

多層制限ボルツマンマシンの実装

Implement For Deep RBM By C

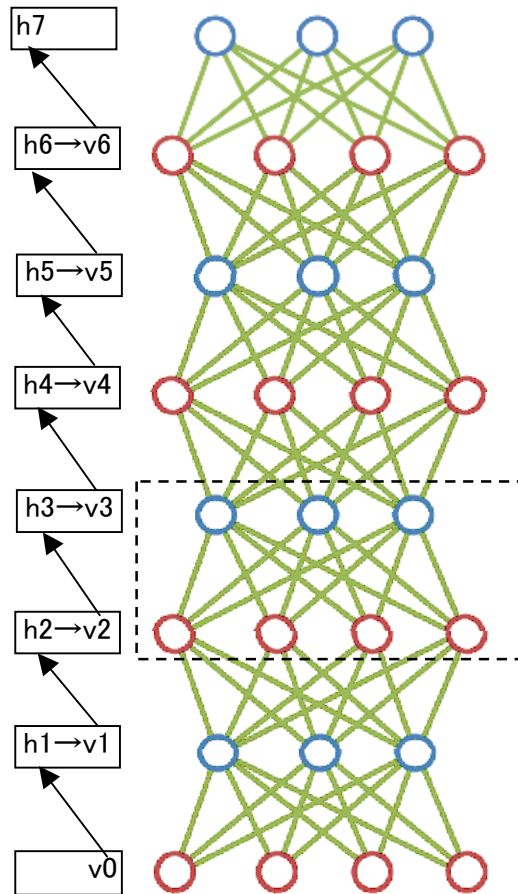
2015/8/5

by @mabonki0725

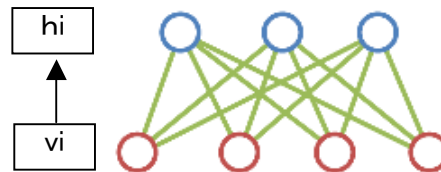
RBMの数理

① 多層制限ボルツマンマシン(Deep RBM)

多層のネットワークは制限ボルツマンモデルの繰返と考え、下層から1段毎に処理していく



制限ボルツマンマシンモデル

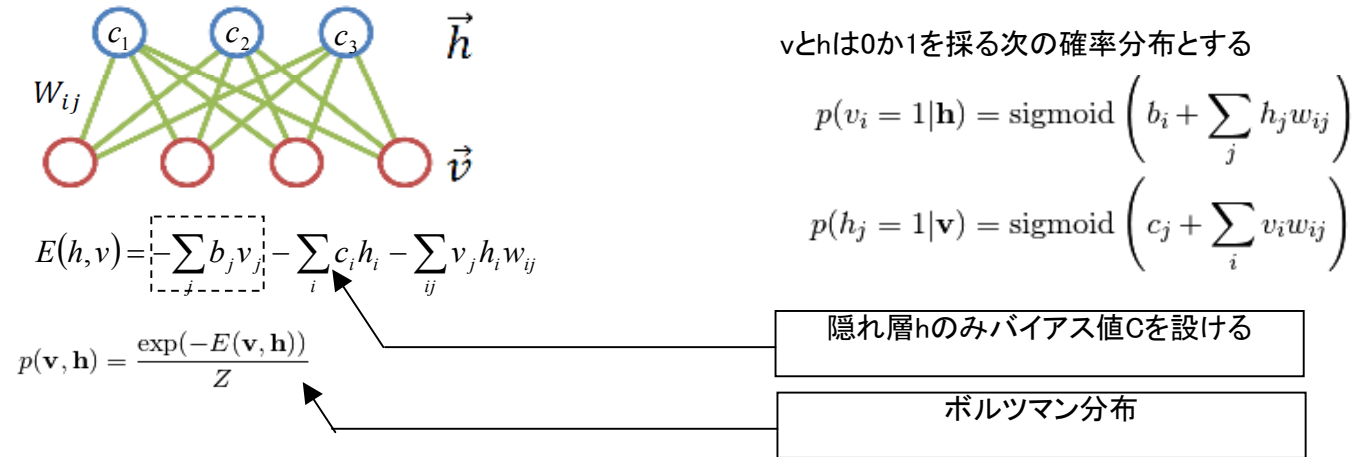


注記
制限なしのボルツマンマシンモデル
は補足資料2 参照

RBMの数理

② 制限ボルツマンマシンモデル

ネットワークが上下にのみ結線しているネットワークモデルを制限ボルツマンマシンモデル(RBM)と云う



\mathbf{v} の観測データがN件数ある場合、 $p(\mathbf{v}, \mathbf{h})$ の最適な連結間の重み \mathbf{W} と バイアス \mathbf{C} を求めるには
最大対数尤度 $\max L$ となる \mathbf{W} と \mathbf{C} を求めればよい

$$L(\mathcal{D}; \theta) \equiv \sum_{\mathbf{v}} q(\mathbf{v}) \log p(\mathbf{v})$$

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

ここで

$$\theta = [\mathbf{C} \ \mathbf{W}]$$

\mathcal{D} : \mathbf{v} のデータ

$q(\mathbf{v})$: 観測データの分布

③ 制限ボルツマンマシンの最尤推定

最大対数尤度を算出するには、対数尤度 L を $\theta=[W_{ij} \ C_i]$ で微分して確率的勾配法で θ を徐々に更新して求めることを考える
(上下層のみの制限ボルツマンマシンなので、対数尤度関数は凹であると考えられる)

