

Értékelőfüggvény közelítése a megerősítéses tanulásban

Kaczúr Flórián

Témavezető: Csáji Balázs Csanád

2023. június 2.

- $X = \{1, \dots, n\}$ állapottér

Markov-döntési folyamatok

- $X = \{1, \dots, n\}$ állapottér
- A véges akciótér, $\mathcal{A} : X \rightarrow \mathcal{P}(A)$

Markov-döntési folyamatok

- $X = \{1, \dots, n\}$ állapottér
- A véges akciótér, $\mathcal{A} : X \rightarrow \mathcal{P}(A)$
- $p_{ij}(a)$ átmenetvalószínűség

Markov-döntési folyamatok

- $X = \{1, \dots, n\}$ állapottér
- A véges akciótér, $\mathcal{A} : X \rightarrow \mathcal{P}(A)$
- $p_{ij}(a)$ átmenetvalószínűség
- $g(i, j)$ költség/jutalom

Markov-döntési folyamatok

- $X = \{1, \dots, n\}$ állapottér
- A véges akciótér, $\mathcal{A} : X \rightarrow \mathcal{P}(A)$
- $p_{ij}(a)$ átmenetvalószínűség
- $g(i, j)$ költség/jutalom

Politika: $\mu : X \rightarrow \Delta(A)$

Markov-döntési folyamatok

- $X = \{1, \dots, n\}$ állapottér
- A véges akciótér, $\mathcal{A} : X \rightarrow \mathcal{P}(A)$
- $p_{ij}(a)$ átmenetvalószínűség
- $g(i, j)$ költség/jutalom

Politika: $\mu : X \rightarrow \Delta(A)$

Átmenetmátrix: $p_{ij} = \sum_{a \in \mathcal{A}(i)} p_{ij}(a) \cdot \mu(i, a)$, feltétel: *irreducibilitás*

Markov-döntési folyamatok

- $X = \{1, \dots, n\}$ állapottér
- A véges akciótér, $\mathcal{A} : X \rightarrow \mathcal{P}(A)$
- $p_{ij}(a)$ átmenetvalószínűség
- $g(i, j)$ költség/jutalom

Politika: $\mu : X \rightarrow \Delta(A)$

Átmenetmátrix: $p_{ij} = \sum_{a \in \mathcal{A}(i)} p_{ij}(a) \cdot \mu(i, a)$, feltétel: *irreducibilitás*

Értékelőfüggvény:

$$J^\mu(i) = \limsup_{N \rightarrow \infty} \mathbb{E} \left\{ \sum_{k=0}^{N-1} \alpha^k g(i_k, i_{k+1}) \mid i_0 = i \right\}$$

Továbbiakban: J^μ közelítése.

Továbbiakban: J^μ közelítése.

Gond: P -t és g -t nem ismerjük.

Továbbiakban: J^μ közelítése.

Gond: P -t és g -t nem ismerjük.

Számításigényes módszerek.

$$\text{span}\{\Phi_1, \dots, \Phi_m\} \leq \mathbb{R}^n$$

$$\text{span}\{\Phi_1, \dots, \Phi_m\} \leq \mathbb{R}^n$$

Fixpont-egyenlet:

$$\Pi_{\xi}^{\Phi} T^{(\lambda)}(\Phi r_{\lambda}^*) = \Phi r_{\lambda}^*,$$

$$\text{span}\{\Phi_1, \dots, \Phi_m\} \leq \mathbb{R}^n$$

Fixpont-egyenlet:

$$\Pi_{\xi}^{\Phi} T^{(\lambda)}(\Phi r_{\lambda}^*) = \Phi r_{\lambda}^*,$$

ahol

- $P^T \xi = \xi$
- $T^{(\lambda)} = (1 - \lambda) \sum_{\ell=0}^{\infty} \lambda^{\ell} T^{\ell+1}$, $0 \leq \lambda < 1$
- $(TJ)(i) = \sum_{j=1}^n p_{ij} \cdot (g(i, j) + \alpha J(j))$, $i = 1, \dots, n$
- $\Phi \in \mathbb{R}^{n \times m}$.

