

Temporal Difference Learning

Reinforcement Learning

September 21, 2021

MSU

Reminder

Reminder: Value functions

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

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

Reminder: Value functions

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

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

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

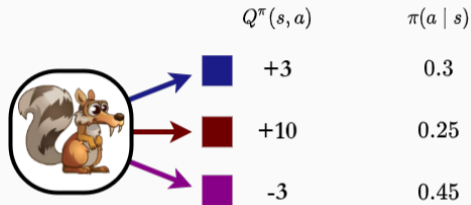
Optimal Q-function: $Q(p; a; q) = \max Q(p; a; q)$

Reminder: Policy Improvement

Policy Improvement:

greedy policy with respect to Q is not worse than π :

$$\pi_{\text{greedy}} = \arg\max_a Q(s; a) \geq \pi^{\pi}$$

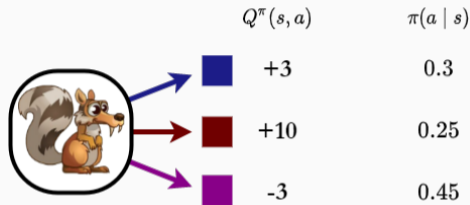


Reminder: Policy Improvement

Policy Improvement:

greedy policy with respect to Q is not worse than π :

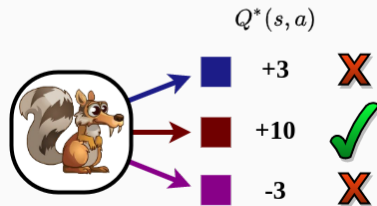
$$\pi_{\text{greedy}}(a|s) = \underset{a}{\operatorname{argmax}} Q(s; a) \geq \pi(a|s)$$



Bellman optimality principle:

greedy policy with respect to optimal Q-function is optimal:

$$\pi_{\text{greedy}}(a|s) = \underset{a}{\operatorname{argmax}} Q(s; a)$$



Bellman (expectation) equations: $@ ; s; a:$

$$Q(p, s; a, q) = r(p, s; a, q) + \gamma \sum_{p^1} p(p^1 | s; a, q) V(p, s^1, q)$$

$$V(p, s, q) = \max_a \sum_{p^1} p(p^1 | s, q) Q(p, s; a, q)$$

Bellman (expectation) equations: $@ ; s ; a :$

$$Q(p, s ; a, q) = r(p, s ; a, q) + \mathbb{E}_{s^1 | p, s ; a, q} V(p, s^1, q)$$
$$V(p, s, q) = \mathbb{E}_a \sum_{p, a | s, q} Q(p, s ; a, q)$$

Bellman (optimality) equations: $@ s ; a :$

$$Q(p, s ; a, q) = r(p, s ; a, q) + \mathbb{E}_{s^1 | p, s ; a, q} V(p, s^1, q)$$
$$V(p, s, q) = \max_a \sum_{p, a | s, q} Q(p, s ; a, q)$$

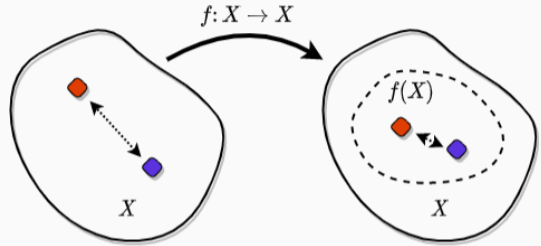
Reminder: Point iteration

Contraction Mapping

$f: X \rightarrow X$ is a contraction mapping if:

$$\forall x_1, x_2: \rho(f(x_1), f(x_2)) \leq \alpha \rho(x_1, x_2)$$

where ρ — some metric, $\alpha < 1$.



