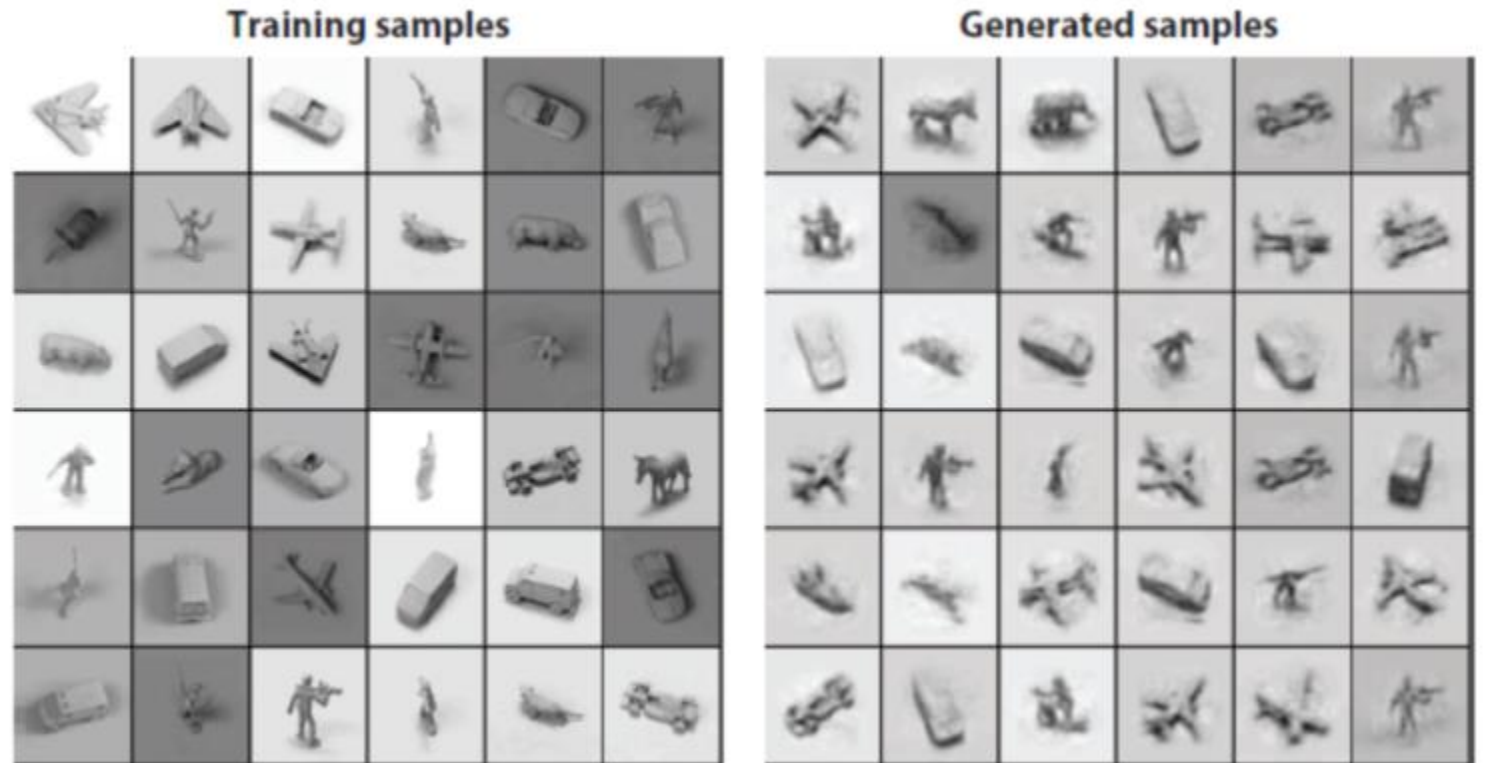


Boltzmann machines



From Hopfield to Boltzmann

- Hopfield networks minimize the quadratic energy function

$$E = -f_{\theta}(\mathbf{x}) = -\left(\sum_{i,j} w_{ij}x_i x_j + \sum_i b_i x_i\right)$$

- Boltzmann machines are stochastic Hopfield networks
- In Boltzmann machines the neuron response on activation a_i is

$$x_i = \begin{cases} +1 & \text{with probability } 1/(1 + \exp(-2a_i)) \\ -1 & \text{otherwise} \end{cases}$$

- Gibbs sampling for pdf $p(\mathbf{x}) = \frac{1}{Z} \exp(\frac{1}{2} \mathbf{x}^T \mathbf{W} \mathbf{x})$

