



# **Strojové učenie II**

prednáška 6 – Aproximácia hodnotových funkcií

Ing. Ján Magyar, PhD.

Katedra kybernetiky a umelej inteligencie

Technická univerzita v Košiciach

2021/2022 letný semester

# Problém škálovania

- najčastejšie škálovanie počtu stavov
  - konečný počet stavov
    - backgammon:  $10^{20}$  stavov
    - go:  $10^{170}$  stavov
  - nekonečný počet stavov – spojité problémy

# Presná reprezentácia hodnotovej funkcie

- každý stav / pár (stav, akcia) má jedinečnú reprezentáciu
  - s rastúcim počtom stavov/akcií sa tabuľka zväčšuje
- určenie hodnoty = vyhládanie príslušného údajov v tabuľke
- problémy
  - príliš veľa miest v pamäti
  - nutnosť uvažovať osobitne každý stav / pár (stav, akcia)

	$s_1$	$s_2$	...	$s_n$
$a_1$		$q(s_2, a_1)$		
...				
$a_m$				

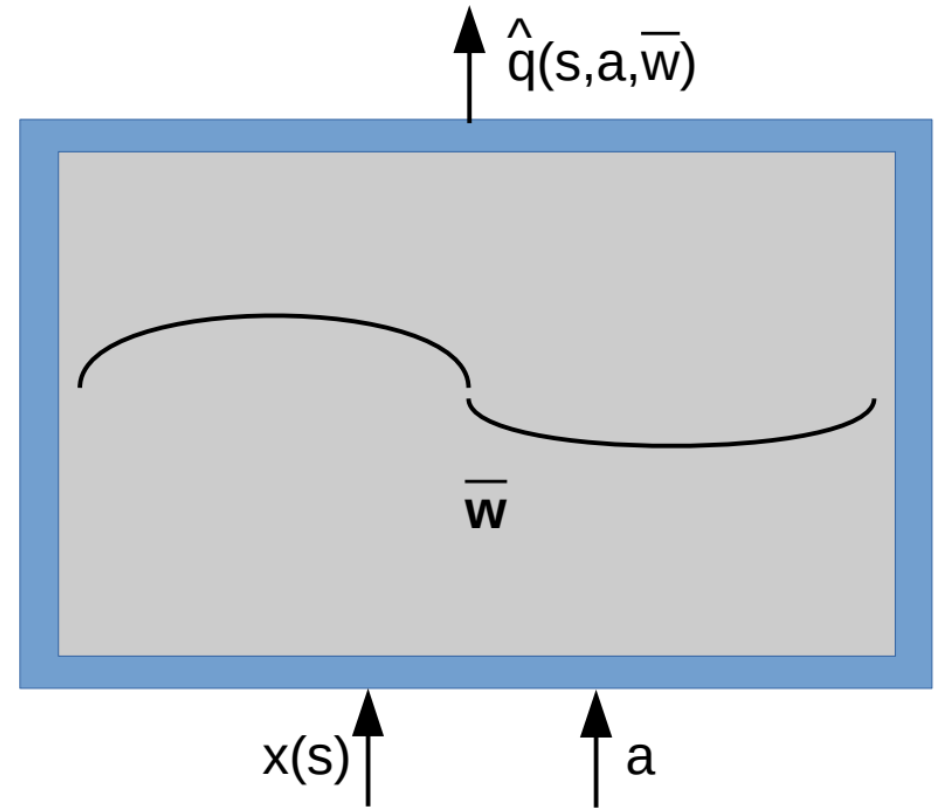
Diagram illustrating a table representing the value function  $q(s, a)$ . The table has states  $s_1, s_2, \dots, s_n$  as columns and actions  $a_1, \dots, a_m$  as rows. The value  $q(s_2, a_1)$  is highlighted in red. Arrows indicate the mapping from state  $s$  and action  $a$  to the value  $q(s, a)$ .