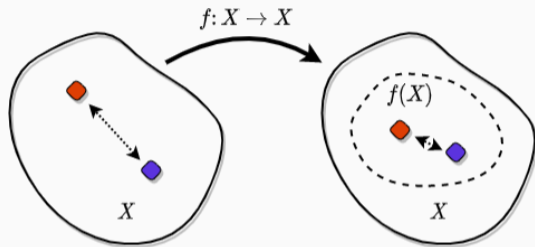
Reminder: Point iteration

Contraction Mapping

$f: X \rightarrow X$ is a contraction mapping if:

$$\forall x_1, x_2: \rho(f(x_1), f(x_2)) \leq \alpha \rho(x_1, x_2)$$

where ρ — some metric, $\alpha < 1$.



If f is a contraction mapping:

- equation $x = f(x)$ has unique solution;
- we can find it using point iteration method

$$x_{k+1} = f(x_k)$$

starting from any x_0 .

Reminder: Policy Iteration

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

Assumption: model is known

Reminder: Policy Iteration

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

Policy Iteration

- **Policy Evaluation:**
compute V^k for current policy π_k ;
- **Policy Improvement:**
 $\pi_{k+1} \in \arg\max_{\pi} Q^k(\pi); a$;
- repeat;

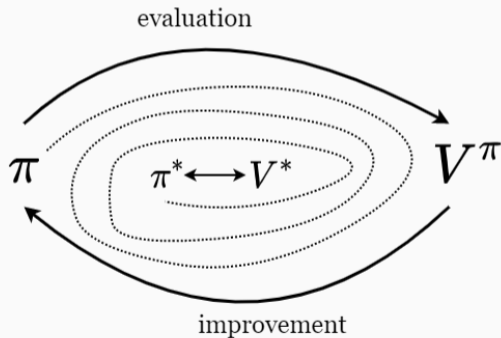
Reminder: Policy Iteration

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

Assumption: model is known

Policy Iteration

- **Policy Evaluation:**
compute V^k for current policy π^k ;
- **Policy Improvement:**
 $\pi^{k+1} = \arg \max_a Q^k(s; a)$;
- repeat;



Reminder: Generalized Policy Iteration

Generalized Policy Iteration

Initialize V ; arbitrary, repeat:

Reminder: Generalized Policy Iteration

Generalized Policy Iteration

Initialize V ; arbitrary, repeat:

- **Policy Evaluation**: move V towards V (*at least a little*)
e.g. perform K iterations of point iteration method for Bellman equation

Reminder: Generalized Policy Iteration

Generalized Policy Iteration

Initialize V ; arbitrary, repeat:

- **Policy Evaluation**: move V towards V (*at least a little*)
e.g. perform K iterations of point iteration method for Bellman equation
- **Policy Improvement**: improve using V (*at least a little*)
e.g. greedy policy improvement

Reminder: Generalized Policy Iteration

Generalized Policy Iteration

Initialize V ; arbitrary, repeat:

- **Policy Evaluation**: move V towards V (*at least a little*)
e.g. perform K iterations of point iteration method for Bellman equation
- **Policy Improvement**: improve using V (*at least a little*)
e.g. greedy policy improvement



This always works!!!

Reminder: Generalized Policy Iteration

Generalized Policy Iteration

Initialize V ; arbitrary, repeat:

- **Policy Evaluation**: move V towards V (*at least a little*)
e.g. perform K iterations of point iteration method for Bellman equation
- **Policy Improvement**: improve using V (*at least a little*)
e.g. greedy policy improvement



This always works!!!

Interesting case: $K = 1$ and greedy policy improvement

Reminder: Value Iteration

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

Assumption: model is known

Value Iteration

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

$$V_k \stackrel{!}{=} \max_a \left[r(s,a) + \gamma \sum_{s'} P(s'|s,a) V_k(s') \right]$$

- $a^* \stackrel{!}{=} \operatorname{argmax}_a Q(s,a)$ is optimal;

