

Temporal Difference Learning

Reinforcement Learning

September 29, 2020

MSU

Reminder: Value functions

State value function (V-function): $V(p, s, q) = E_T \sum_{t \neq 0}^{\infty} \gamma^t r_t \mid s_0 = s$

State-action value function (Q-function): $Q(p, s, a, q) = E_T \sum_{t \neq 0}^{\infty} \gamma^t r_t \mid s_0 = s; a_0 = a$

Reminder: Value functions

State value function (V-function): $V(s) = E_T \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s$

State-action value function (Q-function): $Q(s; a) = E_T \sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s; a_0 = a$

Optimal V-function: $V^*(s) = \max_a V(s; a)$

Optimal Q-function: $Q^*(s; a) = \max_a Q(s; a)$

Reminder: Value functions

State value function (V-function): $V(p; s; q) = E_T \sum_{t \neq 0}^{\infty} r_t \mid s_0 = s$

State-action value function (Q-function): $Q(p; s; a; q) = E_T \sum_{t \neq 0}^{\infty} r_t \mid s_0 = s; a_0 = a$

Optimal V-function: $V^*(p; s; q) = \max V(p; s; q)$

Optimal Q-function: $Q^*(p; s; a; q) = \max Q(p; s; a; q)$

Bellman equations: $@ ; s; a:$

$$Q(p; s; a; q) = r(p; s; a) + E_{s^1} [p(p; s^1; a; q) E_{a^1} [p(a^1; s^1; q) Q(p; s^1; a^1; q)]]$$

$$Q(p; s; a; q) = r(p; s; a) + E_{s^1} [p(p; s^1; a; q) \max_{a^1} Q(p; s^1; a^1; q)]$$

Reminder: Dynamic Programming

Setup: $|S| \neq \infty; |A| \neq \infty$
Assumption: model is known

Reminder: Dynamic Programming

Setup: $|S| \leq 8; |A| \leq 8$

Assumption: model is known

Policy Iteration

Policy Evaluation:

compute V^k for current policy π_k ;

Policy Improvement:

$\pi_{k+1}(s) \in \underset{a}{\operatorname{argmax}} Q^k(s; a)$;

repeat;

Reminder: Dynamic Programming

Setup: $|S| \neq \infty; |A| \neq \infty$
Assumption: model is known

Policy Iteration

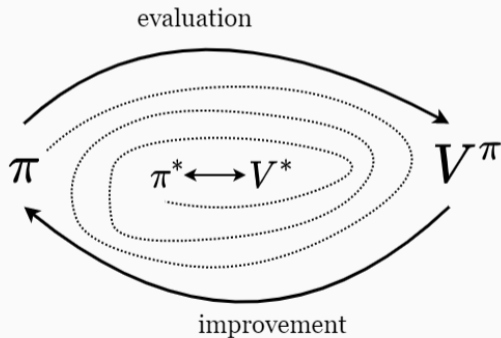
Policy Evaluation:

compute V^k for current policy π^k ;

Policy Improvement:

$\pi^{k+1}(s) \in \arg\max_a Q^k(s; a)$;

repeat;



Reminder: Dynamic Programming

Setup: $|S| \neq \infty; |A| \neq \infty$

Assumption: model is known

Value Iteration

÷ solve Bellman optimality equation
(e. g. for V):

$$V_k(p, s, q) \stackrel{\Delta}{=} \max_a \{ r(p, s; a, q) + \gamma E_{s'} V_k(p, s', q) \}$$

$(p, s, q) \stackrel{\Delta}{=} \operatorname{argmax}_a Q(p, s; a, q)$ is optimal;

Reminder: Dynamic Programming

Setup: $|S| \neq \infty; |A| \neq \infty$

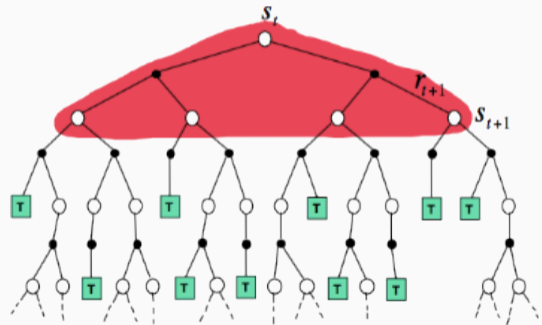
Assumption: model is known

Value Iteration

÷ solve Bellman optimality equation
(e. g. for V):

$$V_k = \max_a \mathbb{E}_{s' \sim p(s'|s,a)} [r + \gamma V_k(s')] - \text{Bellman optimality equation}$$

$a^* = \text{argmax}_a Q(s,a)$ is optimal;



Trial and Error learning

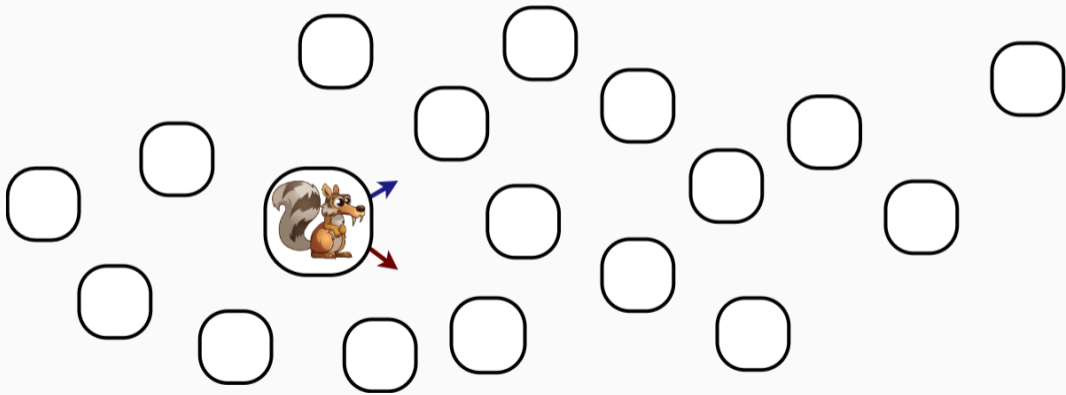
Setup: $|S| = 8; |A| = 8$

~~Assumption: model is known~~

Trial and Error learning

Setup: $|S| = 8; |A| = 8$

Assumption: ~~model is known~~



Idea 1: Model-based RL



Learn model through experience
Use VI / PI with learned model to find optimal policy

