

Log-linear MRFs: Ising, Boltzmann, Deep Belief, Metric

Sargur Srihari

srihari@cedar.buffalo.edu

Topics

- Log-linear MRF Applications
 - Ising Model
 - Boltzmann Distribution
 - Energy Based Model
 - Boltzmann Machine
 - Restricted Boltzmann Machine
 - Deep Belief Networks
 - Metric MRF

General Log-linear model with features

- A distribution P is a log-linear model over \mathcal{H} if

$$P(X_1, \dots, X_n) = \frac{1}{Z} \exp \left[- \sum_{i=1}^k w_i f_i(D_i) \right]$$

Note that k is the
no of features
Not no of subgraphs

- Can have several functions over same scope
- Each term is an energy function
- Equivalent to Gibbs distribution

$$P_\phi(X_1, \dots, X_n) = \frac{1}{Z} \tilde{P}(X_1, \dots, X_n) \text{ where } \tilde{P}(X_1, \dots, X_n) = \prod_{i=1}^m \phi_i(D_i)$$

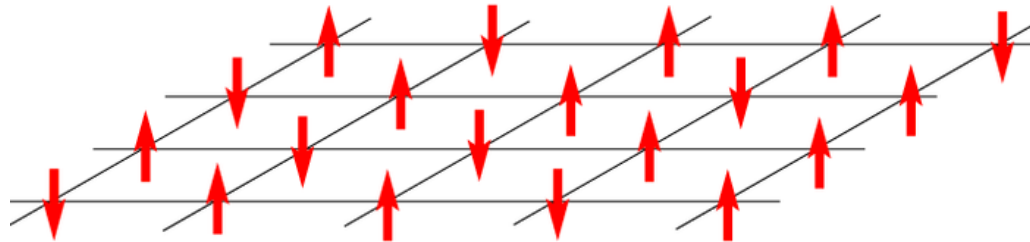
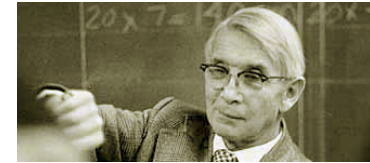
is an unnormalized measure and $Z = \sum_{X_1, \dots, X_n} \tilde{P}(X_1, \dots, X_n)$

- Rewrite factor $\phi(D)$ as $\phi(D) = \exp(-\varepsilon(D))$

Where $\varepsilon(D) = -\ln \phi(D)$ is the *energy* function

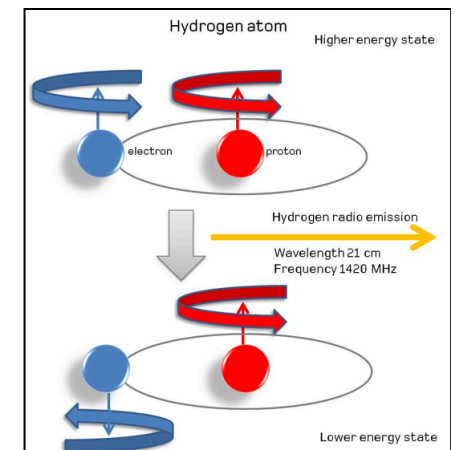
Ising Model

- Earliest Markov Network model
- Model of interacting atoms
 - Square lattice model



- Energy determined by magnetic spins of atoms
 - Atom's spin is sum of its electron spins

One atom going from higher to lower energy releases radio emission and changes electron spin



Ising Model

- Each atom associated with binary random variable
 - $X_i \in \{+1, -1\}$ whose value is direction of atom's spin
- Energy function parametric form
$$\varepsilon_{ij}(x_i, x_j) = -w_{ij}x_ix_j$$
- Symmetric in X_i, X_j : note scope is pairwise
 - Makes contribution w_{ij} to energy when $X_i = X_j$ (same spin)
 - $-w_{ij}$ otherwise

