

# Learning the Parameters of Markov Networks

srihari@buffalo.edu

# Topics

- Generative vs Discriminative Training
- Log-likelihood for Markov Network
- Iterative Method for MN parameters
- Training Energy-based models
- Learning for RBMs

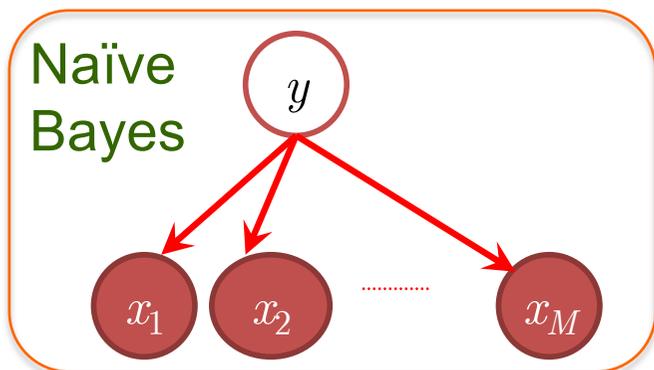
# Overview

- Both BNs and MNs are generative models
- Training can be generative or discriminative
  - Generative: determines joint  $P(\chi)$
  - Discriminative: a conditional  $P(y|\chi)$
- We will look at both discriminative and generative training of MNs

# BN Parameter Learning

Variables  $\mathbf{x} = \{x_1, \dots, x_M\}$  and classifier target  $y \in \{0, 1\}$

- Estimate parameters of variables independently



For classification: Determine joint  $p(y, \mathbf{x})$

$$p(y, \mathbf{x}) = p(y) \prod_{i=1}^M p(x_i | y)$$

From joint get required conditional  $p(y|\mathbf{x})$

## Parameter estimation

independently estimate  $M$  sets of parameters,  $p(x_i | y)$ ,  $i=1, \dots, M$

## But independence is usually false

We can estimate  $M(M+1)/2$  covariance matrix

# MN Parameter Discriminative Learning

Variables  $\mathbf{x} = \{x_1, \dots, x_M\}$  and classifier target  $y \in \{0, 1\}$

Estimate joint parameters  $w_i$

## Joint Distribution

Potential Functions (log-linear)

$$\phi_i(x_i, y) = \exp\{w_i x_i \mathbb{I}\{y=1\}\},$$

$$\phi_0(y) = \exp\{w_0 \mathbb{I}\{y=1\}\}$$

Naïve  
Markov



It has value 1 when  $y=1$ , else 0



## Conditional Distribution

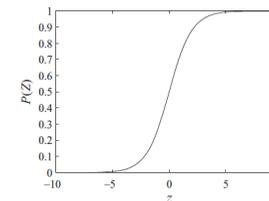
We can directly write the conditional probabilities as unnormalized probs:

$$\tilde{P}(y=1|\mathbf{x}) = \exp\left\{w_0 + \sum_{i=1}^M w_i x_i\right\} \quad \tilde{P}(y=0|\mathbf{x}) = \exp\{0\} = 1$$

Normalized

$$P(y=1|\mathbf{x}) = \text{sigmoid}\left\{w_0 + \sum_{i=1}^M w_i x_i\right\} \quad \text{where } \text{sigmoid}(z) = \frac{e^z}{1+e^z}$$

This is Logistic Regression



Jointly optimize  $M$  parameters

More complex estimation but correlations accounted for

Can use much richer features:

Edges, image patches sharing same pixels

**multiclass**

$$p(y_i|\phi) = y_i(\phi) = \frac{\exp(a_i)}{\sum_j \exp(a_j)}$$

where  $a_j = \mathbf{w}_j^T \phi$

# Logistic Regression Parameter Learning

- For a data set  $(\phi_n, t_n)$  where  $t_n \in \{0,1\}$  and  $\phi_n = \phi(\mathbf{x}_n)$ ,  $n = 1, \dots, N$
- Likelihood function is

$$p(t | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

- Taking negative logarithm we get *Cross-entropy Error*

$$E(\mathbf{w}) = -\ln p(t | \mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

where  $y_n = \sigma(a_n)$  and  $a_n = \mathbf{w}^T \phi_n$

- We need to minimize  $E(\mathbf{w})$

At its minimum, derivative of  $E(\mathbf{w})$  is zero

So we need to solve for  $\mathbf{w}$  in the equation

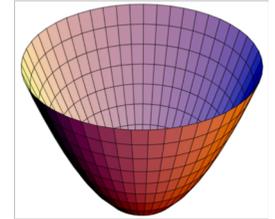
$$\nabla E(\mathbf{w}) = 0$$

# Gradient Descent for Logistic Reg.

- Given Gradient of error function

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$$

where  $y_n = \sigma(\mathbf{w}^T \phi_n)$



- Solve using an iterative approach

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E_n$$

Takes precisely same form as  
Gradient of Sum-of-squares  
error for linear regression

- where  $\nabla E_n = (y_n - t_n) \phi_n$



Error x Feature Vector

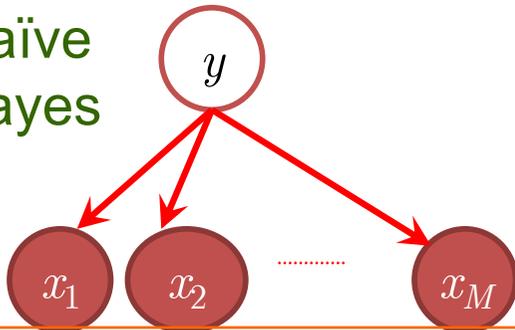
Samples are presented one at a time in  
which each each of the weight vectors is updated

# BN vs MN Training

Variables  $\mathbf{x} = \{x_1, \dots, x_M\}$  and classifier target  $y$

## 1. BN generative: estimate parameters of variables independently

Naïve Bayes



For classification:  
Determine joint:

$$p(y, \mathbf{x}) = p(y) \prod_{i=1}^M p(x_i | y)$$

From joint get required conditional  $p(y | \mathbf{x})$

Simple estimation

independently estimate  $M$  sets of parameters  
But independence is usually false  
We can estimate  $M(M+1)/2$  covariance matrix

## 2. MN: generative/discriminative: estimate joint params $w_i$

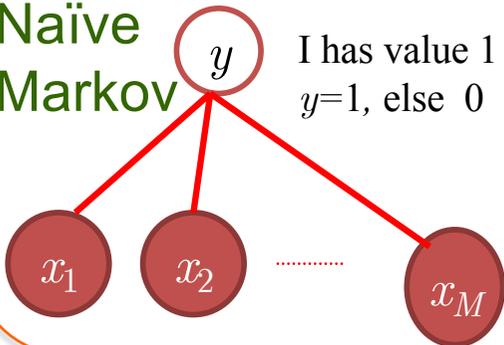
Generative Model (log-linear)

$$\phi_i(x_i, y) = \exp\{w_i x_i I\{y=1\}\},$$

$$\phi_0(y) = \exp\{w_0 I\{y=1\}\}$$

Naïve Markov

I has value 1 when  $y=1$ , else 0



Discriminative Model:

Unnormalized

$$\tilde{P}(y=1 | \mathbf{x}) = \exp\left\{w_0 + \sum_{i=1}^M w_i x_i\right\}$$

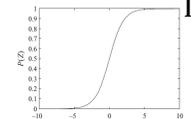
$$\tilde{P}(y=0 | \mathbf{x}) = \exp\{0\} = 1$$

Normalized

$$P(y=1 | \mathbf{x}) = \text{sigmoid}\left\{w_0 + \sum_{i=1}^M w_i x_i\right\}$$

where  $\text{sigmoid}(z) = \frac{e^z}{1 + e^z}$

Logistic Regression



Jointly optimize  $M$  parameters

More complex estimation but correlations accounted for

Can use much richer features:

Edges, image patches sharing same pixels

**multiclass**

$$p(y_i | \phi) = y_i(\phi) = \frac{\exp(a_i)}{\sum_j \exp(a_j)}$$

where  $a_j = \mathbf{w}_j^T \phi$

# Likelihood Function

- Basis for all discussion of learning
- How likelihood function can be optimized to find maximum likelihood parameter estimates
- Begin with form of likelihood function for Markov networks, its properties and their computational implications
- Existence of a global partition function couples the different parameters in a Markov network
  - greatly complicating the estimation problem

# Generative training of MN (Gibbs)

- Markov network  $A—B—C$

- Parameterized by potentials  $\phi_1(A, B)$ ,  $\phi_2(B, C)$ :

- Recall Gibbs:

$$P(a, b, c) = (1/Z) \phi_1(a, b) \phi_2(b, c) \text{ where}$$

- Log-likelihood of instance  $(a, b, c)$  is  $Z = \sum_{a, b, c} \phi_1(a, b) \phi_2(b, c)$

$$\ln P(a, b, c) = \ln \phi_1(a, b) + \ln \phi_2(b, c) - \ln Z$$

$\phi_1[A, B]$		
$a^0$	$b^0$	30
$a^0$	$b^1$	5
$a^1$	$b^0$	1
$a^1$	$b^1$	10

$\phi_2[B, C]$		
$b^0$	$c^0$	100
$b^0$	$c^1$	1
$b^1$	$c^0$	1
$b^1$	$c^1$	100

- Log-likelihood of data set  $\mathcal{D}$  with  $M$  instances:

$$\begin{aligned} \ell(\theta : D) &= \sum_{m=1}^M (\ln \phi_1(a[m], b[m]) + \ln \phi_2(b[m], c[m]) - \ln Z(\theta)) \\ &= \sum_{a, b} M[a, b] \ln \phi_1(a, b) + \sum_{b, c} M[b, c] \ln \phi_2(b, c) - M \ln Z(\theta) \end{aligned}$$

$$M[a, b] = \text{no. with value } (a, b)$$

Parameter  $\theta$  consists of all values of factors  $\phi_1$  and  $\phi_2$

Summing over instances  $\{a[m], b[m], c[m]\}, m=1, \dots, M$

Then summing over different values of  $a \in \{a^0, a^1\}$  and  $b \in \{b^0, b^1\}$

- This likelihood has three terms which we analyze

# Coupled Likelihood of Simple MN

- Simple Markov network  $A—B—C$
- Log-likelihood of data set  $\mathcal{D}$  with  $M$  instances:

$$\begin{aligned}\ell(\theta : \mathcal{D}) &= \sum_{m=1}^M (\ln \phi_1(a[m], b[m]) + \ln \phi_2(b[m], c[m]) - \ln Z(\theta)) \\ &= \sum_{a,b} M[a,b] \ln \phi_1(a,b) + \sum_{b,c} M[b,c] \ln \phi_2(b,c) - M \ln Z(\theta)\end{aligned}$$

$M[a,b]$  = no. with value  $(a,b)$

- This likelihood consists of three terms.

- First term involves only  $\phi_1$
- Second only  $\phi_2$ . But third involves

$$\ln Z(\theta) = \ln \left( \sum_{a,b,c} \phi_1(a,b) \phi_2(b,c) \right)$$

- Thus  $\ln Z(\theta)$  is a function of both  $\phi_1$  and  $\phi_2$ . It couples the two potentials in the likelihood function
- When we change one potential  $\phi_1$ ,  $Z(\theta)$  changes, possibly changing the value of  $\phi_2$  that maximizes  $-\ln Z(\theta)$

