
Riešenie MDP

Juraj Gazda

iislab.kpi.fei.tuke.sk



Odmeny a epizódy

Reward vs return (odmena vs úžitok)

Agent: entita maximalizujúca očakávanú kumulatívnu odmenu:

$$G_t \doteq R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$$

Diskontovanie:

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Základné pravidlo RL (guess from the guess in general):

$$\begin{aligned} G_t &\doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots \\ &= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \gamma^2 R_{t+4} + \dots) \\ &= R_{t+1} + \gamma G_{t+1} \end{aligned}$$



Hodnota stavu, q stav

Hodnota (úžitok) stavu S

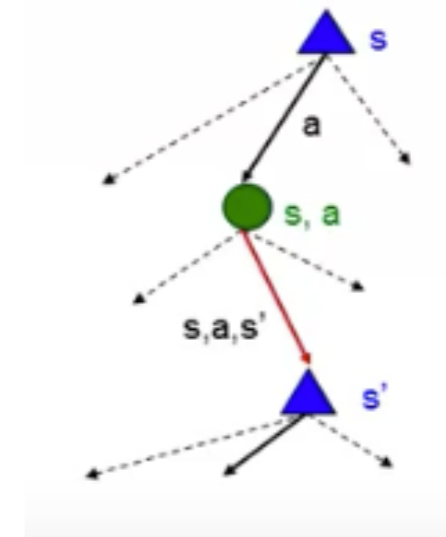
$v_*(s)$ = očakávaný úžitok zo stavu s nasledujúc opt. strat.profil

Hodnota (úžitok) stavu $q(s, a)$

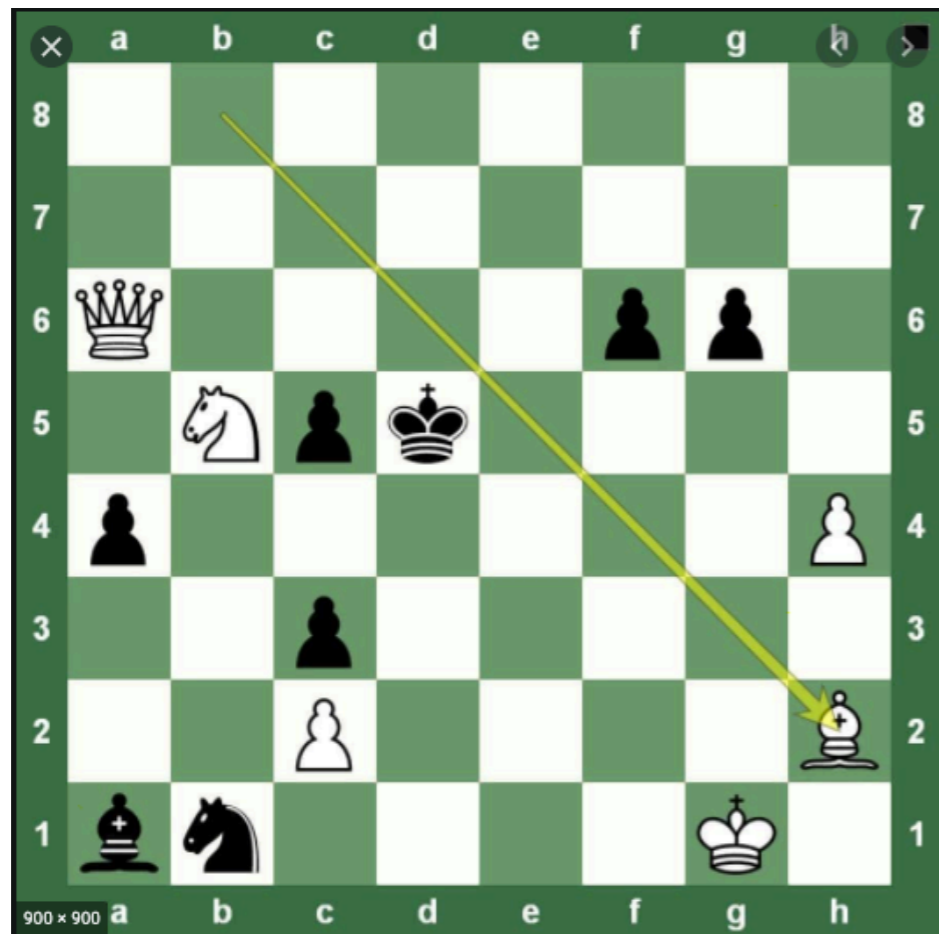
$Q_*(s, a)$ = očakávaný úžitok zo stavu s po vykonaní akcie a nasledujúc opt. strat.profil

