# Sparse Classification RBM

Weidong Liang

## Mathematical Formulation

$$E(y, \mathbf{X}, \mathbf{h}) = E(y, x_1, \dots, x_C, \mathbf{h}) = -\sum_{1 \le i \le C} \left( \mathbf{h}^\mathbf{T} \mathbf{W}^{(i)} \mathbf{e}_{x_i} + \mathbf{b}^{(i)\mathbf{T}} \mathbf{e}_{x_i} \right) - \mathbf{h}^\mathbf{T} \mathbf{c} - dy - \mathbf{h}^\mathbf{T} \mathbf{U} y$$

$$p(y, \mathbf{X}, \mathbf{h}) = p(y, x_1, \dots, x_C, \mathbf{h}) = \frac{\exp(-E(y, x_1, \dots, x_C, \mathbf{h}))}{Z}$$

$$Z = \sum_y \sum_{\mathbf{X}} \sum_{\mathbf{h}} p(y, \mathbf{X}, \mathbf{h})$$

$$\mathbf{e_y} = (1_{i=y})_{1 \le i \le C}$$

$$y \in \{0,1\}$$

$$x_i \in \{1, \dots, k_i\}, 1 \le 1 \le C$$

$$h_i \in \{0,1\}, 1 \le i \le H$$

Parameter Set: $\theta = \{\mathbf{W}, \mathbf{b}, d, \mathbf{c}, \mathbf{U}\}$

## Model Explanation

Y: click (y=1) or no click (y=0).
$x_i$: feature value of the i-th feature class (there are total of C feature classes).
$h_i$: i-th hidden unit value.

## Derivation of Properties

$$p(\mathbf{h}|y, \mathbf{X}) = \prod_j p(h_j|y, \mathbf{X})$$

$$
\begin{aligned}
p(\mathbf{h}|y, \mathbf{X}) &= \frac{p(y, \mathbf{X}, \mathbf{h})}{p(y, \mathbf{X})} = \frac{p(y, \mathbf{X}, \mathbf{h})}{\sum_{\mathbf{h}'} p(y, \mathbf{X}, \mathbf{h}')} = \frac{\exp[-E(y, \mathbf{X}, \mathbf{h})]}{\sum_{\mathbf{h}'} \exp[-E(y, \mathbf{X}, \mathbf{h}')]} \\[2mm]
&= \frac{\exp\{\sum_{1 \le i \le C}(\mathbf{h}^\mathbf{T} \mathbf{W}^{(i)} \mathbf{e}_{x_i} + \mathbf{b}^{(i)\mathbf{T}} \mathbf{e}_{x_i}) + \mathbf{h}^\mathbf{T}\mathbf{c} + dy + \mathbf{h}^\mathbf{T}\mathbf{U}y\}}{\sum_{\mathbf{h}'} \exp\{\sum_{1 \le i \le C}(\mathbf{h'}^\mathbf{T} \mathbf{W}^{(i)} \mathbf{e}_{x_i} + \mathbf{b}^{(i)\mathbf{T}} \mathbf{e}_{x_i}) + \mathbf{h'}^\mathbf{T}\mathbf{c} + dy + \mathbf{h'}^\mathbf{T}\mathbf{U}y\}} \\[2mm]
&= \frac{\exp(dy)}{\exp(dy)} \frac{\exp(\sum_{1 \le i \le C} \mathbf{b}^{(i)\mathbf{T}} \mathbf{e}_{x_i})}{\exp(\sum_{1 \le i \le C} \mathbf{b}^{(i)\mathbf{T}} \mathbf{e}_{x_i})} \frac{\exp\{\mathbf{h}^T[\sum_{1 \le i \le C} \mathbf{W}^{(i)} \mathbf{e}_{x_i} + \mathbf{c} + \mathbf{U}y]\}}{\sum_{\mathbf{h}'} \exp\{\mathbf{h'}^T[\sum_{1 \le i \le C} \mathbf{W}^{(i)} \mathbf{e}_{x_i} + \mathbf{c} + \mathbf{U}y]\}} \\[2mm]
&= \frac{\exp\left\{\sum_{1 \le j \le H}\left[h_j(\sum_{1 \le i \le C} W_{j,x_i}^{(i)} + c_i + U_i y)\right]\right\}}{\sum_{\mathbf{h}'} \exp\left\{\sum_{1 \le j \le H}\left[h'_j(\sum_{1 \le i \le C} W_{j,x_i}^{(i)} + c_i + U_i y)\right]\right\}}
\end{aligned}
$$

$$= \frac{\prod_{1 \leq j \leq H} \exp[h_j(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_i + U_i y)]}{\sum_{\mathbf{h}'} \prod_{1 \leq j \leq H} \exp[h'_j(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_i + U_i y)]}$$

$$= \frac{\prod_{1 \leq j \leq H} \exp[h_j(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_i + U_i y)]}{\sum_{h'_1} \dots \sum_{h'_H} \prod_{1 \leq j \leq H} \exp[h'_j(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_i + U_i y)]}$$

$$= \frac{\prod_{1 \leq j \leq H} \exp[h_j(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_i + U_i y)]}{\left\{\sum_{h'_1 \in \{0,1\}} \exp[h'_1(\sum_{1 \leq i \leq C} W^{(i)}_{1,x_i} + c_1 + U_1 y)]\right\} \dots \left\{\sum_{h'_H \in \{0,1\}} \exp[h'_H(\sum_{1 \leq i \leq C} W^{(i)}_{H,x_i} + c_H + U_H y)]\right\}}$$

