

Contents

1	Introduction	1
2	Lecture two: Markov Decision Process (MDP)	2
2.1	Markov Process	2
2.2	Markov Reward process (MRP)	2
2.3	Markov Decision Process (MDP)	4
2.4	The Optimal Value Function	6
3	Lecture Three: Planning by Dynamic Programming	8
3.1	Introduction	8
3.2	Synchronous Dynamic programming algorithms	8
3.3	Asynchronous dynamic programming	10
4	Lecture Four: Model-Free Prediction	12
4.1	Monte-Carlo Learning	12
4.2	Temporal-Difference Learning	13
4.3	TD(λ)	16
5	Model-Free Control	19
5.1	On and Off-policy learning	19
5.2	ϵ -Greedy Exploration	19
5.3	TD(λ) Control	20
5.4	Off-Policy Learning	23

1 Introduction

Script of the Reinforcement Learning Course by David Silver as uploaded on YouTube from 2015

written by Alex Liebenau

created in November 2022

free to use

2 Lecture two: Markov Decision Process (MDP)

2.1 Markov Process

The current state characterises the process - i.e. we are told the state - i.e. environment is fully observable

- Almost all RL problems can be formalised as MDPs
- i.e. Optimal control primarily deals with continuous MDPs
- i.e. Any partially observable problems can be converted into MDPs
- i.e. Bandits are MDPs with one state

"The future is independent of the past given the present"

$$\mathbb{P}[S_{t+1} | S_t] = \mathbb{P}[S_{t+1} | S_1, \dots, S_t] \quad (1)$$

What happens next only depends on what happened on the state before - you can throw away anything else.

For a Markov state s and successor state s' , the *state transition probability* is defined as

$$\mathcal{P}_{ss'} = \mathbb{P}[S_{t+1} = s' | S_t = s] \quad (2)$$

State transition matrix ρ defines transition probabilities from all states s to all successor states s' .

$$\mathcal{P} = \underset{\text{from}}{\begin{pmatrix} \mathcal{P}_{11} & \dots & \mathcal{P}_{1n} \\ \vdots & & \\ \mathcal{P}_{n1} & \dots & \mathcal{P}_{nn} \end{pmatrix}} \overset{\text{to}}{\quad} \quad (3)$$

Each row of the matrix sums to 1!

A Markov process is a memoryless random process, i.e. a sequence of random states, S_1, S_2, \dots with the Markov property. It is defined as a tuple $\langle \mathcal{S}, \mathcal{P} \rangle$.

- \mathcal{S} is a (finite) set of states
- \mathcal{P} is a state transition probability matrix

2.2 Markov Reward process (MRP)

The Markov Reward Process is defined as a tuple $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma \rangle$.

- \mathcal{S} is a (finite) set of states
- \mathcal{P} is a state transition probability matrix
- \mathcal{R} is a reward function, $\mathcal{R}_s = \mathbb{E}[R_{t+1} | S_t = s]$

- γ is a discount factor, $\gamma \in [0, 1]$

The *return* \mathcal{G}_t is the total discounted reward from time-step t .

$$\mathcal{G}_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (4)$$

- The discount γ is the present value of future rewards - the closer γ is to zero, the less are later rewards accounted (e.g. more 'short-sighted').
- The value of receiving reward \mathcal{R} after $k + 1$ time-steps is $\gamma^k \mathcal{R}$.

Why discount?

- Unless you really trust your model and believe that everything turns out as planned, you need to discount in deviations - Uncertainty about the future may not be fully represented
- Mathematically convenient to discount rewards, avoids infinite returns
- If the reward is financial, immediate rewards may earn more interest than delayed rewards
- Animal/human behaviour shows preference for immediate reward
- It is sometimes possible to use *undiscounted* Markov reward process (i.e. $\gamma = 1$), e.g. if all sequences terminate.

The value function

The value function $v(s)$ gives the long-term value of state s . It defines the expected return in a MRP starting from state s :

$$v(s) = \mathbb{E} [G_t \mid S_t = s] \quad (5)$$

The value function can be decomposed into two parts:

- immediate reward $R + 1$
- discounted value of successor state $\gamma v(S_{t+1})$

This resolves into the **Bellman equation for MRPs**:

$$v(s) = \mathbb{E} [G_t \mid S_t = s] = \mathbb{E} [R_{t+1} + \gamma v(S_{t+1}) \mid S_t = s] \quad (6)$$

By averaging all possible outcomes we get

$$v(s) = \mathcal{R}_s + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'} v(s') \quad (7)$$

The Bellman equation can be concisely using matrices, where v is a column vector with one entry per state

$$v = \mathcal{R} + \gamma \mathcal{P} v \rightarrow \begin{pmatrix} v(1) \\ \vdots \\ v(n) \end{pmatrix} = \begin{pmatrix} \mathcal{R}_1 \\ \vdots \\ \mathcal{R}_n \end{pmatrix} + \gamma \begin{pmatrix} \mathcal{P}_{11} & \dots & \mathcal{P}_{1n} \\ \vdots & & \vdots \\ \mathcal{P}_{n1} & \dots & \mathcal{P}_{nn} \end{pmatrix} \begin{pmatrix} v(1) \\ \vdots \\ v(n) \end{pmatrix} \quad (8)$$

The Bellman equation is linear. It can be solved directly by $v = (I - \gamma \mathcal{P})^{-1} \mathcal{R}$.

- The Computational complexity is $O(n^3)$ for n states.
- Direct solution only possible for small MRPs
- Iterative methods for large MRPs, e.g. Dynamic programming, Monte-Carlo evaluation, Temporal-Difference learning

2.3 Markov Decision Process (MDP)

A MDP is a MRP with decisions. It is an *environment* in which all states are Markov.

A MDP is a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$.