Optimálny strategický profil $\pi^*(s)$

$\pi_*(s)$ = optimálna akcia zo stavu s , t.j. $V_*(s) \leftarrow \max Q_*(s, a)$



Hodnota stavu, q stav



Výhra 100 bodov, prehra 100 bodov

π náhodný strategický profil

Biely hráč $v^*(s) = 99$

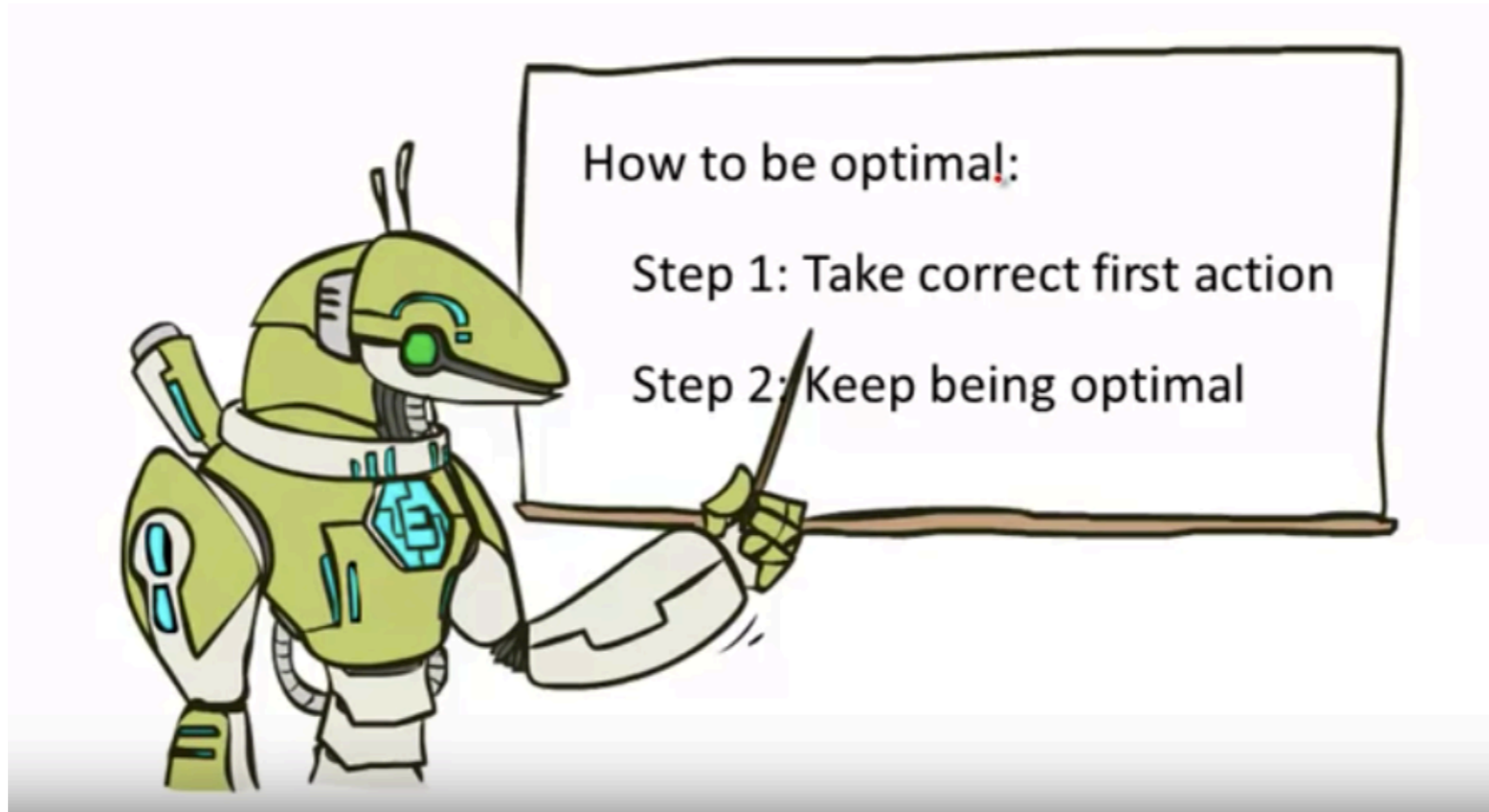
Čierny hráč $v^*(s) = 0,01$

Biely hráč $v_{\pi}(s) = 5$

Čierny hráč $v_{\pi}(s) = 2$



Hodnota stavu, q stav



Hodnota stavu, q stav

Hodnota stavu $v_\pi(s)$ pre ľubovoľný strateg. profil π

$$v_\pi(s) \doteq \mathbb{E}_\pi [G_t | S_t = s] = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right], \forall s \in \mathcal{S}$$

q hodnota pre stav s a akciu a , $q(s,a)$ pre strat. profil π

$$q_\pi(s, a) \doteq \mathbb{E}_\pi [G_t | S_t = s, A_t = a] = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

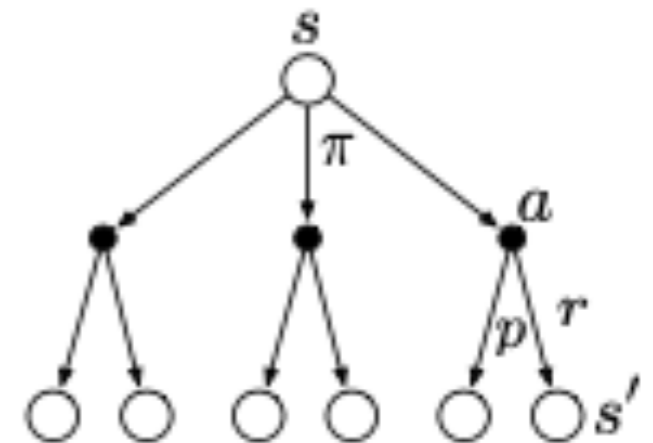


Rekurzívna implementácia výpočtu $v(s)$

Bellmanove rovnice (Bellman equations)

$$\begin{aligned}v_{\pi}(s) &\doteq \mathbb{E}_{\pi} [G_t | S_t = s] \\&= \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s] \\&= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \mathbb{E}_{\pi} [G_{t+1} | S_{t+1} = s']] \\&= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')], \quad \forall s \in \mathcal{S}\end{aligned}$$

Uvažujme N stavov, potom riešenie môžeme získať riešením N lineárnych rovníc



Nie je riešiteľné, potrebujeme vedieť pravú stranu rovnice. Využijeme aproximatívny prístup



Rekurzívna implementácia výpočtu $v(s)$

Probability to take action a from state s following the our policy
(in our case 0.25 since we follow a random policy and we have 4 actions)

multiplied by the sum of:
the reward r (-1 in our case) plus
the expected value of end state s'
multiplied by a discounting factor
gamma

$$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

Value function

Sum for all actions

Sum for all end states s' and reward r

probability to end up in state s' and receive reward r starting in state s and picking action a
(in our case 1, because actions are deterministic and the reward is always -1)



Policy evaluation (predikcia kvality str. profilu)

Motivácia: určenie kvality (hodnoty, resp. využitím hodnotovej funkcie $v_\pi(s), \forall s$) strategického profilu π

V literatúre známe ako *prediction problem* - príklad **policy evaluation**

Control problem - snaha iteratívnym spôsobom zlepšiť strategický profil
- príklad: **policy improvement, value iteration**



Policy evaluation (predikcia kvality str. profilu)