$$= \frac{\prod_{1 \leq j \leq H} \exp[h_j(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_i + U_i y)]}{\left\{1 + \exp[(\sum_{1 \leq i \leq C} W^{(i)}_{1,x_i} + c_1 + U_1 y)]\right\} \dots \left\{1 + \exp[\sum_{1 \leq i \leq C} W^{(i)}_{H,x_i} + c_H + U_H y]\right\}}$$

$$= \prod_{1 \leq j \leq H} \frac{\exp[h_j(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_i + U_i y)]}{1 + \exp[(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_j + U_j y)]}$$

$$= \prod_{1 \leq j \leq H} p(h_j | y, \mathbf{X})$$

$$p(h_j = 1 | y, \mathbf{X}) = \sigma(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_j + U_j y)$$

$$p(h_j = 1 | y, \mathbf{X}) = \sum_{\mathbf{h}' \in \{\dots, h_j = 1, \dots\}} p(\mathbf{h} | y, \mathbf{X})$$

$$= \sum_{\mathbf{h}' \in \{\dots, h_j = 1, \dots\}} \prod_{1 \leq k \leq H} \frac{\exp[h_j(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_i + U_i y)]}{1 + \exp[(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_j + U_j y)]}$$

$$= \frac{\exp[h_j(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_i + U_i y)]}{1 + \exp[(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_j + U_j y)]} \prod_{1 \leq k \leq H, k \neq j} \sum_{h_k \in \{0,1\}} \frac{\exp[h_k(\sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + U_k y)]}{1 + \exp[(\sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + U_k y)]}$$

$$= \frac{\exp[\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_i + U_i y]}{1 + \exp[(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_j + U_j y)]}$$

$$= \sigma(\sum_{1 \leq i \leq C} W^{(i)}_{j,x_i} + c_j + U_j y)$$

$$p(y | \mathbf{X}) = \frac{\exp\{-F(y, X)\}}{\sum_{y''} \exp\{-F(y'', X)\}}$$

Note that:

$$\sum_{\mathbf{h}'} p(y, \mathbf{X}, \mathbf{h}') = \frac{1}{Z} \sum_{\mathbf{h}'} \exp(-E(y, \mathbf{X}, \mathbf{h}')$$

$$= \frac{1}{Z} \sum_{\mathbf{h}'} \exp\left\{\sum_{1 \leq i \leq C} (\mathbf{h}'^{\mathbf{T}} \mathbf{W}^{(i)} \mathbf{e}_{x_i} + \mathbf{b}^{(i)\mathbf{T}} \mathbf{e}_{x_i}) + \mathbf{h}'^{\mathbf{T}} \mathbf{c} + \mathrm{dy} + \mathbf{h}'^{\mathbf{T}} \mathbf{U} y\right\}$$

$$= \frac{1}{Z} \sum_{\mathbf{h}'} \{\exp(dy)\} \{\exp(\sum_{1 \leq i \leq C} \mathbf{b}^{(i)\mathbf{T}} \mathbf{e}_{x_i})\} \left\{ \exp\left[ \mathbf{h}'^{\mathbf{T}} \left( \sum_{1 \leq i \leq C} \mathbf{W}^{(i)} \mathbf{e}_{x_i} + \mathbf{c} + \mathbf{U} y \right) \right] \right\}$$

$$= \frac{1}{Z} \{\exp(\sum_{1 \leq i \leq C} \mathbf{b}^{(i)\mathbf{T}} \mathbf{e}_{x_i})\} \{\exp(dy)\} \sum_{\mathbf{h}'} \left\{ \exp\left[ \mathbf{h}'^{\mathbf{T}} \left( \sum_{1 \leq i \leq C} \mathbf{W}^{(i)} \mathbf{e}_{x_i} + \mathbf{c} + \mathbf{U} y \right) \right] \right\}$$

$$= \frac{1}{Z} \{\exp(\sum_{1 \leq i \leq C} \mathbf{b}^{(i)\mathbf{T}} \mathbf{e}_{x_i})\} \{\exp(dy)\} \prod_{1 \leq k \leq H} \sum_{h_k \in \{0,1\}} \left\{ \exp\left[ h_k \left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k y \right) \right] \right\}$$

$$= \frac{1}{Z} \{\exp(\sum_{1 \leq i \leq C} \mathbf{b}^{(i)\mathbf{T}} \mathbf{e}_{x_i})\} \{\exp(dy)\} \prod_{1 \leq k \leq H} \left\{ 1 + \exp\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k y \right) \right\}$$

Hence:

$$p(y|\mathbf{X}) = \frac{p(y,\mathbf{X})}{p(\mathbf{X})} = \frac{\sum_{\mathbf{h}'} p(y,\mathbf{X},\mathbf{h}')}{\sum_{y''} \sum_{\mathbf{h}''} p(y'',\mathbf{X},\mathbf{h}'')}$$

$$= \frac{\{\exp(dy)\} \prod_{1 \leq k \leq H} \left\{ 1 + \exp\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k y \right) \right\}}{\sum_{y''} \{\exp(dy'')\} \prod_{1 \leq k \leq H} \left\{ 1 + \exp\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k y'' \right) \right\}}$$