Idea 1: Model-based RL

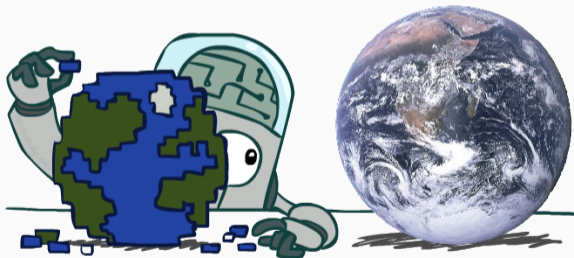


Learn model through experience
Use VI / PI with learned model to find optimal policy

What is easier to create:

a) universe

b) brain



Idea 1: Model-based RL



Learn model through experience
Use VI / PI with learned model to find optimal policy

What is easier to create:

a) **universe**

$S \quad A \quad \tilde{N} \quad PpSq$

× (almost) supervised learning

b) **brain**



Idea 1: Model-based RL



Learn model through experience
Use VI / PI with learned model to find optimal policy

What is easier to create:

a) **universe**

$S \quad A \quad \tilde{N} \quad PpSq$

× (almost) supervised learning

b) **brain**

× $S \quad \tilde{N} \quad A$

dataset is not provided :(



Idea 1: Model-based RL



Learn model through experience
Use VI / PI with learned model to find optimal policy

What is easier to create:

a) **universe**

$S \ A \ \tilde{N} \ PpSq$

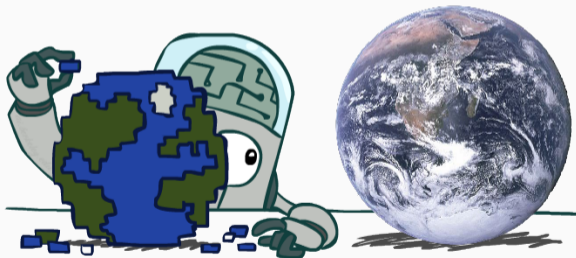
× (almost) supervised learning

b) **brain**

× $S \ \tilde{N} \ A$

dataset is not provided :(

?!?



Idea 2: Model-free RL



Do not try to learn the model
Map states and actions to reward directly

Idea 2: Model-free RL



Do not try to learn the model
Map states and actions to reward directly

Where do we use the model?

Policy Iteration

evaluate policy

e. g. by solving Bellman equation:

$$V^k(p, s, q) = E_a [r(p, s; a, q) + \gamma E_{s'} V^k(p, s', q)]$$

$$k+1(p, s, q) \in \operatorname{argmax}_a Q^k(p, s; a, q);$$

repeat;

Value Iteration

÷ solve Bellman optimality equation
(*e. g. for V*):

$$V_{k+1}(p, s, q) \in \max_a [r(p, s; a, q) + \gamma E_{s'} V_k(p, s', q)]$$

$$p, s, q \in \operatorname{argmax}_a Q(p, s; a, q) \text{ is optimal;}$$

Idea 2: Model-free RL



Do not try to learn the model
Map states and actions to reward directly

Where do we use the model?

Policy Iteration

evaluate policy

e. g. by solving Bellman equation:

$$V^k(p, s, q) = E_a [r(p, s; a, q) + \gamma E_{s'} V^k(p, s', q)]$$

$$k+1(p, s, q) \ni \operatorname{argmax}_a Q^k(p, s; a, q);$$

repeat;

Value Iteration

÷ solve Bellman optimality equation
(*e. g. for V*):

$$V_{k+1}(p, s, q) \ni \max_a [r(p, s; a, q) + \gamma E_{s'} V_k(p, s', q)]$$

$$p, s, q \ni \operatorname{argmax}_a Q^*(p, s; a, q) \text{ is optimal;}$$

Switching to Q-function

Policy Iteration

evaluate policy:

$$Q^k(p, s; a, q) = r(p, s; a, q) + \mathbb{E}_{s'} \mathbb{E}_{a'} Q^k(p, s'; a', q)$$

$$k \leftarrow k + 1; p, s, q \leftarrow \arg\max_a Q^k(p, s; a, q);$$

repeat;

Value Iteration

÷ solve Bellman optimality equation:

$$Q_{k+1}(p, s; a, q) = r(p, s; a, q) + \mathbb{E}_{s'} \max_{a'} Q_k(p, s'; a', q)$$

$$p, s, q \leftarrow \arg\max_a Q(p, s; a, q) \text{ is optimal};$$

Switching to Q-function

Policy Iteration

evaluate policy:

$$Q^k(p, s; a) = r(p, s; a) + \gamma E_{s'} E_{a'} Q^k(p, s'; a')$$

$$k+1(p, s) \in \underset{a}{\operatorname{argmax}} Q^k(p, s; a);$$

repeat;

Value Iteration

÷ solve Bellman optimality equation:

$$Q_{k+1}(p, s; a) \in \underset{a'}{E_{s'}} \max Q_k(p, s'; a')$$

$$p, s \in \underset{a}{\operatorname{argmax}} Q(p, s; a) \text{ is optimal};$$

$$x = E_{s'} f(p, s'; x); \quad x = ?$$

Switching to Q-function

Policy Iteration

evaluate policy:

$$Q^k(p, s; a, q) = r(p, s; a, q) + \mathbb{E}_{s'} \mathbb{E}_{a'} Q^k(p, s'; a', q)$$

$$k+1(p, s, q) \in \underset{a}{\operatorname{argmax}} Q^k(p, s; a, q);$$

repeat;

Value Iteration

÷ solve Bellman optimality equation:

$$Q_{k+1}(p, s; a, q) \in \mathbb{E}_{s'} \max_{a'} Q_k(p, s'; a', q)$$

$$p, s, q \in \underset{a}{\operatorname{argmax}} Q(p, s; a, q) \text{ is optimal};$$

$$x = \mathbb{E}_{s'} f(p, s'; x, q); \quad x = ?$$

Can we derive *some* model-free PI / VI?

