

Markov Definition

Definition (Bayes Theorem)
 $p(AB) = p(A|B) \times p(B) = p(B|A) \times p(A)$

Markov Property

Definition
A state s_t is **Markov** if and only if $P[s_{t+1}|s_t] = P[s_{t+1}|s_1, \dots, s_t]$. The future is independent of the past given the present.

- The present state s_t captures all information in the history of the agent's events.
- Once the state is known, then any data of the history is no longer needed.

Definition (Stationarity)
If the $P[s_{t+1}|s_t]$ do not depend on t , but only on the origin and destination states, we say the Markov chain is **stationary** or **homogeneous**.

A Markov Reward Process (MRP) is a Markov chain which emits rewards.

Definition (Markov Reward Process)
A Markov Reward Process is a tuple $(\mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma)$. \mathcal{S} is a set of states. $\mathcal{P}_{s'}$ is a state transition probability matrix. $\mathcal{R}_{s'} = E[r_{t+1}|s_t = s] = s$ is an expected immediate reward that we collect upon departing state s , this reward collection occurs at time step $t + 1$. $\gamma \in [0, 1]$ is a discount factor.

Why Discounting is a good idea
Most Markov reward processes are discounted with a $\gamma < 1$. Why?

- Mathematically convenient to discount rewards
- Avoids infinite returns in cyclic or infinite processes
- Uncertainty about the future may not be fully represented
- If the reward is financial, immediate rewards may earn more interest than delayed rewards
- Human and animal decision making shows preference for immediate rewards
- It is sometimes useful to adopt undiscounted processes (i.e. $\gamma = 1$), e.g. if all sequences terminate and also when sequences are equally long (why?)

Definition (State value function)
The state value function $v(s)$ of an MRP is the **expected return** R starting from state s at time t .

$$v(s) = E[R_t|S_t = s] \quad (7)$$

Bellman Equation

Forms of the Bellman Equation for MRPs

- Expectation notation:
$$v(s) = E[R_t|S_t = s]$$
- Sum notation (expectation written out):
$$v(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} v(s') \quad (13)$$

We have n of these equations, one for each state.

- Vector notation:
$$v = \mathcal{R} + \gamma P v \quad (14)$$

The vector is n -dimensional.

Direct solution of Bellman Equation

Direct solution

The Bellman equation is a linear, self-consistent equation:

$$v = \mathcal{R} + \gamma P v$$

we can solve for it directly:

$$\begin{aligned} v &= \mathcal{R} + \gamma P v \\ v - \gamma P v &= \mathcal{R} \\ (1 - \gamma P)v &= \mathcal{R} \\ v &= (1 - \gamma P)^{-1} \mathcal{R} \end{aligned} \quad (16) \quad (17) \quad (18) \quad (19)$$

Matrix inversion is computational expensive at $O(n^3)$ for n states (e.g. Backgammon has 10^{30} states), so direct solution only feasible for small MRPs. Fortunately there are many iterative methods for solving large MRPs:

- Dynkin programming
- Value-iteration
- Temporal-Difference learning

These are at the core of Reinforcement learning, we will learn all 3 algorithms. By the way you have met the solution of **self-consistent equations** before, whenever you have to find a solution to an unknown in a unknowns you exploited that the equations and the unknowns have to be self-consistent (i.e. related to each other by the common structure of the problem).

Definition (Policy)
A policy $\pi(a|s) = P[A_t = a|S_t = s]$ is the conditional probability distribution to execute an action $a \in A$ given that one is in state $s \in S$ at time t .

The general form of the policy is called a **probabilistic** or **stochastic** policy, so π is a probability. If for a given state s only a single a is possible, then the policy is **deterministic**: $\pi(a|s) = 1$ and $\pi(a'|s) = 0, \forall a' \neq a$. A shorthand is to write $\pi(s) = a$, implying that the function π returns an action for a given state.

Now we "only" need to work out how to choose an action ...

