

# Stochastic Maximum Likelihood and Contrastive Divergence

Sargur N. Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics in Partition Function

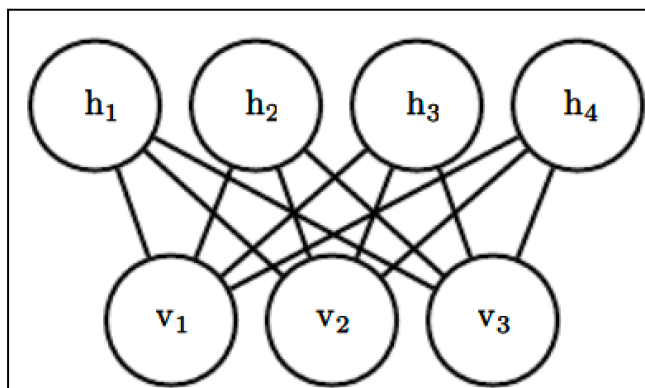
- Definition of Partition Function
  1. The log-likelihood gradient
  2. Stochastic maximum likelihood and contrastive divergence
  3. Pseudolikelihood
  4. Score matching and Ratio matching
  5. Denoising score matching
  6. Noise-contrastive estimation
  7. Estimating the partition function

# Undirected models in deep learning

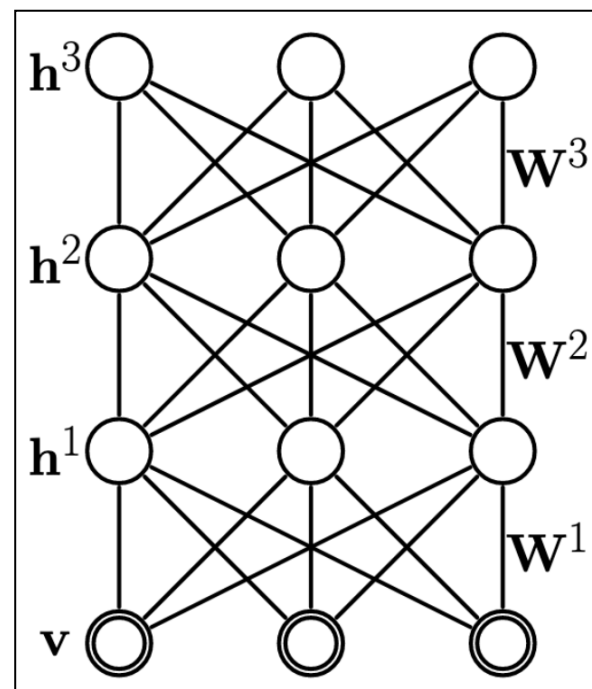
- $p_{\text{model}}(\mathbf{x})$  is an undirected model

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

A Restricted Boltzmann machine



A deep Boltzmann machine



# Stochastic Maximum Likelihood

- Gibbs distribution

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

- Max likelihood estimation of Gibbs parameters

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

- Naive implementation of negative phase

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

- Compute by burning set of Markov chains from a random initialization every time gradient is needed
- When learning is performed by SGD this means chains must be burned in once per gradient step
  - This approach leads to following training procedure

# A Naïve MCMC algorithm

- Algorithm 18.1 Maximizing log-likelihood with an intractable partition function using gradient ascent

Set  $\epsilon$ , the step size, to a small positive number.

Set  $k$ , the number of Gibbs steps, high enough to allow burn in. Perhaps 100 to train an RBM on a small image patch.

**while** not converged **do**

Sample a minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from the training set.

$\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$ . ← **Positive phase: sum over data set**

Initialize a set of  $m$  samples  $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$  to random values (e.g., from a uniform or normal distribution, or possibly a distribution with marginals matched to the model's marginals).

**for**  $i = 1$  to  $k$  **do**

**for**  $j = 1$  to  $m$  **do**

$\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs\_update}(\tilde{\mathbf{x}}^{(j)})$ . ← **Negative phase**

**end for**

**end for**

$\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta})$ .

