

1. はじめに

高機能材料など多品種少量で高品質な製品の製造には、高度な自動制御システムに加えて、熟練オペレータによるプロセス状態の確認と必要に応じた手動操作が欠かせない。しかし近年、団塊の世代と呼ばれる大量の熟練者が定年を迎え、長年の経験による優れた技能を持つ彼らが、製造現場から去りつつある。一方、国際競争の中、製造現場で働く作業者の多くは需要に柔軟に対処できる派遣スタッフとなり、正規スタッフが大幅に減少している職場もある。多くの製造現場で、どのように熟練者の技能を伝承していくのかが大きな課題となっている。従来は、例えば個々の反応工程を一時遅れ関数などの伝達関数によって記述しプロセスモデルを構築してきた。しかし、投入原料の成分変動や運転条件などによって、プロセスの状態は大きく変化してしまう。日々の実データからプロセス特性を得ることができれば、これらの問題は解決する。本稿は、大量の時系列データからプロセス応答モデルを構築し、運転員の操作ガイダンスとなる制御ルールの発見的探索手法を紹介する。そして、実際のバイオプラントの操業データに適用した結果を報告する。モデルの基本的な考え方は、時系列データ間の相互相関係数最大化、MDL(Minimum Description Length)基準¹⁾とタブー探索手法を取り入れた分類子学習である。

図1は、あるバイオプラントのプロセスデータを正規化したものである。図1に示す正規化されたデータを見ただけでは、プロセスデータ間にどのような関連があるのかを読み取ることは不可能に近い。実際のプラント操業を行う運転員は、プロセスの特性を経験的に獲得し、制御設定や手動操作を頻繁に行うことによってプラント全体をコントロールしている。これに対し、本手法を適用すると単純かつ説得力のあるワークフローが得られる。

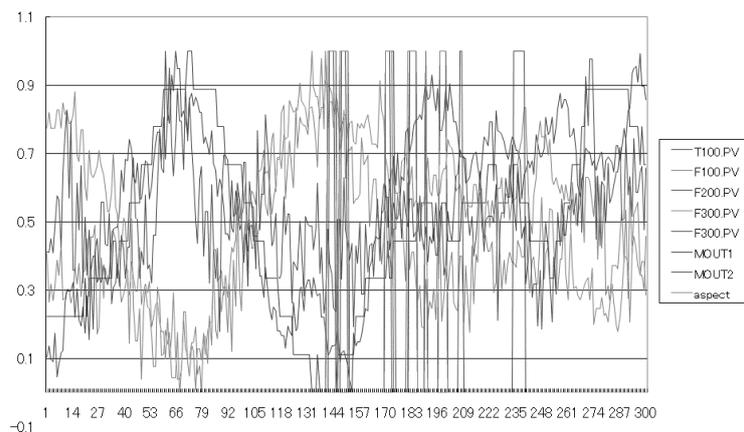


図1 プロセス時系列データ

2. 問題設定

本手法では、このような複雑に見えるプラント操業時系列データから、熟練者の操作情報を取り出すことを目的とする。プロセスデータ解析のためのモデル構築のステップを以下に示す。

