

Úkol 7: Chceme provést MDP, chceme vytvořit policy.

- 1) V init se jednou vytvoří policy, pak už se zavolej nepotřebuje.
- 2) Chceme to počítat prostřednictvím "utility" a optimální policy.

Policy iteration je těžší implementovat, ale počítá rychleji.



C > čas je důležitý v RecastExu

Value iteration má méně času už v některých nejdále.

Vč sloučenosti stále implementovat precompute-probability-policy-trivial

$$\pi(s) = \arg \max_{a \in A} [P(s'|s, a) \cdot u(s')]$$

$$u(s) = R(s) + \max_{a \in A} P(s'|s, a) \cdot u(s')$$

✓ prezentaci!

Value iteration

initialization $u(s) = R(s)$

Bellman update:

$$u_{i+1}(s) = R(s) + \gamma \max_{a \in A} P(s'|s, a) \cdot u_i(s')$$

Co vlastně očekávám počítat?

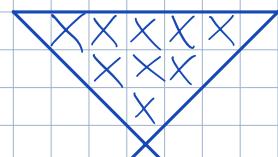
→ úplné mnoho řešení vzdálení stavy.
Tedy $\gamma = 1$

$$R(s) = M[s] + \gamma \min_{a \in A} P(s'|s, a) \cdot u(s')$$

$\{L, R, U, D\}$

Rewards

$$\begin{array}{cccc} 0 & 0 & 0 & +1 \\ 0 & \bullet & 0 & -1 \\ 0 & 0 & 0 & 0 \end{array} \quad \gamma = 0,9 \quad 0,1 \leftarrow \boxed{+1} \rightarrow 0,1$$



$$u_1(0, \epsilon) = 0 + 0,9 \cdot \max \left(0, 0,1 \cdot 1, 0,8 \cdot 1, 0,1 \cdot 1 \right) = 0,9 \cdot 0,8$$

$$u_2(0, \epsilon) = 0 + 0,9 \cdot \max \left(0, 0,1 \cdot 0,72, 0,8 \cdot 0,72, 0,1 \cdot 0,72 \right) = \dots$$

$$U_2(1,2) = 0 + 0,9 \cdot \max \left(0,1 \cdot 0,72, \underbrace{0,8 \cdot 0,72 + 0,1 \cdot (-1)}_{\text{John price receiving money}} \dots \right)$$

John price receiving money
obdrocení

Jak získat optimální akci? $\pi(s) = \arg \max_{a \in A} \left[\sum_{s'} P(s'|s, a) \cdot U(s') \right]$

Policy iteration

Initializujte náhodnou policy $\pi_0(s)$

do

evaluate policy

$$U(s) = R(s) + \gamma \cdot \sum_{s'} P(s'|s, \pi(s)) \cdot U(s')$$

$$\text{if } \max_{a \in A} \sum_{s'} P(s'|s, a) \cdot U(s') > \sum_{s'} P(s'|s, \pi(a)) \cdot U(s)$$

zjistit akci, kterou má výšší hodnotu