

10417/10617

Intermediate Deep Learning: Fall2019

Russ Salakhutdinov

Machine Learning Department
rsalakhu@cs.cmu.edu

<https://deeplearning-cmu-10417.github.io/>

Deep Belief Networks

Neural Networks Online Course

- **Disclaimer:** Much of the material and slides for this lecture were borrowed from Hugo Larochelle's class on Neural Networks:
<https://sites.google.com/site/deeplearningsummerschool2016/>

http://info.usherbrooke.ca/hlarochelle/neural_networks

- Hugo's class covers many other topics: convolutional networks, neural language model, Boltzmann machines, autoencoders, sparse coding, etc.

- We will use his material for some of the other lectures.

Click with the mouse or tablet to draw with pen 2

RESTRICTED BOLTZMANN MACHINE

Topics: RBM, visible layer, hidden layer, energy function

The diagram illustrates the architecture of a Restricted Boltzmann Machine (RBM). It consists of two layers of binary units: a hidden layer (h) and a visible layer (x). The hidden layer units are represented by circles with bias terms b_j and are connected to the visible layer units by weights W . The visible layer units are represented by circles with bias terms c_k . The connections between the layers are labeled W and \mathbf{W} . The bias terms are labeled b_j and c_k . The visible layer units are labeled \mathbf{x} and the hidden layer units are labeled \mathbf{h} .

Energy function:
$$E(\mathbf{x}, \mathbf{h}) = -\mathbf{h}^\top \mathbf{W} \mathbf{x} - \mathbf{c}^\top \mathbf{x} - \mathbf{b}^\top \mathbf{h}$$
$$= -\sum_j \sum_k W_{j,k} h_j x_k - \sum_k c_k x_k - \sum_j b_j h_j$$

Distribution: $p(\mathbf{x}, \mathbf{h}) = \exp(-E(\mathbf{x}, \mathbf{h})) / Z$ ← partition function (intractable)

A small video inset in the bottom right corner shows a person with glasses and a beard, wearing a red and blue shirt, speaking.

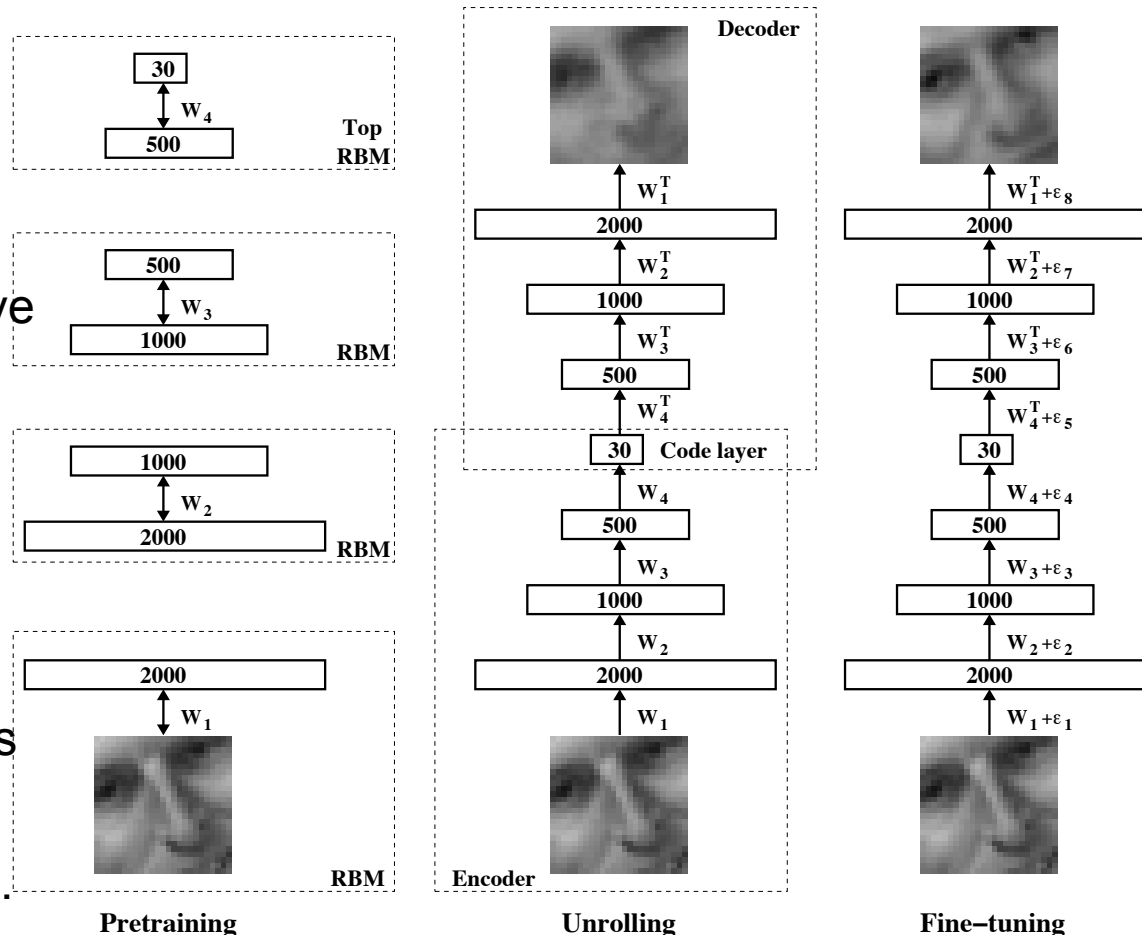
Deep Autoencoder

- Pre-training can be used to initialize a deep autoencoder

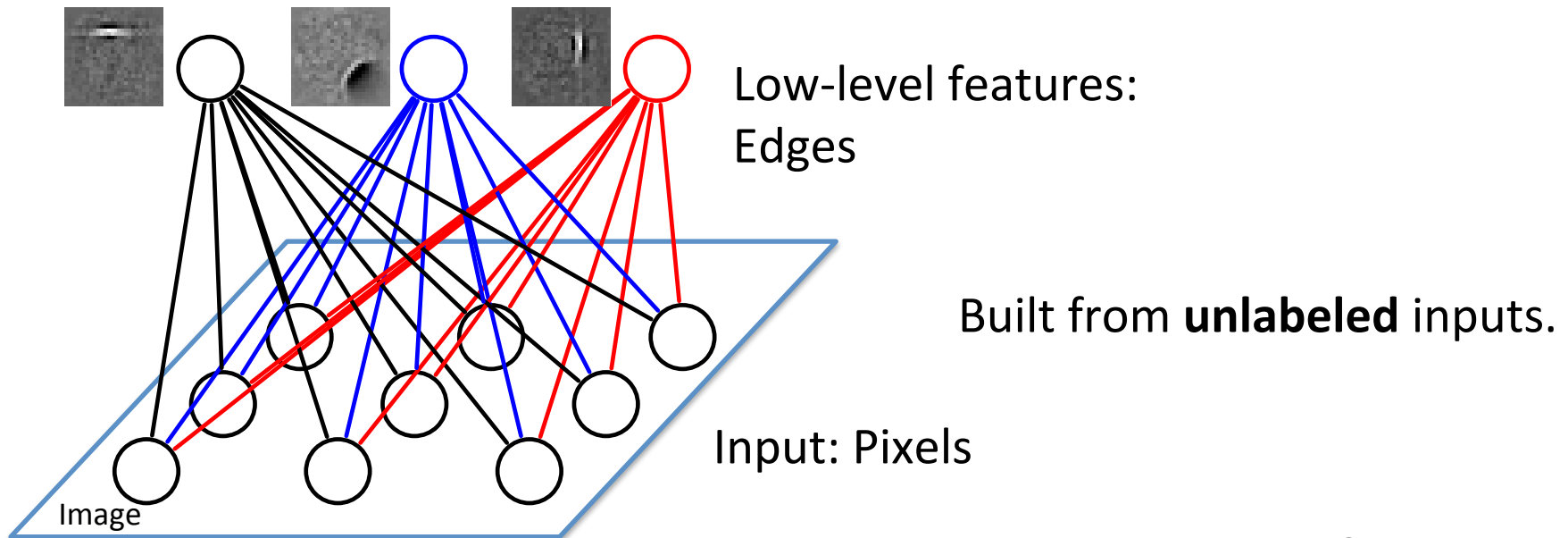
- Pre-training initializes the optimization problem in a region with better local optima of the training objective

- Each RBM used to initialize parameters both in encoder and decoder (“unrolling”)

- Better optimization algorithms can also help: Deep learning via Hessian-free optimization. Martens, 2010



Deep Belief Network



(Hinton et.al. Neural Computation 2006)