Solve bellman equation directly

$$\begin{aligned} V'(s) &= E_\pi[R_t|S_t = s] \\ &= E_\pi[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1}|S_t = s] \\ &= E_\pi[r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2}|S_t = s] \\ &= \sum_{a \in A} \pi(a|s) \left(\sum_{s' \in S} P_{ss'} \left[R_{ss'} + \gamma E_\pi \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+2}|S_{t+1} = s' \right] \right] \right) \\ &= \sum_{a \in A} \pi(a|s) \sum_{s' \in S} P_{ss'} \left[R_{ss'} + \gamma V'(s') \right] \end{aligned}$$

This is a version of Bellmann's equation. A fundamental property of value functions is that they satisfy a recursive consistency equation. Crucially V' has unique solution.

Iterative policy evaluation

```

Input:  $\pi$  (policy),  $R$  (rewards),  $\gamma$  (discount)
Output:  $V$  (value function)
Repeat:
  For each  $s \in S$ :
     $v \leftarrow V(s)$ 
    For each  $a \in A$ :
       $v \leftarrow v + \pi(a|s) \sum_{s' \in S} P_{ss'} [R_{ss'} + \gamma V(s')]$ 
     $\Delta \leftarrow \max_a |V(s) - v|$  ( $\Delta$  is a small positive number)
  Until  $\Delta < \theta$  (a small positive number)

```

In iterative policy evaluation we sweep through all successor states, we call this kind of operation a **full backup**.

To produce such successive approximation V_{t+1} from V_t , iterative policy evaluation applies the same operation to each state s : it replaces the old value of s in place with a new value obtained from the old values of the successor states of s , and the expected immediate rewards, along all the one-step transitions possible under the policy being evaluated. We could also run a code version storing old and new arrays for V . This turns out to converge slower, why?

State-Action Value function as Cost-To-Go

How good is it to be in a given state and take a given action when you follow a policy π ?

Definition (State-Action Value function "Cost to Go")

$$Q^\pi(s, a) = E[R_t|S_t = s, A_t = a] = E[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1}|S_t = s, A_t = a] \quad (20)$$

The relation between (state) value function and the state-action value function is straightforward:

$$V^\pi(s) = \sum_a \pi(a|s) Q^\pi(s, a) \quad (21)$$

Optimal Policy and optimal value and Q function

Value functions define a partial ordering over policies. A policy is defined to be better than or equal to a policy if its expected return is greater than or equal to that of all others. In other words, $\pi \geq \pi'$ if and only if $V^{\pi}(s) \geq V^{\pi'}(s)$ for all $s \in S$.

Definition (Optimal Value function)

$$V^*(s) = \max_\pi V^\pi(s), \forall s \in S \quad (22)$$

Therefore, the policy π^* that maximizes the value function is the **optimal policy**. There is always at least one optimal policy. There may be more than one optimal policy.

Definition (Optimal State-Action Value function)

$$Q^*(s, a) = \max_\pi Q^\pi(s, a), \forall s \in S, a \in A \quad (23)$$

In analogy to above we have

$$Q^*(s, a) = E[r_{t+1} + \gamma V^*(s_{t+1})|S_t = s, A_t = a] \quad (24)$$

V^* is the optimal value function and so conveniently the self-consistency condition can be rewritten in a form without reference to any specific policy π^* : $V^* = V^*$. This yields the Bellman Optimality Equation for an optimal policy:

Definition (Bellman Optimality Equation for V)

$$\begin{aligned} V^*(s) &= \max_{s'} \sum_{a'} P(s', a') (r(s, a', s') + \gamma V^*(s')) \\ &= \max_{s'} \sum_{a'} P_{ss'}^a (R_{ss'} + \gamma V^*(s')) \end{aligned} \quad (25) \quad (26)$$

Intuitively, the **Bellman Optimality Equation** expresses the fact that the value of a state under an optimal policy must equal the expected return for the best action from that state.

Three Assumption for BOE

On solving the Bellman Optimality Equations

Explicitly solving the Bellman optimality equation provides one route to finding an optimal policy, and thus to solving the reinforcement learning problem. This solution approach can often be challenging at best, if not impossible, because it is like an exhaustive search, looking ahead at all possible trades, computing their probabilities of occurrence and their desirabilities in terms of expected rewards. Moreover, this solution relies on at least 3 assumptions that are rarely true in practice:

- we accurately know the dynamics of the environment
- we have computational resources to find the solution
- the Markov property.

