

The Deep Learning Approach to Structured Probabilistic Models

Sargur N. Srihari

srihari@cedar.buffalo.edu

Topics in Structured PGMs for Deep Learning

0. Overview

1. Challenge of Unstructured Modeling

2. Using graphs to describe model structure

3. Sampling from graphical models

4. Advantages of structured modeling

5. Learning about Dependencies

6. Inference and Approximate Inference

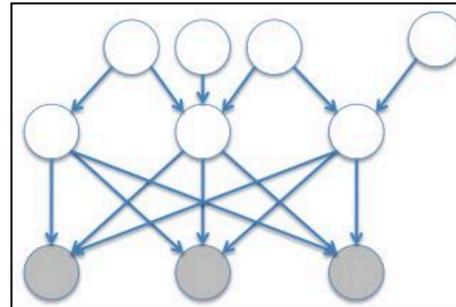
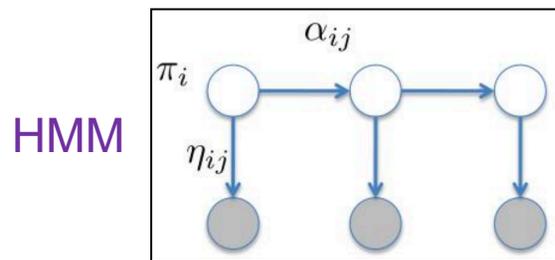
7. Deep learning approach to structured models

1. Ex: The Restricted Boltzmann machine

2. Training an RBM

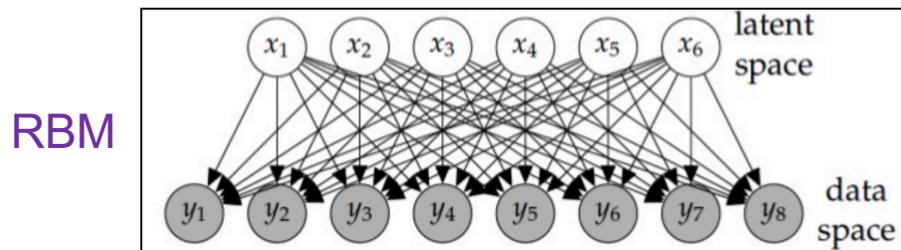
Special nature of PGMs in Deep Learning

- General PGMs



Shaded nodes are observed

- PGMs in Deep Learning



- Same computational tools used

- But different design decisions on combining tools

- Very different flavor from traditional PGMs

Differences in PGM treatment

- Traditional PGMs vs. PGMs in deep learning
 1. Depth
 2. Proportion of observed to latent variables
 3. Latent semantics (meaning of a latent variable)
 4. Connectivity and inference algorithm
 5. Intractability and approximation

1. Depth of PGM

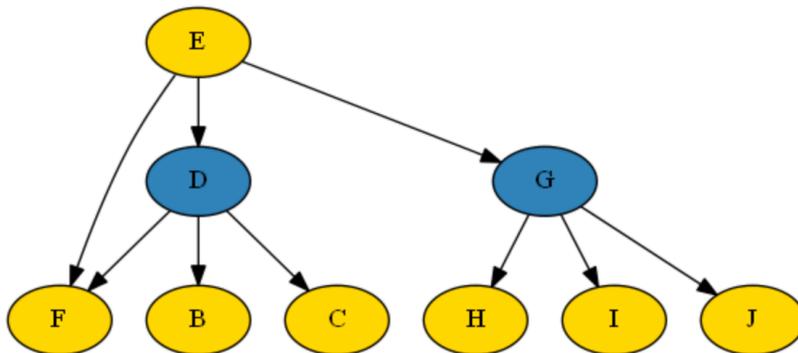
- Depth in a PGM:
 - Latent variable h_i is at depth j if shortest path from h_i to observed variable is j steps
 - Depth of model is the greatest such depth
- Depth in a computational graph (deep learning)
 - Generative models often have no latent variables or only one layer of latent variables
 - But use deep computational graphs to define the conditional distributions within a model
- PGMs in deep learning are not *deep* PGMs

