

Machine Learning - SS 2011

Lecture 6, Deep Learning

Grégoire Montavon

Technische Universität Berlin

May 17, 2011

Deep learning

- Deep learning is the field of research that is concerned with learning automatically hierarchical representations of data.
- By learning a hierarchical model of data, deep learning aims to simultaneously deliver *high generalization* and *low computational cost* at prediction time.

A spectrum of machine learning tasks

Statistics

- There is not much structure in the data
- Simple models are typically sufficient to model the task
- Main challenge: separating the signal from the noise
- Prior knowledge can bring significant improvement

Artificial Intelligence

- There is a huge amount of structure in the data
- We need complex hierarchical models for the task
- Main challenge: discovering the structure of data
- Prior knowledge is already in the data

A spectrum of machine learning tasks

Statistics

Typical applications

- Brain computer interfaces
- Financial data analysis
- Statistical physics

Methods

- Generalized linear models
- Structured losses
- Kernels

Artificial Intelligence

Typical applications

- Computer vision
- Natural language processing
- Sequential reasoning

Methods

- Engineering
- Deep networks
- Unsupervised learning

Representing probability densities

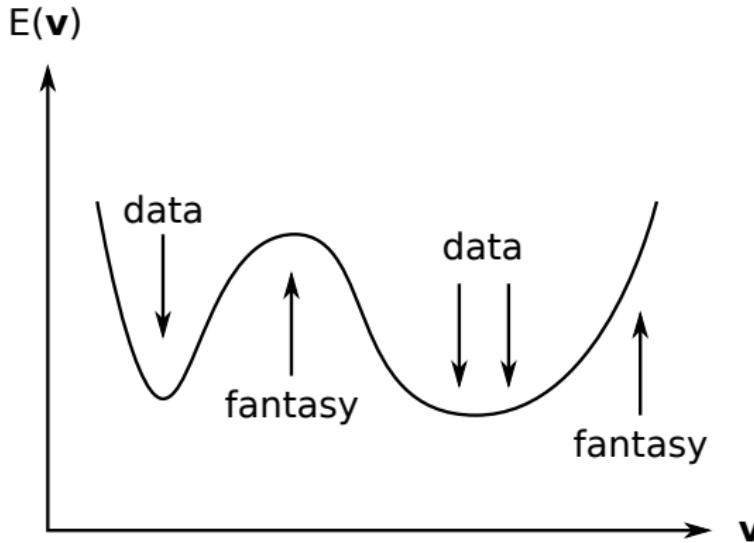
Problem:

- Let $p(\mathbf{v})$, $\mathbf{v} \in \{0, 1\}^d$ be an *unknown* probability distribution
- Estimate $p(\mathbf{v})$ from a finite set of samples

Idea: learn an energy function $E(\mathbf{v}; \theta)$ of such that

- $E(\mathbf{v}; \theta)$ is low when $\mathbf{v} \in \text{data}$
- $E(\mathbf{v}; \theta)$ is high when $\mathbf{v} \notin \text{data}$
- $E(\mathbf{v}; \theta)$ is smooth

Representing probability densities



Find a model such that $E(\mathbf{v}; \theta)$ is low for $\mathbf{v} \in \text{data}$ and high elsewhere

A simple energy-based model

- Define an energy function:

$$E(\mathbf{v}; \theta) = - \sum_{i=1}^d v_i b_i \quad \theta = \mathbf{b}$$

- Associate a probability to the energies:

$$p(\mathbf{v}) = \frac{1}{Z(\theta)} e^{-E(\mathbf{v}; \theta)}$$

Where the partition function $Z(\theta) = \sum_{\mathbf{u}} e^{-E(\mathbf{u}; \theta)}$ forces probabilities to sum to 1.