Thus, in reinforcement learning often we have to (and want to) settle for approximate solutions.

Convergence Rule

Theorem (Bellman Optimality Equation convergence theorem)
For an MDP with a finite state and action space

- The Bellman (Optimality) equations have a unique solution.
- The values produced by value iteration converge to the solution of the Bellman equations.

Dynamic Programming

Two assumption of DP

- MDP to be finite
- Environment model for environment, means we know the transition and reward function

Dynamic Programming – Origins 2

To be able to apply Dynamic Programming requires problems to have

- Optimal substructure**, meaning that the solution to a given optimization problem can be obtained by the combination of optimal solutions to its sub-problems.
- Overlapping sub-problems**, meaning that the space of sub-problems must be small, i.e., any recursive algorithm solving the problem should solve the same sub-problems over and over, rather than generating new sub-problems.

Example
The maximum path sum problem or Dijkstra's algorithm for the shortest path problem are dynamic programming solutions. In contrast, if a problem can be solved by combining optimal solutions to non-overlapping sub-problems, the strategy is called "divide and conquer" instead (e.g. quick sort).

First-visit MC prediction for estimating $V \approx v$

Input: a policy π to be evaluated
Initialize:
 $V(s) \in \mathbb{R}$, arbitrarily, for all $s \in S$
 $RetURNS(s) \leftarrow$ an empty list, for all $s \in S$
Loop for each episode (for each episode):
 Generate an episode following π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$
 $G \leftarrow 0$
 Loop for each step of episode, $t = 1, 2, \dots, T$:
 $G \leftarrow G + R_{t-1}$
 Append G to $RetURNS(S_t)$
 $V(s) \leftarrow \text{average}(RetURNS(S_t))$
First Visit MC vs Every visit MC
Above we have the **First visit MC** algorithm, as we append the return of the episode from the first occurrence of a state s . Another version is **Every visit MC**, where we append the return of the episode (from that point) on every occurrence of state s in the episodes. Batch VS Online Monte-Carlo

Batch & Online Monte-Carlo

Update policy after each episode
Online MC:
Batch MC:
Vanilla MC:
No Need to store samples traces
Performance Learning – Monte Carlo Learning

Incremental Monte-Carlo Updates

We can now update value functions without having to store sample traces

- Update $V(s)$ incrementally after step s, a_t, r_{t+1}, s_{t+1}
- For each state s_t with return R_t (up to this point) and $N(s)$ the visitor counter to this state:
$$N(s_t) \leftarrow N(s_t) + 1$$

$$V(s_t) \leftarrow V(s_t) + \frac{1}{N(s_t)} (R_t - V(s_t))$$

Moreover, if the world is **non-stationary**, it can be useful to track a **running mean**, i.e. by gradually forgetting old episodes.

$$V(s) \leftarrow V(s) + \alpha(R_t - V(s))$$

The parameter α controls the rate of forgetting old episodes (learning rate).

Why should we consider non-stationary conditions?

Temporal Difference Learning

Temporal-Difference (TD) Learning

- TD methods learn directly from episodes of experience (but also works for non-episodic tasks)
- TD is model-free: no knowledge of MDP transitions or rewards needed
- TD learns from incomplete episodes, by **bootstrapping**
- TD updates a guess towards a guess

MC update the value using the actual return R , where dp use the estimated return to update the value

Temporal Difference Learning update rule

Recall that we use the following Monte-Carlo update

$$V(s_t) \leftarrow V(s_t) + (r_{t+1} - V(s_t)) \quad (8)$$

We update the value of $V(s_t)$ towards the **actual return** R_t . Note, how we use only measurements to form our estimates. Temporal Difference methods perform a similar update after every time-step, i.e.

$$V(s_t) \leftarrow V(s_t) + \alpha(r_{t+1} - V(s_t)) \quad (9)$$

We update the value of $V(s_t)$ towards the **estimated return** $r_{t+1} - \gamma V(s_{t+1})$. Note, how we are combining a measurement r_{t+1} with an estimate $V(s_{t+1})$ to produce a better estimate $V(s_t)$.

