

Markov Definition

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

Definition
A state s_t is **Markov** if and only if $P[s_{t+1}|s_t] = P[s_{t+1}|s_1, s_2, \dots, s_t]$. 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)
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 $(S, \mathcal{P}, \mathcal{R}, \gamma)$. S is a set of states. \mathcal{P}_{st} is a state transition probability matrix. $R_{st} = E[r_{t+1} | S_t = s]$ 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 reward
- It is sometimes useful to adopt undiscounted processes (i.e. $\gamma = 1$). 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** starting from state s at time t .

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

Solutions of the Bellman Equation for MRPs

- Expectation notation:

$$v(s) = E[R_t | S_t = s] \quad (16)$$

- Sum notation (expectation written out):

$$v(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} v(s') \quad (17)$$

$$v(s) = R_s + \gamma \sum_{s' \in S} P_{ss'} v(s') \quad (18)$$

$$v = (I - \gamma \mathcal{P})^{-1} R \quad (19)$$

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

- Vector notation:

$$\mathbf{v} = \mathcal{R} + \gamma \mathcal{P} \mathbf{v} \quad (14)$$

The vector \mathbf{v} is n -dimensional.

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

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

we can solve for it directly:

$$\mathbf{v} = \frac{\mathcal{R}}{1 - \gamma \mathcal{P}} \quad (20)$$

$$\mathbf{v} = \frac{\mathcal{R}}{1 - \gamma \mathcal{P}} \quad (21)$$

$$\mathbf{v} = (I - \gamma \mathcal{P})^{-1} \mathcal{R} \quad (19)$$

Matriкс inversion is computational cost of $O(n^3)$ for n states (e.g. Backgammon has 10¹² states), so direct solution only feasible for small MRPs. Fortunately there are many iterative methods for solving large MRPs:

- Dynamical programming
- Monte-Carlo evaluation
- Temporal-difference learning

These are at the core of Reinforcement learning, we will learn all 3 algorithms. By the way we met the solution of **self-consistent equations** before, whenever you solved for a set of n linear equations in unknowns you exploited that the equations and the unknowns had 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 $a \in \mathcal{A}$ given that one is in state $s \in \mathcal{S}$ 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 ...

Value function self-consistency
 $V^\pi(s) = E_{\pi}[R_t | S_t = s] = E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | S_t = s\right] = E_{\pi}\left[r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} | S_t = s\right] = \sum_{s' \in S} \pi_a(s) \sum_{s'' \in S} P_{ss'} \left[R_{s''} + \gamma E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k r_{s''+k+2} | S_{s''} = s'\right]\right] = \sum_{s' \in S} \pi_a(s) \sum_{s'' \in S} P_{ss'} \left[R_{s''} + \gamma V^\pi(s')\right]$

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

Iterative Policy Evaluation Algorithm
 $V^\pi(s) = \text{Initial Value } V(s) = \text{for all } s \in S$
 Repeat:
 $\quad \text{For each } s \in S:$
 $\quad \quad V(s) \leftarrow V(s) + \gamma \sum_{s' \in S} P_{ss'} [\mathbb{E}_{\pi}[r_{s'} + \gamma V(s')] - V(s)]$
 $\quad \quad \Delta \leftarrow \max(|V(s) - \mathbb{E}_{\pi}[r_{s'} + \gamma V(s')]|)$
 $\quad \quad \text{Until } \Delta < \theta \text{ (a small positive number)}$

In Iterative Policy evaluation we sweep through all successor states, we call this kind of operation a **full backup**. To produce each successive approximation V_{k+1} from V_k iterative policy evaluation starts with the initial state s_0 and the true value of the state s_0 plus a new value obtained from the old values of the successor states of s_0 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")
 $Q^\pi(s, a) = E[R_t | S_t = s, A_t = a] = E\left[\sum_{k=0}^{\infty} \gamma^k r_{t+k+1} | S_t = s, A_t = a\right] \quad (20)$

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

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a) Q^\pi(s, a) \quad (21)$$

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 states. In other words, $\pi \geq \pi'$ if and only if $V^\pi(s) \geq V^{\pi'}(s)$ for all $s \in \mathcal{S}$.

