
Reinforcement learning

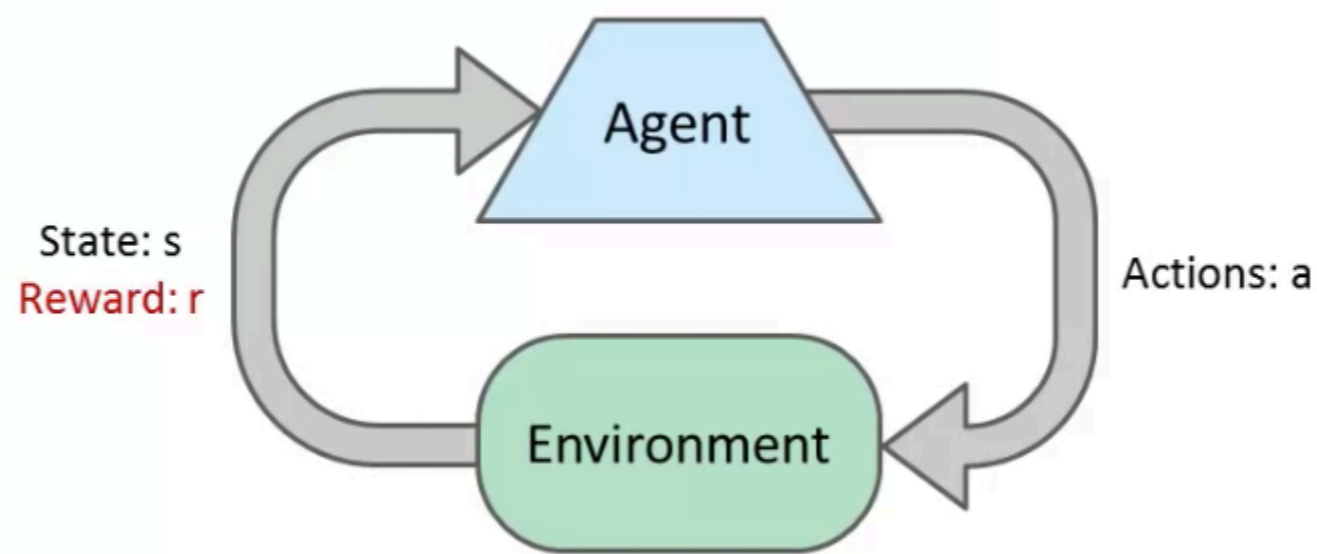
Q learning, SARSA



Reinforcement learning



Reinforcement learning



- Základné predpoklady
- Agent interaguje s prostredím a získava po každej akcii odmenu
- Úžitok agenta je definovaný získanými diskontovanými odmenami
- Cieľ: vykonávať akcie maximalizujúce jeho očakávaný úžitok
- Agent nemá vedomosti o modele. Učí sa na základe vzorkovania (samplingu) prostredia



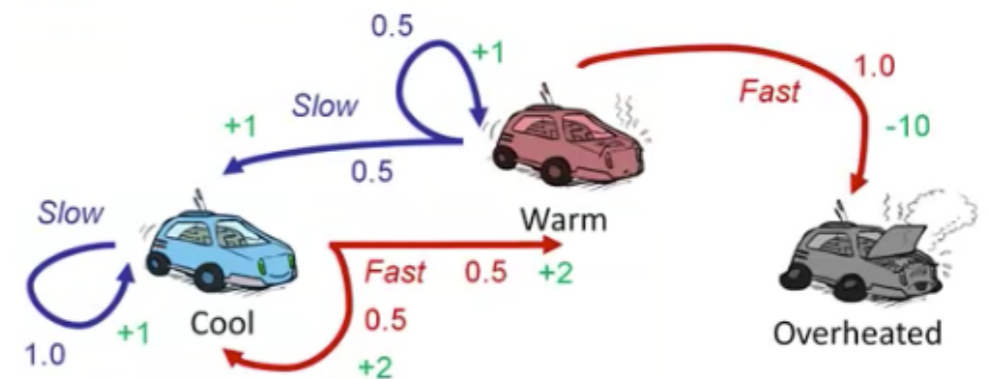
Príklady



Reinforcement learning

- ▶ Uvažujeme stále základný model MDP

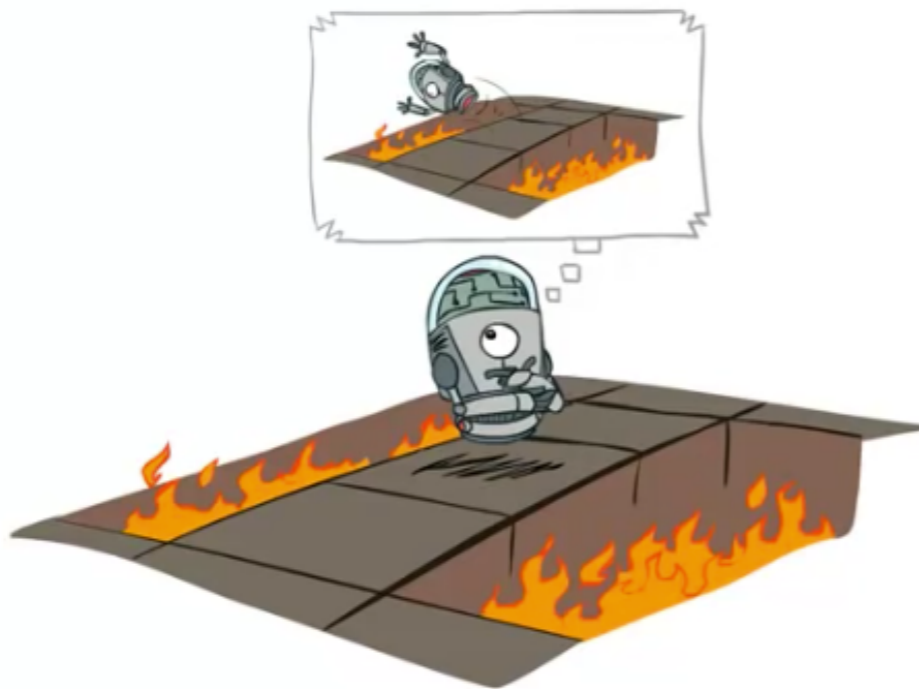
- ▶ Množinu stavov $s \in \mathcal{S}$
- ▶ Množinu akcií nad každým stavom $a, a \in \mathcal{A}$
- ▶ Dynamiku modelu $p(s', r | s, a)$
- ▶ Odmenu, resp odmeňovaciú funkciu $R(s, a, s')$
- ▶ Cieľ: nájsť optimálny strategický profil $\pi_*(s)$



Nepoznáme dynamiku modelu a odmeňovaciú funkciu!!!!!!



