

Zpěťovazební učení

Pasivní učení

Policy je fixní, učíme se jenom Utility stavů

- Direct Utility Estimation:
 - Jendoduše necháme agenta, aby provedl několik průchodů prostorem a stavům přiřadíme průměr z utilit
- Adaptive Dynamic Programming
 - To samé jako DUE, ale utility vždy přepočítáme na základě Bellmanových rovnic
 - $U^\pi(s) = R(s) + \gamma \sum_{s'} P(s'|s, \pi s) U(s')$
- Temporal-difference Learning
 - Nemusím přepočítat úplně všechny utility, jenom ty, které se mohly změnit
 - $U^\pi(s) = U^\pi(s) + \alpha(R(s) + \gamma U^\pi(s') - U^\pi(s))$

Aktivní učení

Policy upravujeme podle toho, jak se nám to zrovna líbí. Musíme udržovat explorační a exploatační abychom nedostali greedy agenta. Preferujeme nové stavy před starými.

Q-Learning a SARSA

Najednou máme matici Q, takovou, že $Q(s, a)$ říká jakou hodnotu má provedení a v s

$$U(s) = \max_a Q(s, a)$$

Update learningu provádíme

$$Q(s, a) = Q(s, a) + \alpha(R(s) + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Pro SARSA update vypadá

$$Q(s, a) = Q(s, a) + \alpha(R(s) + \gamma Q(s', a') - Q(s, a))$$