

Restricted Boltzmann Machine の導出に至る自分用まとめ

齋藤 真樹

2012年8月18日

1 問題設定

可視素子 (visible node) の集合 $\mathbf{v} = (v_1, \dots, v_L), \forall v_i \in \{0, 1\}$, 隠れ素子 (hidden node) の集合 $\mathbf{h} = (h_1, \dots, h_M), \forall h_j \in \{0, 1\}$ から成る結合確率 $p(\mathbf{v}, \mathbf{h})$ を以下のように定義する.

$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} \exp(-E(\mathbf{v}, \mathbf{h})) \quad (1)$$

$$E(\mathbf{v}, \mathbf{h}) = - \sum_i b_i v_i - \sum_j c_j h_j - \sum_i \sum_j v_i W_{ij} h_j \quad (2)$$

このような $p(\mathbf{v}, \mathbf{h})$ を一般に Restricted Boltzmann Machine (RBM) と呼ぶ. b_i, c_j, W_{ij} はそれぞれ可視素子のバイアス, 隠れ素子のバイアス, ウェイトパラメータと呼ばれている. Z は正規化定数あるいは分配関数と呼ばれており,

$$Z = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h})) \quad (3)$$

の関係を持つ.

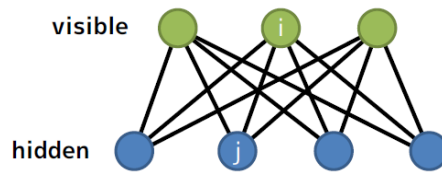


図1 RBM のグラフィカルモデル. 計算上の理由から可視素子間, 隠れ素子間の結合は存在しない.

ここで, 我々は手元に可視素子に関する N 個の観測データ $\{v^1, v^2, \dots, v^N\}$ を手に入れているが, 隠れ素子に関しては手に入れていないものとする. このような制約の中からどうにかして, 我々が求めるべきパラメータである b_i, c_j, W_{ij} を推測したい.

この場合, 典型的な戦略としてはまず $p(\mathbf{v}, \mathbf{h})$ を \mathbf{h} に関して積分を行うことによって, 周辺分布 $p(\mathbf{v}) = \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h})$ を得る. この分布は \mathbf{h} に依存しておらず, \mathbf{v} それ自身によってのみ決定される分布である. この分布を用いて最尤推定を行うことで, パラメータを推測する戦略をとることにしよう.

2 計算の概要

以降は具体的な計算について示す．まず，最尤推定では対数尤度を，パラメータ θ に関する最大化を行うことで表現する．

$$J \equiv \left\langle \ln \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}) \right\rangle_q = \left\langle \ln \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \right\rangle_q - \ln Z \quad (4)$$

ここで， $\langle \cdot \rangle_q$ のような形を確率分布 $q(\mathbf{v})$ の期待値として表現する．つまり

$$\langle f(\mathbf{v}) \rangle_q = \sum_{\mathbf{v}} f(\mathbf{v}) q(\mathbf{v}) \quad (5)$$

である． $q(\mathbf{v})$ は観測データに関する確率分布 $q(\mathbf{v})$

$$q(\mathbf{v}) = \frac{1}{N} \sum_k \delta(\mathbf{v} - \mathbf{v}^k) \quad (6)$$

であり， $\delta(x)$ はディラックのデルタ関数である．このようにして定義された J を最大化するような b_i, c_j, W_{ij} を求める戦略が最尤推定である．

パラメータの推定値を求めるために，勾配法を用いて最適化を行うこととする．そのためにはまず，パラメータに関する微分 $\partial J / \partial \theta$ が必要である．ゆえに，任意のパラメータ θ に関して微分を取ることによって下記を得る．

$$\frac{\partial J}{\partial \theta} = - \left\langle \frac{1}{\sum_{\mathbf{h}} e^{-E}} \sum_{\mathbf{h}} \frac{\partial E}{\partial \theta} e^{-E} \right\rangle_q + \frac{1}{Z} \sum_{\mathbf{v}} \sum_{\mathbf{h}} \frac{\partial E}{\partial \theta} e^{-E} \quad (7)$$

ここで，条件付き確率の定義

$$p(\mathbf{h}|\mathbf{v}) = \frac{e^{-E(\mathbf{v}, \mathbf{h})}}{\sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}} \quad (8)$$