1) プロセスデータ収集

プロセスデータを定周期で収集し、プロセスデータベースに時系列データとして格納する。

2) データの正規化

収集したプロセスデータのスケールはそれぞれ異なるため、0.0~1.0 に正規化する。

3) 相互相関分析

正規化された二つのプロセスデータを選択し、時間を徐々にシフトして最大の相関を示す時間差を探索する。

4) プロセス応答モデル

シフト時間と相関係数から、プロセスデータ間の関係を応答モデルとして記述する。

5) 制御ルールの発見

プロセス応答モデルから、注目するプロセスデータをクラスとして、MDL 基準学習分類子システムを実行し、制御ルールを見出す。

6) ワークフローの発見

運転操作イベントとプロセスデータに対して、タブー学習分類子システムを実行し、それぞれの発生時刻が異なるシーケンシャルな操作に対応したワークフローを見出す。

7) 評価

発見されたワークフローに基づいて運転支援システムを修正し、プラント運転を実行し、その評価を行う。

3. 学習分類子システムによる知識発見手法の提案

プラントなどの制御ルールにおいては、オペレータが対象プロセスの状態を早急に把握し、制御を行うために、簡潔で明快なルールの発見が必要となる。

3.1. MDL 原理

データとそれを記述するモデルは複雑さの概念で測定することができる²⁾。MDL 基準は、モデルの複雑さとデータの複雑さを最小にするものとして提案された^{3,4)}。以下に MDL 基準を示す。ここで、長さ m のデータ列 $y^m = y_1 \cdots y_m$ の中で m_1, m_0 をそれぞれ $y = 1$ および $y = 0$ の生起数として、 $m = m_1 + m_0$ とする。また、 c_i は前件部の各条件がワイルドカード#の場合 0、それ以外は 1、 t_i は各条件のプロセスデータ分割数を表す。そのとき、データの記述長とモデルの記述長は、以下となる。

$$dataLength = mH\left(\frac{m_1}{m}\right) + \frac{1}{2} \log\left(\frac{m\pi}{2}\right) + o(1),$$

$$modelLength = \sum_{i=1}^k c_i(1 + \log t_i),$$

where, $H(x) = -x \log(x) - (1 - x) \log(1 - x).$

3.2. MDL ベース学習分類子システム

ピッツアプローチを変形して DNF(disjunctive normal form)の条件部と、ひとつの結論部を持つ学習分類子システム(LCS)を採用する⁵⁻¹³⁾。ランダムサンプリングによって事象の分布を推定する学習方法を用い、対象事象の多さに対応する。本モデルでの学習分類子システムの概要を図2に示す。

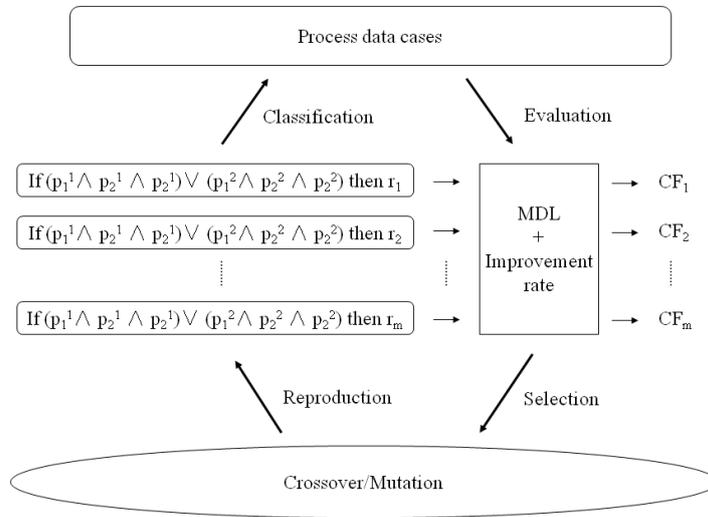


図2 学習分類子システム

はじめに、ランダムに生成したルールを分類子としてプロセスデータを分類する。このルールと分類結果を MDL 基準で評価し、結果を CF_i とする。各分類子に対し、得られた CF_i に基づいてトーナメント選択、交叉および突然変異を実行し、新たな分類子を生成する。分類子数は 200、データ数は 9 タグ × 300min、学習回数は 300、交叉確率は 0.7、突然変異率は 0.5%、数値データは、正規化した後に 25%毎に分割し、それぞれ下降傾向、安定、上昇傾向に分類してルール発見を行った。表 1 に学習分類子のパラメータを示す。ここで、 p は前件部の条件項、 r は後件部の結果項としたとき、分類子の構造は、

$$\mathbf{P} = (p_1^1 \wedge p_2^1 \cdots \wedge p_k^1) \vee (p_1^2 \wedge p_2^2 \cdots \wedge p_k^2) \cdots,$$

$$\mathbf{r} = r_1, r_2, \cdots, r_n.$$

となる。以下に例を示す。

$$((T1 \leq 0.25) \wedge (0.25 < F3 \leq 0.5) \wedge (F4 \text{ is up})) \vee$$

$$((0.5 < T3 \leq 0.75) \wedge (F5 \text{ is down}) \wedge (F1 \text{ is stable}))$$

then $T2 \leq 0.25$.

