

福岡工業大学 機関リポジトリ

FITREPO

Title	確率環境下でのモデルベース学習
Author(s)	田村剛士, 鶴岡久, 山口明宏
Citation	福岡工業大学研究論集 第44巻1号(通巻67号) P9-P12
Issue Date	2011-9
URI	http://hdl.handle.net/11478/1289
Right	
Type	Departmental Bulletin Paper
Textversion	Publisher

Fukuoka Institute of Technology

確率環境下でのモデルベース学習

田 村 剛 士 (綑たけびし)
鶴 岡 久 (情報システム工学科)
山 口 明 宏 (情報システム工学科)

Model Based Learning in Stochastic Environments

Takeshi TAMURA (Takebishi Co.)

Hisashi TSURUOKA (Department of Information and Systems Engineering)

Akihiro YAMAGUCHI (Department of Information and Systems Engineering)

Abstract

Model based learning can reevaluate the utility of every state, according to a measure of urgency. Prioritized sweeping is a typical algorithm for efficient state updating. In a stochastic environment, a probability distribution can be used to represent the uncertainty of the Q-value caused by probabilistic state transitions or probabilistic rewards.

The product of the confidence interval and the Bellman error is used to provide a measure for prioritizing, which takes account of the level of confidence and also yields a measure of urgency. The performance of this approach in the trap domain is examined and compared with that of the ordinary sweeping method. Experimental results indicate that the proposed approach results in a more effective exploration of the state than does the use of conventional sweeping methods.

Key words: *model based learning, prioritized sweeping, Q-value, Bellman error, confidence interval*

1. はじめに

TD 学習に分類されている Q 学習は、強化学習における重要な発展の一つである。有限マルコフ決定過程において状態行動対が全て更新され続ければ、最適行動価値へ収束されることが証明されているからである。しかし Q 学習を含む TD 学習は、動的計画法やモンテカルロ法の考え方を取り入れたものではあるが、動的計画法で取り扱っていた遷移確率集合を直接扱わなくなった。エージェントの遷移が決定論的に行われるならば TD 学習でも問題は無いが、現実の問題では環境とエージェントのやり取りが完全である事は少ない。ゲーム問題等以外の問題は、常に確率的な環境を考えるべきである。未知の環境を不確定な行動をしながら学習するロボット等には確率環境下で学習効率の良い強化学習アルゴリズムが求められる。

強化学習には TD 学習のように実際の行動を通して経験

から価値を推定し、方策を改善するモデルなし学習と、エージェントが観測結果から環境のダイナミクスを構築して状態価値を推定し、方策改善を行なうモデルベース学習がある。モデルベース学習は学習効率をあげる重要な手法の一つであり、Dyna-Q や優先度スイープはその代表的手法であるものの、確率環境下では確率論的モデルに適応させる必要がある。

2. 背景

2.1 プランニング

強化学習では環境モデルをもち、エージェントがモデルの環境と相互作用して、方策を作り出すか、あるいは改善を行う任意の計算過程をプランニングと呼ぶ。実際の経験に対して、プランニング・エージェントには少なくとも 2 つの役割がある。1 つは、モデルの改良、もう 1 つは、直接的に価値関数と方策を改善することである。前者をモデル学習、後者を直接的強化学習 (RL) と呼ぶ。学習手法とプランニング手法の双方において中心となるのは、バック

アップ操作による推定価値である。相違点は、プランニングでは、モデルが生成したシミュレーション上の経験を使うのに対して、学習手法では、環境が作り出す実際の経験を使う点である。Dyna-Q アルゴリズムは各ステップにおいて、行動、モデル学習、RL を行った後に、Q 学習アルゴリズムの N 回繰り返し計算を行うものである (Fig. 1)。

```

Initialize  $Q(s, a)$  and  $Model(s, a)$  for all  $s \in S$  and  $a \in A$ 
Do forever:
  (a)  $s \leftarrow$  current (nonterminal) state
  (b)  $a \leftarrow \epsilon$ -greedy ( $s, Q$ )
  (c) Execute action  $a$ ; observe resultant state,  $s'$  and reward,  $r$ 
  (d)  $Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a' \in A} Q(s', a') - Q(s, a)]$ 
  (e)  $Model(s, a) \leftarrow s', r$ 
  (f) Repeat N times:
     $s \leftarrow$  random previously observed state
     $a \leftarrow$  random action previously taken in  $s$ 
     $s', r \leftarrow Model(s, a)$ 
     $Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a' \in A} Q(s', a') - Q(s, a)]$ 