- Learning rule:

$$b_i = b_i + \epsilon (\langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}})$$

- Problem: doesn't capture dependencies between variables

A simple energy-based model

- Define an energy function:

$$E(\mathbf{v}; \theta) = - \sum_{i=1}^d v_i b_i - \sum_{i,j} v_i v_j w_{ij} \quad \theta = (W, \mathbf{b})$$

- Associate a probability to the energies:

$$p(\mathbf{v}) = \frac{1}{Z(\theta)} e^{-E(\mathbf{v}; \theta)}$$

- Learning rule:

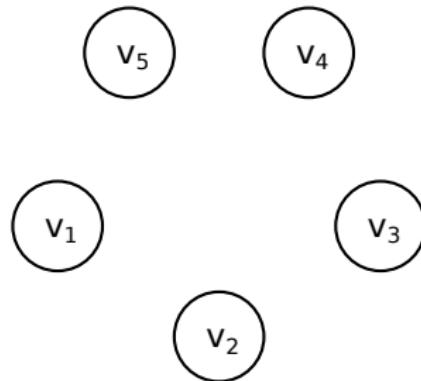
$$b_i = b_i + \epsilon (\langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}})$$

$$w_{ij} = w_{ij} + \epsilon (\langle v_i v_j \rangle_{\text{data}} - \langle v_i v_j \rangle_{\text{model}})$$