Restricted Boltzmann machines

- Boltzmann machines are too parameter heavy
 - For \mathbf{x} with $256 \times 256 = 65536$ the \mathbf{W} has 4.2 billion parameters
- Boltzmann machines learn no features

- Instead, add bottleneck latents \mathbf{v}

$$E = -f_{\theta}(\mathbf{x}) = -\left(\sum_{i,j} w_{ij}x_i v_j + \sum_i b_i x_i + \sum_j c_j v_j\right)$$

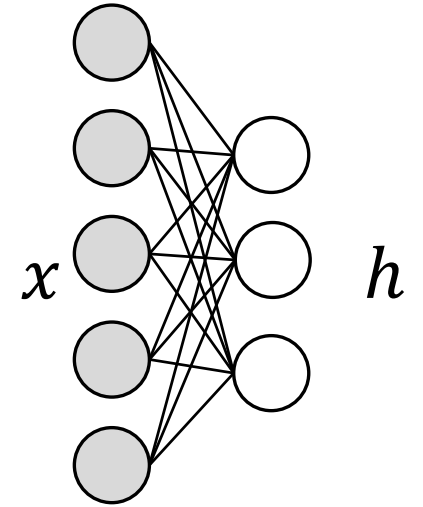
- x_i and v_j are still binary variables in the original model
- The quadratic term captures correlations
- The unary terms capture priors: how likely is a (latent) pixel to be +1 or -1

Restricted Boltzmann Machines

- Energy function: $E(\mathbf{x}) = -\mathbf{x}^T \mathbf{W} \mathbf{v} - \mathbf{b}^T \mathbf{x} - \mathbf{c}^T \mathbf{v}$

$$p(\mathbf{x}) = \frac{1}{Z} \sum_{\mathbf{v}} \exp(-E(\mathbf{x}, \mathbf{v}))$$

- Not in the form $\propto \exp(\mathbf{x})/Z$ because of the \sum



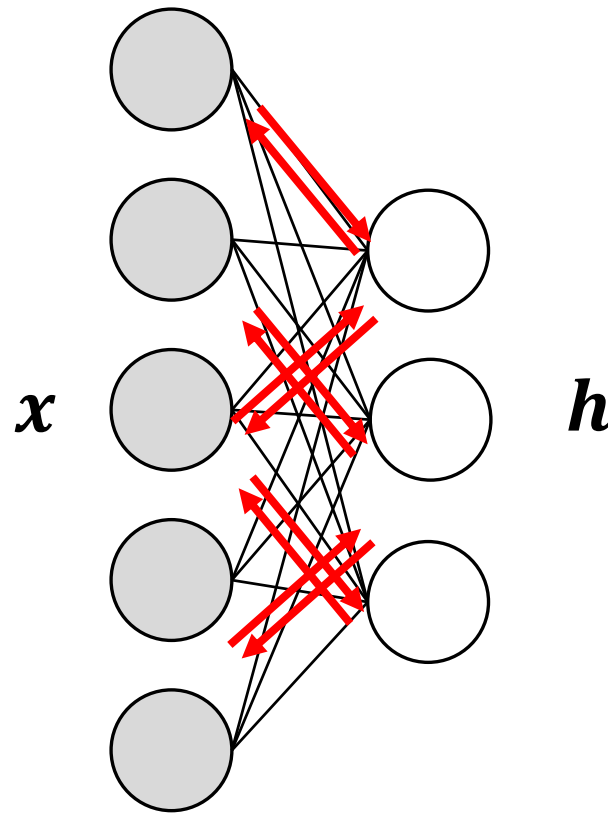
- Free energy function: $F(\mathbf{x}) = -\mathbf{b}^T \mathbf{x} - \sum_i \log \sum_{v_i} \exp(v_i(c_i + \mathbf{W}_i \mathbf{x}))$

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

$$Z = \sum_{\mathbf{x}} \exp(-F(\mathbf{x}))$$

Restricted Boltzmann Machines

- The $F(\mathbf{x})$ defines a bipartite graph with undirected connections
 - Information flows forward and backward