Reminder: Value Iteration

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

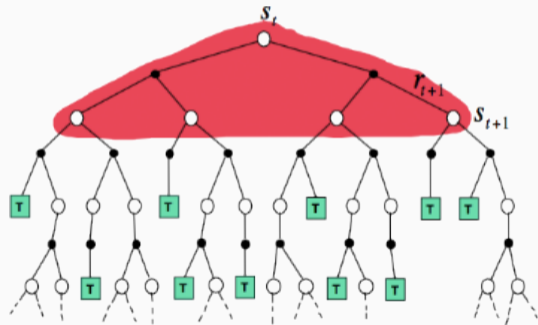
Assumption: model is known

Value Iteration

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

$$V_k = \max_a \mathbb{E}_{S^1} [r + \gamma V_k(S^2) | S^1 = s, a]$$

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



Model-free learning

Trial and Error learning

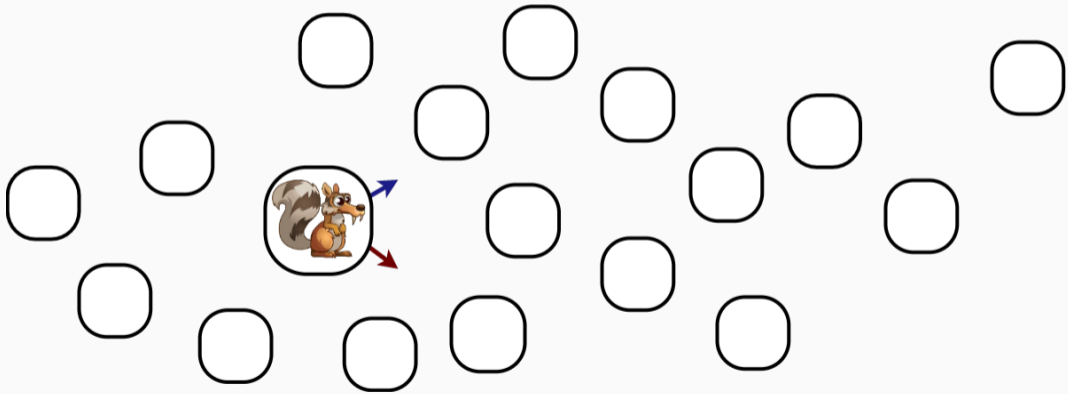
Setup: $|S| = 8; |A| = 8$

~~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

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

× (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**

$S \quad A \quad \tilde{N} \quad P \quad p \quad S \quad q$

× (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**

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

× (almost) supervised learning

b) **brain**

× $S \quad \tilde{N} \quad A$

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:

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

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

Value Iteration

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

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

- $p, s, q \in \arg\max_a Q^*(p, s; a, q)$ 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:

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

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

Value Iteration

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

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

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

Switching to Q-function

Policy Iteration

- evaluate policy:

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

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

Value Iteration

- ÷ solve Bellman optimality equation:

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

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

Switching to Q-function

Policy Iteration

- evaluate policy:

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

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

Value Iteration

- ÷ solve Bellman optimality equation:

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

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

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

Switching to Q-function

Policy Iteration

- evaluate policy:

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

- $p_{k+1} \in \underset{a}{\operatorname{argmax}} Q^k(p; a);$
- repeat;

Value Iteration

- ÷ solve Bellman optimality equation:

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

- $p^* \in \underset{a}{\operatorname{argmax}} Q^*(p; a)$ is optimal;

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

Ok, not clear how to work with this...

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

Monte Carlo Backup

$$Q(p, s; a, q) = E_T \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid s_0 = s, a_0 = a \right]$$

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

Monte Carlo Backup

$$Q(p, s; a, q) = E_T \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid s_0 = s, a_0 = a \right]$$

How to solve this: $x = E_{p, p} \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid p, q \right]$?



Play a lot of games using p to collect data for Monte Carlo estimations

Monte Carlo Backup

$$Q(p, s; a, q) = E_T [R + \gamma V(s') | s_0, s; a_0, a]$$

How to solve this: $x = E_{p, p'} [f(p', q)]$?