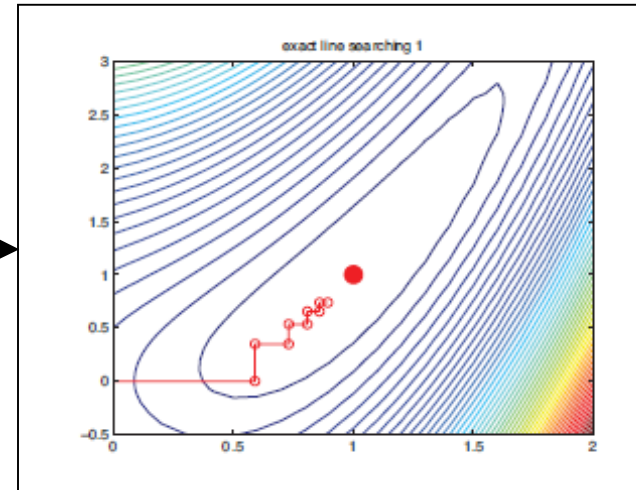
$$\theta' = \theta + \eta * \Delta L(D, \theta)$$

そこで θ で部分した対数尤度は以下で与えらる。

$$\frac{\partial \log p(v)}{\partial \theta} = \frac{\partial}{\partial \theta} \left\{ \log \sum_h \exp(-E(v, h)) \right\} - \frac{\partial}{\partial \theta} \log Z$$

$$\begin{aligned} \frac{\partial}{\partial \theta} \left\{ \log \sum_h \exp(-E(v, h)) \right\} &= \frac{\frac{\partial}{\partial \theta} \{ \sum_h \exp(-E(v, h)) \}}{\sum_h \exp(-E(v, h))} \\ &= \frac{\sum_h \exp(-E(v, h)) \frac{\partial}{\partial \theta} E(v, h)}{\sum_h \exp(-E(v, h))} = \frac{Z \sum_h p(v, h) \frac{\partial}{\partial \theta} E(v, h)}{Z p(v)} \\ &= - \sum_h p(h|v) \frac{\partial E(v, h)}{\partial \theta} \end{aligned}$$

$$\begin{aligned} - \frac{\partial}{\partial \theta} \log Z &= - \frac{1}{Z} \frac{\partial}{\partial \theta} \sum_{v, h} \exp(-E(v, h)) \\ &= - \sum_{v, h} p(v, h) \frac{\partial E(v, h)}{\partial \theta} \end{aligned}$$



確率的勾配法

$$\frac{\partial \mathcal{L}(D, \theta)}{\partial \theta} = - \left\langle \frac{\partial E(v, h)}{\partial \theta} \right\rangle_{data} + \left\langle \frac{\partial E(v, h)}{\partial \theta} \right\rangle_{model}$$

$$- \left\langle \frac{\partial E(v, h)}{\partial \theta} \right\rangle_{data} = - \sum_v q(v) \sum_h p(h|v) \frac{\partial E(v, h)}{\partial \theta} = - \frac{1}{N} \sum_j \sum_h p(h|v_j) \frac{\partial E(v_j, h)}{\partial \theta}$$

$$\left\langle \frac{\partial E(v, h)}{\partial \theta} \right\rangle_{model} = \sum_v q(v) \sum_{vh} p(v, h) \frac{\partial E(v, h)}{\partial \theta} = \frac{1}{N} \sum_j \sum_{jh} p(v_j, h) \frac{\partial E(v_j, h)}{\partial \theta}$$

$\theta = W_{ij}$ の場合

$$\begin{aligned} \frac{\partial \mathcal{L}(D, w_{ij})}{\partial w_{ij}} &= - \frac{1}{N} \sum_h p(h|v) \frac{\partial E(v, h)}{\partial w_{ij}} + \sum_{vh} p(v, h) \frac{\partial E(v, h)}{\partial w_{ij}} \\ &= \frac{1}{N} \left\{ \sum_h p(h|v_j) h_i v_j - \sum_{ih} p(v_j, h) h_i v_j \right\} \\ &= \frac{1}{N} \{ p(h=1|v_j) v_j - p(v_j, h=1) v_j \} \end{aligned}$$

$\theta = C_i$ の場合

$$\frac{\partial \mathcal{L}(D, c_i)}{\partial c_i} = \frac{1}{N} \{ p(h=1|v_j) - p(v_j, h=1) \}$$

④CD-k(k Cycle Contrastive Divergence)法

$$-\left\langle \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} \right\rangle_{\text{data}} = \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} / N$$

$$\left\langle \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} \right\rangle_{\text{model}} = \sum_{\mathbf{v}} \sum_{\mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} / N$$