Ising Model as a Markov Network

- Features are pairwise and single
 - Pair Energy function between X_i, X_j
 - Parametric form $\varepsilon_{ij}(x_i, x_j) = -w_{ij}x_ix_j$
 - Contributes w_{ij} when $X_i = X_j$, same spin, and $-w_{ij}$ otherwise
 - Node potentials are u_i which bias atom spin x_i
- Probability distribution over atoms (energy function)

$$P(\xi) = \frac{1}{Z} \exp\left(-\sum_{i < j} w_{ij} x_i x_j - \sum_i u_i x_i\right)$$

$\xi \in \text{Val}(\mathbf{X})$ is a full assignment of the variables

 - When $w_{ij} > 0$ model prefers aligned spins: ferromagnetism
 - $w_{ij} < 0$: antiferromagnetic
 - $w_{ij} = 0$: non-interacting

Ising Model studies

- To answer a variety of questions
 - Usually as the no. of atoms (variables) goes to infinity

$$P(\xi) = \frac{1}{Z} \exp \left(- \sum_{i < j} w_{ij} x_i x_j - \sum_i u_i x_i \right)$$

- Inference problems, e.g.,
 - Determine probability of a configuration where majority of spins are +1 (or -1) versus more mixed ones
 - Answer depends on strength of interactions w_{ij}
 - e.g., Multiply all weights by temperature parameter
 - Many other problems investigated extensively
 - Answers known--some even analytically



Boltzmann Distribution

- Variant of Ising Model
- Variables X_i have value $\{0, 1\}$ instead of $\{+1, -1\}$
 - Energy function has same parametric form

$$\varepsilon_{ij}(x_i, x_j) = -w_{ij}x_ix_j$$
 - Nonzero contribution $-w_{ij}$ from edge X_i-X_j only when $X_i=X_j=1$
 - Ising model has contribution w_{ij} when variables are same and $-w_{ij}$ when they are different
- Has the same energy function as Ising model

$$P(\xi) = \frac{1}{Z} \exp \left(- \sum_{i < j} w_{ij} x_i x_j - \sum_i u_i x_i \right)$$

Mapping 0 to -1

Boltzmann Distrib. & Statistical Mechanics

- Boltzmann Probability distribution

$$P(\text{state}) \propto \exp[-E / kT]$$

Where E is state energy (varies from state to state)

- kT is a constant of the distribution
 - k = Boltzmann's constant, T = absolute temperature
- Ratio over two states depends on energy difference

$$P(\text{state1}) / P(\text{state2}) = \exp[E2 - E1] / kT$$

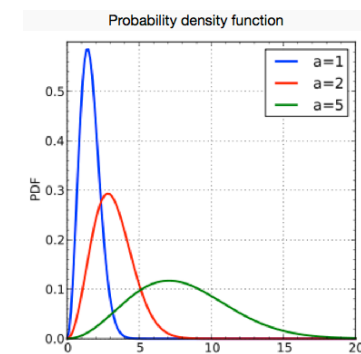
- Later investigated by Josiah Gibbs

- Boltzmann distribution also known as Gibbs measure

- Maxwell-Boltzmann distribution

- Is χ^2 with 3 degrees of freedom

$$f(v) = \sqrt{\left(\frac{m}{2\pi kT}\right)^3} 4\pi v^2 e^{-\frac{mv^2}{2kT}}$$



Boltzmann Distribution with sigmoid

- Probability distribution of each variable X_i given assignment of neighbors X_j is *sigmoid* (z) where

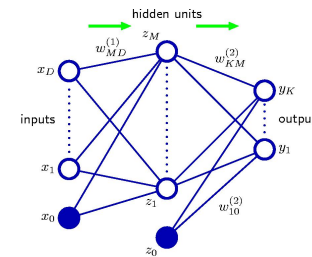
$$z = - \left(\sum_j w_{ij} x_j \right) - w_i$$