- \mathcal{S} is a (finite) set of states
- \mathcal{A} is a (finite) set of actions
- \mathcal{P} is a state transition probability matrix
 $\mathcal{P}_{ss'}^a = \mathbb{P}[S_{t+1} = s' | S_t = s, A_t = a]$
- \mathcal{R} is a reward function, $\mathcal{R}_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$
- γ is a discount factor, $\gamma \in [0, 1]$

A policy π is a distribution over actions given states. It fully defines the behaviour of an agent.

$$\pi(a|s) = \mathbb{P}[S_t = s, A_t = a] \quad (9)$$

In an MDP, the policies depend on the current state (not the history). Policies are stationary: $A_t = \pi(|S_t), \forall t > 0$

Given an MDP $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ and a policy π :

- The state sequence S_1, S_2, \dots is a Markov process $\langle \mathcal{S}, \mathcal{P}^\pi \rangle$
- The state and reward sequence $S_1, R_1, S_2, R_2, \dots$ is a Markov reward process $\langle \mathcal{S}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ where

$$\mathcal{P}_{s,s'}^\pi = \sum_{a \in \mathcal{A}} \pi(a|s) \mathcal{P}_{s,s'}^a \quad \mathcal{R}_s^\pi = \sum_{a \in \mathcal{A}} \pi(a|s) \mathcal{R}_s^a \quad (10)$$

The *state-value function* $v_\pi(s)$ of an MDP is the expected return starting from state s , and then following policy π :

$$v_\pi(s) = \mathbb{E}_\pi[G_t \mid S_t = s] \quad (11)$$

The *action-value function* $q_\pi(s, a)$ of an MDP is the expected return starting from state s , taking action a , and then following policy π :

$$q_\pi(s, a) = \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a] \quad (12)$$

The state-value and action-value functions can again be decomposed into Bellman equations consisting of immediate reward plus discounted value of successor state:

$$v_\pi = \mathbb{E}_\pi [R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s] \quad (13)$$

$$q_\pi = \mathbb{E}_\pi [R_{t+1} + \gamma v_\pi(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a] \quad (14)$$

Basically, the state-value averages over the different actions that can be taken:

$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) q_\pi(s, a) \quad (15)$$

The other way around, by using the probabilities of the transition dynamics we can average through the values of the succeeding states we can evaluate a certain action:

$$q_\pi(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s') \quad (16)$$

- v is telling how good it is to be in a particular state
- q is telling how good it is to take a particular action

Sticking these together in order to solve MDPs end up in a two-step lookahead:

- Consider all action we might take next (v)
- Consider all the things the environment might do to us (q)
- Evaluate the successor state after what the environment did after that point

By double-averaging over the policy as well as the transition probability we get the answer to how good it is to be in a particular state (as the Bellman Equation):

$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_\pi(s') \right) \quad (17)$$

Starting from a particular action, we can also do the same two-step lookahead and see where the wind blows us, and from there consider which action might be taken next. By averaging the same way above we get same recursive relationship:

$$q_{\pi}(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \sum_{a' \in \mathcal{A}} \pi(a'|s') q_{\pi}(s', a') \quad (18)$$

This shows how the value- and action function of the next step relates to itself, we therefore get a recursive relationship. **After all, the value function of the current time step is equal to the immediate reward plus the value function of where you end up.**

To flatten the MDP into a MRP one could average out the matrices as in (7) following policy π and then solve accordingly as in (8). Thereby we can determine the value function.

----- **For more info on undiscounted MDPs:**

Neurodynamic Programming / Dynamic Programming; Dimitri Bertsekas
Dynamic Programming and Optimal Control (2 Vol Set) -----

2.4 The Optimal Value Function

The optimal *state-value* and *action-value* functions are the maximum functions over all policies

$$v_*(s) = \max_{\pi} v_{\pi}(s) \quad q_*(s, a) = \max_{\pi} q_{\pi}(s, a) \quad (19)$$

- The optimal value function specifies the best performance in an MDP
- An MDP is 'solved' when we know the optimal value function

If you know q_* , you're basically done. It tells you under all different ways to behave which one gives the most reward.

There is always one optimal policy. This optimal policy will then achieve the optimal *state-value* and *action-value* function v_{π_*} and q_{π_*} . So we can define a partial ordering over policies:

$$\pi \geq \pi' \quad \text{if} \quad v_{\pi}(s) \geq v_{\pi'}(s), \forall s \quad (20)$$

An optimal policy can be found by maximising over $q_{\pi_*}(s, a)$:

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_{a \in \mathcal{A}} q_{\pi_*}(s, a) \\ 0 & \text{otherwise} \end{cases} \quad (21)$$

The Bellman Optimality Equation

The optimal value functions are recursively related by the Bellman optimality equations:

Instead of taking the average of all actions, you pick the maximum. That determines the optimal value function of a state:

$$v_* = \max_a q_*(s, a) \quad (22)$$

Now looking at the action-value function, we can inductively assume that each of the states we might end up in has a optimal state-value v_* . Therefore, similar as for (16), to determine the action-value we add the immediate reward to the average of the optimal state-values where we could end up in:

$$q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s') \quad (23)$$

As we did before in (17), we can put these two together to create the recursive relationship with the two-step lookahead. For the optimal solution, we pick the maximum instead of the average. This is the Bellman optimality equation for v_* :

$$v_*(s) = \max_{a \in \mathcal{A}} \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s') \right) \quad (24)$$

By reordering the same idea, we can again define the Bellman optimality equation for q_* .

$$q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a \sum_{a' \in \mathcal{A}} \max_{a'} q_*(s', a') \quad (25)$$

Solving the Bellman Optimality Equation

- The Bellman Optimality Equation is non-linear
→ Solving by inverting matrices like in (7) will not work
- No closed form solution (in general)
- Many iterative solution methods
→ Value iteration, policy iteration, **Q-Learning**, Sarsa

3 Lektion Three: Planning by Dynamic Programming

3.1 Introduction

Dynamic: sequential or temporal component to problem

Programming: optimising a 'problem', e.g. a policy

→ Solving complex problems by breaking them down into subproblems, then combine their solutions

Properties of problems for dynamic programming:

- Optimal substructure:
 - Principle of optimality applies
 - Optimal solution can be decomposed into subproblems
- Overlapping substructure
 - Subproblems reoccur many times → Solutions can be cached and reused

MDPs satisfy both properties. The Bellman Equation gives the recursive decomposition. The value function stores and reuses the solutions. Dynamic Programming thereby assumes full knowledge of the MDP, it is used for *planning* in an MDP.