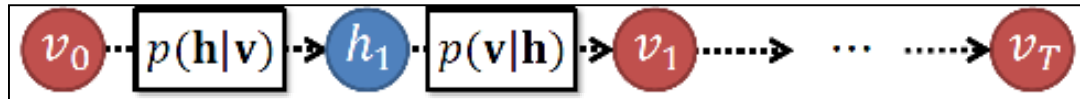
$$p(v_i = 1|\mathbf{h}) = \text{sigmoid} \left(b_i + \sum_j h_j w_{ij} \right)$$

$$p(h_j = 1|\mathbf{v}) = \text{sigmoid} \left(c_j + \sum_i v_i w_{ij} \right)$$

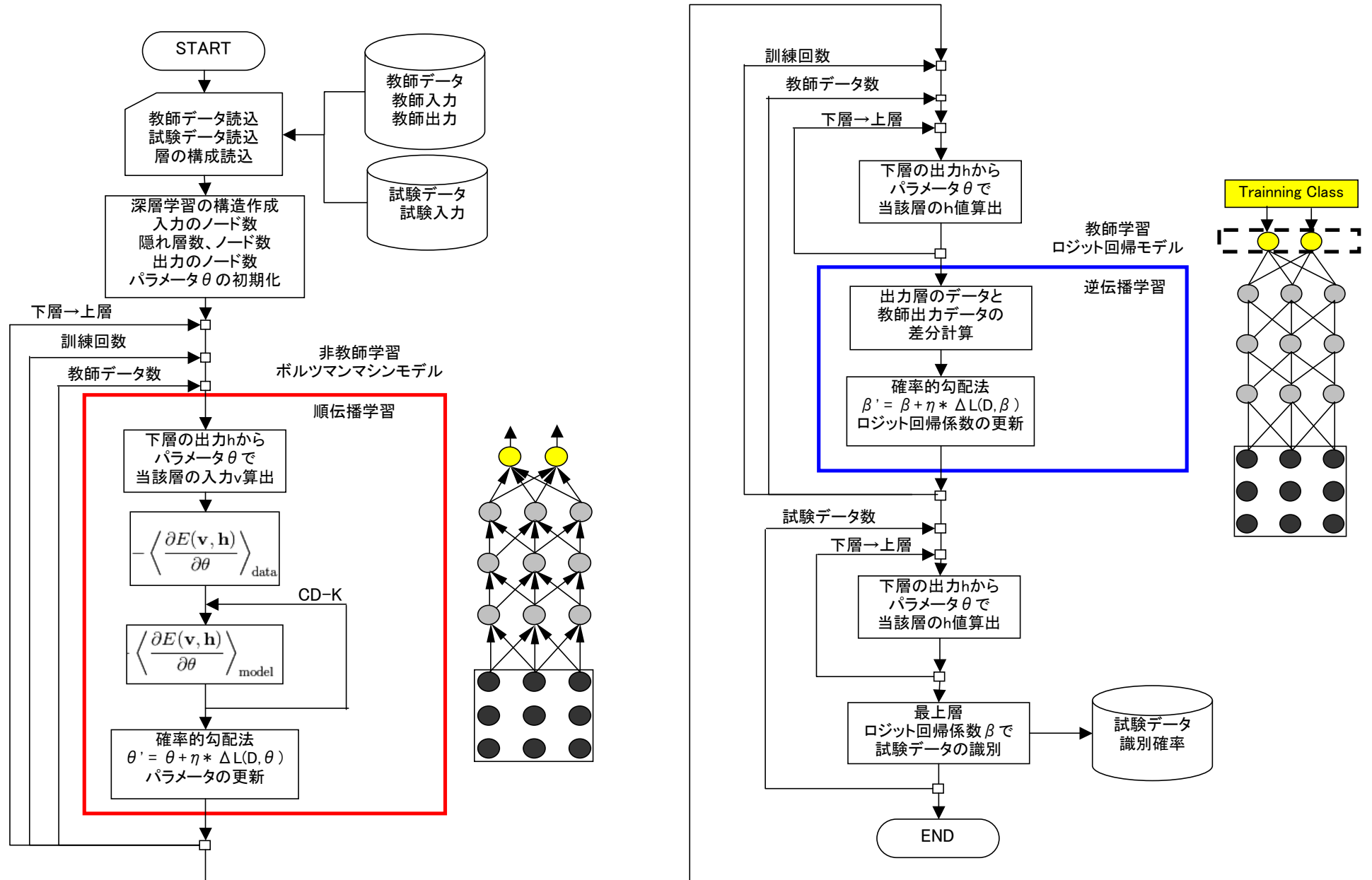
定義より算出できる

ここで $P(\mathbf{v}, \mathbf{h})$ は計算できないので、パラメータを交互に繰り返すギブスサンプラーを行う
これをCD-k(k Cycle Contrastive Divergence)法という。

CD-k



全体フロー図



全体処理図

main(argc,argv)

I 入力データを読み込む

- 1.1 ファイルからデータを読む
- 1.2 引数から訓練データ行数と試験データ行数を得る

II DBMによる訓練と試験データの識別を行う

教師用データ 教師用識別 試験用データ 識別結果 データ数 教師データ数 入力ノード数 出力ノード数 隠れノード数 隠れ層数

test_dbn(train_x, train_y, test_x, test_y, nInst, nTest, nInp, nOut, nHide, dHide);

2 層とノードの構成を作成

ネット 教師データ数 入力ノード数 隠れノード数 出力ノード数 隠れ層数

DBN_construct(&dbn, train_N, n_ins, hidden_layer_sizes, n_outs, n_layers);