A simple energy-based model

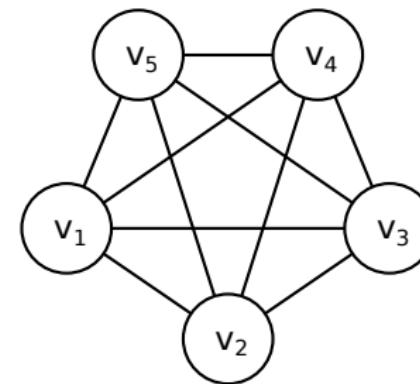
Independent variables

$$\theta = \mathbf{b}$$



Dependent variables

$$\theta = (W, \mathbf{b})$$



$$E(\mathbf{v}; \theta) = - \sum_{i=1}^d v_i b_i$$

$$E(\mathbf{v}; \theta) = - \sum_{i=1}^d v_i b_i - \sum_{i,j} v_i v_j w_{ij}$$

A simple energy-based model

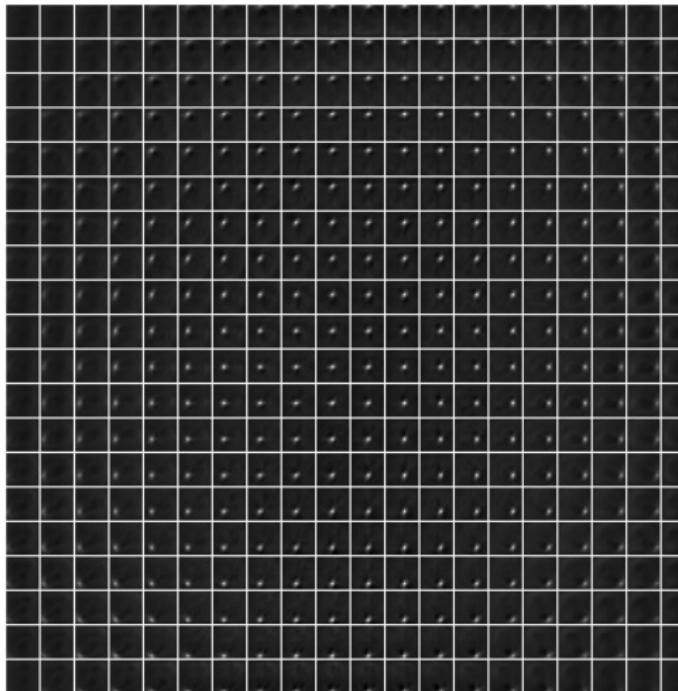
Example: handwritten digits

- Each pixel is a variable
- W models the relation between each pixels



A simple energy-based model

$W =$



$b =$



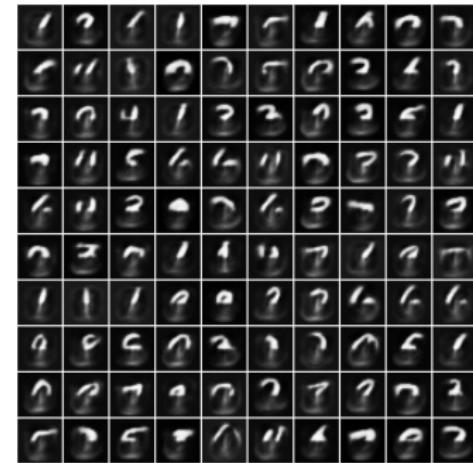
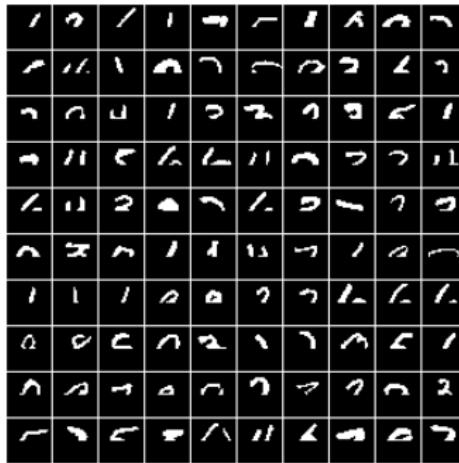
Denoising samples

1	3	6	1	7	5	1	0	9	7
5	4	1	0	2	5	8	3	6	3
7	9	4	1	3	3	2	3	5	1
7	4	5	6	6	4	2	3	2	4
6	4	3	9	2	6	3	7	2	3
9	3	7	1	1	4	7	1	9	7
1	1	1	9	9	2	2	6	6	6
8	8	5	0	3	3	2	0	5	1
0	9	7	9	9	2	3	2	9	3
5	3	5	5	0	4	6	7	9	2



1	3	6	1	7	5	1	0	9	7
5	4	1	0	2	5	8	3	6	3
7	9	4	1	3	3	2	3	5	1
7	4	5	6	6	4	2	3	2	4
6	4	3	9	2	6	3	7	2	3
9	3	7	1	1	4	7	1	9	7
1	1	1	9	9	2	2	6	6	6
8	8	5	0	3	3	2	0	5	1
0	9	7	9	9	2	3	2	9	3
5	3	5	5	0	4	6	7	9	2

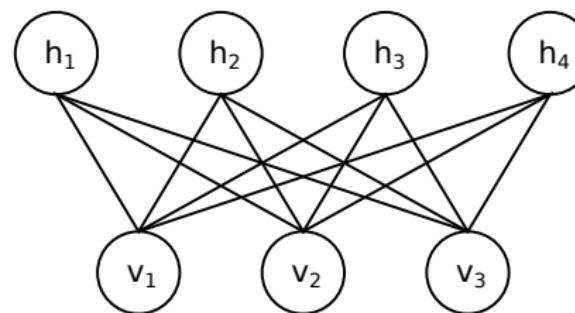
Completing samples



On complicated problems, pairwise interactions between variables are not sufficient to model the distribution. Therefore, we need to introduce hidden variables.

The restricted Boltzmann machine (RBM)

- Let \mathbf{v} be the visible variables (observations)
- Let \mathbf{h} be the hidden variables (causes)



The energy function of the RBM is defined as follows:

$$E(\mathbf{v}, \mathbf{h}; \theta) = - \sum_{i \in \text{visible}} v_i a_i - \sum_{j \in \text{hidden}} h_j b_j - \sum_{i,j} v_i w_{ij} h_j$$

The restricted Boltzmann machine (RBM)

- The energy function of the RBM is defined as follows:

$$E(\mathbf{v}, \mathbf{h}; \theta) = - \sum_{i \in \text{visible}} v_i a_i - \sum_{j \in \text{hidden}} h_j b_j - \sum_{i,j} v_i w_{ij} h_j$$

- A probability $p(\mathbf{v}, \mathbf{h})$ is associated to the energies:

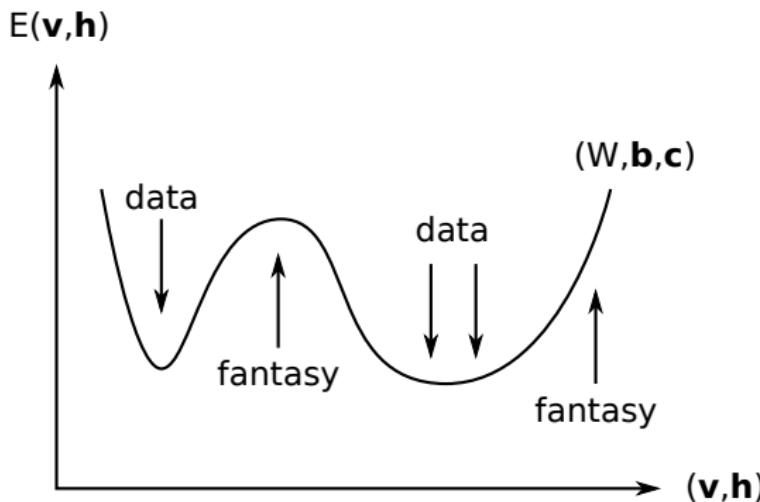
$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z(\theta)} e^{-E(\mathbf{v}, \mathbf{h}; \theta)}$$

- Sampling from the model:

$$p(\mathbf{h}|\mathbf{v}) = \sigma(W^\top \mathbf{v} + \mathbf{b}) \quad \text{and} \quad p(\mathbf{v}|\mathbf{h}) = \sigma(W\mathbf{h} + \mathbf{a})$$

where $\sigma(t) = \frac{1}{1+e^{-t}}$ is the logistic function.

Training an RBM



Idea: find W, b, a such that $E(v, h; \theta)$ is low when $(v, h) \in \text{data}$ and high elsewhere

Training an RBM

- We want to minimize $\langle \log p(\mathbf{v}) \rangle_{\text{data}}$
- Gradient descent:

$$\left\langle \frac{\partial \log p(\mathbf{v})}{\partial w_{ij}} \right\rangle_{\text{data}} = \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}}$$

- Update rule:

$$\Delta w_{ij} = \epsilon (\langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}})$$

- Unfortunately, we don't know $\langle v_i h_j \rangle_{\text{model}}$
- Idea, approximate $\langle v_i h_j \rangle_{\text{model}}$ by Gibbs sampling

Training an RBM

Idea: estimate $\langle v_i h_j \rangle_{\text{model}}$ by running one step of Gibbs sampling from the data.

```

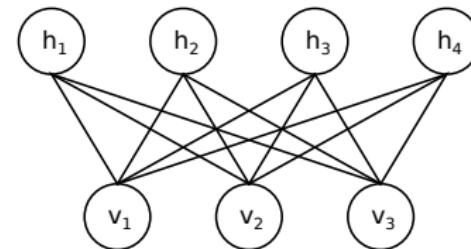
while True:

    # Collect data statistics
     $v_{\text{data}} \sim \text{data}$ 
     $h_{\text{data}} \sim p(h|v_{\text{data}}; (W, a, b))$ 

    # Collect model statistics
     $v_{\text{model}} \sim p(v|h_{\text{data}}; (W, a, b))$ 
     $h_{\text{model}} \sim p(h|v_{\text{model}}; (W, a, b))$ 

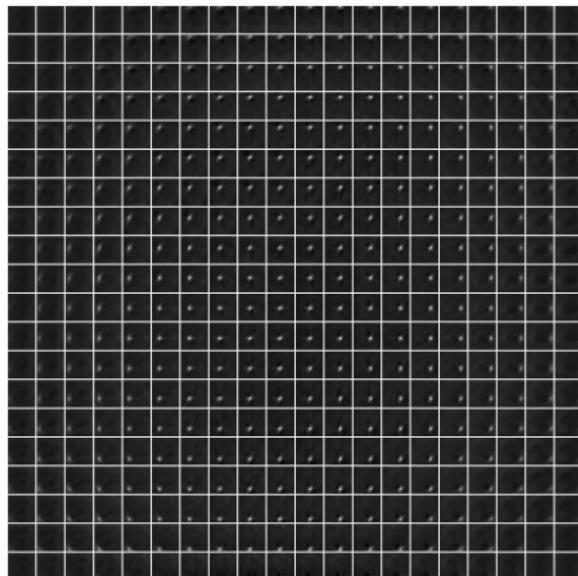
    # Update model parameters
     $W = W + \epsilon(\langle vh^T \rangle_{\text{data}} - \langle vh^T \rangle_{\text{model}})$ 
     $a = a + \epsilon(\langle v \rangle_{\text{data}} - \langle v \rangle_{\text{model}})$ 
     $b = b + \epsilon(\langle h \rangle_{\text{data}} - \langle h \rangle_{\text{model}})$ 