$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}$ . **Parameter update**

**end while**

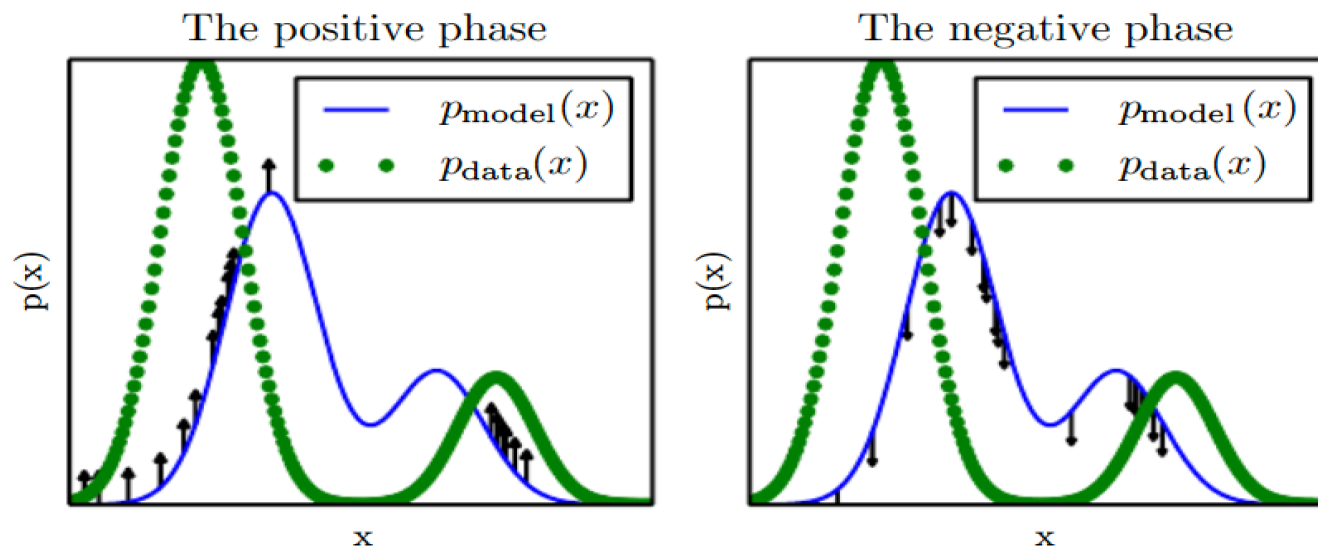
**Negative phase**

With random initialization Gibbs samples from the model are used to determine gradient update

# MCMC: balance between forces

- MCMC approach tries to balance two forces
  - Pushing up on model distribution where data occurs
  - Pushing down on model distribution where model samples occur
- The two forces correspond to
  - Maximizing  $\log \tilde{p}(\mathbf{x})$  and minimizing  $\log Z$
  - For RBM
    - Maximizing  $\tilde{p}(\mathbf{v}) = \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))$  over observed  $\mathbf{v}$
    - Minimizing  $Z = \sum_{\mathbf{v}, \mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}))$  over the model  $\mathbf{h}$

# Phases of MCMC Algorithm



## Positive Phase:

We sample points from the data distribution and push up on their unnormalized probability

## Negative Phase:

We sample points from the model distribution and push down on their unnormalized probability.

This counteracts positive phase's tendency to just add a constant to the unnormalized probability everywhere

When data distribution and model distribution are equal, the positive phase has the same chance to push up at a point as the negative phase has to push down, When this happens, there is no longer any gradient (in expectation) and training must terminate.

# Interpreting the negative phase

- Negative phase involves drawing samples from the model's distribution
  - Can think of it as finding points that the model believes in strongly
- Because the negative phase acts to reduce the probability of those points, they are generally considered to represent the model's incorrect beliefs about the world
- Referred to as *hallucinations* or *fantasy particles*

