

制約ボルツマンマシン

清家大嗣

平成31年11月1日

1 制約ボルツマンマシン (RBM: Restricted Boltzmann Machine) について ([1] の 8 章を参考)

制約ボルツマンマシンとは、その名が示す通りボルツマンマシンの 1 種であり、ユニット間の結合に制約があるものをいう。例えば、図 1 が示す制約ボルツマンマシンは、隠れ変数 (Hidden Variable) を持つユニット、可視変数 (Visible Variable) を持つユニットは、各々同種のユニット同士では結合を持たない。また、可視変数を持つユニットは全ての隠れ変数を持つユニットとの結合を持つという特徴を持っている (この制約ボルツマンマシンを用いたディープビリーフネットワークをベースをきっかけとして、今日の深層学習は成功した)。本稿では、一般のボルツマンマシンや、制約ボルツマンマシンの持つ性質について理解することを目的とする。

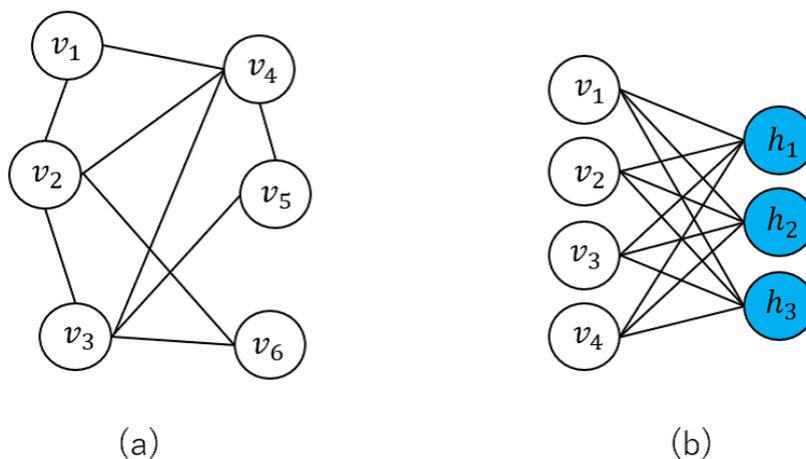


図 1: (a) 一般のボルツマンマシン (b) 制約ボルツマンマシンの一例. h_i は隠れ変数を表し、水色をしたユニットは隠れ変数を持つユニットである.

1.1 ボルツマンマシンについて

ボルツマンマシンでは、複数のユニットが向きを持たない結合により結びついている無向グラフについて考える。各ユニットは 2 種 (0, 1) どちらかの値を持ち i 番目のユニット状態を $x_i (= 0 \text{ or } 1)$

とし、全ユニットの状態を M 次元ベクトル \mathbf{x} により表す。この状態 \mathbf{x} は、ボルツマンマシンでは次式の確率分布により与えられる。

$$p(\mathbf{x}|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp\{-\Phi(\mathbf{x}, \boldsymbol{\theta})\}. \quad (1)$$

ここで、 $\Phi(\mathbf{x}, \boldsymbol{\theta})$ はエネルギー関数と呼ばれる。エネルギー関数 Φ は次式により定義される。

$$\Phi(\mathbf{x}, \boldsymbol{\theta}) = -\mathbf{b}^T \mathbf{x} - \sum_{(i,j) \in \epsilon} w_{ij} x_i x_j \quad (2)$$

ϵ はユニット間の結合であるエッジ集合である。(2) 式より、 $\boldsymbol{\theta}$ はバイアスベクトル \mathbf{b} 、エッジの重み集合 $\{w_{i,j} | (i,j) \in \epsilon\}^1$ により表現される。 $Z(\boldsymbol{\theta})$ は規格化定数²であり、分配関数と呼ばれる。

各ユニットが 0, 1 の 2 種類の状態を持つため、全状態数は 2^M となる。従って、(1, 2) 式により 2^M 通りの状態の生起確率を与えられる。また、一般に (2) 式の形の分布をボルツマン分布 (またはギブス分布) と呼ぶ。