2. Proportion of Observed/Latent Variables

- Deep Learning has more latent variables than observed variables
 - Since it always uses distributed representations
 - Even shallow models (such as for pretraining) have a single *large* layer of latent variables
 - Complicated nonlinear interactions between variables are accomplished via indirect connections that flow through multiple latent variables
- Traditional PGMs contain mostly variables that are observed (i.e., few latent variables)

3. Latent variable semantics

- In DL
 - Not take any specific semantics ahead of time
 - Training algorithm invents concepts to model data
 - Latent variables not easy to interpret after the fact
- In traditional PGMs
 - Specific semantics in mind
 - Document Topic, intelligence, disease causing symptoms



4. Connectivity

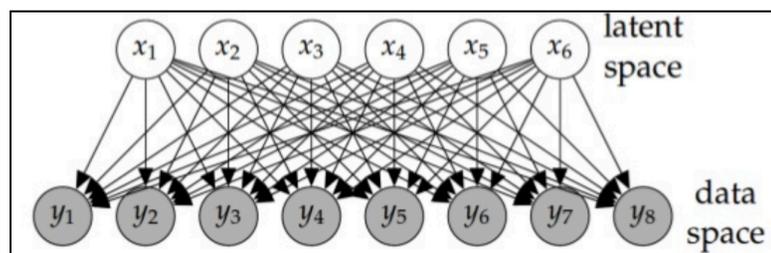
- Deep learning PGMs have large groups of units connected other large groups of units
- So that interactions can be described by a single matrix
- Traditional PGMs have few connections and choice of connections is individually designed
- Design of model structure is tightly linked to choice of inference algorithm

5. Inference

- Traditional PGMs: tractability of exact inference
 - When this is too limiting a popular approximate approach is loopy belief propagation
 - Both approaches work well with sparsely connected graphs
- Deep learning PGMs: graphs are not sparse
 - Use either Gibbs sampling or variational inference
 - Rather than simplifying model until exact inference is feasible, make model complex enough as long as we can compute a gradient

Example: RBMs or Harmonium

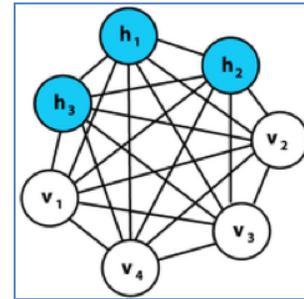
- RBMs are a quintessential example of how graphical models are used for deep learning



- RBM itself is not a deep model
 - It has a single layer of latent units that may be used to learn a representation for the input
 - RBMs can be used to build many deeper models

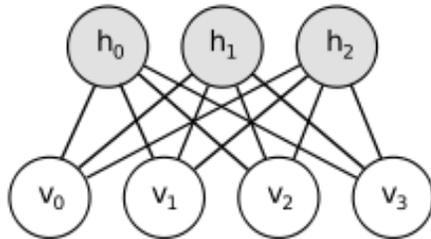
Boltzmann Machine and RBM

- A general Boltzmann machine can have arbitrary connections



General BM

- RBM: No visible-visible and hidden-hidden connections— Bipartite graph



- Used to learn features for input to neural networks in Deep Learning

RBM Characteristics

- Units are organized into large groups called layers
- Connectivity between layers is described by a matrix
- Connectivity is relatively dense
- Allows efficient Gibbs sampling
- Learns latent variables whose semantics is not defined by the designer

Canonical RBM

- A general energy-based model is defined as

$$\tilde{p}(\mathbf{x}) = \exp(-E(\mathbf{x})) \quad p(\mathbf{x}) = \frac{1}{Z} \tilde{p}(\mathbf{x}) \quad Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}$$