For prediction

- **Input:** An MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ and policy π or an MRP $\langle \mathcal{S}, \mathcal{P}^\pi, \mathcal{R}^\pi, \gamma \rangle$
- **Output:** value function v_π

or

For Optimal Control

- **Input:** An MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
- **Output:** Optimal value function v_* and optimal policy π_*

3.2 Synchronous Dynamic programming algorithms

Policy Evaluation and Policy Iteration

Solution: iterative application of Bellman Equation $v_1 \rightarrow v_2 \rightarrow v_3 \rightarrow \dots v_n$
Using *synchronous* backups:

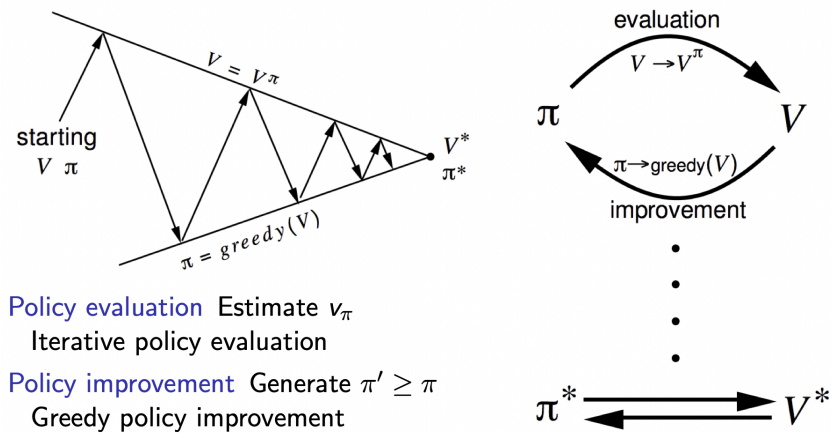
- At each iteration $k + 1$
- For all states $s \in \mathcal{S}$

- Update $v_{k+1}(s)$ from $v_k(s')$ where s' is a successor state of s

Given a policy π , how do you find out if it's optimal?

- *Evaluate* the policy $\pi \rightarrow$ Compute the value function v_π
- *Improve* the policy by acting greedily with respect to v_π

Acting greedily means picking the action that returns the most reward. This will in every case improve the value function. In general, it always needs a lot of improvement and evaluation. This process will *always* converge to π_* , no matter in which state you begin.



What does policy improvement mean?

Consider a deterministic policy $a = \pi(s)$

- We can *improve* the policy by acting greedily: $\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} q_{\pi_*}(s, a)$
- This will improve the value from any state s over one step:
 $q_\pi(s, \pi'(s)) = \max_{a \in \mathcal{A}} q_{\pi_*}(s, a) \geq q_{\pi_*}(s, a) + v_\pi$
- Keep iterating over the coming time-steps
- If improvement stops:
 $q_\pi(s, \pi'(s)) = \max_{a \in \mathcal{A}} q_\pi(s, a) = q_\pi(s, a) = v_\pi$
 \rightarrow Bellman optimality equation is satisfied: $v_\pi(s) = \max_{a \in \mathcal{A}} q_\pi(s, a)$
- Therefore: $v_\pi = v_*$ for all $s \in \mathcal{S} \rightarrow \pi$ is an optimal policy

Value iteration

Any optimal policy can be subdivided into two components:

- An optimal first action A_*

- Followed by an optimal policy from successor state S'

The policy is optimal, if from each state we might end up the policy is optimal from that state onwards. If we know the solution to subproblems $v'_*(s)$, the solution $v_*(s)$ can be found iteratively by one-step lookahead:

$$v_*(s) \leftarrow \max_{a \in \mathcal{A}} \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s') \right) \quad (26)$$

The idea of value iteration is to apply these updates iteratively, thereby moving through the whole state-space. The intuition is that you begin with your final reward and work your way backwards.

Unlike policy iteration, there is no explicit policy. Intermediate value functions may not correspond to any policy.

Problem	Bellman Equation	Algorithm
Prediction	Bellman Expectation Equation	Iterative Policy Evaluation
Control	Bellman Expectation Equation + Greedy Policy Improvement	Policy Iteration
Control	Bellman Optimality Equation	Value Iteration

These algorithms are based on the state-value function $v_\pi(s)$ or $v_*(s)$. They have a complexity of $\mathcal{O}(mn^2)$ per iteration, for m actions and n states. The algorithms could also apply to $q_\pi(s)$ or $q_*(s)$, where the complexity would be $\mathcal{O}(m^2n^2)$.

3.3 Asynchronous dynamic programming

It is not necessary to update every state at each sweep of the algorithm. In *asynchronous* backups, the relationship of updating every state in each iteration is broken. This saves a lot of unnecessary work and as long as all states keep being selected at some time, the algorithm will still converge to the optimal solution.

