

## 関数近似を用いた強化学習の Khepera ロボットへの応用

甲南大学大学院 自然科学研究科 情報・システム科学専攻 \*大前 雅裕 OHMAE Masahiro  
中山 弘隆 NAKAYAMA Hirotaka

Application of the Reinforcement learning using function approximation to Khepera robot

Abstract : Reinforcement learning is made by the interaction with environment. But it is difficult to apply reinforcement learning to robots as it is. Because enormous time of computation is needed. This paper proposes a method which uses function approximation for Reinforcement learning, and employs it to Khepera robot's feature efficiently in an actual environment.

## 1 はじめに

ロボットの学習方法として強化学習の関心が高まっている。強化学習は環境との相互作用で学習を行うが、実際の環境においてはデータが非常に多くなりロボットにそのまま適用するのは困難である。そこで本稿では強化学習に関数近似を用い、さらに Khepera ロボットの特徴を生かし Khepera ロボットが実際の環境で学習できる手法を提案する。

## 2 強化学習

## 2.1 強化学習とは

強化学習とは、試行錯誤を通じて環境に適応する学習制御の枠組みである。強化学習は報酬というスカラーの情報を手掛かりに学習する。Fig.1 で強化学習におけるエージェントと環境の相互作用を示す。

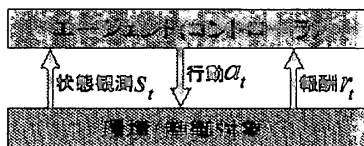


Fig. 1: 強化学習におけるエージェントと環境間の相互作用

強化学習では Fig.1 のように、まずエージェントは環境から状態観測  $S_t$  を受け取り、行動を決定し、行動  $a_t$

を起こす。その行動によって環境は報酬  $r_t$  を渡す。これを繰り返し学習する。

## 2.2 RBF ネットワーク

RBF ネットワーク (Radial Basis Function:放射基底関数) とは、中心からの距離が離れるにつれて、値が単調に減少 (または増加) し、その等高線が超球になる関数のことをいい、代表的なものにガウス関数

$$h(x) = \exp\left(-\frac{\|x - c\|^2}{r^2}\right)$$

がある。本研究ではこの RBF ネットワークを用いて関数近似を行う。

## 2.3 関数近似

強化学習における関数近似の使い方について述べる。強化学習は状態価値関数 ( $V(s)$ ) によって学習する。そこで状態価値関数 ( $V(s)$ ) を

$$V(s) = \sum_{i=1}^n w_i h_i(s)$$

によって近似する。  $n$  は基底の数、  $w$  を重みとする。本場の正しい値を  $v(s)$  とすると、すべての状態  $s$  に対して

$$\frac{1}{2} \{v(s) - V(s)\}^2 \rightarrow \text{Min}$$

となるように重みを調節していく必要がある。そこで最急降下法の考え方をを用いて、勾配ベクトルの逆の方向  $d$  は  $\frac{1}{2} \{v(s) - V(s)\}^2$  を偏微分したものに  $-1$  を掛けると、

$$d = \{v(s) - V(s)\} h_i(s) \quad (1)$$

となる。

重みは最急降下法を利用して、

$$w_i \leftarrow w_i + \alpha d$$

を更新していくことによって求めていく。式(1)を代入すると、

$$w_i \leftarrow w_i + \alpha \{v(s) - V(s)\} h_i$$

強化学習の考え方により、 $v(s) - V(s)$ をTD誤差で近似すると、

$$w_i \leftarrow w_i + \alpha \{r_{t+1} + \gamma V(s') - V(s)\} h_i \quad (2)$$

と導ける。この式(2)を更新することによって最適な重みを求めていく。

## 3 Kheperaロボットへの応用

### 3.1 Kheperaロボット

Kheperaロボットには左に1つ、左斜め前に1つ、前に2つ、右斜め前に1つ、右に1つ、後ろに2つ、計8つの赤外線近接センサ（左斜め後ろ、右斜め後ろにはない）が装備されている。このセンサは0～1023の整数値を返し、障害物が遠くにあるほど小さく、近くにあるほど大きくなる。

### 3.2 応用方法

実際の環境でKheperaロボットの位置を強化学習の状態とすると、Kheperaロボットの正確な位置を判断するのは非常に困難になる。Kheperaロボットの目的は障害物を避ける、目的地に行く、物を運ぶなどなので、Kheperaロボットのセンサによって目標物との距離が測れれば十分である。そこでKheperaロボットのセンサ値を強化学習の状態のパラメータとすると、 $1024^8$ 個の状態を用意することになる。しかしKheperaロボットに $1024^8$ 個のデータを持たせることは不可能である。そこでKheperaロボットのセンサ値を何個かに区切り、センサも8個から4個に減らすことにより、データ数を減らし学習させてみた。しかしこの方法では近い値のセンサ値なら同じ行動をとったり、使用しないセンサ部が死角になったりして正しく学習ができなかった。そこで関数近似を用いることでさらに多くのデータ数を扱え、より学習が可能だと考えた。この手法だとRBF

ネットワークの基底における重みの数だけデータを用意すれば、あとは計算でそれぞれの状態価値関数 $V(s)$ を求めることができる。

### 3.3 実験環境

実験環境として、第1の実験は、迷路の中にKheperaを置きスタートからゴールまでの経路を探索させる。迷路はKheperaが常に壁を感知できるものとし、ゴールから先には壁を置かずに壁を感知させないつくりにする。Kheperaにはセンサしか無いので、センサが壁を感知しなければゴールしたと判断させることができる。どれくらいの学習させればゴールできるようになるか計測する。

第2の実験は、2台のKheperaを使い追いかけっこをさせる。逃げる方に本研究を用い、追う方は学習させずに常にもう1台のKheperaを追うようにプログラムする。そして学習用のKheperaは追跡用のKheperaからどれくらい逃げられるかを計測する。

## 4 結果

実際にKheperaロボットで実験を行ったところ、センサをすべて使えるためKheperaロボットの死角が減り、センサの感知する範囲内において、障害物を回避することができた。また関数近似を使うことによってデータを正確に使えるようになり、関数近似を使わない方法よりも良く学習できた。詳細は当日発表する。

## 5 謝辞

Kheperaロボットについて詳しく教えていただいた、福井大学工学部知能システム工学科の皆さんに深く感謝の意を表します。

## 参考文献

- [1] Richard S.Sutton and Andrew G.Barto(2001):「強化学習」森北出版
- [2] 堀越 俊之(2002):「適格度トレースと関数近似を用いた強化学習」甲南大学理学部卒業論文
- [3] 福井大学リカレント講座(2002):「Khepera入門」