

要 旨

本研究では一方向のみAGV (Autoatically Guided Vicle) が移動できる機械工場において、複数AGVが各AGV間の衝突を回避しながら自律行動するときの行動決定問題を取り扱う。この行動決定法として、各AGVが全体の周回完了時間を最短とするような強化学習法 (Reinforcement Learning) に基づくモデルを提案する。また、数値計算実験より本モデルが有効である事を検証する。

1. はじめに

現在、工場での無人搬送手段としてAGVが使用されている。工場内で複数のAGVの行動を決定する場合、大きく分けて二通りの方法がある。一つは事前に全AGVのスケジュールを決定し与える方法であり、もう一つは各AGVに自律的意志決定をさせる方法である。本研究は後者の方法を採用し、強化学習法を用いて全AGVの行動時間を最小化する行動パターンの決定を強化学習を採用して行う方法を提案する。また、本方法の有効性を、数値シミュレーションにより検証する。

2. AGVモデル

本研究で採用した自動化工場を図1に示す。AGVに關する制約条件を、以下に記述する。

- 1) 各AGVの速度は一定とする。
- 2) 移動時間は、図1の1マス分を移動する時間を1単位時間とし、機械工作に要する時間はその5倍とする。
- 3) レーンはAGV1台が走行可能とし、並走はできないものとする。
- 4) レーンは一方通行とする。
- 5) 各AGVの工程表は乱数により決定する。
- 6) AGVは、工程表に記された機械を全て回り出発点に戻った時点で作業を終了する。また、巡回順は自由とする。
- 7) AGV間の衝突は回避する。

AGVの自律的行動決定方法として強化学習法(RLA)を適用し、確率的にAGVの行動を決定する。個々のAGVは、添え字*i*を用いて以下のように記述される。

$$RLA_i = \{S_i, A_i, r, P_i, Q_i\} \quad (1)$$

- S_i : AGVの走行レーン上の状態
- A_i : 出力
- r : AGVの移動に休する応答
- P_i : 出力関数
- Q_i : 強化値更新アルゴリズム

(a) AGV、工場内等の状態 (S_i)

工場内のAGVは、位置及び作業の進行状況等の状態を8bitsのストリングで表し、これを $S_i = \{S_0, S_1, \dots, S_n\}$ とする。各々のbitの意味は以下のとおりである。

- S_0 : 前方4step内の他のAGVの有無
- S_1 : 前方4step内の作業場の有無
- S_2 : 右方の進路の有無
- S_3 : 進行方向の壁の有無
- S_4 : 右方4step内他のAGVの有無
- S_5 : 右方の作業場の有無
- S_6 : 全作業の終了判定
- S_7 : AGVの状態 (作業中または移動中)

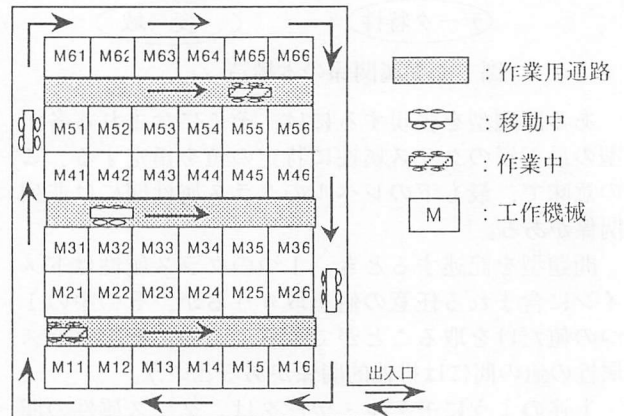


図1 機械工場モデル

(b) 出力 (A_i)

出力 A_i は、AGVの行動パターンを記述する。行動は5種類あり、それを以下に示す。

- ACT1 : 停止
- ACT2 : 直進
- ACT3 : 右折
- ACT4 : 作業中
- ACT5 : 出口に向かう

ただし、いくつかの状態では行動を制限するので、あらかじめ行動を設定し上記の行動を禁止する。また、一方通行のため左折及び後退は行われない。

(c) AGVの行動に対する応答 (r)

時間*n*と時間(*n*-1)における最短の目的地(作業場または出口)までの距離DISを用いて応答*r*を以下の式で決定する。ただし、作業中は応答はないものとする。

$$r = DIS(n-1) - DIS(n) \quad (2)$$

このとき、*r*がプラスであれば成功、0またはマイナスであれば失敗とする。

(d) 出力関数 (P_i)

出力 A_i (AGVの行動パターン)は、状態 x_i における各行動の強化値のマトリックス Q_i から

$$P(S|a) = \text{Prob} \left(\frac{e^{Q_{S,a}(t)}}{T} / \sum_{b=0}^m e^{Q_{S,b}(t)} / T \right) \quad (3)$$

から、確率的に決定する。ここで*b*は行動を示している。ただし、作業場付近、衝突回避状態では各々作業中、停止が絶対的に選択される。

(e) 行動重み更新アルゴリズム (Q_i)

強化学習の重みのマトリックスの学習、更新する。ただし、応答がないならば、学習はしないものとする。

$$V_x(t) = \max_{b=0} Q_{x,b}(t)$$

$$Q_{a,i}(t+1) = (1 - \alpha) Q_{a,i}(t) + \alpha r(t) + \alpha \gamma V_x \quad (4)$$