Deep Belief Network

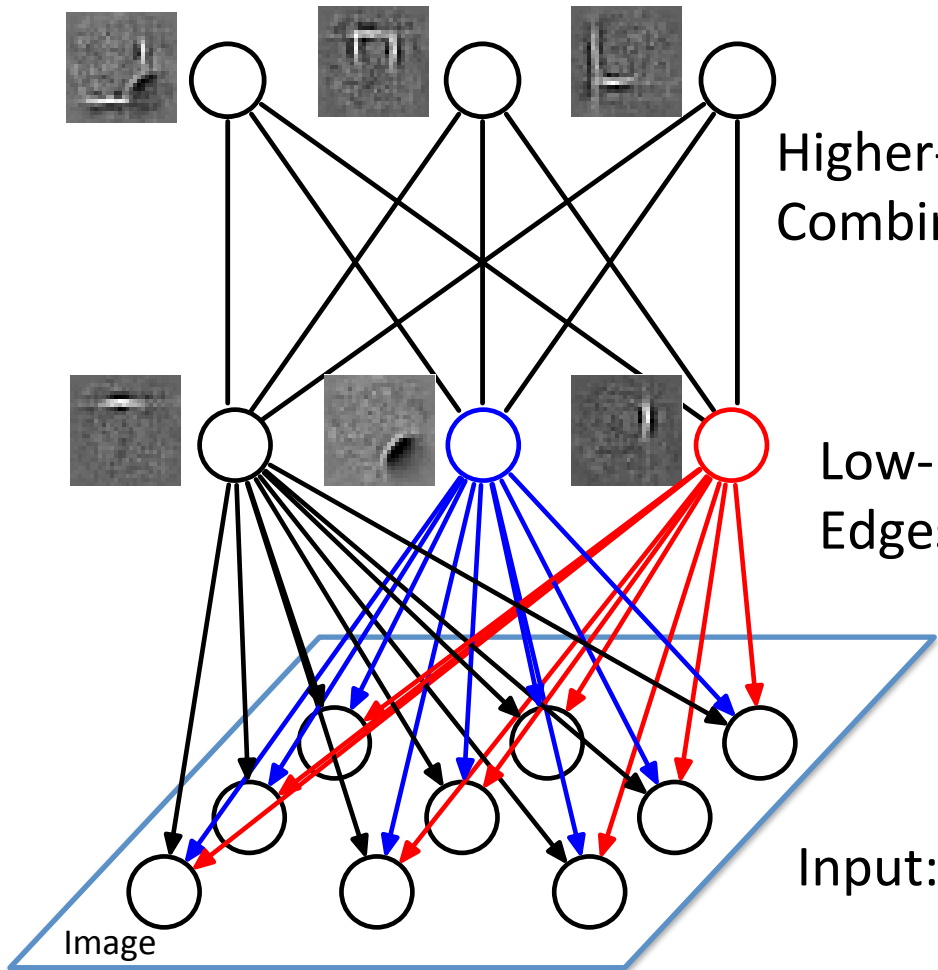
Internal representations capture higher-order statistical structure

Higher-level features:
Combination of edges

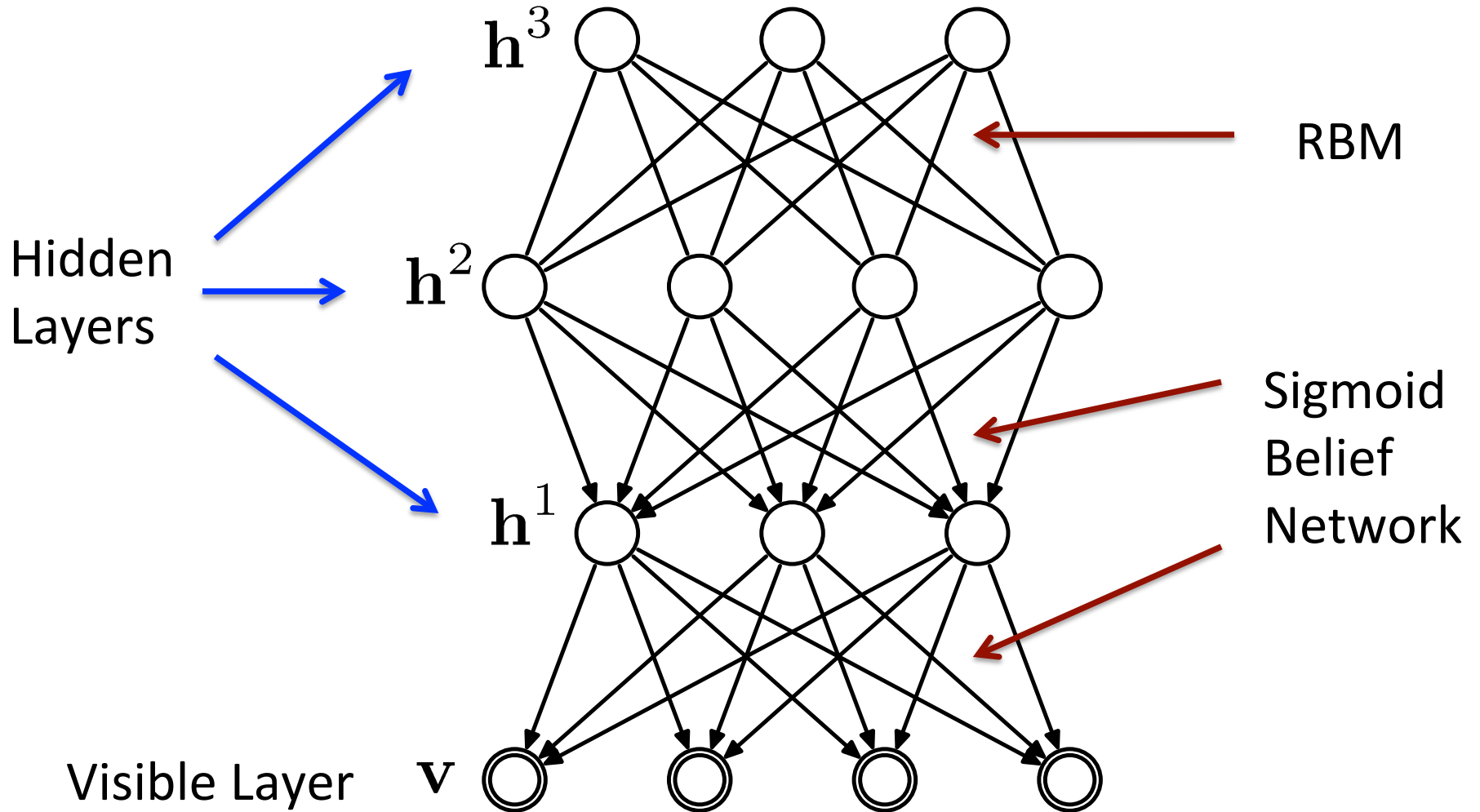
Low-level features:
Edges

Built from **unlabeled** inputs.

Input: Pixels



Deep Belief Network



Deep Belief Network

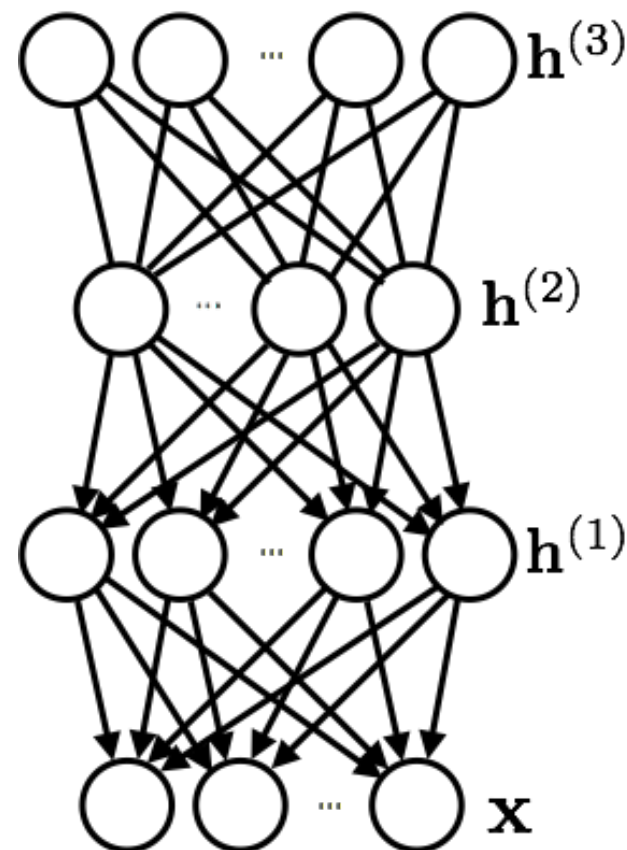
- Deep Belief Networks:

- it is a **generative model** that mixes undirected and directed connections between variables
- top 2 layers' distribution $p(\mathbf{h}^{(2)}, \mathbf{h}^{(3)})$ is an RBM!
- other layers form a **Bayesian network** with conditional distributions:

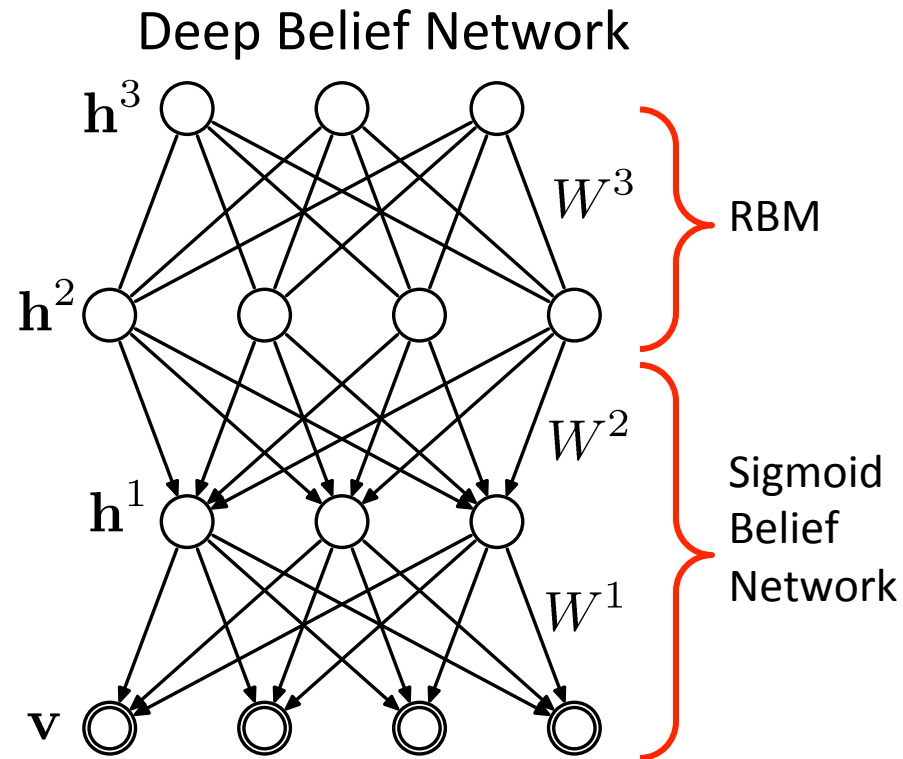
$$p(h_j^{(1)} = 1 | \mathbf{h}^{(2)}) = \text{sigm}(\mathbf{b}^{(1)} + \mathbf{W}^{(2)\top} \mathbf{h}^{(2)})$$

$$p(x_i = 1 | \mathbf{h}^{(1)}) = \text{sigm}(\mathbf{b}^{(0)} + \mathbf{W}^{(1)\top} \mathbf{h}^{(1)})$$