Play a lot of games using x to collect data for Monte Carlo estimations

Monte Carlo estimation:

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

where $T_i = s_0, s; a_0, a$

Monte Carlo Algorithm

Monte Carlo Algorithm

Initialize 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 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 does not utilize MDP structure;
- wastes information about state connections;
- high variance of collected samples;

Monte Carlo Algorithm

Monte Carlo Algorithm

Initialize policy $\pi(s)$ and $Q(s; a)$ 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 does not utilize MDP structure;
- wastes information about state connections;
- high variance of collected samples;



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 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

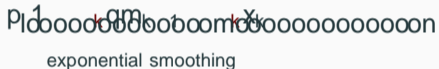
Suppose we want to compute online mean of i. i. d. samples $x_1; x_2; \dots$. On step:

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

Closer look at Monte Carlo

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

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



Closer look at Monte Carlo

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

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

$\frac{1}{k} \sum_{i=1}^k x_i$ exponential smoothing $\frac{k-1}{k} m_{k-1} + \frac{1}{k} x_k$ stoch. gradient descent

Closer look at Monte Carlo

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

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

$\frac{k-1}{k} m_{k-1}$ exponential smoothing
 $\frac{1}{k} x_k$ stoch. gradient descent

Convergence:

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

Ex if learning rate satisfies **Robbins Monro conditions** :

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

Temporal Difference Idea

Solving equations

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

Temporal Difference: intuition 1

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

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

$f(s'; x)$

Temporal Difference: intuition 1

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

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

$f(s^1; x, q)$
 $f(s^1; a^1; x, q)$

Temporal Difference: intuition 1

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

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

To update value for $s; a$ we need only **transition** $s; a; r; s^1; a^1$

$$y = r + \gamma Q(s^1; a^1)$$

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

Requirements for Temporal Difference

For transition $s; a; r; s^1; a^1q$

Bellman target

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

ooooooooooooooooooooomoooooooooooooooooooo

temporal difference

Requirements for Temporal Difference

For transition $s; a; r; s^1; a^1q$

Bellman target

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

looo

temporal difference

$s; a$ considered fixed.

s^1 $p(s^1 | s; a)$ random variable from interaction experience.

Requirements for Temporal Difference

For transition $s; a; r; s^1; a^1q$

Bellman target

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

ooooooooooooooooooooomoooooooooooooooooooo

temporal difference

$s; a$ considered fixed.

s^1 $p(s^1 | s; a)$ random variable from interaction experience.

a^1 $p(a^1 | s^1q)$ random variable from algorithm.

Temporal Difference: intuition 2

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

$$r \quad \overset{1}{|} \overset{2}{|} \overset{2,3}{|} \dots | \text{ooooooooooooooooooooo}$$

$Q_k \quad ps^1; a^1q$

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

$$r \quad \overset{1}{|} \overset{2}{|} \overset{2,3}{|} \dots \overset{\infty}{|} \text{oooooooooooooooooooo}$$

$Q_k \text{ ps}^1; a^1 q$

our own predictions used as targets!

Temporal Difference vs Monte Carlo

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

Which is better?

Temporal Difference

$$y(s; a) = r + \gamma Q_k(s'; a')$$

Monte Carlo

$$y(s; a) = r + \gamma (r^1 + \gamma^2 r^2 + \dots)$$

Temporal Difference vs Monte Carlo

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

Which is better?

Temporal Difference

$$y(s; a) = r + \gamma Q_k(s'; a')$$

X updates after each step
slow reward propagation

Monte Carlo

$$y(s; a) = r + \gamma r^1 + \gamma^2 r^2 + \dots$$

updates at the end of the episode
X fast reward propagation

Temporal Difference vs Monte Carlo

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

Which is better?

Temporal Difference

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

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

Monte Carlo

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

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

Temporal Difference vs Monte Carlo

$$Q_k(s; a) \leftarrow Q_k(s; a) + \alpha (r + \gamma V_k(s') - Q_k(s; a))$$

Which is better?

Temporal Difference

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

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

Monte Carlo

$$Q_k(s; a) \leftarrow r + \gamma V_k(s') - r + r + \gamma V_k(s') - r + r + \gamma V_k(s') - r + \dots$$

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

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

Types of backups

Q-learning

Value Iteration and Policy Iteration connection

What if we improve our policy after every step?