並びにマルコフ確率場の定義を用いることによって，上式は以下のように簡略化できる．

$$\frac{\partial J}{\partial \theta} = - \sum_{\mathbf{v}} \sum_{\mathbf{h}} \frac{\partial E}{\partial \theta} p(\mathbf{h}|\mathbf{v}) q(\mathbf{v}) + \sum_{\mathbf{v}} \sum_{\mathbf{h}} \frac{\partial E}{\partial \theta} p(\mathbf{v}, \mathbf{h}) \quad (9)$$

$$\equiv - \left\langle \frac{\partial E}{\partial \theta} \right\rangle_{\text{data}} + \left\langle \frac{\partial E}{\partial \theta} \right\rangle_{\text{model}} \quad (10)$$

ここで， $\langle \cdot \rangle_{\text{data}}, \langle \cdot \rangle_{\text{model}}$ は，それぞれ確率分布 $p_{\text{data}}(\mathbf{v}, \mathbf{h}) = p(\mathbf{h}|\mathbf{v})q(\mathbf{v})$ と $p_{\text{model}}(\mathbf{v}, \mathbf{h}) = p(\mathbf{v}, \mathbf{h})$ に対する期待値を表す．この表記を用いることで，パラメータ θ に関する微分操作の本質的意味合いをより簡潔に表現できるようになった．

一番初めの定義へと戻ろう．我々が現在持っている観測データは \mathbf{v}^k のみであり，隠れ素子についての観測データは持ち合わせていない．そのため， \mathbf{v} に関する周辺化を行うことによって隠れ素子の変数を消去し，その枠組の中で最大化を行うこととした．これは，本来得られてることが望ましい観測データの結合確率 $\frac{1}{N} \sum_k \delta(\mathbf{v} - \mathbf{v}^k) \delta(\mathbf{h} - \mathbf{h}^k)$ を用いるのではなく，現在のパラメータの元で推測した結合確率 $p(\mathbf{h}|\mathbf{v})q(\mathbf{v})$ を用いて計算を行なっていることと等価である．

幸いにも，Restricted Boltzmann Machine に関しては可視素子間，隠れ素子間についての結合が存在していないため，条件付き確率 $p(\mathbf{v}|\mathbf{h}), p(\mathbf{h}|\mathbf{v})$ の厳密解は容易に計算できる．しかも，条件付き確率下での各引

数は独立であり ($p(\mathbf{v}|\mathbf{h}) = \prod_i p(v_i|\mathbf{h}), p(\mathbf{h}|\mathbf{v}) = \prod_j p(h_j|\mathbf{v})$), このことも RBM の計算の省力化に貢献している.

それぞれの条件付き確率を計算することによって下記を得る.

$$p(v_i = 1|\mathbf{h}) = \sigma(b_i + \sum_j W_{ij}h_j) \quad (11)$$

$$p(h_j = 1|\mathbf{v}) = \sigma(c_j + \sum_i v_i W_{ij}) \quad (12)$$

ここで, $\sigma(x) = 1/(1 + \exp(-x))$ はシグモイド関数である.

3 Contrastive Divergence

上式の第一項 $\langle \cdot \rangle_{\text{data}}$ は比較的簡単に計算できるが, 一般に第二項 $\langle \cdot \rangle_{\text{model}}$ に関する計算は非常に困難である場合が多い. 何故なら, 第二項はあらゆる v, h に関する和を取らなければならないため, ノード数の増加によって状態量が指数的に爆発してしまうためである.

このような直接的に計算不可能な期待値の計算を行う際, 大まかに分けて 2 種類の方法が存在する. まず一つ目はマルコフ連鎖モンテカルロ (MCMC) を用いることによって, 力技で計算を行う手法と, 平均場近似などを用いて近似的に解を求めるという手法である. 一般に前者は計算速度, 後者は近似精度が問題となる場合が多いが, 今回用いる Contrastive Divergence (CD) では MCMC の手法を用いてはいるものの高速に計算を行うことができ, しかもその精度は経験的にも理論的にも良好であることが確認されている.

具体的な手法を下記に示す. まず, MCMC では任意の初期値 v_0, h_0 から, 緩和に至るまでに十分な回数だけ遷移を行うことによって, $p(v, h)$ に従ったサンプルを無数に生成し, その平均を取ることによって近似的に期待値を得る. しかしこの手法は非常に多くの回数の遷移を行う必要があるため, その計算速度は比較的遅いことが欠点であった.