$$= \frac{\exp\left\{ dy + \ln\left[ \prod_{1 \leq k \leq H} \left\{ 1 + \exp\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k y \right) \right\} \right] \right\}}{\sum_{y''} \exp\left\{ dy'' + \ln\left[ \prod_{1 \leq k \leq H} \left\{ 1 + \exp\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k y'' \right) \right\} \right] \right\}}$$

$$= \frac{\exp\left\{ dy + \sum_{1 \leq k \leq H} \ln\left[ 1 + \exp\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k y \right) \right] \right\}}{\sum_{y''} \exp\left\{ dy'' + \sum_{1 \leq k \leq H} \ln\left[ 1 + \exp\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k y'' \right) \right] \right\}}$$

$$= \frac{\exp\left\{ dy + \sum_{1 \leq k \leq H} \text{softplus}\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k y \right) \right\}}{\sum_{y''} \exp\left\{ dy'' + \sum_{1 \leq k \leq H} \text{softplus}\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k y'' \right) \right\}}$$

$$= \frac{\exp\left\{ dy + \sum_{1 \leq k \leq H} \text{softplus}\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k y \right) \right\}}{\exp\left\{ \sum_{1 \leq k \leq H} \text{softplus}\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k \right) \right\} + \exp\left\{ d + \sum_{1 \leq k \leq H} \text{softplus}\left( \sum_{1 \leq i \leq C} W^{(i)}_{k,x_i} + c_k + \mathbf{U}_k \right) \right\}}$$

$$= \frac{\exp\{-F(y,X)\}}{\sum_{y''} \exp\{-F(y'',X)\}}$$

$$p(y|\mathbf{h}) = \frac{\exp\{dy + \mathbf{h}^T \mathbf{U} y\}}{1 + \exp\{d + \mathbf{h}^T \mathbf{U}\}}$$

Note that

$$\sum_{\mathbf{X}} p(y,\mathbf{X},\mathbf{h}) = \frac{1}{Z} \sum_{\mathbf{X}} \exp\left\{ \sum_{1 \leq i \leq C} (\mathbf{h}^{\mathbf{T}} \mathbf{W}^{(i)} + \mathbf{b}^{(i)\mathbf{T}}) \mathbf{e}_{x_i} \right\} \exp\{\mathbf{h}^T \mathbf{c}\} \exp\{dy + \mathbf{h}^T \mathbf{U} y\}$$

$$= \frac{1}{Z} \exp\{\mathbf{h}^T \mathbf{c}\} \exp\{dy + \mathbf{h}^T \mathbf{U} y\} \sum_{\mathbf{X}} \exp\left\{ \sum_{1 \leq i \leq C} (\mathbf{h}^{\mathbf{T}} \mathbf{W}^{(i)} + \mathbf{b}^{(i)\mathbf{T}}) \mathbf{e}_{x_i} \right\}$$

$$= \frac{1}{Z} \exp\{\mathbf{h}^T\mathbf{c}\} \exp\{dy + \mathbf{h}^T\mathbf{U}y\} \prod_{1 \le i \le C} \sum_{x_i} \exp\{(\mathbf{h}^T\mathbf{W}^{(i)} + \mathbf{b}^{(i)T})\mathbf{e}_{x_i}\}$$

$$= \frac{1}{Z} \exp\{\mathbf{h}^T\mathbf{c}\} \exp\{dy + \mathbf{h}^T\mathbf{U}y\} K(\mathbf{h})$$

Hence:

$$p(y|\mathbf{h}) = \frac{p(y,\mathbf{h})}{p(\mathbf{h})} = \frac{\sum_{\mathbf{X}'} p(y,\mathbf{X}',\mathbf{h})}{\sum_{y''}\sum_{\mathbf{X}''} p(y'',\mathbf{X}'',\mathbf{h})}$$

$$= \frac{\exp\{\mathbf{h}^T\mathbf{c}\}\exp\{dy + \mathbf{h}^T\mathbf{U}y\}K(\mathbf{h})}{\sum_{y''}\exp\{\mathbf{h}^T\mathbf{c}\}\exp\{dy'' + \mathbf{h}^T\mathbf{U}y''\}K(\mathbf{h})} = \frac{\exp\{dy + \mathbf{h}^T\mathbf{U}y\}}{\sum_{y''}\exp\{dy'' + \mathbf{h}^T\mathbf{U}y''\}} = \frac{\exp\{dy + \mathbf{h}^T\mathbf{U}y\}}{1 + \exp\{d + \mathbf{h}^T\mathbf{U}\}}$$

$$p(y = 1|\mathbf{h}) = \sigma(d + \mathbf{h}^T\mathbf{U})$$

$$\mathbf{p}(\mathbf{X}|\mathbf{h}) = \prod_{1 \le i \le C} p(x_i|\mathbf{h})$$

Note that:

$$\sum_{y'} p(y',\mathbf{X},\mathbf{h}) = \frac{1}{Z} \sum_{y'} \exp\left\{ \sum_{1 \le i \le C} \left(\mathbf{h}^T\mathbf{W}^{(i)}\mathbf{e}_{x_i} + \mathbf{b}^{(i)T}\mathbf{e}_{x_i}\right) + \mathbf{h}^T\mathbf{c} + dy' + \mathbf{h}^T\mathbf{U}y' \right\}$$

$$= \frac{1}{Z} \sum_{y'} \exp\left\{ \sum_{1 \le i \le C} \left(\mathbf{h}^T\mathbf{W}^{(i)} + \mathbf{b}^{(i)T}\right)\mathbf{e}_{x_i} \right\} \exp\{\mathbf{h}^T\mathbf{c}\} \exp\{dy' + \mathbf{h}^T\mathbf{U}y'\}$$

