

# The Log-likelihood Gradient

Sargur N. Srihari  
srihari@cedar.buffalo.edu

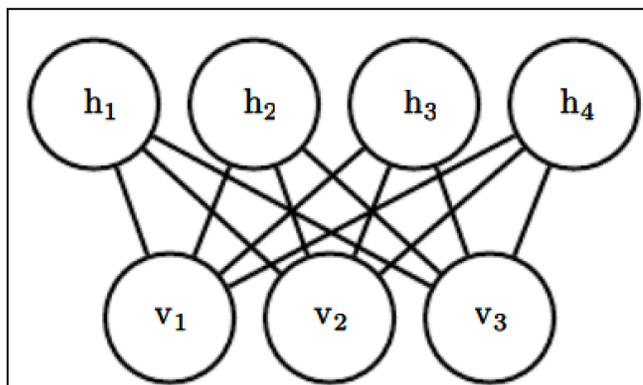
# Topics

- Definition of Partition Function
1. The log-likelihood gradient
  2. Stochastic maximum likelihood and contrastive divergence
  3. Pseudolikelihood
  4. Score matching and Ratio matching
  5. Denoising score matching
  6. Noise-contrastive estimation
  7. Estimating the partition function

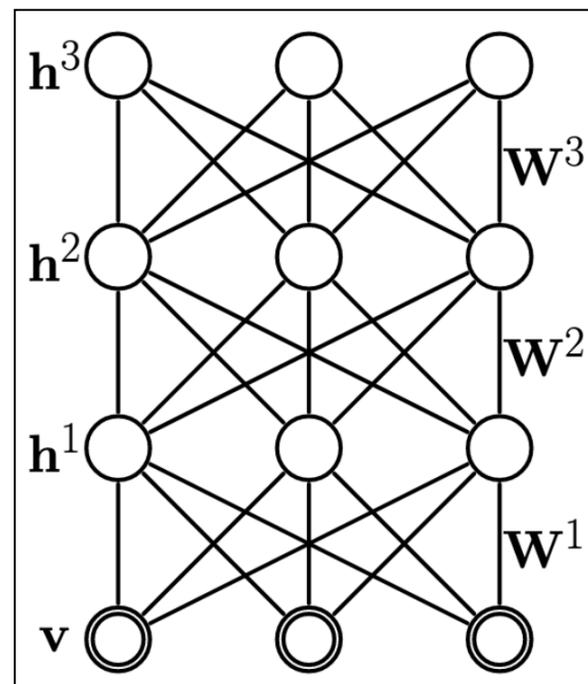
# Undirected models in deep learning

- $p_{\text{model}}(\mathbf{x})$  is an undirected model

A Restricted Boltzmann machine



A deep Boltzmann machine



- We study how the parameters are to be determined

# Finding most likely parameters $\theta$

- Task of interest:
  - Determine parameters  $\theta$  of a Gibbs distribution

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

- where  $Z(\theta) = \sum_x \tilde{p}(\mathbf{x}, \theta)$  is the partition function
- Learning undirected model by MLE is difficult because partition function depends on parameters
- First recall the Maximum Likelihood principle

# Maximum Likelihood Expression

- Given  $m$  i.i.d. examples  $X = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m)}\}$ 
  - From true but unknown distribution  $p_{\text{data}}(\mathbf{x})$
- Let  $p_{\text{model}}(\mathbf{x}; \theta)$  be parametric indexed by  $\theta$ 
  - i.e.,  $p_{\text{model}}(\mathbf{x}; \theta)$  maps any  $\mathbf{x}$  to the true probability  $p_{\text{data}}(\mathbf{x})$

– MLE for  $\theta$  is:

$$\begin{aligned} \theta_{\text{ML}} &= \arg \max_{\theta} p_{\text{model}}(X; \theta) \\ &= \arg \max_{\theta} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \theta) \end{aligned}$$

- Equivalently, by taking logarithms

$$\theta_{\text{ML}} = \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta).$$

- Replacing summation with expectation

$$\theta_{\text{ML}} = \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{x}; \theta).$$

- Which is solved using gradient descent

$$\theta \leftarrow \theta + \varepsilon \mathbf{g}$$

- Where  $\mathbf{g}$  is gradient with terms

$$\mathbf{g}_m = \nabla_{\theta} \log p(\mathbf{x}^{(m)}; \theta) = \nabla_{\theta} \log \tilde{p}(\mathbf{x}^{(m)}; \theta) - \nabla_{\theta} \log Z(\theta)$$

