

方策勾配法による静的局面評価関数の 強化学習についての一考察

五十嵐治一, 森岡祐一, 山本一将

●提案手法の特徴:

- ①将棋の専門知識を必要しない
- ②強化学習として方策勾配法を採用
- ③「指し手評価の期待値」による行動決定
- ④ボルツマン分布による soft-max探索
- ⑤諸探索との関連が明確

将棋のルールを教え、対局させる
だけでコンピュータはプロ棋士に
勝てるようになるだろうか？

強化学習

- ・バックギャモンでは Neural Net + TD(λ)法により実現 [’94 Tesauro]
- ・チェスではTDLeaf(λ)法の有効性が確認 [’98 Baxter et al.]

方策ベースの強化学習

(マルコフ性のない報酬, 方策でもOK)

●本研究で用いた方策勾配法: [Williams92, 石原&五十嵐04]

目的関数 $E(a(t);s(t),\omega)$ [離散時刻 t , 状態 s , 行動 a]: 方策の知識例, ルール重み $E(a(t);s(t),\omega) = -\omega(s(t),a(t)) [\omega \geq 0]$

Boltzmann分布による確率的方策:

$$\pi(a(t);s(t),\omega) \equiv e^{+E(a(t);s(t),\omega)/T} / \sum_{a \in A} e^{+E(a(t);s(t),\omega)/T} \quad \text{仮定}$$

エピソードごとの報酬期待値を極大化する学習則:

$$\Delta\omega = \varepsilon \cdot r \cdot \sum_{t=0}^{L-1} e_{\omega}(t) \quad [L: \text{エピソード長}, r: \text{報酬}, \varepsilon: \text{学習係数}]$$

