Reinforcement Learning for Strategic Board Games with n-Tuple Systems

Wolfgang Konen, Samineh Bagheri and Markus Thill TH Köln – University of Applied Sciences

Connect-4



- Two-player board game
- Players in turn drop their stones into 1 of 7 columns
- Gravity: Stones fall to lowest empty cell in the column Goal: Create a row of four
- connected stones (horizontally, vertically or diagonally) State-Space Complexity: 4.5 · 10¹²

Eligibility Traces



Fig. 3.: Schematic view of different eligibility trace variants.

a) Without eligibility traces: a weight is activated only at isolated time points. Dotted vertical line: random move. b) Resetting traces: eligibility trace with reset on random move.

c) Replacing traces: a new activation replaces the old one.

Summary & Result



Fig. 1.: Connect-4 Position with a win for Yellow.

N-Tuple Systems

- An n-tuple is a sequence of board cells where each cell can have one out
- of P=4 states:
- 0=empty and not reachable, 1=Yellow, 2=Red, 3=empty and reachable

- Example: Board position s₊
 - The system with *m* n-tuples specifies a linear function

3	LUT ₁ LUT ₂					LU	T _m
2 0 2 1	index	weight	index	weight		index	weight
	0	$w_0^{(1)}$	0	$w_0^{(2)}$		0	$w_0^{(m)}$
	1	$w_1^{(1)}$	1	$w_1^{(2)}$		1	$w_1^{(m)}$
abcdefg	:	:	:	÷		÷	:
$T_1 = (a1, a2, b2, c2)$	99	$w_{99}^{(1)}$	102	$w_{102}^{(2)}$		ξ_m	$W_{\xi_m}^{(m)}$
$\xi_1 = 3 \cdot 4^0 + 0 \cdot 4^1 + 2 \cdot 4^2 + 1 \cdot 4^3 = 99$:	÷	:	÷		÷	:
	255	$w_{255}^{(1)}$	255	$W_{255}^{(2)}$		255	$w_{255}^{(m)}$
$T_2 = (d6, d5, e5, e4)$							
$\xi_2 = 2 \cdot 4^0 + 1 \cdot 4^1 + 2 \cdot 4^2 + 1 \cdot 4^3 = 102$	$V(s_t) = \tanh(w_{99}^{(1)} + w_{102}^{(2)} + \dots + w_{\xi_m}^{(m)})$						



Fig. 1.: Initial results for TCL-EXP with eligibility traces, the options replacing traces and resetting traces turned off. For values $\lambda \ge 0.9$ the system breaks down.

Unified Algorithmic Description

Temporal Difference Learning

Goal: learn a value function V(s) by means of Temporal Difference Learning (TDL) as follows:



Step-Size Adaptation

Algorithm 2 General TCL algorithm in pseudo code. The counters A_i and N_i are initialized once at the beginning of the training. In the original algorithm, the individual learning rates α_i are calculated with the identity transfer function g(x) = x, but also other transfer functions may be selected, taking into consideration certain conditions

1: Initialize counters $A_i = 0$ and $N_i = 0 \forall i$	
 Set global learning rate α. 	
 for (every weight index i) do 	
4: $\alpha_i = \begin{cases} g\left(\frac{ N_i }{A_i}\right), & \text{if } A_i \neq 0\\ 1, & \text{otherwise} \end{cases}$	▷ Individual step-size
5: $r_{i,t} = \delta_t e_{i,t}$	▷ recommended weight change
6: $w_i \leftarrow w_i + \alpha \alpha_i r_i$	▷ TD update for each weight