In-place dynamic programming

Synchronous value iteration stores two copies of the value function:

$$v_{new}(s) \leftarrow \max_{a \in \mathcal{A}} \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_{old}(s') \right) \quad (27)$$

$$v_{old} \leftarrow v_{new} \quad (28)$$

In-place value iteration only stores one copy of value function by using the latest version of the value function available. That means if a state has been visited

before, the algorithm uses the already calculated value instead of the old one. This tends to be more efficient, but the order of the states being swept is not relevant for the efficiency of the algorithm.

Prioritized sweeping

By using the magnitude of the Bellman error to guide state selection, the states with the largest difference after an update have a higher priority. These will have the most effect in downstream computing. It is important though to update the Bellman error after each backup, which required knowledge of reverse-dynamics of a system for the predecessor states. The Bellman error is described as following:

$$\left| \max_{a \in \mathcal{A}} \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_{old}(s') \right) - v(s) \right| \quad (29)$$

Real-time Dynamic Programming

The idea of *Real-time programming* is to only update the states which are relevant to the agent. The experience of the agent is used to guide the selection of states. After each time-step S_t, A_t, R_{t+1} , a backup of state S_t is done.

$$v(S_t) \leftarrow \max_{a \in \mathcal{A}} \left(\mathcal{R}_{S_t}^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{S_t s'}^a v(s') \right) \quad (30)$$

Full-Width and Sample backups

Dynamic programming uses *full-width backups*, meaning that every state is backed up. For medium-sized problems consisting of a few millions of states, this is fine. But for large problems, dynamic programming suffer from the *curse of dimensionality*, following from the exponential growth of the number of states with the number of state variables.

This is where sample backups come into play. Instead of backing up the whole state-space (which can be already too expensive for one iteration), only certain samples are backed up. The advantages are:

- Model-free: No advance knowledge of MDP required
- Breaks the curse of dimensionality through sampling
- Cost of backup is constant, independent of $|n| = \mathcal{S}$

4 Lecture Four: Model-Free Prediction

4.1 Monte-Carlo Learning

Monte-Carlo Learning methods learn directly from episodes of experience. It is therefore important that the MDP is applied to *episodic* MDP, where the (complete) episodes do terminate. The value function is assumed to be the simplest possible, the mean return of the samples.

The goal is to learn the value function v_π of a random policy π by looking at some streams $S_1, A_1, R_2, \dots, S_k \sim \pi$ of episode of experience and then evaluate the *return* as the total discounted reward \mathcal{G}_t from (4). Instead of using the *expected* return of the value function, Monte-Carlo Learning uses the *empirical mean* from the sampled episodes.

First-Visit Monte-Carlo policy evaluation

- To evaluate the state s , the first time-step it is visited in an episode, an incremental counter is initialized: $N(s) \leftarrow N(s) + 1$
- The total return is added up: $S(s) \leftarrow S(s) + G_t$
- The value of the state is estimated by mean return: $V(s) = S(s)/N(s)$

By the law of large numbers, the value of the state $V(s)$ converges to the actual value v_π of the state is visited many times so that the estimated mean consists of a large number of visits: $V(s) \rightarrow v_\pi$ as $N(s) \rightarrow \infty$

The **Every-Visit Monte-Carlo policy evaluation** works in a very similar fashion, just that the counter and return is incremented - as the name suggests - every time the episode visits a state.

Incremental Monte-Carlo

The mean μ_1, μ_2, \dots of a sequence x_1, x_2, \dots can be computed incrementally:

$$\begin{aligned}\mu_k &= \frac{1}{k} \sum_{j=1}^k x_j \\ &= \frac{1}{k} \left(x_k + \sum_{j=1}^{k-1} x_j \right) \\ &= \frac{1}{k} \left(x_k + (k-1) \mu_{k-1} \right) \\ &= \mu_{k-1} + \frac{1}{k} \left(x_k - \mu_{k-1} \right)\end{aligned}\tag{31}$$

The difference $x_k - \mu_{k-1}$ basically defines the error term, where μ_{k-1} represents the expected value and x_k the actual value. In the process of (31), the value of the mean will then get an update slightly into the direction of the error term.

This can be applied to the Monte-Carlo Learning. The value $V(s)$ of a state s will be updated incrementally after an episode $S_1, A_1, R_2, \dots, S_t$:

- For each state S_t with return G_t :

$$N(s) \leftarrow N(s) + 1$$

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(s)}(G_t - V(S_t))$$
- In non-stationary problems, it can be useful to track a running mean, e.g. forget old episodes

$$V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$$

4.2 Temporal-Difference Learning

Like with Monte-Carlo Learning, Temporal-Difference Learning learns from episodes of experience and needs no knowledge of the MDP transitions or rewards, they are both *model-free*. In contrast to Monte Carlo Learning however, Temporal-Difference Learning also learns from incomplete episodes, they don't have to terminate. By taking partial episodes and estimating the remaining return, it can learn from incomplete episodes. It iteratively updates the guesses of the return while going through an episode. This is called *bootstrapping*.

The simplest Temporal-Difference Learning algorithm: TD(0)

- The value of V_s is updated towards the *estimated* return: $R_{t+1} + \gamma V(S_{t+1})$

$$V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

- $R_{t+1} + \gamma V(S_{t+1})$ is called the TD target (the estimated return)
- $\delta_t = (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$ is called the TD error

As seen in caption 1, the predicted total travel time gets updated after every state in TD Learning, while in MC learning all the values of the states will be updated after the episode.