- This is **not a feed-forward** neural network



Deep Belief Network



➤ top 2 layers' distribution $p(\mathbf{h}^{(2)}, \mathbf{h}^{(3)})$ is an RBM

➤ other layers form a **Bayesian network** with conditional distributions:

$$p(h_j^{(1)} = 1 | \mathbf{h}^{(2)}) = \text{sigm}(\mathbf{b}^{(1)} + \mathbf{W}^{(2)\top} \mathbf{h}^{(2)})$$

$$p(x_i = 1 | \mathbf{h}^{(1)}) = \text{sigm}(\mathbf{b}^{(0)} + \mathbf{W}^{(1)\top} \mathbf{h}^{(1)})$$

Deep Belief Network

- The **joint distribution** of a DBN is as follows

$$p(\mathbf{x}, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}) = p(\mathbf{h}^{(2)}, \mathbf{h}^{(3)}) p(\mathbf{h}^{(1)} | \mathbf{h}^{(2)}) p(\mathbf{x} | \mathbf{h}^{(1)})$$

where

$$p(\mathbf{h}^{(2)}, \mathbf{h}^{(3)}) = \exp \left(\mathbf{h}^{(2)\top} \mathbf{W}^{(3)} \mathbf{h}^{(3)} + \mathbf{b}^{(2)\top} \mathbf{h}^{(2)} + \mathbf{b}^{(3)\top} \mathbf{h}^{(3)} \right) / Z$$

$$p(\mathbf{h}^{(1)} | \mathbf{h}^{(2)}) = \prod_j p(h_j^{(1)} | \mathbf{h}^{(2)})$$

$$p(\mathbf{x} | \mathbf{h}^{(1)}) = \prod_i p(x_i | \mathbf{h}^{(1)})$$

- As in a deep feed-forward network, **training a DBN is hard**

Layer-wise Pretraining

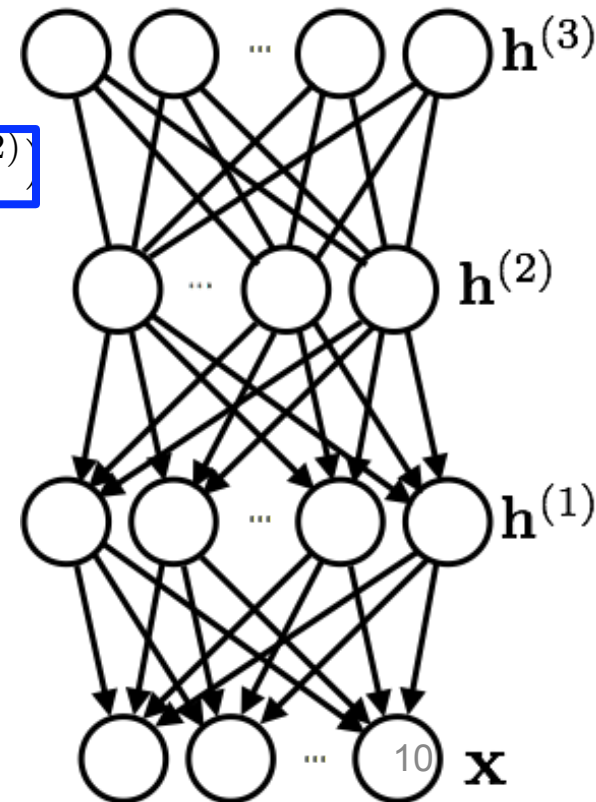
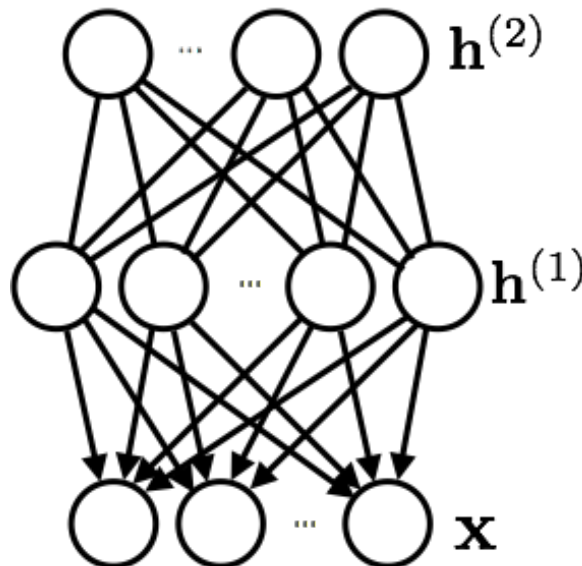
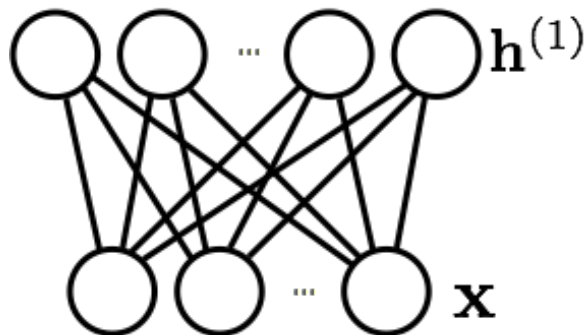
- This is where the RBM stacking procedure comes from:

- **idea:** improve prior on last layer by adding another hidden layer

$$p(\mathbf{h}^{(1)}, \mathbf{h}^{(2)}) = p(\mathbf{h}^{(1)} | \mathbf{h}^{(2)}) \sum_{\mathbf{h}^{(3)}} p(\mathbf{h}^{(2)}, \mathbf{h}^{(3)})$$

$$p(\mathbf{x}, \mathbf{h}^{(1)}) = p(\mathbf{x} | \mathbf{h}^{(1)}) \sum_{\mathbf{h}^{(2)}} p(\mathbf{h}^{(1)}, \mathbf{h}^{(2)})$$

$$p(\mathbf{x}) = \sum_{\mathbf{h}^{(1)}} p(\mathbf{x}, \mathbf{h}^{(1)})$$

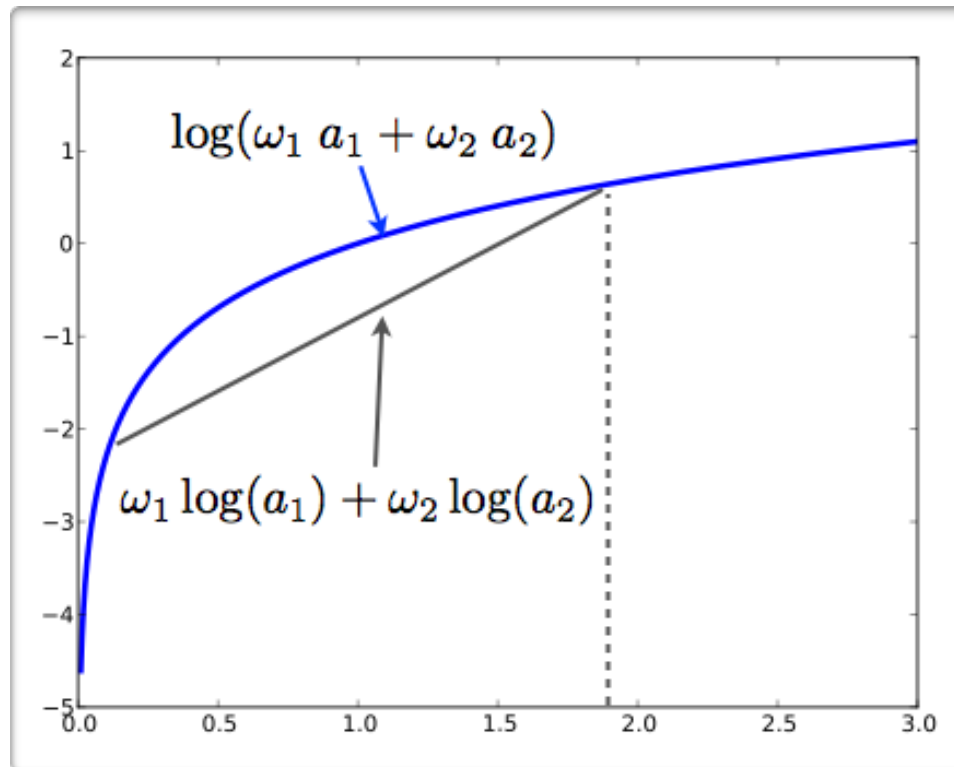


Concavity

- We will use the fact that the logarithm function is **concave**:

$$\log\left(\sum_i \omega_i a_i\right) \geq \sum_i \omega_i \log(a_i)$$

(where $\sum_i \omega_i = 1$ and $\omega_i \geq 0$)



Variational Bound

- For any model $p(\mathbf{x}, \mathbf{h}^{(1)})$ with latent variables $\mathbf{h}^{(1)}$ we can write:

$$\begin{aligned}\log p(\mathbf{x}) &= \log \left(\sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)} | \mathbf{x}) \frac{p(\mathbf{x}, \mathbf{h}^{(1)})}{q(\mathbf{h}^{(1)} | \mathbf{x})} \right) \\ &\geq \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)} | \mathbf{x}) \log \left(\frac{p(\mathbf{x}, \mathbf{h}^{(1)})}{q(\mathbf{h}^{(1)} | \mathbf{x})} \right) \\ &= \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)} | \mathbf{x}) \log p(\mathbf{x}, \mathbf{h}^{(1)}) \\ &\quad - \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)} | \mathbf{x}) \log q(\mathbf{h}^{(1)} | \mathbf{x})\end{aligned}$$

where $q(\mathbf{h}^{(1)} | \mathbf{x})$ is any **approximation** to $p(\mathbf{h}^{(1)} | \mathbf{x})$

Variational Bound

- This is called a **variational bound**

$$\begin{aligned}\log p(\mathbf{x}) &\geq \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{h}^{(1)}) \\ &\quad - \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \log q(\mathbf{h}^{(1)}|\mathbf{x})\end{aligned}$$

- if $q(\mathbf{h}^{(1)}|\mathbf{x})$ is equal to the true conditional $p(\mathbf{h}^{(1)}|\mathbf{x})$, then we have an equality – **the bound is tight!**
- the more $q(\mathbf{h}^{(1)}|\mathbf{x})$ is different from $p(\mathbf{h}^{(1)}|\mathbf{x})$ the less tight the bound is.

