

A Succinct Summary of Reinforcement Learning

Sanjeevan Ahilan*

Abstract

This document is a concise summary of many key results in single-agent reinforcement learning (RL). The intended audience are those who already have some familiarity with RL and are looking to review, reference and/or remind themselves of important ideas in the field.

Contents

1	Acknowledgements	2
2	Fundamentals	2
2.1	The RL paradigm	2
2.2	Agent and environment	2
2.3	Observability	3
2.4	Markov processes and Markov reward processes	3
2.5	Markov decision processes	3
2.6	Policies, values and models	3
2.7	Dynamic programming	5
3	Model-free approaches	6
3.1	Prediction	6
3.2	Control with action-value functions	8
3.3	Value function approximation	9
3.4	Policy gradient methods	10
3.5	Baselines	11
3.6	Compatible function approximation	11
3.7	Deterministic policy gradients	12
4	Model-based Approaches	12
4.1	Model Learning	12
4.2	Combining model-free and model-based approaches	13
5	Latent variables and partial observability	14
5.1	Latent variable models	14
5.2	Partially observable Markov decision processes	14
6	Deep reinforcement learning	15
6.1	Experience replay	15
6.2	Target networks	15

*sanjeevanahilan@gmail.com. Much of this work was done at the Gatsby Unit, UCL.

1 Acknowledgements

I would like to thank Peter Dayan, David Silver, Chris Watkins and ChatGPT for helpful feedback. Much of this work was drawn from David Silver’s UCL course¹ and Sutton and Barto’s textbook (Sutton and Barto, 2018) and formed the introductory chapter of my PhD thesis (Ahilan, 2021).

2 Fundamentals

2.1 The RL paradigm

The field of reinforcement learning (RL) (Sutton and Barto, 2018) concerns itself with the computational principles underlying goal-directed learning through interaction. Although primarily seen as a field of machine learning, it has a rich history spanning multiple fields. In psychology it can be used to model classical (Pavlovian) and operant (instrumental) conditioning. In neuroscience it has been used to model the dopamine system of the brain (Schultz et al., 1997). In economics, it relates to fields such as bounded rationality, and in engineering it has extensive overlap with the field of optimal control (Bellman, 1957). In mathematics, investigation has continued under the guise of operations research. The plethora of perspectives ensures that RL continues to be an exciting and extraordinarily interdisciplinary field.

2.2 Agent and environment

RL problems typically draw a separation between the agent and the environment. The agent receives observation o_t and scalar reward r_t from the environment and emits action a_t , where t indicates the time step. The environment receives action a_t from the agent and then emits a reward r_{t+1} and an observation o_{t+1} . The cycle then begins again with the agent emitting its next action.

How the environment responds to the agent’s action is determined by the environment state s_t , which is updated at every time step. The conditional distribution for the next environment state depends only on the present state and action and therefore satisfies the Markov property:

$$P(s_{t+1}|s_t, a_t) = P(s_{t+1}|s_1, \dots, s_t, a_1, \dots, a_t) \quad (1)$$

The environment state is in general private from the agent, which only receives observations and rewards. The conditional distribution for the next observation given the current observation is not in general Markov, and so it may be beneficial for an agent to construct its own notion of state s_t^α , which it uses to determine its next action. This can be defined as $s_t^\alpha = f(h_t)$, where h_t is the history of the agent’s sequence of observations, actions and rewards:

$$h_t = a_1, o_1, r_1, \dots, a_t, o_t, r_t \quad (2)$$

¹<https://www.davidsilver.uk/teaching/>

2.3 Observability

A special case exists when the observation received by the agent o_t is identical to the environment state s_t (such that there is no need to distinguish between the two). This is the assumption underlying the formalism of Markov decision processes covered in the next section. An environment is partially observable if the agent cannot observe the full environment state, meaning that the conditional distribution for its next observation given its current observation does not satisfy the Markov property. This assumption underlies the formalism of a partially observable Markov decision process which we describe in Section 5.2.

2.4 Markov processes and Markov reward processes

A Markov process (or Markov chain) is a sequence of random states with the Markov property. It is defined in terms of the tuple $\langle \mathcal{S}, \mathcal{P} \rangle$ where \mathcal{S} is a finite set of states and $\mathcal{P} : \mathcal{S} \times \mathcal{S} \rightarrow [0, 1]$ is the state transition probability kernel.

A Markov Reward Process (MRP) $\langle \mathcal{S}, \mathcal{P}, r, \gamma \rangle$ extends the Markov process by including a reward function $r : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$ for each state transition and a discount factor γ . The immediate expected reward in a given state is defined as: $r(s) = \sum_{s'} \mathcal{P}(s, s')r(s, s')$.

The discount factor $\gamma \in [0, 1]$ is used to determine the present value of future rewards. Conventionally, a reward received k steps into the future is of worth γ^k times what it would be worth if received immediately. As we will shortly see, the cumulative sum of discounted rewards is a quantity RL agents often seek to maximise, and so $\gamma < 1$ ensures that this sum is bounded (assuming r is bounded).