Advantages and Disadvantages of Temporal-Difference vs. Monte-Carlo Learning

- TD can learn before knowing the final outcome
 - TD can learn online after every step
 - MC must wait until the end of an episode before return is known
- TD can learn without the final outcome
 - TD can learn from incomplete sequences
 - MC can only learn from complete sequences
 - TD works in continuing environments

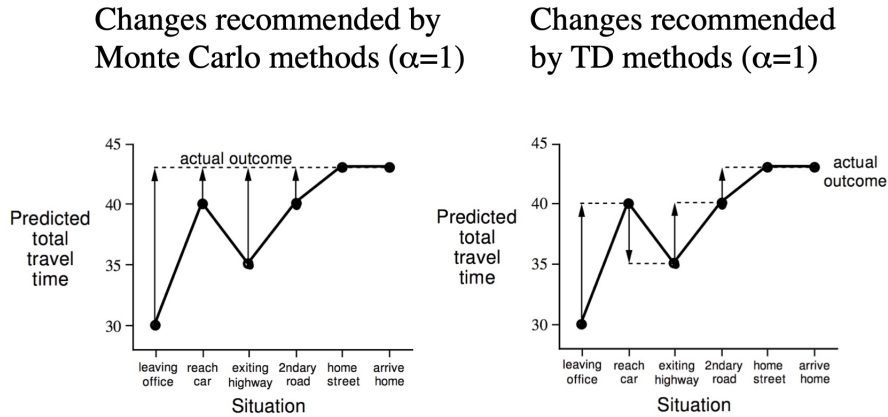


Figure 1: The difference between Monte-Carlo and Temporal-Difference Learning.

- MC only works for episodic (terminating) environments
- Certainty equivalence:
 - MC converges to solution with minimal mean-squared error
→ Best fit to the observed returns
 - TD converges to solution of max likelihood Markov model
→ Solution to the MDP that best fits the data
- TD exploits Markov property by using it, therefore is more efficient in a Markov environment. MC ignores the Markov property, therefore does not exploit it.

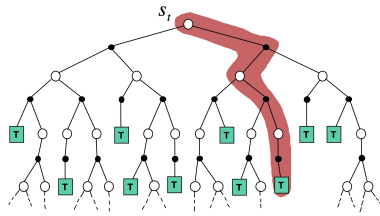
Bias / Variance Trade-Off

- The return G_t and the true TD target $R_{t+1} + \gamma v_\pi(S_{t+1})$ are both *unbiased* estimates of $v_\pi(S_t)$.
- The actual TD target $R_{t+1} + \gamma V(S_{t+1})$ is a *biased* estimate of $v_\pi(S_t)$. It is much lower variance than the return:
 - Return depends on *many* random actions, transitions and rewards.
 - TD target depends on *one* random action, transition and reward.
- MC has high variance, zero bias
 - Good convergence properties, even with function approximation
 - Not very sensitive to initial value

- Very simple to understand and use.
- TD has low variance, some bias
 - Usually more efficient than MC
 - TD(0) converges to $v_\pi(s)$ - but not always with function approximation
 - More sensitive to initial value

Monte-Carlo

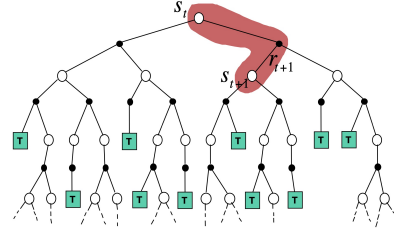
$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$



Backup:

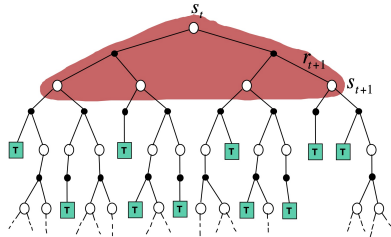
Temporal-Difference Backup:

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$



Dynamic programming

$$V(S_t) \leftarrow \mathbb{E}_\pi [R_{t+1} + \gamma V(S_{t+1})]$$



Pro-Backup:

• Bootstrapping:

update involves an estimate

- MC does not bootstrap
- TD and DP do

• Sampling:

update samples an expectation

- MC and TD sample
- DP does not sample

Note: The 'Exhaustive search' would consider every branch as in Dynamic Programming, but then follow its leafs until the end like Monte-Carlo instead of taking an estimate.

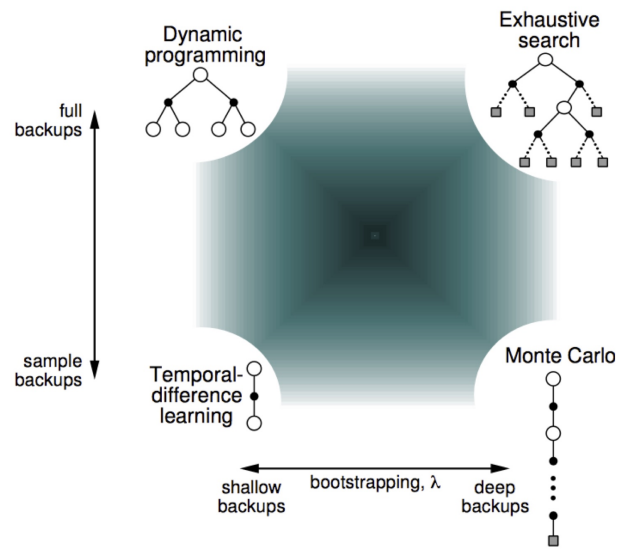
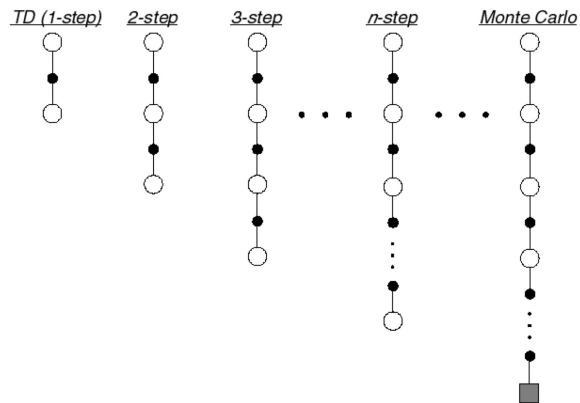