Variational Bound

- This is called a variational bound

$$\begin{aligned}\log p(\mathbf{x}) &\geq \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{h}^{(1)}) \\ &\quad - \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \log q(\mathbf{h}^{(1)}|\mathbf{x})\end{aligned}$$

- In fact, difference between the left and right terms is the **KL divergence** between $q(\mathbf{h}^{(1)}|\mathbf{x})$ and $p(\mathbf{h}^{(1)}|\mathbf{x})$:

$$\text{KL}(q||p) = \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \log \left(\frac{q(\mathbf{h}^{(1)}|\mathbf{x})}{p(\mathbf{h}^{(1)}|\mathbf{x})} \right)$$

Variational Bound

- This is called a variational bound


$$\begin{aligned}\log p(\mathbf{x}) \quad \geq \quad & \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \left(\log p(\mathbf{x}|\mathbf{h}^{(1)}) + \log p(\mathbf{h}^{(1)}) \right) \\ & - \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \log q(\mathbf{h}^{(1)}|\mathbf{x})\end{aligned}$$

- for a single hidden layer DBN (i.e. an RBM), both **the likelihood** $p(\mathbf{x}|\mathbf{h}^{(1)})$ and **the prior** $p(\mathbf{h}^{(1)})$ depend on the parameters of the first layer.
- we can now improve the model by building a better prior $p(\mathbf{h}^{(1)})$

Variational Bound

- This is called a variational bound

adding 2nd layer means
untying the parameters



$$\log p(\mathbf{x}) \geq \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \left(\log p(\mathbf{x}|\mathbf{h}^{(1)}) + \log p(\mathbf{h}^{(1)}) \right) - \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \log q(\mathbf{h}^{(1)}|\mathbf{x})$$


- When adding a second layer, we model $p(\mathbf{h}^{(1)})$ using a separate set of parameters
 - they are the parameters of the RBM involving $\mathbf{h}^{(1)}$ and $\mathbf{h}^{(2)}$
 - $p(\mathbf{h}^{(1)})$ is now the marginalization of the second hidden layer

$$p(\mathbf{h}^{(1)}) = \sum_{\mathbf{h}^{(2)}} p(\mathbf{h}^{(1)}, \mathbf{h}^{(2)})$$

Variational Bound

- This is called a variational bound

adding 2nd layer means
untying the parameters



$$\log p(\mathbf{x}) \geq \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \left(\log p(\mathbf{x}|\mathbf{h}^{(1)}) + \log p(\mathbf{h}^{(1)}) \right) - \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \log q(\mathbf{h}^{(1)}|\mathbf{x})$$

- we can train the parameters of the bound. This is equivalent to training the other terms are constant:

$$- \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \log p(\mathbf{h}^{(1)})$$


Layerwise pretraining
improves variational
lower bound

- this is like training an RBM on data generated from $q(\mathbf{h}^{(1)}|\mathbf{x})$!

Variational Bound

- This is called a variational bound

adding 2nd layer means
untying the parameters



$$\log p(\mathbf{x}) \geq \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \left(\log p(\mathbf{x}|\mathbf{h}^{(1)}) + \log p(\mathbf{h}^{(1)}) \right) - \sum_{\mathbf{h}^{(1)}} q(\mathbf{h}^{(1)}|\mathbf{x}) \log q(\mathbf{h}^{(1)}|\mathbf{x})$$

- for $q(\mathbf{h}^{(1)}|\mathbf{x})$ we use **the posterior of the first layer RBM**. This is equivalent to a feed-forward (sigmoidal) layer, followed by sampling
- by initializing the weights of the second layer RBM as the transpose of the first layer weights, **the bound is initially tight!**
- a 2-layer DBN with tied weights is equivalent to a 1-layer RBM

Layer-wise Pretraining

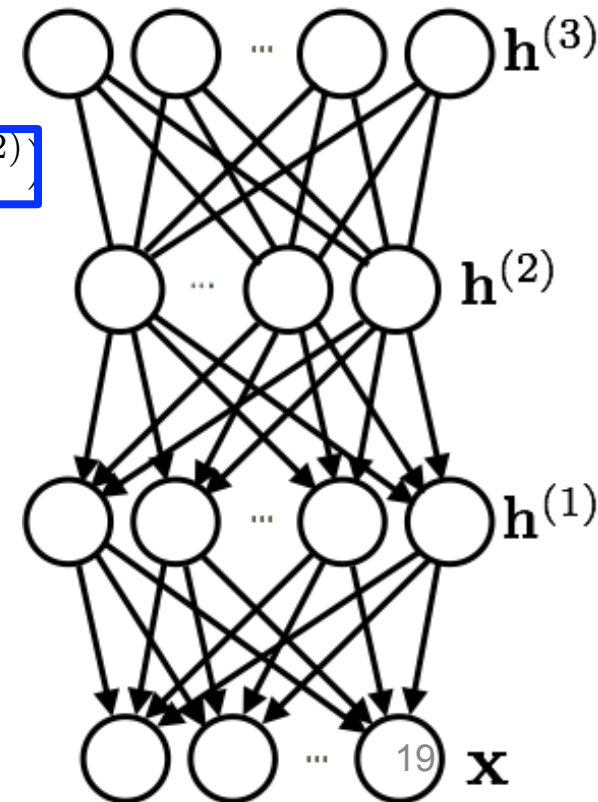
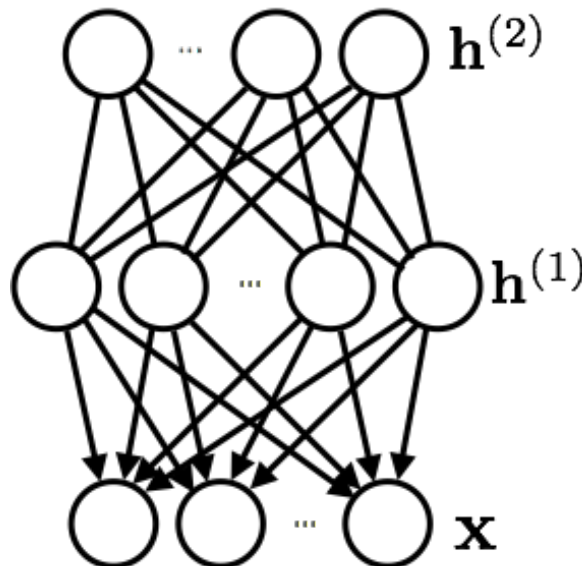
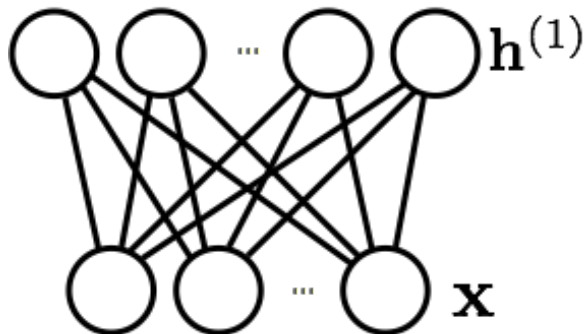
- This is where the RBM stacking procedure comes from:

- **idea:** improve prior on last layer by adding another hidden layer

$$p(\mathbf{h}^{(1)}, \mathbf{h}^{(2)}) = p(\mathbf{h}^{(1)} | \mathbf{h}^{(2)}) \sum_{\mathbf{h}^{(3)}} p(\mathbf{h}^{(2)}, \mathbf{h}^{(3)})$$

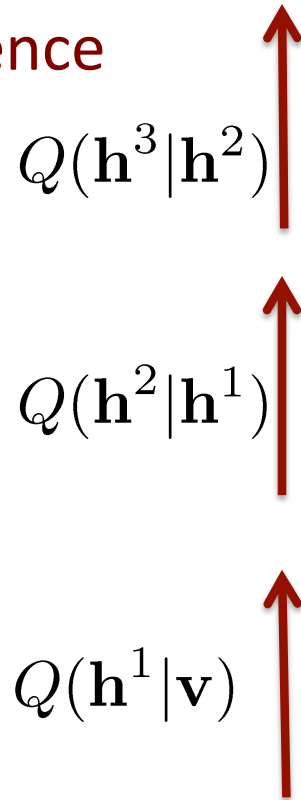
$$p(\mathbf{x}, \mathbf{h}^{(1)}) = p(\mathbf{x} | \mathbf{h}^{(1)}) \sum_{\mathbf{h}^{(2)}} p(\mathbf{h}^{(1)}, \mathbf{h}^{(2)})$$