Definition (Optimal Value function)
 $V^*(s) = \max_{\pi} Q^\pi(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A} \quad (23)$

In analogy to before we also have

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

Bellman optimality equation for V^*
 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: $V^* = V^*$. This yields the Bellman Optimality Equation for an optimal policy:

Definition (Bellman Optimality Equation for V)

$$V^*(s) = \max_{a \in \mathcal{A}} P(s'|s, a) (r(s, a, s') + \gamma V^*(s')) \quad (25)$$

$$= \max_{a \in \mathcal{A}} P_{ss'} (R_{s'} + \gamma V^*(s')) \quad (26)$$

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 SOE
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 traces, 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 assumptions of DP
1. MDP to be finite
2. A perfect 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 optimisation 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 time step t solving the problem must start with 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 problems. 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).

For every states \Rightarrow
Policy Improvement Theorem
Let π and π' be two deterministic policies such that $\forall s \in \mathcal{S}: Q^\pi(s, a) \geq Q^{\pi'}(s, a)$. Then π' must be as good or better than π :
 $V^{\pi'}(s) \geq V^\pi(s), \forall s \in \mathcal{S}$.

In short, update value function using policy until value function converges, then be greedy to update the policy, then go back to step 2 until the policy converges.

Policy Iteration Algorithm
1. Initialization: $V^\pi(s) \in \mathbb{R}$ and $\pi(s, a)$ arbitrarily for all $s \in \mathcal{S}$
2. Policy Evaluation: $\mathbb{E}[r_{t+1} + \gamma \mathbb{E}_{\pi}[r_{t+2} + \gamma \mathbb{E}_{\pi}[r_{t+3} + \dots]] = V^\pi(s)$
3. Policy Improvement: $\pi^{pol_stable} \leftarrow \text{true}$
For each $s \in \mathcal{S}$:
 $\pi(s) \leftarrow \arg \max_a \sum_{s' \in \mathcal{S}} P_{ss'} [R_{s'} + \gamma V^\pi(s')]$
until $\Delta < \theta$ (Δ a small positive number)

Output a deterministic policy, π , such that $\pi(s) = \arg \max_a \sum_{s' \in \mathcal{S}} P_{ss'} [R_{s'} + \gamma V^\pi(s')]$

Unlike policy iteration, there is no explicit policy.

Value Iteration Algorithm
Initialize V^π arbitrarily, e.g., $V(s) = 0$, for all $s \in \mathcal{S}^*$
Repeat:
 $\Delta \leftarrow 0$
For each $s \in \mathcal{S}$:
 $V(s) \leftarrow \max_a \sum_{s' \in \mathcal{S}} P_{ss'} [R_{s'} + \gamma V^\pi(s')]$
 $\Delta \leftarrow \max(|V(s) - \mathbb{E}_{\pi}[r_{s'} + \gamma V^\pi(s')]|)$

TD value function estimation Algorithm
1: procedure TD-ESTIMATION(π)
2: Init
3: $V(s)$ – arbitrary value, for all $s \in \mathcal{S}$.
4: Evaluate
5: repeat (for each episode)
6: Initialize r
7: repeat (for each step of episode)
8: action chosen from π 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 **before** knowing the final outcome ("you can back-up near-death")
TD can learn online after every step
TD must wait until end of episode before return is known
TD can learn **without** the final outcome
TD can learn from incomplete sequences
MC can only learn complete sequences
TD works in continuous (non-terminal) environments MC only works for episodic (terminating) environments

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 be precise TD(0) – converges to $V^*(s)$
convergence not guaranteed with function approximation
more sensitive to initial value

TD exploits Markov property
Usually more efficient in Markov environments

TD uses full-width backups
For each synchronous or asynchronous backup
We need complete knowledge of the MDP transitions and reward function: aside from **optimizing** the policy, our agent is not doing a lot of machine learning
Every success state and action is considered

TD is efficient for solving problems (millions of states)
For example: solving the 100-state gridworld problem

TD is vulnerable to overfitting problems (millions of states)
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).

Let us take a step back from DP
• TD uses full-width backups
• For each synchronous or asynchronous backup
• We need complete knowledge of the MDP transitions and reward function: aside from **optimizing** the policy, our agent is not doing a lot of machine learning
• Every success state and action is considered

TD is efficient for solving problems (millions of states)
For example: solving the 100-state gridworld problem

TD is vulnerable to overfitting problems (millions of states)
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).

Key Characteristics of the Markov Property:
1. 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 .