1.1.1 (1, 2) 式に従うボルツマン分布の学習

$p(\mathbf{x}, \boldsymbol{\theta})$ が真の分布 $p_g(\mathbf{x})$ に近くなるように、データ集合 $\mathbf{x}_1, \dots, \mathbf{x}_N$ から最尤推定により $\boldsymbol{\theta}$ を推定する。(1) 式を用いて全データに関する結合確率 $L(\boldsymbol{\theta})$ の対数を取り、それを最大化するような $\boldsymbol{\theta}$ を求める。つまり、次式を最大化する $\boldsymbol{\theta}$ を求める。

$$\log L(\boldsymbol{\theta}) = \sum_{n=1}^N \log p(\mathbf{x}_n|\boldsymbol{\theta}) = \sum_{n=1}^N \{-\Phi(\mathbf{x}_n, \boldsymbol{\theta}) - \log Z(\boldsymbol{\theta})\} \quad (3)$$

この対数尤度関数 (log-likelihood function) のパラメータ $\boldsymbol{\theta}$ に関する勾配は次式により計算できる³。

$$\begin{aligned} \frac{\partial \log L(\boldsymbol{\theta})}{\partial b_i} &= \sum_{n=1}^N \left\{ x_{ni} - \frac{1}{Z(\boldsymbol{\theta})} \cdot \frac{\partial Z(\boldsymbol{\theta})}{\partial b_i} \right\} = \sum_{n=1}^N \left\{ x_{ni} - \frac{1}{Z(\boldsymbol{\theta})} \cdot \sum_{\mathbf{x}} \exp\{-\Phi(\mathbf{x}, \boldsymbol{\theta})\} x_i \right\} \\ &= \sum_{n=1}^N \left\{ x_{ni} - \sum_{\mathbf{x}} p(\mathbf{x}, \boldsymbol{\theta}) x_i \right\} = \sum_{n=1}^N x_{ni} - N E_{\boldsymbol{\theta}}[x_i] \end{aligned} \quad (4)$$

$$\begin{aligned} \frac{\partial \log L(\boldsymbol{\theta})}{\partial w_{ij}} &= \sum_{n=1}^N \left\{ x_{ni} x_{nj} - \frac{1}{Z(\boldsymbol{\theta})} \cdot \frac{\partial Z(\boldsymbol{\theta})}{\partial w_{ij}} \right\} = \sum_{n=1}^N \left\{ x_{ni} x_{nj} - \frac{1}{Z(\boldsymbol{\theta})} \cdot \sum_{\mathbf{x}} \exp\{-\Phi(\mathbf{x}, \boldsymbol{\theta})\} x_i x_j \right\} \\ &= \sum_{n=1}^N \left\{ x_{ni} x_{nj} - \sum_{\mathbf{x}} p(\mathbf{x}, \boldsymbol{\theta}) x_i x_j \right\} = \sum_{n=1}^N \{x_{ni} x_{nj}\} - N E_{\boldsymbol{\theta}}[x_i x_j] \end{aligned} \quad (5)$$

ここで、 $E_{\boldsymbol{\theta}}[\cdot]$ は $\boldsymbol{\theta}$ が与えられた場合に定まるモデル分布から計算される \cdot の期待値である。従って、勾配法を実行する際、その時点での $\boldsymbol{\theta}$ から $p(\mathbf{x}|\boldsymbol{\theta})$ は計算可能なように思える。しかし、 M が大きくなると期待値計算に必要な計算回数が指数関数的に増加する (例えば $M = 10$ の場合、全ての $p(\mathbf{x}|\boldsymbol{\theta})$ を計算し終えるまでに 1,048,576 回の施行が必要となる)。従って、(4, 5) 式を用いて、素朴に勾配法を実行するのは現実問題として無理である。

1.1.2 ギブスサンプリングについて

前節のように組み合わせ爆発により期待値計算が困難な場合に、近似的に期待値を計算する手法について考える。一般に多変数確率分布からそれに従う変数を生成することは困難であるが、ボル

¹無効グラフであるため $w_{ij} = w_{ji}$ となる。

² $Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \Phi(\mathbf{x}, \boldsymbol{\theta})$ ($\sum_{\mathbf{x}} = \sum_{x_1=0,1} \sum_{x_2=0,1} \dots \sum_{x_M=0,1}$)

