

複数 AGV の Q-学習による衝突回避行動の獲得

旭川高専 ○池田将晴 渡辺美知子 古川正志

要旨

複数 AGV の自律運転技術を確立するためには、他の AGV との衝突を避け円滑に走行する必要がある。本研究では、逆に他の AGV との衝突予想行動を Q 学習によって獲得することを行い、得られた知識により衝突回避を行う方法を提案し、衝突予想のシミュレーション結果を示し有効性を検証した。

1. はじめに

自律分散型生産システムの構築において、各ショップ間を自由に移動できる AGV の自律行動を獲得することは難しい問題である。このような問題には経路獲得問題と衝突回避問題があり、すでに Q 学習を採用した経路獲得法を提案してきた。本研究では、AGV 同士が様々な方向から衝突する可能性を考慮して、自由度の高い衝突回避行動を逆に衝突予想行動を Q 学習により獲得することで実現する。

2. 衝突行動獲得問題

衝突回避行動を獲得するにあたって、走行中の AGV を避ける走行ルートは多数存在する。衝突回避は衝突が予想される行動以外の経路であればよいから、本研究では、衝突行動を獲得する方法を取り扱う。衝突回避問題ではこの衝突行動を行わなければ目的が達成される。

2.1 AGV のモデルと前提条件

実験に用いる AGV のモデルを図 1 に示す。

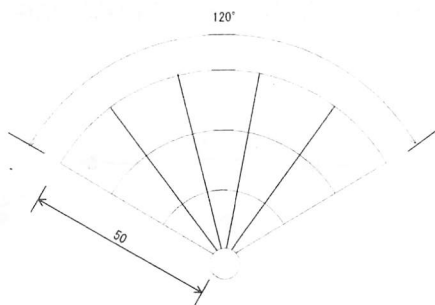


Fig.1 AGV のモデル

AGV には以下のような前提条件を設定する。

- 1) AGV は前方に 5 個のセンサを持つ。視野は等角度の小領域に 1/5 分割されている。
- 2) 前提としたセンサによって近づいた他の AGV の距離・方向を知覚できる。

2.2 問題の記述

前提条件をもとに AGV が最短時間で相手 AGV に到達するような一連の行動を求めることである。

2.3 オートマトンによる学習モデル

AGV を 1 個のオートマトンと見なす。自 AGV と相手 AGV をそれぞれ A_i, A_j とすると、各 AGV, A_i は、

$$A_i = \{ I_i, O_i, S_i, F_i, G_i \} \quad (1)$$

の 5 項組で記述できる。ここで各記号は、

$$\begin{aligned} I_i &: \text{入力} & S_i &: \text{AGV の状態} \\ O_i &: \text{出力} & F_i &: \text{状態遷移関数} \\ G_i &: \text{出力関数} \end{aligned}$$

である。また、 A_j も同様である。

1) 入力 I_i

入力は以下のように構成する。

$$I_i = \{ T_i(t), D_i(t) \}$$

ここで、 t は時刻を表し、

$$T_i(t) : \text{センサーがとらえた } A_j \text{ の情報}$$

$$D_i(t) : A_i \text{ の進行方向}$$

である。

センサから得る情報 $T_i(t)$ は、

$$T_i(t) = \{ DIS_{A_j}, ANG_{A_j} \}$$

とする。ここで DIS_{A_j}, ANG_{A_j} は相手 AGV との距離、方角である。図 1 で示すように、これらの情報は AGV のセンサの小領域に相手 AGV がいるか否かで定まる。

AGV の進行方向 D_i は、図 2 に示すように

$$D_i(t) = \{ D \mid 0, 1, \dots, 11 : \text{時計方向に対応} \}$$

とする。

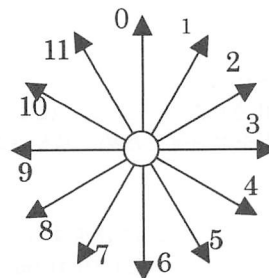


Fig.2 AGV の進行方向

2) 出力 O_i

出力は時刻 t における AGV の行動 $O_i(t)$ とする。

$$O_i(t) = \{ h_left, l_left, forward, l_right, h_right \}$$

: 左急旋回, 左旋回, 前進, 右旋回, 右急旋回

AGV は 5 方向に移動することができる。

3) 状態 S_i

状態は, センサーの情報である。

$$S_i(t) = \{ T_i(t), T_i(t-1), T_i(t-2) \}$$

衝突経路は一連の行動のつながりであること考え, 時刻 $t-2$ までのセンサ情報を時刻 t における状態とする。

4) 状態遷移関数 F_i

AGV の状態は相手 AGV がセンサ内に入るたびに遷移する。

状態の遷移関数としては Q 学習で用いられる Q 値を採用する。

$$Q_{i,t+1}(s_{i,t}, o_{i,t}) = (1-\alpha)Q_{i,t}(s_{i,t}, o_{i,t}) + \alpha \left[f_c + \gamma \max_{b \in O} Q_{i,t}(s_{i,t+1}, b) \right] \quad (2)$$

