

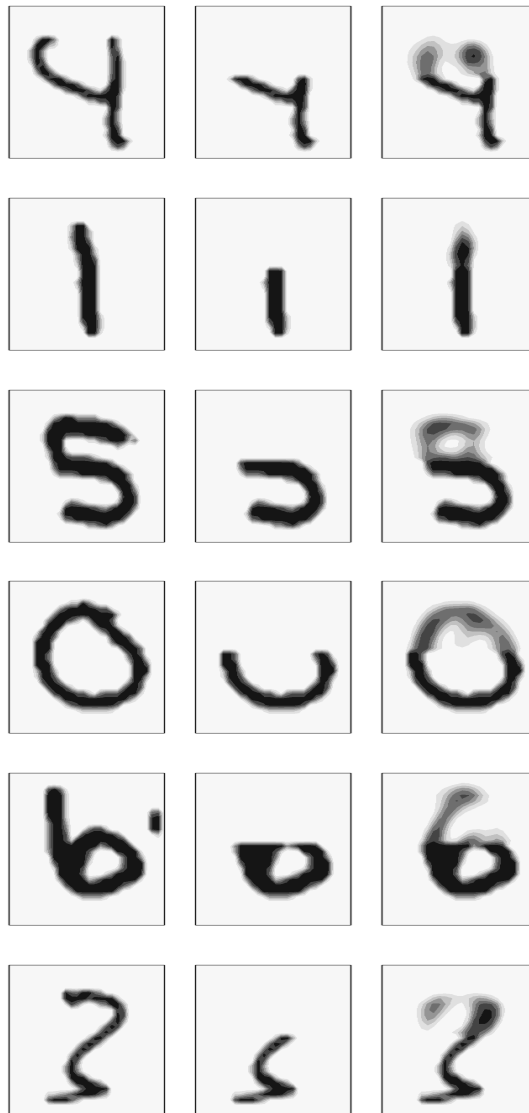
# Implement An Unsupervised Learning Machine and Corrupted Data Reconstruction

Han Chen

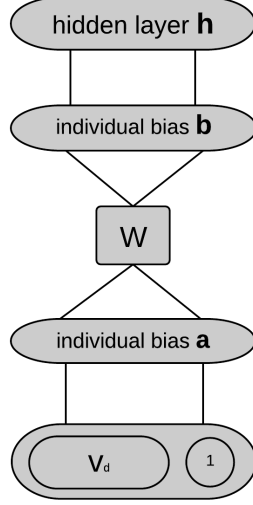
April 1, 2016

## 1 Results

- First column: original image.
- Second column: corrupt image by covering the top half.
- Third column: image inferred from the corrupted image.



## 2 Restricted Boltzmann Machine



The visible and hidden variables are conditional Bernoulli

$$\begin{aligned} P(v_i = 1|h) &= \sigma(a_i + W_{i,\cdot} \cdot h) \\ P(h_j = 1|v) &= \sigma(b_j + W_{\cdot,j} \cdot v) \end{aligned} \quad (1)$$

where

$$\sigma(x) = \frac{1}{1 + e^{-x}}. \quad (2)$$

Gradient of data log-likelihood

$$\begin{aligned} \frac{\partial \log p(v)}{\partial W_{i,j}} &= \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}} \\ \frac{\partial \log p(v)}{\partial a_i} &= \langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}} \\ \frac{\partial \log p(v)}{\partial b_j} &= \langle h_j \rangle_{\text{data}} - \langle h_j \rangle_{\text{model}} \end{aligned} \quad (3)$$