Restricted Boltzmann Machines

- The hidden variables v_j are independent conditioned on the visible variables

$$p(\mathbf{v}|\mathbf{x}) = \prod_j p(v_j|\mathbf{x}, \boldsymbol{\theta})$$

- The visible variables x_i are independent conditioned on the hidden variables

$$p(\mathbf{x}|\mathbf{v}) = \prod_i p(x_i|\mathbf{v}, \boldsymbol{\theta})$$

Training RBM conditional probabilities

- The conditional probabilities are defined as sigmoids

$$p(v_j | \mathbf{x}, \boldsymbol{\theta}) = \sigma(\mathbf{W}_{.j} \mathbf{x} + b_j)$$

$$p(x_i | \mathbf{v}, \boldsymbol{\theta}) = \sigma(\mathbf{v}^T \mathbf{W}_{i.} + c_i)$$

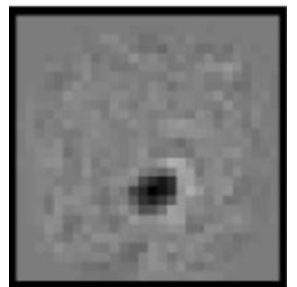
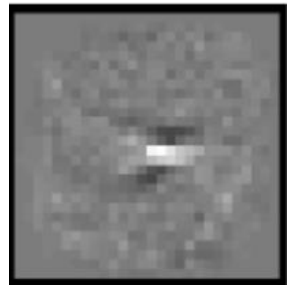
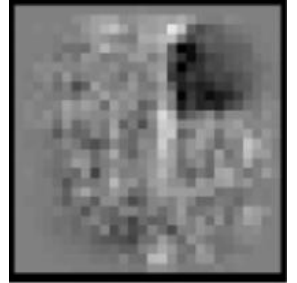
- Since RBMs are bidirectional \Rightarrow “Loop” between visible and latent

$$\mathbf{v}^{(1)} \sim \sigma(\mathbf{W}_{.j} \mathbf{x}^{(0)} + b_j) \Rightarrow$$

$$\mathbf{x}^{(1)} \sim \sigma(\mathbf{W}_{.j} \mathbf{v}^{(2)} + b_j) \Rightarrow$$

$$\mathbf{v}^{(2)} \sim \sigma(\mathbf{W}_{.j} \mathbf{x}^{(1)} + b_j) \Rightarrow \dots$$

Latent activations



Training any energy model

- Maximizing log-likelihood

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \sum_n \log p_{\boldsymbol{\theta}}(\mathbf{x}_n) = \mathbb{E}_{p_0}[\log p_{\boldsymbol{\theta}}(\mathbf{x})]$$

- The expectation w.r.t. a pdf is equivalent to
 - sampling from the pdf and
 - then taking the average

$$\mathbb{E}_{\mathbf{x} \sim p_0}[\log p(\mathbf{x}|\boldsymbol{\theta})] = \mathbb{E}_{\mathbf{x} \sim p_0}[-E_{\boldsymbol{\theta}}(\mathbf{x})] - \log Z(\boldsymbol{\theta})$$

- where $\log Z(\boldsymbol{\theta}) = \log \sum_{\mathbf{x}'} \exp(-E_{\boldsymbol{\theta}}(\mathbf{x}'))$
- and $p_0(\mathbf{x})$ is the data distribution

Taking gradients of any energy model

$$\begin{aligned}\frac{d}{d\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x}) &= -\frac{d}{\partial \boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) - \frac{d}{d\boldsymbol{\theta}} \log Z(\boldsymbol{\theta}) = \\ &= -\frac{d}{\partial \boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) - \frac{1}{Z(\boldsymbol{\theta})} \frac{d}{d\boldsymbol{\theta}} Z(\boldsymbol{\theta}) \\ &= -\frac{d}{\partial \boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) - \sum_{\mathbf{x}'} \frac{1}{Z(\boldsymbol{\theta})} \exp(-E_{\boldsymbol{\theta}}[\mathbf{x}']) \left(-\frac{d}{d\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}') \right) \\ &= -\frac{d}{\partial \boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) + \sum_{\mathbf{x}'} p_{\boldsymbol{\theta}}(\mathbf{x}') \frac{d}{d\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}') \\ &= -\frac{d}{\partial \boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) + \mathbb{E}_{\mathbf{x}' \sim p_{\boldsymbol{\theta}}} \left[\frac{d}{\partial \boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}') \right]\end{aligned}$$

Remember: $\sum p(x) f(x) = \mathbb{E}_{p(x)}[f(x)]$
 $\int_x p(x) f(x) dx = \mathbb{E}_{p(x)}[f(x)]$