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

Value Iteration and Policy Iteration connection

What if we improve our policy after every step?

$$Q_k(p, s, a) \in \arg\max_a Q_k(p, s; a)$$

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

Value Iteration and Policy Iteration connection

What if we improve our policy after every step?

$$Q_k(p,s) \stackrel{\text{D}}{=} \operatorname{argmax}_a Q_k(p,s; a)$$

$$Q_{k+1}(p,s; a) \stackrel{\text{D}}{=} Q_k(p,s; a) + \gamma \left(\sum_{q \in \mathcal{A}} p(q) \left[\max_{a'} Q_k(p,s; a') - Q_k(p,s; a) \right] \right)$$

Value Iteration and Policy Iteration connection

What if we improve our policy after every step?

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

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

$$Q_k(p, s; a, q) + \gamma \left(\sum_{s'} P_{ss'}(a) Q_k(p, s'; \arg\max_{a^1} Q_k(p, s'; a^1, q) - Q_k(p, s; a, q) \right)$$

$$Q_k(p, s; a, q) + \gamma \left(\sum_{s'} P_{ss'}(a) \operatorname{argmax}_{a^1} Q_k(p, s'; a^1, q) - Q_k(p, s; a, q) \right)$$

Value Iteration and Policy Iteration connection

What if we improve our policy after every step?

$$Q_k(p,s) \stackrel{\Delta}{=} \operatorname{argmax}_a Q_k(p,s; a)$$

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

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

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

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

Value Iteration and Policy Iteration connection

What if we improve our policy after every step?

$$Q_k(p,s) \stackrel{\Delta}{=} \operatorname{argmax}_a Q_k(p,s; a)$$

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

$$Q_k(p,s; a) + \gamma \sum_{s'} P_{ss'}(a) [Q_k(p,s'; a) - Q_k(p,s; a)]$$

$$Q_k(p,s; a) + \gamma \sum_{s'} P_{ss'}(a) [Q_k(p,s'; \operatorname{argmax}_{a^1} Q_k(p,s'; a^1)) - Q_k(p,s; a)]$$

$$Q_k(p,s; a) + \gamma \sum_{s'} P_{ss'}(a) [\max_{a^1} Q_k(p,s'; a^1) - Q_k(p,s; a)]$$

Value Iteration and Policy Iteration connection

What if we improve our policy after every step?

$$Q_k(p; a) \stackrel{\Delta}{=} \operatorname{argmax}_a Q_k(p; a)$$

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

$$Q_k(p; a) + \gamma \sum_{q'} P(q' | p, a) [Q_k(p; a') - Q_k(p; a)]$$

$$Q_k(p; a) + \gamma \sum_{q'} P(q' | p, a) [Q_k(p; \operatorname{argmax}_{a'} Q_k(p; a')) - Q_k(p; a)]$$

$$Q_k(p; a) + \gamma \sum_{q'} P(q' | p, a) [\max_{a'} Q_k(p; a') - Q_k(p; a)]$$

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

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

Convergence of temporal difference

Convergence of Q-learning updates

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

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

Then, if:

$$\sum_k \alpha_k = \infty \quad \text{and} \quad \sum_k \alpha_k^2 < \infty$$

Convergence of temporal difference

Convergence of Q-learning updates

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

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

Then, if:

- $\sum_{k=0}^{\infty} \alpha_k = \infty$
- with probability 1 for every s, a a learning rate $\alpha_k(s; a) \in [0, 1]$ satisfies Robbins-Monro conditions:

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

Convergence of temporal difference

Convergence of Q-learning updates

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

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

Then, if:

- $\sum_{k=0}^{\infty} \alpha_k = \infty$
- with probability 1 for every s, a a learning rate $\alpha_k \in (0, 1]$ satisfies Robbins-Monro conditions:

$$\sum_{k=0}^{\infty} \alpha_k^2 < \infty$$

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

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 ϵ -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 $1 - \epsilon$