Factors are of the form $\exp(a)$, $\exp(b)=\exp(a+b)$

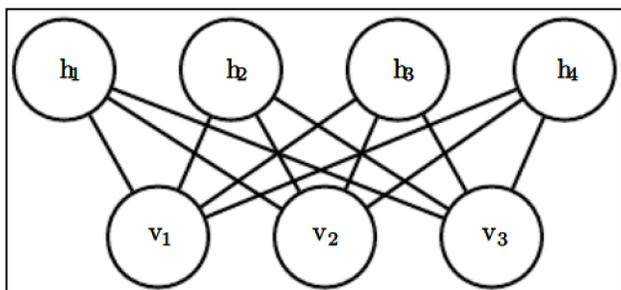
- RBM is an energy based model with binary hidden and visible units, with $\mathbf{x}=\{\mathbf{v},\mathbf{h}\}$

$$E(\mathbf{v},\mathbf{h})= -\mathbf{b}^T \mathbf{v} -\mathbf{c}^T \mathbf{h} -\mathbf{v}^T \mathbf{W} \mathbf{h}$$

- where \mathbf{b} , \mathbf{c} and \mathbf{W} are unconstrained, real-valued learnable parameters
- We can see that the model is divided into two groups of units \mathbf{v} and \mathbf{h} and the interaction between them is described by matrix \mathbf{W}
- Model is graphically depicted next

An RBM drawn as a Markov network

- The model is depicted graphically as



$$\tilde{p}(\mathbf{v}, \mathbf{h}) = \exp(-E(\mathbf{v}, \mathbf{h}))$$

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}$$

$$\mathbf{x} = \{\mathbf{v}, \mathbf{h}\}$$

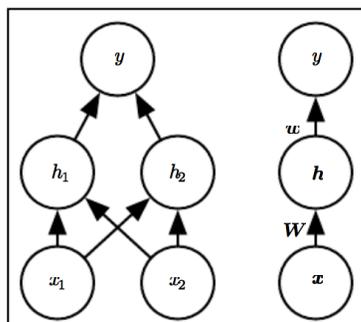
$$p(\mathbf{x}) = \frac{1}{Z} \hat{p}(\mathbf{x})$$

$$Z = \int \tilde{p}(\mathbf{x}) d\mathbf{x}$$

- None between visible units or between hidden units
 - Hence *restricted*, a general Boltzmann has arbitrary connections

- This is different from a neural network

Two representational forms of a neural network:



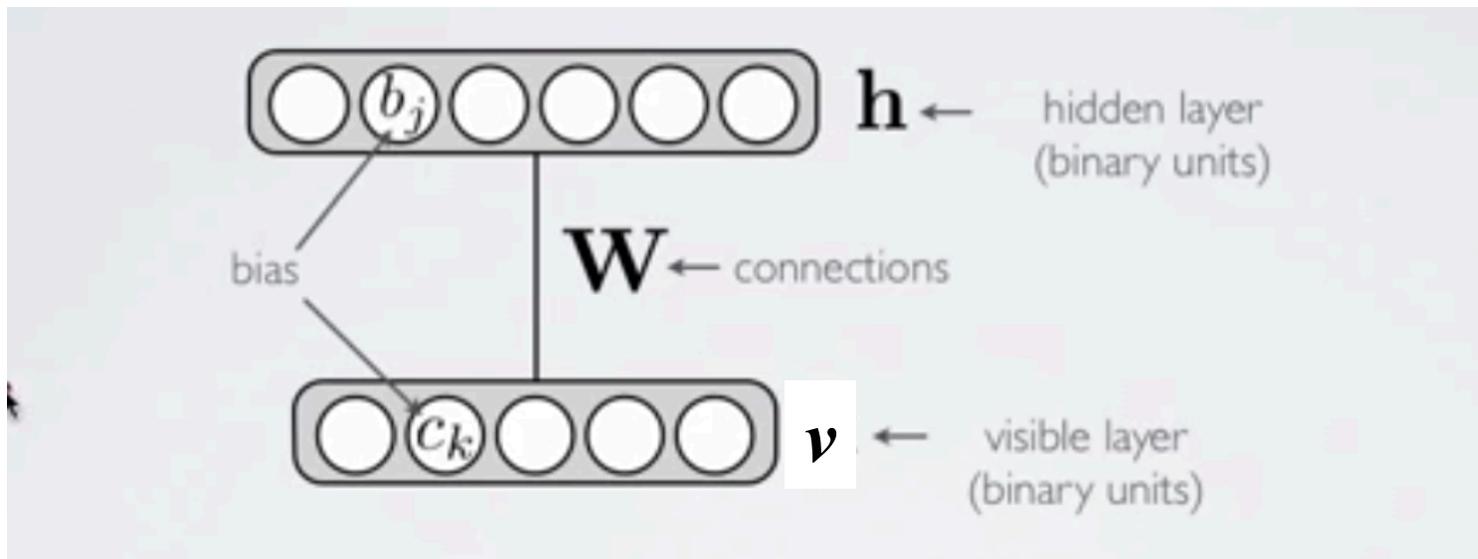
$$\mathbf{h} = g(\mathbf{W}^T \mathbf{x} + \mathbf{c})$$

