

自己組織化マップにおける競合層の動的構成法

日大生産工 ○檜山 昌弘 日大生産工 山内 ゆかり

1 まえがき

情報技術の急速な発達に伴い、世界中でデータの量が複雑かつ膨大になっていることから、データマイニングは重要な課題である。クラスタリングは代表的な分析手法の一つであり、自己組織化マップ (Self-Organizing Map: SOM)[1] の利用が注目されている。SOM は、1982 年に Kohonen によって提案された教師なしニューラルネットワークである。SOM はニューロン間のつながりである位相的順位を保持できるという利点を持ち、類似性のあるデータを分類するのに適している。SOM の改良は、ニューロンの役割による学習率の差別化 [2] や、データの類似度を考慮した近傍関数 [3] など、学習部分に関する研究が主流である。競合層の最適化に関する研究は、入力データから動的に競合層を構成する (Growing Self-Organizing Map: GSOM)[4] などの研究がある。これらの手法の評価には量子化誤差、トポロジー誤差、ニューロン使用率などの指標が用いられるが、各手法によりトレードオフがある。また、用いられる競合層は基本的に格子状である。

長谷川らは、ニューラルネットワークの構造を動的に構成する自己増殖型ニューラルネットワーク (Self-Organizing Incremental Neural Network: SOINN)[5] を提唱し、予めデータの特徴やクラスのないクラスタリングでの有効性を示した。

本研究では、従来の格子状の競合層を用いるのではなく、須藤らの自己増殖の概念に基づき、SOM の学習と共に競合層を動的に構成する動的自己組織化マップを提案する。

2 従来研究

2.1 Self-Organizing Map

Kohonen によって提案された基本の Self-Organizing Map: SOM はニューラルネットワークの一種で、与えられた入力情報の類似度を

マップ上での距離で表現するモデルであり、格子状にニューロンを配置した競合層 (図 1) と、入力そうで構成される。以下に SOM のアルゴリズムを示す。

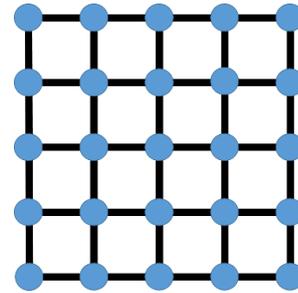


図 1. 格子状の競合層

1. すべての競合層のノード i の参照ベクトル m_i をランダムに初期化する
2. 入力データからランダムに選んだものを入力ベクトル x とする。
3. 式(1)により入力ベクトル x と競合層のノード i の参照ベクトル m_i とのユークリッド距離を求め、その距離が最小となる勝者ノード c を見つける。

$$|x - m_c| = \min |x - m_i| \quad i = 1, 2, \dots, M \quad (1)$$

4. 式(2)を用いて、勝者ノードとその近傍の参照ベクトル m_i を更新する。ここで h_{ci} は近傍関数と呼ばれ、式(3)で定義される。

$$m_i(t+1) = m_i(t) + h_{ci}(t)m_i[x(t) - m_i(t)] \quad (2)$$

$$h_{ci}(t) = \alpha(t) \exp\left(-\frac{|r_c - r_i|^2}{2\sigma^2(t)}\right) \quad (3)$$

2.2 Lazy Self-Organizing Map

従来の SOM の改良としてニューロンの役割による学習率の差別化を行った研究は、原口らによる怠け者を考慮した Lazy Self-Organizing Map: LSOM[2] がある。LSOM は 3 種類のニューロン(働きニューロン, 怠けニューロン, 優柔不断ニューロン)の役割があり、それぞれに異なった学習率係数を用いることにより個性を持たせている。

SOM と同様に競合層と入力層の 2 層で構成され、ランダムに選ばれた $(p \times M)$ 個のニューロ

ンは、怠けニューロンとしての集合 S_{lazy} に分類される。

勝者ニューロン c を決定し、 c が S_{lazy} に含まれている場合、 c は S_{lazy} から取り除かれる。このとき、入力データ x_j から最も遠く、 S_{lazy} の中に存在しないニューロン f が S_{lazy} に入るように選ばれる。つまり、勝者ニューロン c になった怠けニューロンは働きニューロンになり、もう一方のニューロン f が怠けニューロンとなる。概要を図 2 に示す。