# Aproximátor

- vytvorenie na základe príkladov tvaru  $s \rightarrow v(s)$
- parametrický aproximátor  $\hat{v}(s, \bar{w})$  má tvar vhodne zvolenej parametrickej funkcie
  - lineárny
    - lineárna kombinácia príznakov
  - nelineárny
    - umelá neurónová sieť
    - rozhodovací strom
- pamäťový (neparametrický) aproximátor  $\hat{v}(s, \mathcal{M})$  má tvar množiny príkladov
  - $\hat{v}(s, \mathcal{M}) = \sum_{s' \in \mathcal{M}} k(s, s') v(s')$ 
    - metóda najbližšieho suseda
    - metóda váženého priemeru

# Parametrická aproximácia funkcie

- namiesto enumeračnej tabuľky použitý aproximátor
- stav je reprezentovaný príznakmi
- namiesto hodnôt  $v / q$  sa učí súbor váh  $\bar{w}$ 
  - váh je menej, ako počet stavov
  - hodnoty  $v / q$  sa odvodzujú výpočtom
- zovšeobecnenie skúmaných stavov na neskúmané



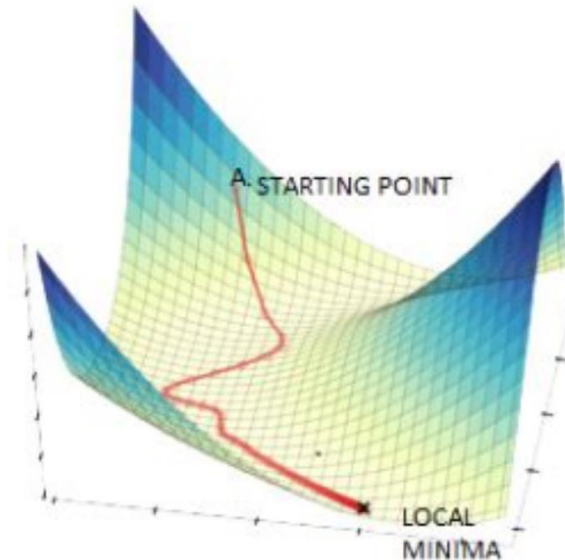
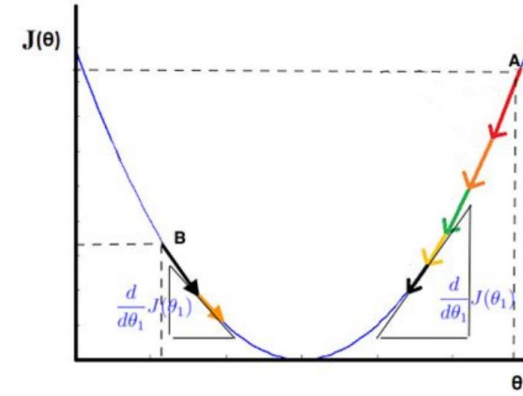
$$\begin{aligned}\hat{v}(s, \bar{w}) &\approx v_{\pi}(s) \\ \hat{q}(s, a, \bar{w}) &\approx q_{\pi}(s, a)\end{aligned}$$

# Aproximačný rámec

- individuálna aktualizácia:  $v(S_t) \mapsto U_t$ 
  - dynamické programovanie:
$$v(S_t) \mapsto E_{\pi}[R_{t+1} + \gamma \hat{v}(S_{t+1}, \bar{w}_t) | S_t = s]$$
  - Monte Carlo:  $v(S_t) \mapsto G_t$
  - time-difference:  $v(S_t) \mapsto R_{t+1} + \gamma \hat{v}(S_{t+1}, \bar{w}_t)$
- aktualizácia hodnoty stavu musí byť urobená ako aktualizácia aproximátora
  - zmenou váhového vektora
  - aktualizácia zasiahne súčasne aj iné stavy (generalizácia)
- požadovaný posun sa vyjadrí ako príklad žiadaného vstupno-výstupného chovania aproximátora
  - aktualizácia sa realizuje pomocou metódou kontrolovaného učenia, ktorá podporuje
    - inkrementálne učenie
    - nestacionárne príklady

# Gradientová minimalizácia

- $J(\bar{w})$  je diferencovateľná funkcia parametra  $\bar{w}$
- gradient
$$\nabla J(\bar{w}) = \left( \frac{\partial J(\bar{w})}{\partial w_1}, \dots, \frac{\partial J(\bar{w})}{\partial w_n} \right)^T$$
- iteračné hľadanie lokálneho minima funkcie  $J(\bar{w})$
- aktualizácia parametra  $\bar{w}$  v smere gradientu  $\Delta \bar{w} = -\alpha \nabla J(\bar{w})$
- metóda GD