2.5 Markov decision processes

Single-agent RL can be formalised in terms of Markov decision processes (MDPs). The idea of an MDP is to capture the key components available to the learning agent; the agent's sensation of the state of its environment, the actions it takes which can affect the state, and the rewards associated with states and actions. An MDP extends the formalism of an MRP to include a finite set of actions on which both \mathcal{P} and r depend. Discrete-time, infinite-horizon MDPs are described in terms of the 5-tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma \rangle$ where \mathcal{S} is the set of states, \mathcal{A} is the set of actions, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the state transition probability kernel, $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the immediate reward function and $\gamma \in [0, 1]$ is the discount factor. The expected immediate reward for a given state and action is defined as $r(s, a) = \sum_{s'} \mathcal{P}(s, a, s')r(s, a, s')$, which we use for convenience subsequently.

2.6 Policies, values and models

Common components of a reinforcement learning agent are a policy, value function and a model. The policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the agent's behaviour function which denotes the probability of taking action a in state s . Agents may also act according to a deterministic policy $\mu : \mathcal{S} \rightarrow \mathcal{A}$. We will assume that policies are stochastic unless otherwise noted.

Given an MDP and a policy π , the observed state sequence is a Markov process $\langle \mathcal{S}, \mathcal{P}^\pi \rangle$.

$$\mathcal{P}^\pi(s, s') = \sum_{a \in \mathcal{A}} \pi(s, a) \mathcal{P}(s, a, s') \quad (3)$$

Similarly, the state and reward sequence is a MRP $\langle \mathcal{S}, \mathcal{P}^\pi, r_\pi, \gamma \rangle$ in which:

$$r_\pi(s) = \sum_{a \in \mathcal{A}} \pi(s, a) r(s, a) \quad (4)$$

Starting from any particular state s at time step $t = 0$, the value function $v_\pi(s)$ is a prediction of the expected discounted future reward given that the agent starts in state s and follows policy π :

$$v_\pi(s) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} | s_0 = s \right] \quad (5)$$

where $r_{t+1} = r(s_t, a_t, s_{t+1})$

which is the solution of an associated Bellman expectation equation:

$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(s, a) \left[r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') v_\pi(s') \right] \quad (6)$$

In matrix form the Bellman expectation equation can be expressed in terms of the induced MRP:

$$\mathbf{v}_\pi = \mathbf{r}_\pi + \gamma \mathcal{P}^\pi \mathbf{v}_\pi = (\mathcal{I} - \gamma \mathcal{P}^\pi)^{-1} \mathbf{r}_\pi \quad (7)$$

where $\mathbf{v}_\pi \in \mathbb{R}^{|\mathcal{S}|}$ and $\mathbf{r}_\pi \in \mathbb{R}^{|\mathcal{S}|}$ are the vector of values and expected immediate rewards respectively for each state under policy π . We can also define a Bellman expectation backup operator:

$$T^\pi(\mathbf{v}) = \mathbf{r}_\pi + \gamma \mathcal{P}^\pi \mathbf{v} \quad (8)$$

which has a fixed point of \mathbf{v}^π .

An action-value for a policy π can also be defined, which is the expected discounted future reward for executing action a and subsequently following policy π .

$$\begin{aligned} q_\pi(s, a) &= r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') v_\pi(s') \\ &= r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') \sum_{a' \in \mathcal{A}} \pi(s', a') q_\pi(s', a') \end{aligned} \quad (9)$$

The process of estimating v_π or q_π is known as policy evaluation. Policies can be evaluated without directly knowing or estimating a model, using instead the directly sampled experience of the environment, an approach which is known as ‘model-free’. However a ‘model-based’ approach is also possible in which a model is used to predict what the environment will do next. A key component of a model is an estimate of $\mathcal{P}(s, a, s')$, the probability of the next state given the current state and action. Another is an estimate of $r(s, a)$, the expected immediate reward.

Policy evaluation enables a value function to be learned for a given policy. However, we often wish to learn the best possible policy. The value function for this is known as the optimal value function and corresponds to the maximum value function over all policies:

$$v_*(s) = \max_{\pi} v_{\pi}(s) \quad (10)$$

The definition of the optimal action-value function (which evaluates the immediate action a in state s) is similarly:

$$q_*(s, a) = \max_{\pi} q_{\pi}(s, a) \quad (11)$$

A partial ordering over policies can be defined according to:

$$\pi \geq \pi' \text{ if } v_{\pi}(s) \geq v_{\pi'}(s), \forall s \quad (12)$$

For any MDP there exists an optimal policy π_* that is better than or equal to all other policies. All optimal policies achieve the optimal value function and optimal action-value function and there is always a deterministic optimal policy for any MDP. The latter is achieved by selecting:

$$a = \arg \max_{a \in \mathcal{A}} q_*(s, a) \quad (13)$$

If there are many possible actions which satisfy this, any of these may be chosen to constitute an optimal policy (of which there may be many). The optimal value and state-value functions satisfy Bellman optimality equations:

$$\begin{aligned} v_*(s) &= \max_{a \in \mathcal{A}} q_*(s, a) \\ v_*(s) &= \max_{a \in \mathcal{A}} \left[r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') v_*(s') \right] \\ q_*(s, a) &= r(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s, a, s') \max_{a'} q_*(s', a') \end{aligned} \quad (14)$$

The Bellman optimality equation is non-linear with no closed form solution (in general). Solving it therefore requires iterative solution methods.