表 1 学習分類子パラメータ

パラメータ	設定値
プロセスデータアイテム数	9
サンプリング時間	1min
収集期間	300min
前件部データ分割数	4
前件部データ変化傾向	上昇, 安定, 下降
後件部データ分割数	4
分類子数	200
選択手法	トーナメント
交叉率	0.7
突然変異率/遺伝子座	0.005
タブーリスト数	10
近傍距離	ハミング距離

結果 r_i はその事象のとりうる全種類の結果を表し、前件部にヒットしたすべての r_i 毎のヒット数をカウントする。これにより、前件部に一致した事象の後件部事象の信頼度の推定値を得ることができる。図 3 にテスト関数を用いた MDL ベース学習分類子システムの探索結果を示す。Tag1 の時系列データに対して、図に示すように Tag3 の関連をデータの変化として捉え、そのルール発見が行われている。

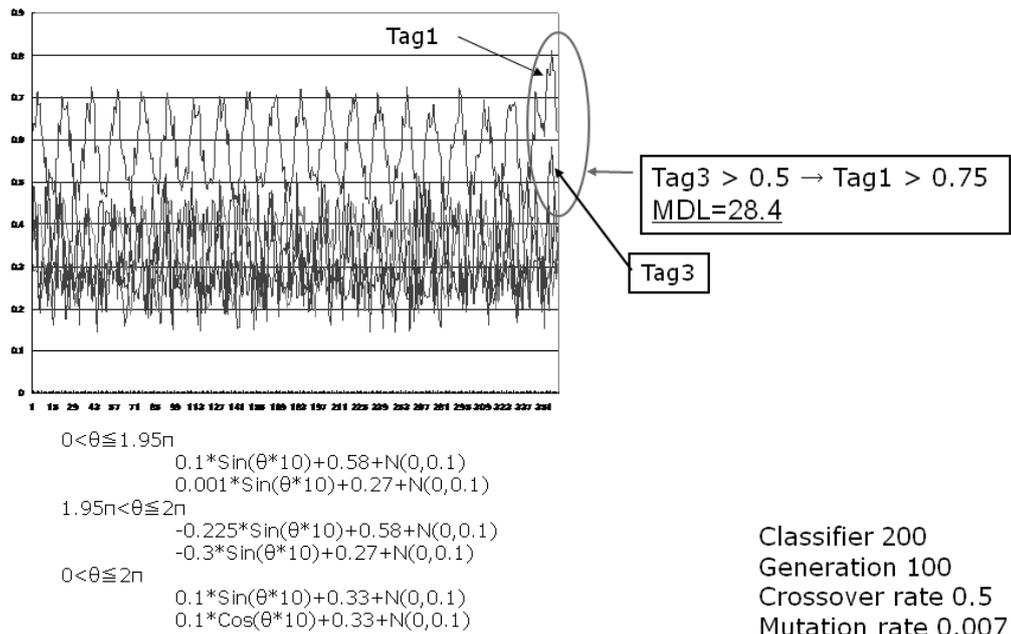


図 3 MDL テスト関数

4. モデル構築と提案手法の適用

2章で示した手順に従ってモデルを構築し、提案手法を適用する。

4.1. 応答モデル

連続プロセスデータの場合、時系列データであってもそれぞれのデータ間に広く相関が認められる。そこで、データ間の関係に着目し正規化された時系列データの相互相関係数を以下の操作により求める。

1対となるプロセスタグの時系列データ x 、 y を選択する。 k をそれぞれのタグの時間シフト量とし、 \bar{x} 、 \bar{y} を平均としたとき、 k 次相互相関係数 $r_{xy}(k)$ は次式にて求まる。

$$r_{xy} = \frac{\sum_{t=k+1}^T (x_{t-k} - \bar{x})(y_t - \bar{y})}{\sqrt{\sum_{t=k+1}^T (x_{t-k} - \bar{x})^2} \sqrt{\sum_{t=k+1}^T (y_t - \bar{y})^2}}$$

$$\max_k r_{xy}(k).$$

図 4 に $\pm 60\text{min}$ の時間シフトを行ったときの相互相関係数の変化を示す。この操作を全てのデータの組み合わせで実行して得られるのが、最大相関係数表とシフト時間表である。その一部を表 2 と表 3 に示す。