式(2)における環境情報からの報酬は, AGV のスタートから衝突するまでの行動で評価する。さらに衝突するまでの経路は一連の行動であることを考えて, バケットブリゲードにより報酬値 f_c を割引係数 Γ , 割引率 d_c を用いて以下のように与える。

$$\Gamma_{i+1} = \Gamma_i \cdot d_c \quad (3)$$

直接評価 $f_c = R(\text{成功}) \cdot \Gamma_i$

$f_c = P(\text{失敗}) \cdot \Gamma_i$

ただし, $i = 1, 2, \dots, T$

$R > 0, P < 0$

ここで, T は過去に状態を遡る数である。

5) 出力関数 G_i

G_i は時刻 t における AGV の状態 S_i から行動 O_i 確率的に決める関数である。本実験では, G_i をボルツマン分布式(4)によって確率的に決定する。

$$O(s, o) = \frac{\exp\{Q(s, o) T\}}{\sum_{b \in O} \exp\{Q(s, o) T\}} \quad (4)$$

3. 衝突回避問題

衝突問題より理想的な衝突行動が獲得されれば, 衝突回避ではそれ以外の行動をとればよい。

しかし, 衝突回避後は本来の目的地を目指す行動に移行する必要がある。

4. 実験結果

AGV が走行する環境を図 3 に示す。ターゲットの経路は一定方向に直線移動し, 外周を障害物とする。Q 学習パラメータは, 学習率 $\alpha=0.8$, 減衰率 $\gamma=0.8$, 温度係数=1.0 とした。オフラインの Q 学習では 17000 回の学習で最適解に収束した。図 4 に学習終了間際の AGV の行動を示す。また図 5 に Q 学習の学習曲線を示す。

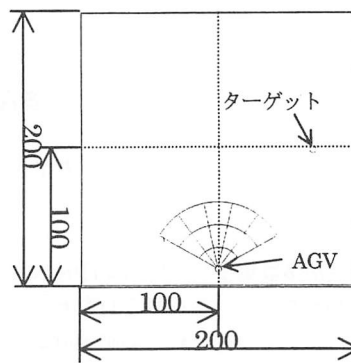


Fig.3 実験モデル

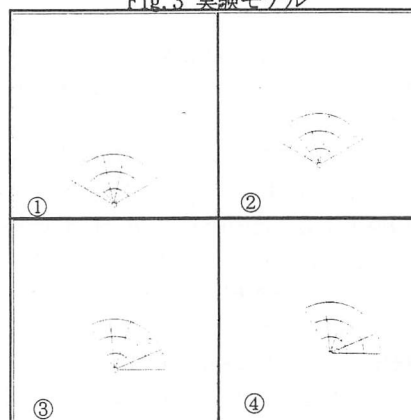


Fig.4 シミュレーションの様子

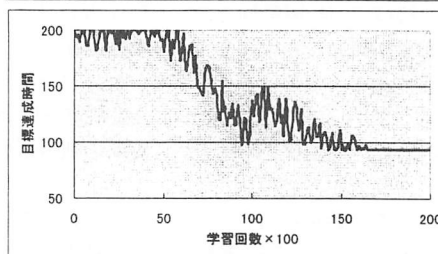
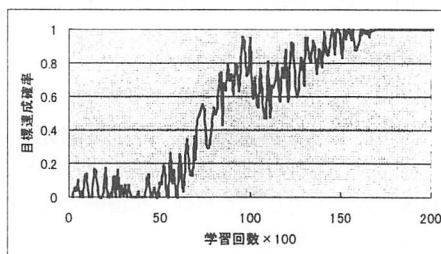


Fig.5 学習曲線

(上:達成時間収束グラフ, 下:達成確率収束グラフ)

5. おわりに

本研究では, AGV に単純なセンサーを持たせているのにもかかわらず, ターゲットの行動を予測し適確に行動させることができた。ただし, この強化学習に対するパラメータに関しては最適値であるとは限らず, 何らかの方法で最適値を求める必要がある。

参考文献

1) 池田将晴, 渡辺美知子, 古川正志 : 仮想工場における多数 AGV の自律運転に関する研究, 1999 年度精密工学会北海道支部講演会論文集(1999)44