2. Simplicity Modeling:
• Assuming 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.

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 states. In other words, $\pi \geq \pi'$ if and only if $V^\pi(s) \geq V^{\pi'}(s)$ for all $s \in \mathcal{S}$.

Definition (Optimal Value function)
 $V^*(s) = \max_{\pi} Q^\pi(s, a), \forall s \in \mathcal{S}$

Therefore, the policy π^* maximises 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)
 $Q^*(s, a) = \max_{\pi} Q^\pi(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A} \quad (23)$

In analogy to before we also have

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

Contracting Mapping Property:
Value iteration repeatedly applies the Bellman update:

$$V_{k+1}(s) = \max_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} P(s'|s, a) [R(s, a, s') + \gamma V_k(s')] \quad (25)$$

This update is a **contracting mapping** with respect to the max-norm $\|\cdot\|_\infty$:

$$\|V_{k+1} - V_k\|_\infty \leq \gamma \|V_k\|_\infty$$

The contraction mapping property arises because $\gamma < 1$, which ensures that each iteration brings V closer to V^* .

Summary:
Value iteration is guaranteed to converge because:

1. The Bellman update is a contraction mapping.
2. The Banach fixed point theorem ensures convergence to the unique optimal value function V^* .
3. The discount factor $\gamma < 1$ ensures bounded rewards and stability.
4. A finite state and action space limits the computational complexity, ensuring a solution in finite time.

Model-Free Reinforcement Learning: Monte Carlo (MC)
MC methods learn directly from episodes of experience

MC is **model-free**: no knowledge of MDP transitions or rewards needed

MC learns from complete episodes (of sample traces): **no bootstrapping**

MC uses the simplest possible idea: value of state = mean return

But this can only be applied to **episodic** MDPs that have terminal states.

Model-Free 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 (Cause of Dimensionality), except by sampling
- MDP model is unknown, but experience can be sampled.

→ policy vs 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'

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

Epison Greedy

c-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 ϵ probability (that allows us to explore the non-greedy / optimal action).

On-policy first-visit MC control algorithm
On-policy first-visit MC control (for c-greedy policies), estimates $\pi \approx \pi_c$

Algorithmic pseudocode: small $\epsilon > 0$

Initialize: $\pi(s, a)$ arbitrarily, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$
 $Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$
 $Return(s, a) \leftarrow \emptyset$, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat forever (for each episode):
1. Choose s_0 randomly (or $s_0 \in Return(s_0, a_0)$)
2. Loop for $t = 0$ to end of episode, $t = T - 1, 2, \dots, 0$:
 a. Take action a_t from s_t
 b. $R_t \leftarrow r_t$
 c. $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [R_t + \max_a Q(s_{t+1}, a)] - Q(s_t, a_t)$
 d. $s_{t+1} \leftarrow s_t$
 e. $Return(s_t, a_t) \leftarrow Return(s_t, a_t) \cup (s_{t+1}, a_{t+1})$

until s_T is terminal

Convergence of exploratory MC algorithms
Definition

Greedy in the Limit with Infinite Exploration (GLIE)
All state-action pairs are explored infinitely many times, $\lim_{t \rightarrow \infty} N_t(s, a) = \infty$

The policy converges on a greedy policy, $\lim_{t \rightarrow \infty} \pi_t(s, a) = \arg \max_{a \in \mathcal{A}} Q_t(s, a)$ where the infix comparison operator == evaluates to 1 if true and 0 else.

GLIE explores us when a schedule for adapting the exploration parameter is sufficient to ensure convergence.

Example
The c-greedy policy is GLIE if c reduces to zero with $c_k = \frac{1}{k}$.

Summary MC control I
In designing Monte Carlo control methods we have followed the overall schema of generalised policy iteration (GPI). Rather than use a model to compute 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.

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.