表 2 最大相互相関係数

	F4	F2	F3	T2	F1
F4	1.00	0.41	-0.32	0.32	-0.28
F2	0.41	1.00	-0.57	-0.63	-0.46
F3	-0.32	-0.57	1.00	-0.80	-0.53
T2	0.32	-0.63	-0.80	1.00	0.66
F1	-0.28	-0.46	-0.53	0.66	1.00

表 3 最大相関シフト時間

	F4	F2	F3	T2	F1
F4	0	-10	-14	-22	-26
F2	10	0	-5	53	-7
F3	14	5	0	-5	-60
T2	22	-53	5	0	-55
F1	26	7	60	55	0

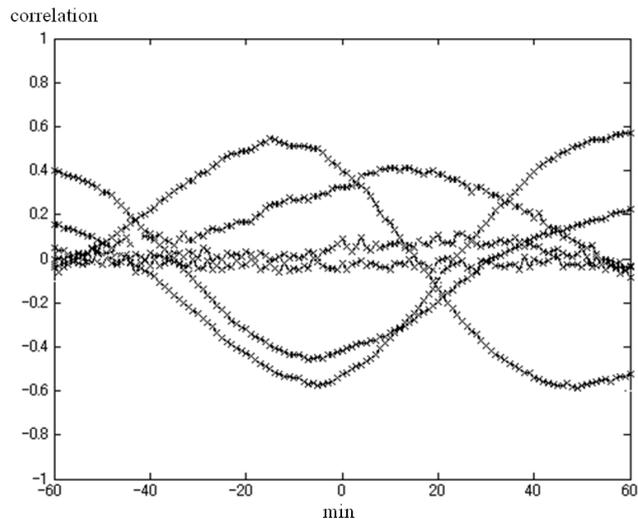


図 4 相互相関係数分析

このようにして多数の時系列タグデータから相関の高いタグを抽出し、プロセス応答モデルを構築することができる。図 5 にその例を示す。このプロセス応答モデル図にあるように、T2 のデータに対して、F2 は 53 分前に -0.6 の相関係数で変化をしており、また F1 は 55 分前に 0.5 の相関係数で変化をしていることが明らかとなっている。これは、逆に見れば、約 50 分前に T2 の変化を予測できることを示している。プロセスタグ間の時系列相関と時間シフトの情報により、プロセス応答の構造を簡潔に示すことができている。

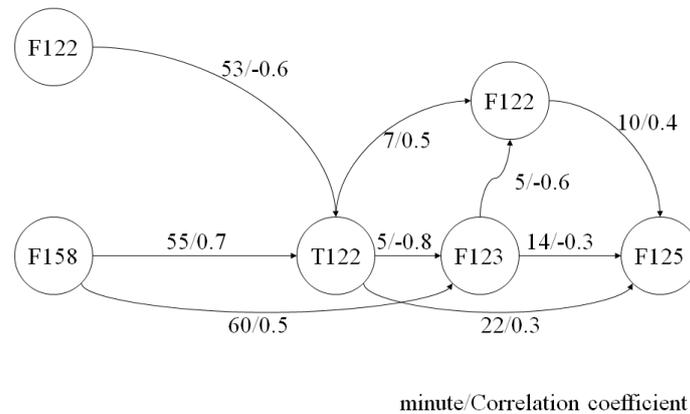


図5 プロセス反応モデル

4.2. 制御ルールの発見的探索

実際の操業では最終品質を安定させる制御ポイントの発見が重要となる。プロセス応答モデルで得られた相関の高いタグのデータを対象に MDL 基準に基づく学習分類子システムによる制御ルールの探索を行う。得られた分類子の例を以下に示す。

$(75\% < F2) \text{ and } (75\% < F3) \text{ then } 50\% < T2$

このときの MDL 値は 32.9 ビットとなっている。

5. 既存手法との比較

時系列データの分析と時系列モデルは統計学分野で発展してきた¹⁴⁾。経済指標の予測を行ったモデルも多く報告されており¹⁵⁾、大量のデータから意味のある情報を見出す方法として多くのデータマイニング手法が研究されている¹⁶⁻²³⁾。