```

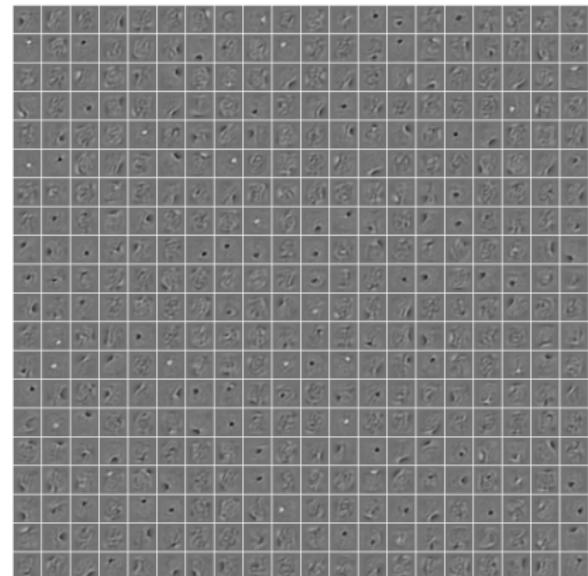


Training an RBM on handwritten digits

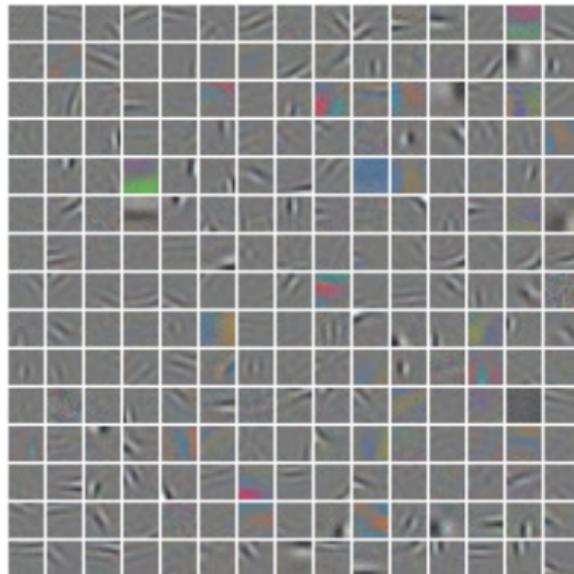
Without latent variables



With latent variables

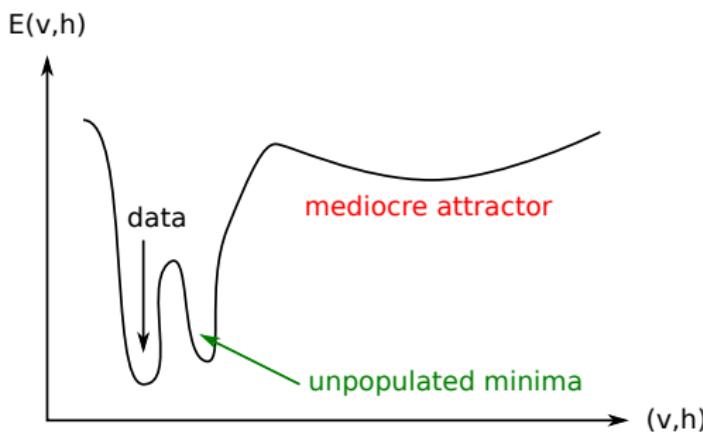


Training an RBM on natural image patches



- High resolution black and white filters
- Low-resolution color filters
- Similar to those observed in the monkey visual cortex