Figure 2: Unified view of Reinforcement Learning

4.3 $TD(\lambda)$

n-step Prediction

- Let TD target look n steps into the future



n-step Return

$$\begin{array}{ll}
 n = 1 & (TD) \quad G_t^{(1)} = R_{t+1} \gamma V(S_{t+1}) \\
 n = 2 & G_t^{(2)} = R_{t+1} \gamma R_{t+2} + \gamma^2 V(S_{t+2}) \\
 \vdots & \vdots \\
 n = \infty & (MC) \quad G_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \gamma^{T-1} R_T
 \end{array}$$

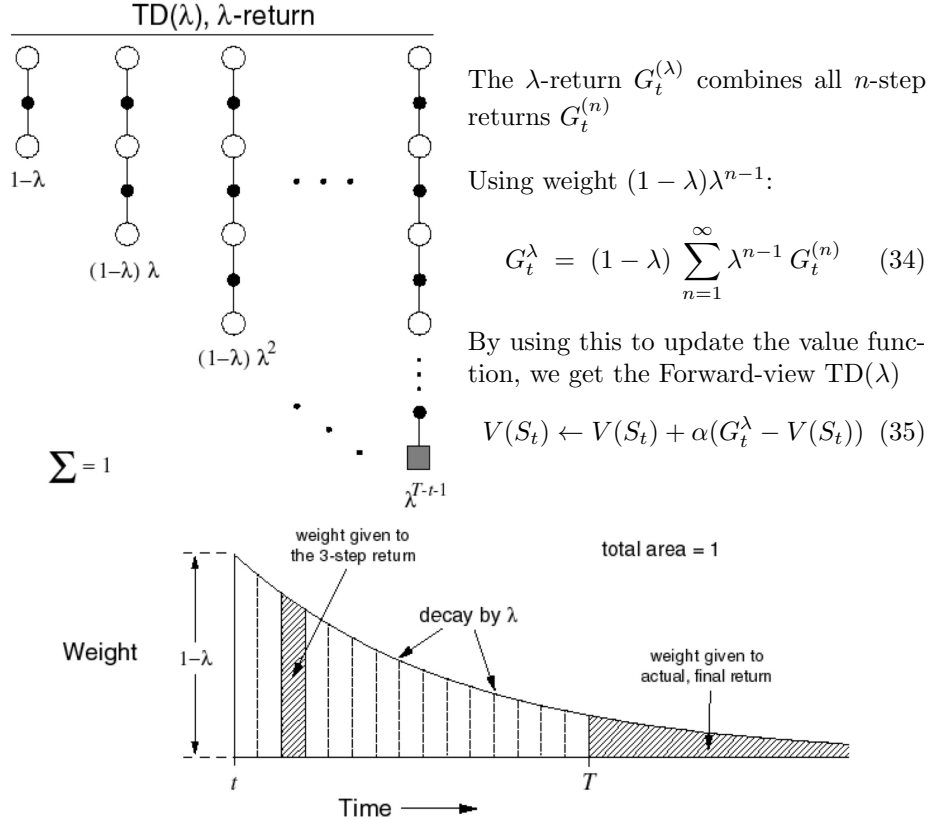
Define the n -step return for the estimator:

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n}) \quad (32)$$

n -step Temporal-Difference Learning:

$$V(S_t) \leftarrow V(S_t) + \alpha(G_t^{(n)} - V(S_t)) \quad (33)$$

Forward-view of TD(λ)



The Forward-view looks into the future to compute $G_t^{(\lambda)}$. Like Monte-Carlo Learning, it can only be computed from complete episodes.

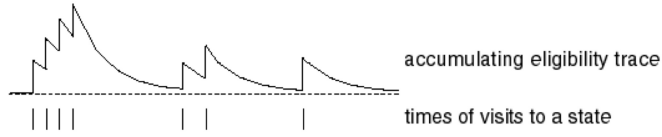
Backward-view of TD(λ)

The Forward-view provides theory, the Backward-view provides mechanism

It updates online, every step, from incomplete sequences.

- **Frequency heuristic:** Assign credit to most frequent states
- **Recency heuristic:** Assign credit to most recent states

→ *Eligibility traces* combine both heuristics!



$$E_0(s) = 0$$

$$E_t(s) = \gamma\lambda E_{t-1}(s) + \mathbf{1}(S_t = s)$$

Backward-view TD(λ)-Algorithm:

- Keep an eligibility trace for every state s
- Update value $V(s)$ for every state s
- In proportion to TD-error δ_t and eligibility trace $E_t(s)$

$$\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \quad (36)$$

$$V(s) \leftarrow V(s) + \alpha \delta_t E_t(s) \quad (37)$$

Relationship between Forward and Backward TD(λ)

The sum of the offline¹ updates is the same identical for forward-view and backward-view TD(λ)

$$\sum_{t=1}^T \alpha \delta_t E_t(s) = \sum_{t=1}^T \alpha (G_t^\lambda - V(S_t)) \mathbf{1}(S_t = s) \quad (38)$$

Offline updates	$\lambda = 0$	$\lambda \in (0, 1)$	$\lambda = 1$
Backward view	TD(0)	TD(λ)	TD(1)
	\parallel	\parallel	\parallel
Forward view	TD(0)	Forward TD(λ)	MC
Online updates	$\lambda = 0$	$\lambda \in (0, 1)$	$\lambda = 1$
Backward view	TD(0)	TD(λ)	TD(1)
	\parallel	\nparallel	\nparallel
Forward view	TD(0)	Forward TD(λ)	MC
	\parallel	\parallel	\parallel
Exact Online	TD(0)	Exact Online TD(λ)	Exact Online TD(1)

¹Even for online updating it is now possible to achieve equivalence for forward and backward view TD(λ) - see

5 Model-Free Control

Model-free prediction: *Estimate* the value function of an *unknown* MDP

Model-free control: *Optimise* the value function of an *unknown* MDP

Model-free control can solve problems whether either the MDP model is unknown, but experience can be sampled or where the MDP model is known but too big to use - except by samples (e.g. when the underlying real-world dynamics are very complicated).