Unlike policy iteration, there is no explicit policy.

If value of all states are updated same time or individually

Synchronous vs Asynchronous backups

- DP methods can use **synchronous** backups i.e. all states are backed up in parallel (this requires two copies of the value function)
- Asynchronous** DP backups states target only states individually – in any order (one copy of value function)
 - For each selected state we apply the appropriate backup function
 - This can significantly reduce computation and
 - we are still guaranteed to converge if over time all states continue to be selected

TDerminology

- $r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$ is the **Temporal Difference Error**.
- $r_{t+1} + \gamma V(s_{t+1})$ is the **Temporal Difference Target**

TD function value estimation Algorithm

```

1: procedure TD-ESTIMATION( $\pi$ )
2:   Init
3:    $V(s) \leftarrow$  arbitrary value, for all  $s \in S$ .
4:    $\mathcal{R} \leftarrow$  empty list
5:   repeat (For each episode)
6:     Initialize  $s$ 
7:     repeat (For each step of episode)
8:       action chosen from  $\pi$  at  $s$ 
9:       Take action  $a$ ; observe  $r$ , and next state,  $s'$ 
10:       $\delta \leftarrow r + \gamma V(s') - V(s)$ 
11:       $V(s) \leftarrow V(s) + \alpha \delta$ 
12:       $s \leftarrow s'$ 
13:      until  $s$  is absorbing state
14:   until Done
15: end procedure

```

Advantages & Disadvantages of MC & TD

- TD** can learn **full-width backups**
 - For each synchronous or asynchronous backup
 - We must complete knowledge of the MDP transitions and reward function
 - aside from **optimizing** the policy, the agent is not doing a lot of machine learning
 - Every transition and action is considered
 - DP is good for medium-sized problems (millions of states)
 - For large problems DP suffers the **Curse of Dimensionality**: Number of states $N = |\mathcal{S}|$ grows exponentially with number of (continuous) state variables
 - Even one backup can be too expensive (e.g. consider control of an n -joint robot arm).
- TD** can learn **without** the final outcome
 - MC can only learn from complete sequences
 - TD works in continuing (non-terminating) environments MC only works for episodic (terminating) environments
- TD** can learn **incomplete sequences**
 - MC has high variance, zero bias
 - Good convergence properties
 - even with function approximation (to be discussed later)
 - not very sensitive to initial value
 - very simple to understand and use
 - MC does not exploit Markov property
 - ⇒ Usually more effective in non-Markov environments
 - TD has low variance, some bias
 - Usually more efficient than MC
 - TD → to precise TD(0) – converges to $V^0(s)$
 - convergence not guaranteed with function approximation
 - more sensitive to initial value
 - TD exploits Markov property
 - ⇒ Usually more efficient in Markov environments

What is Markov Property

Key Characteristics of the Markov Property:

- Dependence on Current State:
 - Only the current state s_t and action a_t influence the future, regardless of how the system arrived at s_t .
- Simplifies Modeling:
 - We assume the Markov property, the transition dynamics $P(s_{t+1}|s_t, a_t)$ can be described without considering the full history of the process.

Summary DP vs TD vs MC

Bootstrapping: update involves an **estimate**

- MC does not bootstrap
- DP bootstraps
- TD bootstraps

Sampling: update does not involve an **expected value**

- MC samples
- DP does not sample
- TD samples

MC Control

Why use model free control

We can use Model-Free Control in two important scenarios:

- MDP model is known, but is too big to use (Curse of Dimensionality), except by sampling
- MDP model is unknown, but experience can be sampled.

On-policy vs off-policy

How can we avoid the unlikely assumption of exploring starts? The only general way to ensure that all actions are selected infinitely often is for the agent to continue to select them:

- On-policy** methods, which attempt to evaluate or improve the policy that is used to make decisions
- Off-policy** methods, that evaluate or improve a policy different from that used to generate the data.