$g(z) = \max\{0, z\}$, is the ReLU activation

$$y = \mathbf{h}^T \mathbf{w} + b$$

Second form used for RBM in next slide

Network view of RBM



$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}$$

$$= -\sum_j \sum_k W_{j,k} h_j v_k - \sum_k c_k v_k - \sum_j b_j h_j$$

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

partition function
(intractable)

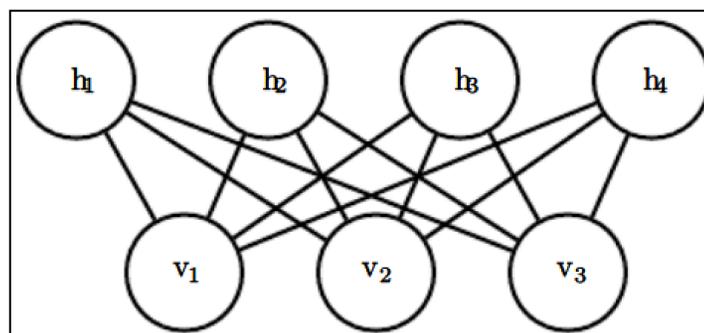
Properties of RBMs

- Restrictions of RBM structure yields nice properties:

$$p(\mathbf{h}|\mathbf{v}) = \prod_i p(h_i|\mathbf{v}) \text{ and}$$

$$p(\mathbf{v}|\mathbf{h}) = \prod_i p(v_i|\mathbf{h})$$

Since nodes at same level are independent



- Individual conditionals are simple to compute
 - For binary RBM $v_j, h_i \in \{0,1\}$ we obtain

$$p(h_i=1|\mathbf{v}) = \sigma(\mathbf{v}^T \mathbf{W}_{:,i} + b_i)$$

$$p(h_i=0|\mathbf{v}) = 1 - \sigma(\mathbf{v}^T \mathbf{W}_{:,i} + b_i)$$

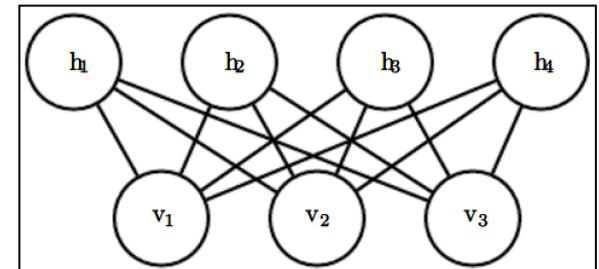
$$\tilde{p}(\mathbf{v}, \mathbf{h}) = \exp(-E(\mathbf{v}, \mathbf{h}))$$

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}$$

Block Gibbs Sampling with *RBM*

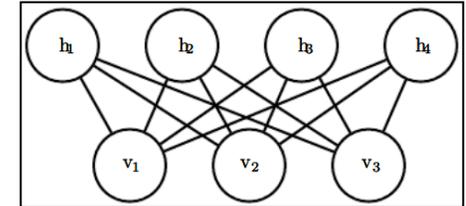
- RBM properties allow for *block Gibbs sampling*
 - Alternate between sampling all \mathbf{h} simultaneously and all \mathbf{v} simultaneously
 - Gibbs sampling and samples obtained shown next