ノード構成を設定する。

入力ノード数、隠れ層の数、隠れ層のノード数、出力ノード数
層間の連結重み $w[i][j]$ を初期値を一様分布で設定する

3 順伝播学習

ネット 教師データ 学習率 ギブス回数 トレーニング回数

DBN_pretrain(&dbn, *train_X, pretrain_lr, k, pretraining_epochs);

- loop 入力層から出力層まで 当該層 i
 - loop トレーニング回数
 - loop 教師データレコード数 N
 - 3.1 入力層から当該層の前層まで、層間の連結重みで、当該層への入力値 layer_input を計算する
 - 3.2 次の関数で当該層の重みを確率的勾配法(学習率 lr)で更新する

当該層 当該層への入力データ 学習率 ギブス回数

RBM_contrastive_divergence(&(this->rbm_layers[i]), layer_input, lr, k);

4 逆伝播学習

ネット 教師用データ 教師用識別 学習率 トレーニング回数

DBN_finetune(&dbn, *train_X, *train_Y, finetune_lr, finetune_epochs);

4.1 多クラスロジスティック回帰係数を確率勾配法で算出

- loop トレーニング回数
 - loop 訓練データレコード数 N
 - loop 入力層から出力層まで連結間の重みでノード値を計算
- 確率勾配法でロジット回帰係数を算出

5 試験データの予測

ネット 試験用データ 試験識別結果

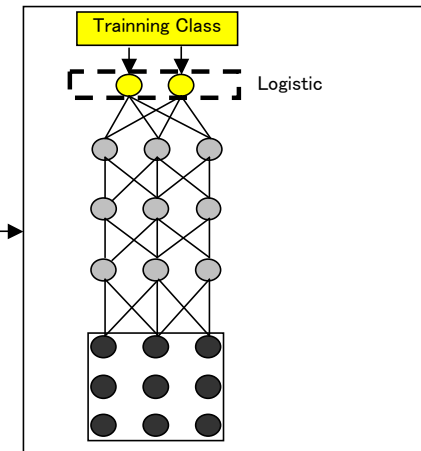
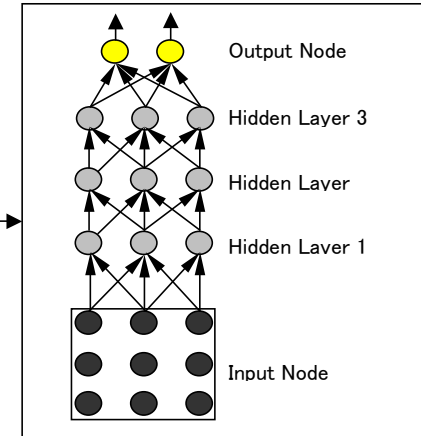
DBN_predict(&dbn, test_X[i], test_Y[i]);

- loop 試験データ数
 - loop 入力層から出力層まで
 - 5.1 入力層から最上層まで、連結間の連結重みで、出力層への入力値 layer_input を計算する
- 試験データの識別確率を計算する