Úloha 1: Výpočet hodnôt stavov pre ľubovoľný strategický profil π

V literatúre známe ako POLICY EVALUATION, ktorý rieši tzv. prediction problém

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi} [G_t | S_t = s] \\ &= \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \mathbb{E}_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s] \\ &= \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')] \end{aligned}$$

Iteratívny prístup: Na začiatku algoritmu priradíme stavom náhodné hodnoty stavov a v ďalších iteráciách hodnoty stavov obnovujeme podľa

$$\begin{aligned} v_{k+1}(s) &\doteq \mathbb{E}_{\pi} [R_{t+1} + \gamma v_k(S_{t+1}) | S_t = s] \\ &= \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_k(s')] \end{aligned}$$



Policy evaluation (predikcia kvality str. profilu)

Iterative Policy Evaluation, for estimating $V \approx v_\pi$

Input π , the policy to be evaluated

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation

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

Loop:

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_a \pi(a|s) \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$



Príklad 3.2 (+ implementácia v Pythone)

Example 3.5: Gridworld Figure 3.2 (left) shows a rectangular gridworld representation of a simple finite MDP. The cells of the grid correspond to the states of the environment. At each cell, four actions are possible: **north**, **south**, **east**, and **west**, which deterministically cause the agent to move one cell in the respective direction on the grid. Actions that would take the agent off the grid leave its location unchanged, but also result in a reward of -1 . Other actions result in a reward of 0 , except those that move the agent out of the special states **A** and **B**. From state **A**, all four actions yield a reward of $+10$ and take the agent to **A'**. From state **B**, all actions yield a reward of $+5$ and take the agent to **B'**.



Bellman optimality equation (control problem)

Tvrdenie: Hodnotu stavu $v_{\pi^*}(s)$ pri opt. strategickom profile π_* je možné určiť ako akciu a prinášajúcu maximálny očakávaný úžitok danom stave s

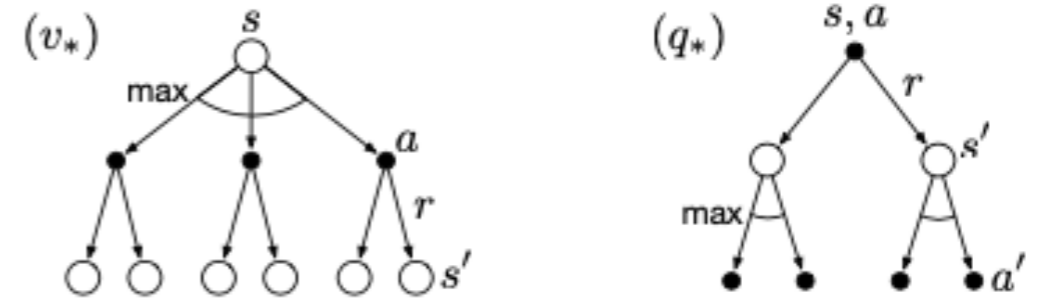
$$v_*(s) = \max_a q_*(s, a)$$

$$q_*(s, a) = \mathbb{E} \left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right]$$

$$= \sum_{s', r} p(s', r \mid s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right]$$

$$q_*(s, a) = \sum_{s'} p(s', r \mid s, a) [R(s, a, s') + \gamma v_*(s')]$$

Uvažujme N stavov, potom riešenie môžeme získať riešením N NElineárnych rovníc



$$v_*(s) = \max_{a \in \mathcal{A}(s)} q_*(s, a)$$

$$= \max_a \mathbb{E}_{\pi_*} [G_t \mid S_t = s, A_t = a]$$

$$= \max_a \mathbb{E}_{\pi_*} [R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a]$$

$$= \max_a \mathbb{E} [R_{t+1} + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a]$$

$$= \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_*(s')]$$



Policy iteration (iterácia strat. profilu)

$$\pi_0 \xrightarrow{\text{E}} v_{\pi_0} \xrightarrow{\text{I}} \pi_1 \xrightarrow{\text{E}} v_{\pi_1} \xrightarrow{\text{I}} \pi_2 \xrightarrow{\text{E}} \dots \xrightarrow{\text{I}} \pi_* \xrightarrow{\text{E}} v_*$$