Offline (MDP) vs online (RL)



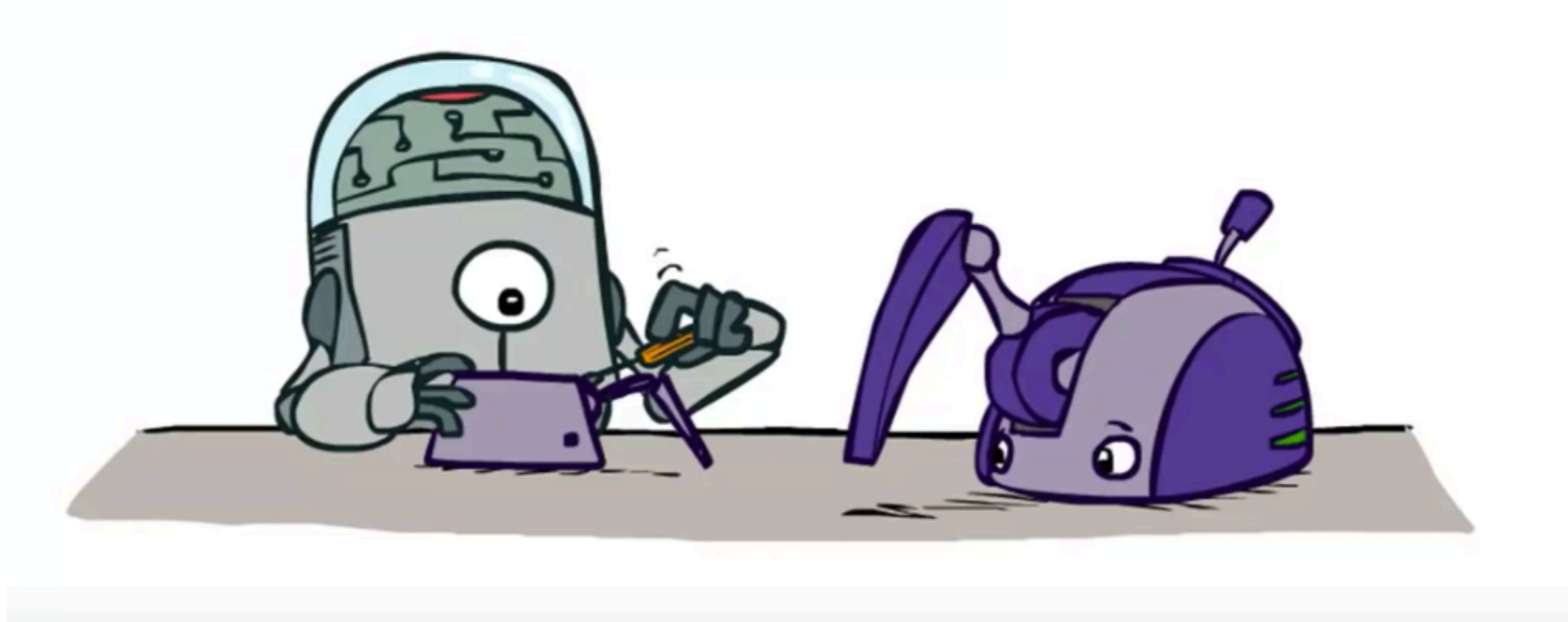
Riešenie offline (dynamika je známa)



Riešenie online (zistujeme parametre modelu)



Model-based learning (učenie na báze modelu)



Model-based learning (učenie na báze modelu)

- ▶ Učenie na báze modelu:

Postupné učenie parametrov modelu na základe interakcie

Aplikácia dosiaľ uvedených algoritmov na naučenom modeli

- ▶ Krok 1. Učenie

Určenie dynamiky modelu, rátaj pravdepodobnosť výskytu stavu s' pre každé s a a

Normovanie vedie k získaniu aproximovaných hodnôt $p(s', r | s, a)$ a $\hat{r}(s, a, s')$



Model-based learning (učenie na báze modelu)

- ▶ Učenie na báze modelu:

Postupné učenie parametrov modelu na základe interakcie

Aplikácia dosiaľ uvedených algoritmov na naučenom modeli

- ▶ Krok 1. Učenie

Určenie dynamiky modelu, vrátaj pravdepodobnosť výskytu stavu s' pre každé s a a

Normovanie vedie k získaniu aproximovaných hodnôt $p(s', r | s, a)$ a $\hat{r}(s, a, s')$

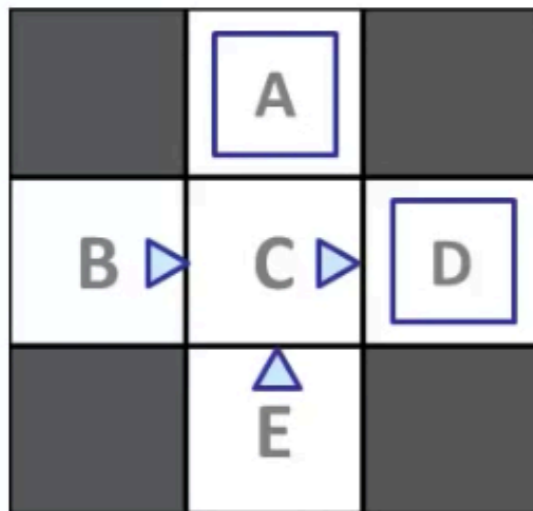
- ▶ Krok 2. Implementácia algoritmov

Value iteration, Policy iteration, resp. Policy evaluation



Model-based learning (učenie na báze modelu)

Input Policy π



Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

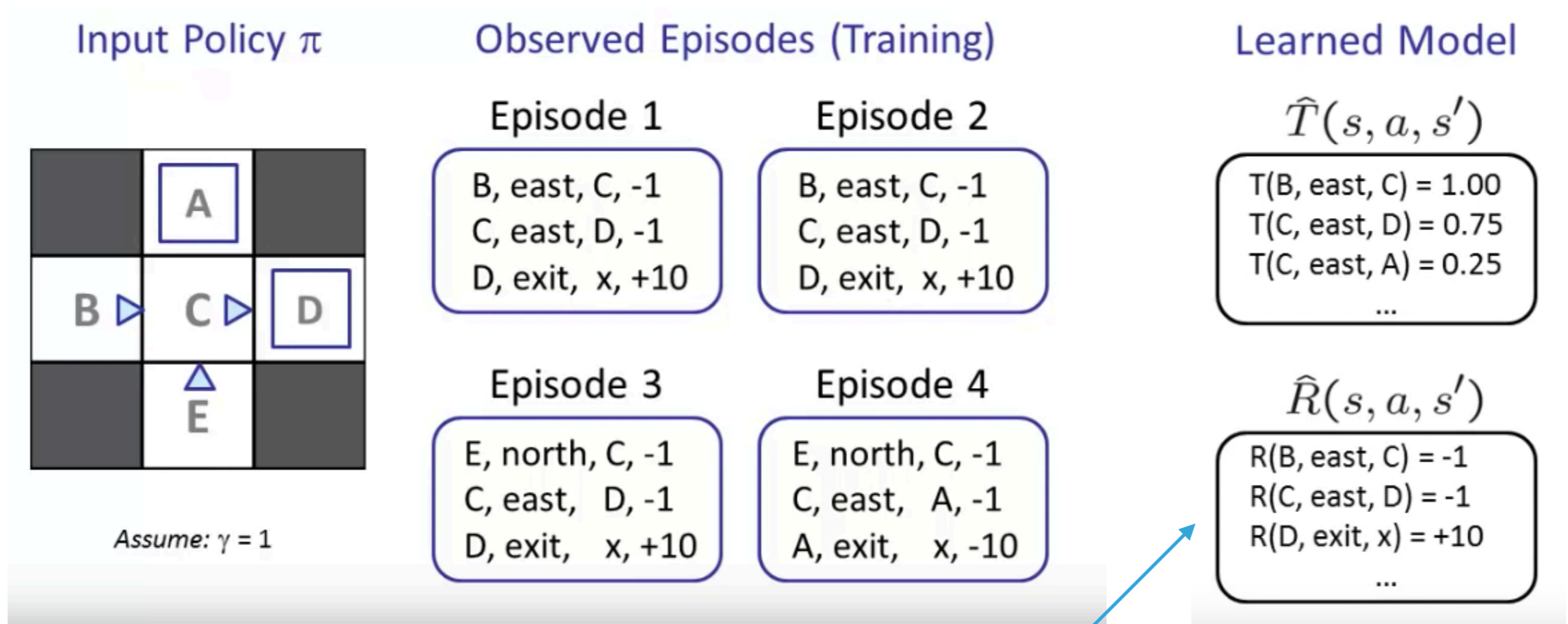
E, north, C, -1
C, east, D, -1
D, exit, x, +10