On-policy learning is "Learn on the job"

- Learn about policy π from experience sampled from π
- Off-policy learning is "Look over someone's shoulder"
- Learn about policy π from experience sampled from π'

Means we have a policy that gives a non-zero probability to all possible actions

Definition
Soft policies have in general $\pi(a|s) > 0 \forall s \in S, \forall a \in A$. i.e. we have a finite probability to explore all actions.

Epsilon Greedy
 ϵ -greedy policies

ϵ -greedy policies are a form of soft policy, where the greedy action a^* (as selected from being greedy on the value or action-value function and choosing the arg max action) has a high probability of being selected, while all other actions available in the state, have an equal share of an ϵ probability (that allows us to explore the non-greedy/ optimal action).

Definition
 ϵ -greedy policy with $\epsilon \in [0, 1]$.
$$\pi(s, a) = \begin{cases} 1 - \epsilon / |\mathcal{A}(s)|, & \text{if } a^* = \operatorname{argmax}_a Q(s, a) \\ \epsilon / |\mathcal{A}(s)|, & \text{if } a \neq a^* \end{cases} \quad (5)$$

On-policy first-visit MC control (for soft policies), estimates $\pi \approx \pi$

Algorithm parameters: small $\epsilon > 0$
Initialize:
 $Q(s) \in \mathbb{R}$, arbitrarily, for all $s \in S$, $a \in A(s)$
 $RetURNS(s) \leftarrow$ an empty list, for all $s \in S$, $a \in A(s)$
Loop for each episode (for each episode):

- Generate an episode following π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$
- Append G to $RetURNS(S_t)$
- $Q(S_t) \leftarrow Q(S_t) + \alpha [R_t + \gamma \max_a Q(S'_t, a) - Q(S_t, a)]$
- $S \leftarrow S_t$

until S is terminal

MC control

The Monte Carlo methods presented in this chapter learn value functions and optimal policies from experience in the form of sample episodes. They give them three kinds of advantages over DP methods:

- designed Monte Carlo control methods we have followed the overall goal of generalized policy iteration (GPI)
- Rather than use a π to compare the value of each state, they simply average many returns that start in the state.
- Because a state's value is the expected return, this average can become a good approximation to the value.

Maintaining sufficient exploration is an issue in Monte Carlo control methods. It is not enough just to select the actions currently estimated to be best, because then no returns will be obtained for alternative actions, and it may never be learned that they are actually better.

One approach is to ignore this problem by assuming that episodes begin with state-action pairs randomly selected to cover all possibilities. Such exploring starts can sometimes be arranged in applications with simulated episodes, but are unlikely in learning from real experience.

In on-policy methods, the agent commits to always exploring and tries to find the best policy that still explores. In off-policy methods, the agent also explores, but learns a deterministic optimal policy that may be unrelated to the policy followed.

Off-policy Monte Carlo prediction refers to learning the value function of a target policy from data generated by a different behaviour policy. Such learning methods are all based on some form of importance sampling, that is, on weighting returns by the ratio of the probabilities of taking the observed actions under the two policies.

Stochastic Gradient Descent

Goal: find parameter vector w minimizing mean-squared error between approximate value function $\hat{V}(s, w)$ and true value function $V^*(s, w)$.

$$J(w) = E[(\hat{V}(s, w) - V^*(s, w))^2] \quad (1)$$

Gradient descent finds a local minimum by sliding down the gradient

$$\Delta w = -\frac{1}{2} \nabla J(w) \quad (2)$$

$$\Delta w = -\alpha (\hat{V}(s, w) - V^*(s, w)) \nabla \hat{V}(s, w) \quad (3)$$

The learning factor α is chosen so it's small enough so the derivative of the squared error. The term α is negative as we want to perform gradient descent (the gradient points upwards otherwise).