5.1 On and Off-policy learning

On-policy learning:

- "Learn on the job"
- "Learn about policy π from experience sampled from π "

Off-policy learning:

- "Look over someone's shoulder"
- "Learn about policy π from experience sampled from μ "

Model-Free Policy iteration using action-value function

Greedy policy improvement over the state value function $V(s)$ requires a model of the MDP:

$$\pi'(s) = \underset{a \in \mathcal{A}}{\operatorname{argmax}} \mathcal{R}_s^a + \mathcal{P}_{ss}^a V(s') \quad (39)$$

...whereas greedy policy improvement over the action value function $Q(s, a)$ is model-free

$$\pi'(s) = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q(s, a) \quad (40)$$

5.2 ϵ -Greedy Exploration

This is the simplest idea for ensuring continual exploration.

- All m actions are tried with non-zero probability
- With probability $1-\epsilon$ choose the greedy action
- With probability ϵ choose an action at random

$$\pi(a|s) = \begin{cases} \epsilon/m + 1 - \epsilon & \text{if } a = \underset{a \in \mathcal{A}}{\operatorname{argmax}} Q(s, a) \\ \epsilon/m & \text{otherwise} \end{cases} \quad (41)$$

GLIE: Greedy in the Limit with Infinite Exploration

- All state-action pairs are explored infinitely many times:
 $\lim_{k \rightarrow \infty} N_k(s, a) = \infty$

- The policy converges on a greedy policy
 $\lim_{k \rightarrow \infty} \pi_k(a|s) = \mathbf{1}(a = \operatorname{argmax}_{a' \in \mathcal{A}} Q_k(s, a'))$

→ for example, ϵ -greedy is GLIE if ϵ reduces to zero at $\epsilon_k = \frac{1}{k}$

GLIE Monte-Carlo Control

- Sample k -th episode using $\pi : S_1, A_1, R_2, \dots, S_T \sim \pi$
- For each state S_t and action A_t in the episode,

$$\begin{aligned} N(S_t, A_t) &\leftarrow N(S_t, A_t) + 1 \\ Q(S_t, A_t) &\leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)} (G_t - Q(S_t, A_t)) \end{aligned}$$

- Improve policy based on new action-value function

$$\begin{aligned} \epsilon &\leftarrow 1/l \\ \pi &\leftarrow \epsilon - \text{greedy}(Q) \end{aligned}$$

5.3 TD(λ) Control

MC vs. TD control

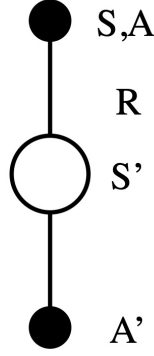
Temporal difference (TD) Learning has several advantages over Monte-Carlo (MC):

- Lower variance
- online
- Incomplete sequences

Natural idea: use TD instead of MC in our control loop

- Apply TD to $Q(S, A)$
- Use ϵ -greedy policy improvement
- Update every time step

Updating Action-Value Functions with *Sarsa*



$$Q^{new}(s_t, a_t) = \underbrace{Q(s_t, a_t)}_{\text{current value}} + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\left(\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_{a \in \mathcal{A}} Q(s_t, a_{t+1})}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{current value}} \right)}_{\text{new value (temporal difference target)}} \quad (42)$$

The *Sarsa*-algorithm for on-policy learning would be the following:

```

Initialize  $Q(s, a), \forall s \in \mathcal{S}, a \in \mathcal{A}(s)$ , arbitrarily, and  $Q(\text{terminal-state}, \cdot) = 0$ 
Repeat (for each episode):
  Initialize  $S$ 
  Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
  Repeat (for each step of episode):
    Take action  $A$ , observe  $R, S'$ 
    Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)
     $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$ 
     $S \leftarrow S'; A \leftarrow A'$ 
  until  $S$  is terminal
    
```

→ ***n*-Step Sarsa**

$$\begin{array}{ll}
 n = 1 & (Sarsa) \quad q_t^{(1)} = R_{t+1} + \gamma Q(S_{t+1}) \\
 n = 2 & q_t^{(2)} = R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(S_{t+2}) \\
 & \vdots \\
 n = \infty & (MC) \quad q_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \gamma^{T-1} R_T
 \end{array}$$

Define the n -step Q-return:

$$q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n}) \quad (43)$$

n -step Sarsa updates $Q(s, a)$ towards the n -step Q-return

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(q_t^{(n)} - Q(S_t, A_t)) \quad (44)$$

By following the same logic as for TD(λ) in (34) and using weight $(1-\lambda)\lambda^{n-1}$

$$q_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} q_t^{(n)} \quad (45)$$

... we get the **Forward-view Sarsa(λ)**

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(q_t^{(\lambda)} - Q(S_t, A_t)) \quad (46)$$

Just like TD(λ), we use *eligibility traces* in an online algorithm for **Backward-View Sarsa(λ)**. In this case, Sarsa(λ) has one eligibility trace for each state-action pair:

$$E_0(s, a) = 0 \quad (47)$$

$$E_t(s, a) = \gamma \lambda E_{t-1}(s, a) + \underbrace{\mathbf{1}(S_t = s, A_t)}_{\text{increment by 1 each time the state-action pair is visited}} \quad (48)$$

