi_*$, all $s \in \mathcal{S}$, and all $a \in \mathcal{A}$, these equations hold:

$$v_*(s) = \max_{a \in \mathcal{A}(s)} q_*(s, a) \tag{12}$$

This simply states that the optimal state value is the maximum over all optimal action values.

$$q_*(s, a) = \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r | s, a)[r + \gamma v_*(s')] \tag{13}$$

Performing the same substitutions as above results in the Bellman optimality equations.

## 2.2 Policy Evaluation

Policy evaluation is a very important subroutine in RL, and refers to computing $v_\pi(s)$ or $q_\pi(s, a)$ for a given fixed policy $\pi$. One algorithm for policy evaluation involves starting from random values and iteratively applying the Bellman equation to update each state value from the values of the next states. One of these iterations is called a *sweep*:

$$V(s) \leftarrow \sum_{a \in \mathcal{A}(s)} \pi(a | s) \sum_{s' \in \mathcal{S}, r \in \mathcal{R}} p(s', r | s, a)[r + \gamma v_\pi(s')] \tag{14}$$

This process always converges to the true state values. This process can be applied in an almost identical manner to calculate action values.

## 2.3 Policy Improvement

Another important subroutine is policy improvement, where given $v_\pi(s)$, the goal is to find a better policy than $\pi$.

Recall that for any optimal policy $\pi_*$:

$$\pi_*(s) = \operatorname*{argmax}_a q_*(s, a) \tag{15}$$

Given this, we can define a new policy $\pi'$ that simply chooses the action with the highest value under the current policy $\pi$:

$$\pi'(s) = \operatorname*{argmax}_a q_\pi(s, a) \tag{16}$$

And due to the Policy Improvement Theorem, we can be sure $\pi'$ is always at least as good as $\pi$:

$$v_\pi(s) \leq v_{\pi'}(s) \qquad \forall s \in \mathcal{S} \tag{17}$$

## 2.4 Generalized Policy Iteration (GPI)

GPI is a class of algorithms that takes advantage of the fact that an arbitrary number of iterations of policy evaluation and improvement always converges on the optimal policy $\pi_*$. Almost all RL methods are well described as GPI.

The simplest GPI algorithm is policy iteration. This algorithm starts with some arbitrary initial policy $\pi_0$ and simply repeats policy evaluation and policy improvement until the optimal policy and optimal value function are reached:

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \cdots \xrightarrow{I} \pi_* \xrightarrow{E} v_* \tag{18}$$

3

# 3    Monte Carlo and Off-Policy Methods

At the root, Monte Carlo methods are algorithms for approximating expectations by sampling batches of data from the probability distribution and averaging them:

$$\mathbb{E}_p[f(\mathbf{x})] \approx \frac{1}{N} \sum_{i=1}^{N} f(\mathbf{x}_i) \qquad \mathbf{x}_i \sim p(\mathbf{x}) \tag{19}$$

Unlike in the assumption of dynamic programming, in MC, we don't have complete knowledge of the MDP. The goal is to use GPI to obtain $\pi \approx \pi_*$.

## 3.1    Model-Free vs Model-Based Methods

In RL, a model is defined as something the agent uses to predict the environment's response to its actions. In model-based methods, a model is used to plan actions before they are taken, such as a chess engine. In model-free methods, the agent learns associations between actions and rewards. Monte Carlo methods are model-free.

## 3.2    Monte Carlo Evaluation

The goal is to estimate $q_\pi$ given sample trajectories under $\pi$. We can express $q_\pi$-estimation as $v_\pi$-estimation by translating to a Markov Reward Process (MRP) where the state-action pairs of the original MDP are the states of the new MRP.

$$S_t^{new} = (S_t, A_t) \tag{20}$$

This effectively bakes in the policy of the original MDP. Estimating $v(s^{new})$ is isomorphic to estimating $q_\pi(s, a)$. This is a useful generalization, as we can now apply any state value estimation method to action value estimation. We use averages to approximate $v_\pi(s)$:

$$v_\pi(s) \approx \frac{1}{C(s)} \sum_{m=1}^{M} \sum_{t=0}^{T_m - 1} \mathbb{I}[s_t^m = s] g_t^m \tag{21}$$

We sum over all $M$ trajectories and again over all timesteps in each trajectory. $\mathbb{I}$ is an indicator function. This works as a filter that only includes rewards gathered from state $s$. $g_t^m$ is the sampled return at time $t$ in the $m$th trajectory. The double sum can then be summarized as the sum of all rewards that followed from being in state $s$. We then divide by the number of occurences of $s$, $C(s)$, to get the average reward.