$$p(\mathbf{x}) = \sum_{\mathbf{h}^{(1)}} p(\mathbf{x}, \mathbf{h}^{(1)})$$

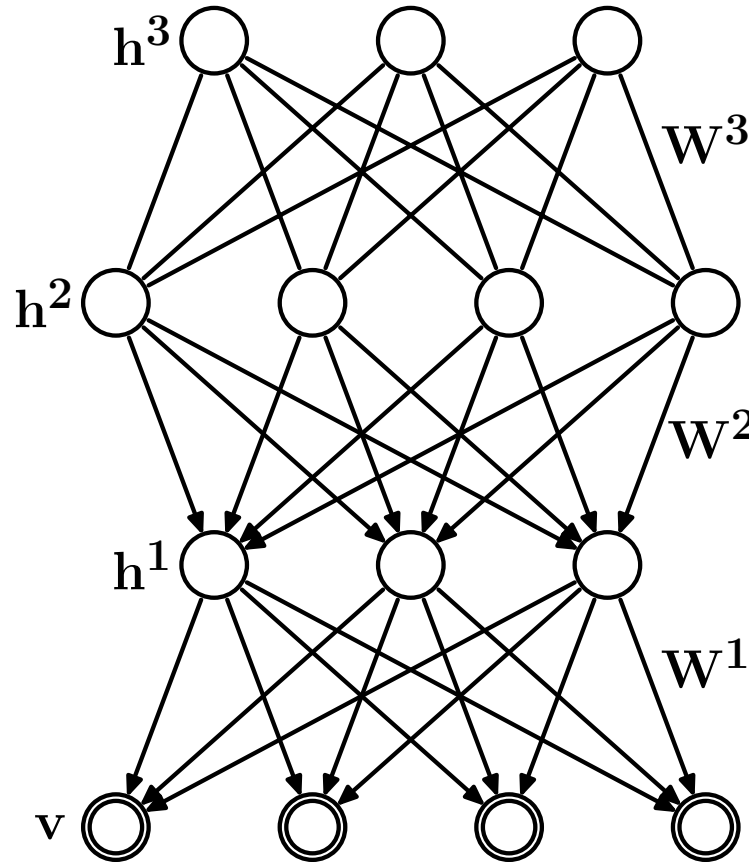
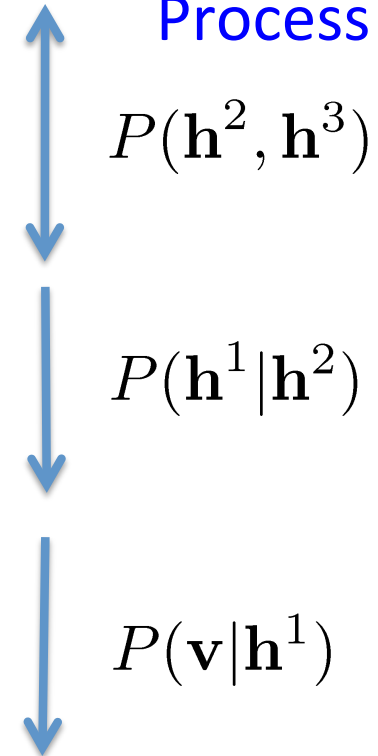


Deep Belief Network

Approximate
Inference



Generative
Process

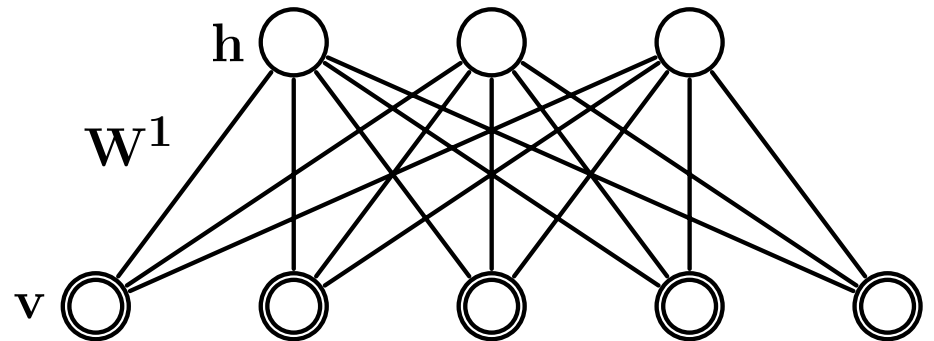


$$Q(\mathbf{h}^t | \mathbf{h}^{t-1}) = \prod_j \sigma \left(\sum_i W^t h_i^{t-1} \right)$$

$$P(\mathbf{h}^{t-1} | \mathbf{h}^t) = \prod_j \sigma \left(\sum_i W^t h_i^t \right)$$

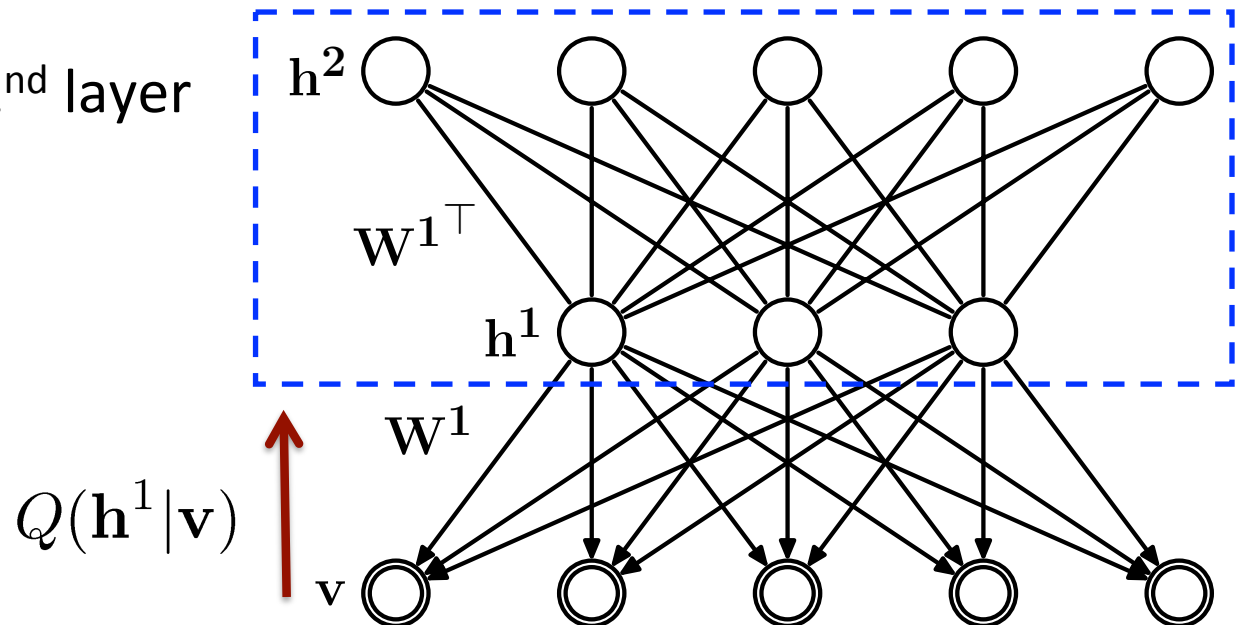
DBN Layer-wise Training

- Learn an RBM with an input layer $v=x$ and a hidden layer h .



DBN Layer-wise Training

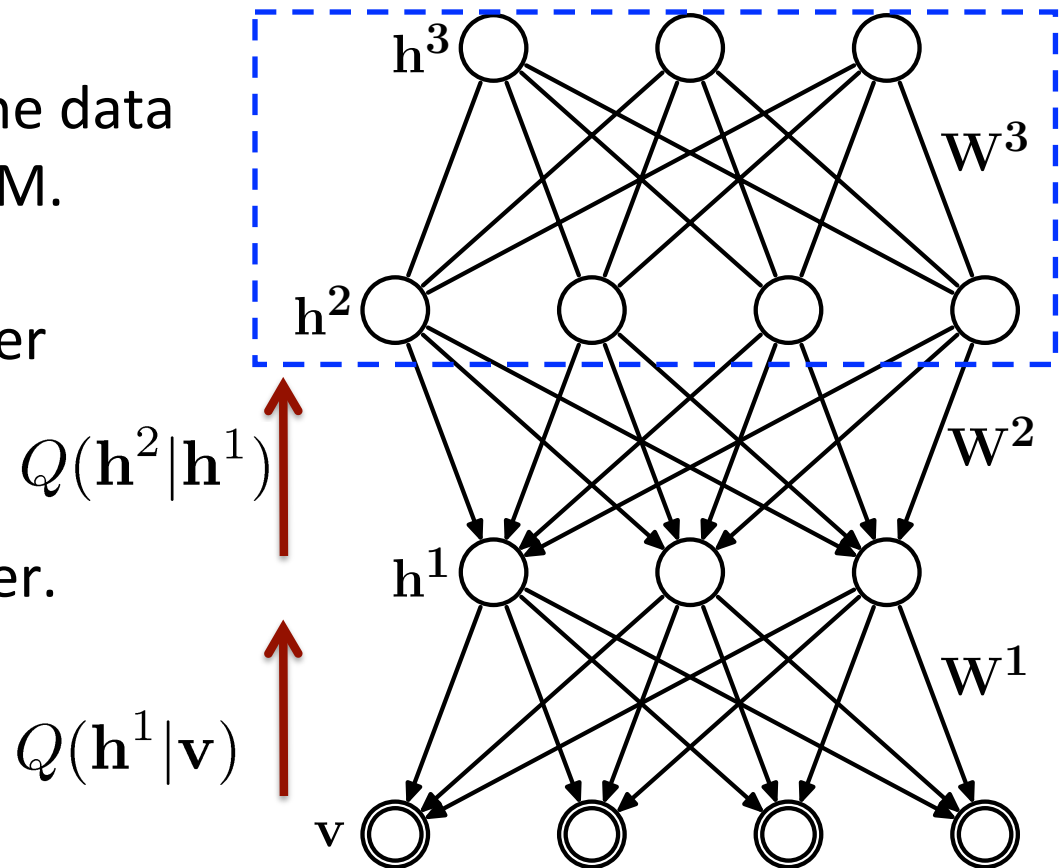
- Learn an RBM with an input layer $\mathbf{v}=\mathbf{x}$ and a hidden layer \mathbf{h} .
- Treat inferred values $Q(\mathbf{h}^1|\mathbf{v}) = P(\mathbf{h}^1|\mathbf{v})$ as the data for training 2nd-layer RBM.
- Learn and freeze 2nd layer RBM.



DBN Layer-wise Training

- Learn an RBM with an input layer $v=x$ and a hidden layer h .
- Treat inferred values $Q(h^1|v) = P(h^1|v)$ as the data for training 2nd-layer RBM.
- Learn and freeze 2nd layer RBM.
- Proceed to the next layer.