Monte-Carlo Backup

How to solve this: $x = E_{s^1} f_{ps^1} q$

How to solve this: $x = E_{s^1} f(p, s^1, q)$

Monte-Carlo estimation:

$$Q(p, s; a, q) = \frac{1}{N} \sum_{i=0}^N R(p, T_i, q)$$

where $T_i \mid s_0, s; a_0, a$

Monte-Carlo Backup

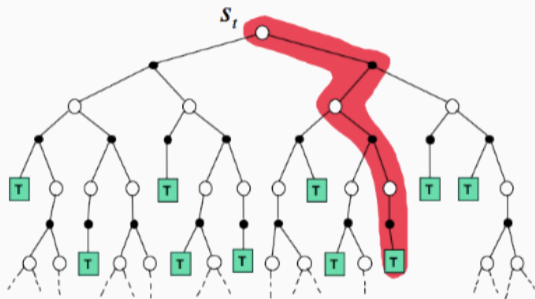
How to solve this: $x = E_{S^1} f(p, s^1, q)$

Monte-Carlo estimation:

$$Q(p, s; a, q) = \frac{1}{N} \sum_{i=0}^N R(p, T_i, q)$$

where $T_i \mid s_0 = s; a_0 = a$

Can we have i. i. d. samples?



Monte Carlo Algorithm

Monte Carlo Algorithm

Initialize $Q(s; a)$ and policy $\pi(s)$ arbitrarily.

for $k = 0; 1; 2; \dots$:

 play several games using π ;

 estimate $Q(s; a)$ using Monte Carlo;

$\pi(s) \leftarrow \underset{a}{\operatorname{argmax}} Q(s; a)$

! requires infinite samples for each pair $s; a$;

Monte Carlo Algorithm

Monte Carlo Algorithm

Initialize $Q(s; a)$ and policy $\pi(s)$ arbitrarily.

for $k = 0; 1; 2; \dots$:

 play several games using π ;

 estimate $Q(s; a)$ using Monte Carlo;

$\pi(s) \leftarrow \underset{a}{\operatorname{argmax}} Q(s; a)$

- ! requires infinite samples for each pair $s; a$;
- still needs full episodes;
- still do not utilize MDP structure;
- wastes information about state connections;
- high variance of collected samples;
- no clear theoretical guarantees;

Monte Carlo Algorithm

Monte Carlo Algorithm

Initialize $Q(s; a)$ and policy $\pi(s)$ arbitrarily.

for $k = 0; 1; 2; \dots$:

 play several games using π ;

update $Q(s; a)$ with new samples;

$\pi(s) \leftarrow \underset{a}{\operatorname{argmax}} Q(s; a)$

- ! requires infinite samples for each pair $s; a$;
- still needs full episodes;
- still do not utilize MDP structure;
- wastes information about state connections;
- high variance of collected samples;
- no clear theoretical guarantees;



Reuse samples from previous policy
(for example, with smaller weight)

Closer look at Monte Carlo

Suppose we want to compute online mean m of i. i. d. samples $x_1; x_2; \dots$. On step k :

$$m_k = \frac{1}{k} \sum_{i=1}^k x_i$$

Closer look at Monte Carlo

Suppose we want to compute online mean m of i. i. d. samples $x_1; x_2; \dots$. On step k :

$$m_k = \frac{1}{k} \sum_{i=1}^k x_i = \frac{k-1}{k} m_{k-1} + \frac{1}{k} x_k$$

Closer look at Monte Carlo

Suppose we want to compute online mean m of i. i. d. samples $x_1; x_2; \dots$. On step k :

$$m_k = \frac{1}{k} \sum_{i=1}^k x_i = \frac{k-1}{k} m_{k-1} + \frac{1}{k} x_k$$

Closer look at Monte Carlo

Suppose we want to compute online mean m of i. i. d. samples $x_1; x_2; \dots$. On step k :

$$m_k = \frac{1}{k} \sum_{i=1}^k x_i = \frac{k-1}{k} m_{k-1} + \frac{1}{k} x_k$$

! $k : \frac{1}{k}$

exponential smoothing

Closer look at Monte Carlo

Suppose we want to compute online mean m of i. i. d. samples $x_1; x_2; \dots$. On step k :

$$m_k = \frac{1}{k} \sum_{i=1}^k x_i$$

!

$$k : \frac{1}{k} \left(\frac{k-1}{k} m_{k-1} + \frac{1}{k} x_k \right)$$

exponential smoothing

$$k : \frac{1}{k} \left(m_{k-1} + \frac{1}{k} (x_k - m_{k-1}) \right)$$

stoch. gradient descent

Closer look at Monte Carlo

Suppose we want to compute online mean m of i. i. d. samples $x_1; x_2; \dots$. On step k :

$$m_k = \frac{1}{k} \sum_{i=1}^k x_i$$

!