³ただし、 x_{ni} とは n 番目のデータの第 i 成分要素とする

ツマンマシンは局所マルコフ性という性質があり、ギブスサンプリング (Gibbs Sampling) と呼ばれる手法で効率的に確率変数ベクトル \mathbf{x} をサンプリングできる。

ユニット i を除いた全ユニットの変数を並べたベクトルを \mathbf{x}_{-i} と表現する。あるパラメータ θ から生成されるモデル分布 $p(\mathbf{x}|\theta)$ から \mathbf{x} を効率的にサンプルするため、各ユニットがベクトル \mathbf{x}_{-i} の値を取る場合のユニット i に関する事後分布について考える。ユニット i の値が x_i となる確率は次式で表される。

$$\begin{aligned} p(x_i|\mathbf{x}_{-i}, \theta) &= \frac{p(\mathbf{x}, \theta)}{\sum_{x_i=0,1} p(\mathbf{x}, \theta)} \\ &= \frac{\exp\{(b_i + \sum_{j \in \mathcal{N}_i} w_{ij} x_j) x_i\}}{1 + \exp\{(b_i + \sum_{j \in \mathcal{N}_i} w_{ij} x_j)\}} \end{aligned} \quad (6)$$

従って、モデルのパラメータベクトル θ とユニット i 以外のユニットの値さえ与えられていれば、ユニット i を (6) 式の分布に従う確率変数としてサンプリングすることができる。これを $i = 1, \dots, M$ に対して順番に繰り返してサンプリングする手法がギブスサンプリングと呼ばれる。例えば、このステップを t 純目まで繰り返した場合、その時のユニット i は次式の分布に従う確率変数としてサンプリングされる。

$$p(x_i|\mathbf{x}_{-i}^{(t)}, \theta) = p(x_i|x_1^{(t)}, \dots, x_{i-1}^{(t)}, x_{i+1}^{(t-1)}, \dots, x_M^{(t-1)}, \theta) \quad (7)$$

このサンプリングした値は t を十分に大きくとれば、モデル分布 $p(\mathbf{x}, \theta)$ に従う [2]。また、このサンプリング列から複数のサンプル $\mathbf{x}^{(t_1)}, \mathbf{x}^{(t_2)}$ を得たい場合、サンプリングする間隔 $|t_2 - t_1|$ は十分に離れている必要がある (独立にするため)。

(6, 7) 式に従うプロセスを繰り返すことで、(4, 5) 式の期待値を近似的に計算できる。つまり $O(a^M)$ ($a > 0$) にかかる計算コストを $O(MT)$ (T はギブスサンプリングの繰り返し回数) に減らすことができる⁴。

1.2 隠れ変数を持つボルツマンマシン

図 1 (b) にあるような、全変数がシステム外部に公開されていないボルツマンマシンもある。つまり、可視変数ベクトル \mathbf{v} 、隠れ変数ベクトル \mathbf{h} が与えられた場合、前節の \mathbf{x} に当たるものが \mathbf{v}, \mathbf{h} であるが、 \mathbf{h} は観測不可能というケースである。

1.2.1 隠れ変数を持つボルツマンマシンの学習

最初にモデル分布 $p(\mathbf{v}, \mathbf{h}|\theta)$ を最大にする最尤法を観測可能な可視変数に対して用いるため、観測していない隠れ変数 \mathbf{h} を周辺化する⁵。

$$p(\mathbf{v}|\theta) = \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}|\theta) \quad (8)$$

従って、尤度関数の対数を取ることで次式を最大化する問題ととらえることができる。

$$\log L(\theta) \propto \frac{1}{N} \sum_{n=1}^N \left[\log \sum_{\mathbf{h}} \exp\{-\Phi(\mathbf{v}_n, \mathbf{h}, \theta)\} - \log Z(\theta) \right] \quad (9)$$