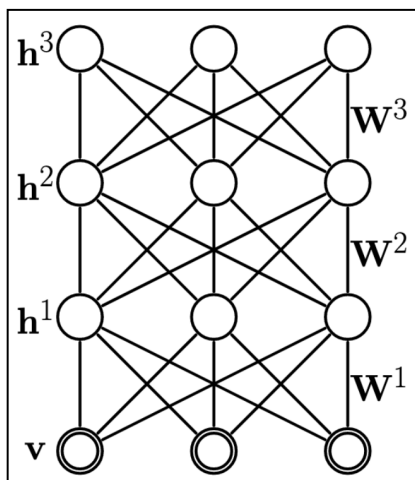


# Wake-Sleep phases of deep learning

Brain maintains a probabilistic model of the world

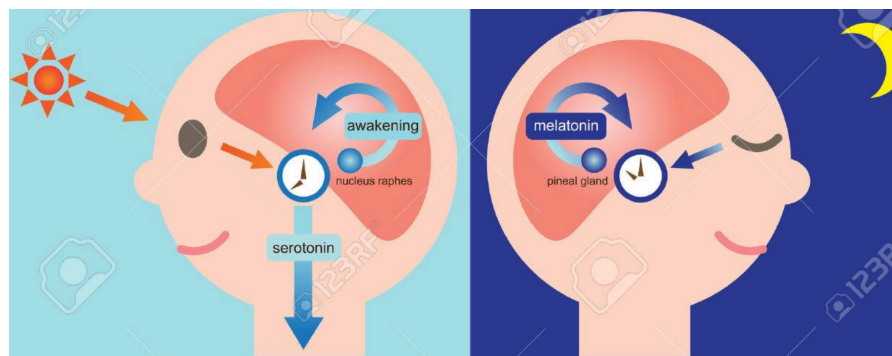
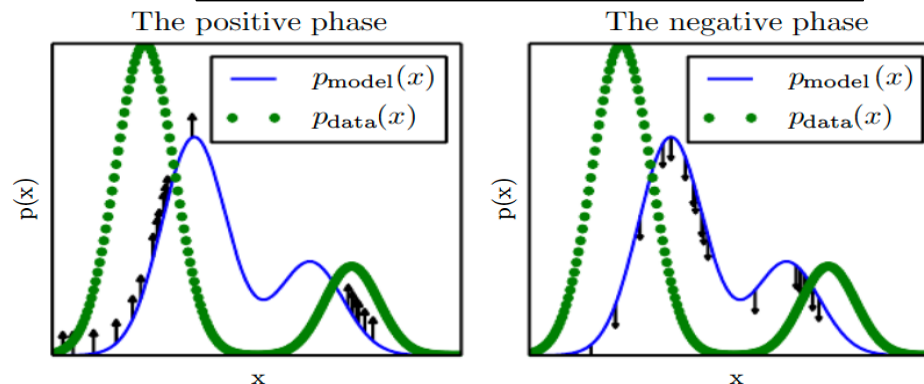
$$p(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \tilde{p}(\mathbf{x}, \boldsymbol{\theta})$$

A deep Boltzmann machine to learn a distribution



Contrastive Divergence Learning Algorithm

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



Follow gradient of  $\log \tilde{p}$  while experiencing real events while awake

Follow negative gradient of  $\log \tilde{p}$  to minimize  $\log Z$  while sleeping and experiencing events sampled from the present model

# Neuroscientific analogy

- Negative phase proposed as explanation of dreaming in humans and other animals
  - Brain maintains a probabilistic model of the world
    - Follows the gradient of  $\log \tilde{p}$  while experiencing real events while awake
    - Follows the negative gradient of  $\log \tilde{p}$  to minimize  $\log Z$  while sleeping and experiencing events sampled from the present model
- Not proven with neuroscientific experiments
  - In ML it is necessary to use positive and negative phases simultaneously
    - Rather than separate periods of wakefulness and REM sleep

# A design less expensive than MCMC

- Given positive/negative phases of learning
  - We can design a less expensive alternative to naive MCMC
- Main cost of naive MCMC:
  - Cost of burning-in the Markov chains from a random initialization at each step
- Natural solution:
  - Initialize Markov chains from a distribution that is very close to the model distribution
    - So that burn-in operation does not take many steps

# Contrastive Divergence algorithm

- CD or CD- $k$  to indicate CD with  $k$  Gibbs steps
- Algorithm initializes the Markov chain at each step with samples from the data distribution
  - This approach presented as Algorithm 18.2
- Obtaining samples from data distribution is free
  - Because they are already in the data set

# CD algorithm intuition

- Initially data distribution is not close to model distribution, so the negative phase is inaccurate
  - Positive phase can increase model's data probability
- After positive phase has time to act, the model distribution is closer to the data distribution
  - And the negative phase starts to become accurate

# Contrastive Divergence Algorithm

- Algorithm 18.2 CD Using gradient ascent as optimization procedure

---

Set  $\epsilon$ , the step size, to a small positive number.

Set  $k$ , the number of Gibbs steps, high enough to allow a Markov chain sampling from  $p(\mathbf{x}; \boldsymbol{\theta})$  to mix when initialized from  $p_{\text{data}}$ . Perhaps 1-20 to train an RBM on a small image patch.

**while** not converged **do**

Sample a minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from the training set.

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

Positive phase

**for**  $i = 1$  to  $m$  **do**

$$\tilde{\mathbf{x}}^{(i)} \leftarrow \mathbf{x}^{(i)}.$$

**end for**

**for**  $i = 1$  to  $k$  **do**

**for**  $j = 1$  to  $m$  **do**

$$\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs\_update}(\tilde{\mathbf{x}}^{(j)}).$$

**end for**

**end for**

$$\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta}).$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}.$$