# Aktualizácia aproximátora

- predpokladajme, že chceme aktualizáciu v zmysle  $\hat{v}(S_t, \bar{w}_t) \mapsto U_t$
- chyba aproximátora pre stav  $S_t$  je  $\frac{1}{2} (U_t - \hat{v}(S_t, \bar{w}_t))^2$
- aktualizácia váhového vektora

$$\begin{aligned}\bar{w}_{t+1} &= \bar{w}_t - \alpha \nabla \left( \frac{1}{2} (U_t - \hat{v}(S_t, \bar{w}_t))^2 \right) \\ &= \bar{w}_t + \alpha (U_t - \hat{v}(S_t, \bar{w}_t)) \nabla \hat{v}(S_t, \bar{w}_t)\end{aligned}$$

- použitie metódy SGD
  - aktualizácia robená iba podľa jedného príkladu, nie viacerých naraz
  - príklad získavaný z interakcie s prostredím (použitím metód MC alebo TD)



# Algoritmus TD odhadu $v_\pi$

## Semi-gradient TD(0) for estimating $\hat{v} \approx v_\pi$

Input: the policy  $\pi$  to be evaluated

Input: a differentiable function  $\hat{v} : \mathcal{S}^+ \times \mathbb{R}^d \rightarrow \mathbb{R}$  such that  $\hat{v}(\text{terminal}, \cdot) = 0$

Algorithm parameter: step size  $\alpha > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

        Choose  $A \sim \pi(\cdot | S)$

        Take action  $A$ , observe  $R, S'$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})] \nabla \hat{v}(S, \mathbf{w})$

$S \leftarrow S'$

    until  $S$  is terminal

# Konvergenca aproximátora

- neminimalizujeme chybovú funkciu
  - iba malý posun smerom k minimu chybovej funkcie
    - minimalizácia chyby pre nejaký stav by znamenala zväčšenie chyby pre iné stavy
  - potrebné pre vybalansovanie chýb pre rôzne stavy
- cieľová hodnota  $U_t$  pre stav  $S_t$  nie je správna hodnota  $v_\pi(S_t)$  ale iba jej náhodný odhad
  - ak  $E[U_t | S_t = s] = v_\pi(S_t)$ , tak  $\bar{w}_t$  konverguje k lokálnemu optimu ak  $\alpha$  sa postupne znižuje tak, že platia vzťahy
$$\sum_{n=1}^{\infty} \alpha_n = \infty \text{ a } \sum_{n=1}^{\infty} \alpha_n^2 < \infty$$
  - ak  $U_t$  závisí na  $\bar{w}$  (lebo závisí na odhade  $\hat{v}(S_t, \bar{w})$ ), tak konvergenca nie je garantovaná
    - semi-gradientové metódy

# Lineárny aproximátor

- aproximátor má tvar lineárnej kombinácie príznačkov reprezentujúcich stav

$$\hat{v}(s, \bar{w}) = \bar{w}^T \bar{x}(s) = \sum_{i=1}^n w_i x_i(s)$$

kde  $\bar{x}(s)$  je príznačkový vektor stavu  $s$

- gradient s ohľadom na  $\bar{w}$

$$\nabla \hat{v}(s, \bar{w}) = \bar{x}(s)$$

- lineárna funkcia nemá lokálne extrémym
- aktualizácia parametrov podľa

$$\bar{w}_{t+1} = \bar{w}_t + \alpha (U_t - \bar{w}^T \bar{x}(S_t)) \bar{x}(S_t)$$