Control problem

E - Policy evaluation

I - Policy improvement

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Loop:

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement

policy-stable \leftarrow true

For each $s \in \mathcal{S}$:

old-action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

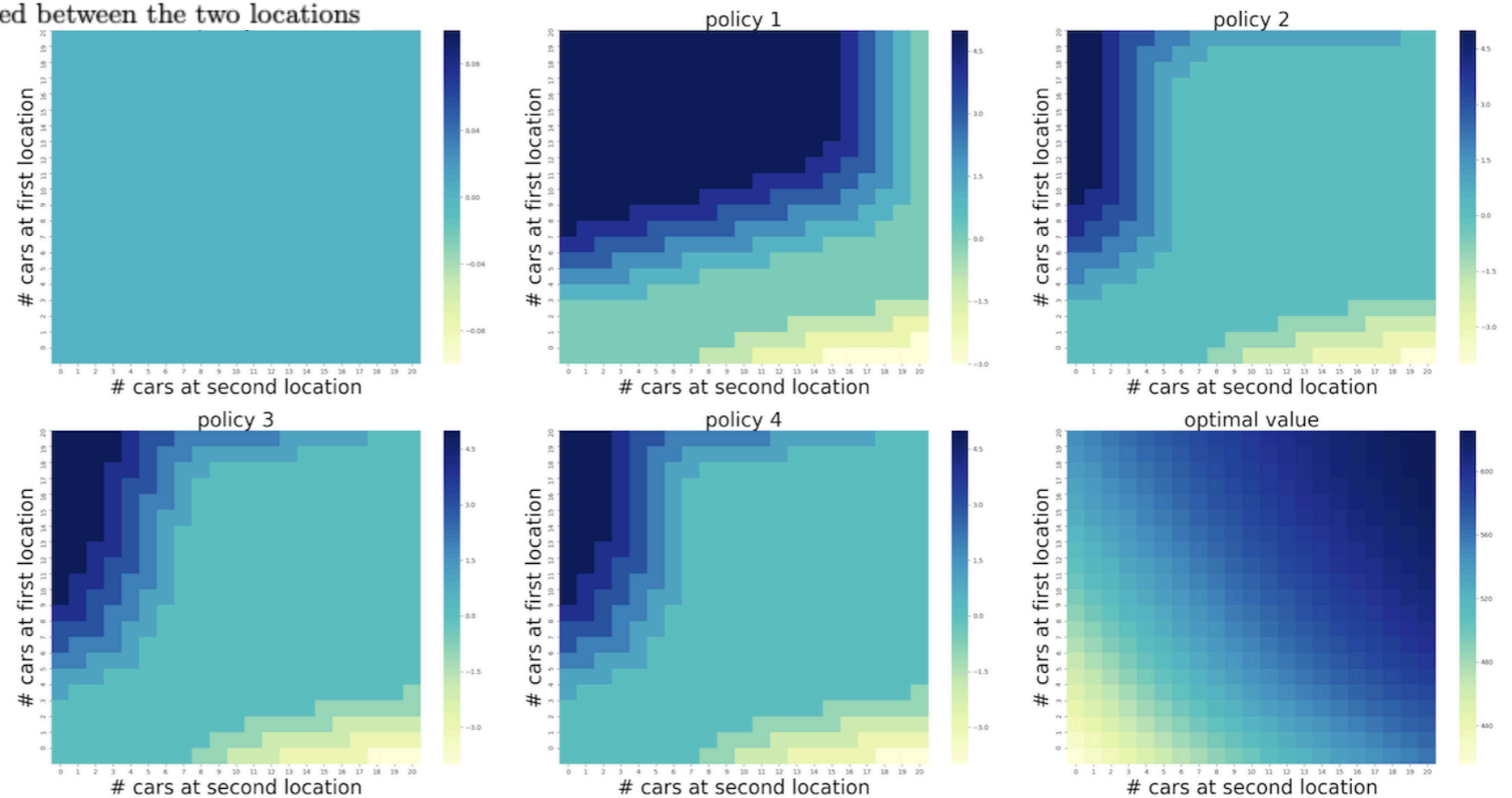
If *old-action* $\neq \pi(s)$, then *policy-stable* \leftarrow false