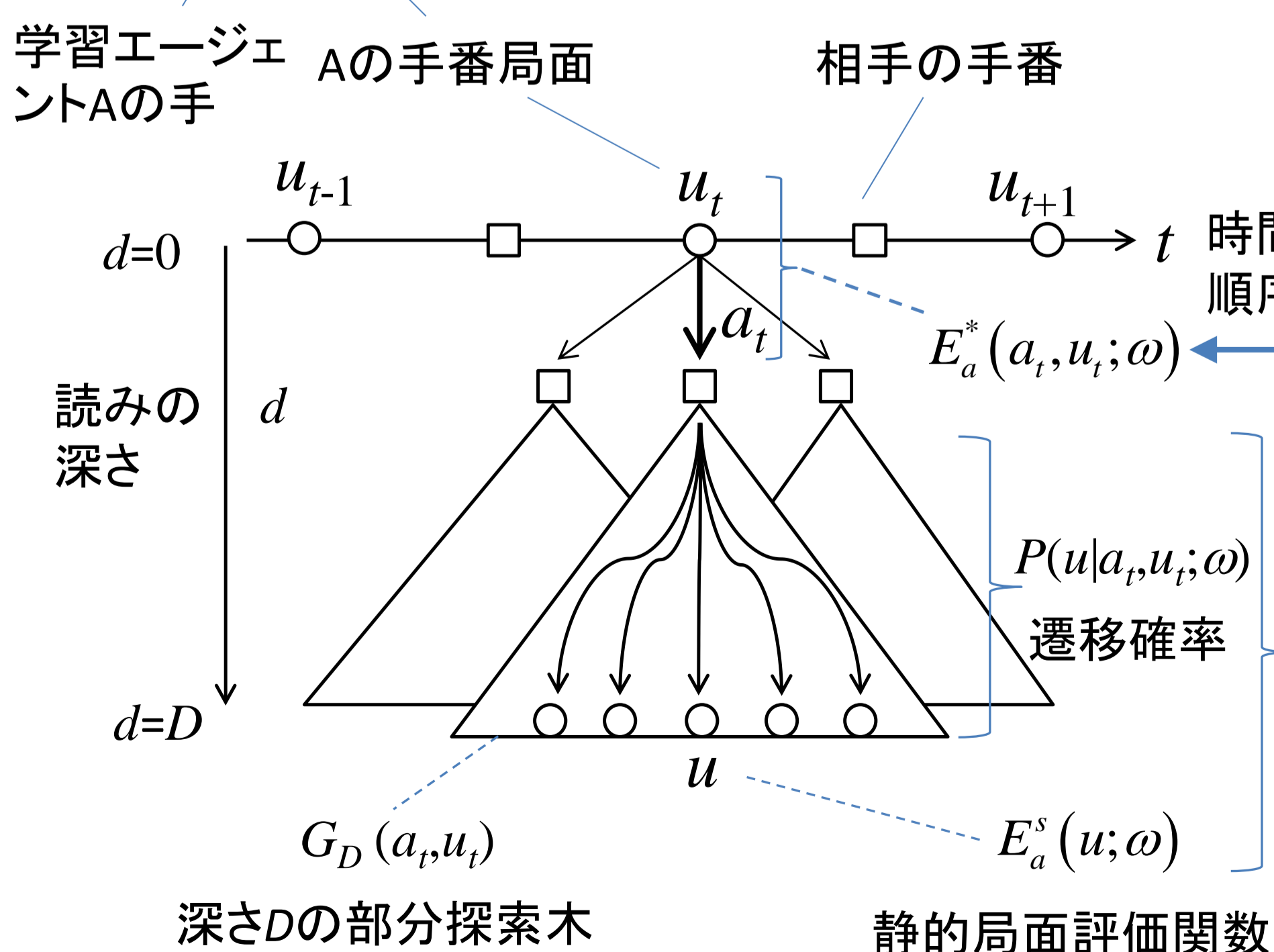
特徴的適正度 $t=0$

$$e_{\omega}(t) \equiv \partial \ln \pi / \partial \omega = \frac{1}{T} \left[\frac{\partial E(a(t);s(t),\omega)}{\partial \omega} - \sum_{a \in A} \frac{\partial E(a; s(t), \omega)}{\partial \omega} \pi(a; s(t), \omega) \right]$$