Function Approximation
Approximating functions

Generalise from seen states to unseen states Update parameter w using MC or TD learning

- Estimate value function with **function approximation**

$V^\pi(s) \approx \hat{V}(s, \mathbf{w})$
 $Q^\pi(s, a) \approx \hat{Q}(s, \mathbf{w}, a)$

- Generalise** from seen states to unseen states (fundamental ability of any good machine learning system)
- Update** function approximation parameter w using MC or TD learning

Stochastic Gradient Descent
Goal: find parameter vector w minimising the mean-squared error between approximate value function $\hat{V}(s, w)$ and true value function $V(s)$

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

Gradient descent finds a local minimum by sliding down the gradient

$$\Delta w = -\frac{1}{N} \sum_{i=1}^N (\hat{V}(s_i, w) - V(s_i))^2 \quad (2)$$

$$= -\frac{1}{N} \sum_{i=1}^N (\hat{V}(s_i, w) - V(s_i)) \nabla_w \hat{V}(s_i, w) \quad (3)$$

The learning factor – η , where the learning rate controls the step size. The term η is chosen to be that cancels out with the derivative of the squared error. The term is negative as we want to perform gradient descent (the gradient points upwards otherwise).

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

$$\Delta w = \frac{1}{N} \sum_{i=1}^N (\hat{V}(s_i, w) - V(s_i)) \nabla_w \hat{V}(s_i, w) \quad (4)$$

Linear Value Function Approximation II
Update rule is particularly simple

$$\nabla_w \hat{V}(s, w) = \nabla_w \mathbf{x}(s)^T \mathbf{w} = x(s)$$

$$\Delta w = A(V^\pi(s) - \hat{V}(s, w)) \nabla_w \hat{V}(s, w) \quad (5)$$

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 $V^\pi(s_t)$
- 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 Monte-Carlo policy evaluation

$$\Delta w = (A_r (R_t - \hat{V}(s, w)) \nabla_w \hat{V}(s, w)) \quad (8)$$

$$= (A_r (R_t - \hat{V}(s, w)) \mathbf{x}(s)) \quad (9)$$

Monte-Carlo evaluation converges to a local optimum, even when using non-linear value function approximation (provable)

TD Learning with Value Function Approximation I
The TD-target $R_{t+1} + \gamma \hat{V}(s_{t+1}, w)$ is a biased (single) sample of the true value $V^\pi(s_t)$

We still perform supervised learning on "digested" training data:

$$(s_1, r_1 + \gamma \hat{V}(s_2, w)), (s_2, r_2 + \gamma \hat{V}(s_3, w)), \dots, (s_T, r_T) \quad (10)$$

The estimates make the estimates closer to the real $V(s)$ value, but not reach the true $V(s)$ value. The estimate will get closer to the true $V(s)$ value.

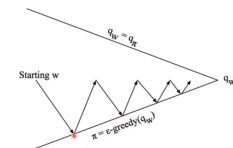
SARSA – On-Policy learning TD control
Initialization: $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
Repeat (for each episode):
1. Choose s from π using policy derived from Q (e.g., c-greedy)
2. Take action a , observe r, s'
3. Choose a' from s' using policy derived from Q (e.g., c-greedy)
 $Q(s', a') = Q(s, a) + \alpha [r + \gamma Q(s', a') - Q(s, a)]$
 $s \leftarrow s'$
until s is terminal

SARSA converges to the optimal action-value function $Q(s, a) = Q^*(s, a)$, under the following conditions

- GLIE sequence of policies $\pi^t(s, a)$
- Robins-Munro sequence of step-sizes α_t :

$$\sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\sum_{t=1}^{\infty} \alpha_t^2 < \infty$$



- Policy evaluation: Approximate policy evaluation $Q(s, a, \pi) \approx Q^{\pi}(s, a)$
- Policy improvement: ϵ -greedy policy improvement

DQN

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:

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

we want to learn $Q(s_t, a; w)$ 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.

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

Taking the gradient of E wrt w we obtain:

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

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

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

DQN: Experience Replay I

The CNN is easily overfitting the latest experienced episodes and the network becomes worse at dealing with different experiences of the game world:

- Complication 1: Inefficient use of interactive experience
 - Training deep neural networks requires many samples, each has its own transition
 - But now each state transition is used only once, then discarded
 - So we need to constantly revisit the same state transitions to train
- This is inefficient and slow.

Complication 2: Experience distribution

- Interactive learning produces training samples that are highly correlated (agents 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 "extinction")
- 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 Replay? (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 pair as they occur during simulation or actual experience, the experience replay buffer 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.

The learning phase is separated from gaining experience, and based on random samples from the mini-batch table.

DQN interleaves the two processes - acting and learning. Improving the policy will lead to a different behaviour (and different state action pairs), resulting too large replay buffers getting stale.