Obtaining unbiased samples from $p(\mathbf{v}, \mathbf{h})$



- Starting sampling from $(\mathbf{v}, \mathbf{h})_{\text{data}}$, it is difficult to move to other low-energy configurations of (\mathbf{v}, \mathbf{h})
- Starting sampling from $(\mathbf{v}, \mathbf{h}) \sim \text{Bernoulli}(p=0.5)$, it is difficult find any low-energy configuration of (\mathbf{v}, \mathbf{h})

Obtaining unbiased samples from $p(\mathbf{v}, \mathbf{h})$

Contrastive divergence work well for some applications, but this is a hack. How can we approximate well $\langle v_i h_j \rangle_{\text{model}}$?

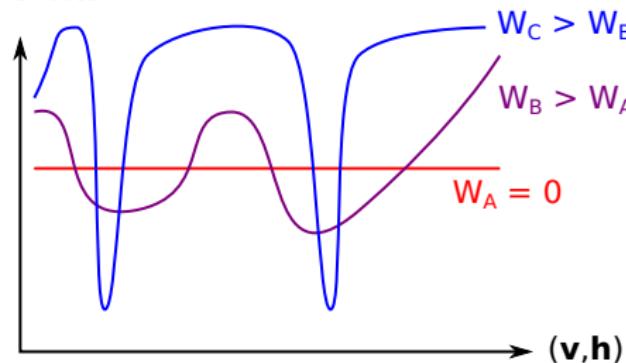
Simulated annealing

- 1 Start with a RBM with parameters $W, \mathbf{b}, \mathbf{a}$ set to zero.
- 2 Sample $(\mathbf{v}, \mathbf{h}) \sim \text{Bernouilli}(p = 0.5)$.
- 3 Update (\mathbf{v}, \mathbf{h}) by running k steps of Gibbs sampling while progressively raising the parameters $(W, \mathbf{b}, \mathbf{a})$ from zero to their original values.

If the process is slow enough, unbiased samples of $p(\mathbf{v}, \mathbf{h})$ are obtained.

Obtaining unbiased samples from $p(\mathbf{v}, \mathbf{h})$

$$E = -\mathbf{v}^T \mathbf{W} \mathbf{h}$$



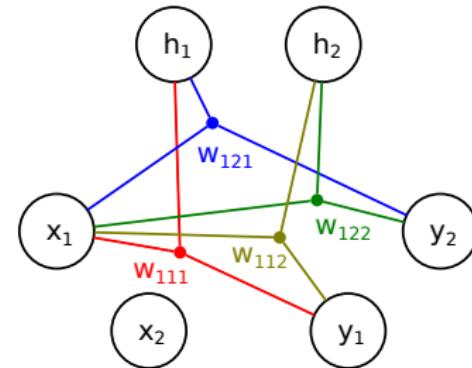
- When W is small, the RBM is “hot” and Gibbs sampling jumps randomly from one state to another.
- When W is large, the RBM is “cold” and Gibbs sampling can be stuck in the same state.
- By “cooling” progressively the RBM, we ensure that unbiased samples of $p(\mathbf{v}, \mathbf{h})$ are obtained.

Third-order restricted Boltzmann machines

- Visible units \mathbf{x}
- Visible units \mathbf{y}
- Hidden units \mathbf{h}
- Parameters $\theta = (W, \mathbf{a}, \mathbf{b}, \mathbf{c})$

Energy function is defined as

$$E(\mathbf{x}, \mathbf{y}, \mathbf{h}) = -\sum_i x_i a_i - \sum_j y_j b_j - \sum_k h_k c_k - \sum_{ijk} x_i y_j h_k w_{ijk}$$



Third-order restricted Boltzmann machines

Energy function is defined as

$$E(\mathbf{x}, \mathbf{y}, \mathbf{h}) = -\sum_i x_i a_i - \sum_j y_j b_j - \sum_k h_k b_k - \sum_{ijk} x_i y_j h_k w_{ijk}$$