その後は、ニューロンの性格と怠けニューロンの割合を考慮して式(3)の学習率に異なった係数を用いて、参照ベクトルを更新していく。

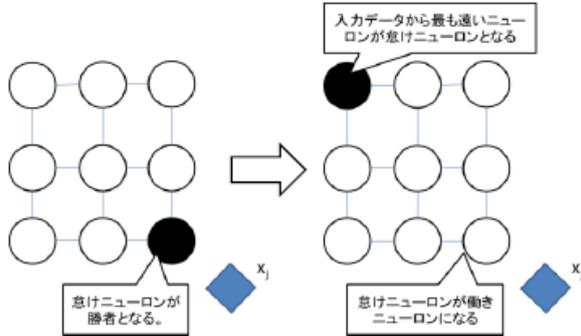


図 2. 怠けニューロンが勝者になったときの働きニューロンとの交代方法

2.3 類似度を考慮した自己組織化マップ

データの類似度を考慮した近傍関数を用いる研究に、竹内らが提案した自己組織化マップ (TSOM) [3]がある。従来の近傍の概念は競合層上の距離に基づいているが、TSOM では参照ベクトルの距離も考慮して近傍関数を定義する。に基づいて学習を行うものである。概念を図 3 に、近傍関数を式(4)に示す。

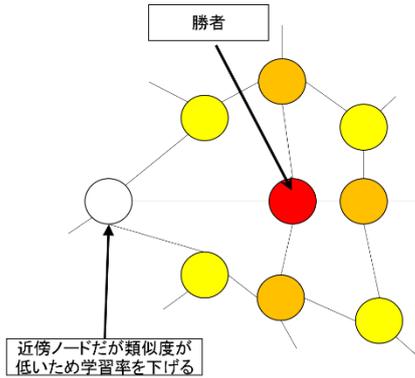


図 3. TSOM

$$h_{ci}(t) = \alpha(t) \exp\left(-\frac{|r_c - r_i|^2}{2\sigma^2(t)}\right) \times \exp\left(-\frac{|x_c - w_i|^2}{2\sigma^2(t)}\right) \quad (4)$$

$|x_c - w_i|^2$ は入力ベクトルと参照ベクトルの特徴差を表している。ただし、勝者ノードの学習は従来の近傍関数による式(3)を用いる。

2.4 自己増殖型ニューラルネットワーク

長谷川らに提唱された自己増殖型ニューラルネットワーク (Self-Organizing Incremental Neural Network: SOINN)[5]は、入力された情報を元にノードの追加や削除、エッジの追加や削除を繰り返し行うことによって、情報数や分類数が不明な状況でも適切なクラスタリングを行う。SOINNのアルゴリズムを下記に示す。

1. 入力データをランダムに 2 つ選択し、その値を参照ベクトル W_j, W_k とするノードを生成する。
2. ランダムに選んだ入力データを入力ベクトル I_c とし、式(1)により 2 つの勝者ノード r, q を見つける。
3. 2 つ勝者ノードの閾値 d_r, d_q をそれぞれ式(5)により算出する。

$$d_i = \begin{cases} \max & k - \text{th node} \in N_i \|W_i - W_k\| (N_i \neq \emptyset) \\ \min & k - \text{th node} \in A \|W_i - W_k\| (N_i \neq \emptyset) \end{cases} \quad (5)$$

4. $\|I_c - W_r\| > d_r, \|I_c - W_q\| > d_q$ のどちらかが成立しない場合は入力ベクトルをノードとして追加する。どちらも成り立つ場合はノード生成せず、第 1, 第 2 勝者をエッジで結合する。ただし元々エッジが引かれていた場合は、年齢を更新する。
5. 第 1 勝者が勝者に選ばれた回数 μ_r に基づき、第 1 勝者は式(7)、第 1 勝者の隣接ノードは式(8)を用いて参照ベクトルを更新する。