Taking gradients in an RBM

- For an RBM we must integrate out the latent variables

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{N} \sum_n \log p_{\boldsymbol{\theta}}(\mathbf{x}_n) = \frac{1}{N} \sum_n \log \sum_v p_{\boldsymbol{\theta}}(\mathbf{x}_n, \mathbf{v})$$

$$\frac{\partial}{\partial \boldsymbol{\theta}} \log \sum_v p_{\boldsymbol{\theta}}(\mathbf{x}_n, \mathbf{v}) = -\mathbb{E}_{\mathbf{v} \sim p_{\boldsymbol{\theta}}(\mathbf{v}|\mathbf{x}_n)} \left[\frac{d}{d\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}_n, \mathbf{v}) \right] + \mathbb{E}_{\mathbf{x}', \mathbf{v} \sim p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{v})} \left[\frac{d}{d\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}', \mathbf{v}) \right]$$

Taking gradients in an RBM

- And since for RBM $E_{\theta}(\mathbf{x}, \mathbf{v}) = -\mathbf{v}^T \mathbf{W} \mathbf{x} - \mathbf{b}^T \mathbf{x} - \mathbf{c}^T \mathbf{v}$

$$\frac{d}{dW_{ij}} E_{\theta}(x_i, v_j) = -x_i v_j \Rightarrow$$

$$\frac{d\mathcal{L}}{dW_{ij}} = \mathbb{E}_{\mathbf{v} \sim p_{\theta}(\mathbf{v} | \mathbf{x}_n)} [x_i v_j] - \mathbb{E}_{\mathbf{x}', \mathbf{v} \sim p_{\theta}(\mathbf{x}, \mathbf{v})} [x_i v_j]$$

- **Easy:** substitute \mathbf{x}_n and sum over \mathbf{v}
- **Hard (normalization):** sum over all 2^{m+d} combinations of images & latents
 - Intractable due to exponential complexity w.r.t. $m + d$
 - Evaluating and optimizing $p_{\theta}(\mathbf{x}, \mathbf{v})$ takes a long time
 - If we had only the unnormalized part we would have no problem

Tackling intractability by sampling

- $\mathbb{E}_{\mathbf{x}', \mathbf{v} \sim p_{\theta}(\mathbf{x}, \mathbf{v})} \left[\frac{d}{d\theta} E_{\theta}(\mathbf{x}', \mathbf{v}) \right]$ stands for an expectation
 - One can sample very many \mathbf{x}', \mathbf{v} from $p_{\theta}(\mathbf{x}, \mathbf{v})$
 - Take average instead of computing analytically (Monte Carlo sampling)
- Question: how to even sample from a hard pdf?
 - Markov Chain Monte Carlo with Gibbs sampling
 - Convergence after many rounds

Initialization: Initialize $\mathbf{x}^{(0)} \in \mathcal{R}^D$ and number of samples N

```
• for  $i = 0$  to  $N - 1$  do
•    $x_1^{(i+1)} \sim p(x_1 | x_2^{(i)}, x_3^{(i)}, \dots, x_D^{(i)})$ 
•    $x_2^{(i+1)} \sim p(x_2 | x_1^{(i+1)}, x_3^{(i)}, \dots, x_D^{(i)})$ 
•    $\vdots$ 
•    $x_j^{(i+1)} \sim p(x_j | x_1^{(i+1)}, x_2^{(i+1)}, \dots, x_{j-1}^{(i+1)}, x_{j+1}^{(i)}, \dots, x_D^{(i)})$ 
•    $\vdots$ 
•    $x_D^{(i+1)} \sim p(x_D | x_1^{(i+1)}, x_2^{(i+1)}, \dots, x_{D-1}^{(i+1)})$ 
return  $(\{\mathbf{x}^{(i)}\}_{i=0}^{N-1})$ 
```

Sampling the normalizing constant

- We can rewrite the gradient as

$$\frac{d}{d\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = -\mathbb{E}_0 \left[\frac{d}{d\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) \right] + \mathbb{E}_{\infty} \left[\frac{d}{d\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}') \right]$$

- $\mathbb{E}_0 \equiv E_{\mathbf{x} \sim p_0}$ means sampling from training data and average gradients
 - $\mathbb{E}_{\infty} \equiv E_{\mathbf{x}, \mathbf{v} \sim p_{\boldsymbol{\theta}}}$ means sampling from the model and average gradients
- Unfortunately, MCMC can be very slow \rightarrow 2nd source of intractability

Ergo, contrastive diverge learning

- To motivate contrastive divergence, we revisit maximum likelihood learning

$$\text{KL}(p_0 \parallel p_\infty) = \int p_0 \log p_0 - \int p_0 \log p_\infty \propto - \int p_0 \log p_\infty$$