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

"-greedy exploration

Do something random sometimes

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

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

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

Decaying exploration:

$$\epsilon_k \sim \frac{1}{k}; \quad \epsilon_k \geq 0$$

Persistent exploration:

$$\epsilon_k = \text{const} \cdot \frac{1}{k}$$

Q-learning (online)

Initialize Q ρ ; aq arbitrarily;

observe s_0 ;

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

\hat{a}_k take action a_k " -greedy Q ρ ; aq

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

Q-learning (online)

Initialize Q ρ ; a_q arbitrarily;

observe s_0 ;

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

\hat{a}_k take action a_k "greedy" Q ρ ; a_q

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

$y := r_k + \max_{a_{k+1}} Q$ ρ ; $s_{k+1}; a_{k+1}$

Q-learning (online)

Initialize Q ρ s; a q arbitrarily;

observe s_0 ;

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

\hat{a}_k take action a_k " -greedy Q ρ s $_k$; a q

 observe r_k ; s_{k+1} ;

$y := r_k + \max_{a_{k+1}} Q$ ρ s $_{k+1}$; a_{k+1} q

Q ρ s $_k$; a_k q $\leftarrow \rho + (1 - \rho) (Q$ ρ s $_k$; a_k q + y)

Q-learning

Q-learning (online)

Initialize Q ψ ; a_q arbitrarily;

observe s_0 ;

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

$\hat{a}_k = \text{argmax}_a Q \psi_k; a_q$ " -greedy;

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

$y = r_k + \max_{a_{k+1}} Q \psi_{k+1}; a_{k+1} - Q \psi_k; \hat{a}_k$

$Q \psi_k; \hat{a}_k \leftarrow Q \psi_k; \hat{a}_k + \alpha (y - Q \psi_k; \hat{a}_k)$

Q-learning (with replay buffer)

Initialize Q ψ ; a_q arbitrarily, D H ;

Q-learning

Q-learning (online)

```
Initialize  $Q$   $\psi$ ;  $a_q$  arbitrarily;  
  
observe  $s_0$ ;  
for  $k = 0; 1; 2; \dots$   
  ^ take action  $a_k$  " -greedy  $Q$   $\psi_k; a_q$   
  ^ observe  $r_k; s_{k+1}$ ;  
  ^  $y := r_k + \max_{a_{k+1}} Q \psi_{k+1}; a_{k+1} q$   
  ^  $Q \psi_k; a_k q \leftarrow \beta + (1 - \beta) y$ 
```

Q-learning (with replay buffer)

```
Initialize  $Q$   $\psi$ ;  $a_q$  arbitrarily,  $D = H$  ;  
  
observe  $s_0$ ;  
for  $k = 0; 1; 2; \dots$   
  ^ take action  $a_k$  " -greedy  $Q$   $\psi_k; a_q$   
  ^ observe  $r_k; s_{k+1}$ ;  
  ^ store  $\psi_k; a_k; r_k; s_{k+1} q$  in  $D$ ;
```


Q-learning

Q-learning (online)

```

Initialize  $Q$   $\psi$ ;  $a_q$  arbitrarily;

observe  $s_0$ ;
for  $k = 0; 1; 2; \dots$ 
    ^ take action  $a_k$  " -greedy  $Q \psi_k; a_q$ 
    ^ observe  $r_k; s_{k+1}$ ;
    ^  $y = r_k + \max_{a_{k+1}} Q \psi_{k+1}; a_{k+1} q$ 
    ^  $Q \psi_k; a_k q \leftarrow \beta Q \psi_k; a_k q + (1 - \beta) y$ 

```

Q-learning (with replay buffer)

```

Initialize  $Q \psi$ ;  $a_q$  arbitrarily,  $D \sim \mathcal{H}$ ;

observe  $s_0$ ;
for  $k = 0; 1; 2; \dots$ 
    ^ take action  $a_k$  " -greedy  $Q \psi_k; a_q$ 
    ^ observe  $r_k; s_{k+1}$ ;
    ^ store  $\psi_k; a_k; r_k; s_{k+1} q$  in  $D$ ;
    ^ sample  $\psi; a; r; s^1 q$  from  $D$ ;
    ^  $y = r + \max_{a^1} Q \psi^1; a^1 q$ 
    ^  $Q \psi; a q \leftarrow \beta Q \psi; a q + (1 - \beta) y$ 