If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2



Príklad 3.2 (+ implementácia v Pythone)

Example 4.2: Jack's Car Rental Jack manages two locations for a nationwide car rental company. Each day, some number of customers arrive at each location to rent cars. If Jack has a car available, he rents it out and is credited \$10 by the national company. If he is out of cars at that location, then the business is lost. Cars become available for renting the day after they are returned. To help ensure that cars are available where they are needed, Jack can move them between the two locations overnight, at a cost of \$2 per car moved. We assume that the number of cars requested and returned at each location are Poisson random variables, meaning that the probability that the number is n is $\frac{\lambda^n}{n!}e^{-\lambda}$, where λ is the expected number. Suppose λ is 3 and 4 for rental requests at the first and second locations and 3 and 2 for returns. To simplify the problem slightly, we assume that there can be no more than 20 cars at each location (any additional cars are returned to the nationwide company, and thus disappear from the problem) and a maximum of five cars can be moved from one location to the other in one night. We take the discount rate to be $\gamma = 0.9$ and formulate this as a continuing finite MDP, where the time steps are days, the state is the number of cars at each location at the end of the day, and the actions are the net numbers of cars moved between the two locations



Value iteration (hodnotová iterácia)

Motivácia: eliminovať nutnosť volania policy evaluation pri hľadani optim. strateg. profilu.

Iteratívne zvyšovanie hodnôt jednotlivých stavov je možné vykonať bez nutnej evaluácie strategického profilu, priamo v jednom kroku. Garancia konvergenencie je zaručená (neuvádzame ju)

Control problem

$$\begin{aligned}v_{k+1}(s) &\doteq \max_a \mathbb{E} \left[R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s, A_t = a \right] \\ &= \max_a \sum_{s',r} p(s', r \mid s, a) [r + \gamma v_k(s')]\end{aligned}$$



Value iteration (hodnotová iterácia)

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

```
|  $\Delta \leftarrow 0$   
| Loop for each  $s \in \mathcal{S}$ :  
|    $v \leftarrow V(s)$   
|    $V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$   
|    $\Delta \leftarrow \max(\Delta, |v - V(s)|)$ 
```

until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that

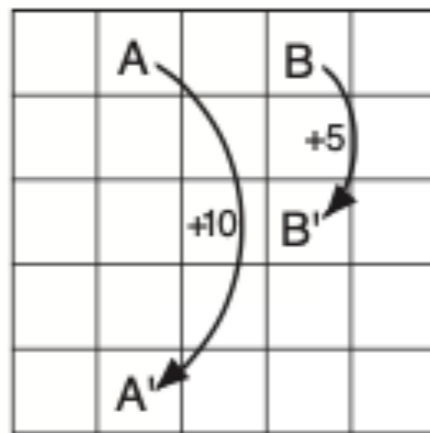
$$\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$$

Control problem



Príklad 3.8 (+ implementácia v Pythone)

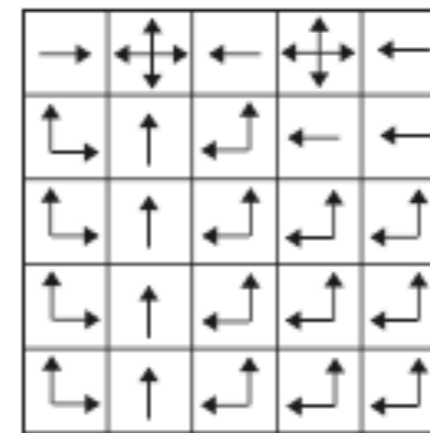
Example 3.8: Solving the Gridworld Suppose we solve the Bellman equation for v_* for the simple grid task introduced in Example 3.5 and shown again in Figure 3.5 (left). Recall that state A is followed by a reward of +10 and transition to state A', while state B is followed by a reward of +5 and transition to state B'. Figure 3.5 (middle) shows the optimal value function, and Figure 3.5 (right) shows the corresponding optimal policies. Where there are multiple arrows in a cell, all of the corresponding actions are optimal.