- Gibbs sampling with M variables
 - Initialize first sample: $\{z_i, i = 1, \dots, M\}$
 - For $t=1, \dots, T$, $T = \text{no of samples}$
 - Sample $z_1^{(\tau+1)} \sim p(z_1 | z_2^{(\tau)}, z_3^{(\tau)}, \dots, z_M^{(\tau)})$
 - Sample $z_2^{(\tau+1)} \sim p(z_2 | z_1^{(\tau+1)}, z_3^{(\tau)}, \dots, z_M^{(\tau)})$
 -
 - Sample $z_j^{(\tau+1)} \sim p(z_j | z_1^{(\tau+1)}, \dots, z_{j-1}^{(\tau+1)}, z_{j+1}^{(\tau)}, \dots, z_M^{(\tau)})$
 -
 - Sample $z_M^{(\tau+1)} \sim p(z_M | z_1^{(\tau+1)}, z_2^{(\tau+1)}, \dots, z_{M-1}^{(\tau+1)})$

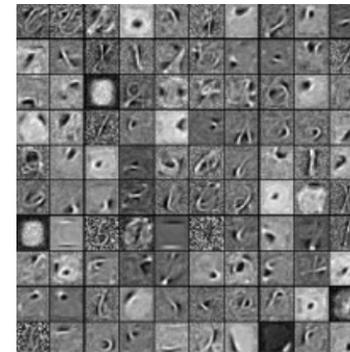


Samples from an RBM

- RBM Gibbs samples
 - Model trained on MNIST data



Corresponding weight vectors:



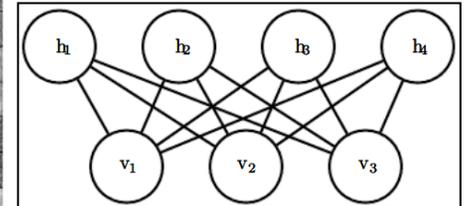
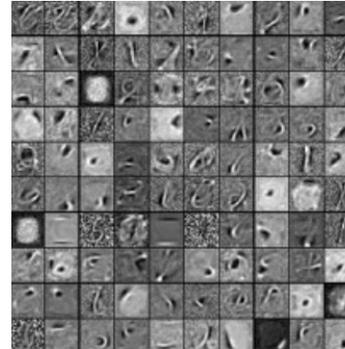
- Each column is a separate Gibbs process
- Each row represents the output of another 1000 steps of Gibbs sampling, Successive samples are highly correlated

RBM: Comparison with Sparse Coding

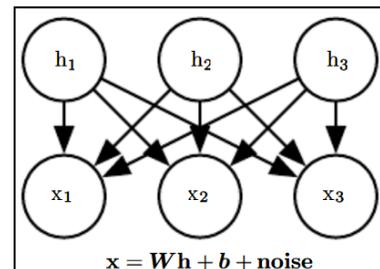
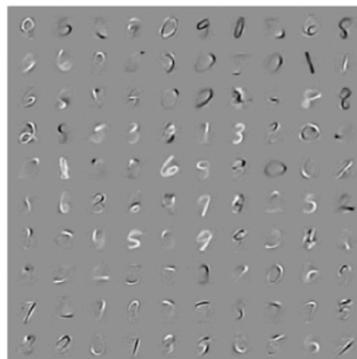
- RBM Gibbs samples (model trained on MNIST data)



Corresponding weight vectors:



- Linear factor model (e.g., PCA) samples
 - Sparse Coding Samples are blurry



Derivatives of Energy Function

- Energy function: $E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}$
 - where \mathbf{b} , \mathbf{c} and \mathbf{W} are unconstrained, real-valued learnable parameters
- Since the energy function is a linear function of its parameters, it is easy to take derivatives
 - E.g., $\frac{\partial}{\partial W_{i,j}} E(\mathbf{v}, \mathbf{h}) = -v_i h_j$
- These two properties, efficient Gibbs sampling and efficient derivatives make training convenient
 - Described next