$$= \frac{1}{Z} \left[ \exp\left\{ \sum_{1 \le i \le C} \left(\mathbf{h}^T\mathbf{W}^{(i)} + \mathbf{b}^{(i)T}\right)\mathbf{e}_{x_i} \right\} \exp\{\mathbf{h}^T\mathbf{c}\} \right] \sum_{y'} \exp\{dy' + \mathbf{h}^T\mathbf{U}y'\}$$

$$= \frac{1}{Z} [1 + \exp(d + \mathbf{h}^T\mathbf{U})] \exp\{\mathbf{h}^T\mathbf{c}\} \exp\left\{ \sum_{1 \le i \le C} \left(\mathbf{h}^T\mathbf{W}^{(i)} + \mathbf{b}^{(i)T}\right)\mathbf{e}_{x_i} \right\}$$

Hence:

$$p(\mathbf{X}|\mathbf{h}) = \frac{p(\mathbf{X},\mathbf{h})}{p(\mathbf{h})} = \frac{\sum_{y'} p(y',\mathbf{X},\mathbf{h})}{\sum_{\mathbf{X}''}\sum_{y''} p(y'',\mathbf{X}'',\mathbf{h})}$$

$$= \frac{\exp\left\{\sum_{1 \le i \le C}\left(\mathbf{h}^T\mathbf{W}^{(i)} + \mathbf{b}^{(i)T}\right)\mathbf{e}_{x_i}\right\}}{\sum_{\mathbf{X}''}\exp\left\{\sum_{1 \le i \le C}\left(\mathbf{h}^T\mathbf{W}^{(i)} + \mathbf{b}^{(i)T}\right)\mathbf{e}_{x_i''}\right\}}$$

$$= \frac{\prod_{1 \le i \le C}\exp\{\left(\mathbf{h}^T\mathbf{W}^{(i)} + \mathbf{b}^{(i)T}\right)\mathbf{e}_{x_i}\}}{\prod_{1 \le i \le C}\sum_{x_i'' \in \{1,\dots,C_i\}}\exp\{\left(\mathbf{h}^T\mathbf{W}^{(i)} + \mathbf{b}^{(i)T}\right)\mathbf{e}_{x_i}\}}$$

$$= \prod_{1 \le i \le C} p(x_i|h)$$

$$p(x_i = k | \mathbf{h}) = \frac{\exp(\sum_{1 \le j \le H} h_j W_{j,k}^{(i)})}{\sum_{1 \le q \le C_i} \exp(\sum_{1 \le j \le H} h_j W_{j,q}^{(i)})}$$

$$p(x_i = k | \mathbf{h}) = \frac{\exp\{(\mathbf{h^T W^{(i)} + b^{(i)T}})\mathbf{e_{x_i}}\}}{\sum_{x_i'' \in \{1,\dots,C_i\}} \exp\{(\mathbf{h^T W^{(i)} + b^{(i)T}})\mathbf{e_{x_i}}\}}$$

$$= \frac{\exp(\sum_{1 \le j \le H} h_j W_{j,k}^{(i)})}{\sum_{1 \le q \le C_i} \exp(\sum_{1 \le j \le H} h_j W_{j,q}^{(i)})}$$

# Derivation of Learning Algorithm

Generative Learning

$$\frac{\partial \ln p(y, \mathbf{X})}{\partial \theta} = -\mathbb{E}_{p(\mathbf{h}|y,\mathbf{X})}\left[\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial \theta}\right] + \mathbb{E}_{p(y,\mathbf{X},\mathbf{h})}\left[\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial \theta}\right]$$

$$\frac{\partial \ln p(y, \mathbf{X})}{\partial \theta} = \frac{\partial}{\partial \theta}\left\{\ln \sum_{\mathbf{h'}} \exp(-E(y, \mathbf{X}, \mathbf{h'}) - \ln Z\right\}$$

$$= -\sum_{\mathbf{h'}}\left[\frac{exp(-E(y, \mathbf{X}, \mathbf{h'})}{\sum_{\mathbf{h''}} exp(-E(y, \mathbf{X}, \mathbf{h''})} \frac{\partial E(y, \mathbf{X}, \mathbf{h'})}{\partial \theta}\right] + \sum_{y',\mathbf{X},\mathbf{h'}}\left[\frac{\exp(-E(y', \mathbf{X'}, \mathbf{h'})}{Z} \frac{\partial E(y', \mathbf{X'}, \mathbf{h'})}{\partial \theta}\right]$$

$$= -\sum_{\mathbf{h'}} p(\mathbf{h'}|y, \mathbf{X}) \frac{\partial E(y, \mathbf{X}, \mathbf{h'})}{\partial \theta} + \sum_{y',\mathbf{X'},\mathbf{h'}} p(y', \mathbf{X'}, \mathbf{h'}) \frac{\partial E(y', \mathbf{X'}, \mathbf{h'})}{\partial \theta}$$