2.7 Dynamic programming

Dynamic programming (DP) (Bertsekas et al., 1995) refers to a collection of algorithms that can be used to compute optimal policies given a perfect model of the environment as an MDP. In general, DP solves complex problems by breaking them down into subproblems and then combining the solutions. It is particularly useful for overlapping subproblems, the solutions to which reoccur many times when solving the overall problem, making it more computationally efficient to cache and reuse them.

When applied to MDPs, the recursive decomposition of DP corresponds to the Bellman equation and the cached solution to the value function. DP assumes that the MDP is fully known and therefore does not address the full RL problem but instead addresses the problem of planning. By planning, the prediction problem can be addressed by finding the value function v_{π} of a given policy

π . This can be evaluated by iterative application of the Bellman Expectation Backup (Equation 8).

This leads to convergence to a unique fixed point v_π , which can be shown using the contraction mapping theorem (also known as the Banach fixed-point theorem) (Banach, 1922). When a Bellman expectation backup operator T^π is applied to two value functions \mathbf{u} and \mathbf{v} over states, we find that it is a γ -contraction:

$$\begin{aligned} \|T^\pi(\mathbf{u}) - T^\pi(\mathbf{v})\|_\infty &= \|(r_\pi + \gamma\mathcal{P}^\pi\mathbf{u}) - (r_\pi + \gamma\mathcal{P}^\pi\mathbf{v})\|_\infty \\ &= \|\gamma\mathcal{P}^\pi(\mathbf{u} - \mathbf{v})\|_\infty \\ &\leq \|\gamma\mathcal{P}^\pi\mathbf{1}\|\|\mathbf{u} - \mathbf{v}\|_\infty \\ &\leq \gamma\|\mathbf{u} - \mathbf{v}\|_\infty \end{aligned} \tag{15}$$

where $\mathbf{1}$ is a vector of ones and the infinity norm of a vector \mathbf{a} is denoted $\|\mathbf{a}\|_\infty$ and is defined as the maximum value of its components. This contraction ensures that both \mathbf{u} and \mathbf{v} converge to the unique fixed point of T^π which is \mathbf{v}_π .

For control, DP can be used to find the optimal value function v_* and in turn the optimal policy π_* . One possibility is *policy iteration* in which the current policy π is first evaluated as described and then subsequently improved to π' such that:

$$\pi'(s) = \arg \max_{a \in \mathcal{A}} q_\pi(s, a) \tag{16}$$

This improves the value from any state s over one step:

$$q_\pi(s, \pi'(s)) = \max_{a \in \mathcal{A}} q_\pi(s, a) \geq \sum_{a \in \mathcal{A}} \pi(s, a) q_\pi(s, a) = v_\pi(s) \tag{17}$$

It can be shown that this improves the value function such that that $v_{\pi'}(s) \geq v_\pi(s)$ (Silver, 2015). This process is then repeated, with improvements ending when the Bellman optimality equation (14) has been satisfied and convergence to π_* achieved. A generalisation of policy iteration is also possible in which, instead of waiting for policy evaluation to converge, only n steps of evaluation are taken before policy improvement occurs and the process is repeated. If $n = 1$ this is known as *value iteration*, as the policy is no longer explicit (being a direct consequence of the value function). Like policy iteration, value iteration is also guaranteed to converge to the optimal value function and policy. This can be demonstrated using the contraction mapping theorem.

3 Model-free approaches

3.1 Prediction

As has been outlined, dynamic programming can be used to solve known MDPs enabling optimal value functions and policies to be found. However, in many cases the MDP is not directly known - instead an agent taking actions in the MDP must learn directly from its experiences, as it transitions from state to state and receives rewards accordingly. One approach, known as ‘model-free’, seeks to solve MDPs without learning transitions or rewards. For prediction, a

key quantity to estimate in this setting is the expected discounted future reward. A sampled estimate of this, starting from state s_t , is known as the return:

$$R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \quad (18)$$

which depends on the actions sampled from the policy, and states from transitions.

Monte-Carlo (MC) methods seek to estimate this directly using complete episodes of experience. Introducing a learning rate α_t , the agent's value function can therefore be updated according to²:

$$v(s_t) \leftarrow v(s_t) + \alpha_t [R_t - v(s_t)] \quad (19)$$

The value function updated in this way will converge to a solution with minimum mean-square error (best fit to the observed returns), assuming a suitable sequential decrease in the learning rate.

Temporal-difference (TD) learning methods learn from incomplete episodes by bootstrapping. For example, if learning occurs after a single step, this is known as TD(0), which has the following update:

$$v(s_t) \leftarrow v(s_t) + \alpha_t [r_{t+1} + \gamma v(s_{t+1}) - v(s_t)] \quad (20)$$

where $r_{t+1} + \gamma v(s_{t+1})$ is known as the *target*. This approximates the full-width Bellman expectation backup (Equation 8) in which every successor state and action is considered, with experiences instead being sampled. TD(0) will converge to the solution of the maximum likelihood Markov model which best fits the data (again assuming a suitable sequential decrease in the learning rate). This solution may be different from the minimum mean-square error solution of MC methods, which do not assume the Markov property.