また、学習分類子のルールや遺伝的プログラミングの枝の膨張を押さえながら、一般性を失わないための手法として、MDL 基準を適用した研究が行われている^{24,25)}。これらの手法では、正確さを損なわないようにルールの数を減らすことに成功している。しかし、意外性のあるルールの発見を考慮していないため、事例数が少なく気づきにくい、確実なルールを見出すことが難しい。

一方、興味深さを測るものとして、J-Measure、i-Measure、I-Measures、IShannon-Measure など、いくつかの指標も提案されている²⁶⁻²⁷⁾。しかしこれらは、ルールの持つ情報量や分類階層の深さなどを対象とするものであり、本提案のようなモデルそのものの記述長と分類されたデータ長の両者の変化を同時に考慮するものではなく、また具体的な操作をオペレータへ示すような制御ルールを発見することはできない。決定木などでは大規模な枝が発生してしまうことが多い。標準的な C4.5²⁸⁾を用いて、実際のプロセ

スデータを分析した結果、枝刈り前でノード数が 87、枝刈りを行ってもノード数が 43 となり、大規模なツリーが出現している。プロセスモデルの場合、より確実な制御応答性能を求められ、簡潔なモデルで、確実な情報を導出することが必要となる。表 4 に C4.5 決定木と MDL 基準による学習分類子システムとの比較表を示す。

分類エラー率は枝刈り前後の C4.5 決定木よりも MDL 基準が低くなっている。またルール記述長は、MDL 基準分類子が C4.5 に比べて大幅に小さくなっている。

表 4 決定木 C4.5 と MDL 基準学習分類子の比較

	MDL	MDL-length	Error
C4.5	410.4	297.2bit	12.0%
C4.5rule/pruning	177.0	50.8bit	14.3%
MDL-LCS	121.3	9.8bit	6.7%

6. 熟練者からのワークフローの抽出

プラントの熟練者運転員は自動制御操作以外に、自らの判断で手動で設定値や操作量を随時変更しているのが一般的である。しかし、従来これらの操作は明示的に記録されているわけではなく、一般化され形式化された知識となっていなかった。そこで、本章ではこれらの熟練者の操作を発見するために、これまでの手法に加えて、操作ワークフローを生成する。

はじめに、バルブの開閉やスイッチのオンオフのようなイベントデータと、プロセス制御のための温度や圧力、流量などの設定値を操作したタイムスタンプ付データを収集する。これらは、プロセス制御コンピュータからプロセスデータベースに記録される。次に、原料組成が変化する前後のデータから、上記の手法を用いて制御ルールを探索する。発見されたこれらの制御ルールは、その時系列データに時間シフト以前のタイムスタンプが付いていることから、操作時間順に並べ替えることができる。ここで用いる MDL 基準の学習分類子システムは、時間シフト後のデータに対しては、同一時刻に操作を行ったようになるため、複数のルールを同時に発見する必要がある。ところが、学習分類子システムは、あるひとつのプロセス状態に至るための制御操作に対し、最適化の結果としてもっとも適応度の高い制御ルールを発見することになる。これでは、もともとは異なる時刻に操作したはずのルールを発見することができない。例えば、次のようなひとつのルールとなってしまう。

$$0.5 < F2 \leq 0.75 \text{ then } T2 \leq 0.25 .$$

これは、MDL 原理自身が冗長なルールを除外するように働くことによる。オペレータが実行した複数の操作や、異なるオペレータによる異なる操作が、この方法では発見できない。そこで、タブー探索の手法を取り入れることで、この問題を解決する。タブー探索は、タブーリストと呼ばれる記憶部に、過去の操作を記録し、ある期間同じ操作は行わないようにするヒューリスティック手法である²⁹⁾。図 6 にタブー探索を用いたタブー学習分類子の概要を示す。タブーリストは、親分類子集合と候補集合の間に位置し、タブーリストの格納されている遺伝子に近似した親分類子は、候補集合への通過を拒否され、他の分類子が選択される。一方、タブーリストにある分類子と近似していても、より適応度の高い分類子は、以前の分類子と入れ替えに、タブーリストに格納される。