Unsupervised Feature Learning.



DBN Layer-wise Training

- Learn an RBM with an input layer $\mathbf{v}=\mathbf{x}$ and a hidden layer \mathbf{h} .

Unsupervised Feature Learning.

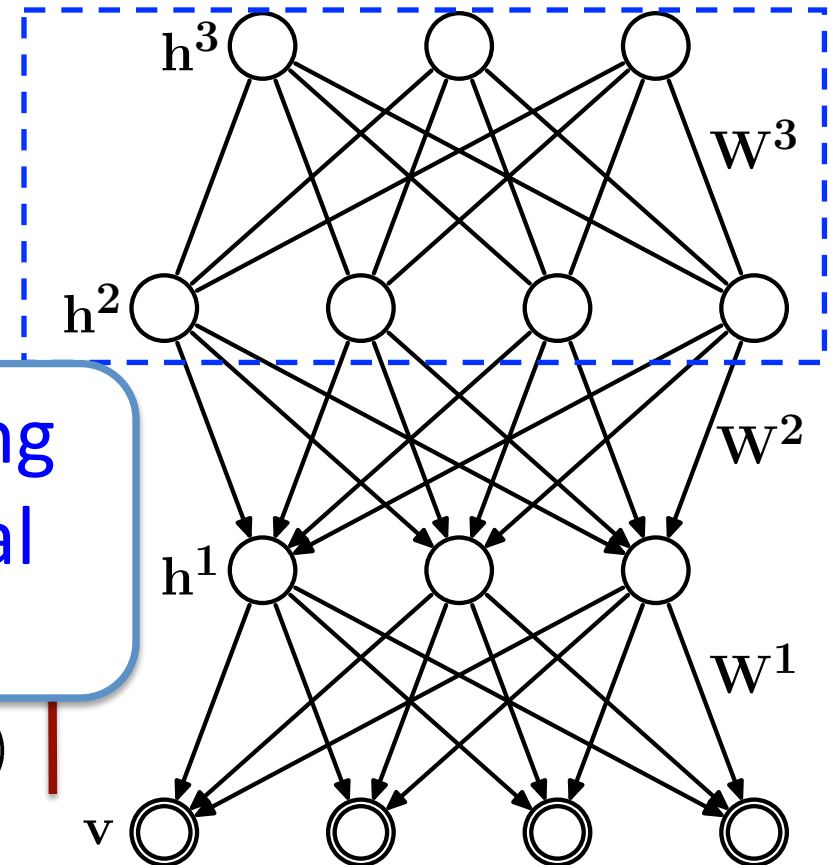
- Treat inferred values $Q(\mathbf{h}^1|\mathbf{v}) = P(\mathbf{h}^1|\mathbf{v})$ as the data for training 2nd-layer RBM.

- Learn and freeze 2nd layer RBM

- Proc

Layerwise pretraining improves variational lower bound

$$Q(\mathbf{h}^1|\mathbf{v})$$



Deep Belief Networks

- This process of adding layers can be repeated recursively
 - we obtain **the greedy layer-wise pre-training** procedure for neural networks
- We now see that this procedure corresponds to **maximizing a bound on the likelihood of the data** in a DBN
 - in theory, if our approximation $q(\mathbf{h}^{(1)}|\mathbf{x})$ is very far from the true posterior, the bound might be very loose
 - this only means we might not be improving the true likelihood
 - we might still be extracting better features!
- Fine-tuning is done by the Up-Down algorithm
 - A fast learning algorithm for deep belief nets. Hinton, Teh, Osindero, 2006.

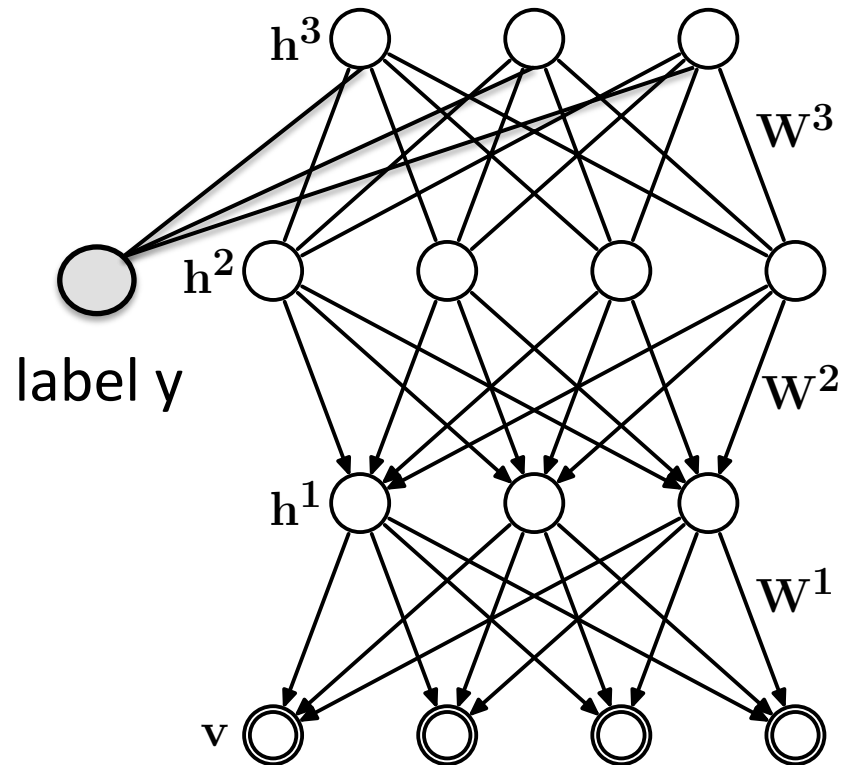
Supervised Learning with DBNs

- If we have access to label information, we can train **the joint generative model** by maximizing the joint log-likelihood of data and labels

$$\log P(\mathbf{y}, \mathbf{v})$$

- Discriminative fine-tuning:
 - Use DBN to initialize a multilayer neural network.
 - Maximize **the conditional distribution**:

$$\log P(\mathbf{y}|\mathbf{v})$$

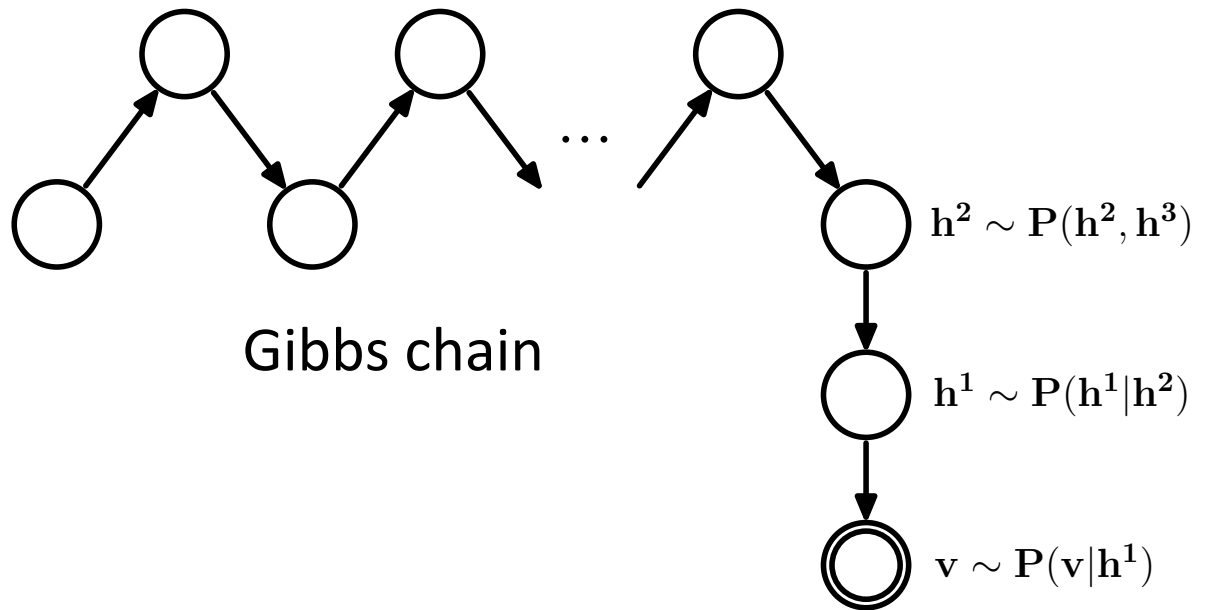
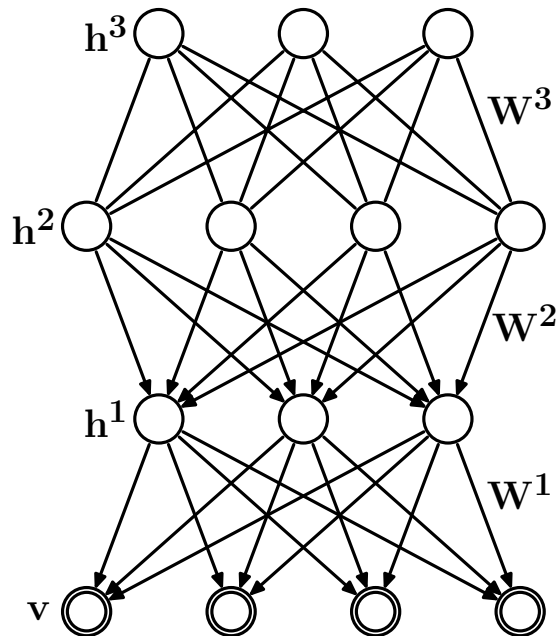


Sampling from DBNs

- To sample from the DBN model:

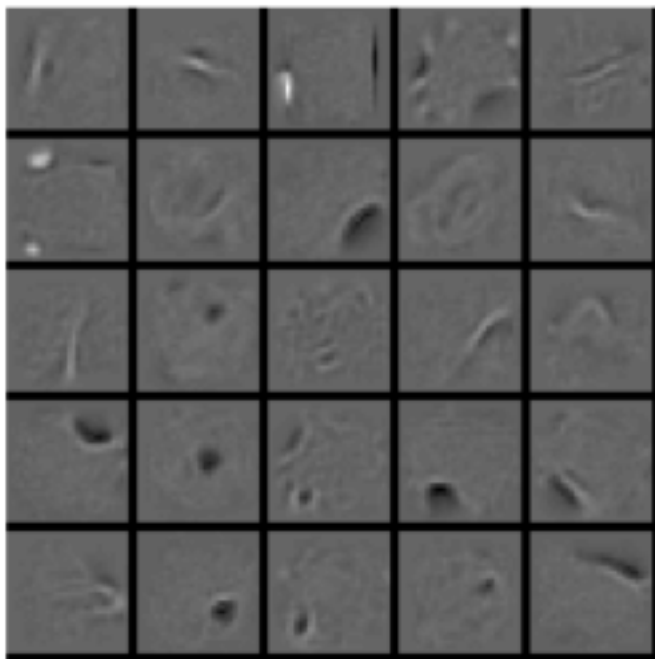
$$P(\mathbf{v}, \mathbf{h}^1, \mathbf{h}^2, \mathbf{h}^3) = P(\mathbf{v}|\mathbf{h}^1)P(\mathbf{h}^1|\mathbf{h}^2)P(\mathbf{h}^2, \mathbf{h}^3)$$

- Sample \mathbf{h}^2 using alternating Gibbs sampling from RBM.
- Sample lower layers using sigmoid belief network.

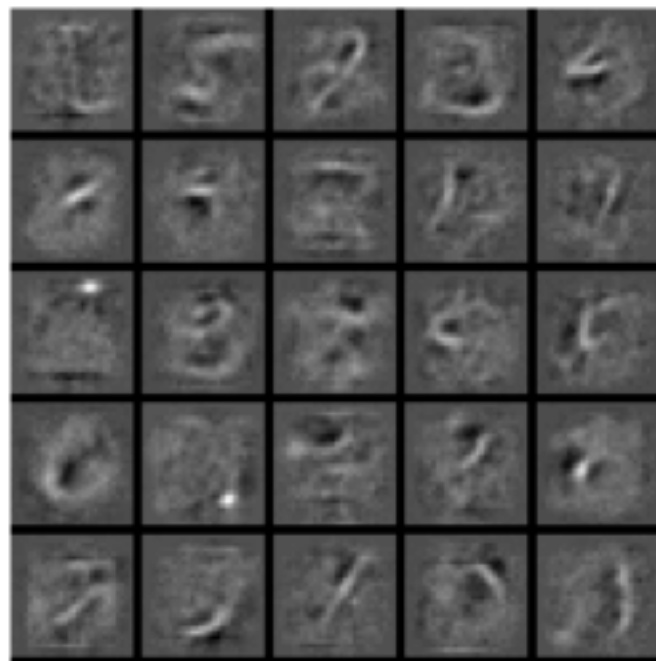


Learned Features

1st-layer features

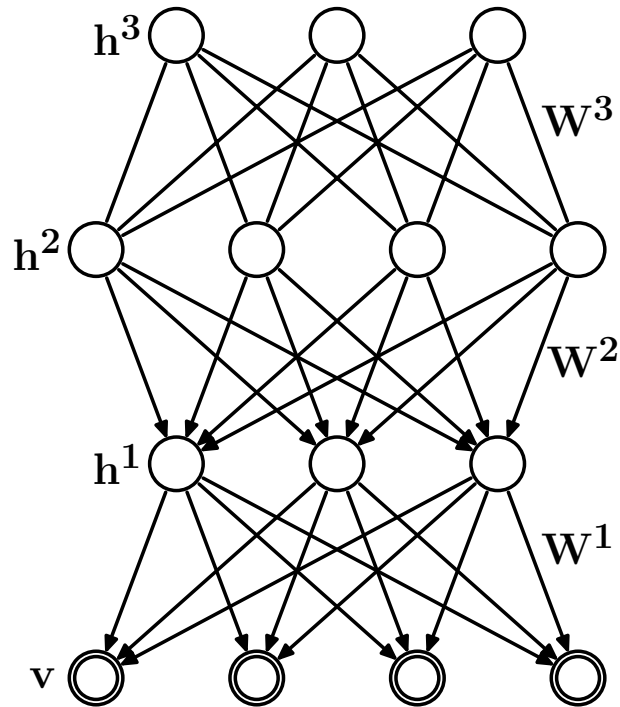


2nd-layer features

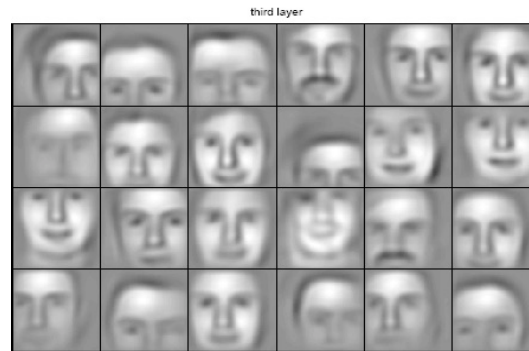


Learning Part-based Representation

Convolutional DBN



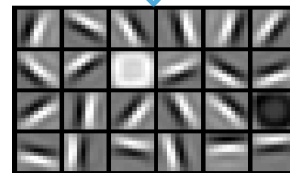
Faces



Groups of parts.



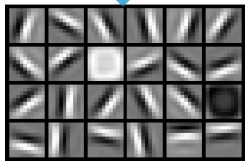
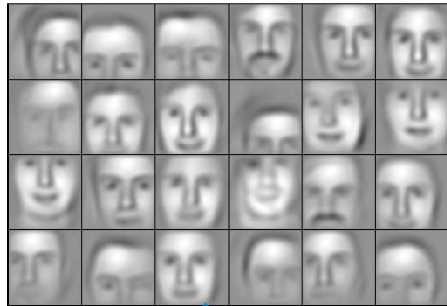
Object Parts



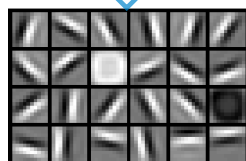
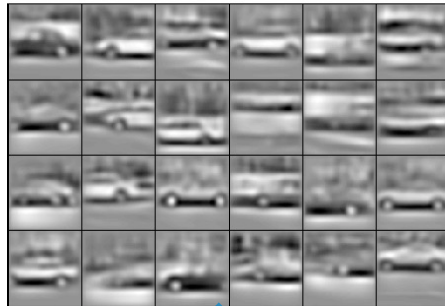
Trained on face images.

Learning Part-based Representation

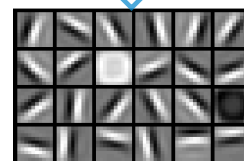
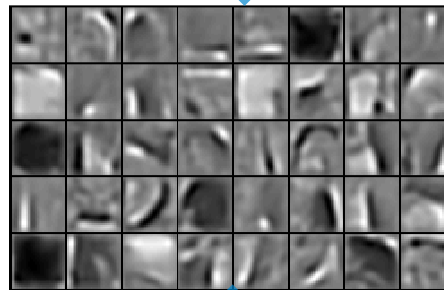
Faces



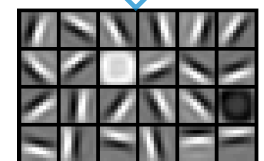
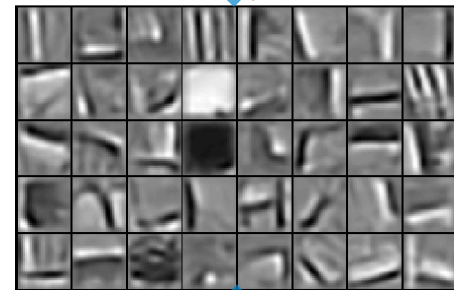
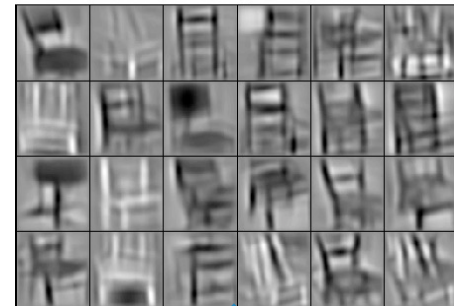
Cars



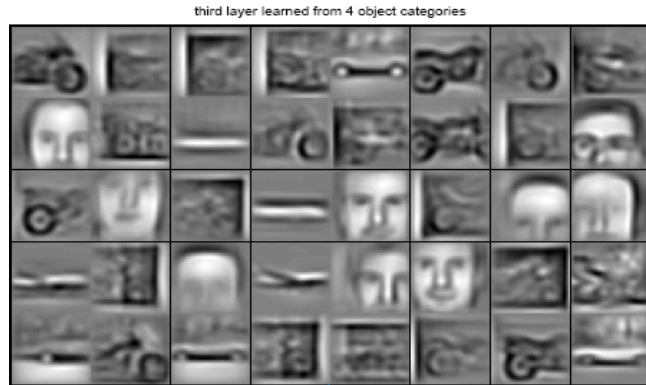
Elephants



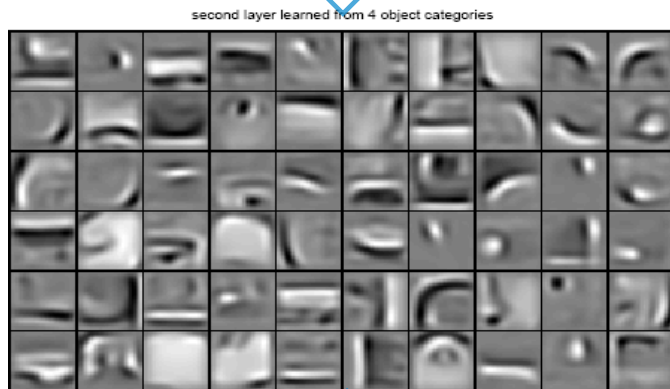
Chairs



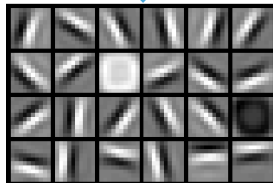
Learning Part-based Representation



Groups of parts.