⁴それでもギブスサンプリングのコストは高いことには注目しておきたい。

⁵こうすることで、 \mathbf{v} のみの分布を最大化する問題として置き換える

右辺第一項を θ で微分すると次式のようなになる (ここで z_i, z_j は \mathbf{v}, \mathbf{h} を連結したベクトル \mathbf{z} の成分である).

$$\begin{aligned} \frac{\partial}{\partial w_{ij}} \log \sum_{\mathbf{h}} \exp\{-\Phi(\mathbf{v}_n, \mathbf{h}, \theta)\} &= \frac{\sum_{\mathbf{h}} z_i z_j \exp\{-\Phi(\mathbf{v}_n, \mathbf{h}, \theta)\}}{\sum_{\mathbf{h}} \exp\{-\Phi(\mathbf{v}_n, \mathbf{h}, \theta)\}} \\ &= \sum_{\mathbf{h}} z_i z_j p(\mathbf{h}|\mathbf{v}_n, \theta) \end{aligned} \quad (10)$$

従って, (9) 式は次のように書き換えることができる.

$$\log L(\theta) \propto E_{\theta}[z_i z_j]_{data} - E_{\theta}[z_i z_j]_{model} \quad (11)$$

ただし, $E_{\theta}[z_i z_j]_{data} = 1/N \sum_{n=1}^N \sum_{\mathbf{h}} z_i z_j p(\mathbf{h}|\mathbf{v}_n, \theta)$ である. \mathbf{h} の次元数を H とした場合の第一項の計算量を考慮すると, 上式の勾配計算のための期待値計算の計算量は $O(a^H + b^M)$ となる. 従って, (4, 5) 式より更に計算が困難となる.

1.3 制約ボルツマンマシン (RBM: Restricted Boltzmann Machine)

制約ボルツマンマシンは, 隠れ変数を持つボルツマンマシンの一種である. 図 1 にあるように, 可視変数間, 隠れ変数間の結合はなく, 任意の可視変数と隠れ変数の間には必ず結合がある. RBM ではこの特別な構造により, 効率的なギブスサンプリングが可能となる.

1.3.1 RBM のグラフ構造を利用したギブスサンプリング

RBM の分布関数で用いる新たな記法について説明する. RBM のエネルギー関数と状態の分布関数 $p(\mathbf{v}, \mathbf{h}, \theta)$ は次式により定義される (i, j は各々, 可視変数と隠れ変数のインデックスである. a_i, b_j は各々のユニットのバイアスである).

$$\Phi(\mathbf{v}, \mathbf{h}, \theta) = -\sum_i a_i v_i - \sum_j b_j h_j - \sum_i \sum_j w_{ij} v_i h_j \quad (12)$$

$$p(\mathbf{v}, \mathbf{h}|\theta) = \frac{1}{Z(\theta)} \exp\{-\Phi(\mathbf{v}, \mathbf{h}, \theta)\} \quad (13)$$

従って, 可視変数ベクトル \mathbf{v} による隠れ変数ベクトルの条件付き分布は次式で定義される⁶.

$$\begin{aligned} p(\mathbf{h}|\mathbf{v}, \theta) &= \prod_j p(h_j|\mathbf{v}, \theta) \\ &= \frac{\exp\{(b_j + \sum_i w_{ij} v_i) h_j\}}{1 + \exp(b_j + \sum_i w_{ij} v_i)} \end{aligned} \quad (14)$$

つまり, 各ユニット変数 h_j, v_i は次のように, \mathbf{v} により定まるベルヌーイ分布に従うと考えられる ($\sigma(x) = 1/(1 + \exp(-x))$ はロジスティック関数である).

$$\begin{aligned} \prod_j p(h_j = 1|\mathbf{v}, \theta) &= \frac{\exp(b_j + \sum_i w_{ij} v_i)}{1 + \exp(b_j + \sum_i w_{ij} v_i)} \\ &= \sigma\left(1 + \exp(b_j + \sum_i w_{ij} v_i)\right) \end{aligned} \quad (15)$$

可視変数の場合も同様に, 次式で表せる.

$$\prod_j p(v_i = 1|\mathbf{h}, \theta) = \sigma\left(1 + \exp(a_i + \sum_j w_{ij} h_j)\right) \quad (16)$$

⁶隠れ変数同士に結合がない (独立) ため, \mathbf{v} が与えられていることと合わせて 2 段目のような変形が可能である. 分母の第一項は $h_j = 0$ のケースを考え, 第二項は $h_j = 1$ のケースについて考えている項である.