```

Fig. 1 Dyna-Q Algorithm

状態数が多い大規模問題に対してはバックアップを緊急度の優先順位、通常は実際の経験から得られた価値関数とモデルで得られる価値関数の差 (Bellman error) の大きいものを優先する、優先度スイープ (Prioritized Sweeping) が有利とされている¹⁾。

2.2 確率環境

状態遷移や報酬獲得が確率分布で与えられる確率環境では、価値関数も学習過程では不確実性を有し、Bellman error そのものが不確実性をもつため、決定論的環境と同じ意味で優先度を考えることができない。状態遷移や報酬が確率分布で表される確率環境下での価値関数の分布は R. Darden 等によって考察され²⁾、行動選択手法に応用されているが、優先度スイープへの直接的な応用はない。

3. 提案手法

4 組 (S, A, PT, PR) (S: 状態, A: 行動, PT: 状態遷移確率 (状態 s で行動 a を起こしたときの状態 s' に到達する確率), PR: 報酬獲得確率 (状態 s で行動 a を起こしたときの報酬 r を受け取る確率)) で表現されるマルコフ決定過程 (MDP) を考える。最適価値関数 V^* 、最適行動価値関数 Q^* 、は下記の Bellman 方程式に従う。

$$V^*(s) = \max_{a \in A} Q^*(s, a) \quad (3.1)$$

$$Q^*(s, a) = \sum r \Pr(r|s, a) + \gamma \sum_{PT(s, s')} V^*(s') \quad (3.2)$$

状態 s で行動 a を起こしたとき、状態遷移先 s' がランダムに揺らいだり、獲得報酬がランダムなノイズを受ける確率環境下では価値関数の値や Q 値が不確定な値になる。その分布関数は、学習過程において複雑であるが、訪問回数の増加にともない、価値関数の分布は正規分布に近づくことが示されている³⁾。従って学習初期の訪問回数が少ない場合は、実現する価値関数の値は正規分布に従う母集団からのランダム抽出と等価と考えられるので、価値関数の分布は訪問回数を自由度とする t 分布で近似できる。M. Wiering⁴⁾、A.L. Strehl⁵⁾ 等は状態遷移確率分布の信頼区間

$$\left(\bar{X} - t_{n-1}(\alpha) \frac{S}{\sqrt{n-1}}, \bar{X} + t_{n-1}(\alpha) \frac{S}{\sqrt{n-1}} \right) \quad (3.3)$$

の上限を利用し、行動選択の探査効率の向上を図っている。しかし優先度スイープの優先度に直接関係するのは状態価値であるから、確率環境下での優先度を次のように考える。

それまで訪問した価値関数の分布における信頼区間は、未知環境に対する新しい知見が得られる長期的尺度を表すと考えられる。また状態価値の信頼区間のみを優先度に反映させたのでは即応的な学習改善が図れないと考えられる。そこでモデル価値と実測価値の誤差を表す従来の Bellman error は報酬改善が期待される短期的尺度と考え、双方を考慮した優先度を採用する。本論文では両者の積に応じた優先度に従ってスイープする。Fig. 2 に提案手法の処理手順を示す。s1~sn は状態 s に対する現時点から過去に遡る n 個のサンプリング点を示し、これから t 分布近似による信頼区間が計算されている。

4. 実験結果

Fig. 3 に提案手法を評価するためのフィールドを示す。エージェントは S から出発し、フラグを集めてゴールに入るとプラス 1 の報酬を得て 1 エピソードが終了するが、途中トラップに入るとマイナス 10 の報酬が加算される。

行動 (up, down, left, right) は進行方向が空いていれば、確率 0.9 で移動が成功するが、確率 0.1 で希望する方向と任意の直角方向へ動く。学習収束判定は続けて 2 回最短ステップでゴールしたこととする。Table 1 に学習収束するまでのステップ数の 50 回平均値を示す。1 エピソード毎に行なう優先度スイープはどの方法でも状態数に等しく 9 回とした。t 分布の信頼区間は両側 95% を採用した。

```

Initialize  $V(s), Q(s, a), Model(s, a)$ 
Initialize  $s$ 
Repeat (for each episode)
  Repeat for each step of episode
     $Q(s, a) \leftarrow \sum P(R + \gamma V(s'))$ 

```

```

V ← max Q
Choose a from s using policy from Q
Next state s'
Model(s, a) ← s', r
Update V(s1) ... V(Sn)
Planning
  Calculation of Confidence Interval
  Decision of Prioritized sequence
  Repeat Ns times
    Q, V update
until s is terminal
    
```

Fig. 2 Proposed Algorithm

S	F	
		T
		G

Fig. 3 Domain used in the experiment

S predicates the start space, G is the goal state, T is the trap state and F is the flag.