Unlike MC methods, TD methods introduce bias into the estimated return as the currently estimated value function may be different from the true value function. However, they generally have reduced variance relative to MC methods, as in MC the estimated return depends on a potentially long sequence of random actions, transitions and rewards.

The distinction between MC and TD methods can be blurred by considering multi-step TD methods (rather than only TD(0)), in which rewards are sampled for a number of steps before the value function is used to compute an estimate of future rewards. The n -step return is defined as:

$$R_t^{(n)} = r_{t+1} + \gamma r_{t+2} + \dots + \gamma^{n-1} r_{t+n} + \gamma^n v(s_{t+n}) \quad (21)$$

As $n \rightarrow \infty$ it tends towards the unbiased MC return. An algorithm may seek to find a good bias-variance tradeoff by estimating a weighted combination of n -step returns; one popular method to do this is known as TD(λ):

$$R_t^\lambda = (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} R_t^{(n)} \quad (22)$$

where $\lambda \in [0, 1]$.

²assuming a table-based representation rather than use of a function approximator

3.2 Control with action-value functions

Model free control concerns itself with optimising rather than evaluating the RL objective. Policies may be evaluated according to various objectives. In the case of continuing environments, the objective can be the average value or the average reward per time-step. We focus instead on episodic environments, assuming an initial distribution over starting states $p_0(s) : \mathcal{S} \rightarrow [0, 1]$. The objective is thus:

$$J(\pi) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_{t+1} | p_0(s) \right] \quad (23)$$

Note that if the domain of the starting state distribution is only over a single starting state, the objective is simply the value function (Equation 5) in that starting state. This objective can equivalently be expressed as:

$$J(\pi) = \mathbb{E}_{s \sim \rho^\pi, a \sim \pi} [r(s, a)] \quad (24)$$

where:

$$\rho^\pi(s) := \sum_{s'} \sum_{t=0}^{\infty} \gamma^t p(s_t = s | s', \pi) p_0(s') \quad (25)$$

is the improper discounted state distribution induced by policy π starting from an initial state distribution $p_0(s')$. In Section 3.4 we describe policy gradient methods which seek to optimise this objective directly.

However, we first consider model-free approaches which rely on an action-value function $q(s, a)$ to achieve control (a value function $v(s)$ alone is insufficient for model-free control). The optimal action-value function $q_*(s, a)$ must be learned, with MC and TD methods both viable. Once it has been learned, an optimal policy may be achieved by selecting the best action in each state (Equation 13).

However, unlike dynamic programming, full-width backups are not used and so if actions are selected greedily (meaning those with highest action-values are always chosen) then certain states and actions may never be correctly evaluated. Model-free RL methods must therefore allow for enough exploration during learning before ultimately exploiting this learning to achieve near-optimal cumulative reward.

One simple approach, known as ϵ -greedy is to take a random action with probability ϵ but otherwise act greedily according to the current estimate of the action-value function. The value of ϵ can be decreased with the number of episodes. This can satisfy a condition known as greedy in the limit of infinite exploration in which all state-action pairs are explored infinitely many times and the policy converges to the greedy policy.

One popular algorithm for model-free control is known as Q-learning (Watkins and Dayan, 1992), which seeks to learn the optimal action-value function whilst using a policy which also takes exploratory actions (such as epsilon greedy). This learning is termed *off-policy* as the policy used to sample experience is different from the policy being learned (the optimal policy). The resulting update is:

$$q(s_t, a_t) \leftarrow q(s_t, a_t) + \alpha [r_{t+1} + \gamma \max_{a' \in \mathcal{A}} q(s_{t+1}, a') - q(s_t, a_t)] \quad (26)$$

An alternative to off-policy Q-learning is on-policy SARSA (Rummery and Niranjan, 1994). This uses the sampled state s_t , action a_t , reward r_{t+1} , next state s_{t+1} , and next action a_{t+1} for updates³:

$$q(s_t, a_t) \leftarrow q(s_t, a_t) + \alpha(r_{t+1} + \gamma q(s_{t+1}, a_{t+1}) - q(s_t, a_t)) \quad (27)$$