Episode 4

E, north, C, -1
C, east, A, -1
A, exit, x, -10



Model-based learning (učenie na báze modelu)



Value Iteration, Policy Improvement, etc

Nevýhoda: zistenie parametrov modelu trvá čas



Model based vs model free učenie

Ciel': zistenie vekovej skupiny nášho predmetu

Predpoklad: perfektná vedomosť hustoty pravdepodobnosti:

$$E(A) = \sum_a P(a)a$$

Ak predpoklad nie je splnený, zbieraj vzorky študentov $a_1, a_2, a_3, \dots, a_N$

Výpočet (je nutné vypočítať $P(a)$ pre Model Based))

$\hat{P}(a) = \frac{\text{num}(a)}{N}$, kde N je počet vzoriek študentov množiny A

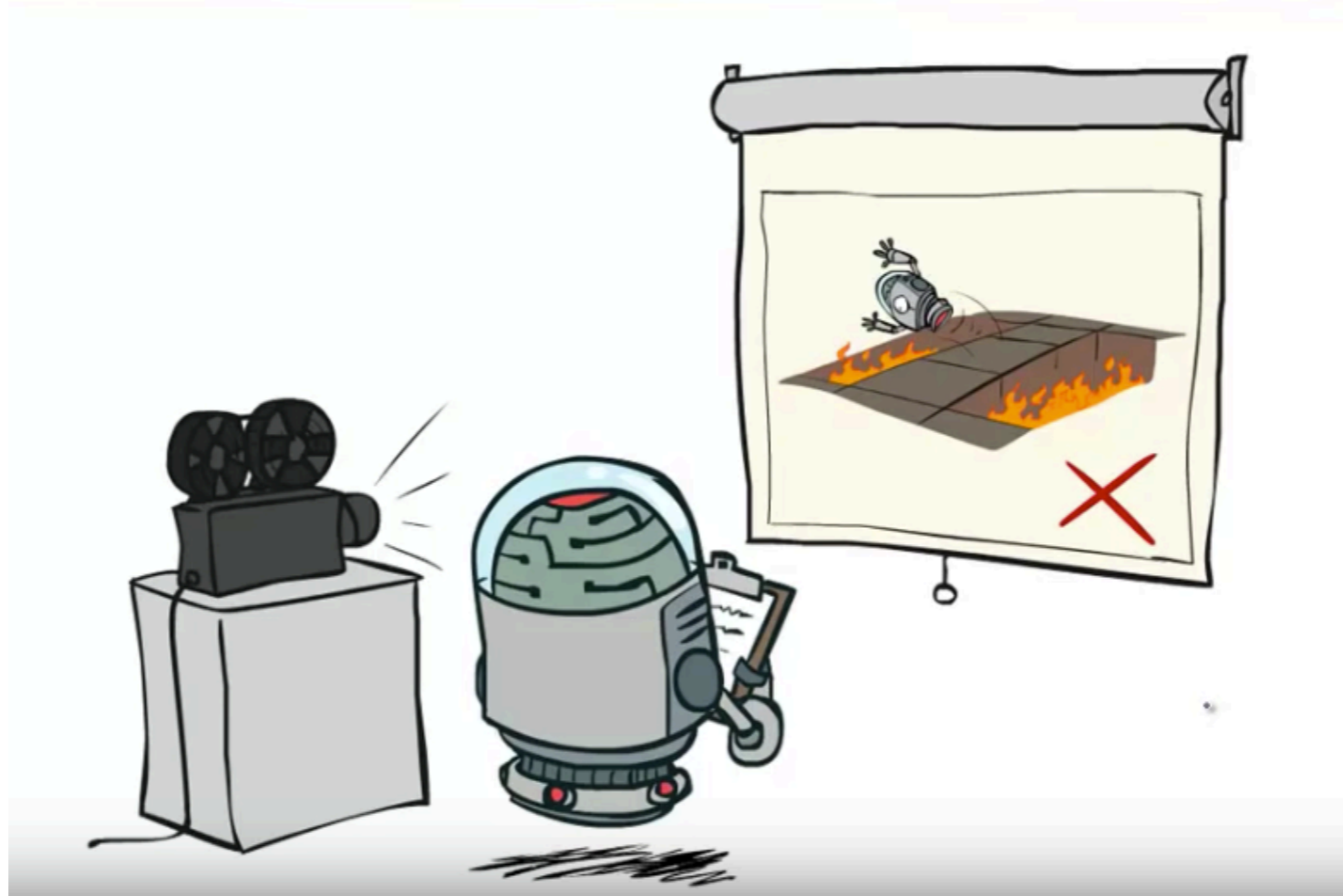
$$E(A) \approx \sum_a \hat{P}(a)a$$

Model free

$$E(A) \approx \frac{1}{N} \sum_i a_i$$



Pasívny reinforcement learning (prediction problem)



Pasívny reinforcement learning (prediction problem)

Zjednodušená úloha: evaluácia daného strategického profilu

Vstup: fixný strategický profil $\pi(s)$

$R(s, a, s')$ a $p(s', r | s, a)$ sú neznáme

Cieľ: určenie hodnôt jednotlivých stavov (i.e. $V(s), \forall s$)

V tomto prípade:

Agent vykonáva akcie free ride na základe definovaného strat. profilu

Strategický profil je konštantný

Vykonáva akcie na základe strategického profilu a ohodnocuje jednotlivé stavy



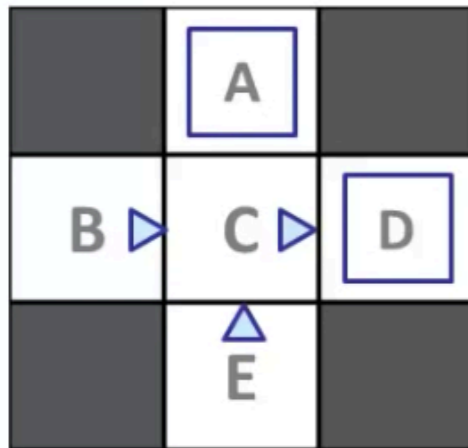
Priama evaluácia

- **Cieľ:** určenie hodnoty jednotlivých stavov na základe $\pi(s)$
- **Myšlienka:** Aritmetický priemer pozorovaných hodnôt (model free)
 - Vykonávajú akcie na základe $\pi(s)$
 - V každom časovom okamihu zapíš sumu diskotovaných odmien pre daný stav
 - Spriemerni (priama evaluácia)



Priama evaluácia

Input Policy π



Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, x, +10

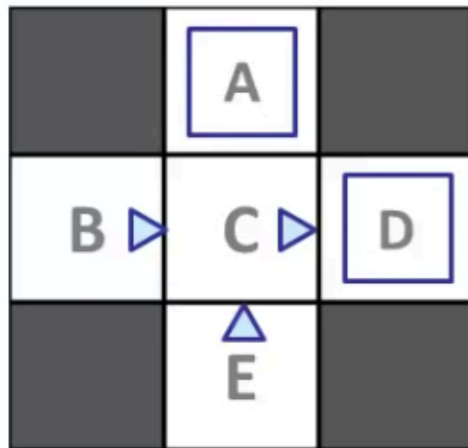
Episode 4

E, north, C, -1
C, east, A, -1
A, exit, x, -10



Priama evaluácia

Input Policy π



Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, x, +10

Episode 4

E, north, C, -1
C, east, A, -1
A, exit, x, -10

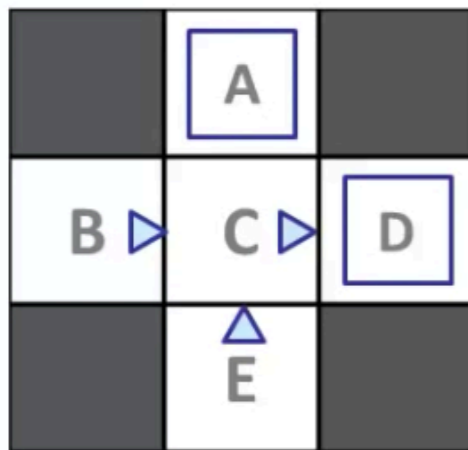
Output Values

$$V_{\pi}(C) = \frac{9 + 9 + 9 - 11}{4}$$



Priama evaluácia

Input Policy π



Assume: $\gamma = 1$

Observed Episodes (Training)

Episode 1

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 2

B, east, C, -1
C, east, D, -1
D, exit, x, +10

Episode 3

E, north, C, -1
C, east, D, -1
D, exit, x, +10

Episode 4

E, north, C, -1
C, east, A, -1
A, exit, x, -10

Output Values

$$V_{\pi}(C) = \frac{9 + 9 + 9 - 11}{4}$$

	-10	
	A	
+8	+4	+10
B	C	D
	-2	
	E	

Strácame súvislosť medzi stavmi! E nereflektuje na hodnotu C

Neuplatňujeme Bellmana, ale je to jednoduché no neefektívne



Policy evaluation?

Úvaha: Využitie policy evaluation pre konštantný strategický profil

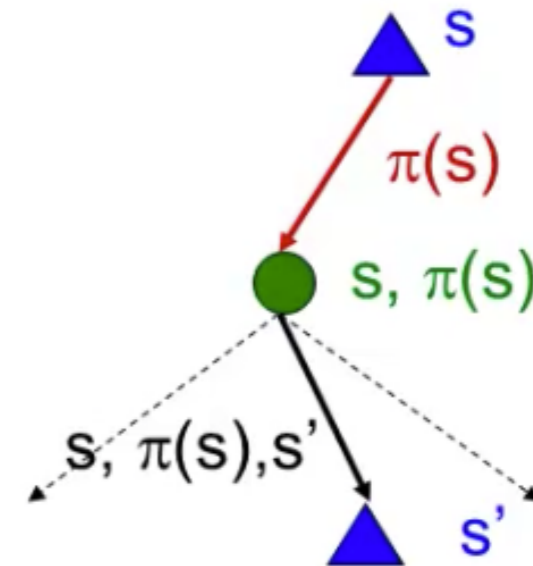
Pre každý časový okamih interakcie agenta s prostredím urobme update hodnoty stavu

$$V_{\pi}^0(s) = 0$$

$$V_{\pi}^{k+1}(s) \leftarrow \sum_{s'} T(s, \pi(s), s') \left[R(s, \pi(s), s') + \gamma V_{\pi}^k(s') \right]$$



Nepoznáme to, riešenie: vzorkovanie



Otázka: Nepoznáme váhy $T(s, \pi(s), s')$, ako môžeme uplatniť daný rekurzívny vzťah?



Policy evaluation by sampling?

Ciel': Snaha zlepšiť presnosť ohodnotenia stavu $V(s)$ spriemerovaním jednotlivých pozorovaní

$$V_{\pi}^{k+1}(s) \leftarrow \sum_{s'} T(s, \pi(s), s') \left[R(s, \pi(s), s') + \gamma V_{\pi}^k(s') \right]$$

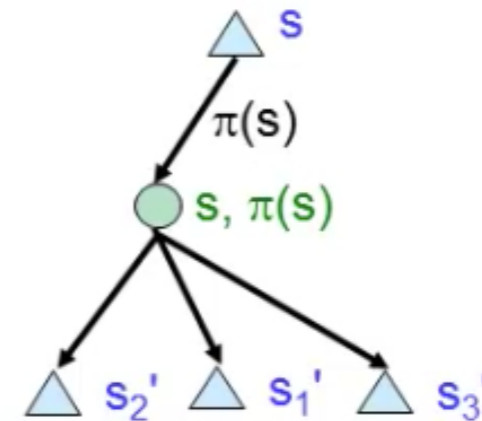
Riešenie: Vzorkovanie (sampling) výstupov zo stavov s' a spriemerovanie

$\pi(s)$, resp alternatívne akcia a v prípade determ. str.profilu

$$sample_1 = R(s, \pi(s), s'_1) + \gamma V_k^{\pi}(s'_1)$$

$$sample_2 = R(s, \pi(s), s'_2) + \gamma V_k^{\pi}(s'_2)$$

$$sample_3 = R(s, \pi(s), s'_n) + \gamma V_k^{\pi}(s'_n)$$



Policy evaluation by sampling?

