

スパース中間層を持つ制限ボルツマンマシン分類器

勝亦 利宗[†]

山形大学大学院理工学研究科[†]

安田 宗樹[‡]

山形大学大学院理工学研究科[‡]

1 はじめに

機械学習の分野において過学習 (overfitting) の問題は避けることのできない問題であり、より高い汎化能力を持つ学習モデルの研究が必要とされている。本稿では制限ボルツマンマシン (restricted Boltzmann machine; RBM)[1] を分類問題に特化したモデルである制限ボルツマンマシン分類器 (discriminative RBM; DRBM)[2] に着目し、DRBM の中間層を連続値に拡張することによって汎化能力の向上が見られた先行研究 [3] を踏まえ、新たに中間層にスパース性を持たせた学習モデルを提案する。スパース中間層を持った DRBM は問題に応じて自動的に中間層の素子数を調整できるようになり、問題に柔軟に対応することが可能となる。また、提案したモデルで数値実験を行い、提案モデルの性能を測定する。

2 制限ボルツマンマシン分類器

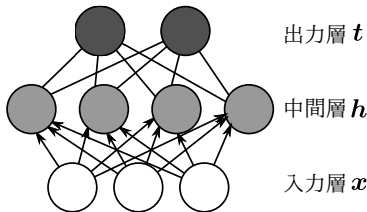


図1 DRBM のグラフ構造。

DRBM のグラフ構造を図1に示す。DRBM は入力層 $\mathbf{x} = \{x_i \in (-\infty, \infty) \mid 1 \leq i \leq X\}$, 中間層 $\mathbf{h} = \{h_j \in \mathcal{X}_h \mid 1 \leq j \leq H\}$, 出力層 $\mathbf{t} = \{t_k \in \{0, 1\} \mid 1 \leq k \leq T\}$ の3層からなる構造を持っている。元論文 [2] では中間層の取る値は2値変数であるが、本稿では文献 [3] を参考にし、 $\mathcal{X}_h = [-1, +1]$ とする。また、出力層 \mathbf{t} は 1-of-K 表現ベクトルとする。1-of-K 表現ベクトルとはベクトルの要素のうち1つが1を取り、その他の要素は0を取るベクトルである。k番目の要素のみが1である 1-of-K ベクトルを $\mathbf{1}_k$ のように表す。このベクトルを使って、分類問題における教師データを表す。

入力層と出力層はデータの説明変数と目的変数に対応する層であり、その素子数 X と T はそれぞれ説明変数の次元数と目的変数のカテゴリ数によって定まる。また、中間層はデータに直接対応しない層であり、素子数 H は問題に応じて適切な数を設定する必要がある。中間層の素子数が多いほどモデルの表現能力は増し、入力変数と出力変数の複雑な関係性を表すことができるようになる一方で、より過学習を増長させる原因

にもなる。

DRBM の確率分布関数は、エネルギー関数

$$E(\mathbf{x}, \mathbf{h}, \mathbf{t}; \theta) = - \sum_{j=1}^H b_j^{(1)} h_j - \sum_{k=1}^T b_k^{(2)} t_k - \sum_{j=1}^H \sum_{i=1}^X w_{i,j}^{(1)} x_i h_j - \sum_{k=1}^T \sum_{j=1}^H w_{j,k}^{(2)} h_j t_k \quad (1)$$

を使って、

$$P(\mathbf{t}, \mathbf{h} \mid \mathbf{x}, \theta) = \frac{1}{Z(\theta)} \exp(-E(\mathbf{x}, \mathbf{h}, \mathbf{t}; \theta)) \quad (2)$$

と定義される。ここで、 $\theta = \{\mathbf{b}^{(1)}, \mathbf{b}^{(2)}, \mathbf{w}^{(1)}, \mathbf{w}^{(2)}\}$ はモデルのパラメータである。それぞれ、 $\mathbf{b}^{(1)}, \mathbf{b}^{(2)}$ はバイアスパラメータ、 $\mathbf{w}^{(1)}, \mathbf{w}^{(2)}$ は層間の結合パラメータを表す。また、 $Z(\theta)$ は分配関数である。