$$m_k = \frac{k-1}{k} m_{k-1} + \frac{1}{k} x_k$$

exponential smoothing

$$m_k = m_{k-1} - \eta \nabla \ell(m_{k-1}; x_k)$$

stoch. gradient descent

Closer look at Monte Carlo

Suppose we want to compute online mean m of i. i. d. samples $x_1; x_2; \dots$. On step k :

$$m_k = \frac{1}{k} \sum_{i=1}^k x_i = \frac{k-1}{k} m_{k-1} + \frac{1}{k} x_k$$

exponential smoothing
stoch. gradient descent

Convergence:

m_k $\xrightarrow[k \rightarrow \infty]{\text{w.p. 1}}$ m Ex if learning rate satisfies Robbins–Monro conditions:

$$\sum_{k=1}^{\infty} \eta_k = \infty; \quad \sum_{k=1}^{\infty} \eta_k^2 < \infty$$

Solving equations

Equation	Method	Update
$x \quad f(x)$	Point Iteration	$x_{k+1} : f(x_k)$
$x \quad E_{s^{-1}} f(p s^{-1} q)$	Exponential smoothing	$x_{k+1} : x_k + \frac{1}{s} (f(p s^{-1} q) - x_k)$
$x \quad E_{s^{-1}} f(p s^{-1} ; x) q$	Stochastic approximation	

Solving equations

Equation	Method	Update
$x = f(x)$	Point Iteration	$x_{k+1} = f(x_k)$
$x = E_{s^1} f(x)$	Exponential smoothing	$x_{k+1} = x_k + \frac{1}{s^1} (f(x_k) - x_k)$
$x = E_{s^1} f(x; x)$	Stochastic approximation	$x_{k+1} = x_k + \frac{1}{s^1} (f(x_k; x_k) - x_k)$

Temporal Difference: intuition

View 1: solving Bellman expectation equation using Robbins-Monro scheme:

$$Q(p,s; a,q) = E_{s^1} [r(p,s; a) + \gamma \sum_{q'} P(p,s^1; a, q') Q(p,s^1; a, q)]$$

Temporal Difference: intuition

View 1: solving Bellman expectation equation using Robbins-Monro scheme:

$$Q(s; a) = E_{s^1} [r(s; a) + \gamma \sum_{x^1} P(x^1 | s^1, a^1) Q(s^1; a^1)]$$

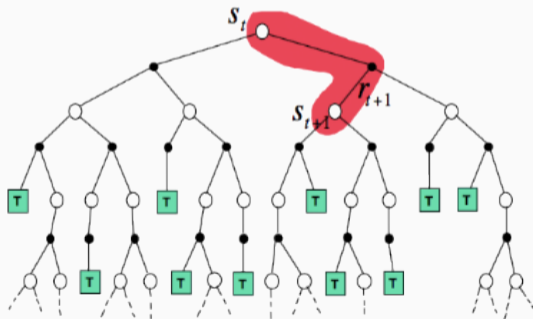
$$Q(s; a) = E_{s^1} E_{a^1} [r(s; a) + \gamma \sum_{x^1} P(x^1 | s^1, a^1) Q(s^1; a^1)]$$

Temporal Difference: intuition

View 1: solving Bellman expectation equation using Robbins-Monro scheme:

$$Q(p; a) = E_{s^1} [r(p; a) + \gamma Q(p; a^1 | s)] \quad E_{s^1} E_{a^1} [r(p; a) + \gamma Q(p; a^1 | s)]$$

$f(p; s^1; x)$
 $f(p; s^1; a^1; x)$



Temporal Difference: intuition

View 1: solving Bellman expectation equation using Robbins-Monro scheme:

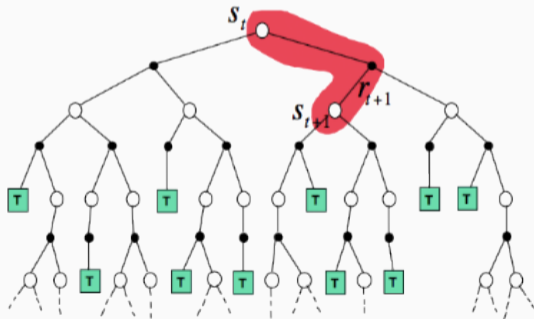
$$Q(s^1; a^1) = E_{s^1} [r + \gamma E_{a^1} [Q(s^1; a^1)]]$$

$$Q(s^1; a^1) = E_{s^1} [r + \gamma Q(s^1; a^1)]$$

View 2: one-step bootstrapping:
approximating future rewards using current approximation.

$$r = r^1 + \gamma r^2 + \gamma^2 r^3 + \dots$$

$$Q_k(s^1; a^1)$$



Temporal Difference: intuition

View 1: solving Bellman expectation equation using Robbins-Monro scheme:

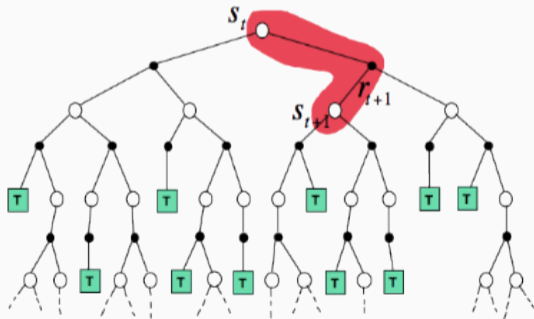
$$Q(s^1; a^1) = E_{s^1} [r + \gamma Q(s^1; a^1) | s^1, a^1] = E_{s^1} [r + \gamma E_{a^1} [r + \gamma Q(s^1; a^1) | s^1, a^1]]$$

View 2: one-step bootstrapping:
approximating future rewards using current approximation.

$$r = r^1 + \gamma r^2 + \gamma^2 r^3 + \dots$$

$$Q_k(s^1; a^1)$$

our own predictions used as targets!



Temporal Difference vs Monte Carlo

$$Q_{k-1}(s; a) \leftarrow Q_k(s; a) + \alpha (y_k - Q_{k-1}(s; a))$$

Which is better?

