

学習モデルの情報理論的分析に基づく学習早期終了タイミングの検知

栗林諒[†]関本快士[‡]安田宗樹[§]山形大学大学院理工学研究科[†]山形大学大学院理工学研究科[‡]山形大学大学院理工学研究科[§]

1. はじめに

過学習は機械学習における重要な問題の1つである。過学習を避けるための手法として、検証データ（またはテストデータ）に対する性能の評価が挙げられる。しかしこの手法は少データ学習を行う際、検証データを用意する必要があるため、十分な数の学習データを用意できないなどの問題が起こってしまう。そこで、検証データなしで性能の評価を行うための案としてモデルの内部情報からモデルの状態を推定することが考えられる。この方法が実現すれば、究極的には検証データが不要となる。先行研究[1]より、モデル内部の隠れ素子とデータに対応する可視層との相互情報量が有用な内部情報であることが示唆される。そこで本研究では、そのような相互情報量の解析計算が可能である分類器型制限ボルツマンマシン (Discriminative Restricted Boltzmann Machine (DRBM))[2]を用いて、相互情報量に基づく学習早期終了タイミングの検知指標の模索を行った。

2. 分類器型制限ボルツマンマシン

DRBMは、三層構造の分類器型確率モデルで、入力層 $\mathbf{x} = \{x_i \mid i = 1, 2, \dots, N\}$ 、中間層 $\mathbf{h} = \{h_j \in \{0, 1\} \mid j = 1, 2, \dots, M\}$ 、出力層 $\mathbf{t} = \{\mathbf{1}_k \mid k = 1, 2, \dots, K\}$ の三層で構成されている。ここで、 $\mathbf{1}_k$ は $\mathbf{t} = \{t_1, t_2, \dots, t_K\}$ の第 k 番目の要素のみが1、他の要素が0となるような K 次元のone-hotベクトルを表している。この時、DRBMは以下で定義される。

$$P_{\theta}(\mathbf{t}, \mathbf{h} \mid \mathbf{x}) := \frac{1}{Z_{\theta}(\mathbf{x})} \exp(-E_{\theta}(\mathbf{t}, \mathbf{h}; \mathbf{x})), \quad (1)$$

ここで $Z_{\theta}(\mathbf{x})$ は規格化定数である。また、 $E_{\theta}(\mathbf{t}, \mathbf{h}; \mathbf{x})$ は

エネルギー関数であり以下で定義される。

$$E_{\theta}(\mathbf{t}, \mathbf{h}; \mathbf{x}) := - \sum_{k=1}^K b_k t_k - \sum_{j=1}^M c_j h_j - \sum_{k=1}^K \sum_{j=1}^M w_{k,j} t_k h_j - \sum_{j=1}^M \sum_{i=1}^N v_{j,i} h_j x_i, \quad (2)$$

ここで θ はパラメータ \mathbf{b} 、 \mathbf{c} 、 \mathbf{w} 、 \mathbf{v} をまとめて表記している。通常DRBMを用いてクラス分類を行う際は式(1)を \mathbf{h} について周辺化した式を用いる。

$$\begin{aligned} P_{\theta}(\mathbf{t} \mid \mathbf{x}) &= \sum_{\mathbf{h}} P_{\theta}(\mathbf{t}, \mathbf{h} \mid \mathbf{x}) \\ &= \frac{1}{Z_{\theta}(\mathbf{x})} \exp \left(\sum_{k=1}^K b_k t_k + \sum_{j=1}^M \text{sfp}(\alpha_j(\mathbf{t}, \mathbf{x})) \right) \end{aligned} \quad (3)$$

ここで、

$$\alpha_j(\mathbf{t}, \mathbf{x}) := c_j + \sum_{k=1}^K \sum_{j=1}^M w_{k,j} t_k + \sum_{j=1}^M \sum_{i=1}^N v_{j,i} x_i, \quad (4)$$

である。また、 $\text{sfp}(x) = \ln(1 + \exp(x))$ はソフトプラス関数である。

3. 相互情報量

モデルの学習早期終了タイミングの検知に用いる指標は、中間素子と入力層の相互情報量 (Mutual information between Hidden layers and Input layers (HIMI))

$$M_{\text{in}}^{(j)} := H[h_j] - H[h_j \mid \mathbf{x}], \quad (5)$$

と、中間素子と出力層の相互情報量 (Mutual information between Hidden layers and Output layers (HOMI))

$$M_{\text{out}}^{(j)} := H[h_j] - H[h_j \mid \mathbf{t}], \quad (6)$$

に基づく。ここで $H[h_j]$ 、 $H[h_j \mid \mathbf{x}]$ 、 $H[h_j \mid \mathbf{t}]$ はそれぞれ各確率変数のエントロピーを表しており、以下のように定義される。

$$H[h_j] := - \sum_{h_j \in \{0, 1\}} P_{\theta}(h_j) \ln P_{\theta}(h_j), \quad (7)$$

$$H[h_j \mid \mathbf{x}] := - \sum_{h_j \in \{0, 1\}} \int d\mathbf{x} P_{\theta}(h_j, \mathbf{x}) \ln P_{\theta}(h_j \mid \mathbf{x}), \quad (8)$$

$$H[h_j \mid \mathbf{t}] := - \sum_{h_j \in \{0, 1\}} \sum_{\mathbf{t}} P_{\theta}(h_j, \mathbf{t}) \ln P_{\theta}(h_j \mid \mathbf{t}), \quad (9)$$

式中の確率分布を求めるには $P_{\theta}(\mathbf{x})$ が必要となるが、DRBMの定義式から導出不可能なため、学習データ分布により近似する。相互情報量の観点から、HIMIは入力 \mathbf{x} に対して中間素子 h_j が持つ情報量、HOMIは出力