### 2.1 Sampling by Contrastive-Divergence

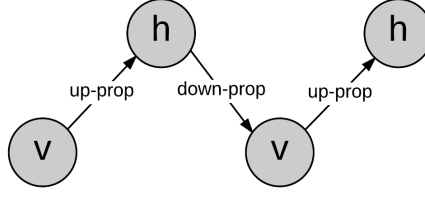
Let  $N$  be the cardinality of the training data. The integrals in (3) can be approximated by

$$\begin{aligned} \langle v_i h_j \rangle_{\text{data}} &\approx \frac{1}{N} \sum_{n=1}^N v_i^{(n)} h_j^{(n)} \\ \langle v_i \rangle_{\text{data}} &\approx \frac{1}{N} \sum_{n=1}^N v_i^{(n)} \quad , \quad h_j^{(n)} \sim P(h_j = 1|v^{(n)}) \\ \langle h_j \rangle_{\text{data}} &\approx \frac{1}{N} \sum_{n=1}^N \sigma(b_j + W_{\cdot,j} v^{(n)}) \end{aligned} \quad (4)$$

and

$$\begin{aligned} \langle v_i h_j \rangle_{\text{model}} &\approx \frac{1}{N} \sum_{n=1}^N \sigma(a_i + W_{i,\cdot} \cdot h^{(n)}) h_j^{(n)} \\ \langle v_i \rangle_{\text{model}} &\approx \frac{1}{N} \sum_{n=1}^N \sigma(a_i + W_{i,\cdot} \cdot h^{(n)}) \quad , \quad \begin{aligned} h_j^{(n)} &\sim P(h_j = 1|v^{(n)}) \\ \hat{v}^{(n)} &\sim P(\hat{v}_i = 1|h^{(n)}) \end{aligned} \\ \langle h_j \rangle_{\text{model}} &\approx \frac{1}{N} \sum_{n=1}^N \sigma(b_j + W_{\cdot,j} \cdot \hat{v}^{(n)}) \end{aligned} \quad (5)$$

(4) and (5) can be expressed by the graphical model



In the implementation, the up-propagation and down-propagation methods produce both the probability and a binary sampling.

## 2.2 Stochastic Gradient Descent

Given a learning rate, (4) and (5) drives the maximization of data log-likelihood by stochastic gradient descent according to (3).

## 2.3 Regularization

An L-1 regularizer for the weights is implemented into the log-likelihood, which modifies (3) to

$$\begin{aligned}\frac{\partial \log p(v)}{\partial W_{i,j}} &= \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}} - \lambda \delta_{W_{i,j}} \\ \frac{\partial \log p(v)}{\partial a_i} &= \langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}} \\ \frac{\partial \log p(v)}{\partial b_j} &= \langle h_j \rangle_{\text{data}} - \langle h_j \rangle_{\text{model}}\end{aligned}\quad , \quad (6)$$

where

$$\delta_x = \begin{cases} 1, & x > 0 \\ 0, & x = 0 \\ -1, & x < 0 \end{cases} \quad (7)$$

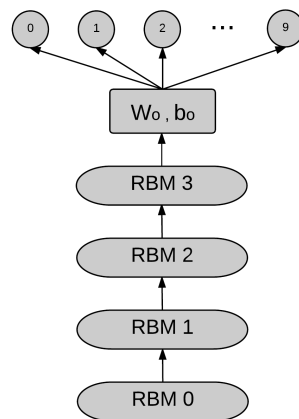
and  $\lambda > 0$  is a small number.

## 2.4 Termination Condition

The training terminates when the free energy of the validation dataset exceeds the free energy of the training data by a certain amount. The free energy is defined by

$$\mathcal{F}(x) = -a \cdot v - \sum_{j=1}^H \log \left( 1 + \exp(b_j + W_{\cdot,j} \cdot v) \right) \quad (8)$$

### 3 Stacked Restricted Boltzmann Machine



#### 3.1 Fine Tuning

After the tuning of each RBM, the RBMs are stacked into a deep neural network. The neural network is implemented using automatic differentiation whose weights and biases are initialized by the unsupervised trained results. The network is tuned by back propagation, where gradients are computed by automatic differentiation.

### 4 Trained RBM Weights



## 5 Gibbs Sampling Transition