# Log-Likelihood Surface for $A-B-C$

- Two factors  $\phi_1(A, B)$ ,  $\phi_2(B, C)$ 
  - With binary variables each factor has four values
  - Thus, total of 8 parameters

$\phi_1[A, B]$		
$a^0$	$b^0$	30
$a^0$	$b^1$	5
$a^1$	$b^0$	1
$a^1$	$b^1$	10

$\phi_2[B, C]$		
$b^0$	$c^0$	100
$b^0$	$c^1$	1
$b^1$	$c^0$	1
$b^1$	$c^1$	100

- We need parameters (values of  $\phi_i$ ) that maximize

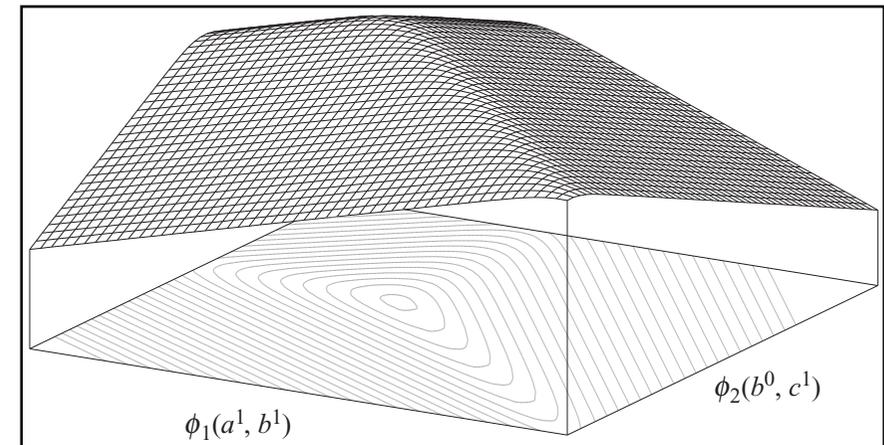
$$\ell(\theta : D) = \sum_{a,b} M[a,b] \ln \phi_1(a,b) + \sum_{b,c} M[b,c] \ln \phi_2(b,c) - M \ln Z(\theta)$$

- Ex: Log-likelihood surface

- X-axis:  $\phi_1(a^1, b^1)$
- Y-axis:  $\phi_2(b^0, c^1)$
- Other parameters set to 1
- Data set has  $M=100$  with  $M[a^1, b^1]=40$ ,  $M[b^0, c^1]=40$

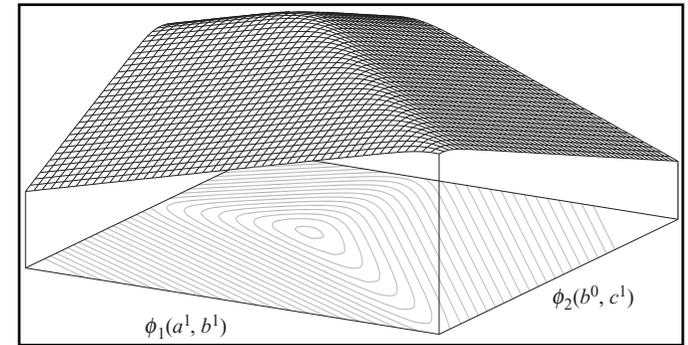
- Coupling problem:

- When  $\phi_1$  changes,  $\phi_2$  that maximizes  $-\ln Z(\theta)$  also changes



# Equivalent Bayesian Network

- Log-likelihood  $A—B—C$  wrt two factors  $\phi_1(A,B)$ ,  $\phi_2(B,C)$ 
  - With binary variables we would have 8 parameters
- In this case, parameter coupling is avoidable
  - Since  $A—B—C$  has an equivalent BN:  
 $A \rightarrow B \rightarrow C$
  - We can estimate parameters using
$$\phi_1(A,B) = P(A)P(B|A)$$
$$\phi_2(B,C) = P(C|B)$$
- In general, cannot convert learned BN parameters into equivalent MN
  - Optimal likelihood achievable by the two representations is not the same