一方, Contrastive Divergence では同じように MCMC を用いて遷移を行っているが, 2 つの点が MCMC とは異なっている. まず, 遷移させる回数は無数ではなく, A 回の非常に少ない遷移で構わない. 大抵の場合は $A = 1$ 回で十分である. これは隠れ素子, 可視素子間の結合が存在していないため, MCMC における緩和時間が比較的高速であるからである. 次に, 初期値には N 個のサンプルである $\{v^1, \dots, v^N\}$ を用いる. この僅かなルール付け足しによって, CD は高速かつ高精度に RBM の勾配を求めることができる.

なお, 条件付き確率を解析的に求められるという利点から, MCMC による遷移は一般に Gibbs Sampling を用いて行われる. また, 遷移によって用いるサンプルは v についてのみ用うことで確率分布 $\hat{q}(v) = \frac{1}{N} \sum_k \delta(v - \hat{v}^k)$ を生成し, これと h についての条件付き確率 $p(h|v)$ の積によって結合確率 $\hat{q}(v)p(h|v)$ を表現している.

Contrastive Divergence の更新式

$A = 1$ の場合 , Contrastive Divergence を用いた第二項の更新は下記の通りとなる .

1. 初期値に $\{v_1, \dots, v_N\}$ をセットする .
2. 条件付き確率 $p(\mathbf{h}|\mathbf{v})$ を用いて , 隠れ素子の状態 $\{\hat{h}^1, \dots, \hat{h}^N\}$ を求める .
3. 2. で得た隠れ素子の状態を用いて , 条件付確率 $p(\mathbf{v}|\mathbf{h})$ を用いて可視素子の状態 $\{\hat{v}^1, \dots, \hat{v}^N\}$ を求める .
4. 3. で得た可視素子の状態と条件付確率から結合確率

$$p_{\text{model}}(\mathbf{v}, \mathbf{h}) = \frac{1}{N} \sum_k \delta(\mathbf{v} - \hat{\mathbf{v}}^k) p(\mathbf{h}|\hat{\mathbf{v}}^k) \quad (13)$$

を求め , 計算を行うことによって近似的に第二項を得る .

4 パラメータの計算

以降の章では具体的なパラメータ b_i, c_j, W_{ij} についての計算を行い , 実際に導出された式について書き下す .
まず , W_{ij} についての計算を行う . $\partial E / \partial W_{ij} = -v_i h_j$ より ,

$$\frac{\partial J}{\partial W_{ij}} = \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}} \quad (14)$$

ここで , 第一項について ,

$$\langle v_i h_j \rangle_{\text{data}} = \sum_{\mathbf{v}} \sum_{\mathbf{h}} v_i h_j q(\mathbf{v}) p(\mathbf{h}|\mathbf{v}) \quad (15)$$

$$= \sum_{\mathbf{h}} \left\{ \frac{1}{N} \sum_k v_i^k h_j p(\mathbf{h}|\mathbf{v}^k) \right\} = \frac{1}{N} \sum_k v_i^k \sum_{h_j} h_j p(\mathbf{h}|\mathbf{v}^k) \quad (16)$$

$$= \frac{1}{N} \sum_k v_i^k \sigma(c_j + \sum_{i'} v_{i'}^k W_{i'j}) \quad (17)$$

第二項も同様にして ,

$$\langle v_i h_j \rangle_{\text{model}} = \frac{1}{N} \sum_k \hat{v}_i^k \sigma(c_j + \sum_{i'} \hat{v}_{i'}^k W_{i'j}) \quad (18)$$

以上より , W_{ij} についての勾配は

$$\frac{\partial J}{\partial W_{ij}} = \frac{1}{N} \sum_k v_i^k \sigma(c_j + \sum_{i'} v_{i'}^k W_{i'j}) - \frac{1}{N} \sum_k \hat{v}_i^k \sigma(c_j + \sum_{i'} \hat{v}_{i'}^k W_{i'j}) \quad (19)$$

によって表現できる .

次に , b_i について . $\partial E / \partial b_i = -v_i$ より ,

$$\frac{\partial J}{\partial b_i} = \langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}} \quad (20)$$

ここで ,

$$\langle v_i \rangle_{\text{data}} = \sum_{\mathbf{v}} \sum_{\mathbf{h}} v_i q(\mathbf{v}) p(\mathbf{h}|\mathbf{v}) = \sum_{\mathbf{v}} v_i q(\mathbf{v}) = \frac{1}{N} \sum_k v_i^k \quad (21)$$

ゆえに, b_i についての勾配は下記の通りとなる.

$$\frac{\partial J}{\partial b_i} = \frac{1}{N} \sum_k v_i^k - \frac{1}{N} \sum_k \hat{v}_i^k \quad (22)$$

