

Restricted Boltzmann Machines

Sargur N. Srihari
srihari@cedar.buffalo.edu

RBM or Harmonium

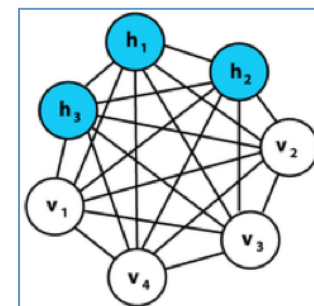
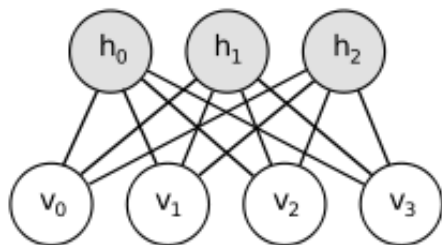
- RBMs are a quintessential example of how graphical models are used for deep learning
- RBM itself is not a deep model
- It has a single layer of latent units that may be used to learn a representation for the input
- RBMs can be used to build many deeper models

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

Restricted Boltzmann Machine

- RBM is a special case of Boltzmann machines and Markov networks
- No visible-visible and hidden-hidden connections— Bipartite graph



General BM

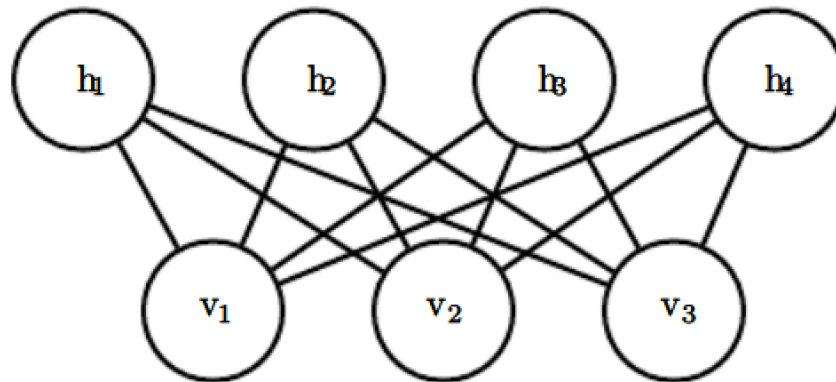
- Used to learn features for input to neural networks in Deep Learning

Canonical RBM

- Energy-based model with binary visible / hidden units
- Its energy function is
$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}$$
 - 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



- No direct interactions between any two visible units or between any two hidden units
 - Hence the “restricted”, a general BM may arbitrary connections

Properties of RBMs

- Restrictions on RBM structure yields properties

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

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

- Individual conditionals are simple to compute

- For binary RBM we obtain

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

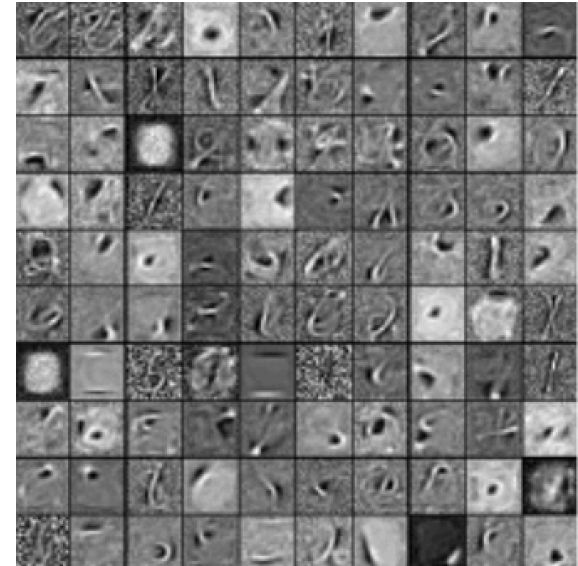
$$p(h_i=0|\mathbf{v}) = 1 - \sigma(\mathbf{v}^T \mathbf{W}_{:,i} + b_i)$$