# TD konvergencia lineárneho aproximátora

$$\begin{aligned}\bar{w}_{t+1} &= \bar{w}_t + \alpha(R_{t+1} + \gamma \bar{w}^T \bar{x}(S_{t+1}) - \bar{w}^T \bar{x}(S_t)) \bar{x}(S_t) \\ &= \bar{w}_t + \alpha(R_{t+1} \bar{x}(S_t) - \bar{x}(S_t)(\bar{x}(S_t) - \gamma \bar{x}(S_{t+1}))^T \bar{w}_t)\end{aligned}$$

$$E[\bar{w}_{t+1} | \bar{w}_t] = \bar{w}_t + \alpha(E[R_{t+1} \bar{x}(S_t)] - E[\bar{x}(S_t)(\bar{x}(S_t) - \gamma \bar{x}(S_{t+1}))^T] \bar{w}_t)$$

$$E[\bar{w}_{t+1} | \bar{w}_t] = \bar{w}_t$$

$$0 = E[R_{t+1} \bar{x}(S_t)] - E[\bar{x}(S_t)(\bar{x}(S_t) - \gamma \bar{x}(S_{t+1}))^T] \bar{w}_{TD}$$

$$\bar{w}_{TD} = (E[\bar{x}(S_t)(\bar{x}(S_t) - \gamma \bar{x}(S_{t+1}))^T])^{-1} E[R_{t+1} \bar{x}(S_t)]$$

kde  $\bar{w}_{TD}$  je fixný bod

doporučené nastavenie:

$$\alpha = 1/(\tau E[\bar{x}^T \bar{x}])$$

# Epizodické učenie politiky

- aproximovaná funkcia  $\hat{q}(S_t, A_t, \bar{w})$ 
  - aproximátor  $x(s, a) \mapsto \hat{q}(s, a, \bar{w})$  na základe  $\bar{w}$
- aktualizácia v zmysle  $\hat{q}(S_t, A_t, \bar{w}) \mapsto U_t$
- chyba aproximátora pre dvojicu  $S_t, A_t$  je  $\frac{1}{2} (U_t - \hat{q}(S_t, A_t, \bar{w}_t))^2$
- aktualizácia váhového vektora
$$\bar{w}_{t+1} = \bar{w}_t - \alpha (U_t - \hat{q}(S_t, A_t, \bar{w}_t)) \nabla \hat{q}(S_t, A_t, \bar{w}_t)$$
- aktualizácia váhového vektora pre Sarsu
$$\bar{w}_{t+1} = \bar{w}_t - \alpha (R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, \bar{w}_t) - \hat{q}(S_t, A_t, \bar{w}_t)) \nabla \hat{q}(S_t, A_t, \bar{w}_t)$$
- zlepšenie politiky  $A_t^* = \operatorname{argmax}_a \hat{q}(S_t, a, \bar{w}_{t-1})$

# Algorithmus semi-gradient Sarsa

## Episodic Semi-gradient Sarsa for Estimating $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step size  $\alpha > 0$ , small  $\varepsilon > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Loop for each episode:

$S, A \leftarrow$  initial state and action of episode (e.g.,  $\varepsilon$ -greedy)

    Loop for each step of episode:

        Take action  $A$ , observe  $R, S'$

        If  $S'$  is terminal:

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$

            Go to next episode

        Choose  $A'$  as a function of  $\hat{q}(S', \cdot, \mathbf{w})$  (e.g.,  $\varepsilon$ -greedy)

$\mathbf{w} \leftarrow \mathbf{w} + \alpha [R + \gamma \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})] \nabla \hat{q}(S, A, \mathbf{w})$

$S \leftarrow S'$

$A \leftarrow A'$

# Kontinuálne učenie politiky

- prístup založený na  $\gamma$  je problematický
- použitie aproximácie hodnotovej funkcie spôsobuje
  - zmena politiky zlepšujúca diskontovanú hodnotu nejakého stavu negarantuje zlepšenie politiky ako celku
  - **teoréma zlepšovania politiky neplatí**
  - $\epsilon$ -greedifikácia môže niekedy vyústiť do menej kvalitnej politiky
- diskontný prístup sa nahrádza prístupom založeným na priemernej odmene

# Priemerná odmena

- priemerná odmena pri použití politiky  $\pi$

$$\begin{aligned} r(\pi) &= \lim_{h \rightarrow \infty} \frac{1}{h} \sum_{t=1}^h E[R_t | S_0, A_{0:t-1} \sim \pi] \\ &= \sum_s \mu_\pi(s) \sum_a \pi(a|s) \sum_r r \sum_{s'} p(s', r | s, a) \\ \mu_\pi(s) &= \lim_{t \rightarrow \infty} P[S_t = s | A_{0:t-1} \sim \pi] \end{aligned}$$