- Contrastive divergence minimizes

$$\text{CD}_n = \text{KL}(p_0 \parallel p_\infty) - \text{KL}(p_n \parallel p_\infty)$$

- Updates weights using CD_n gradients instead of ML gradients

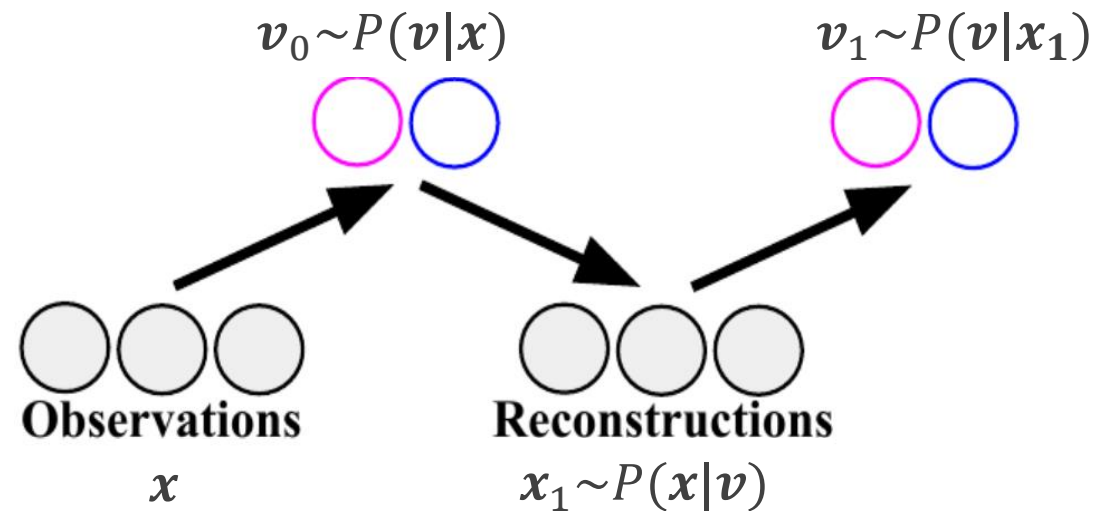
$$\frac{d}{d\boldsymbol{\theta}} \text{CD}_n = -\mathbb{E}_0 \left[\frac{d}{d\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}) \right] + \mathbb{E}_n \left[\frac{d}{d\boldsymbol{\theta}} E_{\boldsymbol{\theta}}(\mathbf{x}') \right] + \frac{d}{d\boldsymbol{\theta}} [\dots]$$

- where \mathbb{E}_n is computed by sampling after n steps in the Markov Chain
- The last term is small and can be ignored

Hinton, *Training Products of Experts by Minimizing Contrastive Divergence*, Neural Computation, 2002

Contrastive diverge learning: intuition

- Make sure after n sampling step not far from data distribution
 - Usually, one step only ($n=1$) is enough
 - Something similar to ‘minimizing reconstruction error’
- Because of conditional independence of $\mathbf{x}|\mathbf{v}$ and $\mathbf{v}|\mathbf{x} \rightarrow$ parallel computations
 - Sample a data point \mathbf{x}
 - Compute the posterior $\mathbf{p}(\mathbf{v}|\mathbf{x})$
 - Take sample of latents $\mathbf{v} \sim \mathbf{p}(\mathbf{v}|\mathbf{x})$
 - Compute the conditional $\mathbf{p}(\mathbf{x}|\mathbf{v})$
 - Sample from $\mathbf{x}' \sim \mathbf{p}(\mathbf{x}|\mathbf{v})$
 - Minimize difference using \mathbf{x}, \mathbf{x}'