本提案

「指し手評価の期待値」を目的関数とする

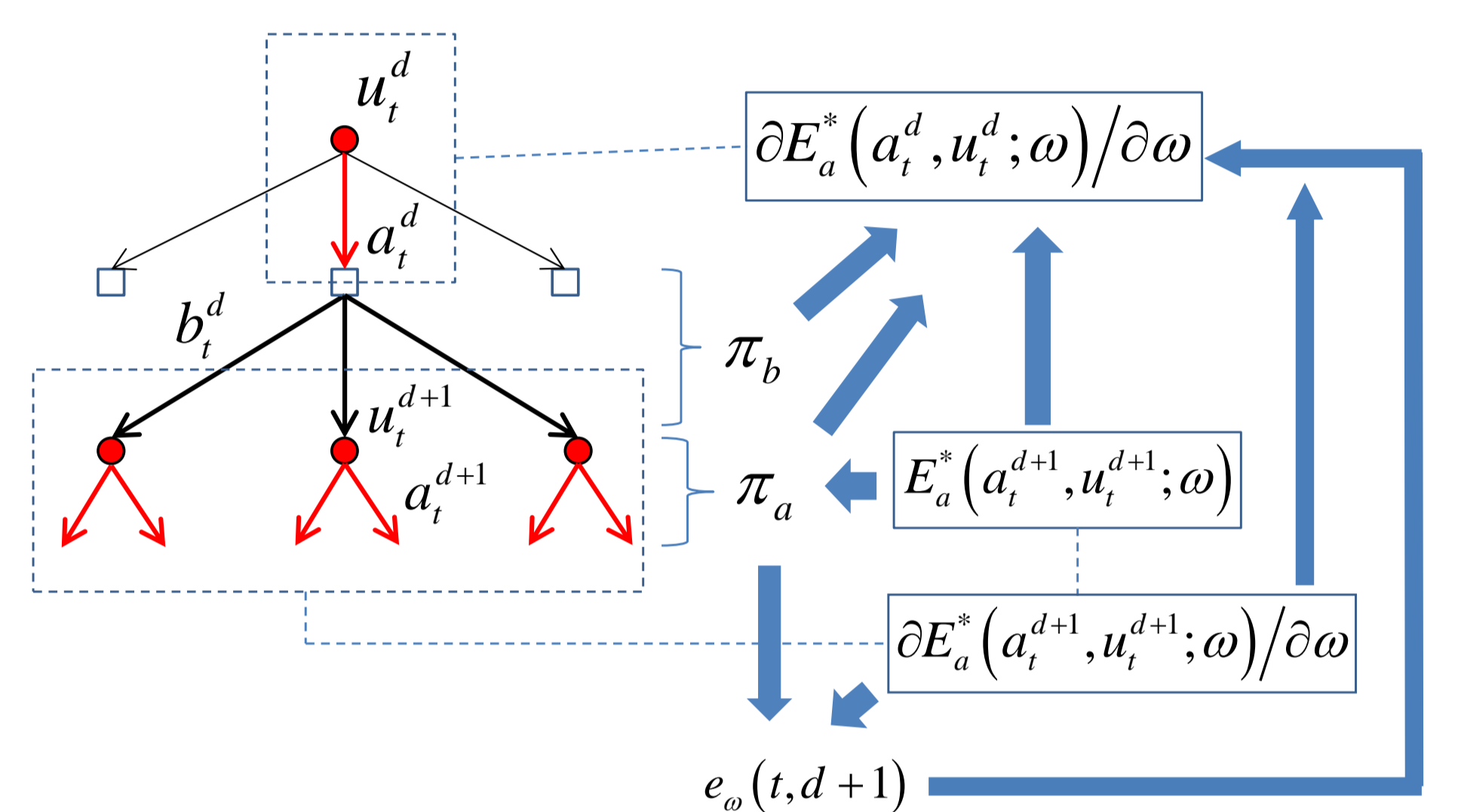
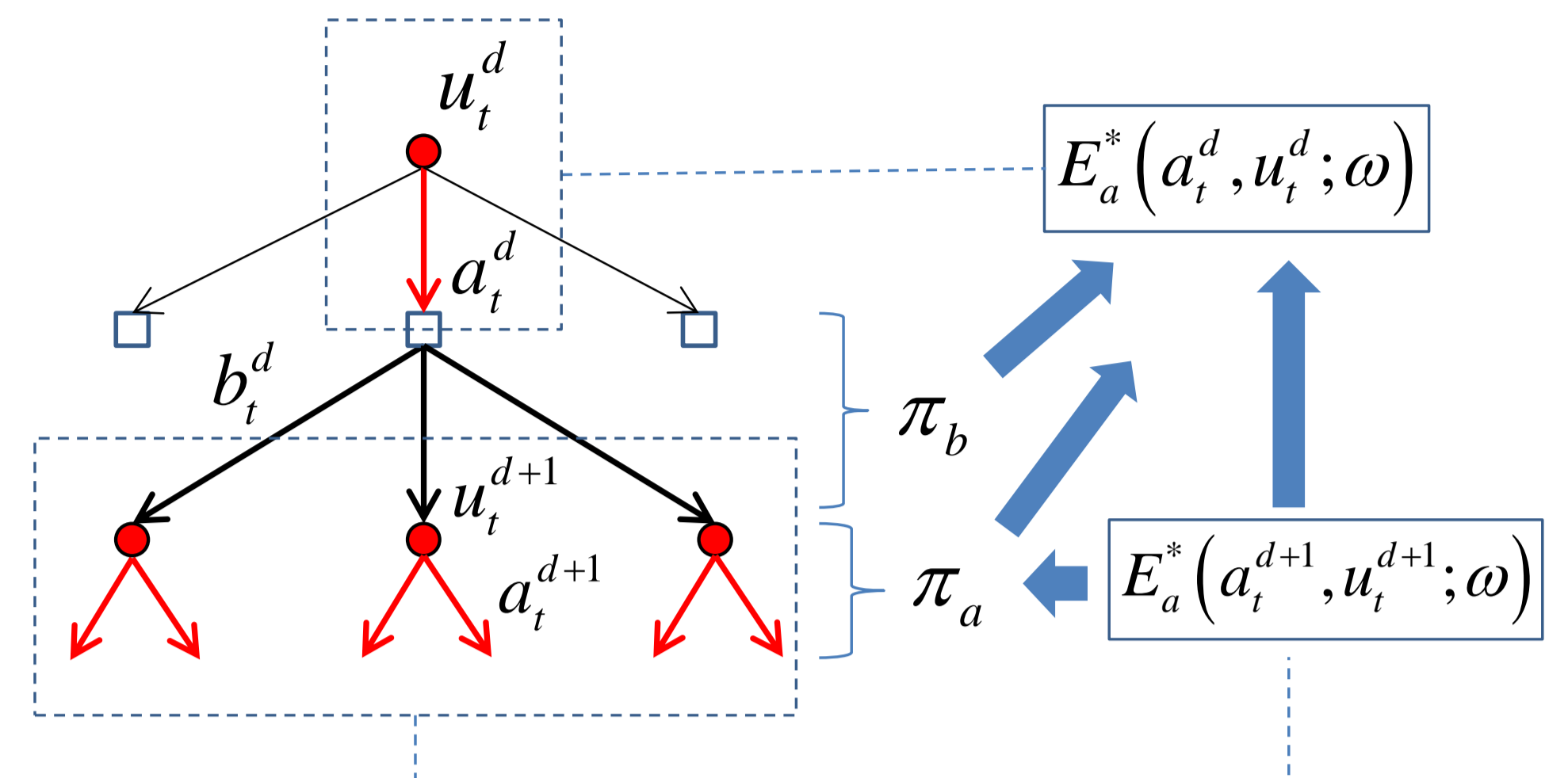
$$E_a^*(a_t, u_t; \omega) \equiv \sum_{u \in U_D(a_t, u_t)} P(u | a_t, u_t; \omega) E_a^s(u; \omega)$$



Boltzmann
分布による
確率的探索

- 従来
- ・min-max探索
 - ・選択探索
 - ・モンテカルロ探索

学習則に必要な再帰的計算



近似的計算法のアイデア