### 3.2.1    Update Rule

In practice, instead of calculating the value function using a batch of trajectories as in (9), we apply an update rule after each trajectory:

$$V(s_t^m) \leftarrow V(s_t^m) + \frac{1}{C(s_t^m)}(g_t^m - V(s_t^m)) \tag{22}$$

Where $C(s_t^m)$ denotes the number of occurences of state $s$ up to time $t$ in trajectory $m$. This means we make smaller updates as we process more data.

### 3.2.2    Constant-$\alpha$ MC

We can change $C(s_t^m)$ to a constant and it will still converge to the value function.

$$V(s_t^m) \leftarrow V(s_t^m) + \alpha(g_t^m - V(s_t^m)) \tag{23}$$

Where $\alpha$ is called the step size or learning rate. For smaller values of $\alpha$, the update rule converges more slowly, but more accurately on the true values. For larger values, it converges more quickly, but has more noisy and erratic behavior over time.

## 3.3  Monte Carlo Control

In MC control, we are now modifying the policy, so we can no longer model our environment as a Markov reward process. To execute MC control, we apply the same update rule, but instead of updating state values, we update action values:

$$Q(s_t^m, a_t^m) \leftarrow Q(s_t^m, a_t^m) + \alpha(g_t^m - Q(s_t^m, a_t^m)) \tag{24}$$

We then apply policy improvement, argmaxing over the actions in each state.

### 3.3.1  Exploration-Exploitation Trade-off

The outstanding issue with the above method is that the initial policy may be biased towards actions that it knows about. It could be leaving truly optimal actions on the table simply because it has never tried them. The exploration-exploitation trade-off states that in order to discover optimal policies, we must *explore* all state-action pairs. However, in order to get high returns, we must *exploit* known high-value pairs.

With infinite data, $\pi_*$ is always discoverable if the policy is *soft*. A policy is soft if it gives all actions in all states a positive probability.

$$\pi(a|s) > 0 \qquad \forall s \in \mathcal{S} \quad \forall a \in \mathcal{A}(s) \tag{25}$$

In other words, in the limit, all state-action pairs are visited an infinite number of times, meaning GPI will find the optimal policy.

### 3.3.2  $\epsilon$-Greedy Policy

With probability $\epsilon$, take an action selected uniformly from $\mathcal{A}(s)$, otherwise take $\text{argmax}_a Q(s, a)$.

### 3.3.3  Off-Policy Methods

Off-policy methods are a generalization of on-policy methods, where the same policy serves to generate data by choosing actions and to be updated and returned as an estimate of the optimal policy. In off-policy methods, these roles are separated into two policies, the behavior policy $b(a|s)$ which generates the data and the target policy $\pi(a|s)$ which is evaluated and improved.

Our goal is the same. We want to find $q_\pi(s, a)$, but since our data comes from $b$ and not $\pi$, we are actually estimating the expected return under $b$:

$$\mathbb{E}_b[G_t | S_t = s, A_t = a] \tag{26}$$

In order to estimate $q_\pi(s, a)$, we use *importance sampling*, which enables estimating an expectation of one distribution using samples from another:

$$q_\pi(s, a) = \mathbb{E}_b[\rho G_t | S_t = s, A_t = a] \tag{27}$$

$$\rho = \frac{p_\pi(G_t)}{p_b(G_t)} = \prod_{\tau=t+1}^{T-1} \frac{\pi(A_\tau | S_\tau)}{b(A_\tau | S_\tau)} \tag{28}$$

We also must assert *coverage*, meaning if the target policy might take a particular action in a particular state, then so must the behavior policy:

$$\pi(a|s) > 0 \Rightarrow b(a|s) > 0 \tag{29}$$

This ensures that in the limit of infinite data, we won't have zero data in places where the target policy might end up.

# 4 Temporal Difference Learning

The MC approach in the last section has a major disabling feature. Each episode must complete before values can be updated. If episodes are long, learning can be very slow. The individual transitions and rewards within an episode contain useful information. If we could update the policy from those transitions during an episode, which would likely improve how the agent explores in that episode.

## 4.1 $n$-step Temporal Difference Learning

Recall the update rule from 3.2.2:

$$V(s_t^m) \leftarrow V(s_t^m) + \alpha(g_t^m - V(s_t^m)) \tag{30}$$

Here, $g_t^m$ is called the *target*. It is what we move our estimate towards on each step. In $n$-step TD, instead of using this target, which is only available at the end of an episode, we add the discounted rewards of the next $n$ steps, and then use our current estimate of the value of the state $n$ steps away $V(s_{t+n})$ to account for the unobserved steps.

$$g_{t:t+n}^m = r_{t+1}^m + \gamma r_{t+2}^m + \cdots + \gamma^{n-1} r_{t+n}^m + \gamma^n V(s_{t+n}^m) \tag{31}$$

This provides an update rule that is first observable at time $t+n$ and can be applied at every time step until $T$, even when $T = \infty$. Using a value estimate in the target like this is called *bootstrapping*.