Conditional probabilities can be computed as:

$$p(x_i | \mathbf{y}, \mathbf{h}) = \sigma\left(\sum_j y_j h_k w_{ijk} + a_i\right)$$

$$p(y_j | \mathbf{x}, \mathbf{h}) = \sigma\left(\sum_i x_i h_k w_{ijk} + b_j\right)$$

$$p(h_k | \mathbf{x}, \mathbf{y}) = \sigma\left(\sum_{ij} x_i y_j w_{ijk} + c_k\right)$$

To reduce capacity, W can be factored: $w_{ijk} = \sum_f w_{if}^x w_{jf}^y w_{kf}^h$

Third-order restricted Boltzmann machines

Key features:

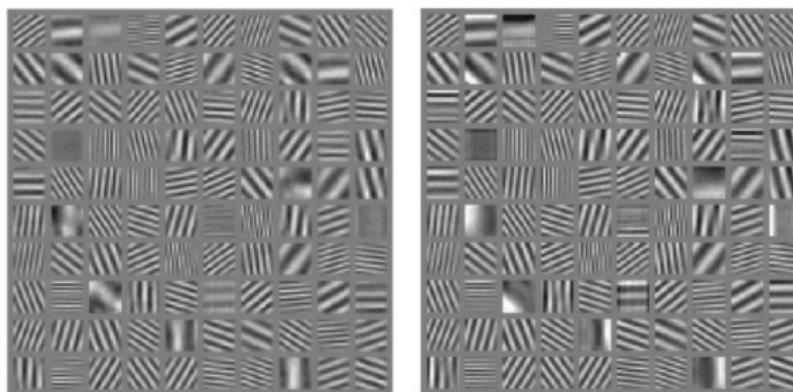
- This architecture models the interaction between x and y rather than x and y themselves.
- Hidden units h acts as a switch to determine whether x and y communicate.

Typical applications:

- Modeling system dynamics: ($x = s(t)$ and $y = s(t + 1)$)
- Modeling covariance matrices: ($x = v$ and $y = v$)
- Transducers

Third-order restricted Boltzmann machines

Modeling image transformations (Memisevic & Hinton, 2010)



When trained on image translations, the filters learned by RBM resemble the Fourier transform.

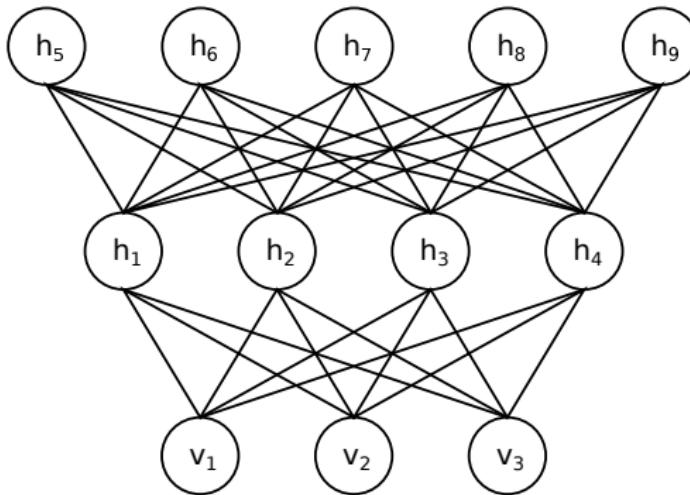
Third-order restricted Boltzmann machines

Modeling image transformations (Memisevic & Hinton, 2010)



Even if the RBM hasn't been trained on the faces, it can still rotate them.

Deep Boltzmann machines



- Theoretically higher representational power as RBM
- Difficult to optimize

Strategy:

- Initialize parameters with a greedy layer-wise pretraining

Does deep learning work?

MNIST (digit classification) - permutation invariant task

SVM	1.4 % error
Deep architecture	0.95% error

NORB (object categorization) - permutation invariant task

SVM	11.6 % error
Deep architecture	10.8% error

TIMIT (phone recognition)

GMM	33 % error
Deep architecture	20.5% error

Selected readings

Tutorial on RBMs

G. Hinton, A practical guide to training restricted Boltzmann machines, TR, 2010

A review on deep learning

Y. Bengio, Learning deep architectures for AI, Foundations and Trends in Machine Learning, 2009