$Q(s, a)$ is updated for every state s and a in proportion to TD-error δ_t and eligibility trace $E_t(s, a)$:

$$\delta_t = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t) \quad (49)$$

$$Q^{new}(s, a) \leftarrow Q(s, a) + \underbrace{\alpha}_{\text{learning rate}} \cdot \underbrace{\delta_t}_{\text{TD-error}} \cdot \underbrace{E_t}_{\text{eligibility trace}} \quad (50)$$

The algorithm looks like this:

```

Initialize  $Q(s, a)$  arbitrarily, for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
Repeat (for each episode):
     $E(s, a) = 0$ , for all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ 
    Initialize  $S, A$ 
    Repeat (for each step of episode):
        Take action  $A$ , observe  $R, S'$ 
        Choose  $A'$  from  $S'$  using policy derived from  $Q$  (e.g.,  $\varepsilon$ -greedy)
         $\delta \leftarrow R + \gamma Q(S', A') - Q(S, A)$ 
         $E(S, A) \leftarrow E(S, A) + \delta$ 
        For all  $s \in \mathcal{S}, a \in \mathcal{A}(s)$ :
             $Q(s, a) \leftarrow Q(s, a) + \alpha \delta E(s, a)$ 
             $E(s, a) \leftarrow \gamma \lambda E(s, a)$ 
         $S \leftarrow S'; A \leftarrow A'$ 
    until  $S$  is terminal

```

Figure 3: Sarsa(λ)-Algorithm

5.4 Off-Policy Learning

Evaluate *target policy* $\pi(a|s)$ to compute $v_\pi(s)$ or $q_\pi(s, a)$ while following *behaviour policy* $\mu(a|s)$: $\{S_1, A_1, R_2, \dots, S_T\} \sim \mu$

Why is this important?

- Learn from observing humans or other agents
- Re-use experience gained from old policies $\pi_1, \pi_2, \dots, \pi_{t-1}$
- Learn about *optimal* policy while follow *exploratory* policy
- Learn about *multiple* policies while following *one* policy

Importance Sampling

$$\begin{aligned}
 \mathbb{E}_{X \sim P}[f(X)] &= \sum P(X) f(x) \\
 &= \sum Q(X) \frac{P(X)}{Q(X)} f(X) \\
 &= \mathbb{E}_{X \sim Q} \left[\frac{P(X)}{Q(X)} f(X) \right]
 \end{aligned}$$

Importance Sampling for Off-Policy Monte Carlo:

- Use return from μ to evaluate π
- Weight return G_t according to similarities between policies

- Multiply importance sampling corrections along whole episode:

$$G_t^{\pi/\mu} = \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} \frac{\pi(A_{t+1}|S_{t+1})}{\mu(A_{t+1}|S_{t+1})} \dots \frac{\pi(A_T|S_T)}{\mu(A_T|S_T)} G_t \quad (51)$$

- Update value towards *corrected* return

$$V(S_t) \leftarrow V(S_t) + \alpha \left(G_t^{\pi/\mu} - V(S_t) \right) \quad (52)$$

- Cannot use if μ is zero when π is non-zero
- Importance sampling can dramatically increase variance
→ **in reality unusable because of multiplied probabilities**

For importance sampling, **bootstrapping becomes necessary** to gain useful information to avoid the multiplied probabilities.

Importance Sampling for Off-Policy TD

- Use TD targets generate from μ to evaluate π
- Weight TD target $R + \gamma V(S')$ by importance sampling
- Only need a single importance sampling correction (because of bootstrapping):

$$V(S_t) \leftarrow V(S_t) + \alpha \left(\frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} ((R_{t+1} + \gamma V(S_{t+1})) - V(S_t)) \right) \quad (53)$$

- Much lower variance than Monte Carlo importance sampling
- Policies only need to be similar over a single time step

Q-Learning

For off-policy learning of Q-values $Q(s, a)$. For this **no importance sampling** is required:

- Next action is chosen using behaviour policy: $A_{t+1} \sim \mu(\cdot|S_t)$
- An alternative successor is considered: $A' \sim \pi(\cdot|S_t)$
- Update $Q(S_t, A_t)$ towards value of alternative action:

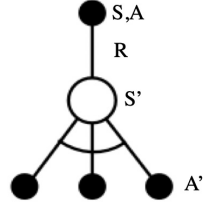
$$Q(S_t, A_t) \rightarrow Q(S_t, A_t) + \alpha ((R_{t+1} + \gamma Q(S_{t+1}, A') - Q(S_t, A_t)) \quad (54)$$

We now allow both policies to improve:

- The target policy π is greedy: $\pi(S_{t+1}) = \operatorname{argmax}_{a'} Q(S_{t+1}, a')$
- The behaviour policy is e.g. ϵ -greedy

- The Q-learning target then simplifies:

$$\begin{aligned} R_{t+1} + \gamma Q(S_{t+1}, A') &= R_{t+1} + \gamma Q(S_{t+1}, \underset{a'}{\operatorname{argmax}} Q(S_{t+1}, a')) \\ &= R_{t+1} + \max_{a'} \gamma Q(S_{t+1}, a') \end{aligned}$$



$$Q(S, A) \leftarrow Q(S, A) + \alpha (R + \gamma \max_{a'} Q(S', a') - Q(S, A)) \quad (55)$$

Initialize $Q(s, a), \forall s \in \mathbb{S}, a \in \mathcal{A}(s)$, arbitrarily, and $Q(\text{terminal-state}, \cdot) = 0$
Repeat (for each episode):
 Initialize S
 Repeat (for each step of episode):
 Choose A from S using policy derived from Q (e.g., ϵ -greedy)
 Take action A , observe R, S'
 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$
 $S \leftarrow S'$;
 until S is terminal

Figure 4: Q-Learning Algorithm