3.3 Value function approximation

So far we have assumed a tabular representation of states and actions such that each state is separately updated. However, in practice we would like value functions and policies to generalise to new states and actions, and so it is beneficial to use function approximators such as deep neural networks. A common approach is to approximate the value function or action-value function:

$$\begin{aligned} v_w(s) &= \hat{v}(s; w) \approx v_\pi(s) \\ q_w(s, a) &= \hat{q}(s, a; w) \approx q_\pi(s, a) \end{aligned} \quad (28)$$

where w are the parameters we wish to learn. If we start by assuming we know the true value function v_π , we can define a mean square error between the approximate value function and the true function:

$$\mathcal{L}(w) = \mathbb{E}_\pi[(v_\pi(s) - v_w(s))^2] \quad (29)$$

Given a distribution of states $s \sim p(s)$ ⁴, we can minimise this iteratively using stochastic gradient descent:

$$w \leftarrow w + \alpha(v_\pi(s_t) - v_w(s_t))\nabla_w v_w(s_t) \quad (30)$$

In reality we can only use a better estimate of v_π provided by the sampled reward(s). For example, if we use the TD(0) target the update is:

$$w \leftarrow w + \alpha(r_{t+1} + \gamma v_w(s_{t+1}) - v_w(s_t))\nabla_w v_w(s_t) \quad (31)$$

Updates like this are known as ‘semi-gradient’ as the gradient of the value function used to define the target is ignored.

If we use a linear function approximator $v_w(s) = x(s)^T w$ (where features $x(s)$ and w are vectors), then we find:

$$w \leftarrow w + \alpha(r_{t+1} + \gamma v_w(s_{t+1}) - v_w(s_t))x(s_t) \quad (32)$$

indicating that the linear weights are updated in proportion to the activity of their corresponding features. Non-linear function approximators can also be used, but typically have weaker convergence guarantees than linear function approximators. Nevertheless, due to their flexibility such approximators have enabled impressive performance in a number of challenging domains, such as Atari games (Mnih et al., 2015) and Go (Silver et al., 2016).

³and also gives SARSA its name

⁴we later discuss a method for sampling states

3.4 Policy gradient methods

Parameterised stochastic policies π_θ may be improved using the *policy gradient theorem* (Sutton et al., 2000). This can be derived for any of the common RL objectives. To demonstrate a derivation of this result we use a starting state objective $J(\theta) = v_{\pi_\theta}(s_0)$ with a single starting state s_0 :

$$\begin{aligned}
\nabla_\theta J(\theta) &= \nabla_\theta v_\pi(s_0) \\
&= \nabla_\theta \sum_a \pi(s_0, a) q_\pi(s_0, a) \\
&= \sum_a \nabla_\theta \pi(s_0, a) q_\pi(s_0, a) + \pi(s_0, a) \nabla_\theta q_\pi(s_0, a) \\
&= \sum_a \nabla_\theta \pi(s_0, a) q_\pi(s_0, a) + \pi(s_0, a) \nabla_\theta \left[r(s_0, a) + \sum_{s'} \gamma \mathcal{P}(s_0, a, s') v_\pi(s') \right] \\
&= \sum_a \nabla_\theta \pi(s_0, a) q_\pi(s_0, a) + \pi(s_0, a) \sum_{s'} \gamma \mathcal{P}(s_0, a, s') \nabla_\theta v_\pi(s')
\end{aligned} \tag{33}$$

We note that we could continue to unroll $\nabla_\theta v_\pi(s')$ on the R.H.S in the same way as we have already done. Considering now transitions from starting state s_0 to arbitrary state s we therefore find:

$$\nabla_\theta v_\pi(s_0) = \sum_s \sum_{t=0}^{\infty} \gamma^t p(s_t = s | s_0, \pi) \sum_a \nabla_\theta \pi(s, a) q_\pi(s, a) \tag{34}$$

where $\sum_{t=0}^{\infty} \gamma^t p(s_t = s | s_0, \pi)$ is the discounted state distribution $\rho^\pi(s)$ from a fixed starting state s_0 (Equation 25). This derivation holds even when there is a distribution over starting states, and gives us the policy gradient theorem:

$$\nabla_\theta J(\theta) = \sum_s \rho^\pi(s) \sum_a \nabla_\theta \pi(s, a) q_\pi(s, a) \tag{35}$$

Using the likelihood ratio trick:

$$\begin{aligned}
\nabla_\theta \pi(s, a) &= \pi(s, a) \frac{\nabla_\theta \pi(s, a)}{\pi(s, a)} \\
&= \pi(s, a) \nabla_\theta \log \pi(s, a)
\end{aligned} \tag{36}$$

this can be equivalently expressed as:

$$\begin{aligned}
\nabla_\theta J(\theta) &= \sum_s \rho^\pi(s) \sum_a \pi(s, a) q_\pi(s, a) \nabla_\theta \log \pi(s, a) \\
&= \mathbb{E}_\pi [q_\pi(s, a) \nabla_\theta \log \pi(s, a)]
\end{aligned} \tag{37}$$

The policy gradient theorem result enables model-free learning as gradients need only be determined for the policy rather than for properties of the environment. There are a variety of approaches for determining q_π . If q_π is approximated using the sample return (Equation 18), this leads to the algorithm known as REINFORCE (Williams, 1992):

$$\theta \leftarrow \theta + \alpha R_t \nabla_\theta \log \pi(s_t, a_t) \tag{38}$$

As there is no bootstrapping here, this is also known as MC policy gradient. An alternative approach is to separately approximate q_π with a ‘critic’ q_w giving rise to what are commonly known as ‘actor-critic’ methods. These introduce two sets of parameter updates; the critic parameters w are updated to approximate q_π , and the policy (actor) parameters θ are updated according to the policy gradient as indicated by the critic. The critic itself can be updated according to the TD error. An example of this approach is SARSA actor-critic:

$$\begin{aligned} w &\leftarrow w + \alpha_1(r_{t+1} + \gamma q_w(s_{t+1}, a_{t+1}) - q_w(s_t, a_t)) \nabla_w q_w(s_t, a_t) \\ \theta &\leftarrow \theta + \alpha_2 q_w(s_t, a_t) \nabla_\theta \log \pi(s_t, a_t) \end{aligned} \quad (39)$$

where different learning rates α_1 and α_2 may be used for the actor and the critic.