# Log-linear form for MN

- Instead of Gibbs, use log-linear framework

- Joint distribution of  $n$  variables  $X_1, \dots, X_n$
- $k$  features  $\mathcal{F} = \{ f_i(D_i) \}_{i=1, \dots, k}$   $k$  depends on no. of values  $D_i$  takes
  - where  $D_i$  is a sub-graph and  $f_i$  maps  $D_i$  to  $\mathcal{R}$

$$P(X_1, \dots, X_n; \theta) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{i=1}^k \theta_i f_i(D_i) \right\}$$

- Parameters  $\theta_i$  are weights we put on features
  - If we have a sample  $\xi$  then its features are  $f_i(\xi(D_i))$  which has the shorthand  $f_i(\xi)$ .
- Representation is general
  - can capture Markov networks with global structure and local structure

# Generative training of MN (log-linear)

- With log-linear form for probability distribution:

$$P(X_1, \dots, X_n; \theta) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{i=1}^k \theta_i f_i(D_i) \right\}$$

$\theta = \{\theta_1, \dots, \theta_k\}$  are table entries  
 $f_i$  are features over instances of  $D_i$

- Let  $\mathcal{D}$  be a data set of  $M$  samples  $\xi [m]_{m=1, \dots, M}$
- Log-likelihood (product of indep probabilities)  
– log gets rid of exp, converts product to sum:

$$\ell(\theta : D) = \sum_i \theta_i \left( \sum_m f_i(\xi[m]) \right) - M \ln Z(\theta)$$

- Sufficient statistics (likelihood depends only on this)

Dividing  
by no.  
of samples

$$\frac{1}{M} \ell(\theta : D) = \sum_i \theta_i \left( E_D [f_i(d_i)] \right) - \ln Z(\theta)$$

$E_{\mathcal{D}}[f_i(d_i)]$  is empirical  
expectation: average  
in the data set

# Maximum Likelihood Parameter Estimation

- Task: Estimate parameters of a MN with a fixed structure given a fully observable data set  $\mathcal{D}$
- Simplest variant of the problem is maximum likelihood parameter estimation
- Log-likelihood with features  $\mathcal{F} = \{f_i, i = 1, \dots, k\}$  is

$$\ell(\theta : \mathcal{D}) = \sum_{i=1} \theta_i \left( \sum_m f_i(\xi[m]) \right) - M \ln Z(\theta)$$

- Although log-likelihood function
  - is concave, no analytical form for its maximum
    - Can use iterative methods, e.g. gradient ascent

# Gradient of Log-likelihood

- Log-likelihood:  $\ell(\theta : D) = \sum_{i=1} \theta_i \left( \sum_m f_i(\xi[m]) \right) - M \ln Z(\theta)$
- Average log-likelihood:  $\frac{1}{M} \ell(\theta : D) = \sum_i \theta_i (E_D[f_i(d_i)]) - \ln Z(\theta)$  Dividing by no. of samples  $M$
- We have seen that gradient of second term  $\ln Z(\theta)$  is
 
$$\frac{\partial}{\partial \theta_i} \ln Z(\theta) = \frac{1}{Z(\theta)} \sum_{\xi} \frac{\partial}{\partial \theta_i} \exp \left\{ \sum_j \theta_j f_j(\xi) \right\} = \frac{1}{Z(\theta)} \sum_{\xi} f_i(\xi) \exp \left\{ \sum_j \theta_j f_j(\xi) \right\} = E_{\theta}[f_i]$$
- We can compute gradient of average log-likelihood as
 
$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : D) = E_D[f_i(\chi)] - E_{\theta}[f_i]$$
First term is average value of  $f_i$  in data  $\mathcal{D}$ . Second term is expected value from distribution
  - Provides a precise characterization of m.l. parameters  $\theta$
- *Theorem:* Let  $\mathcal{F}$  be a feature set. Then  $\theta$  is a m.l. assignment if and only if  $E_D[f_i(\chi)] = E_{\hat{\theta}}[f_i]$  for all  $i$ 
  - i.e., expected value of each feature relative to  $P_{\theta}$  matches its empirical expectation in  $\mathcal{D}$