最後に, c_j について, $\partial E/\partial c_j = -h_j$ より,

$$\frac{\partial J}{\partial c_j} = \langle h_j \rangle_{\text{data}} - \langle h_j \rangle_{\text{model}} \quad (23)$$

ここで,

$$\langle h_j \rangle_{\text{data}} = \sum_{\mathbf{v}} \sum_{\mathbf{h}} h_j q(\mathbf{v}) p(\mathbf{h}|\mathbf{v}) = \sum_{\mathbf{h}} \left\{ \frac{1}{N} h_j \sum_k p(\mathbf{h}|\mathbf{v}^k) \right\} \quad (24)$$

$$= \frac{1}{N} \sum_k \sum_{h_j} h_j p(h_j|\mathbf{v}^k) = \frac{1}{N} \sum_k \sigma(c_j + \sum_{i'} v_{i'}^k W_{i'j}) \quad (25)$$

ゆえに, c_j についての勾配は下記の通りとなる.

$$\frac{\partial J}{\partial c_j} = \frac{1}{N} \sum_k \sigma(c_j + \sum_{i'} v_{i'}^k W_{i'j}) - \frac{1}{N} \sum_k \sigma(c_j + \sum_{i'} \hat{v}_{i'}^k W_{i'j}) \quad (26)$$

最後に, これらの計算結果をまとめることにしよう.

二値の Restricted Boltzmann Machine についての更新式

可視素子 v_i が二値 $\{0, 1\}$ しか取り得ない場合の, RBM のパラメータに関する更新式は下記の通りとなる.

$$\frac{\partial J}{\partial W_{ij}} = \frac{1}{N} \sum_k v_i^k \sigma(c_j + \sum_{i'} v_{i'}^k W_{i'j}) - \frac{1}{N} \sum_k \hat{v}_i^k \sigma(c_j + \sum_{i'} \hat{v}_{i'}^k W_{i'j}) \quad (27)$$

$$\frac{\partial J}{\partial b_i} = \frac{1}{N} \sum_k v_i^k - \frac{1}{N} \sum_k \hat{v}_i^k \quad (28)$$

$$\frac{\partial J}{\partial c_j} = \frac{1}{N} \sum_k \sigma(c_j + \sum_{i'} v_{i'}^k W_{i'j}) - \frac{1}{N} \sum_k \sigma(c_j + \sum_{i'} \hat{v}_{i'}^k W_{i'j}) \quad (29)$$

ここで, N はサンプル数, \hat{v}_i^k は可視素子のサンプル v_i^k を Gibbs Sampling によって A 回だけ遷移させた場合についてのサンプルを表している. 通常の場合, A は 1 回だけで良い.

5 可視素子が連続的な値を取る場合の RBM

前節では可視素子が二値という離散的な場合における RBM の更新式を導出したが, 可視素子が連続的な値を取る場合においては以前のエネルギー関数をそのまま用いてしまうと不都合が生じる. そのため, 可視素子が連続的な値を取っても問題のないよう修正した RBM も存在する. その場合のエネルギー関数は下記の通りとなる.

$$E(\mathbf{v}, \mathbf{h}) = \frac{1}{2} \sum_i v_i^2 - \sum_i b_i v_i - \sum_j c_j h_j - \sum_i \sum_j v_i W_{ij} h_j \quad (30)$$

この場合，条件付き確率 $p(\mathbf{v}|\mathbf{h}), p(\mathbf{h}|\mathbf{v})$ は下記の通りである．

$$p(v_i|\mathbf{h}) = N(v_i; b_i + \sum_j W_{ij}h_j, 1) \quad (31)$$

$$p(h_j = 1|\mathbf{v}) = \sigma(c_j + \sum_i v_i W_{ij}) \quad (32)$$

ここで， $N(x; \mu, \sigma^2)$ はガウス関数である． $p(\mathbf{v}|\mathbf{h})$ がシグモイド関数でなく，正規分布の確率密度関数として変化することに注意されたい*¹．

幸いにもこの場合におけるパラメータの更新式は，Contrastive Divergence を行う際に用いる条件付き確率が正規分布に変化したことを除いて，二値のRBMと全く同等である．ただし，仮定している生成モデルが異なるため，連続値の場合において推測されたパラメータは二値のそれとは本質的に異なる．

*¹ これは観測したデータの分布がガウシアンであることを強制するものではない．観測データが従う分布は周辺分布 $p(\mathbf{v}) = \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h})$ それ自身であり，条件付き確率ではない．