```

Behavior and target policies

The policy interacting with environment (collecting data) is called **behavior policy**;

Behavior and target policies

The policy interacting with environment (collecting data) is called **behavior policy**;

The trained policy that will be outputted by algorithm is called **target policy** ;

Behavior and target policies

The policy interacting with environment (collecting data) is called **behavior policy**;

The trained policy that will be outputted by algorithm is called **target policy** ;

If we learn Q using samples from π , we are **on-policy**. If $\pi \neq \pi'$, we are **off-policy**.

Behavior and target policies

The policy interacting with environment (collecting data) is called **behavior policy**;

The trained policy that will be outputted by algorithm is called **target policy**;

If we learn Q using samples from π , we are **on-policy**. If $\pi \neq \pi'$, we are **off-policy**.

Q-learning

off-policy

(arbitrary π' can be used)

Vanilla version: ϵ -greedy Q

Monte Carlo

on-policy

($\pi = \pi'$ is required)

Experience replay is useless!

Q-learning is ϵ -policy

Data source	Behavior policy	Convergence
Online experience	On-line visitation	Q
Expert data	-	Q
Experience replay	On-line visitation	Q

SARSA

Expectation

Expectation

Reality

Expectation

Reality

Consider the family of ϵ -soft policies :

$$\pi_\epsilon(a|s) = \frac{\pi(a|s) + \epsilon}{|A|}$$

Family of ϵ -soft policies : what changes

Policy Evaluation : nothing changes, Q is the same.

Family of ϵ -soft policies : what changes

Policy Evaluation : nothing changes Q^* is the same.

Optimal value function:

$$Q^*_{\epsilon\text{-greedy}}(p; a) = \max_{P^{\epsilon\text{-soft}}} Q(p; a)$$

Family of ϵ -soft policies : what changes

Policy Evaluation : nothing changes, Q is the same.

Optimal value function:

$$Q_{\epsilon\text{-greedy}}(p; a) = \max_{P \text{ } \epsilon\text{-soft}} Q(p; a)$$

Optimal policy:

$$p_a | s \leftarrow \epsilon\text{-greedy}(Q_{\epsilon\text{-greedy}}(p; a))$$

Family of ϵ -soft policies : what changes

Policy Evaluation : nothing changes, Q is the same.

Optimal value function:

$$Q_{\epsilon\text{-greedy}}(p; aq) = \max_{P \text{ -soft}} Q(p; aq)$$

Optimal policy:

$$p_a | s_q = \epsilon\text{-greedy}(Q_{\epsilon\text{-greedy}}(p; aq))$$

ϵ -soft Policy Improvement

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

$$\tilde{\pi} : \epsilon\text{-greedy}(Q_{\tilde{\pi}}(p; aq)) \quad \tilde{\pi} \sim \pi \quad \tilde{\pi} \succ \pi$$

Q-learning vs SARSA

$$Q_k(s; a) \leftarrow Q_k(s; a) + \alpha [r + \max_{a'} Q_k(s'; a') - Q_k(s; a)]$$

a^1 must be taken from target policy

Q-learning vs SARSA

$$Q_k(s; a) \leftarrow Q_k(s; a) + \alpha [r + \max_{a'} Q_k(s; a') - Q_k(s; a)]$$

a^1 must be taken from target policy

Q-learning

SARSA

Target policy:

$$\arg \max_a Q_k(s; a)$$

"-greedy $Q_k(s; a)$

Behavior policy:

"-greedy $Q_k(s; a)$

"-greedy $Q_k(s; a)$

Q-learning vs SARSA

$$Q_k(s; a) \leftarrow Q_k(s; a) + \alpha [r + \max_{a'} Q_k(s; a') - Q_k(s; a)]$$

a^1 must be taken from target policy

Q-learning

SARSA

Target policy:

$$\arg \max_a Q_k(s; a)$$

"-greedy $Q_k(s; a)$

Behavior policy:

"-greedy $Q_k(s; a)$

"-greedy $Q_k(s; a)$

Q-learning vs SARSA

$$Q_k(s; a) \leftarrow Q_k(s; a) + \alpha [r + \max_{a'} Q_k(s'; a') - Q_k(s; a)]$$

a^1 must be taken from target policy

Q-learning

SARSA

Target policy:

$$\arg \max_a Q_k(s; a)$$

"-greedy $Q_k(s; a)$

Behavior policy:

"-greedy $Q_k(s; a)$

"-greedy $Q_k(s; a)$

Convergence:

$$Q_k(s; a) \rightarrow \arg \max_a Q(s; a)$$

"-greedy $Q(s; a)$

SARSA

SARSA

Initialize Q ρ ; a 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 $a_{k+1} \sim \pi(\cdot | s_{k+1}; a_k)$

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

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

^ on-policy by construction

SARSA

SARSA

Initialize Q ρ ; a 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 $a_{k+1} \sim \pi$ -greedy Q ρ ; a

^ $y = r_k + \gamma Q$ ρ ; a_{k+1}

^ Q ρ ; $a_k \leftarrow (1 - \alpha) Q$ ρ ; $a_k + \alpha y$

^ on-policy by construction

^ same as Q-learning if $\epsilon = 0$

^ converges to Q^* if $\sum_k \alpha_k = \infty$

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 $a_{k+1} \sim \pi$ -greedy $Q(p_{s_k}; a_q)$

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

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

^ on-policy by construction

^ same as Q-learning if $\pi = \pi^*$

^ converges to Q^* if $\sum_k \alpha_k = \infty$

Expected-SARSA:

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

SARSA is on-policy

For transition $p(s, a; r; s^1; a^1q$

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

Data source	Target policy	Convergence
Online experience	Decaying exploration	Q
	Persistent exploration	Q _n -greedy
Expert data (a ¹ taken from bu er)	-	

SARSA is on-policy

For transition $p(s, a; r; s^1; a^1q$

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

Data source	Target 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 $p(s, a, r, s')$; a^1, q

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

Data source	Target 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

O-policy SARSA (not really sarsa now)

Can we learn Q for given o-policy?

O-policy SARSA (not really sarsa now)

Can we learn Q for given π -policy?

For transition $p(s', a; r; s, a)$ generate $a^1 \sim \pi(a^1 | s^1)$

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

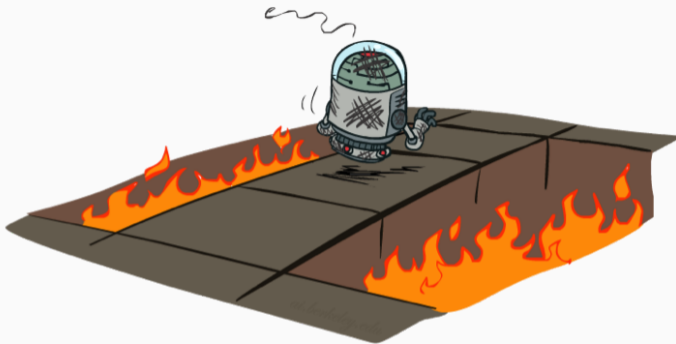
O-policy SARSA (not really sarsa now)

Can we learn Q for given π -policy?

For transition $p(s, a; r; s', q)$ generate $a^1 \sim p(a^1 | s^1, q)$

$$Q_k(s, a; q) \leftarrow Q_k(s, a; q) + \alpha (r + \sum_{s'} p(s', a; q) Q_k(s', a; q) - Q_k(s, a; q))$$

Data source	Target policy	Convergence
(any, but a^1 generated online)	Decaying exploration	Q
	Persistent exploration	Q_{π} -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;*