$$\Delta W_r = \frac{1}{\mu_r} (I_c - W_r) \quad (6)$$

$$\Delta W_i = \frac{1}{100\mu_r} (I_c - W_i) \quad (7)$$

6. 2.-5.を学習回数分繰り返し、ネットワークの更新周期 λ 毎に、少ないエッジを持つノードと勝者ノード間に選ばれないエッジの削除を行う。

2.5 自己組織化マップの評価指標

SOM の学習及びデータ構造の表現能力を定量的に比較するために、量子化誤差、トポロジー誤差、ニューロン使用率の 3 つの指標を用いる。

量子化誤差 Q_e はそれぞれの入力ベクトルとその勝者との距離の平均を計算した値である。

$$Q_e = \frac{1}{N} \sum_{j=1}^N \|x_j - w_j\| \quad (8)$$

w_j は入力データ x_j に対する勝者ニューロンの重みベクトルである。従って、 Q_e は 0 に近いほど入力状態に近いことを示す。

トポロジー誤差 Te は SOM がどのくらい入力データのトポロジーを保存できているかを示す値である。

$$Te = \frac{1}{N} \sum_{j=1}^N u(x_j) \quad (9)$$

N は入力データの総数である。また、入力データ x_j に対する 1 番目の勝者と 2 番目の勝者が互いに 1 近傍以内なら $u(x_j)$ は 0、それ以外なら 1 となる。つまり、 Te は 0 に近いほど入力データのトポロジーを保存できていることを示す。

ニューロン利用率 U は 1 度以上勝者になったニューロンの割合を示す値である。

$$U = \frac{1}{nm} \sum_{i=1}^{nm} u_i \quad (10)$$

もし、ニューロン i が 1 度以上勝者になったなら、 u_i は 1 となり 1 度も勝者にならなかったなら u_i は 0 となる。つまり、 U は 1 に近いほどより多くのニューロンを有効利用できている、不活性ニューロンが少ないことを示す。

3 提案手法

従来の改善手法において、LSOM は Q_e が改善し、TSOM は U が改善されるが、どちらもトポロジー誤差が悪化する。3つの評価指標がすべてにおいて改善することは難しい。本研究では、競合層を入力データの特徴に基づきノードの配置やエッジの作成を行うことにより、従来の SOM より Q_e と U の 2 つの指標を同時に改善することを目指す。そこで、SOINN の自己増殖アルゴリズムを取り入れた、動的に競合層を構成する自己組織化マップを提案する。以下に提案手法のアルゴリズムを示す。

1. 入力データをランダムに 2 つ選択し、その値を参照ベクトル W_j, W_k とするノードを生成する。
2. ランダムに選んだ入力データを入力ベクトル I_c とし、式(1)により 2 つの勝者ノード r, q を見つける。

3. 2 つ勝者ノードの閾値 d_r, d_q をそれぞれ式(5)により算出する。
4. $\|I_c - W_r\| > d_r, \|I_c - W_q\| > d_q$ のどちらかが成立しない場合は入力ベクトルをノードとして追加し、新しく生成されたノードと第 1, 第 2 勝者間をそれぞれエッジで結合する。どちらも成り立つ場合はノード生成せず、第 1, 第 2 勝者をエッジで結合する。ただし、1 つのノードがもつエッジの最大数は 6 とした。
5. 式(2)を用いて、勝者ノードとその近傍の参照ベクトル m_i を更新する。近傍関数は DSOM では式(11)、DTSOM では式(12)を用いる。

$$h_{ci}(t) = \alpha(t) \exp\left(-\frac{|r_c - r_i|^2}{2l^2(t)}\right) \quad (11)$$

$$h_{ci}(t) = \alpha(t) \exp\left(-\frac{|r_c - r_i|^2}{2l^2(t)}\right) \times \exp\left(-\frac{|x_c - w_i|^2}{2\sigma^2(t)}\right) \quad (12)$$

ここで、 $l(t)$ はネットワークの最長経路長に基づき決定する。

6. 2-5. を学習回数分繰り返す

4 実験環境

本実験では、図 4 に示すターゲットデータ[6]を用いた。データ数は 770 で、中央、周辺円、四隅に偏りのあるデータ配置となっている。

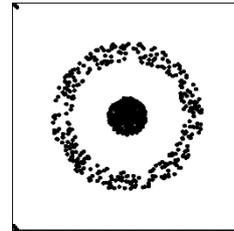


図 4. ターゲットデータ

ノード数 $N=100$ 、学習回数 $T=15400$ とし、それぞれの手法のシミュレーションを 100 回行った平均で各指標を比較する。

従来の格子状の競合層は、格子の対角に第 1 勝者と第 2 勝者が選ばれた場合は、近傍にあるにも関わらず Te が悪するという欠点がある。そこで、正方格子ではなく、三角格子のハニカム構造(図 5)で競合層を構成することにより、従来手法の SOM のトポロジー誤差が改善できるか実験を行った。