# Gradient has two phases

- Positive and Negative phases of learning
  - Gradient of the log-likelihood wrt parameters has a term corresponding to gradient of partition function

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

$$\nabla_{\boldsymbol{\theta}} \log p(\mathbf{x}; \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}; \boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} \log Z(\boldsymbol{\theta})$$

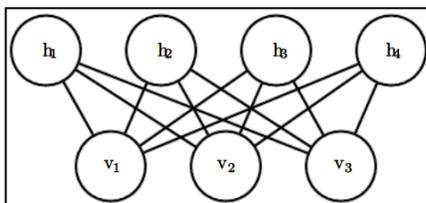
# Tractability: Positive, Negative phases

- For most undirected models: negative phase

$\nabla_{\theta} \log Z(\theta)$  is difficult

- Models with no latent variables or few interactions between latent variables have a tractable positive phase

– RBM: straight-forward positive phase, difficult negative phase



- This chapter is about difficulties with the negative phase

# Computing Gradient for Negative Phase

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} \log Z &= \frac{\nabla_{\boldsymbol{\theta}} Z}{Z} \\ &= \frac{\nabla_{\boldsymbol{\theta}} \sum_{\mathbf{x}} \tilde{p}(\mathbf{x})}{Z} \\ &= \frac{\sum_{\mathbf{x}} \nabla_{\boldsymbol{\theta}} \tilde{p}(\mathbf{x})}{Z}.\end{aligned}$$

For models that guarantee  $p(\mathbf{x}) > 0$  for all  $\mathbf{x}$  we can substitute  $\exp(\log \tilde{p}(\mathbf{x}))$  for  $\tilde{p}(\mathbf{x})$ .

$$\begin{aligned}& \frac{\sum_{\mathbf{x}} \nabla_{\boldsymbol{\theta}} \exp(\log \tilde{p}(\mathbf{x}))}{Z} \\ &= \frac{\sum_{\mathbf{x}} \exp(\log \tilde{p}(\mathbf{x})) \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x})}{Z} \\ &= \frac{\sum_{\mathbf{x}} \tilde{p}(\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x})}{Z} \\ &= \sum_{\mathbf{x}} p(\mathbf{x}) \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x})\end{aligned}$$

Derivation made use of summation over discrete  $\mathbf{x}$   
Similar result applies using integration over continuous  $\mathbf{x}$   
In the continuous version we use Leibniz rule for differentiation

$$\nabla_{\boldsymbol{\theta}} \int \tilde{p}(\mathbf{x}) d\mathbf{x} = \int \nabla_{\boldsymbol{\theta}} \tilde{p}(\mathbf{x}) d\mathbf{x}$$

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

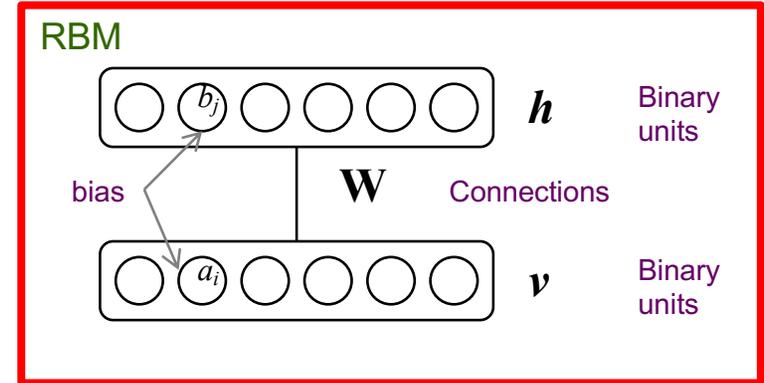
# MLE for RBM

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:

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

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

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

# Monte Carlo methods

- **The identity**  $\nabla_{\theta} \log Z(\theta) = E_{\mathbf{x} \sim p(\mathbf{x})} \nabla_{\theta} \log \tilde{p}(\mathbf{x})$ 
  - is the basis for Monte Carlo methods for MLE of models with intractable partition functions
- MC approach provides intuitive framework for both positive and negative phases
  - Positive phase: increase  $\log \tilde{p}(\mathbf{x})$  for  $\mathbf{x}$  drawn from data
  - Negative phase: decrease partition function by decreasing  $\log \tilde{p}(\mathbf{x})$  drawn from the model distribution