Stochastic gradient descent samples the gradient and the average update is equal to the full gradient update.

$$\Delta w = \alpha (\hat{V}(s, w) - V^*(s, w)) \nabla \hat{V}(s, w) \quad (4)$$

Linear Value Function Approximation II

- Update rule is particularly simple

$\nabla_w \hat{V}(s, w) = \nabla(s, w)^\top \hat{V}(s, w) = x(s)$

$$\Delta w = \alpha (\hat{V}'(s, w) \nabla(s, w)) \nabla \hat{V}(s, w) \quad (5)$$

Definition
Update = learning step size \times prediction error \times feature value

Monte-Carlo with Value Function Approximation

- Return R_t is an unbiased, noisy sample of true value $V^*(s_t)$
 - We apply supervised learning to "training data" of state return trace:
- $$(s_1, r_1), (s_2, r_2), \dots, (s_T, r_T) \quad (7)$$
- For example, using linear Monte-Carlo policy evaluation
- $$\Delta w = a(r_t - \hat{V}(s, w))\nabla_w \hat{V}(s, w) \quad (8)$$
- $$= a(R_t - \hat{V}(s, w))x(s_t) \quad (9)$$
- Monte-Carlo evaluation converges to a local optimum, even when using non-linear value function approximation (probable)

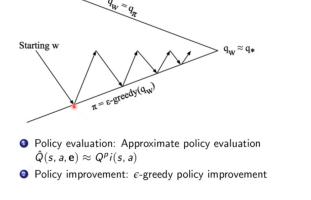
TD Learning with Value Function Approximation I

- The TD-target $R_{t+1} + \gamma V(s_{t+1}, w)$ is a biased (single) sample of the true value $V^*(s_t)$
 - We still perform supervised learning on "digested" training data:
- $$(s_1, r_1 + \gamma \hat{V}(s_2, w)), (s_2, r_2 + \gamma \hat{V}(s_3, w)), \dots, (s_T, r_T) \quad (10)$$
- For example, using linear TD
- $$\Delta w = a(r + \gamma \hat{V}(s', w) - \hat{V}(s, w))\nabla_w \hat{V}(s, w) \quad (11)$$
- $$= a(R_t + \gamma \hat{V}(s', w) - \hat{V}(s, w))x(s_t) \quad (12)$$

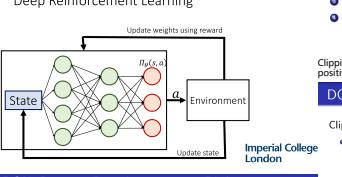
Linear TD converges "close" to the global optimum (probable). This does not extend to non-linear TD (see Sutton & Barto, 2018).

The estimates make the estimates closer to the real Q/Q value, but not reach the real Q/Q value, in the end, it will get close enough to the true Q/Q value.

From evaluation to control: FA in GPI



Deep Reinforcement Learning



DQN: Bringing Deep Learning into RL I

In principle it sounds easy: use convolutional neural networks to link screen shots to values.

Given our TD-Q-learning update:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + a[r_t + \gamma \max_{a_{t+1}} Q(s_{t+1}, a') - Q(s_t, a_t)] \quad (1)$$

We want to learn $Q(s_t, a_t, w)$ as a parametrised function (a neural network) with parameters w .

We can define the TD error as our learning target that we want to reduce to zero.

$$TD \text{ error}(w) = r_{t+1} + \gamma \max_{a_{t+1}} Q(s_{t+1}, a') - Q(s_t, a; w) \quad (2)$$

Taking the gradient of E wrt w we obtain:

$$\Delta w = a[r + \gamma \max Q(s_{t+1}, a_{t+1}; w) - Q(s_t, a; w)]\nabla_w Q(s_t, a; w) \quad (3)$$

However, the Atari playing problem required more engineering to solve properly

- Experience Replay
- Target Network
- Clipping of Rewards
- Skipping of Frames

This is inefficient and slow.

- Complication 1: Experience distribution
- Interactive learning produces training samples that are highly correlated (recent actions reflect recent policy)
- The probability distribution of the input data changes

The network forgets transitions that were in the not so recent past (overwriting past memory, in neuroscience known as "inhibition")

Moreover, because networks are monolithic, changes in one part of the network can have side-effects on other parts of the network [Why is this not a problem for table-based RL?]

Solution: Experience Replay² (alternative use Hierarchical RL).

Experiences (traces) are stored and replayed in mini-batches to train the neural networks on more than just the last episode.

Instead of running Q-learning on each state-action pairs as they occur during simulation or actual experience, the experience replay buffer system stores the traces sampled.