Gridworld

22.0	24.4	22.0	19.4	17.5
19.8	22.0	19.8	17.8	16.0
17.8	19.8	17.8	16.0	14.4
16.0	17.8	16.0	14.4	13.0
14.4	16.0	14.4	13.0	11.7

v_*

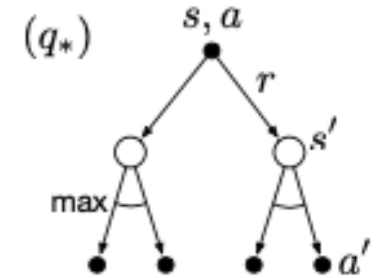
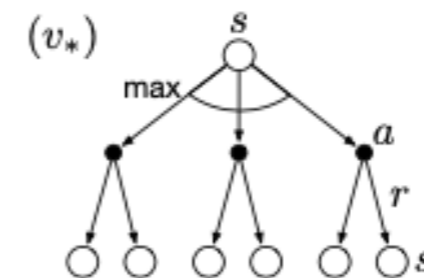


π_*

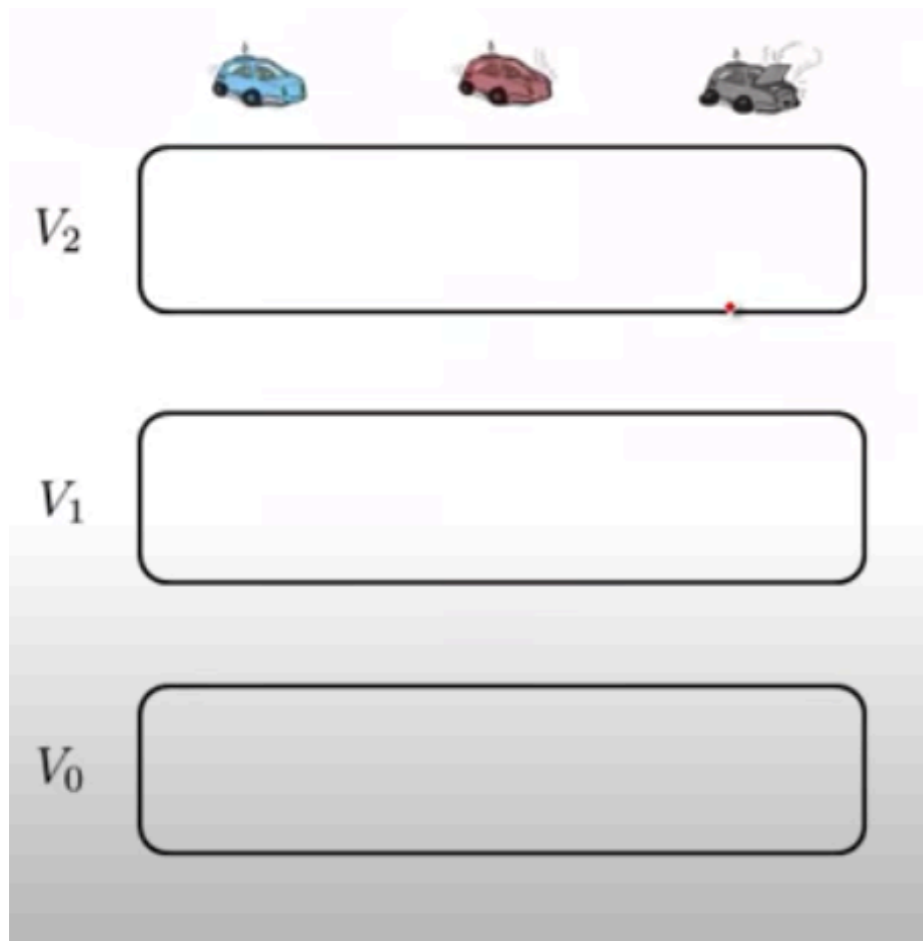
Policy extraction



$$\pi_*(s) = \arg \max_a \sum_{s'} p(s', r | s, a) [R(s, a, s') + \gamma V_*(s')]$$