Ciel': Snaha zlepšiť presnosť ohodnotenia stavu $V(s)$ spriemerovaním jednotlivých pozorovaní

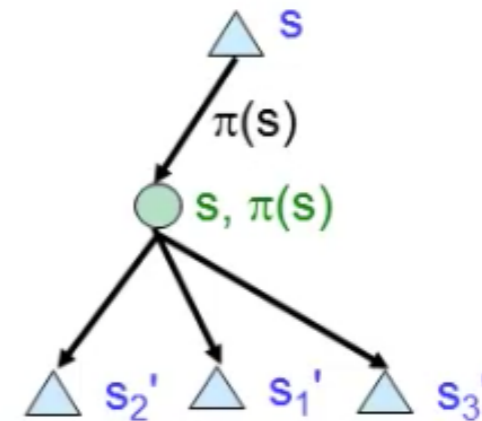
$$V_{\pi}^{k+1}(s) \leftarrow \sum_{s'} T(s, \pi(s), s') \left[R(s, \pi(s), s') + \gamma V_{\pi}^k(s') \right]$$

Riešenie: Vzorkovanie (sampling) výstupov zo stavov s' a spriemerovanie

$$sample_1 = R(s, \pi(s), s'_1) + \gamma V_k^{\pi}(s'_1)$$

$$sample_2 = R(s, \pi(s), s'_2) + \gamma V_k^{\pi}(s'_2)$$

$$sample_3 = R(s, \pi(s), s'_n) + \gamma V_k^{\pi}(s'_n)$$



$$V_{\pi}^{k+1} = \frac{1}{N} \sum_i sample_i$$

Problém: MDP: nevieme sa vracat' v čase

$$V_k^{\pi}(s)$$



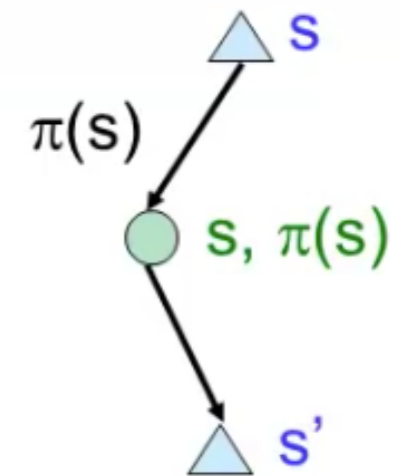
Temporal difference learning

- ▶ **Ciel:** učenie na základe každej interakcie s modelom
 - ▶ Update $V_{\pi}^k(s)$ po jednotlivých prechodoch (s, a, s', r)
 - ▶ Stavy s' , ktorých pravdepodobnosť výskytu $p(s', r | s, a)$, resp. $T(s', \pi(s), s)$ je vyššia budú ovplyvňovať $V_{\pi}^k(s)$ viac

$$\text{Vzorka } V(s): \quad \text{sample}_1 \leftarrow R(s, \pi(s), s'_1) + \gamma V_{\pi}^k(s'_1)$$

$$\text{Vzorka } V(s): \quad V_{\pi}(s) \leftarrow (1 - \alpha)V_{\pi}(s) + \alpha * \text{sample}$$

$$\text{resp:} \quad V_{\pi}(s) \leftarrow V_{\pi}(s) + \alpha(\text{sample} - V_{\pi}(s))$$



Kĺzavý priemer (exponential moving average)

- ▶ Kĺzavý priemer pozorovaní:

$$\bar{x}(n) = (1 - \alpha) \cdot \bar{x}_{n-1} + \alpha \cdot x_n$$

- ▶ Posledné pozorované vzorky majú väčšiu váhu, ako tie ktoré boli pozorované v minulosti

$$\bar{x}(n) = \frac{x_n + (1 - \alpha) \cdot x_{n-1} + (1 - \alpha)^2 \cdot x_{n-2} + \dots}{1 + (1 - \alpha) + (1 - \alpha)^2 + \dots}$$

- ▶ **Snaha:** Prvé pozorovania boli značne vychýlené, ich vplyv sa bude s časom zmeňovať.

- ▶ **Cieľ:** Znižovanie parametra α nám zaručuje konvergenciu uvedených algoritmov

$$\alpha(n) = \frac{1}{n}$$

$$\alpha(n) \leftarrow \alpha(n) * 0.99$$



Príklad TD

States

	A	
B	C	D
	E	

Assume: $\gamma = 1$, $\alpha = 1/2$



Príklad TD

States

	A	
B	C	D
	E	

Observed Transitions

B, east, C, -2

	0	
0	0	8
	0	

	0	
-1	0	0
	0	

Assume: $\gamma = 1$, $\alpha = 1/2$

$$V_{\pi}(s) \leftarrow (1 - \alpha)V_{\pi}(s) + \alpha * (R(s, \pi(s), s') + \gamma \cdot V_{\pi}(s'))$$



Príklad TD

Observed Transitions

B, east, C, -2

	0	
0	0	8
	0	

	0	
-1	0	8
	0	

C, east, D, -2

	0	
-1	3	8
	0	

$$V_{\pi}(s) \leftarrow (1 - \alpha)V_{\pi}(s) + \alpha * (R(s, \pi(s), s') + \gamma \cdot V_{\pi}(s'))$$




Temporal difference: problémy

TD je vynikajúci prístup na ohodnocovanie konkrétneho strategického profilu $\pi(s)$ využitím základných Bellman rovníc a ich iteratívneho riešenia (na báze model free prístupu)

Problém: ako zlepšiť strategický profil (hodnotové funkcie stavov nám nehovoria nič o akciách)?

$$\pi(s) = \arg \max Q(s, a), \forall a \in \mathcal{A}$$

Nemožné: POLICY EXTRACTION sa použiť nedá

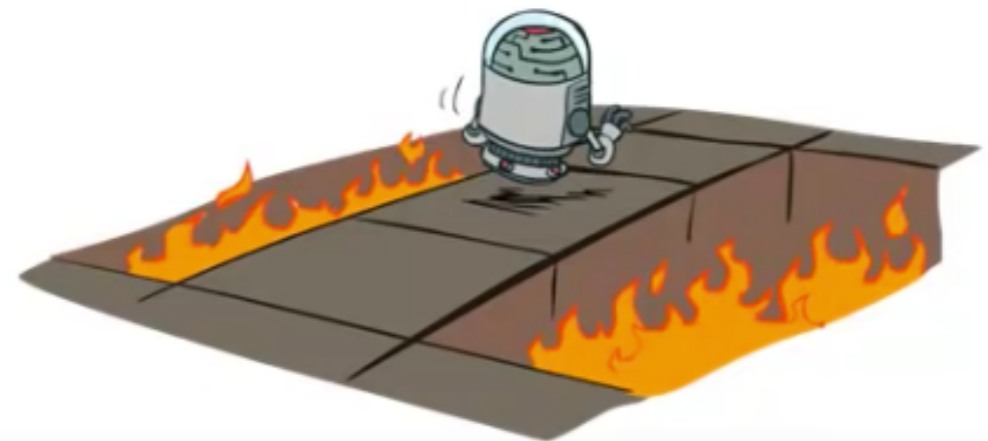
$$Q(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V(s')]$$