Temporal Difference
 $y_k = r + \gamma V(s')$

Monte Carlo
 $y_k = r + \gamma V(s')$

Temporal Difference vs Monte Carlo

$$Q_k(s; a) \leftarrow Q_k(s; a) + \alpha (y_k - Q_k(s; a))$$

Which is better?

Temporal Difference

$$y_k : r + \gamma Q_k(s^1; a^1)$$

- × updates after each step
- slow reward propagation

Monte Carlo

$$y_k : r + \gamma (r^1 + \gamma^2 r^2 + \dots)$$

- updates at the end of the episode
- × fast reward propagation

Temporal Difference vs Monte Carlo

$$Q_{k+1}(s; a) \leftarrow Q_k(s; a) + \alpha (y_k - Q_k(s; a))$$

Which is better?

Temporal Difference

$$y_k = r + \gamma V(s_{k+1}) - V(s_k)$$

- × updates after each step
- slow reward propagation
- × low variance
- introduces bias

Monte Carlo

$$y_k = r + \gamma V(s_{k+1}) - V(s_k)$$

- updates at the end of the episode
- × fast reward propagation
- high variance
- × unbiased

Temporal Difference vs Monte Carlo

$$Q_{k+1}(s; a) \leftarrow Q_k(s; a) + \alpha (r_k + \gamma V_k(s') - V_k(s; a))$$

Which is better?

Temporal Difference

$$Q_k(s; a) : r + \gamma V_k(s')$$

- × updates after each step
- slow reward propagation
- × low variance
- introduces bias

Monte Carlo

$$Q_k(s; a) : r + r^1 + \gamma^2 r^2 + \dots$$

- updates at the end of the episode
- × fast reward propagation
- high variance
- × unbiased

Best estimation is somewhere in between (see TD(0)).

Value Iteration and Policy Iteration connection



What if we improve our policy after every step?

$$Q_k(p, s; a) \in \operatorname{argmax}_a Q_k(p, s; a)$$

Value Iteration and Policy Iteration connection



What if we improve our policy after every step?

$$Q_k(p, s; a) \in \operatorname{argmax}_a Q_k(p, s; a)$$

$$Q_{k+1}(p, s; a) \in Q_k(p, s; a) + r + \gamma \max_{a'} Q_k(p, s; a')$$

Value Iteration and Policy Iteration connection



What if we improve our policy after every step?

$$Q_k(p, s; a) \in \operatorname{argmax}_a Q_k(p, s; a)$$

$$Q_{k+1}(p, s; a) \in Q_k(p, s; a) + \gamma \left(\sum_{s'} P_{ss'}(a) Q_k(p, s'; a) - Q_k(p, s; a) \right)$$

Value Iteration and Policy Iteration connection



What if we improve our policy after every step?

$$Q_k(p, s, a) \in \operatorname{argmax}_a Q_k(p, s; a, q)$$

$$Q_k(p, s; a, q) \in Q_k(p, s; a, q) \quad k, r \quad Q_k(p, s^1; a^1, q) \quad Q_k(p, s; a, q)$$

$$Q_k(p, s; a, q) \quad k, r \quad Q_k(p, s^1; k, p, s^1, q, q) \quad Q_k(p, s; a, q)$$

$$Q_k(p, s; a, q) \quad k, r \quad Q_k(p, s^1; \operatorname{argmax}_{a^1} Q_k(p, s^1; a^1, q), q) \quad Q_k(p, s; a, q)$$

Value Iteration and Policy Iteration connection



What if we improve our policy after every step?

$$Q_k(p, s, q) \in \operatorname{argmax}_a Q_k(p, s, a, q)$$

$$Q_{k-1}(p, s, a, q) \in Q_k(p, s, a, q) \quad k, r \quad Q_k(p, s^1, a^1, q) \quad Q_k(p, s, a, q)$$

$$Q_k(p, s, a, q) \quad k, r \quad Q_k(p, s^1, a^1, q) \quad Q_k(p, s, a, q)$$

$$Q_k(p, s, a, q) \quad k, r \quad \operatorname{argmax}_{a^1} Q_k(p, s^1, a^1, q) \quad Q_k(p, s, a, q)$$

$$Q_k(p, s, a, q) \quad k, r \quad \max_{a^1} Q_k(p, s^1, a^1, q) \quad Q_k(p, s, a, q)$$

Value Iteration and Policy Iteration connection



What if we improve our policy after every step?

$$Q_k(p, s; a) \stackrel{\Delta}{=} \arg\max_a Q_k(p, s; a)$$

$$Q_{k+1}(p, s; a) \stackrel{\Delta}{=} Q_k(p, s; a) + \gamma \sum_{s'} P(s' | p, s, a) [Q_k(p, s'; a) - Q_k(p, s; a)]$$

$$Q_{k+1}(p, s; a) \stackrel{\Delta}{=} Q_k(p, s; a) + \gamma \sum_{s'} P(s' | p, s, a) [Q_k(p, s'; a) - Q_k(p, s; a)]$$

$$Q_{k+1}(p, s; a) \stackrel{\Delta}{=} Q_k(p, s; a) + \gamma \sum_{s'} P(s' | p, s, a) [Q_k(p, s'; a) - Q_k(p, s; a)]$$

$$Q_{k+1}(p, s; a) \stackrel{\Delta}{=} Q_k(p, s; a) + \gamma \sum_{s'} P(s' | p, s, a) [Q_k(p, s'; a) - Q_k(p, s; a)]$$

Value Iteration and Policy Iteration connection



What if we improve our policy after every step?