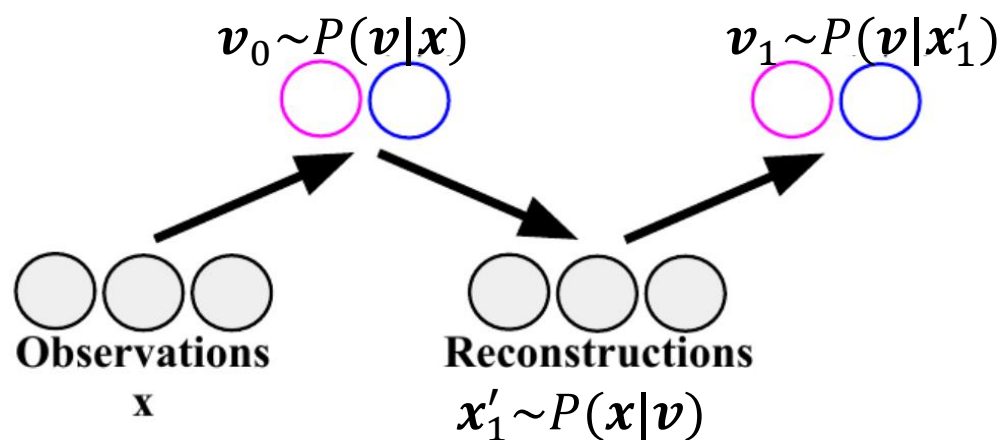


Contrastive divergence for RBMs

- Contrastive divergence approximates gradient by k-steps Gibbs sampler

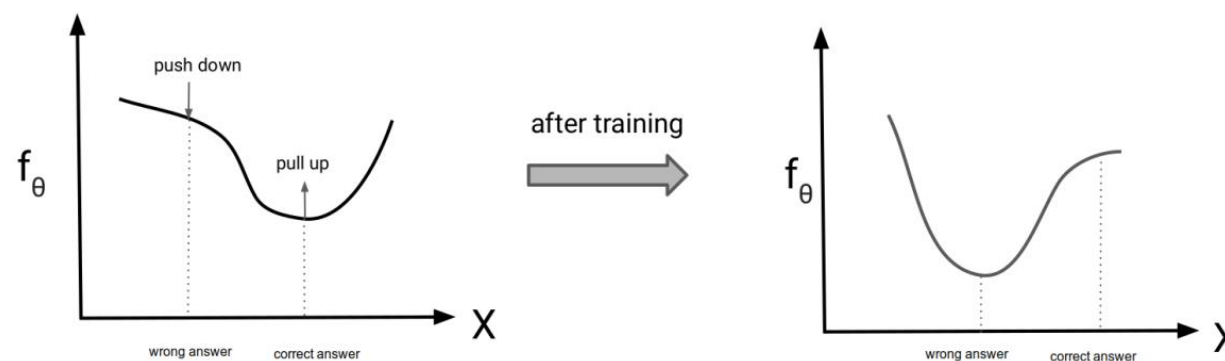
$$\frac{d}{d\theta} \log p(\mathbf{x}_n | \theta) = -\frac{d}{d\theta} E_{\theta}(\mathbf{x}_n, \mathbf{v}_0) - \frac{d}{d\theta} E_{\theta}(\mathbf{x}'_k, \mathbf{v}_k)$$

- Pushing the nominator up while pushing the denominator down



Hinton, 2002

Carreira-Perpinan and Hinton, 2005



[Ermon and Grover, Deep Generative Models](#)