3.5 Baselines

Whether we use REINFORCE or an actor-critic based approach to policy gradients, it is possible to reduce the variance further by the introduction of baselines. If this baseline depends only on the state s , then we find it introduces no bias:

$$\begin{aligned} \sum_s \rho^\pi(s) \sum_a \nabla_\theta \pi(s, a) b(s) &= \sum_s \rho^\pi(s) b(s) \nabla_\theta \sum_a \pi(s, a) \\ &= \sum_s \rho^\pi(s) b(s) \nabla_\theta 1 \\ &= 0 \end{aligned} \quad (40)$$

A natural choice for the state-dependent baseline is the value function:

$$\begin{aligned} \nabla_\theta J(\theta) &= \mathbb{E}_\pi[(q_\pi(s, a) - v_\pi(s)) \nabla_\theta \log \pi(s, a)] \\ &= \mathbb{E}_\pi[A_\pi(s, a) \nabla_\theta \log \pi(s, a)] \end{aligned} \quad (41)$$

where A_π is known as the advantage, which may in some algorithms be approximated directly (rather than approximating both q_π and v_π).

3.6 Compatible function approximation

In the general case, our choice to approximate q_π with q_w introduces bias such that there are no guarantees of convergence to a local optimum. However, in the special case of a compatible function approximator we can introduce no bias and take steps in the direction of the true policy gradient. This becomes possible when the critic’s function approximator reaches a minimum in the mean-squared error:

$$\begin{aligned} 0 &= \mathbb{E}_\pi[\nabla_w (q_\pi(s, a) - q_w(s, a))^2] \\ &= \mathbb{E}_\pi[(q_\pi(s, a) - q_w(s, a)) \nabla_w q_w(s, a)] \end{aligned} \quad (42)$$

If we choose $q_w(s, a)$ such that $\nabla_w q_w(s, a) = \nabla_\theta \log \pi(s, a)$ we find:

$$\mathbb{E}_\pi[q_\pi(s, a) \nabla_\theta \log \pi(s, a)] = \mathbb{E}_\pi[q_w(s, a) \nabla_\theta \log \pi(s, a)] \quad (43)$$

where the L.H.S is equal to the true policy gradient and so our function approximation has introduced no bias. For example, if the policy is a Boltzmann policy with a linear combination of features, of the form:

$$\pi(s, a) = \frac{e^{\theta^T \phi(s, a)}}{\sum_{a'} e^{\theta^T \phi(s, a')}} \quad (44)$$

then a compatible value function must be linear in the same features as the policy except normalised to zero mean for each state using a subtractive baseline (Sutton et al., 2000).

$$q_w(s, a) = \mathbf{w}^T [\phi(s, a) - \sum_{a'} \phi(s, a') \pi(s, a')] \quad (45)$$

3.7 Deterministic policy gradients

Rather than have a policy specify a probability for certain actions in certain states we can instead have it simply be a function mapping states to actions $\mu_\theta : \mathcal{S} \rightarrow \mathcal{A}$ and, in the case of continuous actions, seek to find the gradient of the objective with respect to the policy parameters. An example of an algorithm which uses such an approach is Deterministic Policy Gradients (DPG) (Silver et al., 2014). The DPG algorithm builds on the deterministic policy gradient theorem:

$$\nabla_\theta J(\theta) = \mathbb{E}_{s \sim \rho^\mu} [\nabla_\theta \mu_\theta(s) \nabla_a q_\mu(s, a)|_{a=\mu_\theta(s)}]. \quad (46)$$

where the parameters of the policy are adjusted in an off-policy fashion using an exploratory behavioural policy (which is a noisy version of the deterministic policy). In practice q_μ is approximated by the critic q_w , which is differentiable in the action and updated using Q-learning:

$$\begin{aligned} \delta_t &= r_{t+1} + \gamma q_w(s_{t+1}, \mu_\theta(s_{t+1})) - q_w(s_t, a_t) \\ w &\leftarrow w + \alpha_1 \delta_t \nabla_w q_w(s_t, a_t) \end{aligned}$$

The parameters of the policy are then updated according to:

$$\theta \leftarrow \theta + \alpha_2 \nabla_\theta \mu_\theta(s_t) \nabla_a q_w(s_t, a_t)|_{a=\mu_\theta(s_t)} \quad (47)$$

4 Model-based Approaches

In model-free RL agents learn to take actions directly from experiences, without ever modelling transitions in the environment or reward functions, whereas in model-based RL the agent attempts to learn these. The key benefit is that if the agent can perfectly predict the environment ‘in its head’, then it no longer needs to interact directly with the environment in order to learn an optimal policy.