- predpokladá sa ergodický MDP
  - z ľubovoľného stavu možno prejsť do ľubovoľného stavu (nemusí bezprostredne v jednom kroku)



# Priemerná odmena – použitie

- $r(\pi)$  reprezentuje kvalitu politiky  $\pi$
- zotriedenie politík podľa dosiahnutej  $r(\pi)$
- za optimálnu politiku bude považovaná tá, ktorá dosahuje maximálne  $r(\pi)$

# TD pre kontinuálne úlohy

- diferenčná odmena  $G_t = R_{t+1} - r(\pi) + R_{t+2} - r(\pi) + \dots$
- diferenčné hodnotové funkcie  
 $v_\pi(s) = E_\pi[G_t | S_t = s], q_\pi(s, a) = E_\pi[G_t | S_t = s, A_t = a]$
- diferenčná forma TD chyby
$$\delta_t = (R_{t+1} - \bar{R}_t + \hat{v}(S_{t+1}, \bar{w}_t)) - \hat{v}(S_t, \bar{w}_t)$$
$$\delta_t = (R_{t+1} - \bar{R}_t + \hat{q}(S_{t+1}, A_{t+1}, \bar{w}_t)) - \hat{q}(S_t, A_t, \bar{w}_t)$$
kde  $\bar{R}_t$  je odhad  $r(\pi)$  v čase  $t$
- aktualizácia parametrického vektora aproximátora
$$\bar{w}_{t+1} = \bar{w}_t + \alpha \delta_t \nabla \hat{q}(S_t, A_t, \bar{w}_t)$$

# Algoritmus diferencálna Sarsa

Differential semi-gradient Sarsa for estimating  $\hat{q} \approx q_*$

Input: a differentiable action-value function parameterization  $\hat{q} : \mathcal{S} \times \mathcal{A} \times \mathbb{R}^d \rightarrow \mathbb{R}$

Algorithm parameters: step sizes  $\alpha, \beta > 0$

Initialize value-function weights  $\mathbf{w} \in \mathbb{R}^d$  arbitrarily (e.g.,  $\mathbf{w} = \mathbf{0}$ )

Initialize average reward estimate  $\bar{R} \in \mathbb{R}$  arbitrarily (e.g.,  $\bar{R} = 0$ )

Initialize state  $S$ , and action  $A$

Loop for each step:

Take action  $A$ , observe  $R, S'$

Choose  $A'$  as a function of  $\hat{q}(S', \cdot, \mathbf{w})$  (e.g.,  $\varepsilon$ -greedy)

$\delta \leftarrow R - \bar{R} + \hat{q}(S', A', \mathbf{w}) - \hat{q}(S, A, \mathbf{w})$

$\bar{R} \leftarrow \bar{R} + \beta \delta$

$\mathbf{w} \leftarrow \mathbf{w} + \alpha \delta \nabla \hat{q}(S, A, \mathbf{w})$

$S \leftarrow S'$

$A \leftarrow A'$

# Triáda nestability

- aproximácia hodnotovej funkcie
  - zovšeobecnenie v priestore stavov (neuvažovanie stavov oddelene, ich vzájomné ovplyvňovanie)
- bootstrapping
  - aktualizácia hodnôt založená na aktuálnom stave týchto hodnôt
- off-policy učenie
  - rozdielnosť medzi cieľovou a exploračnou politikou

# Dávkový vs inkrementálny prístup

- minimalizácia chyby pre jeden stav alebo pár (stav, akcia) – inkrementálny prístup
  - SGD
    - oprava pre jeden stav znamená pokazenie pre iný stav
    - preto pohyb iba malý kúsok v smere opravy
    - skúsenosť sa po jednorazovom použití zahodí
- minimalizácia chyby pre viac stavov alebo párov (stav, akcia) – dávkový prístup
  - LS (least squares) algoritmy
  - SGD + opakované prehrávanie
    - skúsenosť sa nezahadzuje ale sa pamätá
    - z pamätnej skúsenosti sa vyberá vzorka (dávka)
    - skúsenosť môže byť opakovane vzorkovaná do rôznych dávok