Nápad: učenie Q hodnôt



Aktívny RL (control based RL)

- Schopnosť zlepšovať strategický profil správanie agenta v neznámom prostredí
 - Nevieme prechodové pravdepodobnosti $p(s', r | s, a)$
 - Nevieme tvar odmeňovacej funkcie $R(s, a, s')$
 - Agent vyberá akcie a snaží sa ich zlepšiť za účelom maximalizácie očakávaného úžitku
- Agent rozhoduje o akciách v prostredí, nie daný strategický profil ako pri TD
 - Explorácia vs. Exploitation - problém ako pri banditoch



Aktívny RL (control based RL)

Recap: Value iteration

Inicializujme $V^0(s) = 0, \forall s$

Iteratívnym prístupom updatujme jednotlivé hodnoty pre všetky stavy s

$$V^{k+1}(s) = \max_a \sum_{s'} \pi(s, a, s') [R(s, a, s') + \gamma(V^k(s'))]$$

Iteratívnym prístupom updatujme jednotlivé hodnoty pre všetky stavy s

Problém: Našou snahou je zistiť Q hodnoty

Inicializujme $Q^0(s, a) = 0, \forall s, a$

Iteratívnym prístupom updatujme jednotlivé hodnoty $Q^k(s, a)$ postupne, na základe interakcií s modelom



Aktívny RL (control based RL)

Recap: Value iteration

Inicializujeme $V^0(s) = 0, \forall s$

Iteratívnym prístupom updatujeme jednotlivé hodnoty pre všetky stavy s

$$V^{k+1}(s) = \max_a \sum_{s'} \pi(s, a, s') [R(s, a, s') + \gamma(V^k(s'))] \longrightarrow \begin{array}{l} \text{jednoduchý sampling ako pri TD nestačí,} \\ \text{máme MAX} \end{array}$$

Iteratívnym prístupom updatujeme jednotlivé hodnoty pre všetky stavy s

Problém: Našou snahou je zistiť Q hodnoty

Inicializujeme $Q^0(s, a) = 0, \forall s, a$

Iteratívnym prístupom updatujeme jednotlivé hodnoty $Q^k(s, a)$ postupne, na základe interakcií s modelom



Aktívny RL (control based RL)

Recap: Value iteration

Inicializujeme $V^0(s) = 0, \forall s$

Iteratívnym prístupom updatujeme jednotlivé hodnoty pre všetky stavy s

$$V^{k+1}(s) = \max_a \sum_{s'} \pi(s, a, s') [R(s, a, s') + \gamma V^k(s')] \longrightarrow \begin{array}{l} \text{jednoduchý sampling ako pri TD nestačí,} \\ \text{máme MAX} \end{array}$$

Iteratívnym prístupom updatujeme jednotlivé hodnoty pre všetky stavy s

Problém: Našou snahou je zistiť Q hodnoty

Inicializujeme $Q^0(s, a) = 0, \forall s, a$

Iteratívnym prístupom updatujeme jednotlivé hodnoty $Q^k(s, a)$ postupne, na základe interakcií s modelom

$$\text{nevidujeme max} \quad Q^{k+1}(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V^k(s')]$$

$$Q^{k+1}(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q^k(s', a')]$$



Q learning

Q learning: updatovanie Q hodnôt jednotlivých stavov na základe interakcie s modelom

$$Q^{k+1} \leftarrow \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma \max_{a'} Q_k(s', a')]$$

Algoritmus: učenie $Q(s,a)$ postupne, ako prechádzame modelom

- 1) získanie vzorky (s,a,s',r)
- 2) Načítanie odhadu hodnoty $Q(s,a)$
- 3) Výpočet nového samplu v tvare $sample = R(s, a, s') + \gamma \max_{a'} Q(s', a')$
- 4) Aplikácia klzavého priemeru na nový update Q pre vykonanú akciu a v stave s hodnoty v tvare

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha(sample)$$



Q learning

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Loop 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

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \arg \max_{A'} Q(S_{t+1}, A'_{t+1}) - Q(S_t, A_t) \right]$$

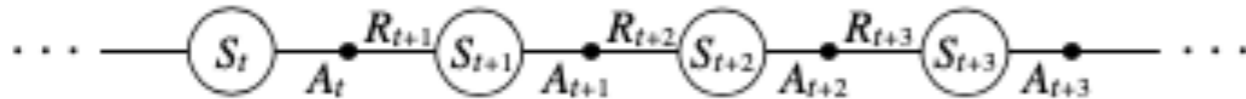


DEMO Q learning

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma \arg \max_{A'} Q(S_{t+1}, A'_{t+1}) - Q(S_t, A_t) \right]$$



SARSA learning



Ciel': transformácia TD na SARSA, náhradou $V(s)$ za $Q(s,a)$ podľa $\pi(s)$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t) \right]$$



pozor, nerobíme max (neaproximujeme $Q_*(s, a)$)

Tvrdenie: SARSA umožňuje nájsť optimálny strategický profil $\pi_*(s)$, za predpokladu, že všetky stavy v MDP sú navštívene dostatočne veľa krát a parameter ϵ klesá k nule dostatočne pomaly.



SARSA learning

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\varepsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

 Initialize S

 Choose A from S using policy derived from Q (e.g., ε -greedy)

 Loop for each step of episode:

 Take action A , observe R, S'

 Choose A' from S' using policy derived from Q (e.g., ε -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



Kvíz

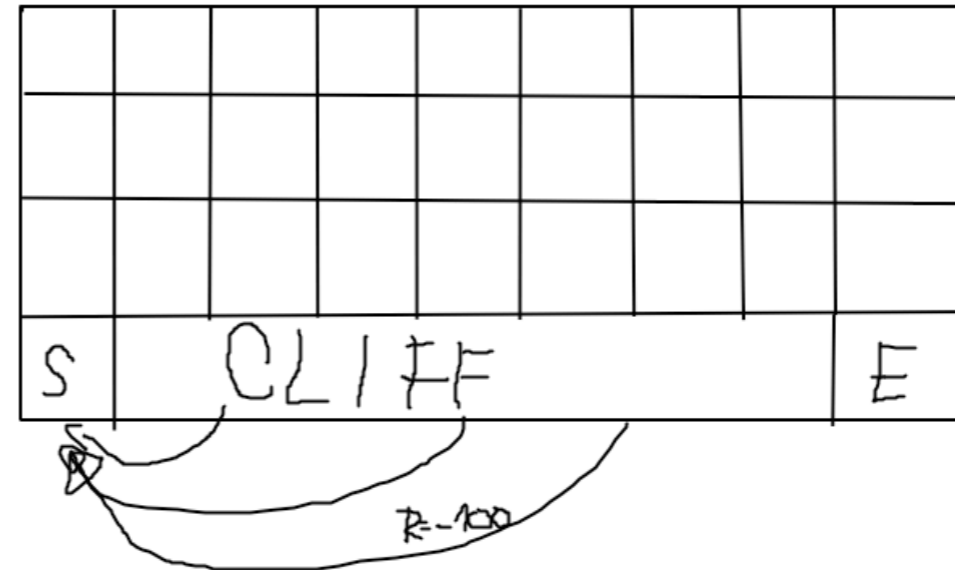
Optimálny strategický profil pre

konštantné $\epsilon = 0.2$

SARSA vs Q learning

Aký je očakávaný užitočnosť pre optimálne strategické profily?