- Together these properties allow for *block Gibbs sampling* which alternate between sampling all \mathbf{h} simultaneously and all \mathbf{v} simultaneously

- Shown next

Samples from a trained RBM

- Samples from a trained RBM on MNIST drawn using Gibbs sampling



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

Derivatives of Energy Function

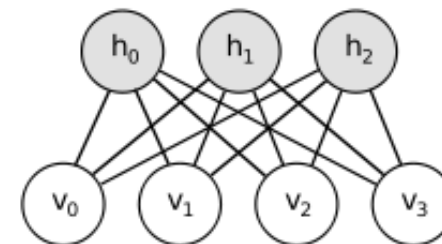
- Energy function: $E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^T \mathbf{v} - \mathbf{c}^T \mathbf{h} - \mathbf{v}^T \mathbf{W} \mathbf{h}$
 - where \mathbf{b} , \mathbf{c} and \mathbf{W} are unconstrained, real-valued learnable parameters
- Since the energy function is a linear function of its parameters, 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
 - Undirected models can be trained by computing such derivatives applied to samples from the model

Energy function of a RBM

- 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



- Defining free energy as

$$\mathcal{F}(\mathbf{v}) = -\mathbf{b}'\mathbf{v} - \sum_i \log \sum_{h_i} e^{h_i(\mathbf{c}_i + \mathbf{W}_i \mathbf{v})}$$

- Due to structure of RBM

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

RBM with binary units

- Using $v_j, h_i \in \{0,1\}$

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

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

- Free energy simplifies to

$$\mathcal{F}(\mathbf{v}) = -b'v - \sum_i \log(1 + e^{(c_i + W_i v)}).$$

- Update equations

$$-\frac{\partial \log p(\mathbf{v})}{\partial W_{ij}} = E_v[p(h_i|\mathbf{v}) \cdot v_j] - v_j^{(i)} \cdot \text{sigm}(W_i \cdot \mathbf{v}^{(i)} + c_i)$$

$$-\frac{\partial \log p(\mathbf{v})}{\partial c_i} = E_v[p(h_i|\mathbf{v})] - \text{sigm}(W_i \cdot \mathbf{v}^{(i)})$$

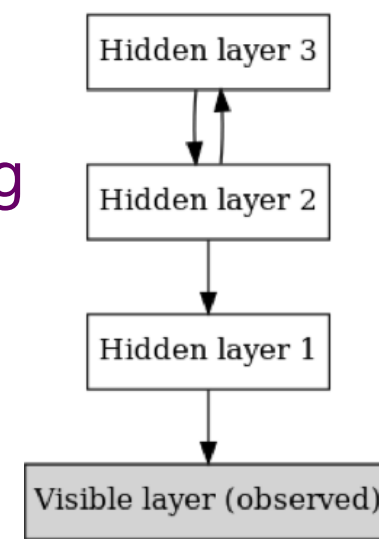
$$-\frac{\partial \log p(\mathbf{v})}{\partial b_j} = E_v[p(v_j|h)] - v_j^{(i)}$$

Training RBMs

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

Deep Belief Networks (DBNs)

- Consist of several layers of RBMs
 - Stacking RBMs
 - Fine tuning resulting deep network using gradient descent and back-propagation
- DBNs are Generative Models
 - Provide estimates of both
$$p(x | C_k) \text{ and } p(C_k | x)$$
 - Conventional neural networks are discriminative
 - Directly estimate $p(C_k | x)$



Training several layers of RBMs

- Let X be a matrix of input feature vectors
 1. Train an RBM on X to obtain weight matrix W
 - Between lower two layers (input and hidden)
 2. Transform X by RBM to produce new data X'
 - by sampling or by computing mean activation of hidden units
 3. Repeat procedure with $X \leftarrow X'$ for next layer pair
 - Until top two layers of network are reached (output and hidden)