# Need for Iterative Method

- Although log-likelihood function

$$\ell(\theta : D) = \sum_i \theta_i \left( \sum_m f_i(\xi[m]) \right) - M \ln Z(\theta)$$

- is concave, there is no analytical form for the maximum
- Since no closed-form solution
  - Can use iterative methods, e.g. gradient ascent as shown next
- Fortunately, exact form of gradient is known:

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : D) = E_D[f_i(\chi)] - E_\theta[f_i]$$

# Computing the gradient

- Gradient wrt  $\theta_i$  of log-likelihood  $\ell(\theta; \mathcal{D})$  is

$$E_{\mathcal{D}}[f_i(\chi)] - E_{\theta}[f_i]$$

– It is the difference between the feature's

1. empirical expectation in data  $\mathcal{D}$  and
2. Its expected value relative to current parameterization

- Computing empirical expectation  $E_{\mathcal{D}}[f_i(\chi)]$  is easy

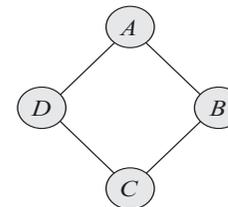
– Ex: for feature

$$f_{a^0 b^0}(a, b) = I\{a = a^0\} I\{b = b^0\}$$

it is the empirical frequency in  $\mathcal{D}$  of the event  $a^0, b^0$

- At a particular parameterization  $\theta$ , the expected  $f_i$  is simply  $P_{\theta}(a^0, b^0)$  since  $E_{\theta}[f_i] = \frac{1}{Z(\theta)} \sum_{\xi} f_i(\xi) \exp\left\{\sum_j \theta_j f_j(\xi)\right\}$