$R = -1$



Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

Initialize S

Loop 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

Sarsa (on-policy TD control) for estimating $Q \approx q_*$

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$

Initialize $Q(s, a)$, for all $s \in \mathcal{S}^+$, $a \in \mathcal{A}(s)$, arbitrarily except that $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

Initialize S

Choose A from S using policy derived from Q (e.g., ϵ -greedy)

Loop for each step of episode:

Take action A , observe R, S'

Choose A' from S' using policy derived from Q (e.g., ϵ -greedy)

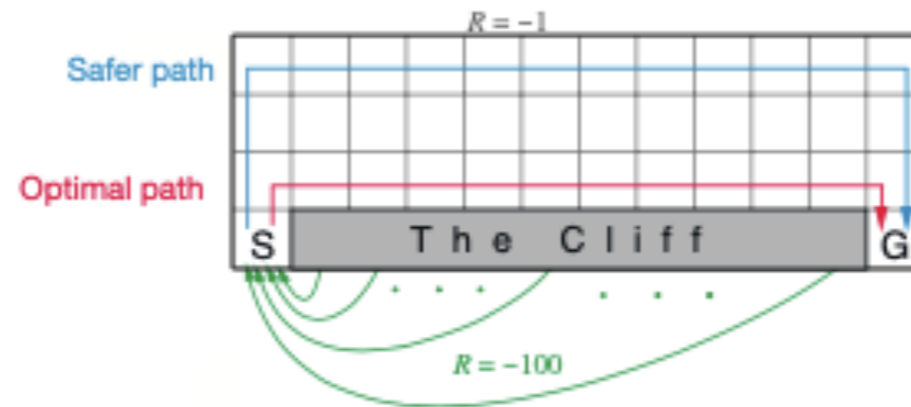
$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma Q(S', A') - Q(S, A)]$

$S \leftarrow S'; A \leftarrow A'$

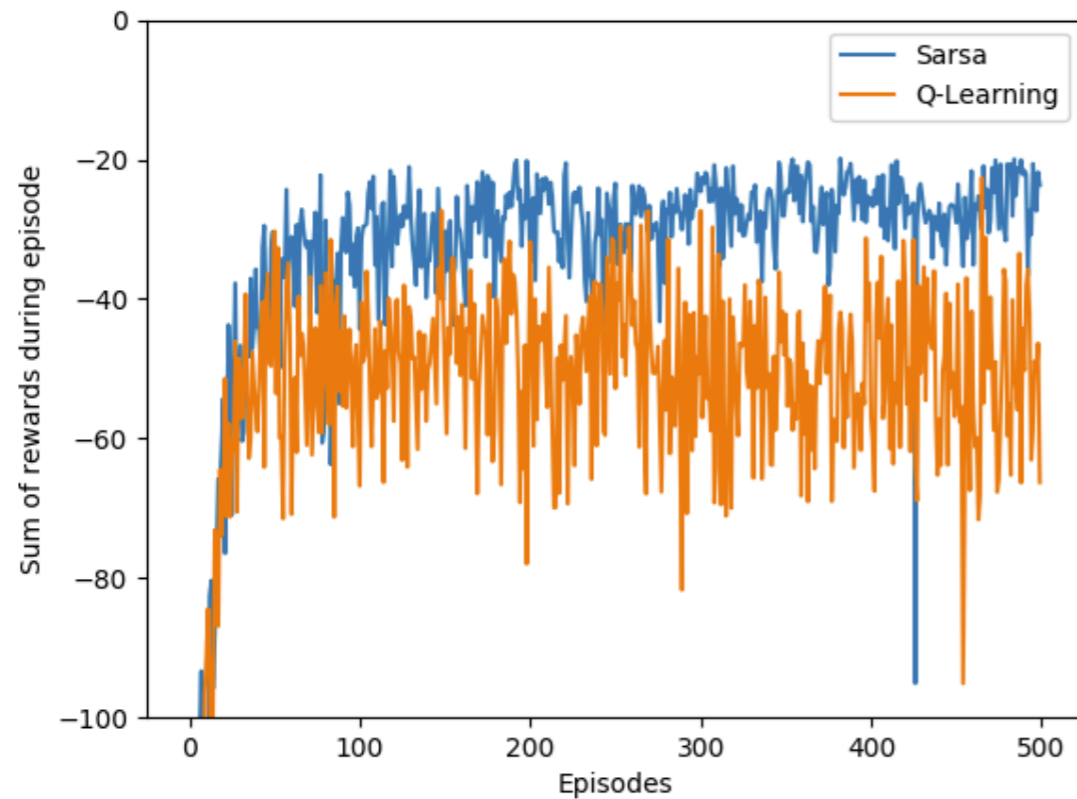
until S is terminal



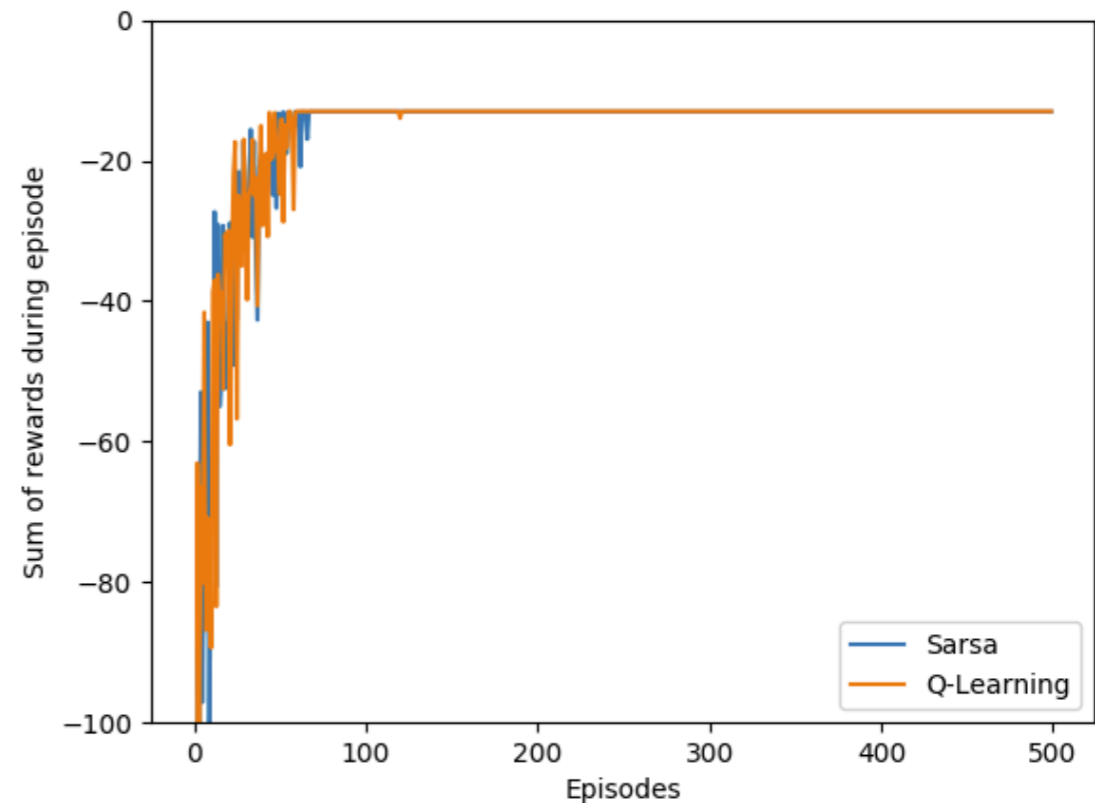
SARSA vs Q learning (príklad Python, 6.6)