- *sigmoid* (z) is a weighted combination of X_i 's neighbors, weighted by strength and direction of connection
 - Where *sigmoid*(z) = $[1/1 + \exp(-z)]$ is a value in $[0, 1]$
- Simplest mathematical approximation of the function employed by a neuron in the brain

Boltzmann Distribution & Neuron

- Boltzmann distribution resembles a neuron
 - Output node of an artificial neural network

$$y_k(x, w) = \sigma \left(\sum_{j=1}^M w_{kj}^{(2)} h \left(\sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right)$$



- Neuron output is a stochastic function of its connected neighbors
- Boltzmann distribution is a type of energy model

Energy-based Model (EBM)

- Probability distribution: associates a scalar energy with each configuration of its variables
- Learning corresponds to modifying energy function so its shape has desirable properties
 - E.g., plausible configurations have low energy
- Energy-based probability distribution

$$p(x) = \frac{1}{Z} \exp(-E(x))$$

- Where Z is the partition function

$$Z = \sum_x \exp(-E(x))$$

Learning the energy-based model

- To determine parameters of

$$p(x) = \frac{1}{Z} \exp(-E(x))$$

- Perform stochastic gradient-descent on negative log-likelihood

- Log-likelihood $\mathcal{L}(\theta, \mathcal{D}) = \frac{1}{N} \sum_{x^{(i)} \in \mathcal{D}} \log p(x^{(i)})$

- Loss function $\ell(\theta, \mathcal{D}) = -\mathcal{L}(\theta, \mathcal{D})$

– Gradient is $-\frac{\partial \log p(x^{(i)})}{\partial \theta}$ where θ are parameters

$$\theta^{(\tau+1)} = \theta^{(\tau)} - \eta \nabla \ell$$

EBMs with hidden units

- Want to include non-observed variables to increase expressive power of model

$$P(x) = \sum_h P(x, h) = \sum_h \frac{e^{-E(x, h)}}{Z}$$

- Introducing free-energy $\mathcal{F}(x) = -\log \sum_h e^{-E(x, h)}$

$$P(x) = \frac{e^{-\mathcal{F}(x)}}{Z} \text{ with } Z = \sum_x e^{-\mathcal{F}(x)}$$

- Data negative log-likelihood gradient

$$-\frac{\partial \log p(x)}{\partial \theta} = \frac{\partial \mathcal{F}(x)}{\partial \theta} - \sum_{\tilde{x}} p(\tilde{x}) \frac{\partial \mathcal{F}(\tilde{x})}{\partial \theta}$$

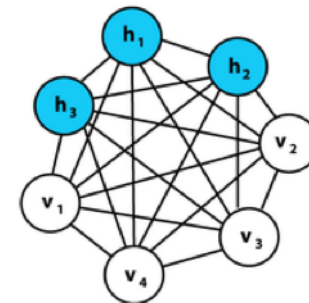
First term increases probability of training data.
Second term decreases probability of samples generated by model

- Sampling version (with samples from P)

$$-\frac{\partial \log p(x)}{\partial \theta} \approx \frac{\partial \mathcal{F}(x)}{\partial \theta} - \frac{1}{|\mathcal{N}|} \sum_{\tilde{x} \in \mathcal{N}} \frac{\partial \mathcal{F}(\tilde{x})}{\partial \theta}$$

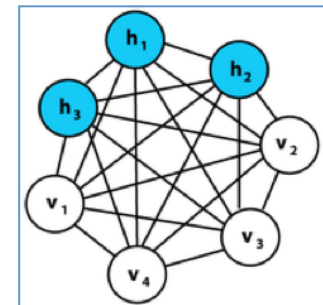
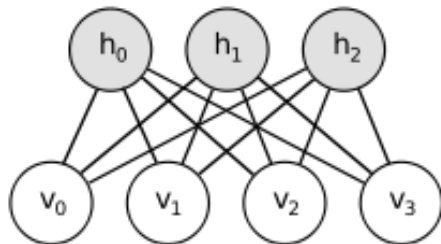
Boltzmann Machine