Állítás

$\Pi_{\xi}^{\Phi} T^{(\lambda)}$ *kontrakció.* [1]

Állítás

$\Pi_{\xi}^{\Phi} T^{(\lambda)}$ kontrakció. [1]

Állítás

$$\|J_{\mu} - \Phi r_{\lambda}^*\|_{\xi} \leq \left(\frac{1}{\sqrt{1 - \alpha_{\lambda}^2}} \right) \|J_{\mu} - \Pi_{\xi}^{\Phi} J_{\mu}\|_{\xi}, \text{ ahol } \alpha_{\lambda} = \frac{\alpha(1 - \lambda)}{1 - \alpha\lambda}.$$

[1]

$$\Pi_{\xi}^{\Phi} T^{(\lambda)}(\Phi r_{\lambda}^*) = \Phi r_{\lambda}^*$$

$$\Pi_{\xi}^{\Phi} T^{(\lambda)}(\Phi r_{\lambda}^*) = \Phi r_{\lambda}^* \quad \rightarrow \quad Ar_{\lambda}^* = b$$

$$\begin{array}{lcl} \Pi_{\xi}^{\Phi} T^{(\lambda)}(\Phi r_{\lambda}^*) = \Phi r_{\lambda}^* & \rightarrow & A r_{\lambda}^* = b \\ & & \downarrow \quad \downarrow \\ & & A_k \quad b_k \end{array}$$

$$\Pi_{\xi}^{\Phi} T^{(\lambda)}(\Phi r_{\lambda}^*) = \Phi r_{\lambda}^* \quad \rightarrow \quad A r_{\lambda}^* = b$$
$$\downarrow \quad \downarrow$$
$$A_k \quad b_k$$

① LSTD(λ): $\hat{r}_k = A_k^{-1} b_k$

$$\begin{array}{ccc} \Pi_{\xi}^{\Phi} T^{(\lambda)}(\Phi r_{\lambda}^*) = \Phi r_{\lambda}^* & \rightarrow & A r_{\lambda}^* = b \\ & & \downarrow \quad \downarrow \\ & & A_k \quad b_k \end{array}$$

① LSTD(λ): $\hat{r}_k = A_k^{-1} b_k$

② LSPE(λ): $r_{k+1} = r_k - \gamma G_k (A_k r_k - b_k)$

$$\Pi_{\xi}^{\Phi} T^{(\lambda)}(\Phi r_{\lambda}^*) = \Phi r_{\lambda}^* \quad \rightarrow \quad A r_{\lambda}^* = b$$

$$\downarrow \quad \downarrow$$

$$A_k \quad b_k$$

① LSTD(λ): $\hat{r}_k = A_k^{-1} b_k$

② LSPE(λ): $r_{k+1} = r_k - \gamma G_k (A_k r_k - b_k)$

③ TD(λ): $r_{k+1} = r_k - \gamma_k z_k q_{k,k}$, ahol

$$q_{k,k} = \Phi(i_k)^T r_k - \alpha \Phi(i_{k+1})^T r_k - g(i_k, i_{k+1}),$$

$$z_k = \sum_{h=0}^k (\alpha \lambda)^{k-h} \Phi(i_h).$$

- kezdeti ζ eloszlás

- kezdeti ζ eloszlás

- $$c_{l,s}(r_k) = \alpha^{N_s-l} \Phi(i_{N_s-1,s})^T r_k + \sum_{h=l}^{N_s-1} \alpha^{h-l} g(i_{h,s}, i_{h+1,s})$$

- kezdeti ζ eloszlás

- $$c_{\ell,s}(r_k) = \alpha^{N_s-\ell} \Phi(i_{N_s-1,s})^T r_k + \sum_{h=\ell}^{N_s-1} \alpha^{h-\ell} g(i_{h,s}, i_{h+1,s})$$

- $$r_{k+1} = \arg \min_{r \in \mathbb{R}^m} \sum_{s=0}^T \sum_{\ell=0}^{N_s-1} (\Phi(i_{\ell,s})^T r - c_{\ell,s}(r_k))^2$$

- kezdeti ζ eloszlás

- $$c_{\ell,s}(r_k) = \alpha^{N_s-\ell} \Phi(i_{N_s-1,s})^T r_k + \sum_{h=\ell}^{N_s-1} \alpha^{h-\ell} g(i_{h,s}, i_{h+1,s})$$

- $$r_{k+1} = \arg \min_{r \in \mathbb{R}^m} \sum_{s=0}^T \sum_{\ell=0}^{N_s-1} (\Phi(i_{\ell,s})^T r - c_{\ell,s}(r_k))^2$$

- $$\Phi r = \Pi_{\zeta}^{\Phi} \mathcal{T}^{(\lambda)} \Phi r.$$



Bertsekas, Dimitri P. "Dynamic programming and optimal control 3rd edition, volume ii." Belmont, MA: Athena Scientific (2011).