It is important to note that the criteria of MC and TD are different. MC tries to minimize the mean-squared error between the estimate and the data. TD tries to maximize the likelihood of the MRP. It more closely models the actual data source.

## 4.2 On-Policy TD Control: $n$-step Sarsa

This algorithm is almost identical to MC control (3.3). Just like in MC control, the same update rule as above is applied, but the state value estimates $V(s_t^m)$ are replaced with action value estimates $Q(s_t^m, a_t^m)$.

$$g_{t:t+n}^m = r_{t+1}^m + \gamma r_{t+2}^m + \cdots + \gamma^{n-1} r_{t+n}^m + \gamma^n Q(s_{t+n}^m, a_{t+n}^m) \tag{32}$$

$$Q(s_t^m, a_t^m) \leftarrow Q(s_t^m, a_t^m) + \alpha(g_{t:t+n}^m - Q(s_t^m, a_t^m)) \tag{33}$$

The main differences between this and MC control are (1) $g_t^m$ is replaced with $g_{t:t+n}^m$ and (2) updates are now applied *during* the episodes with an $n$-step delay. The name Sarsa comes from 1-step Sarsa, where every episodes operates on some $s_t^m, a_t^m, r_{t+1}^m, s_{t+1}^m, a_{t+1}^m$.

## 4.3 Q-Learning

This algorithm is very similar to $n$-step TD control, except we replace the last term of the target (32) with the maximum action value from that step. For simplification, this is in context of 1-step TD, but can easily be generalized to $n$-step TD.

$$r_{t+1}^m + \gamma \max_a Q(s_{t+1}^m, a) \tag{34}$$

Note that the use of max instead of using the action value directly makes this an off-policy method, since that last term is now a different policy functioning as the target. Another difference is the timing of the update. Instead of updating at the end of each Sarsa tuple, a Q-learning algorithm updates before the next action.

## 4.4 Expected Sarsa

In this algorithm, instead of maximizing the next action value, the target uses the policy weighted sum of all the next action values, or the expectation of the next action value.

$$r_{t+1}^m + \gamma \sum_a \pi(a|s_{t+1}^m) Q(s_{t+1}^m, a) \equiv r_{t+1}^m + \gamma \mathbb{E}_\pi[Q(s_{t+1}^m, a)] \tag{35}$$

# 5   Function Approximation

So far, all the algorithms and environments mentioned have had small state spaces. Of course, in real world problems, state spaces are not small, but arbitrarily large. It is therefore infeasible to move forward with this tabular approach. In order to solve RL problems with large state spaces, it is necessary for the agent to generalize its experiences in a miniscule subset of the state space to select high reward actions in the rest of the state space. The solution is function approximation.

The goal is to approximate $v_\pi(s)$ where data is generated under a fixed policy $\pi$ and the state space is arbitrarily large. It is assumed that the true value function can be approximated by some function of that state and a fixed-length parameter vector $\mathbf{w}$.

$$v_\pi(s) \approx \hat{v}(s, \mathbf{w}) \qquad \mathbf{w} \in \mathbb{R}^d \tag{36}$$

The additional constraint on this approach is that since $d \ll |S|$, any change to $\mathbf{w}$ results in changes to the estimated value for many states. In other words, unlike the tabular case, the parameters do not provide enough degrees of freedom to adjust the value of one state at a time. It follows from this constraint that it is impossible to arrive at the true value function, even with infinite data.

## 5.1   Selecting w

The goal is to select a vector $\mathbf{w}$ that minimizes the *value error* (loss), which is the weighted mean-squared error between the approximate value function and the true value function.

$$\overline{VE}(\mathbf{w}) = \sum_{s \in \mathcal{S}} \mu(s)[v_\pi(s) - \hat{v}(s, \mathbf{w})]^2 \tag{37}$$

Here, $\mu(\cdot)$ is a distribution over states, weighing how much we care about each state. It is often defined as the proportion of time spent in each state.

Clearly, this is impossible to calculate. It requires summing over an arbitrarily large space and knowledge of the true value function, which is an unknown. Instead of trying to calculate this metric, we approximate it.

## 5.2   Stochastic Gradient Descent

This is one common algorithm for approximately minimizing $\overline{VE}(\mathbf{w})$ by finding a local minimum. It requires three additional assumptions:

- States are visited in proportion to $\mu(\cdot)$
- $\hat{v}(s, \mathbf{w})$ is differentiable
- We observe a surrogate $U_t$ for $v_\pi(S_t)$

Given this, SGD applies an update rule to $\mathbf{w}$ for each visit to a state:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha[U_t - \hat{v}(S_t, \mathbf{w})]\nabla\hat{v}(S_t, \mathbf{w}) \tag{38}$$

Here, each update nudges $\mathbf{w}$ in the direction that maximally reduces the error from the target.