**end while**

---

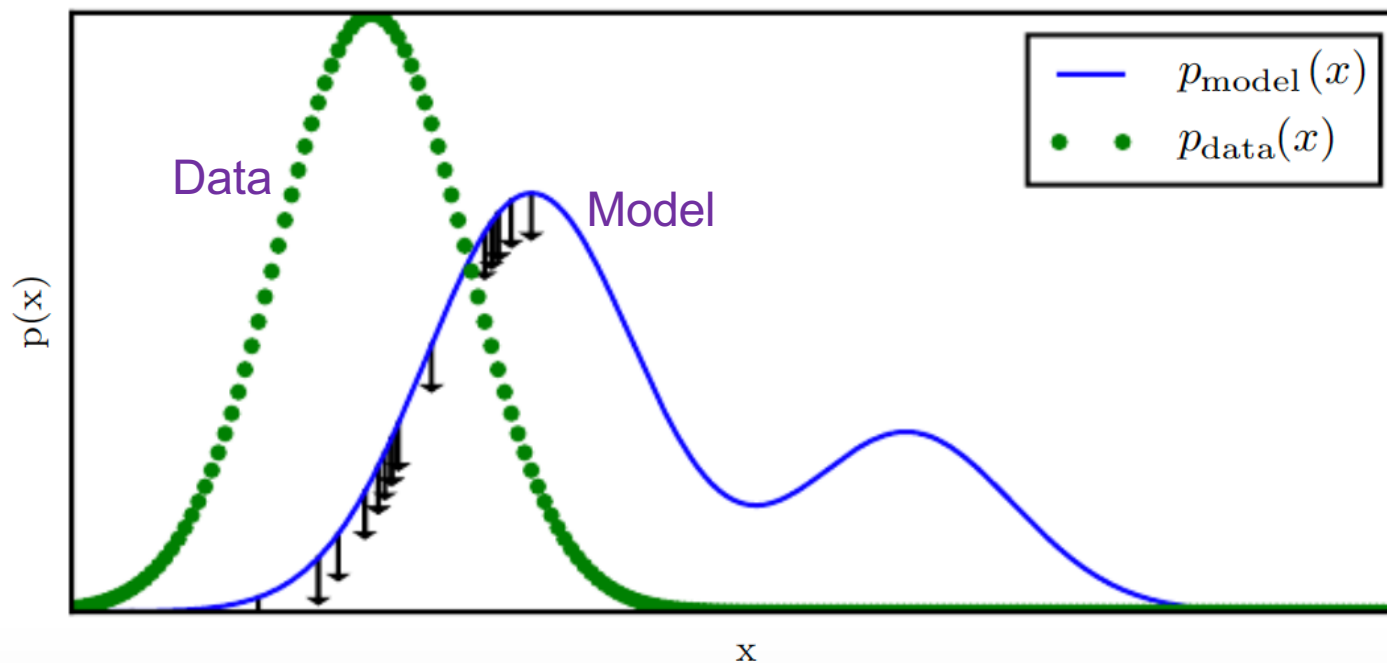
Negative Phase:

Gibbs samples from the model initialized with training data are used to determine gradient update

# CD is still an approximation

- CD is still an approximation to the correct negative phase
- The main way that CD qualitatively fails to implement the correct negative phase is:
  - It fails to suppress regions of high probability that are far from actual training examples
    - These regions that have high probability under the model but low probability under the data generating distribution are called spurious modes
    - See figure as to why this happens

# How negative phase of CD fails to suppress spurious modes



Essentially, it is because modes in the model distribution that are far from the data distribution will not be visited by Markov chains initialized at training points, unless  $k$  is very large.



# Use of Contrastive Divergence

- CD useful for training shallow models like RBMs
- These can in turn be stacked to initialize deeper models like DBNs or DBMs
- However, CD does not provide much help for training deeper models directly
  - Because it is difficult to obtain samples of the hidden units given samples of the visible units

# Stochastic Maximum Likelihood

- Different strategy
- Initialize the Markov chains at each gradient step with their states from the previous gradient step
- Basic idea:
  - As long as the steps taken by the stochastic gradient algorithm are small, then the model from the previous step will be similar to the model from the current step
- Also known as Persistent Contrastive Divergence, or PCD- $k$

# SML Algorithm

- Algorithm 18.3 The stochastic maximum likelihood / persistent contrastive divergence algorithm using gradient ascent as the optimization procedure

Set  $\epsilon$ , the step size, to a small positive number.

Set  $k$ , the number of Gibbs steps, high enough to allow a Markov chain sampling from  $p(\mathbf{x}; \boldsymbol{\theta} + \epsilon \mathbf{g})$  to burn in, starting from samples from  $p(\mathbf{x}; \boldsymbol{\theta})$ . Perhaps 1 for RBM on a small image patch, or 5-50 for a more complicated model like a DBM. Initialize a set of  $m$  samples  $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$  to random values (e.g., from a uniform or normal distribution, or possibly a distribution with marginals matched to the model's marginals).

**while** not converged **do**

Sample a minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from the training set.

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

**for**  $i = 1$  to  $k$  **do**

**for**  $j = 1$  to  $m$  **do**

$\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs\_update}(\tilde{\mathbf{x}}^{(j)})$ .

**end for**

**end for**

$\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta})$ .

$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}$ .

**end while**