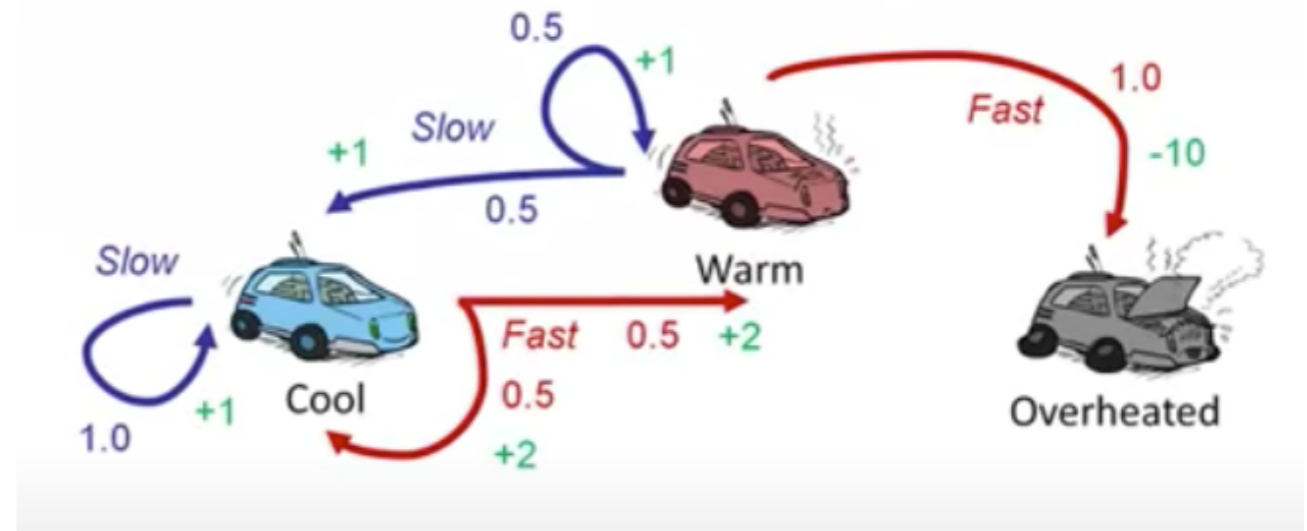
Value iteration (hodnotová iterácia)



V_2

V_1

V_0



$$v_{k+1}(s) = \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma v_k(s')]$$



Value iteration (hodnotová iterácia)

https://cs.stanford.edu/people/karpathy/reinforcejs/gridworld_dp.html

GridWorld: Dynamic Programming Demo

Policy Evaluation (one sweep) Policy Update Toggle Value Iteration Reset

-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00
-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00
-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00
-0.00	-0.00	-0.00	-1.00	-0.00	-0.00	-0.00	-0.00	-0.00	-0.00
-0.01	-0.01	-0.01	-0.01	-0.10	-1.00	-0.00	-0.00	-0.00	-0.00
-0.01	-0.01	-0.01	-0.01	1.00	-0.10	-0.00	-1.00	-0.00	-0.00
-0.01	-0.01	-0.01	-0.01	0.90	0.81	-0.01	-1.01	-0.01	-0.01
-0.01	-0.01	-0.01	-1.01	-0.01	-1.01	-0.01	-0.01	-0.01	-0.01
-0.01	-0.01	-0.01	-0.01	-0.01	-0.01	-0.01	-0.01	-0.01	-0.01
-0.01	-0.01	-0.01	-0.01	-0.01	-0.01	-0.01	-0.01	-0.01	-0.01

Cell rewards (select a cell)



Zhrnutie

☆ Naučili sme sa:

- Výpočet hodnotovej funkcie pre daný strategický profil (**policy evaluation**)
- Nájdenie optimálneho strategického profilu iteratívnym zlepšením hodnotovej funkcie stavov (**value iteration**), resp. implementovaním **policy evaluation** v kombinácií s **policy extraction**