- A form of Energy Based Model
- Structure of a recurrent neural network (RNN):
 - one where there are directed cycles
 - Unlike feed-forward neural networks RNN can use internal memory to process arbitrary sequences
 - Can process time-varying real-valued inputs
 - Have nodes which are inputs, hidden and outputs
- Boltzmann machines are a type of RNN



Restricted Boltzmann Machine

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



Not an RBM

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

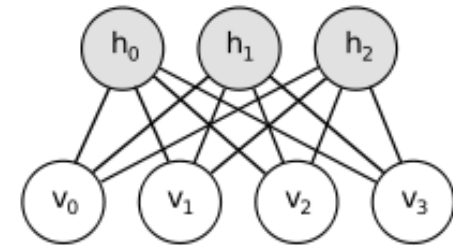
Energy function of a RBM

- Energy function

$$E(v, h) = -b'v - c'h - h'Wv$$

where W is weight matrix connecting hidden and visible units

$v = [v_0, v_1, \dots], h = [h_0, h_1, \dots]$, with offset vectors b, c



- Defining free energy as

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

- Due to structure of RBM

$$p(h|v) = \prod_i p(h_i|v)$$

$$p(v|h) = \prod_j p(v_j|h)$$

RBM with binary units

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

$$P(h_i = 1|v) = \text{sigm}(c_i + W_i v)$$

$$P(v_j = 1|h) = \text{sigm}(b_j + W'_j h)$$

- Free energy simplifies to

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

- Update equations

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

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

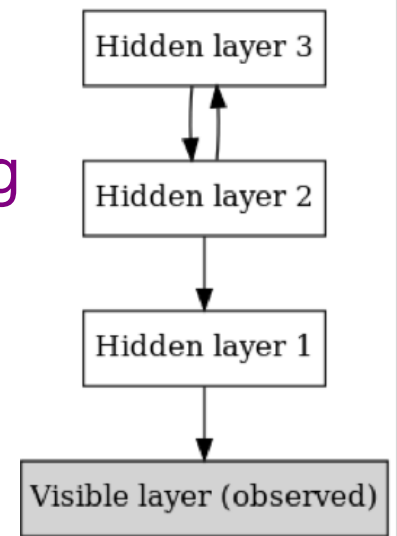
$$-\frac{\partial \log p(v)}{\partial b_j} = E_v[p(v_j|h)] - v_j^{(i)}$$