4.1 Model Learning

Recall that MDPs are defined in terms of the 5-tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma \rangle$. Although models can be predictions about anything, a natural starting point is to approximate the state transition function $\mathcal{P}_\eta \approx \mathcal{P}$ and immediate reward function

$r_\eta \approx r$. We can then use dynamic programming to learn the optimal policy for an approximate MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}_\eta, r_\eta, \gamma \rangle$, the performance of which may be worse than for the true MDP.

Given a fixed set of experiences, a model can be learned using supervised methods. For predicting immediate expected scalar rewards, this is a regression problem whereas for predicting the distribution over next states this is a density estimation problem. Given the simplicity of this framing, a range of function approximators may be employed, including neural networks and Gaussian processes.

4.2 Combining model-free and model-based approaches

Once a model is learned it can be used for planning. However, in many situations it is computationally infeasible to do the full-width backups of dynamic programming as the state space is too large. Instead, experiences can be sampled from the model and used as data by a model-free algorithm.

A well known architecture which combines model-based and model-free RL is the Dyna architecture (Sutton, 1991). Dyna treats samples of simulated and real experience similarly, using both to learn a value function. Simulated experience is generated by the model which is itself learned from real experience. In Dyna, model-free based updates depend on the state the agent is currently in, whereas for the model-based component starting states can be sampled randomly and then rolled forwards using the model to update the value function using e.g. TD learning.

One potential disadvantage of Dyna is that it does not preferentially treat the state the agent is currently in. In many cases, such as deciding on the next move in chess, it is useful to start all rollouts from the current state (the board position) when choosing the next move. This is known as forward search, where a search tree is built with the current state as the root. Forward based search often uses sample based rollouts rather than full-width ones so as to be computationally tractable and this is known as simulation-based search.

An effective algorithm for simulation-based search is Monte-Carlo Tree search (Coulom, 2007). It uses the MC return to estimate the action-value function for all nodes in the search tree using the current policy. It then improves the policy, for example by being ϵ -greedy with respect to the new action-value function (or more commonly handling exploration-exploitation using Upper Confidence Trees, see Kocsis and Szepesvári (2006) for a more detailed discussion). MC Tree Search is equivalent to MC control applied to simulated experience and therefore is guaranteed to converge on the optimal search tree. Instead of using MC control for search it is also possible to use TD-based control, which will increase bias but reduce variance.

Model-based RL is a highly active area of research. Recent advances include MuZero (Schrittwieser et al., 2020), which extends model-based predictions to value functions and policies, and Dreamer which plans using latent variable models (Hafner et al., 2019).

5 Latent variables and partial observability

5.1 Latent variable models

Hidden or ‘latent’ variables correspond to variables which are not directly observed but nevertheless influence observed variables and thus may be inferred from observation. In reinforcement learning, it can be beneficial for agents to infer latent variables as these often provide a simpler and more parsimonious description of the world, enabling better predictions of future states and thus more effective control.

Latent variable models are common in the field of unsupervised learning. Given data $p(x)$ we may describe a probability distribution over x according to:

$$p(x; \theta_{x|z}, \theta_z) = \int dz p(x|z; \theta_{x|z})p(z; \theta_z) \quad (48)$$

where $\theta_{x|z}$ parameterises the conditional distribution $x|z$ and θ_z parameterises the distribution over z .

Key aims in unsupervised learning include capturing high-dimensional correlations with fewer parameters (as in probabilistic principal components analysis), generating samples from a data distribution, describing an underlying generative process z which describes causes of x , and flexibly modelling complex distributions even when the underlying components are simple (e.g. belonging to an exponential family).

5.2 Partially observable Markov decision processes

A partially observable Markov decision process (POMDP) (Kaelbling et al., 1998) is a generalisation of an MDP in which the agent cannot directly observe the true state of the system, the dynamics of which is determined by an MDP. Formally, a POMDP is a 7-tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, r, \mathcal{O}, \Omega, \gamma \rangle$ where \mathcal{S} is the set of states, \mathcal{A} is the set of actions, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is the state transition probability kernel, $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the reward function, \mathcal{O} is the set of observations, $\Omega : \mathcal{S} \times \mathcal{A} \times \mathcal{O} \rightarrow [0, 1]$ is the observation probability kernel and $\gamma \in [0, 1)$ is the discount factor. As with MDPs, agents in POMDPs seek to learn a policy $\pi(s_t^\alpha)$ which maximises some notion of cumulative reward, commonly $\mathbb{E}_\pi[\sum_{t=0}^{\infty} \gamma^t r_{t+1}]$. This policy depends on the agent’s representation of state $s_t^\alpha = f(h_t)$, which is a function of its history.