Table 1 Comparison of sweeping methods

Sweeping	Steps to Convergence
Dyna-Q	14.66
Prioritized	10.58
Confidence Interval	10.76
Proposed	10.12

Confidence Interval は信頼区間の大きさのみを優先度の判断としたことを示す。Table 1から提案方式が Bellman error のみを指標とした従来の Dyna-Q スイープ, 優先度スイープより学習効率が良いことが分かる。

次に20回学習を繰り返したときの, 平均獲得報酬の変化を Fig.4に示す。優先度スイープは学習初期でマイナス報酬が大きく, トラップを回避できていないことを示している。4ステップ以降を拡大したのが Fig. 5である。DYNA や優先度スイープは獲得報酬が容易に安定して正の値に向かわないことが分かる。

5. 結言

学習変化していくモデルを活用して学習効率を向上させる方向には行動選択とプランニングがある。本論文は後者に注目し, 確率環境下での優先度スイープに従来の1ポイントの時間に着目する Bellman error による優先度決定よ

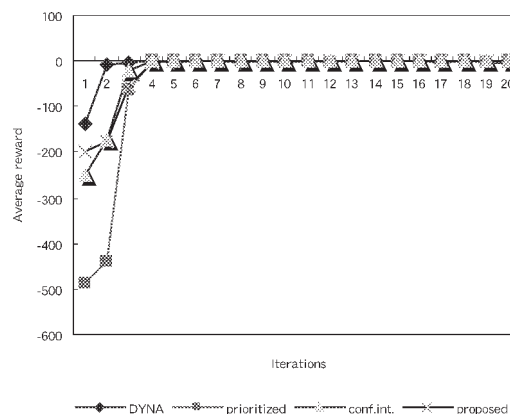


Fig. 4 Plots of average reward as a function of number of iterations

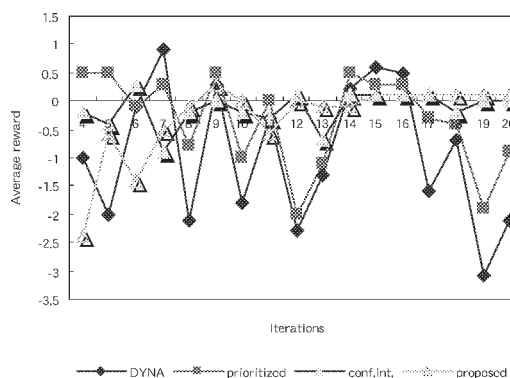


Fig. 5 Enlarged plots from 4 to 20 iterations

り, 時間スケールの長い信頼区間も考慮した効率の良い優先度を提案し, 実験的にこれを確認した。

行動から得られる経験から状態遷移分布や報酬分布を学習していく手段にはベイジアンモデルや母数モデルがあり, 正確には信頼区間はこのような個々の経験によって個別に求められた分布から決めるべきと考えられるが, 今後の課題である。

6. 参考文献

- 1) A.Moor and C.G.Atkeson: Prioritized sweeping; Reinforcement learning with less data and less time, Machine Learning, 13, pp.103-130 (1993)
- 2) R.Dearden, N.Friedman, S.Russell: Bayesian Q-learning, Fifteenth National Conf. on Artificial Intelligence (AAAI), pp.761-768 (1998)
- 3) R.Dearden, N.Friedman, D.Andre: Model based Bayesian Exploration, Proc. Fifteenth Conf. on Uncertainty in Artificial Intelligence, (1999)
- 4) M.Wiering, J.Schmidhuber: Model-Based Exploration, From Animals to Animats 5, pp.223-228 (1998)

- 5) A.L.Strehl, M.L.Littman: An Empirical Evaluation of Interval Estimation for Markov Decision Process, The 16th IEEE Int. Conf. on Tools with A.I. (ICTAI-2004), pp.128-135 (2004)