Discriminative Learning

$$\frac{\partial \ln p(y|\mathbf{X})}{\partial \theta} = -\mathbb{E}_{p(\mathbf{h}|y,\mathbf{X})}\left[\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial \theta}\right] + \mathbb{E}_{p(y,\mathbf{h}|\mathbf{X})}\left[\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial \theta}\right]$$

$$\frac{\partial \ln p(y|\mathbf{X})}{\partial \theta} = \frac{\partial}{\partial \theta}\ln\left[\frac{\frac{1}{Z}\sum_{\mathbf{h'}} p(y, \mathbf{X}, \mathbf{h'})}{\frac{1}{Z}\sum_{y'}\sum_{\mathbf{h'}} p(y', \mathbf{X}, \mathbf{h'})}\right]$$

$$= \frac{\partial}{\partial \theta}\left\{\ln\left[\sum_{\mathbf{h'}} p(y, \mathbf{X}, \mathbf{h'})\right] - \ln\left[\sum_{y'}\sum_{\mathbf{h'}} p(y', \mathbf{X}, \mathbf{h'})\right]\right\}$$

$$= -\sum_{\mathbf{h'}}\left\{\frac{p(y, \mathbf{X}, \mathbf{h'})}{\sum_{\mathbf{h''}} p(y, \mathbf{X}, \mathbf{h''})} \frac{\partial E(y, \mathbf{X}, \mathbf{h'})}{\partial \theta}\right\} + \sum_{y'}\sum_{\mathbf{h'}}\left\{\frac{p(y', \mathbf{X}, \mathbf{h'})}{\sum_{y''}\sum_{\mathbf{h''}} p(y'', \mathbf{X}, \mathbf{h''})} \frac{\partial E(y', \mathbf{X}, \mathbf{h'})}{\partial \theta}\right\}$$

$$= -\sum_{\mathbf{h'}} p(\mathbf{h'}|y, \mathbf{X}) \frac{\partial E(y, \mathbf{X}, \mathbf{h'})}{\partial \theta} + \sum_{y'}\sum_{\mathbf{h'}} p(y', \mathbf{h'}|\mathbf{X}) \frac{\partial E(y', \mathbf{X}, \mathbf{h'})}{\partial \theta}$$

Hybrid Learning

$$\frac{\partial \ln Hybrid(\alpha, y, \mathbf{X})}{\partial \theta} = \frac{\partial \ln p(y|\mathbf{X})}{\partial \theta} + \alpha \frac{\partial \ln p(y, \mathbf{X})}{\partial \theta}$$

$$= -(1 + \alpha)\mathbb{E}_{p(\mathbf{h}|y,\mathbf{X})}\left[\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial \theta}\right] + \mathbb{E}_{p(y,\mathbf{h}|\mathbf{X})}\left[\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial \theta}\right] + \alpha\mathbb{E}_{p(y,\mathbf{X},\mathbf{h})}\left[\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial \theta}\right]$$

$$= -(1 + \alpha)A(y, \mathbf{X}, \mathbf{h}, \theta) + B(y, \mathbf{X}, \mathbf{h}, \theta) + \alpha C(y, \mathbf{X}, \mathbf{h}, \theta)$$

Derivatives of Gradients

$$E(y, \mathbf{X}, \mathbf{h}) = E(y, x_1, \dots, x_C, \mathbf{h}) = - \sum_{1 \leq i \leq C} \left(\mathbf{h}^\mathsf{T} \mathbf{W}^{(i)} \mathbf{e}_{x_i} + \mathbf{b}^{(i)\mathsf{T}} \mathbf{e}_{x_i}\right) - \mathbf{h}^\mathsf{T} \mathbf{c} - dy - \mathbf{h}^\mathsf{T} \mathbf{U} y$$

Parameter Set: $\theta = \{\mathbf{W}, \mathbf{b}, d, \mathbf{c}, \mathbf{U}\}$

$$\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial W_{j,k}^{(i)}} = -h_j 1_{(x_i=k)}$$

$$\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial b_j^i} = -1_{(x_i=j)}$$

$$\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial d} = -y$$

$$\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial c_i} = -h_i$$

$$\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial U_i} = -h_i y$$

Parameter Updates

$$A(y, \mathbf{X}, \mathbf{h}, \boldsymbol{\theta})$$

$$A\left(y, \mathbf{X}, \mathbf{h}, W_{j,k}^{(i)}\right) = -p(h_j = 1|y, X)1_{(x_i=k)}$$

$$A\left(y, \mathbf{X}, \mathbf{h}, b_k^i\right) = -1_{(x_i=k)}$$
$$A(y, \mathbf{X}, \mathbf{h}, d) = -y$$
$$A\left(y, \mathbf{X}, \mathbf{h}, c_j\right) = -p(h_j = 1|y, X)$$
$$A\left(y, \mathbf{X}, \mathbf{h}, U_j\right) = -p(h_j = 1|y, X)y$$

$$A(y, \mathbf{X}, \mathbf{h}, \boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{h}|y,\mathbf{X})}\left[\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial \theta}\right] = \sum_{\mathbf{h}'} p(\mathbf{h}'|y, \mathbf{X})\frac{\partial E(y, \mathbf{X}, \mathbf{h}')}{\partial \theta}$$

$$A\left(y, \mathbf{X}, \mathbf{h}, W_{j,k}^{(i)}\right) = -\sum_{\mathbf{h}'} p(\mathbf{h}'|y, \mathbf{X})h_j 1_{(x_i=k)} = -p(h_j = 1|y, X)1_{(x_i=k)} \prod_{q,q \neq j} \sum_{h_q \in \{0,1\}} p(h_q|y, X)$$