ソフトマック関数で多クラス毎の確率を算出

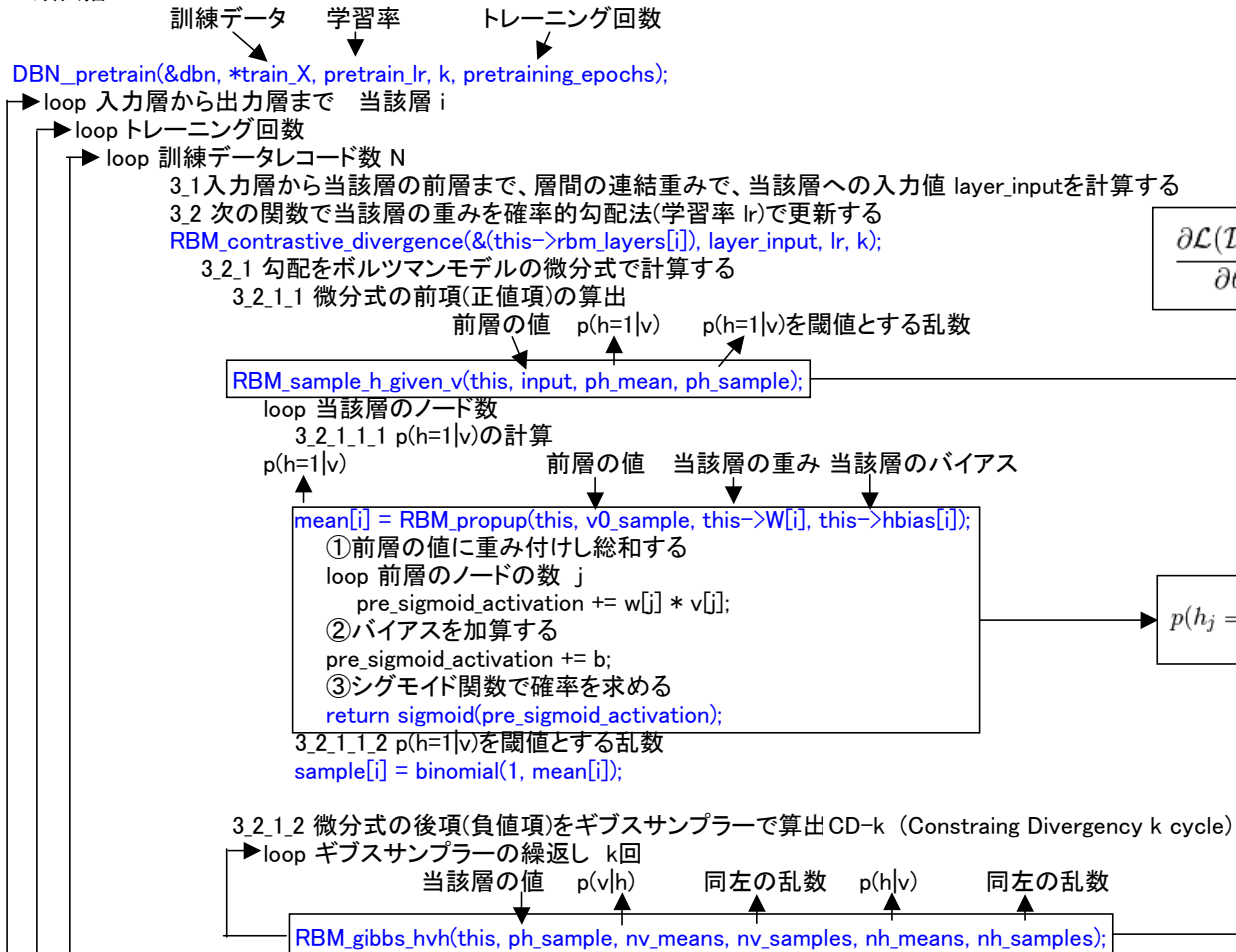
当該層 識別確率

LogisticRegression_softmax(&(this->log_layer), y);



3 順伝播学習

3 順伝播



$$\frac{\partial \mathcal{L}(\mathcal{D}, \theta)}{\partial \theta} = - \left\langle \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} \right\rangle_{\text{data}} + \left\langle \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} \right\rangle_{\text{model}}$$

$$p(h_j = 1|\mathbf{v}) = \text{sigmoid} \left(c_j + \sum_i v_i w_{ij} \right)$$

3 順伝播学習

3.2.1.2.1 当該層の値 $h \rightarrow$ 前層の値 v の逆算出

h の値 $P(v|h)$ $p(v|h)$ の乱数

RBM_sample_v_given_h(this, nh_sample, nv_means, nv_samples);

3.2.1.2.1.1 $p(h=1|v)$ の計算

$p(v=1|h)$

当該層の値

当該層のバイアス

```
mean[i]=RBM_propdown(this, h0_sample, i, this->vbias[i]);
①当該層の値に重み付けし総和する
loop 当該層のノードの数 j
    pre_sigmoid_activation += this->W[j][i] * h[j];
②バイアスを加算する
pre_sigmoid_activation += b;
③シグモイド関数で確率を求める
return sigmoid(pre_sigmoid_activation);
```

$$p(v_i = 1|h) = \text{sigmoid} \left(b_i + \sum_j h_j w_{ij} \right)$$

3.2.1.2.1.2 $p(h=1|v)$ を閾値とする乱数

sample[i] = binomial(1, mean[i]);

3.2.1.2.2 層の入力 $v \rightarrow$ 層の出力 h の順算出

v の値 $P(v|h)$ $p(v|h)$ の乱数

RBM_sample_h_given_v(this, nv_samples, nh_means, nh_samples);

3.2.1.2.2.1 $p(v=1|h)$ の計算

$p(h=1|v)$

前層の値

当該層の重み

当該層のバイアス

```
mean[i] = RBM_propup(this, v0_sample, this->W[i], this->hbias[i]);
①前層の値に重み付けし総和する
loop 前層のノードの数 j
    pre_sigmoid_activation += w[j] * v[j];
②バイアスを加算する
pre_sigmoid_activation += b;
③シグモイド関数で確率を求める
return sigmoid(pre_sigmoid_activation);
```

$$p(h_j = 1|v) = \text{sigmoid} \left(c_j + \sum_i v_i w_{ij} \right)$$

3.2.1.2.2.2 $p(v=1|h)$ を閾値とする乱数

sample[i] = binomial(1, mean[i]);

3.2.1.2 確率的勾配法での重みの更新

loop i 層のノードの数

loop j 前層のノード数

3.2.1.2.1 重みの更新

$p(h|v)$

v 値

$p(v,h)$

v 値(Gibbs)

this->W[j][i] += lr * (ph_mean[i] * input[j] - nh_means[i] * nv_samples[j]) / this->N;

3.2.1.2.2 バイアス項の更新

$p(h|v)$

$p(v,h)$

this->hbias[i] += lr * (ph_means[i] - nh_means[i]) / this->N;

$$\begin{aligned} \frac{\partial L(D, w_{ij})}{\partial w_{ij}} &= -\frac{1}{N} \sum_h p(h|v) \frac{\partial E(v, h)}{\partial w_{ij}} + \sum_{v, h} p(v, h) \frac{\partial E(v, h)}{\partial w_{ij}} \\ &= \frac{1}{N} \left\{ \sum_h p(h|v_j) h_i v_j - \sum_{i, h} p(v_j, h) h_i v_j \right\} \\ &= \frac{1}{N} \{ p(h=1|v_j) v_j - p(v_j, h=1) v_j \} \end{aligned}$$

$$\frac{\partial L(D, c)}{\partial c_i} = \frac{1}{N} \{ p(h=1|v_j) - p(v_j, h=1) \}$$

