

Az UCT és a MoGo algoritmusokról

Szörényi Balázs

1 UCT: UCB-1 for trees

Algorithm 1 UCT-KIÉRTÉKELŐ($s_0 \in \mathcal{S}$)

```
1: INICIALIZÁLÁS
2: minden  $(s, a)$  párra  $n(s, a) := 0$ ,  $Q_{\text{UCT}}(s, a) := 0$ 
3: TANULÁS (ON-LINE)
4: for  $\text{szim} := 1$  to  $\text{SzimulációSzám}$  do
5:   EPIZÓD GENERÁLÁSA
6:   Generáljunk egy  $s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T$  epizódot (avagy játszmát) a
7:    $\pi_{\text{UCT}} : \mathcal{S} \mapsto \begin{cases} \text{egy véletlen választott lépés az } s\text{-ből még nem próbáltak közül,} & \text{ha van ilyen} \\ \operatorname{argmax}_a \left( Q_{\text{UCT}}(s, a) + c \cdot \sqrt{\frac{n(s)}{n(s, a)}} \right) & \text{egyébként} \end{cases}$ 
8:   stratégiával self-play módban.
9:   UPDATE (MONTE-CARLO)
10:  for  $t := 1$  to  $T$  do
11:     $n(s_t, a_t) := n(s_t, a_t) + 1$ 
12:     $n(s_t) := \sum_a n(s_t, a_t)$ 
13:     $Q_{\text{UCT}}(s_t, a_t) := Q_{\text{UCT}}(s_t, a_t) + (1/n(s_t, a_t))(R_t - Q_{\text{UCT}}(s_t, a_t))$ 
14:     $\{R_t = r_t + r_{t+1} + \dots + r_T = a \text{ } t\text{-edik lépés után szerzett összjuttalom}\}$ 
15:  end for
16: end for
17: TÉNYLEGES LÉPÉS
18: return  $\operatorname{argmax}_a Q_{\text{UCT}}(s_0, a)$ 
```

Pl. $T = 3000$ vagy $T = 70000$.

2 MoGo: Monte-Carlo Go

π_{MoGo} : offline (azaz még éles alkalmazás előtt) kézzel összeállított véletlen stratégia, beépített (emberi) Go háttértudással. (Persze más probléma esetén használható egy ilyen helyett is valamilyen offline módon tanított stratégia.)

Q_{RLGO} : offline (azaz még éles alkalmazás előtt) megtanult cselekvés-érték függvény. Konkrétan: ℓ darab, kézzel összeválogatott bináris $\phi_1, \dots, \phi_\ell : \mathcal{S} \times \mathcal{A} \rightarrow \{0, 1\}$ jellemző megfelelő $(\theta_1, \dots, \theta_\ell)$ súlyozására tanultak rá TD(0) algoritmussal (gradiens módszerrel) self-play-t alkalmazva. A súlyokból a cselekvés-érték függvény $\sigma(\sum_{i=1}^{\ell} \theta_i \cdot \phi_i(s, a))$ képlettel került kiszámításra, ahol $\sigma(y) = 1/(1 + e^{-y})$ a szigmoid függvény.

Algorithm 2 UCT_{RAVE}-KIÉRTÉKELŐ($s_0 \in \mathcal{S}$)

```
1: INICIALIZÁLÁS
2: minden  $(s, a)$  párra  $n(s, a) := m(s, a) := 50$ ,  $Q_{\text{UCT}}(s, a) := Q_{\text{RLGO}}$ 
3: TANULÁS (ON-LINE)
4: for szim := 1 to SzimulációSzám do
5:   EPIZÓD GENERÁLÁSA
6:   Generáljunk egy  $s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_T, a_T, r_T$  epizódot (avagy játszmát) az
7:    $s \mapsto \begin{cases} \pi_{\text{MoGo}} \text{ szerint, ha van olyan lépés, amit az } s\text{-ből még nem próbáltunk} \\ \operatorname{argmax}_a (1 - \beta(s, a)) \left( Q_{\text{UCT}}(s, a) + c \cdot \sqrt{\frac{n(s)}{n(s, a)}} \right) + \beta(s, a) \left( Q_{\text{RAVE}}(s, a) + c \cdot \sqrt{\frac{m(s)}{m(s, a)}} \right) \end{cases}$ 
   egyébként
8:   stratégiával self-play módban, ahol  $\beta(s, a) = \sqrt{k/(k + 3n(s))}$ .
9:   UPDATE (MONTE-CARLO)
10:  for  $t := 1$  to  $T$  do
11:     $n(s_t, a_t) := n(s_t, a_t) + 1$ 
12:     $n(s_t) := \sum_a n(s_t, a_t)$ 
13:     $Q_{\text{UCT}}(s_t, a_t) := Q_{\text{UCT}}(s_t, a_t) + (1/n(s_t, a_t))(R_t - Q_{\text{UCT}}(s_t, a_t))$ 
14:  end for
15:  UPDATE (RAVE: Rapid Action Value Estimation)
16:  {Heurisztika: "a későbbi lépések akár hamarabb is ugyanolyan jók lennének".}
17:  for  $t := 1$  to  $T$  do
18:    for  $\tau := t$  to  $T$  do
19:       $m(s_t, a_\tau) := n(s_t, a_\tau) + 1$ 
20:       $m(s_t) := \sum_a n(s_t, a_\tau)$ 
21:       $Q_{\text{RAVE}}(s_t, a_\tau) := Q_{\text{RAVE}}(s_t, a_\tau) + (1/n(s_t, a_\tau))(R_t - Q_{\text{RAVE}}(s_t, a_\tau))$ 
22:    end for
23:  end for
24: end for
25: TÉNYLEGES LÉPÉS
26: return  $\operatorname{argmax}_a Q_{\text{UCT}}(s_0, a)$  vagy  $Q_{\text{RAVE}}(s_0, a)$ 
```

Pl. $k = 1000$.

References

- [1] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Mach. Learn.*, 47(2-3):235–256, 2002.
- [2] S. Gelly and D. Silver. Combining online and offline knowledge in UCT. In *ICML '07: Proceedings of the 24th international conference on Machine learning*, pages 273–280, New York, NY, USA, 2007. ACM.
- [3] L. Kocsis and C. Szepesvári. Bandit based monte-carlo planning. In *ECML*, pages 282–293, 2006.