$$Q_k(p, s; a) \in \operatorname{argmax}_a Q_k(p, s; a)$$

$$Q_{k+1}(p, s; a) \in Q_k(p, s; a) + \gamma \left(\sum_{q'} P_{k+1}(a, s, q) Q_k(p, s; a) - Q_k(p, s; a) \right)$$

$$Q_k(p, s; a) + \gamma \left(\sum_{q'} P_k(a, s, q) Q_k(p, s; a) - Q_k(p, s; a) \right)$$

$$Q_k(p, s; a) + \gamma \left(\sum_{q'} P_k(a, s, q) \max_{a'} Q_k(p, s; a') - Q_k(p, s; a) \right)$$

$$Q_k(p, s; a) + \gamma \left(\max_{a'} \sum_{q'} P_k(a, s, q) Q_k(p, s; a') - Q_k(p, s; a) \right)$$

View 1: Policy Iteration with policy improvement after *each* policy evaluation update;

View 2: solving Bellman optimality equation (model-free Value Iteration);

Convergence of temporal difference

Convergence of Q-learning updates

Let $Q_0(p, s; a, q)$ be initialized arbitrary, and for every $s; a$ the following update is used:

$$Q_{k+1}(p, s; a, q) \leftarrow Q_k(p, s; a, q) + \alpha (r(p, s; a, q) + \max_{a'} Q_k(p, s; a', q) - Q_k(p, s; a, q))$$

Convergence of temporal difference

Convergence of Q-learning updates

Let $Q_0(p,s;a; q)$ be initialized arbitrary, and for every $s; a$ the following update is used:

$$Q_{k+1}(p,s;a; q) \leftarrow Q_k(p,s;a; q) + \alpha (r(p,s;a; q) + \max_{a'} Q_k(p,s;a'; q) - Q_k(p,s;a; q))$$

Then, if:

$$s_{k;s;a}^1 \rightarrow pps^1 | s; a; q;$$

Convergence of temporal difference

Convergence of Q-learning updates

Let $Q_0(p, s; a, q)$ be initialized arbitrary, and for every $s; a$ the following update is used:

$$Q_{k+1}(p, s; a, q) \leftarrow Q_k(p, s; a, q) + \alpha_k(p, s; a, q) [r(p, s; a, q) + \max_{a'} Q_k(p, s; a', q) - Q_k(p, s; a, q)]$$

Then, if:

$$\sum_{k=0}^{\infty} \alpha_k(p, s; a, q) = \infty$$

with probability 1 for every $s; a$ learning rate $\alpha_k(p, s; a, q) \in [0, 1]$ satisfies Robbins-Monro conditions:

$$\sum_{k=0}^{\infty} \alpha_k(p, s; a, q) = \infty \quad \sum_{k=0}^{\infty} \alpha_k^2(p, s; a, q) < \infty$$

Convergence of temporal difference

Convergence of Q-learning updates

Let $Q_0(p, s; a)$ be initialized arbitrary, and for every $s; a$ the following update is used:

$$Q_{k+1}(p, s; a) \leftarrow Q_k(p, s; a) + \alpha_k(p, s; a) (r(p, s; a) + \max_{a'} Q_k(p, s; a') - Q_k(p, s; a))$$

Then, if:

$$s_{k+1}^1 \sim p(p, s^1 | s; a);$$

with probability 1 for every $s; a$ learning rate $\alpha_k(p, s; a) \in (0, 1]$ satisfies Robbins-Monro conditions:

$$\sum_{k=0}^{\infty} \alpha_k(p, s; a) = \infty \quad \text{and} \quad \sum_{k=0}^{\infty} \alpha_k^2(p, s; a) < \infty$$

then $Q_k(p, s; a) \xrightarrow[k \rightarrow \infty]{\text{w.p. 1}} Q^*(p, s; a)$

Analogy 1

Global optimization

Exploration vs Exploitation

Analogy 1

Global optimization

Exploration:
random search

Exploration vs Exploitation

Analogy 1

Global optimization

Exploration:
random search

Exploitation:
local optimization



Exploration vs Exploitation

Analogy 1

Global optimization

Exploration:
random search

Exploitation:
local optimization

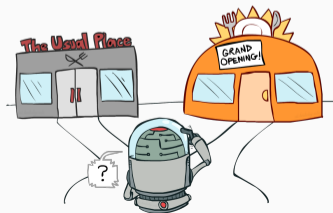


Analogy 2

Multi-armed bandits
(regret minimization)

Exploration:
try something new

Exploitation:
try your favourite



Exploration vs Exploitation

Analogy 1

Global optimization

Exploration:
random search

Exploitation:
local optimization

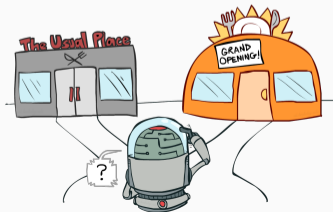


Analogy 2

Multi-armed bandits
(regret minimization)

Exploration:
try something new

Exploitation:
try your favourite



Analogy 3

General MDP

Exploration:
explore the environment

Exploitation:
solve the task



"-greedy exploration



Do something random *sometimes*

Let's call policy "-greedy with respect to some Q-function, if it behaves:
 randomly with probability " ϵ "
 greedy with probability $1 - \epsilon$ "

"-greedy exploration



Do something random *sometimes*

Let's call policy π -greedy with respect to some Q-function, if it behaves:

randomly with probability ϵ
greedy with probability $1 - \epsilon$

$$p(a|s) = \begin{cases} \epsilon & \text{when } a = \underset{a}{\operatorname{argmax}} Q(s; a) \\ \frac{1 - \epsilon}{|A|} & \text{otherwise} \end{cases}$$

"-greedy exploration



Do something random *sometimes*

Let's call policy ϵ -greedy with respect to some Q-function, if it behaves:

randomly with probability ϵ
 greedy with probability $1 - \epsilon$

$$p(a | s) = \begin{cases} \epsilon & \text{when } a = \underset{a}{\operatorname{argmax}} Q(s; a) \\ \frac{1 - \epsilon}{|A|} & \text{otherwise} \end{cases}$$