4 逆伝播学習

4 逆伝播学習

訓練データ 教師データ 学習率 学習繰返回数

```
DBN_finetune(&dbn, *train_X, *train_Y, finetune_lr, finetune_epochs);
```

4.1 多クラスロジステック回帰係数を確率勾配法で算出

→ loop トレーニング回数

→ loop 訓練データレコード数 N

→ loop 入力層から出力層まで

→ 4.1 入力層から当該層の前層まで、層間の連結重みで、当該層への入力値 layer_inputを計算する

4.2 次の関数で出力層のロジット回帰係数を確率的勾配法(学習率 lr)で更新する
前層の値 出力層の教師データ 学習率

```
LogisticRegression_train(&(this->log_layer), layer_input, train_Y, lr);
```

4.3 多クラスロジステック関数で確率計算

前層の値 指数化

```
LogisticRegression_softmax(this, p_y_given_x);
```

→ loop 出力ノード数

dy[i] = y[i] - p_y_given_x[i];

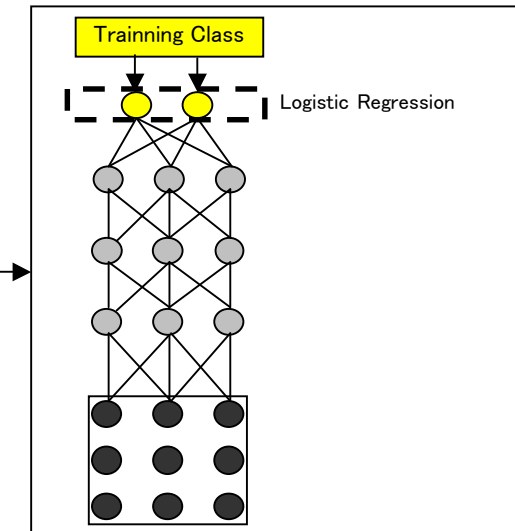
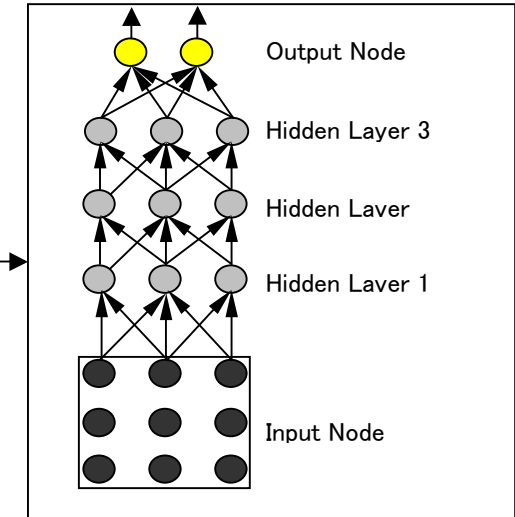
→ loop 出力層の前層のノード数

ロジット回帰係数の更新

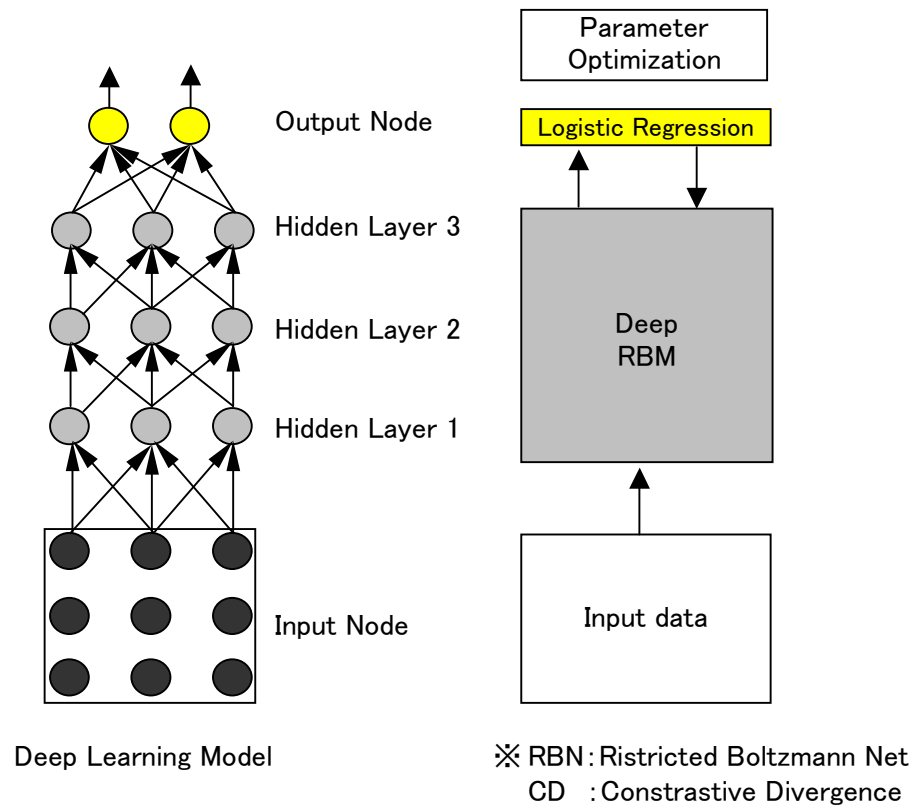
this->W[i][j] += lr * dy[i] * x[j] / this->N;