One approach to solving POMDPs is by maintaining a belief state over the latent environment state - transitions for which satisfy the Markov property. Maintaining a belief over states only requires knowledge of the previous belief state, the action taken and the current observation. Beliefs may then be updated according to:

$$b'(s') = \eta \Omega(o'|s', a) \sum_{s \in \mathcal{S}} \mathcal{P}(s'|s, a) b(s) \quad (49)$$

where $\eta = 1 / \sum_{s'} \Omega(o'|s', a) \sum_{s \in \mathcal{S}} \mathcal{P}(s'|s, a) b(s)$ is a normalising constant.

A Markovian belief state allows a POMDP to be formulated as an MDP where every belief is a state. However, in practice, maintaining belief states in POMDPs will be computationally intractable for any reasonably sized problem. In order to address this, approximate solutions may be used. Alternatively,

agents learning using function approximators which condition on the past can construct their own state representations, which may in turn enable relevant aspects of the state to be approximately Markov.

6 Deep reinforcement learning

The policies and value functions used in reinforcement learning can be learned using artificial neural network function approximators. When such networks have many layers they are conventionally denoted as ‘deep’, and are typically trained on large amounts of data using stochastic gradient descent (LeCun et al., 2015). The application of deep networks in model-free reinforcement learning garnered extensive attention when they were successfully used to learn a variety of Atari games from scratch (Mnih et al., 2013). For the particular problem of learning from pixels a convolutional neural network architecture was used (LeCun et al., 1998), which are highly effective at extracting useful features from images. They have been extensively used on supervised image classification tasks due to their ability to scale to large and complex datasets (LeCun et al., 2015).

A deep analysis of deep reinforcement learning (DRL) is beyond the scope of this summary. However we review two key techniques used to overcome the technical challenge of stabilising training.

6.1 Experience replay

As an agent interacts with its environment it receives experiences that can be used for learning. However, rather than using those experiences immediately, it is possible to store such experience in a ‘replay buffer’ and sample them at a later point in time for learning. The benefits of such an approach were introduced by Mnih et al. (2013) for their ‘deep Q-learning’ algorithm. At each timestep, this method stores experiences $e_t = (s_t, a_t, r_{t+1}, s_{t+1})$ in a replay buffer over many episodes. After sufficient experience has been collected, Q-learning updates are then applied to randomly sampled experiences from the buffer. This breaks the correlation between samples, reducing the variance of updates and the potential to overfit to recent experience. Further improvements to the method can be made by prioritised (as opposed to random) sampling of experiences according to their importance, determined using the temporal-difference error (Schaul et al., 2015).

6.2 Target networks

When using temporal difference learning with deep function approximators a common challenge is stability of learning. A source of instability arises when the same function approximator is used to evaluate both the value of the current state and the value of the target state for the temporal difference update. After such updates, the approximated value of both current and target state change (unlike tabular methods), which can lead to a runaway target. To address this, deep RL algorithms often make use of a separate target network that remains stable even whilst the standard network is updated. As it is not desirable for the target network to diverge too far from the standard network’s improved

predictions, at fixed intervals the parameters of the standard network can be copied to the target network. Alternatively, this transition is made more slowly using Polyak averaging:

$$\phi_{targ} \leftarrow \rho\phi_{targ} + (1 - \rho)\phi \tag{50}$$

where ϕ are the parameters of the standard network and ρ is a hyperparameter typically close to 1.

References

- Sanjeevan Ahilan. *Structures for Sophisticated Behaviour: Feudal Hierarchies and World Models*. PhD thesis, UCL (University College London), 2021.
- Stefan Banach. Sur les opérations dans les ensembles abstraits et leur application aux équations intégrales. *Fund. math*, 3(1):133–181, 1922.
- Richard Bellman. A markovian decision process. *Journal of mathematics and mechanics*, pages 679–684, 1957.
- Dimitri P Bertsekas, Dimitri P Bertsekas, Dimitri P Bertsekas, and Dimitri P Bertsekas. *Dynamic programming and optimal control*, volume 1. Athena scientific Belmont, MA, 1995.
- Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In *International conference on computers and games*, pages 72–83. Springer, 2007.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- Leslie Pack Kaelbling, Michael L Littman, and Anthony R Cassandra. Planning and acting in partially observable stochastic domains. *Artificial intelligence*, 101(1-2):99–134, 1998.
- Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *European conference on machine learning*, pages 282–293. Springer, 2006.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444, 2015.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- Gavin A Rummery and Mahesan Niranjan. *On-line Q-learning using connectionist systems*, volume 37. University of Cambridge, Department of Engineering Cambridge, UK, 1994.
- Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *arXiv preprint arXiv:1511.05952*, 2015.

- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839):604–609, 2020.
- Wolfram Schultz, Peter Dayan, and P Read Montague. A neural substrate of prediction and reward. *Science*, 275(5306):1593–1599, 1997.
- David Silver. Lecture 3: Planning by dynamic programming. *Google DeepMind*, 2015.
- David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic policy gradient algorithms. In *ICML*, 2014.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- Richard S Sutton. Dyna, an integrated architecture for learning, planning, and reacting. *ACM Sigart Bulletin*, 2(4):160–163, 1991.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. 2018.
- Richard S Sutton, David A McAllester, Satinder P Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, pages 1057–1063, 2000.
- Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
- Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.