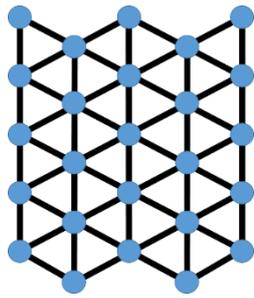


図 5. ハニカム構造の競合層

5 実験結果

通常四角格子で競合層を構成した SOM、LSOM、TSOM の実験結果を表 1 に示す。

表 1. 通常競合層での各指標

	SOM	LSOM	TSOM
Qe	0.018	0.015	0.018
Te	0.094	0.313	0.138
U	0.804	0.804	0.858

表 1 より LSOM は Qe が、TSOM は U が最も良い結果となった。しかし両手法とも Te が従来の SOM より悪化した。

表 2 にハニカム構造での競合層を用いた、実験結果を示す。

表 2. ハニカム構造の競合層での各指標

	SOM	LSOM	TSOM
Qe	0.021	0.020	0.021
Te	0.020	0.022	0.028
U	0.822	0.827	0.871

表 2 より、競合層をハニカム構造で構成することにより、Te の数値が大幅に改善されたことがわかる。これは四角格子では対角にあったノード間も隣接ノードとなることで、データの距離と競合層の距離の概念の乖離が改善されたからだと考えられる。

表 3 に提案手法により動的に競合層を構成した、通常 SOM と TSOM の実験結果を示す。LSOM は競合層の変化には適さないので、提案手法の適用は行わない。

表 3. 提案手法の結果

	DSOM	DTSOM
Qe	0.014	0.014
Te	0.206	0.202
U	0.905	0.917

表 3 より Qe 及び U は従来およびハニカム構造の競合層を用いた結果より改善することができた。この結果より、提案手法では、入力データの特徴に合わせて適切にノードの配置

を行ったことで、不活性ニューロンの割合が減少し、データの偏りを反映した競合層の構成が行えたと考えられる。

しかし Te は従来手法より悪化する結果となった。この原因は、ノードやエッジの削除を行っていないので、学習が不十分な初期ノードやエッジの影響が残っているからだと考えられ、SOINN のようにノードやエッジの削除といった概念を取り入れれば、改善できるのではないかと考えられる。

6 まとめ

本研究では、入力データに基づき、動的に競合層を構成する自己組織化マップのアルゴリズムを提案した。提案手法により、従来の自己組織化マップの各手法より量子化誤差及びニューロン使用率を改善することができた。しかし、トポロジー誤差は従来手法より悪化する結果となった。今後の課題として、ノードやエッジの削除を取り入れることで、トポロジー誤差の改善を試みる事が挙げられる。

「参考文献」

- 1) T.Kohonen, *Self-Organizing Maps*, Springer, vol.30, (1995)
- 2) 原口卓, 松下春奈, 西尾芳文: 「怠け者を考慮した Lazy Self-Organizing Map とその振る舞い」, 電子情報通信学会技術研究報告. NLP, 非線形問題 108(174), 5-10, 2008
- 3) 竹内良, 山内ゆかり: 「自己組織化マップと交差点特徴を用いた手書き文字認識」日本大学生産工学部第 46 回学術講演会 P-13 (2013)
- 4) Daminda Alahakoon, Saman K. Halgamuge, Bala Srinivasan: Dynamic Self-Organizing Maps with Controlled Growth for Knowledge Discovery, IEEE Trans Neural Networks, 11(3), 601-614, (2000)
- 5) 須藤明人, 佐藤彰祥, 長谷川修: 「自己増殖型ニューラルネットワークを用いたノイズのある環境下での追加学習が可能な連想記憶システム」日本神経回路学会誌 Vol.15, No.2, 98-109, (2008)
- 6) Data - Philipps-Universität Marburg - Datenbionik (AG Ultsch): http://www.uni-marburg.de/fb12/datenbionik/data?language_sync=1, (最終アクセス日時:2014 年 10 月 30 日 9 時 0 分)