定数項の更新

this->b[i] += lr * dy[i] / this->N;



Deep Learning for recognition about close or open figure

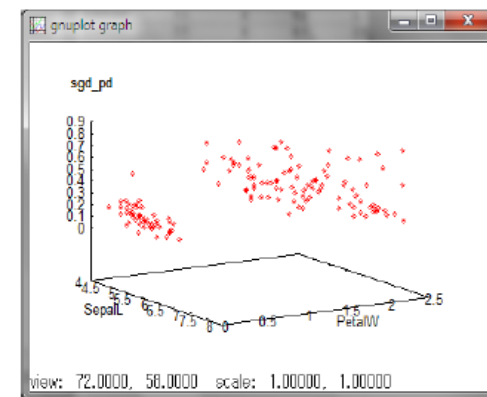
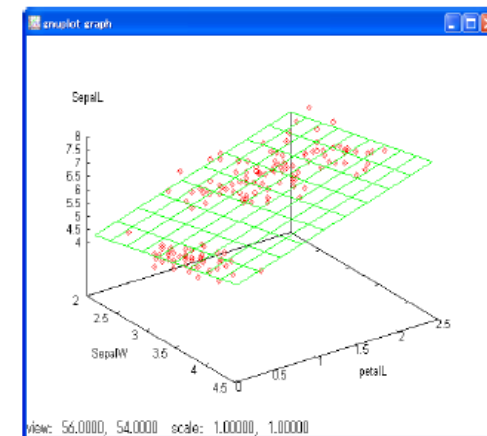
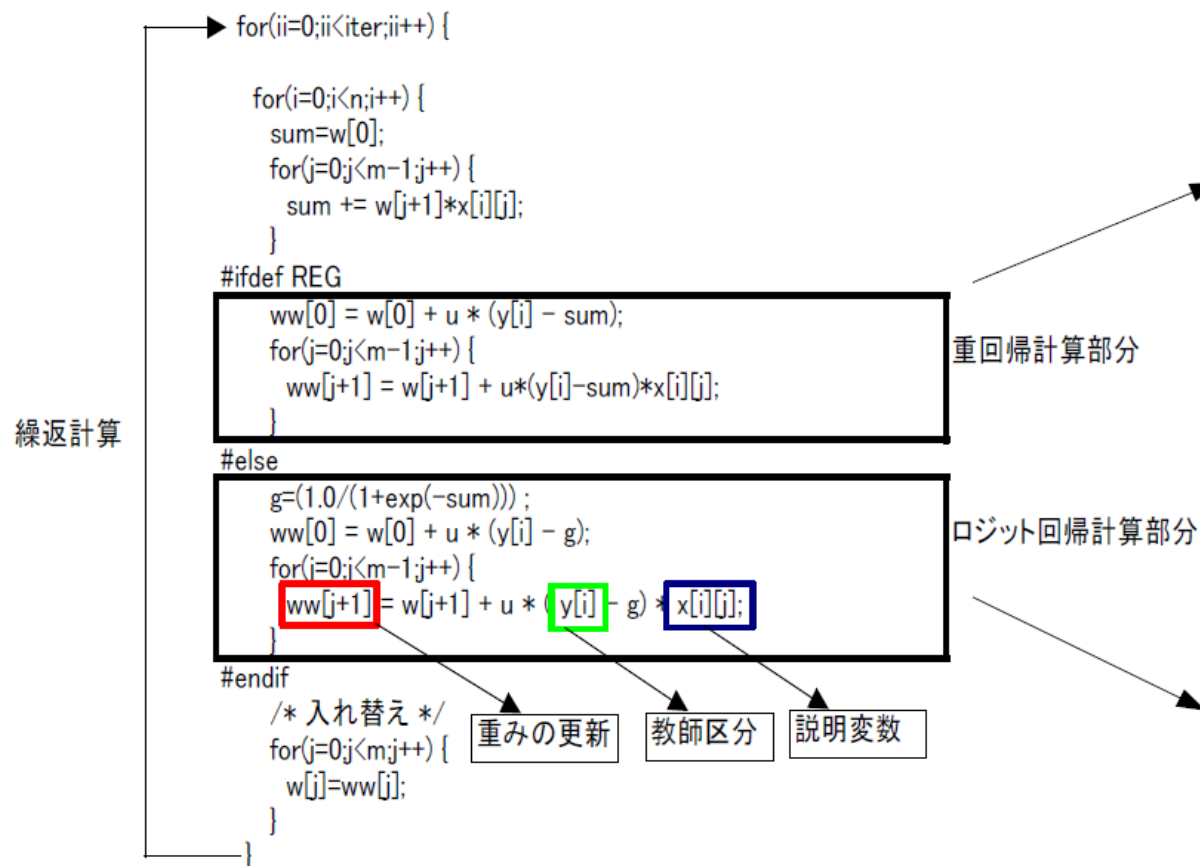