Algorithm 3 Incremental $TD(\lambda)$ algorithm for board games.	Prior to the first game, the weight
ector \vec{w} is initialized with random values. Then the follow	ing algorithm is executed for each
omplete board game. During self-play the player is toggled	between $+1$ (first player) and -1
second player).	
1: Set $REP = true$, if using replacing traces	
2: Set $RES = true$, if resetting elig. traces on random moves	
3: Set the initial state s_0 (usually the empty board) and $p = 1$	
4: Use partially trained weights w from previous games	
5: function TDETRAIN(s_0, w)	. T. (4) - 1 - 1) - (1 - (1 - (1 - (1 - (1 - (1
$e_0 \leftarrow \bigvee_{\vec{w}} w, x(s_0)$	\triangleright Initial eligibility traces
$ (: \text{ for } (t \leftarrow 0; s_t \notin S_{Final}; t \leftarrow t+1, p \leftarrow (-p)) \text{ do} $	Value for
S: $V_{old} \leftarrow w_t, x(s_t)$	\triangleright value for s_t
$f:$ generate randomly $q \in [0, 1]$	
$\begin{array}{ccc} \text{II} & (q < \epsilon) \text{ then} \\ \text{Bendemby colort } e \end{array}$	Ermlanative mare
$\begin{array}{c} \text{I:} \qquad \text{Randomly select } s_{t+1} \\ \text{if } (DEC) \text{ there} \end{array}$	▷ Explorative move
\vec{x} : If (<i>RES</i>) then $\vec{z} \neq 0$	
$e_t \leftarrow 0$	
: else	t Create man
Select after-state s_{t+1} , which maximizes	▷ Greedy move
$7: \qquad p \cdot \begin{cases} R(s_{t+1}), & \text{if } s_{t+1} \in S_{Final} \\ \hline \end{array}$	
$\left(\vec{w}_t, \vec{x}(s_{t+1}) + R(s_{t+1}), \text{ otherwise}\right)$	
end if	
$V(s_{t+1}) \leftarrow \vec{w}_t, \vec{x}(s_{t+1})$	
$\delta_t \leftarrow R(s_{t+1}) + \gamma V(s_{t+1}) - V_{old}$	▷ TD error-signal
: if $(q \ge \epsilon s_{t+1} \in S_{Final})$ then	
$: \qquad \vec{w}_{t+1} \leftarrow \vec{w}_t + \alpha \delta_t \vec{e}_t$	▷ Weight-update
end if	
for (every weight index i) do	\triangleright Update elig. traces
$\Delta e_i \leftarrow \nabla_{w_i} \vec{w}_{t+1}, \vec{x}(s_{t+1})$	
$\int \Delta e_i, \qquad \text{if } x_i(s_{t+1}) \neq 0 \land REP$	
$\gamma \lambda e_{i,t+1} = \left\{ \gamma \lambda e_{i,t} + \Delta e_i, \text{ otherwise} \right\}$	
i: end for	
end for	
end function	\triangleright End of $TD(\lambda)$ self-play algorithm

Experimental Setup

Connect-4 Learning Framework



http://github.com/MarkusThill/Connect-Four

7: $A_i \leftarrow A_i + |r_{i,t}|$ $N_i \leftarrow N_i + r_{i,t}$ 9: end for

> - - PieceLinear Standard TCI

> > 0.25

0.50

0.75

×0.50-

0.25

0.00

10: with a transfer function g(x), 11: the approximation error δ_t and, 12: the eligibility traces $e_{i,t}$.

 \triangleright update accumulating counter A \triangleright update accumulating counter N

All agents are trained solely by self-play

- Exploration rate ϵ : 0.1, discount factor γ : 1.0
- Eligibility traces factor: $0 \le \lambda \le 1$ •
- α_{init} : 0.004 , α_{final} : 0.002
- We used 70 n-tuples all of length 8 and generated by random walks on the board
- For a system with 70×8-tuple, $2 \cdot 70 \cdot 4^8 \approx 6 \cdot 10^6$ weights are created
- Weights are initialized with random values uniformly drawn from $\left[-\frac{\xi}{2},+\frac{\xi}{2}\right]$, with $\xi = 0.001$.
- Eligibility trace variants: [et] for standard eligibility traces ۲ without any further options, [res] for resetting traces, [rep] for replacing traces, and [rr] for resetting and replacing traces.
- Benchmark: Each evaluation point based on 200 matches against a perfect-playing Minimax agent from the empty board.







The TCL algorithm developed by Beal and Smith is an extension of TDL. It has an adjustable learning rate α_i for every weight w_i and a global constant α_{init} . For each weight two counters N_i and A_i accumulate the sum of weight changes and sum of absolute weight changes. If all weight changes have the same sign, then the learning rate stays at its upper bound, otherwise the learning rate will be largely reduced.