Training RBMs

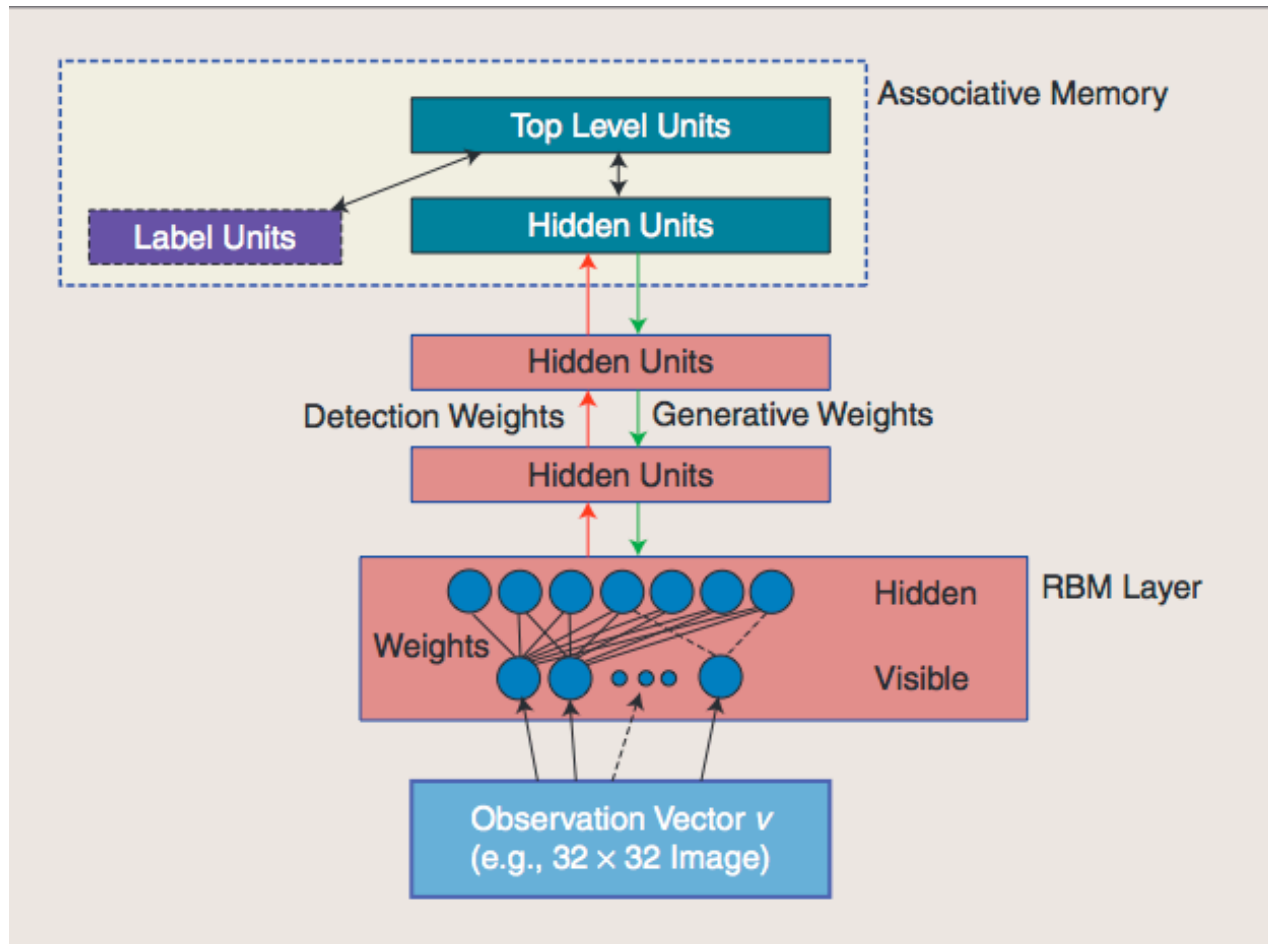
- Contrastive Divergence
- A method to overcome exponential complexity in dealing with the partition function

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 $p(x/C_k)$ and $p(C_k/x)$
 - Conventional neural networks are discriminative
 - Directly estimate $p(C_k/x)$