In summary

- Reduction of correlation between experiences in updating DNN - mini-batches effectively make the samples more independent and identically distributed.

- Increased learning speed with mini-batches due to reuse/multi-use of experience

- Reusing/Replacing past transitions to avoid catastrophic forgetting of associated rewards

DQN: Experience Replay V

Initialise replay memory \mathcal{D} to capacity N

Initialise action-value function Q with random weights

for episode = 1..M do

Initialise state s_t

for t = 1..T do

 Action a_t = ϵ -greedy select a random action a_t

 Execute action a_t and observe reward r and state s_{t+1}

 Store transition (s_t, a_t, r, s_{t+1}) in \mathcal{D}

 Set $\theta_t = \theta$

 Sample random minibatch of transitions (s_t, a_t, r, s_{t+1}) from \mathcal{D}

 Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_a Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(s_t, a_t; \theta))^2$

 end for

end for

DQN: Target Network I

Whenever we run an update step on a Q-network's state we also update "near-by" states (monolithic architecture + generalisation ability of CNNs).

$$Q(s_t, a) \leftarrow Q(s_t, a) + \alpha[r_{t+1} + \gamma \max Q(s_{t+1}, a') - Q(s_t, a)] \quad (4)$$

Issue: Unstable training resulting from bootstrapping a continuous state space representation

- We may have an unstable learning process, because changes to $Q(s'_t, a)$ change $Q(s'_t, a')$ (with $s = s' \Rightarrow s'$) which changes $Q(s_t, a)$ on the next CNN and update in turn $Q(s'_t, a')$ forth.

- The effect may result in a run-away bias from bootstrapping and subsequently dominating the system numerically, causing the estimated Q values to diverge.

- In calculating the TD-error, the target function is changing too frequently when using convolutional neural networks.

Solution: Slow things down (as resonance damper) Using a separate target network Q' , updated every few timer (e.g. every 1000 steps) with a copy of the latest learned weight parameters, controls this stability (relaxation time).

- Initialising two Q networks: a main Q-network (Q), and a target network (Q')

When calculating the TD-error, use Q' , not Q

- Infrequently set $Q = Q'$

This gives the highly fluctuating Q-time to settle (we can measure this so called "relaxation time" empirically and set it accordingly) before updating Q'

Closing Remarks: Variation in rewards makes the training unstable → Clip positive reward to 1, and negative reward to -1

DQN: Clipping rewards I

Clipping rewards:

- Each ATARI game has different score scales. For example, in Pong, players get 1 point when winning the play.

Otherwise, players get -1 point. However, in Space Invaders, players get 10-30 points when defeating invaders. This difference would make training unstable.

Clipping Rewards: clips the rewards, with all positive rewards are set +1 and all negative rewards are set -1.

DQN: Skipping frames I

Skipping frames:

- Computer can simulate games (e.g. the ATARI simulator does this at 60 Hz) much faster than a human player would be able to react to the game, and this implies that the video game design and the required actions can be slower.
- Frame skipping uses only every 4 video game frames (i.e. 15 Hz...) (and frame stacking uses the past 4 frames as inputs), thus reducing computational cost and accelerating training times. It also makes the game run at a speed comparable to human reaction time.

Double Q network

Normal netowrk will produce a maximisation bias

Max at the Casino II

- So when update the value for $Q(Start, Left)$ we perform a max over actions in successor state Casino (i.e. from traces yielding r in Casino, A , Hotel) = -1, r (Casino, B , Hotel) = +1, r (Casino, A , Hotel) = -1) and immediately note down the +1 reward from choosing B .
- In the limit of large number of trials all actions will yield an average reward of -0.02, but that is only asymptotically reached.
- This results in a maximisation bias, namely that when taking a max over all actions with (very) finite data we may always overestimate values. This illustrates a general relationship (that we do not prove here), but we can see empirically.

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,
- Obtain 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 there is an unusually high value from the target network Q' . The problem with Deep-Q-Learning 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 updates become independent from each other (different networks) so we become less susceptible to variance 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 was the crucial first step at popularising Policy Gradients (PGs)

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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

REINFORCE algorithm: sample from $\pi(a_t | s_t)$ (run the policy)

$$\nabla_\theta J(\theta) \approx \sum_t \left[\nabla_\theta \log \pi(a_t | s_t) \right] \left[\sum_t r(s_t, a_t) \right]$$

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