Detection of early-stopping timing based on information-theoretic analysis of learning model

[†] Ryo Kurabayashi; Graduate School of Science and Engineering, Yamagata University

[‡] Kaiji Sekimoto; Graduate School of Science and Engineering, Yamagata University

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

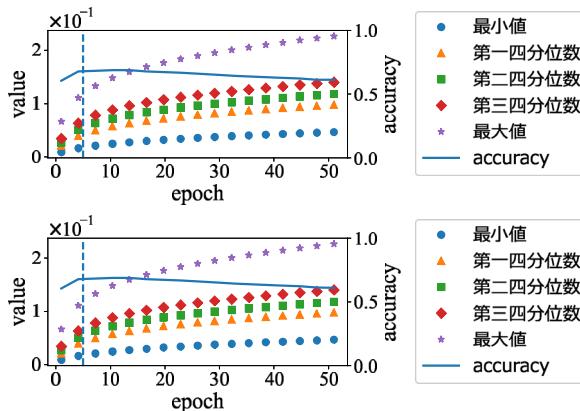


図 1: 学習経過による代表値の推移。上図が HIMI のグラフで、下図が HOMI のグラフである。正答率の最大値との誤差が 0.01 以下となった初めてのエポックを破線で表している。

を決定するうえでの中間素子 h_j の重要度となっている。ここで隠れ変数の重要度を表す HIMI と HOMI の学習経過による推移を観察していく。 $M = 500$ の中間素子数を持つ DRBM に対し、手書き文字のデータセットである MNIST を学習させたときの推移は図 1 のようになる。図 1 より HIMI, HOMI ともに未学習時と過学習時では変化量に違いが観察できる。このことから、過学習を検知するための指標として HIMI と HOMI の変化量が有用であると考察した。

4. 相互情報量を利用した学習早期終了タイミングの検知指標

本研究では過学習を検知するための指標として、HIMI と HOMI の平均値の変化量と、HIMI と HOMI に対するジニ係数の値の変化量を用いた。ジニ係数とは分布の偏り度合いをあらわした指標で、 $[0,1]$ の値をとる。0 の時に完全平等となり、分布が一様分布であることを示し、1 の時に完全不平等となり、分布がデルタ関数的であることを示す。HIMI と HOMI のジニ係数の変化量は分布の形状変化を捉えることができるが分布の平行移動は捉えられない。一方、平均値の変化量は分布の平行移動は捉えられるが分布の形状の変化をとらえられない。よって二つの指標は相補的な関係となっている。

前節と同様の実験を行ったときの HIMI と HOMI の平均値の変化量と、ジニ係数の変化量の推移は図 2 のようになった。平均値の変化量とジニ係数の変化量は、未学習時に大きな推移が起こり、学習完了以降は徐々に 0 に近い値をとることがわかる。同様の実験をファッション商品のデータセットである，Fashion MNIST (F-

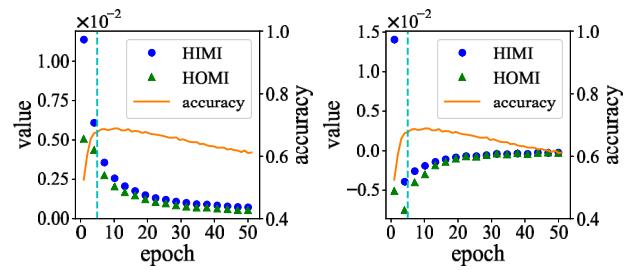


図 2: 学習経過による指標の推移 (MNIST)。左図が平均値の変化量、右図がジニ係数の変化量である。

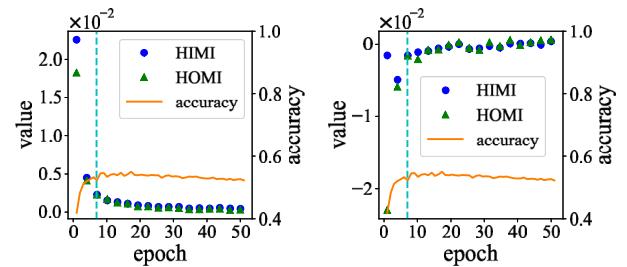


図 3: 学習経過による指標の推移 (F-MNIST)。左図が平均値の変化量、右図がジニ係数の変化量である。

MNIST) を用いて行ったときの結果を図 3 に示す。図 3 においても同様の振る舞いが見られた。これらのことから、学習早期終了のタイミングを平均値の変化量と、ジニ係数の変化量の推移から検知できることが期待される。

5. まとめ

本研究では、過学習を回避し検証データを必要としない手法として、HIMI と HOMI に基づく学習早期終了タイミングの検知指標を提案した。指標は HIMI と HOMI の平均値およびジニ係数の変化量から構成され、数値実験より未学習と過学習の兆候を反映する有用な指標であることが示された。今後は実験を増やすことによって今回の結果がどの程度普遍的であるかを確認することこのような変化が起こる数理的な解釈付けてしていく。

謝辞

本研究は JSPS 科研費 JP24KJ0452 の助成を受けたものである。

文献

- [1] 熊中仁, 安田宗樹. 制限ボルツマンマシンを用いた特徴量抽出と特徴重要度分析. 第 6 回情報処理学会東北支部研究会, 2023.
- [2] Hugo Larochelle and Yoshua Bengio. Classification using discriminative restricted boltzmann machines. In *Proceedings of the 25th International Conference on Machine Learning*, p. 536–543, 2008.