Training Undirected Models

- Training undirected models is accomplished by computing derivatives $\frac{\partial}{\partial W_{i,j}} E(\mathbf{v}, \mathbf{h}) = -v_i h_j$
 - And applying to samples from the model
 - Training induces a representation \mathbf{h} of data \mathbf{v}
 - We can use $E_{\mathbf{h} \sim p(\mathbf{h}|\mathbf{v})}[\mathbf{h}]$ as a set of features to describe \mathbf{v}
- RBM demonstrates the typical deep learning approach to graphical models:
 - Representation learning accomplished via layers of latent variables combined with efficient interactions between layers parameterized by matrices

RBM with visible and hidden units

- Joint configuration (\mathbf{v}, \mathbf{h})

- visible and hidden units has an energy (Hopfield 1982)

$$E(\mathbf{v}, \mathbf{h}) = - \sum_{i \in \text{visible}} a_i v_i - \sum_{j \in \text{hidden}} b_j h_j - \sum_{i,j} v_i h_j w_{ij}$$

- Network assigns a probability to every pair of hidden and visible vectors

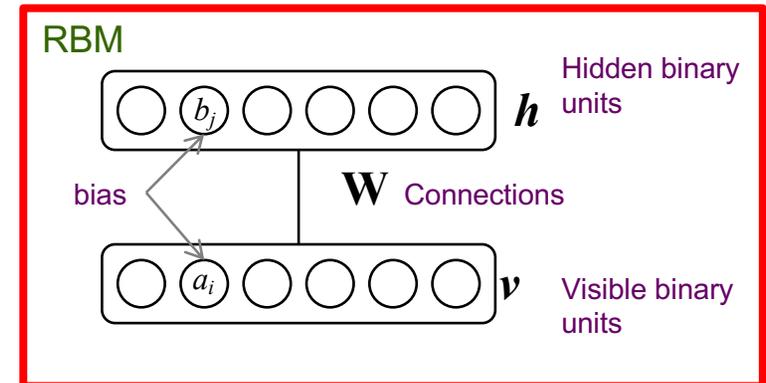
$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})}$$

- where partition function Z is a sum over all possible pairs of visible/hidden vectors

$$Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$$

- Probability that network assigns to a visible vector \mathbf{v} is

$$p(\mathbf{v}) = \frac{1}{Z} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$$



Stochastic binary pixels \mathbf{v} connected to stochastic binary feature detectors \mathbf{h} using symmetrically weighted connections

Changing probability of \mathbf{v}

- Probability network assigns to training data \mathbf{v} is raised by adjusting weights and biases
 - Lower the energy of that image & raise energy of other images
 - Especially those that have low energies and make a high contribution to the partition function
- Maximum likelihood approach to determine W, a, b

$$\text{Likelihood: } P(\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(M)}\}) = \prod_m p(\mathbf{v}^{(m)})$$

Log-likelihood:

$$\ln P(\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(M)}\}) = \sum_m \ln p(\mathbf{v}^{(m)}) = \sum_m \ln \left(\frac{1}{Z} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})^{(m)}} \right) = \sum_m \ln \left(\sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})^{(m)}} \right) - \sum_m \ln \left(\sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \right)$$

Derivative of the log-probability of a training vector wrt a weight:

$$\frac{\partial \ln p(\mathbf{v})}{\partial w_{ij}} = \mathbb{E}_{\text{data}}(v_i h_j) - \mathbb{E}_{\text{model}}(v_i h_j)$$

Learning rule for stochastic steepest ascent

$$\Delta w_{ij} = \varepsilon \left(\mathbb{E}_{\text{data}}(v_i h_j) - \mathbb{E}_{\text{model}}(v_i h_j) \right), \text{ where } \varepsilon \text{ is the learning rate}$$

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

$$p(\mathbf{v}) = \frac{1}{Z} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$$

$$E(\mathbf{v}, \mathbf{h}) = - \sum_{i \in \text{visible}} a_i v_i - \sum_{i \in \text{hidden}} b_j h_j - \sum_{i,j} v_i h_j w_{ij}$$

$$\frac{\partial}{\partial w_{i,j}} E(\mathbf{v}, \mathbf{h}) = -v_i h_j$$

$$\frac{d}{dx} \ln x = \frac{1}{x}$$

Error for gradient descent step is $E_{\text{data}}(v_i h_j) - E_{\text{model}}(v_i h_j)$