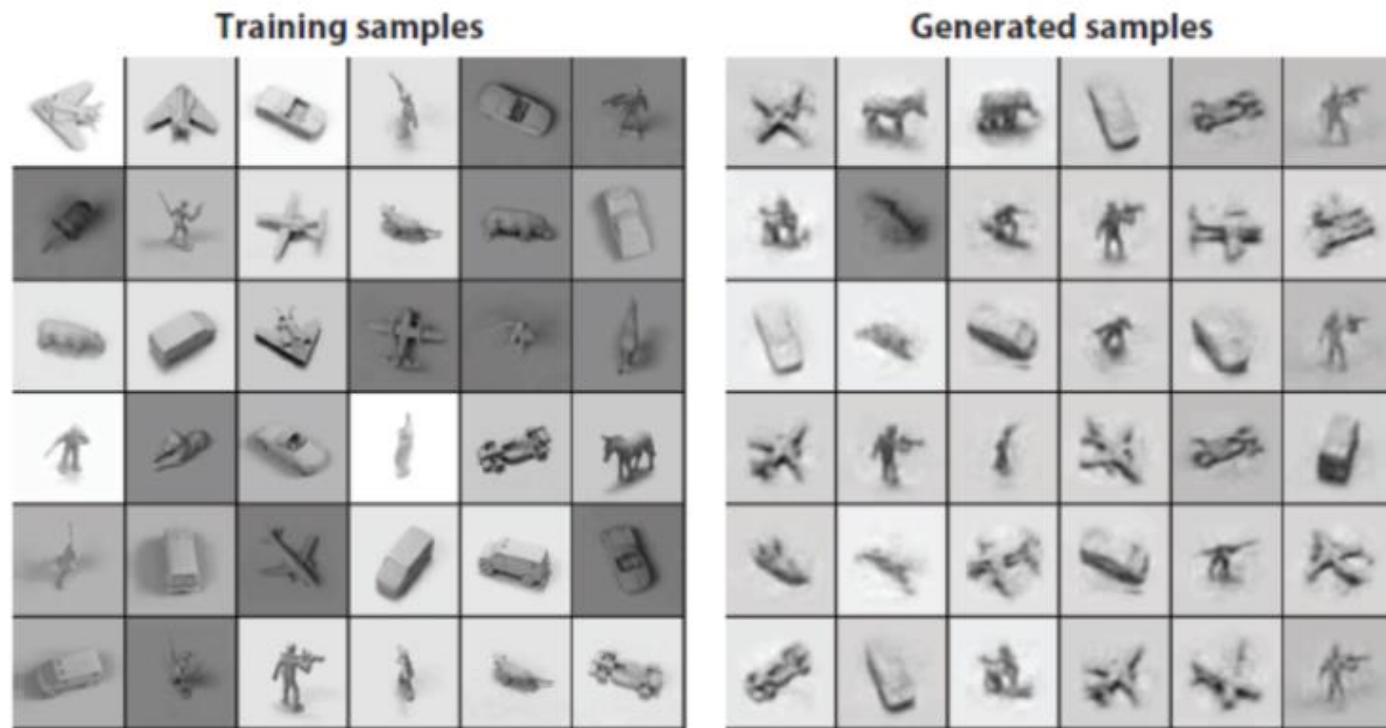
How to sample? Markov Chain Monte Carlo

We want to sample an \mathbf{x} from a pdf $p_{\boldsymbol{\theta}}(\mathbf{x})$ with MCMC with Gibbs sampler

- Step 1. Initialize \mathbf{x}^0 randomly
- Step 2. Let $\hat{\mathbf{x}} = \mathbf{x}^t + \text{noise}$
 - If $f_{\boldsymbol{\theta}}(\hat{\mathbf{x}}) > f_{\boldsymbol{\theta}}(\mathbf{x}^t)$, set $\mathbf{x}^{t+1} = \hat{\mathbf{x}}$
 - Otherwise $\mathbf{x}^{t+1} = \mathbf{x}^t$ with probability $\frac{p(\hat{\mathbf{x}})}{p(\mathbf{x}^t)} = \exp(f_{\boldsymbol{\theta}}(\hat{\mathbf{x}}) - f_{\boldsymbol{\theta}}(\mathbf{x}^t))$
- Go to step 2
- Because of the ratio of likelihoods \rightarrow no $Z(\boldsymbol{\theta})$

Using RBMs

- Some of the first models to show nice generations of images
- Use RBMs to pretrain networks for classification afterward