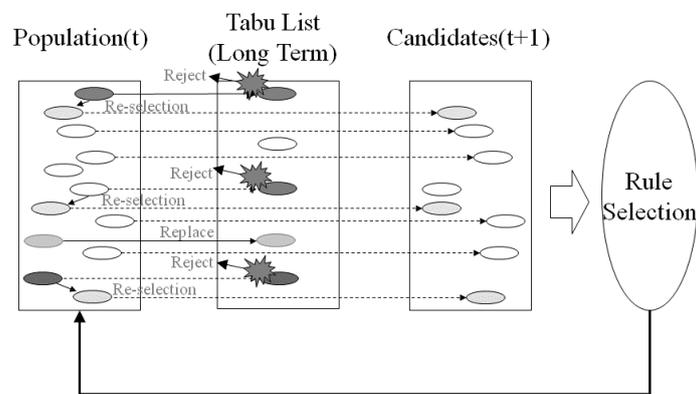


図 6 タブー学習分類子

図 7 は、タブー学習分類子システムを実行した結果である。

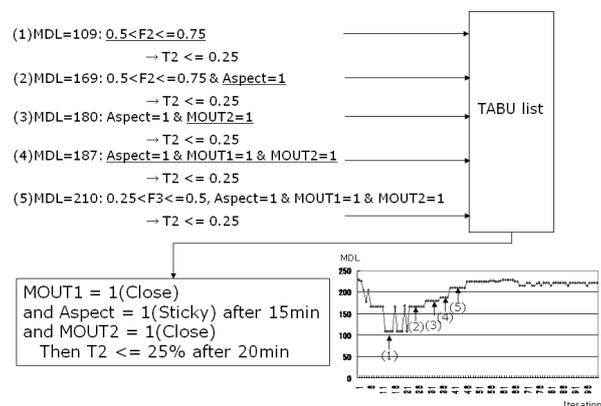


図 7 タブー学習分類子実験結果

最初に発見したルールの MDL は 109 であり、これをタブーリストに格納した後、次に発見されたルールの MDL は 169 であった。最終的に、5 件のルールを発見したところで、MDL が低下しなくなったため、ここで探索を停止した。これらのルールを時系列順に並び替え再構成したワークフローを図 8 に示す。

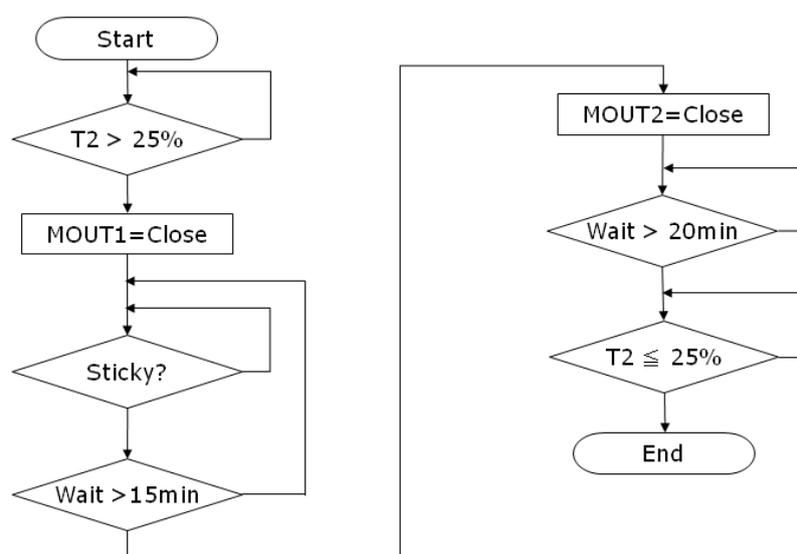


図 8 抽出された熟練者ワークフロー

このワークフローは小規模ではあるが、手動操作記録とプロセスデータからオペレータの操作を抽出したものとなっている。これらの操作は、普段は意識されていないものであり、熟練運転員の暗黙知の発見であると考えられる。本手法の採用により、これらの知識のように、ベテランの運転員から未熟練運転員へ熟練技能を伝承していくことが可能となる。