Decaying exploration:

$$\epsilon_k \sim 0; \quad \epsilon_k \geq 0$$

Persistent exploration:

$$\epsilon_k = \text{const} \cdot k^{-\alpha}$$

Q-learning (online)

Initialize Q ψ ; q arbitrarily;

observe s_0 ;

for $k = 0; 1; 2; \dots$:

 take action $a_k = \text{argmax}_{a \in \mathcal{A}} Q(s_k, a)$;

 observe $r_k; s_{k+1}$;

Q-learning (online)

Initialize Q $p_s; a_q$ arbitrarily;

observe s_0 ;

for $k = 0; 1; 2; \dots$:

take action a_k ϵ -greedy $p_{s_k}; a_q$;

observe $r_k; s_{k+1}$;

$y = r_k + \max_{a_{k+1}} Q(p_{s_{k+1}}; a_{k+1})$;

Q-learning (online)

Initialize Q $p_s; a_q$ arbitrarily;

observe s_0 ;

for $k = 0; 1; 2; \dots$:

take action a_k ϵ -greedy $p_{s_k; a_q}$;

observe $r_k; s_{k+1}$;

$y = r_k + \max_{a_{k+1}} Q(p_{s_{k+1}}; a_{k+1})$;

$Q(p_{s_k; a_k}) \leftarrow \alpha y + (1 - \alpha) Q(p_{s_k; a_k})$

Q-learning

Q-learning (online)

Initialize Q $s; a$ arbitrarily;

observe s_0 ;

for $k = 0; 1; 2; \dots$:

take action $a_k = \text{argmax}_{a \in \mathcal{A}(s_k)} Q(s_k, a)$;

observe $r_k; s_{k+1}$;

$y = r_k + \gamma \max_{a \in \mathcal{A}(s_{k+1})} Q(s_{k+1}, a)$;

$Q(s_k, a_k) \leftarrow Q(s_k, a_k) + \alpha (y - Q(s_k, a_k))$

Q-learning (with replay buffer)

Initialize Q $s; a$ arbitrarily, $D = \mathcal{H}$;

Q-learning

Q-learning (online)

Initialize $Q(p_s; a_q)$ arbitrarily;

observe s_0 ;

for $k = 0; 1; 2; \dots$:

take action $a_k = \text{argmax}_{a \in \mathcal{A}} Q(p_{s_k}; a)$;

observe $r_k; s_{k+1}$;

$y = r_k + \gamma \max_{a \in \mathcal{A}} Q(p_{s_{k+1}}; a)$;

$Q(p_{s_k}; a_k) \leftarrow Q(p_{s_k}; a_k) + \alpha (y - Q(p_{s_k}; a_k))$

Q-learning (with replay buffer)

Initialize $Q(p_s; a_q)$ arbitrarily, $D = \emptyset$;

observe s_0 ;

for $k = 0; 1; 2; \dots$:

take action $a_k = \text{argmax}_{a \in \mathcal{A}} Q(p_{s_k}; a)$;

observe $r_k; s_{k+1}$;

store $(p_{s_k}; a_k; r_k; s_{k+1})$ in D ;

Q-learning

Q-learning (online)

Initialize $Q(s; a)$ arbitrarily;

observe s_0 ;

for $k = 0; 1; 2; \dots$:

take action $a_k = \text{argmax}_a Q(s_k; a)$;

observe $r_k; s_{k+1}$;

$y = r_k + \gamma \max_{a_{k+1}} Q(s_{k+1}; a_{k+1})$;

$Q(s_k; a_k) \leftarrow \gamma Q(s_k; a_k) + (1 - \gamma)y$

Q-learning (with replay buffer)

Initialize $Q(s; a)$ arbitrarily, $D = \emptyset$;

observe s_0 ;

for $k = 0; 1; 2; \dots$:

take action $a_k = \text{argmax}_a Q(s_k; a)$;

observe $r_k; s_{k+1}$;

store $(s_k; a_k; r_k; s_{k+1})$ in D ;

sample $(s; a; r; s')$ from D ;

$y = r + \gamma \max_{a'} Q(s'; a')$;

$Q(s; a) \leftarrow \gamma Q(s; a) + (1 - \gamma)y$

Q-learning is ϵ -policy

Data source	Interaction policy	Convergence
Online experience	Infinite visitation	Q
Expert data	-	Q
Experience replay	Infinite visitation	Q



Expectation

Expectation

Reality

Expectation

Reality



Consider the family of « ϵ -soft policies » :

$$\pi(a|s) \geq \frac{\epsilon}{|A|}$$

Q-learning vs SARSA

ε-soft Policy Improvement

In the family of « ε-soft policies » ε-greedy is policy improvement:

~ : ε-greedy Q ps; aqq ñ ~ ©

Q-learning vs SARSA

"-soft Policy Improvement

In the family of « "-soft policies» "-greedy is policy improvement:

$$\pi_k : \text{"-greedy} \pi_{k-1} \text{ } Q_{k-1} \text{ } \pi_k \text{ } \tilde{r} \text{ } \pi_k \text{ } \odot$$

$$Q_{k-1}(s,a) \in Q_k(s,a) \quad k \text{ } r(s,a) \quad Q_k(s^1; a^1) \quad Q_k(s,a)$$

Q-learning

a^1 taken from target policy

$$k \text{ } \pi_k : \arg \max_a Q_k(s,a)$$

SARSA

a^1 taken from behavior policy

$$k : \text{"-greedy} \pi_k \text{ } Q_k \text{ } \pi_k \text{ } \odot$$

Q-learning vs SARSA

"-soft Policy Improvement

In the family of « "-soft policies» "-greedy is policy improvement:

$$\pi_k : \text{"-greedy}_{\pi_k} Q_{\pi_k}; a, q \quad \tilde{\pi} \quad \sim \odot$$

$$Q_{\pi_{k-1}}(s; a, q) \in Q_{\pi_k}(s; a, q) \quad k \quad r(s; a, q) \quad Q_{\pi_{k-1}}(s; a^1, q) \quad Q_{\pi_k}(s; a, q)$$