$\epsilon = 0.2$



$\epsilon \leftarrow \epsilon * 0.99$



SARSA a Q learning

- **Tvrdenie:** Q learning umožňuje získanie optimálneho strategického profilu bezohľadu aký je strategický profil agenta (náhodný, pseudonáhodný, greedy, atď).
- Táto skutočnosť podmieňuje nazývať Q learning ako OFF-POLICY algoritmus
- **Predpoklady:**
 - Musíme dostatočne dlho spoznávať model (fáza exploration)
 - Správna rovnováha v parametri learning rate. Nesmie klesať príliš rýchlo, ale ani veľmi pomaly
 - Teoreticky pri Q learningu nie je podstatné ako navštevujeme stavy, ani akým strat. profilom sa riadime.



SARSA a Q learning

SARSA algoritmus: on policy algoritmus, update Q hodnôt jednotlivých stavov je vykonávaný podľa aktuálneho strategického profilu

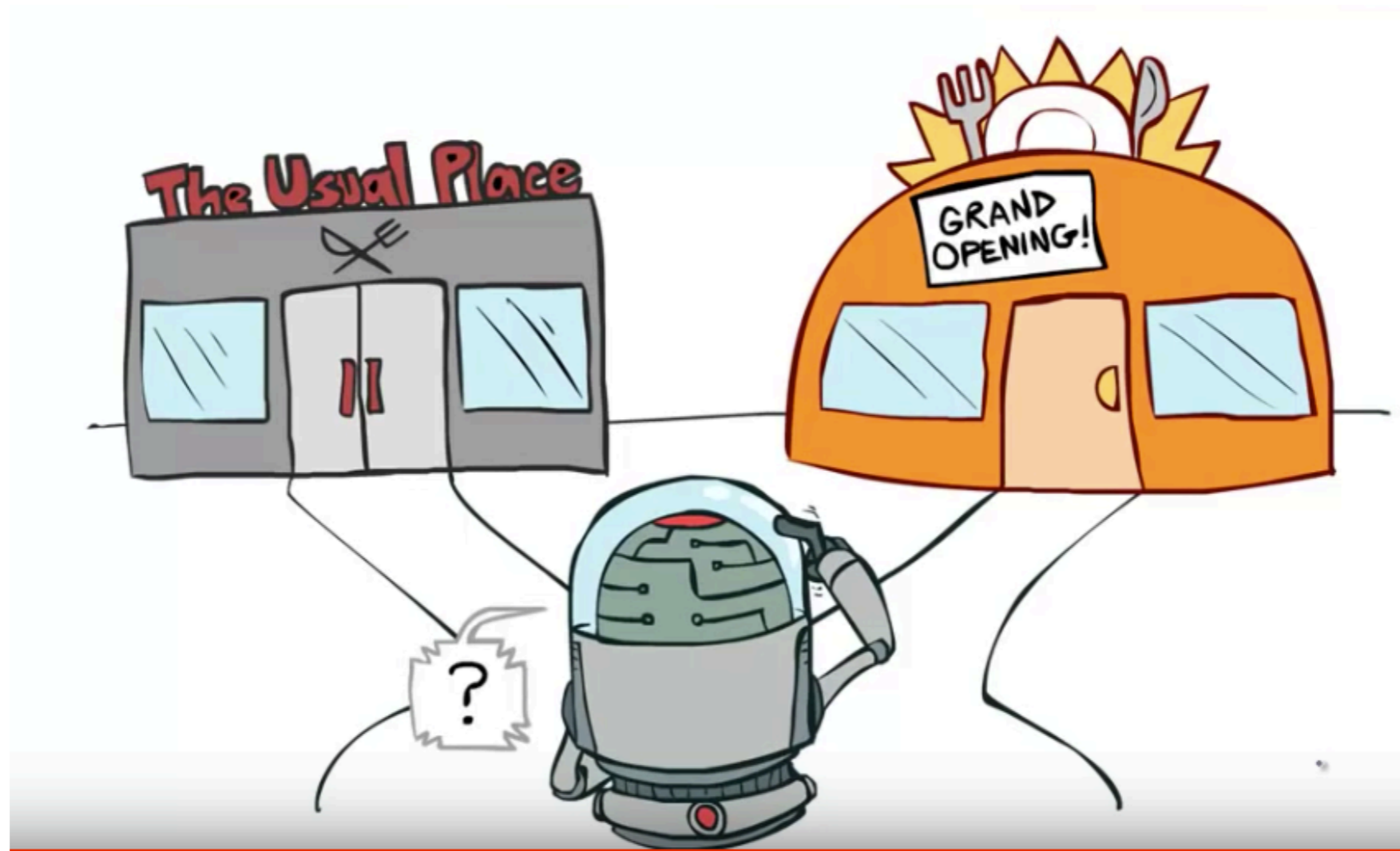
Je viac SAFE ako Q learning.

Pri vhodne klesajúcej miere explorácie dosahujú v tabuľkových formách oba algoritmy rovnaké výsledky



Exploration vs Exploitation

Bádanie v modeli, resp využívanie modelu



Exploration vs Exploitation

Vo všeobecnosti existuje niekoľko prístupov podporujúcich exploráciu

Jednoduché riešenie: ϵ greedy algoritmus

v každom časovom okamihu si hod' mincou

S malou $p()$ ϵ vykonaj náhodnú akciu

S veľkou $p()$ $1-\epsilon$, vykonaj akciu podľa súčasného strategického profilu

Úvaha: náhodné akcie spôsobujú po čase zníženie výkonnosti agenta

Po natrénovaní modelu môžu náhodné akcie agenta dostať do stavov pre neho nevýhodných

Riešenie: znižujeme parameter ϵ (analógia so simulovaným žiňaním)