Test Analysis				
data	Input	pattern	close figure	open figure
Train	0 1 0		1	0
	1 0 1			
	0 1 0			
	0 1 0		1	0
	1 0 1			
	1 1 1			
	1 0 1		0	1
	0 1 0			
	1 0 1			
	1 1 1		1	0
	1 0 1			
	1 1 1			
	1 0 0		0	1
	1 0 0			
	1 1 1			
Prediction	0 1 1		1	0
	1 0 1			
	0 1 1			
	1 1 1		0	1
	0 0 1			
	0 0 1			
	1 1 0		0.846612	0.153388
	1 0 1		0.316327	0.683673
	1 1 0		0.237903	0.762097

補足資料1 ロジット回帰(SGD版)

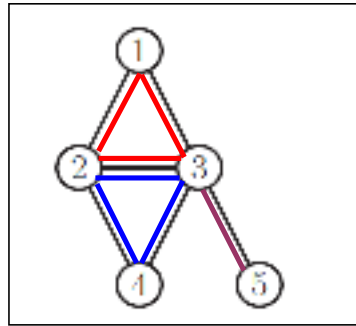
download <http://www1.m.jcnnet.jp/mabonki/download.htm>

一般にロジット回帰は多次元ニュートン法で解くが、SGD法だと5行でできる。



一般ボルツマンマシン

ノードとその連結(無方向)があり、ノード間の結合を**クリーク**単位で分割します。
クリーク内はボルツマン分布に従うと仮定した場合、全ノードの確率分布を計算します。



$$p(\mathbf{y}|\theta) = \frac{1}{Z(\theta)} \psi_{123}(y_1, y_2, y_3) \psi_{234}(y_2, y_3, y_4) \psi_{35}(y_3, y_5)$$

$$Z = \sum_{\mathbf{y}} \psi_{123}(y_1, y_2, y_3) \psi_{234}(y_2, y_3, y_4) \psi_{35}(y_3, y_5)$$

クリーク内は $\Phi(\mathbf{y})$ とパラメータ θ の線形結合 注) Φ は \mathbf{y} の特徴関数(カーネル等)

$$\log \psi_c(\mathbf{y}_c) \triangleq \phi_c(\mathbf{y}_c)^T \theta_c$$

全ノードのボルツマン確率分布

$$p(\mathbf{y}|\theta) = \frac{1}{Z(\theta)} \exp \left(\sum_c \theta_c^T \phi_c(\mathbf{y}) \right)$$

隠れ値 \mathbf{h} がある場合のボルツマン確率分布

$$p(\mathbf{y}, \mathbf{h}|\theta) = \frac{1}{Z(\theta)} \exp \left(\sum_c \theta_c^T \phi_c(\mathbf{h}, \mathbf{y}) \right)$$

一般ボルツマンマシンのパラメータ推定

各ノード値のデータが N 件ある場合、そのデータに合うパラメータを推定します。
推定には最大対数尤度を計算します。
最大値の算出には、対数尤度をパラメータで微分して、最大解を求めます。
ここでは確率的降下法を使ったアルゴリズムを示します。

$$\ell(\theta) \triangleq \frac{1}{N} \sum_i \log p(\mathbf{y}_i|\theta) = \frac{1}{N} \sum_i \left[\sum_c \theta_c^T \phi_c(\mathbf{y}_i) - \log Z(\theta) \right]$$

$$\frac{\partial \ell}{\partial \theta_c} = \frac{1}{N} \sum_i \left[\phi_c(\mathbf{y}_i) - \frac{\partial}{\partial \theta_c} \log Z(\theta) \right]$$

$$\frac{\partial \log Z(\theta)}{\partial \theta_c} = \mathbb{E}[\phi_c(\mathbf{y})|\theta] = \sum_{\mathbf{y}} \phi_c(\mathbf{y}) p(\mathbf{y}|\theta)$$

$$\frac{\partial \ell}{\partial \theta_c} = \left[\frac{1}{N} \sum_i \phi_c(\mathbf{y}_i) \right] - \mathbb{E}[\phi_c(\mathbf{y})]$$

隠れ値 \mathbf{h} がある場合の対数尤度の微分式

$$\frac{\partial \ell}{\partial \theta_c} = \frac{1}{N} \sum_i \{ \mathbb{E}[\phi_c(\mathbf{h}, \mathbf{y}_i)|\theta] - \mathbb{E}[\phi_c(\mathbf{h}, \mathbf{y})|\theta] \}$$

確率的降下法によるボルツマンマシンの学習

Algorithm 19.1: Stochastic maximum likelihood for fitting an MRF

```

1 Initialize weights  $\theta$  randomly;
2  $k = 0, \eta = 1$ ;
3 for each epoch do
4   for each minibatch of size  $B$  do
5     for each sample  $s = 1 : S$  do
6       Sample  $\mathbf{y}^{s,k} \sim p(\mathbf{y}|\theta_k)$ ;
7        $\hat{E}(\phi(\mathbf{y})) = \frac{1}{S} \sum_{s=1}^S \phi(\mathbf{y}^{s,k})$ ;
8       for each training case  $i$  in minibatch do
9          $\mathbf{g}_{ik} = \phi(\mathbf{y}_i) - \hat{E}(\phi(\mathbf{y}))$ ;
10       $\mathbf{g}_k = \frac{1}{B} \sum_{i \in B} \mathbf{g}_{ik}$ ;
11       $\theta_{k+1} = \theta_k - \eta \mathbf{g}_k$ ;
12       $k = k + 1$ ;
13   Decrease step size  $\eta$ ;
```