Mini-batches are only feasible when running multiple passes with the same data is somewhat stable with respect to the samples. I.e. the transitions should show low variance/entropy for the next immediate outcomes (reward, next state) given the same state, action pair.

Replaying past data is also a more efficient use of previous experience, by learning (training the neural network) with it multiple times. This is key when gaining real-world experience is costly (e.g. simulation time) as the Q-learning updates are incremental and do not converge quickly.

Target Network

Use the main network instead of the target network to do action selection (compared with normal target network)

Double Q Learning (van Hasselt et al., 2017, AAAI) I

Situation: In regular Deep Q-learning we use the target network for two tasks:

- Identify action with highest Q value,
- Determine Q value of this action

Complication: The maximum Q value may be overestimated (variance-bias problem), because an unusually high value from the main network Q does not mean that this is an unusually high value from the target network Q' .

The reason for this overestimation is that the same samples (i.e. the Q network) are being used to decide which action is the best (highest expected reward), and the same samples are also being used to estimate that specific action-value.

Solution: In Double Q-Learning (van Hasselt et al., 2017, AAAI) we separate the two estimates: We use a target network Q' for 1. but

use the regular Q for 2. (or the converse). Thus, we make the selection of the action with highest Q value, and selection of the Q value used in Bellman update based on samples from each other (different networks) so we become less susceptible to bias in each estimate. This reduces the frequency by which the maximum Q value may be overestimated, as it is less likely that both the networks are overestimating the same action.

REINFORCE algorithms

Before we had to learn a value function through function approximation and then derive a corresponding policy

Often learning a function value can be intractable (more unstable / was, at the time, intractable for large state spaces).

REINFORCE provided a way to directly optimize policies to get around this problem.

REINFORCE algorithm:

1. $\pi(a_t | s_t) \sim \pi_\theta(a_t | s_t)$ (run the policy)
2. $\nabla_\theta J(\theta) \approx \sum_{i=1}^T (\sum_{a_i} \nabla_\theta \log \pi_\theta(a_i | s_i)) (\sum_i r(s_i, a_i))$
3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

The method suffer from high variance in the sampled trajectories, thus stabilising model parameters is difficult.

REINFORCE was the crucial first step at popularising Policy Gradients (PGs).

What was the impact of REINFORCE?

Before: we had to learn a value function through function approximation and then derive a corresponding policy

But, often learning a value function can be intractable (more unstable, computationally demanding for large state spaces, observability).

REINFORCE: provided a way to directly optimize policies to get around this problem. It was the direct result of rendering the PG theorem useful after its derivation.

REINFORCE was the crucial first step at popularising Policy Gradients (PGs).

REINFORCE algorithms suffer greatly from variance, a single erratic trajectory can cause a suboptimal shift into wired area of optimisation.

Asynchronous Advantage Actor-Critic (A3C) II

The agents (or workers) are trained in parallel and update periodically a global network, which holds shared parameters.

Agents (workers) are trained in parallel and update periodically a global network, which holds shared parameters.

Updates are not happening simultaneously and that's where the asynchronous comes from.

DDPG (Lillicrap et al., 2015) I

Deep Deterministic Policy Gradients (DDPG) is a deterministic policy, model-free, off-policy, actor-critic algorithm which

- learns the value model (Q -function) and a policy π .

learns by off-policy methods and uses the Bellman equation to learn the Q -function (bootstrapping part 1)

uses the Q -function to learn the policy π (bootstrapping part 2)

We will next look first at the algorithm (as re-printed from Lillicrap et al., 2015) and then what ideas it contains. Note, Lillicrap 2015's notation differs from ours: θ_Q is the parameters of the value function approximation (we used w) and θ_π is the parameters of the policy function approximation (we used θ). R is the replay buffer (we used it for the total return).

DDPG (Lillicrap et al., 2015) I

To get DDPG to work well a few special elements were introduced:

Exploration in continuous action spaces

Use Gaussian distributed noise that perturbs a mean action $\mu'(s) = \mu(s) + N(0, \sigma)$ (i.e. Gaussian Policies).