– Note: for Bernoulli( $x|p$ ), expectation is  $p$



All four binary-valued  $a \in \{a^0, a^1\}$  etc.  
Features are indicator fns

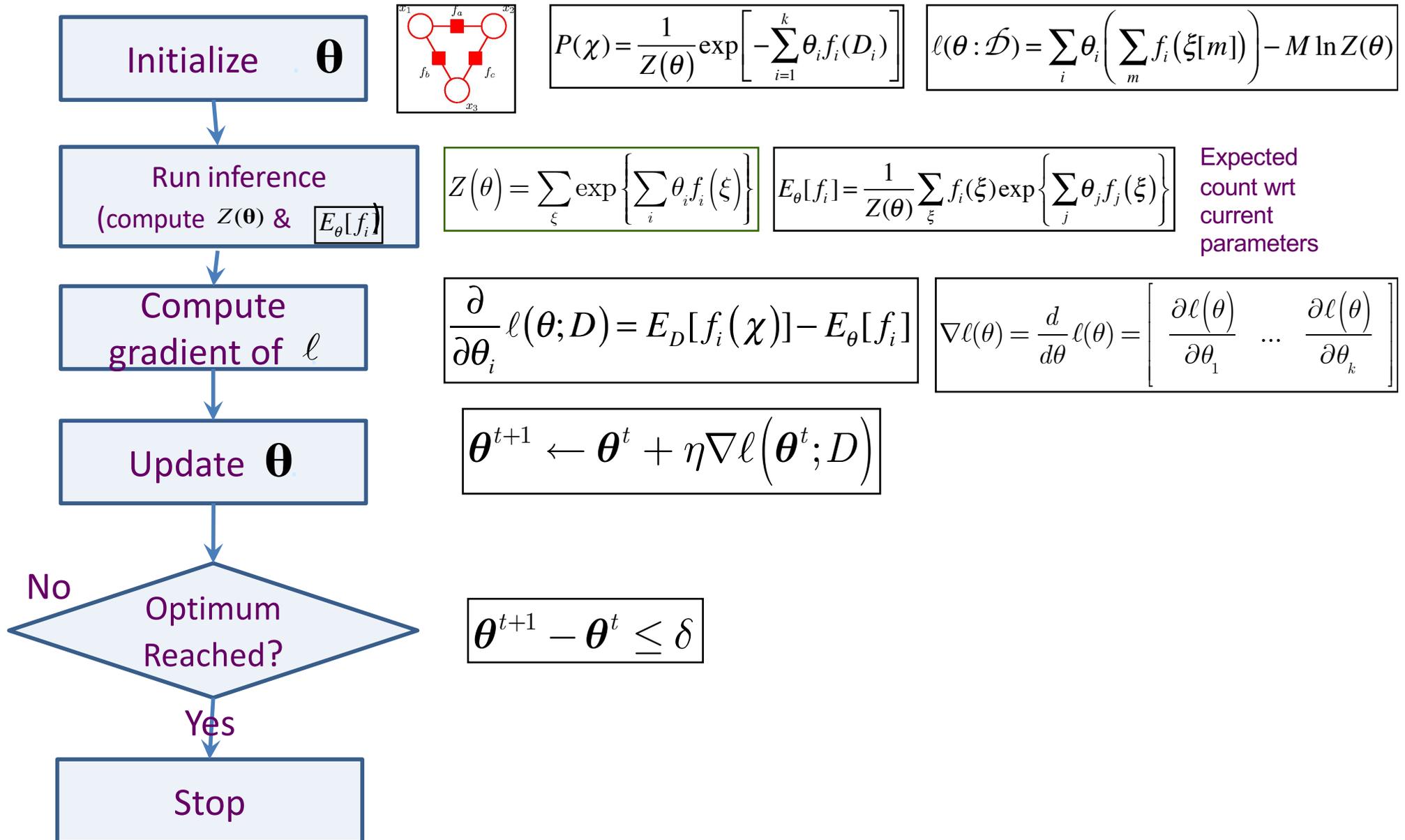
# Computing the expected count

- Second term in the gradient is

$$E_{\theta}[f_i] = \frac{1}{Z(\theta)} \sum_{\xi} f_i(\xi) \exp \left\{ \sum_j \theta_j f_j(\xi) \right\}$$

- We need to compute the different probabilities of the form  $P_{\theta^t}(a, b)$ 
  - Since expectation is a probability-weighted average
- Computing probability requires running inference over the network

# Iterative solution for MN parameters $\theta$



# Difficulty with Iterative method

- Gradient ascent over parameter space
- Good news:
  - likelihood function is concave
    - Guaranteed to converge to global optimum
- Bad news:
  - each step needs inference
  - Simple parameter estimation is intractable
  - Bayesian parameter estimation even harder
    - Integration done using MCMC

# Energy-based Models (EBMs)

- Boltzmann distribution is an energy model
  - Probability distribution: associates a scalar energy with each configuration of its variables
- 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 corresponds to modifying energy function so its shape has desirable properties
  - E.g., plausible configurations have low energy

# Learning EBM parameters

- 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  $\theta$  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  $P$ )

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

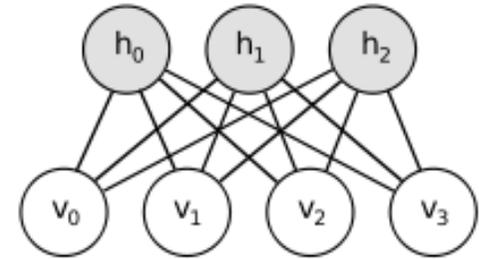
# Learning with RBMs

- 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

# Summary of MN Param Learning

- MN parameter learning can be performed generatively
- Maximize log-likelihood using an iterative algorithm
- Applicable to energy-based models
- Computational complexity of partition function handled using contrastive divergence algorithm