$$= -p(h_j = 1|y, X)1_{(x_i=k)}$$

$$A\left(y, \mathbf{X}, \mathbf{h}, b_k^i\right) = -\sum_{\mathbf{h}'} p(\mathbf{h}'|y, \mathbf{X})1_{(x_i=k)} = -1_{(x_i=k)}$$

$$A(y, \mathbf{X}, \mathbf{h}, d) = -\sum_{\mathbf{h}'} p(\mathbf{h}'|y, \mathbf{X})y = -y$$

$$A\left(y, \mathbf{X}, \mathbf{h}, c_j\right) = -\sum_{\mathbf{h}'} p(\mathbf{h}'|y, \mathbf{X})h_j = -p(h_j = 1|y, X)$$

$$A\left(y, \mathbf{X}, \mathbf{h}, U_j\right) = -\sum_{\mathbf{h}'} p(\mathbf{h}'|y, \mathbf{X})h_j y = -p(h_j = 1|y, X)y$$

$$B(y, \mathbf{X}, \mathbf{h}, \boldsymbol{\theta})$$

$$B\left(y, \mathbf{X}, \mathbf{h}, W_{j,k}^{(i)}\right) = -1_{(x_i=k)} \sum_{y'} p(y'|\mathbf{X})p(h_j = 1|y', \mathbf{X})$$

$$B\left(y, \mathbf{X}, \mathbf{h}, b_k^i\right) = -1_{(x_i=k)}$$

$$B(y, \mathbf{X}, \mathbf{h}, d) = -\sum_{y'} p(y|\mathbf{X})y'$$

$$B\left(y, \mathbf{X}, \mathbf{h}, c_j\right) = -\sum_{y'} p(y'|\mathbf{X})p(h_j = 1|y', \mathbf{X})$$

$$B\left(y, \mathbf{X}, \mathbf{h}, U_j\right) = -p(y = 1|\mathbf{X})p(h_j = 1|y = 1, \mathbf{X})$$

$$B(y, \mathbf{X}, \mathbf{h}, \boldsymbol{\theta}) = \mathbb{E}_{p(y,\mathbf{h}|\mathbf{X})}\left[\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial \theta}\right] = \sum_{y'}\sum_{\mathbf{h}'} p(y', \mathbf{h}'|\mathbf{X})\frac{\partial E(y', \mathbf{X}, \mathbf{h}')}{\partial \theta}$$

$$B\left(y, \mathbf{X}, \mathbf{h}, W_{j,k}^{(i)}\right) = -\sum_{y'}\sum_{\mathbf{h}'} p(y', \mathbf{h}'|\mathbf{X}) h_j 1_{(x_i=k)} = -1_{(x_i=k)} \sum_{y'} p(y|\mathbf{X}) \sum_{\mathbf{h}'} p(\mathbf{h}'|y, \mathbf{X})h_j$$

$$= -1_{(x_i=k)} \sum_{y'} p(y'|\mathbf{X})p(h_j = 1|y', \mathbf{X})$$

$$B\left(y, \mathbf{X}, \mathbf{h}, b_k^i\right) = -\sum_{y'}\sum_{\mathbf{h}'} p(y', \mathbf{h}'|\mathbf{X}) h_j 1_{(x_i=k)} = -1_{(x_i=k)}$$

$$B(y, \mathbf{X}, \mathbf{h}, d) = -\sum_{y'}\sum_{\mathbf{h}'} p(y', \mathbf{h}'|\mathbf{X})y = -\sum_{y'} p(y'|\mathbf{X})y' \sum_{\mathbf{h}'} p(\mathbf{h}'|y', \mathbf{X}) = -\sum_{y'} p(y|\mathbf{X})y'$$

$$B\left(y, \mathbf{X}, \mathbf{h}, c_j\right) = -\sum_{y'} p(y|\mathbf{X}) \sum_{\mathbf{h}'} p(\mathbf{h}'|y, \mathbf{X})h_j = -\sum_{y'} p(y'|\mathbf{X})p(h_j = 1|y', \mathbf{X})$$

$$B\left(y, \mathbf{X}, \mathbf{h}, U_j\right) = -\sum_{y'} y' p(y'|\mathbf{X}) \sum_{\mathbf{h}'} p(\mathbf{h}'|y', \mathbf{X})h_j = -p(y = 1|\mathbf{X})p(h_j = 1|y = 1, \mathbf{X})$$

$$C(y, \mathbf{X}, \mathbf{h}, \boldsymbol{\theta}) \text{ Using CD-k Approximation}$$

$$C\left(y, \mathbf{X}, \mathbf{h}, W_{j,k}^{(i)}\right) = -h_j 1_{(x_i=k)}$$

$$C\left(y, \mathbf{X}, \mathbf{h}, b_k^i\right) = -1_{(\hat{x}_i=k)}$$

$$C(y, \mathbf{X}, \mathbf{h}, d) = -\hat{y}$$

$$C\left(y, \mathbf{X}, \mathbf{h}, c_j\right) = -\widehat{h_j}$$

$$C\left(y, \mathbf{X}, \mathbf{h}, U_j\right) = -\widehat{h_j}\hat{y}$$

$$C(y, \mathbf{X}, \mathbf{h}, \boldsymbol{\theta}) = \mathbb{E}_{p(y,\mathbf{X},\mathbf{h})}\left[\frac{\partial E(y, \mathbf{X}, \mathbf{h})}{\partial \theta}\right] = \sum_{y',\mathbf{X}',\mathbf{h}'} p(y', \mathbf{X}', \mathbf{h}')\frac{\partial E(y', \mathbf{X}', \mathbf{h}')}{\partial \theta} \approx \frac{\partial E(\hat{y}, \widehat{\mathbf{X}}, \hat{\mathbf{h}})}{\partial \theta}$$