従って, $\mathbf{v}^{(0)}$ を初期値として, (15, 16) 式を用いて逐次的に次のようにサンプリングしていくことができる.

$$\mathbf{v}^{(0)} \rightarrow \mathbf{h}^{(0)} \rightarrow \mathbf{v}^{(1)} \rightarrow \mathbf{h}^{(1)} \rightarrow \dots \quad (17)$$

このサンプリング手法では, $\mathbf{v}^{(t)}, \mathbf{h}^{(t)}$ の値が変わらないことから並列計算が可能であり, 幾分か計算効率を改善することが可能である. しかし, (6) 式と (15, 16) 式を比較しても分かるように結合の数だけ積の計算をする必要があるのは変わらないため,それほど大きな改善にはならないと思われる.

1.3.2 RBM の学習

(17) 式を用いて次のように RBM の勾配を計算することができる (\mathbf{h} による周辺化に注意. 隠れ変数同士の結合がないため, $p(\mathbf{h}|\mathbf{v}, \boldsymbol{\theta}) = \prod_j p(h_j|\mathbf{v}_n, \boldsymbol{\theta})$. (10) 式を参照).

$$\begin{aligned} \frac{1}{N} \frac{\partial \log L}{\partial w_{ij}} &= \frac{1}{N} \sum_{n=1}^N \left\{ \sum_{h_j=0,1} v_{ni} h_j p(h_j|\mathbf{v}_n, \boldsymbol{\theta}) \right\} - \sum_{\mathbf{v}, \mathbf{h}} v_i h_j p(\mathbf{v}, \mathbf{h}|\boldsymbol{\theta}) \\ &= \frac{1}{N} \sum_{n=1}^N \left\{ v_{ni} p(h_j = 1|\mathbf{v}_n, \boldsymbol{\theta}) \right\} - \sum_{\mathbf{v}, \mathbf{h}} v_i h_j p(\mathbf{v}, \mathbf{h}|\boldsymbol{\theta}) \end{aligned} \quad (18)$$

同様に, バイアスも次のように計算できる.

$$\frac{1}{N} \frac{\partial \log L}{\partial a_i} = \frac{1}{N} \sum_{n=1}^N v_{ni} - \sum_{\mathbf{v}, \mathbf{h}} v_i p(\mathbf{v}, \mathbf{h}|\boldsymbol{\theta}) \quad (19)$$

$$\frac{1}{N} \frac{\partial \log L}{\partial b_j} = \frac{1}{N} \sum_{n=1}^N \sum_{h_j=0,1} h_j p(h_j|\mathbf{v}_n, \boldsymbol{\theta}) - \sum_{\mathbf{v}, \mathbf{h}} h_j p(\mathbf{v}, \mathbf{h}|\boldsymbol{\theta}) \quad (20)$$

つまり, 各重みとバイアスの更新式は次式により表せるので, (17) 式を用いて学習が実行可能である.

$$\Delta w_{ij} \propto \mathbf{E}_{\boldsymbol{\theta}}[v_i h_j]_{data} - \mathbf{E}_{\boldsymbol{\theta}}[v_i h_j]_{model} \quad (21)$$

$$\Delta a_i \propto \mathbf{E}_{\boldsymbol{\theta}}[v_i]_{data} - \mathbf{E}_{\boldsymbol{\theta}}[v_i]_{model} \quad (22)$$

$$\Delta b_j \propto \mathbf{E}_{\boldsymbol{\theta}}[h_j]_{data} - \mathbf{E}_{\boldsymbol{\theta}}[h_j]_{model} \quad (23)$$

参考文献

- [1] 岡谷貴之. 2015. 『深層学習』 講談社.
- [2] K. P. Murphy. *Machine Learning: A probabilistic perspective*. MIT Press, 2012.