このときの*b*は行動で V_x はその状態のなかで一番大きい

Qの値である。また、 α は学習率、 γ は割引率とする。

3. 計算機実験

図1のような工場に対しAGVの台数を5台、一つのAGVが受け持つ作業数を5ヶ所、 $\alpha=0.0003$ 、 $\gamma=0.8$ 、 $T=0.05$ として数値シミュレーションを行った。全体の作業数は25ヶ所である。学習回数1000回とした結果を図2に示す。学習を重ねると全てのAGVが作業を終える時間は160前後に収束するが、140から200の間で変動する。また、各AGVの学習状況を図3に、ある状態での強化値の変化を図4に示す。尚、1000世代での最短の作業所要時間は111であった。

4. 考察

本実験での学習が収束している状態での作業時間が変化する理由として、以下のことが考えられる。

- 1) 行動が確率的に決定されるので、必ずしも良い行動が選択されるとは限らない。このことは、逆に新しい環境に対して学習適応性を持つ。
- 2) AGVに与えられる作業場が複数ヶ所に接している場合、AGVは移動時間にロスや殆ど被ることなく、複数の作業を実行できる。
- 3) AGVに与えられる作業場は、ランダムに選んでいるために作業場が近くなったり遠くなったり毎回変化する。

また、状態数は256通り用意したが、実際に使用されたのは62通りである。これは、62通りの規則生成を行えたことを意味する。

5. おわりに

本実験ではAGVの自律意志決定問題に強化学習法を適用し、良好な結果を得ることができた。しかし、パラメータの調整が収束に影響することが分かった。今後の課題はこのパラメータの調整と、自動化工場のモデルの大きい場合、AGVや工作機械の性能の違いを持たせた場合についての検討をする必要がある。

参考文献

- 1) Jeffery A. Clouse and Paul E. Utgoff; A Teaching Method for Reinforcement Learning, Machine Learning (ML92)
- 2) 畝見達夫; 実例に基づく強化学習法、人工知能学会誌 Vol.7 No.4 July 1992
- 3) 渡辺他; 群AGVのSLAによる自律的意志決定、1994年度精密工学会北海道支部学術講演会論文集、1994

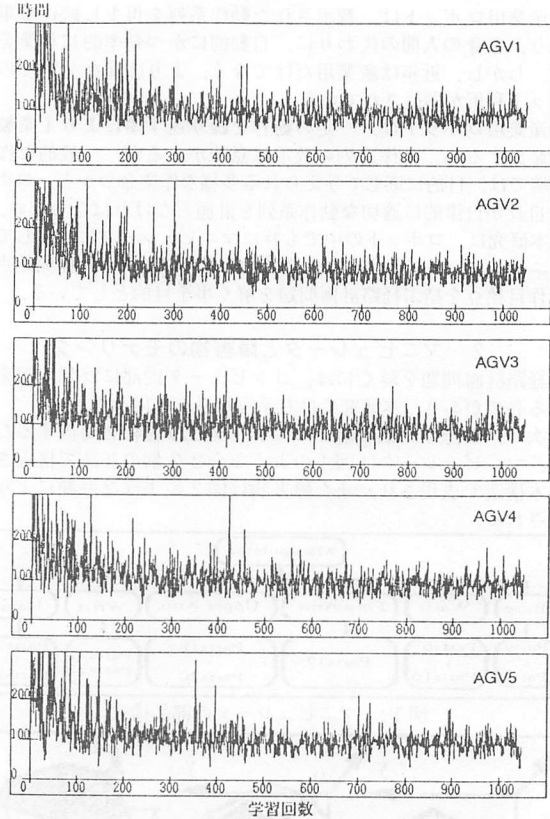


図3 各AGVの周回完了時間の収束

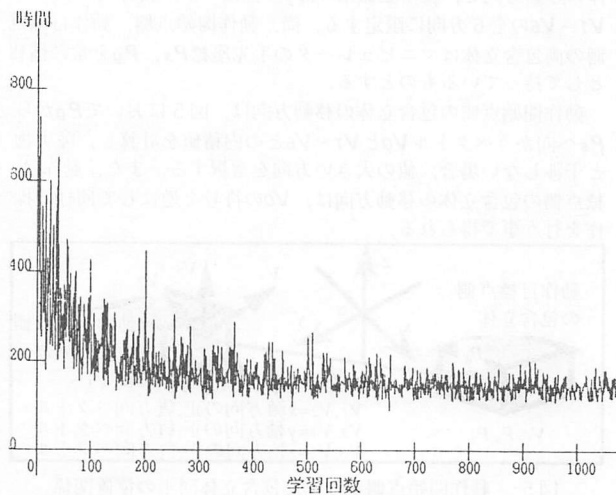


図2 AGV全体の周回完了時間の収束

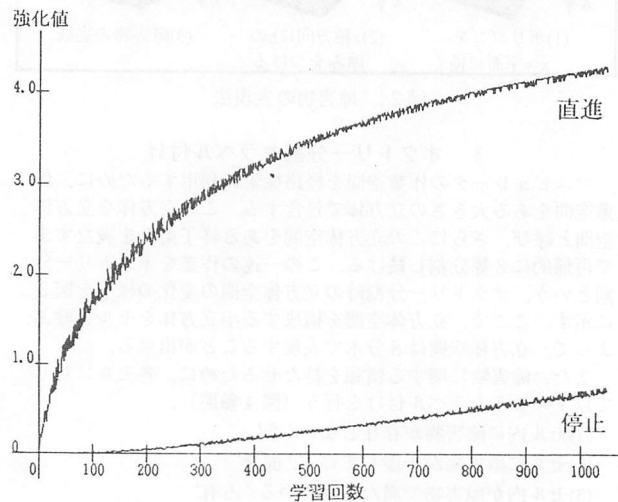


図4 01100110の状態におけるACT1とACT2の行動の強化値の様子