Q-learning

a^1 taken from target policy

$$\pi_k(s, q) : \underset{a}{\operatorname{argmax}} Q_{\pi_k}(s; a, q)$$

SARSA

a^1 taken from behavior policy

$$\pi_k : \text{"-greedy}_{\pi_k} Q_{\pi_k}; a, q$$

SARSA convergence properties

SARSA converges to optimal "-soft policy, satisfying

$$\pi_k : \text{"-greedy}_{\pi_k} Q_{\pi_k}; a, q$$

SARSA

Initialize Q $p_s; a_q$ arbitrarily.

observe s_0 , select a_0 randomly;

for $k = 0; 1; 2; \dots$:

take action a_k , observe $r_k; s_{k+1}$;

sample $a_{k+1} \sim \text{"-greedy } p_{s_{k+1}; a_q}$

$y = r_k + \gamma Q(p_{s_{k+1}; a_{k+1}})$;

$Q(p_{s_k; a_k}) \leftarrow p_1 + \alpha (y - Q(p_{s_k; a_k}))$

SARSA

SARSA

Initialize $Q(p_s; a_q)$ arbitrarily.

observe s_0 , select a_0 randomly;

for $k = 0; 1; 2; \dots$:

take action a_k , observe $r_k; s_{k+1}$;

sample $a_{k+1} \sim \text{greedy}(Q(p_{s_{k+1}}; a_q))$

$y = r_k + \gamma Q(p_{s_{k+1}}; a_{k+1})$;

$Q(p_{s_k}; a_k) \leftarrow \gamma Q(p_{s_k}; a_k) + (1 - \gamma)y$

Expected-SARSA

Initialize $Q(p_s; a_q)$ arbitrarily.

observe s_0 ;

for $k = 0; 1; 2; \dots$:

sample $a_k \sim \text{greedy}(Q(p_{s_k}; a_q))$;

take action a_k , observe $r_k; s_{k+1}$;

$y = r_k + \gamma E_{a_{k+1}} Q(p_{s_{k+1}}; a_{k+1})$;

$Q(p_{s_k}; a_k) \leftarrow \gamma Q(p_{s_k}; a_k) + (1 - \gamma)y$;

SARSA is on-policy

For transition $p(s; a; r; s^1; a^1|q)$:

$$Q_k(s; a) \leftarrow Q_k(s; a) + \alpha (r + \gamma Q_k(s^1; a^1) - Q_k(s; a))$$

Data source	Interaction policy	Convergence
Online experience	Decaying exploration	Q
	Persistent exploration	Q_{ϵ} -greedy
Expert data (a^1 taken from buffer)	-	

SARSA is on-policy

For transition $p(s; a; r; s^1; a^1|q)$:

$$Q_k(s; a) \leftarrow Q_k(s; a) + \alpha [r + \gamma Q_k(s^1; a^1) - Q_k(s; a)]$$

Data source	Interaction policy	Convergence
Online experience	Decaying exploration	Q
	Persistent exploration	Q_{ϵ} -greedy
Expert data (a^1 taken from buffer)	-	Q_{expert}
Experience replay (a^1 taken from buffer)		

SARSA is on-policy

For transition $p(s; a; r; s^1; a^1|q)$:

$$Q_{k+1}(s; a) \leftarrow Q_k(s; a) + \alpha [r + \sum_{s'} p(s'; a|q) Q_k(s'; a) - Q_k(s; a)]$$

Data source	Interaction policy	Convergence
Online experience	Decaying exploration	Q
	Persistent exploration	Q_{ϵ} -greedy
Expert data (a^1 taken from buffer)	-	Q_{expert}
Experience replay (a^1 taken from buffer)	(arbitrary)	divergence

« Fixing » SARSA (not really sarsa now)

For transition $p_s; a; r; s^1q$ generate $a^1 \quad p a^1 \mid s^1q$:

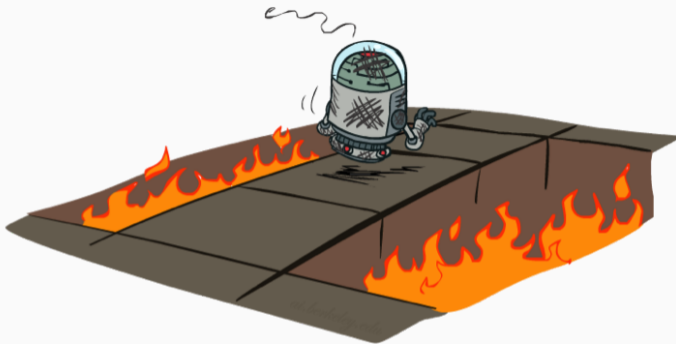
$$Q_{k+1}(p_s; aq) \leftarrow Q_k(p_s; aq) + \alpha (r + Q_k(p_{s^1}; a^1q) - Q_k(p_s; aq))$$

« Fixing » SARSA (not really sarsa now)

For transition $p(s; a; r; s^1|q)$ generate $a^1 \quad p(a^1 | s^1|q)$:

$$Q_k(s^1|ps; aq) \leftarrow Q_k(ps; aq) + \alpha (r + \sum_{q'} p(s^1|ps; aq) Q_k(ps; aq') - Q_k(s^1|ps; aq))$$

Data source	Interaction policy	Convergence
(any, but a^1 generated online)	Decaying exploration	Q
	Persistent exploration	Q_{ϵ} -greedy



Literature

Sutton, Barto — Reinforcement Learning, an Introduction, ch. 5-6;

Watkins, Dayan — Technical Note, Q-learning;

Pictures from Berkeley CS 188 | Introduction to Artificial Intelligence;