7 まとめ

本稿では、プロセスデータから運転操作ワークフローを生成するために、MDL 基準とタブー探索を組み込んだ学習分類子システムを用いて、時系列データから、制御および操作ルールを探索する手法を紹介させていただいた。そして、実際のプラントのデータから、熟練者の運転技能を抽出できることを示した。提案手法は以下のフェーズから構成される。(1)プロセスデータ収集 (2)データの正規化 (3)相互相関分析 (4)プロセス応答モデル (5)制御ルールの発見 (6)ワークフローの発見 (7)評価。この手法は比較的単純な方法ではあるが、プロセス時系列データから暗黙知としての熟練者の運転操作を抽出することができ、非熟練者への技能伝承に役立つことを願っている。

参考文献

- 1) J. Rissanen : *Automatica*, 14 (1978), 465.
- 2) C. Adami : *Introduction to Artificial Life*, (1998).
- 3) 山西健司 : *計測自動制御学会誌*, 38 (1999) 7, 420.
- 4) M. Mehta, J. Rissanen, and R. Agrawal : In *Proceedings of the First International Conference on Knowledge Discovery and Data Mining (KDD'95)*, (1995), 216.
- 5) J. H. Holland and J. S. Reitman : *SIGART Bull.*, 63 (1977), 49.
- 6) S. Smith : Ph.D thesis, University of Pittsburgh, (1980).
- 7) S. Smith : In *Proceedings 8th International Joint Conference on Artificial Intelligence*, (1983).
- 8) M. V. Butz, M. Pelikan, X. Llor`a, and D. E. Goldberg : In *GECCO '05: Proceedings of the 2005 conference on Genetic and evolutionary computation*, (2005), 655.
- 9) A. Orriols-Puig, A. Llor`a, and D. E. Goldberg : In *GECCO '10: Proceedings of the 2010 conference on Genetic and evolutionary computation*, (2010), 1023.
- 10) M. V. Butz and O. Herborg : In *GECCO '08: Proceedings of the 2008 conference on Genetic and evolutionary computation*, (2008), 1357.
- 11) M. V. Butz, P. L. Lanzi, and S. W. Wilson : In *GECCO '06: Proceedings of the 8th annual conference on Genetic and evolutionary computation*, (2006), 1457.
- 12) A. Knittel : In *GECCO '10: Proceedings of the 2010 conference on Genetic and evolutionary computation*, (2010), 1871.
- 13) G. Ene and M. Peroumalnaik : In *GECCO '08: Proceedings of the 2008 conference on Genetic and evolutionary computation*, (2008), 2001.
- 14) A. C. Harvey : *Time Series Models*, (1993).
- 15) J. H. Stock and M. W. Watson : *National Bureau of Economic Research*, 2772 (1988).
- 16) P. Adriaans and D. Zantinge : *Data Mining*, (1996).
- 17) A. Barry, J. Holme, and X. Llor`a : *Applications of Learning Classifier Systems*, (2004), 15.
- 18) D. J. Berndt and J. Clifford : *A dynamic programming approach*, (1996), 229.
- 19) A. A. Freitas : *Data mining and knowledge discovery with evolutionary algorithms*, (2002).
- 20) D. E. Goldberg : *Genetic Algorithms in Search, Optimization, and Machine Learning*, (1989).
- 21) M. L. Hetland and P. Saetrom : *Machine Learning*, 58 (2005) 23, 107.
- 22) J. H. Stock and M. W. Watson : *National Bureau of Economic Research*, 2772 (1988).
- 23) S. M. Weiss and N. Indurkha : *Predictive Data Mining, A Practical Guide*, (1997).

- 24) J. Bacardit and J. M. Garrell : In Sixth International Workshop on Learning Classifier Systems (IWLCS-2003), (2003).
- 25) H. Iba, H. de Garis, and T. Sato : Advances in Genetic Programming, (1994), 265.
- 26) J. R. Quinlan : C4.5:Programs for Machine Learning, (1993).
- 27) F. Glover and M. Laguna : Tabu Search, (1997).