3 スパース正則化

制限ボルツマンマシン分類器

スパース正則化 DRBM (Sparse DRBM; SDRBM) を、式 (3) にスパース正則化項を加え、

$$P(\mathbf{t}, \mathbf{h} \mid \mathbf{x}, \hat{\theta}) = \frac{1}{Z(\hat{\theta})} \exp \left\{ -E(\mathbf{x}, \mathbf{h}, \mathbf{t}; \theta) - \sum_{j=1}^H \sigma(\gamma_j) |h_j| \right\} \quad (3)$$

と定義する。ここで、関数 $\sigma(\gamma_j)$ を

$$\sigma(\gamma_j) := \ln(\exp(\gamma_j) + 1)$$

とする。この関数はソフトプラス関数として知られている関数であり、常に正の値を取る。また、 γ_j はモデルに新たに加わるパラメータである。このパラメータと前節のパラメータ θ を合わせ、 $\hat{\theta} = \{\mathbf{b}^{(1)}, \mathbf{b}^{(2)}, \mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \gamma\}$ とする。新たに加わった項は常に負の値を取り、 γ_j が大きいとき j 番目の中間素子の値が確率分布関数の出力に与える影響を減少させる。適切な γ_j の値は学習によって調節するため、中間素子の取る値をどれだけ軽視するかをモデル自体が行うことができ、それが不要な素子を自動的に削減する仕組みとなっている。さらに、 $\gamma_j \rightarrow -\infty$ のとき、確率分布関数は2節で説明した DRBM と同じものとなる。

4 数値実験

4.1 人工データを用いた数値実験

この節ではバイアスパラメータを0、結合パラメータをXavierの方法 [4] で初期化したパラメータ θ_{gen} から定められる2節の DRBM $P_{\text{gen}}(\mathbf{1}_k \mid \mathbf{x}, \theta_{\text{gen}})$ を生成モデルとし、この生成モデルからサンプリングしたデータを、元論文 [2] の中間

Discriminative restricted Boltzmann machine with sparse hidden layer

[†] Tomu Katsumata, Graduate School of Science and Engineering, Yamagata University

[‡] Muneki Yasuda, Graduate School of Science and Engineering, Yamagata University

層が2値変数である従来のDRBM, 中間層が連続値である2節のDRBM[3], スパース正則化DRBMの3つの学習モデルで学習し, 比較を行う. 学習モデルと生成モデルはどちらも確率分布関数であるので, その差はカルバックライブラー情報量 (Kullback-Leibler divergence; KLD) で表すことができる. 学習モデル $P_{\text{fit}}(\mathbf{1}_k | \mathbf{x}, \theta_{\text{fit}})$ と生成モデル $P_{\text{gen}}(\mathbf{1}_k | \mathbf{x}, \theta_{\text{gen}})$ 間の, 入力データ \mathbf{x} に対する KLD は

$$K(\mathbf{x}) = \sum_{k=1}^T P_{\text{gen}}(\mathbf{1}_k | \mathbf{x}, \theta_{\text{gen}}) \ln \frac{P_{\text{gen}}(\mathbf{1}_k | \mathbf{x}, \theta_{\text{gen}})}{P_{\text{fit}}(\mathbf{1}_k | \mathbf{x}, \theta_{\text{fit}})} \quad (4)$$

となる. この実験では $K(\mathbf{x})$ の平均を, モンテカルロ積分により近似する. また, 生成モデルの構造を $X = 20, H = 100, T = 10$, 学習モデルの構造を $X = 20, H = 200, T = 10$ とし, より過学習が起こりやすい状況で比較を行う. 学習モデルのパラメータ $\mathbf{b}^{(1)}, \mathbf{b}^{(2)}, \mathbf{w}^{(1)}, \mathbf{w}^{(2)}$ の初期値は生成モデルと同様の方法で決定し, スパース正則化DRBMのパラメータ γ の初期値は10とする. γ が大きな状態から学習を始めることで, 必要最低限の素子のみが有効になると考えられるためである. また, γ 以外のパラメータの更新は Adamax[6] を使用し, γ のみ確率的勾配法を使用した. 実験の結果を図2に示す.