Samples for Computing Expectations

- Getting unbiased samples for $E_{\text{data}}(v_i h_j)$

- h_j : Given random training image \mathbf{v} , the binary state h_j for each hidden unit is set to 1 with probability

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

- v_i : Given a random training image \mathbf{v} , the binary state v_i for a visible unit is set to 1 with probability

$$p(v_i = 1 | \mathbf{v}) = \sigma \left(a_i + \sum_j h_j w_{ij} \right)$$

- Getting unbiased samples for $E_{\text{model}}(v_i h_j)$

- Can be done by starting at a random state of visible units and performing Gibbs sampling for a long time
 - One iteration of alternating Gibbs sampling consists of updating all hidden units in parallel followed by updating all visible units

Summary of RBM training

Probability Distribution of Undirected model (Gibbs)

$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \tilde{p}(\mathbf{x}, \boldsymbol{\theta})$$

$$Z(\boldsymbol{\theta}) = \sum_{\mathbf{x}} \tilde{p}(\mathbf{x}, \boldsymbol{\theta})$$

Intractable
Partition function

For an RBM: $\mathbf{x} = \{\mathbf{v}, \mathbf{h}\}$ $\boldsymbol{\theta} = \{\mathbf{W}, \mathbf{a}, \mathbf{b}\}$

Determine parameters $\boldsymbol{\theta}$ that maximize log-likelihood (negative loss)

$$\max_{\boldsymbol{\theta}} L(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}; \boldsymbol{\theta}) = \sum_m \log p(\mathbf{x}^{(m)}; \boldsymbol{\theta})$$

$$L(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(M)}\}; \boldsymbol{\theta}) = \sum_m \log \tilde{p}(\mathbf{x}^{(m)}; \boldsymbol{\theta}) - \sum_m \log Z(\boldsymbol{\theta})$$

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{h}^T \mathbf{W} \mathbf{v} - \mathbf{a}^T \mathbf{v} - \mathbf{b}^T \mathbf{h} = \sum_i \sum_j W_{i,j} v_i h_j - \sum_i a_i v_i - \sum_j b_j h_j$$

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

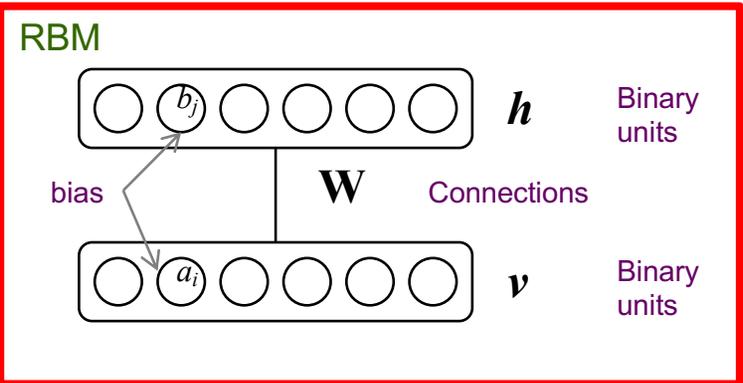
$$\tilde{p}(\mathbf{x}) = \exp(-E(\mathbf{v}, \mathbf{h}))$$

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

For stochastic gradient ascent, take derivatives:

$$g_m = \nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}^{(m)}; \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(m)}; \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta}) \quad \boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \varepsilon g$$

$$\frac{\partial}{\partial W_{i,j}} E(\mathbf{v}, \mathbf{h}) = -v_i h_j$$



Derivative of positive phase:

$$\frac{1}{M} \sum_{m=1}^M \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(m)}; \boldsymbol{\theta})$$

Summation is over samples
from the training set

Since it is summed m times $1/m$ has no effect

Derivative of negative phase:

$$\nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}) \quad \text{An identity}$$

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(m)}; \boldsymbol{\theta})$$

Summation is over samples from the RBM

RBM Training Algorithm

- Contrastive Divergence
- A method to overcome exponential complexity in dealing with the partition function