$$C\left(y, \mathbf{X}, \mathbf{h}, W_{j,k}^{(i)}\right) = -\widehat{h_j} 1_{(\hat{x}_i=k)}$$

$$C\left(y, \mathbf{X}, \mathbf{h}, b_k^i\right) = -1_{(\hat{x}_i=k)}$$

$$C(y, \mathbf{X}, \mathbf{h}, d) = -\hat{y}$$
$$C(y, \mathbf{X}, \mathbf{h}, c_j) = -\widehat{h}_j$$
$$C(y, \mathbf{X}, \mathbf{h}, U_j) = -\widehat{h}_j \hat{y}$$

Combine All

$$\frac{\ln Hybrid(\alpha, y, \mathbf{X})}{\partial \theta} = -(1 + \alpha)A(y, \mathbf{X}, \mathbf{h}, \theta) + B(y, \mathbf{X}, \mathbf{h}, \theta) + \alpha C(y, \mathbf{X}, \mathbf{h}, \theta)$$

$$\frac{\partial \ln Hybrid(\alpha, y, \mathbf{X})}{\partial W_{j,k}^{(i)}} = 1_{(x_i = k)} \left\{ (1 + \alpha)p(h_j = 1|y, X) - \sum_{y'} p(y'|\mathbf{X})p(h_j = 1|y', \mathbf{X}) - \alpha \widehat{h}_j \right\}$$

$$\frac{\partial \ln Hybrid(\alpha, y, \mathbf{X})}{\partial b_j^i} = \alpha \left[ 1_{(x_i = k)} - 1_{(\hat{x}_i = k)} \right]$$

$$\frac{\partial \ln Hybrid(\alpha, y, \mathbf{X})}{\partial d} = (1 + \alpha)y - p(y = 1|\mathbf{X}) - \alpha \hat{y}$$

$$\frac{\partial \ln Hybrid(\alpha, y, \mathbf{X})}{\partial c_i} = (1 + \alpha)p(h_j = 1|y, X) - \sum_{y'} p(y'|\mathbf{X})p(h_j = 1|y', \mathbf{X}) - \alpha \widehat{h}_j$$

$$\frac{\partial \ln Hybrid(\alpha, y, \mathbf{X})}{\partial U_i} = (1 + \alpha)p(h_j = 1|y, \mathbf{X})y - p(y = 1|\mathbf{X})p(h_j = 1|y = 1, \mathbf{X}) - \alpha \widehat{h}_j \hat{y}$$

CD-K Updates

$$X^{(0)} \leftarrow X, y^{(0)} \leftarrow y$$
For t = 0, …, k − 1 do

    For i = 1, …, H do sample $h_i^{(t)} \sim p(h_i = 1|y^{(t)}, X^{(t)})$

    For i = 1, …, C do sample $x_i^{(t+1)} \sim p(x_i|y^{(t)}, h^{(t)})$

    Do Sample $y^{(t+1)} \sim p(y = 1|h^{(t)})$

For i = 1,…H do $h_i^{(k)} \leftarrow p(h_i = 1|y^{(k)}, X^{(k)})$

# Maximum Likelihood and the Delta Rule

- Maximum Likelihood

$$\ln L(\vec{\theta}|\vec{S}) = \ln \prod_{i=1}^{l} p(\overrightarrow{v^i}|\vec{\theta}) = \sum_{i=1}^{l} \ln p(\overrightarrow{v^i}|\vec{\theta})$$

- Mini-batch Gradient Ascent (Delta Rule)

$$\overrightarrow{\theta^{t+1}} = \overrightarrow{\theta^t} + \eta \underbrace{\frac{\partial}{\partial \overrightarrow{\theta^t}} \left[ \sum_{i=1}^{N} \ln L(\overrightarrow{\theta^t}|\overrightarrow{v^i}) \right] - \lambda \overrightarrow{\theta^t} + v \Delta \overrightarrow{\theta^{t-1}}}_{:= \Delta \overrightarrow{\theta^t}}$$

# Learning Algorithm

//CD-k, only if $\alpha \neq 0$

$X^{(0)} \leftarrow X, y^{(0)} \leftarrow y$

For t = 0, …, k − 1 do

    For j = 1, …, H do sample $h_j^{(t)} \sim p\left(h_j = 1 | y^{(t)}, \mathbf{X^{(t)}}\right) = \sigma(\sum_{1 \leq i \leq C} W_{j,x_i^{(t)}}^{(i)} + c_j + U_j y^{(t)})$ [O(C+H)]

    For i = 1, …, C do sample $x_i^{(t+1)} \sim \boldsymbol{p(x_i = k | \mathbf{h^{(t)}})} = \dfrac{\exp(\sum_{1 \leq j \leq H} h_j^{(t)} W_{j,k}^{(i)})}{\sum_{1 \leq q \leq C_i} \exp(\sum_{1 \leq j \leq H} h_j^{(t)} W_{j,q}^{(i)})}$ (Computation Intensive! **O(H\*V)**)

    Do Sample $y^{(t+1)} \sim p\left(y = 1 | \mathbf{h^{(t)}}\right) = \sigma(d + \mathbf{h^{(t)T}U})$ [O(H)]