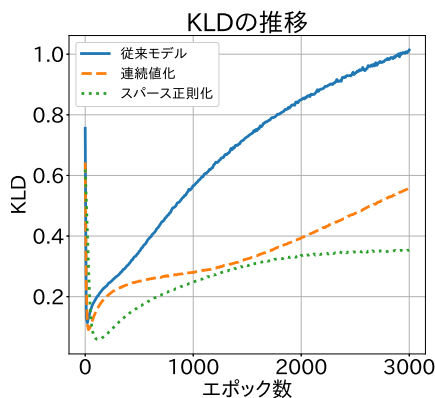


図2 DRBM(実線)[2], 連続値化DRBM(破線)[3], SDRBM(点線)のKLDの推移. プロットは100回の試行の平均である.

図2から, 提案モデルは最もKLDの小さい学習モデルであり, 従来のモデルよりも汎化能力が向上していると言える.

4.2 実データを用いた数値実験

この節では, MNISTの手書き数字画像のデータセットを使用し, 実データに対する数値実験を行う. MNISTは 28×28 のサイズを持つ256階調の画像データであり, 学習用のデータが60,000個, テスト用のデータが10,000個用意されている. この実験では過学習の傾向がより顕著に現れるよう, 学習データを1,000個に減らし, さらにテストデータに $\mathcal{N}(0, 120)$ のガウスノイズを加算した後に標準化し学習を行う. また, 中間層の素子数は200とし, 学習モデルのパラメータは前節と同様に初期化し, 更新する. 実験では学習データに対する誤認識率を訓練誤差, テストデータに対する誤認識率をテスト誤差とし, この2つの値を比較する. 実験の結果を図3に示す.

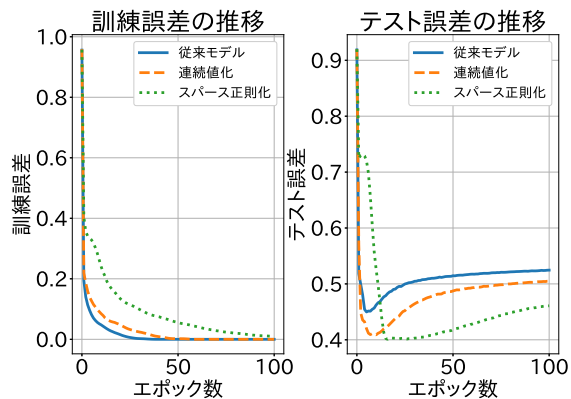


図3 DRBM(実線)[2], 連続値化DRBM(破線)[3], SDRBM(点線)の訓練誤差・テスト誤差の推移. プロットは100回の試行の平均である.

図3より, SDRBMは訓練誤差が上昇し, テスト誤差は減少していることが分かる. 過学習が起こるほど訓練誤差は減少し, テスト誤差は増加するため, SDRBMは過学習を抑制できていると考えられる.

5 まとめ

本稿ではDRBMの確率分布にスパース正則化項を加えた新たな学習モデルを提案し, 従来のDRBMとの比較を行った. 数値実験の結果, SDRBMは総じて従来のDRBMよりも優れた性能を示した. 今回の手法は隠れ素子を持つボルツマンマシンであればいずれも適用可能であると考えられるため, より隠れ層の数を増やしたモデルであるディープボルツマンマシン [5] 等においても, 同様に性能の向上が示されるかどうか今後の課題として挙げられる.

参考文献

- [1] D. Ackley, G. Hinton and T. Sejnowski: A learning algorithm for boltzmann machines, *Cognitive science*, 9(1):147–169, 1985.
- [2] H. Larochelle and Y. Bengio: Classification using discriminative restricted boltzmann machines, *International Conference on Machine Learning*, pp.536–543, 2008.
- [3] Y. Yokoyama, T. Katsumata and M. Yasuda: Restricted Boltzmann Machine with Multivalued Hidden Variables: a model suppressing over-fitting, *The Review of Socionetwork Strategies*, DOI:10.1007/s12626-019-00042-4, pp.1–14, 2018.
- [4] X. Glorot and Y. Bengio, Understanding the difficulty of training deep feedforward neural networks, *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics*, pp.249–256, 2010.
- [5] R. Salakhutdinov and G. Hinton: Deep Boltzmann Machines, *Artificial intelligence and statistics*, pp.448–455, 2009.
- [6] D. Kingma and J. Ba: Adam: a method for stochastic optimization, *Proceedings of the 3rd International Conference on Learning Representations (ICLR 2015)*, 2015.