This is almost equivalent to discrete action ϵ -greedy (only that we explore more often and much closer to the intended action μ)

To control learning variability DDPG uses replay buffers/minibatches & target networks (just as in DQN)

Error-based learning on the mean-squared Bellman error (MSBE)

MSBE (here denoted by L , tells us how closely our approximation of Q comes to satisfying the Bellman equation (compare this to function approximation Q-learning/DQN)).

Policy gradient based learning on the actor network

Here we are using a deterministic action version of the PG theorem, which proof (not examinable here) was published in Silver et al. (2014).

DDPG therefore uses "smooth" updates on the parameter vector θ (that is shared by both actor and critic) which updates slowly but steadily over time steps:

$$\theta^{t+1} = \theta^t + (1 - \beta)\theta^t, 0 < \beta < 1$$

The parameter β controls the degree of blending (or forgetting) between the current parameter estimate θ^t and thus θ^t to determine the parameter vector at the next $t+1$.

Just like the Policy Gradient theorem can be thought of a policy-centered equivalent of the Bellman Theorem, DDPG can consider the continuous action version of the discrete action DQN. Why?

Desired qualities for Reinforcement Learning

Stable learning & little sensitivity to hyperparameters: we want to avoid parameter tuning which is (computationally) expensive and more an art than a science. Ideally, RL algorithms work out of the box.

Sample Efficiency: Good sample complexity is the prerequisite for efficient (and in some cases effective) skill acquisition.

Most algorithms we covered so far are relatively poor at sample efficiency (although we have seen how introducing features improved their stability of learning). This is ok for simulation settings or "toy" problems such as video games.

Many real-world applications, e.g. robots require us to put the second aspect, good exploration at the centre of our attention, as interactive learning is costly and does not scale well.

An operational definition of the sample complexity of exploration for a reinforcement learning algorithm is the number of time steps in which the algorithm does not select near-optimal actions.

SAC - Objective Function I

In the Soft Actor-Critic we use the entropy concept to boost its sample complexity. The primary goal of SAC is to maximize a modified expected cumulative reward, which includes an entropy term to encourage exploration:

$$J(\pi) = E_{(s_t, a_t) \sim p} \left[\sum_t \gamma^t (r(s_t, a_t) + \lambda H(\pi(s_t))) \right] \quad (7)$$

Here,

$J(\pi)$ is the cost-to-go objective function,

(s_t, a_t) is the immediate reward function,

$p_{\pi} = p_{\pi}(s_t, a_t)$ is the state-action distribution under policy π ,

γ is the discount factor,

λ is a temperature parameter controlling the importance of the entropy term H ,

$\pi(s_t)$ is the policy distribution given state s_t .

SAC - Critic Update I

SAC employs two Q-value functions, Q_{μ} and Q_{σ} (cf. target networks in DDQN), to mitigate maximisation bias in policy improvement. In this entropy-regularized reinforcement learning, the agent gets an bonus reward at each time step proportional to the entropy of the policy at that timestep.

The critic (Q -functions) are updated to minimize the Bellman residual:

$$J_Q(\theta) = E_{(s_t, a_t) \sim p} \left[\left(Q_{\mu}(s_t, a_t) - (r(s_t, a_t) + \gamma E_{s_{t+1} \sim p} [\max_a Q_{\mu}(s_{t+1}, a)]) \right)^2 \right]$$

where D is the replay buffer and p_{θ_Q} represents the environment dynamics. Note a' comes from the current policy π , not the replay buffer

SAC - Critic Update II

$J_Q(\theta) = E_{(s_t, a_t) \sim p} \left[\left(Q_{\mu}(s_t, a_t) - (r(s_t, a_t) + \gamma E_{s_{t+1} \sim p} [\max_a Q_{\mu}(s_{t+1}, a') - \lambda \log \pi(a'|s_{t+1})]) \right)^2 \right]$

The Eq. can be interpreted as follows: Across all sampled states from our experience replay buffer, decrease the squared difference between the prediction of our Q -value network and the expected prediction of the Q function plus the entropy of the policy function π measured here by the negative log of the policy function (note, not to be confused with the $\log \pi$ as a Policy Gradient term).