Deep Belief Network Framework



Training DBNs

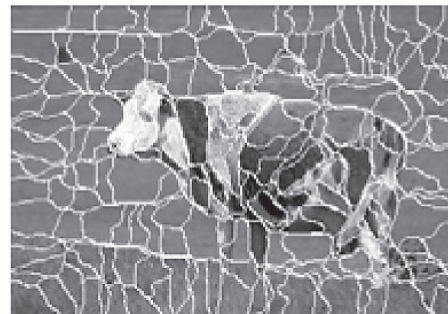
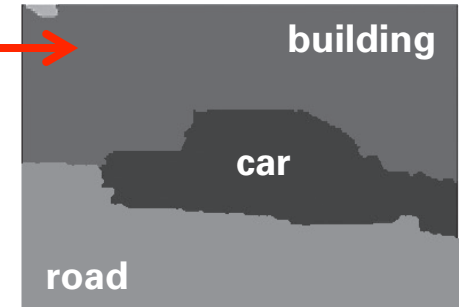
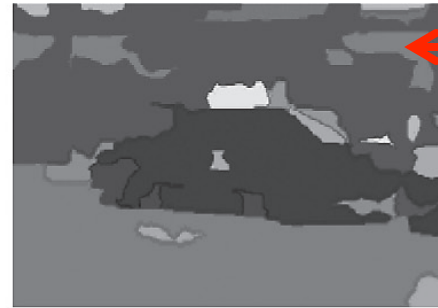
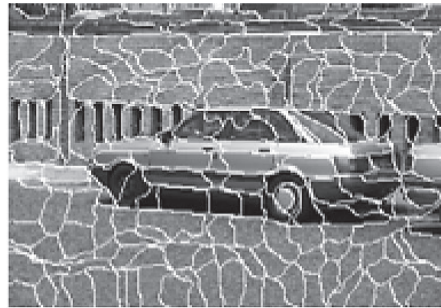
- Let X be a matrix of input feature vectors
 1. Train an RBM on X to obtain weight matrix W
 - Between lower two layers (input and hidden)
 2. Transform X by RBM to produce new data X'
 - by sampling or by computing mean activation of hidden units
 3. Repeat procedure with $X \leftarrow X'$ for next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

Metric MRF for Labeling

- Task:
 - Graph with nodes X_1, \dots, X_n , edges E
 - Assign to each X_i a label in $V = \{v_1, \dots, v_k\}$
 - E.g., labeling super-pixels in image
 - Each node, in isolation, has a preferred label
 - E.g., color specifies a label
 - However, we want *smoothness* constraint over neighbors
 - Neighboring nodes should have “similar” values

Importance of Modeling

Correlations between superpixels



Original image

Oversegmented
image-superpixels
Each superpixel is
a random variable

Classification using
node potentials
alone-each
superpixel classified
independently

Segmentation using
pairwise Markov
Network encoding
interactions
between adjacent
superpixels

Solution for Labeling

- Solution:
 - Encode node preferences as edge potentials
 - Smoothness preferences as edge potentials
 - Encode model in negative log-space, using energy functions

- Energy function

$$E(x_1, \dots, x_n) = \sum_i \varepsilon_i(x_i) + \sum_{i,j \in E} \varepsilon_{ij}(x_i, x_j)$$

- For MAP objective, ignore partition function
- Goal: Minimize the energy (MAP objective)

$$\arg \min_{x_1, \dots, x_n} E(x_1, \dots, x_n)$$

- How to define smoothness? Next.

Smoothness for Metric MRF

- Many variants
- Simplest one is a variant of Ising model

$$\mathcal{E}_{i,j}(x_i, x_j) = \begin{cases} 0 & x_i = x_j \\ \lambda_{i,j} & x_i \neq x_j \end{cases}$$

- for $\lambda_{ij} \geq 0$

- In this model:
 - Lowest pairwise energy (0) when neighbors have same value
 - Higher energy otherwise λ_{ij}

Generalizations of Smoothness for Metric MRF

1. Potts model (when there are more than two labels)
2. Distance Function on labels
 - Prefer neighboring nodes to have labels smaller distance apart
 - Metric MRF
 - Need a metric $\mu(v_k, v_l)$ on labels

Metric Requirement

- Function $\mu: V \times V \rightarrow [0, \infty)$
 - Reflexivity, symmetry and triangle inequality
 - Semi-metric if triangle inequality is violated
- Metric MRF
 - Define $\varepsilon_{i,j}(v_k, v_l) = \mu(v_k, v_l)$
 - Where μ is a metric (or semi-metric)
 - Assume same for all variables
 - » Simplifies no. of parameters needed
 - » Usually holds in practice
 - Example metric p -norm: $\varepsilon(x_i, x_j) = \min(c \|x_i - x_j\|_p, \text{dist}_{\max})$
- Metric interactions arise frequently
 - Plays important role in computer vision