For i = 1,…H do $h_i^{(k)} \leftarrow p(h_i = 1 | y^{(k)}, X^{(k)}) = \sigma(\sum_{1 \leq i \leq C} W_{j,x_i^{(k)}}^{(i)} + c_j + U_j y^{(k)})$   [O(H)]


//Gradient Calculation

$\hat{y} \leftarrow y^{(k)}, \hat{X} \leftarrow X^{(k)}, \hat{h} \leftarrow h^{(k)}$

$$\boldsymbol{p(h_j = 1 | y, X) = \sigma(} \sum_{1 \leq i \leq C} \boldsymbol{W_{j,x_i}^{(i)} + c_j + U_j y)}$$

$$\mathbf{p(y = 1 | X)} = \frac{\exp\{dy + \sum_{1 \leq k \leq H} \text{softplus}\left(\sum_{1 \leq i \leq C} W_{k,x_i}^{(i)} + c_k + U_k y\right)\}}{\exp\{\sum_{1 \leq k \leq H} \text{softplus}\left(\sum_{1 \leq i \leq C} W_{k,x_i}^{(i)} + c_k\right)\} + \exp\{d + \sum_{1 \leq k \leq H} \text{softplus}\left(\sum_{1 \leq i \leq C} W_{k,x_i}^{(i)} + c_k + U_k\right)\}}$$

[O(H\*C)]


$$\frac{\partial \ln Hybrid(\alpha, y, \mathbf{X})}{\partial W_{j,k}^{(i)}} = \mathbf{1}_{(x_i = k)} \left\{ (1+\alpha) \boldsymbol{p(h_j = 1 | y, X)} - \sum_{y'} \mathbf{p(y' | X) p(h_j = 1 | y', X)} - \alpha \widehat{\boldsymbol{h_j}} \right\}$$

$$\frac{\partial \ln Hybrid(\alpha, y, \mathbf{X})}{\partial b_j^i} = \alpha \left[ \mathbf{1}_{(x_i = k)} - \mathbf{1}_{(\hat{x}_i = k)} \right]$$

$$\frac{\partial \ln Hybrid(\alpha, y, \mathbf{X})}{\partial d} = (1+\alpha) y - \mathbf{p(y = 1 | X)} - \alpha \hat{\boldsymbol{y}}$$

$$\frac{\partial \ln Hybrid(\alpha, y, \mathbf{X})}{\partial c_i} = (1+\alpha) \mathbf{p(h_j = 1 | y, X)} - \sum_{y'} \mathbf{p(y' | X) p(h_j = 1 | y', X)} - \alpha \widehat{\boldsymbol{h_j}}$$

$$\frac{\partial \ln Hybrid(\alpha, y, \mathbf{X})}{\partial U_i} = (1+\alpha) \mathbf{p(h_j = 1 | y, X)} y - \mathbf{p(y = 1 | X) p(h_j = 1 | y = 1, X)} - \alpha \widehat{\boldsymbol{h_j}} \hat{\boldsymbol{y}}$$


//Parameter Update

$$\overrightarrow{\theta^{t+1}} = \overrightarrow{\theta^t} + \eta \underbrace{\frac{\partial}{\partial \overrightarrow{\theta^t}} \left[ \sum_{i=1}^{N} \ln L(\overrightarrow{\theta^t} | \overrightarrow{v^i}) \right] - \lambda \overrightarrow{\theta^t} + v \Delta \overrightarrow{\theta^{t-1}}}_{:= \Delta \overrightarrow{\theta^t}}$$

[Note, storing two sets of parameters would be memory hungry, store the differences only??]

Note:

Efficient calculation of the following is the key to the performance of CD-k!

Naïve implementation requires $O(|X| * |H|)$

$$x_i^{(t+1)} \sim p(x_i = k | \mathbf{h}^{(t)}) = \frac{\exp\left(\sum_{1 \le j \le H} h_j^{(t)} W_{j,k}^{(i)}\right)}{\sum_{1 \le q \le C_i} \exp\left(\sum_{1 \le j \le H} h_j^{(t)} W_{j,q}^{(i)}\right)}, 1 \le k \le C_i$$

Use Mini-batch + Cache Strategy??

Nonexact sampling??MCMC???Importance Sampling, rejection sampling etc to avoid the normalization constant? ????

What else??

# SparseClassRBM Verses Logistic Regression

Note that in logisitc regression, we have:

$$p(y|X) = \frac{\exp\{y(W^T x + c)\}}{1 + \exp\{W^T x + c\}}$$

And for SparseClassRBM, we have

$$p(y|X) = \frac{\exp\{dy\} \prod_{1 \le k \le H} \sum_{h_k \in \{0,1\}} \exp\{h_k[\sum_{1 \le i \le C} W_{k,x_i}^{(i)} + c_k + U_k y]\}}{\sum_{y'} \exp\{dy'\} \prod_{1 \le k \le H} \sum_{h_k \in \{0,1\}} \exp\{h_k[\sum_{1 \le i \le C} W_{k,x_i}^{(i)} + c_k + U_k y']\}}$$

Setting H=1, $h_1 = 1$, d=1, U = 0, we have

$$p(y|X) = \frac{\exp\{y\left(\sum_{1 \le i \le C} W_{x_i}^{(i)} + c\right)\}}{1 + \exp\{\sum_{1 \le i \le C} W_{x_i}^{(i)} + c\}}$$

which is a form of logistic regression.

Therefore, we can view logistic regression as a special form of RBM with less variables.