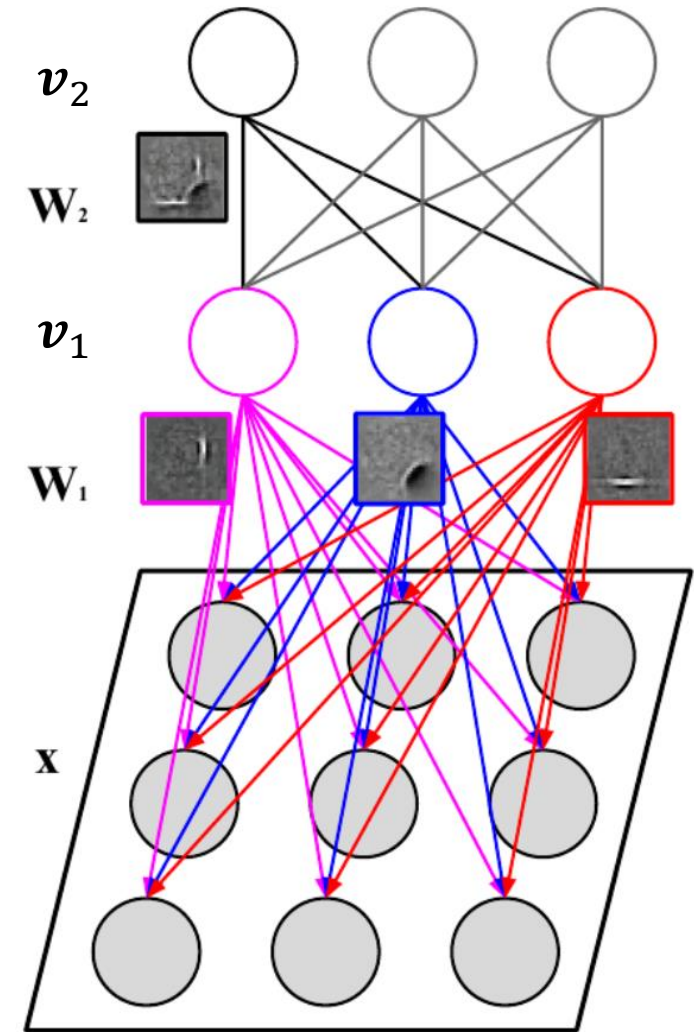


Deep Belief Network

- Stack RBM layers assuming conditional independence

$$p(\mathbf{x}, \mathbf{v}_1, \mathbf{v}_2) = p(\mathbf{x}|\mathbf{v}_1) \cdot p(\mathbf{v}_1|\mathbf{v}_2)$$

- Deep Belief Networks are directed models
- Dense layers with single forward flow
 - As RBM is directional: $p(x_i|\mathbf{v}, \boldsymbol{\theta}) = \sigma(\mathbf{W}_{.i}\mathbf{x} + c_i)$



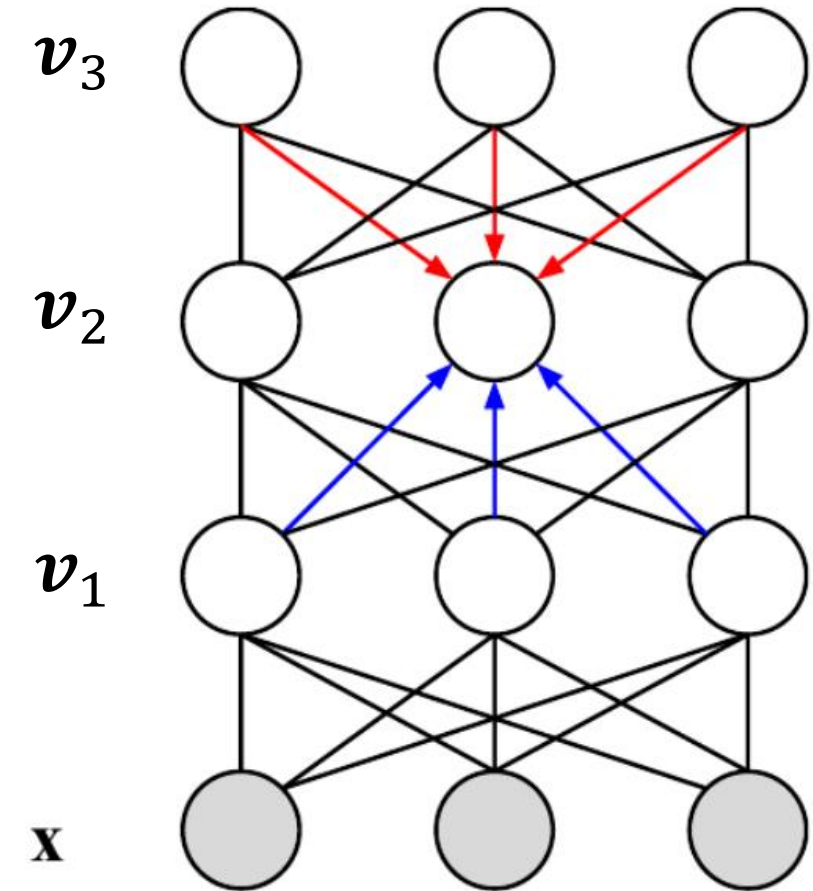
Deep Boltzmann machines

- Stacking RBM layers from **above** and **below** layers
 - Markov model

- Energy function

$$E(\mathbf{x}, \mathbf{v}_1, \mathbf{v}_2 | \boldsymbol{\theta}) = \mathbf{x}^T \mathbf{W}_1 \mathbf{v}_1 + \mathbf{v}_1^T \mathbf{W}_2 \mathbf{v}_2 + \mathbf{v}_2^T \mathbf{W}_3 \mathbf{v}_3$$

$$p(\mathbf{v}_2^k | \mathbf{v}_1, \mathbf{v}_3) = \sigma\left(\sum_j \mathbf{w}_1^{jk} \mathbf{v}_1^j + \sum_l \mathbf{w}_3^{kl} \mathbf{v}_3^l\right)$$



Training deep Boltzmann machines

- Computing gradients is intractable
- Instead, variational methods (mean-field) or sampling methods are used