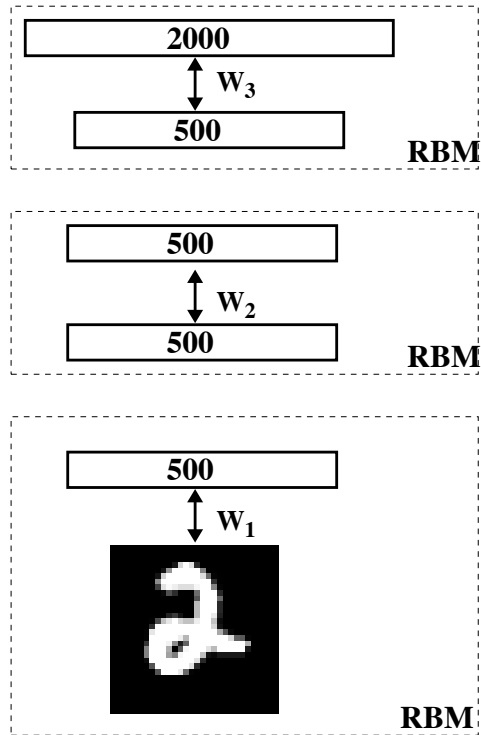


Class-specific object parts

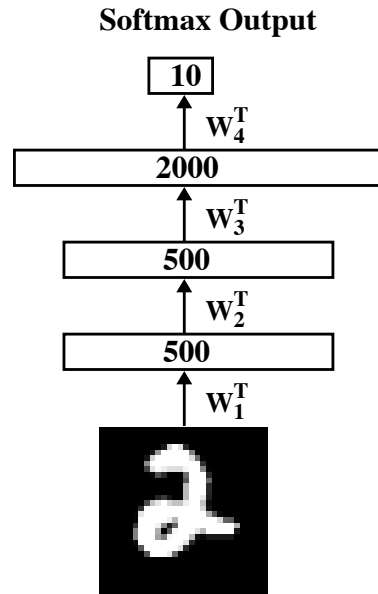


Trained from multiple classes (cars, faces, motorbikes, airplanes).

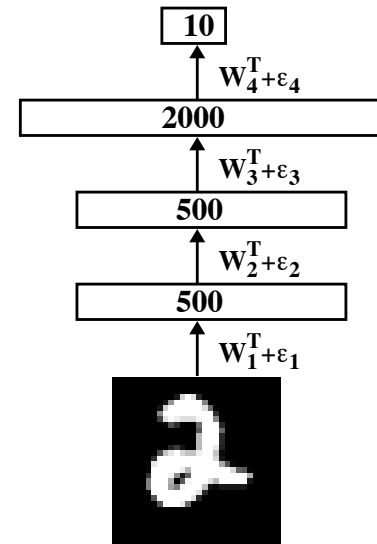
DBNs for Classification



Pretraining



Unrolling



Fine-tuning

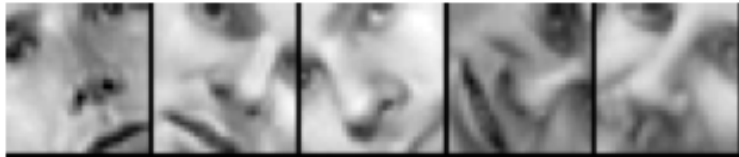
- After layer-by-layer **unsupervised pretraining**, discriminative fine-tuning by backpropagation achieves an error rate of 1.2% on MNIST. SVM's get 1.4% and randomly initialized backprop gets 1.6%.
- Clearly unsupervised learning helps generalization. It ensures that most of the information in the weights comes from modeling the input data.

DBNs for Regression

Predicting the orientation of a face patch

Training Data

-22.07 32.99 -41.15 66.38 27.49



Test Data



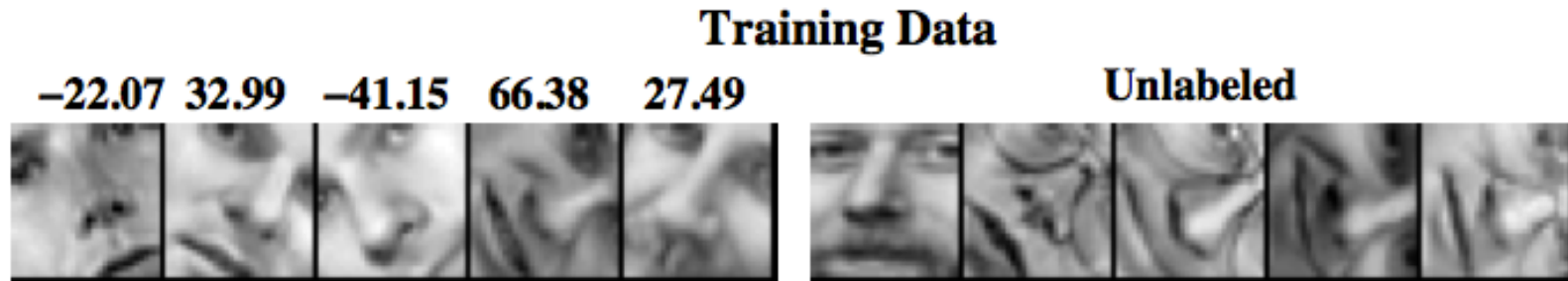
Training Data: 1000 face patches of 30 training people.

Test Data: 1000 face patches of 10 new people.

Regression Task: predict orientation of a new face.

Gaussian Processes with spherical Gaussian kernel achieves a RMSE (root mean squared error) of 16.33 degree.

DBNs for Regression

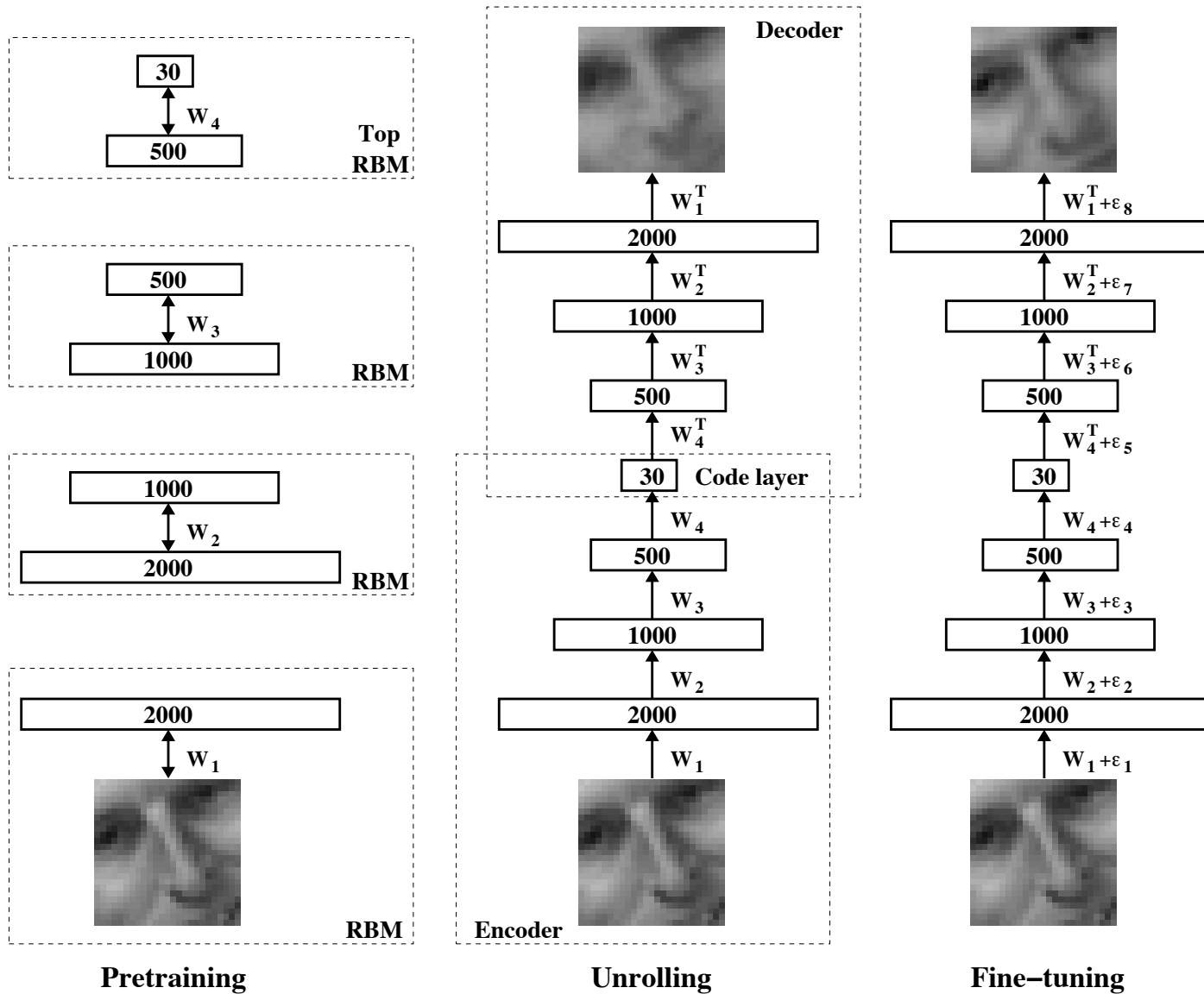


Additional Unlabeled Training Data: 12000 face patches from 30 training people.

- Pretrain a stack of RBMs: 784-1000-1000-1000.
- **Features were extracted with no idea of the final task.**

The same GP on the top-level features:	RMSE: 11.22
GP with fine-tuned covariance Gaussian kernel:	RMSE: 6.42
Standard GP without using DBNs:	RMSE: 16.33

Deep Autoencoders



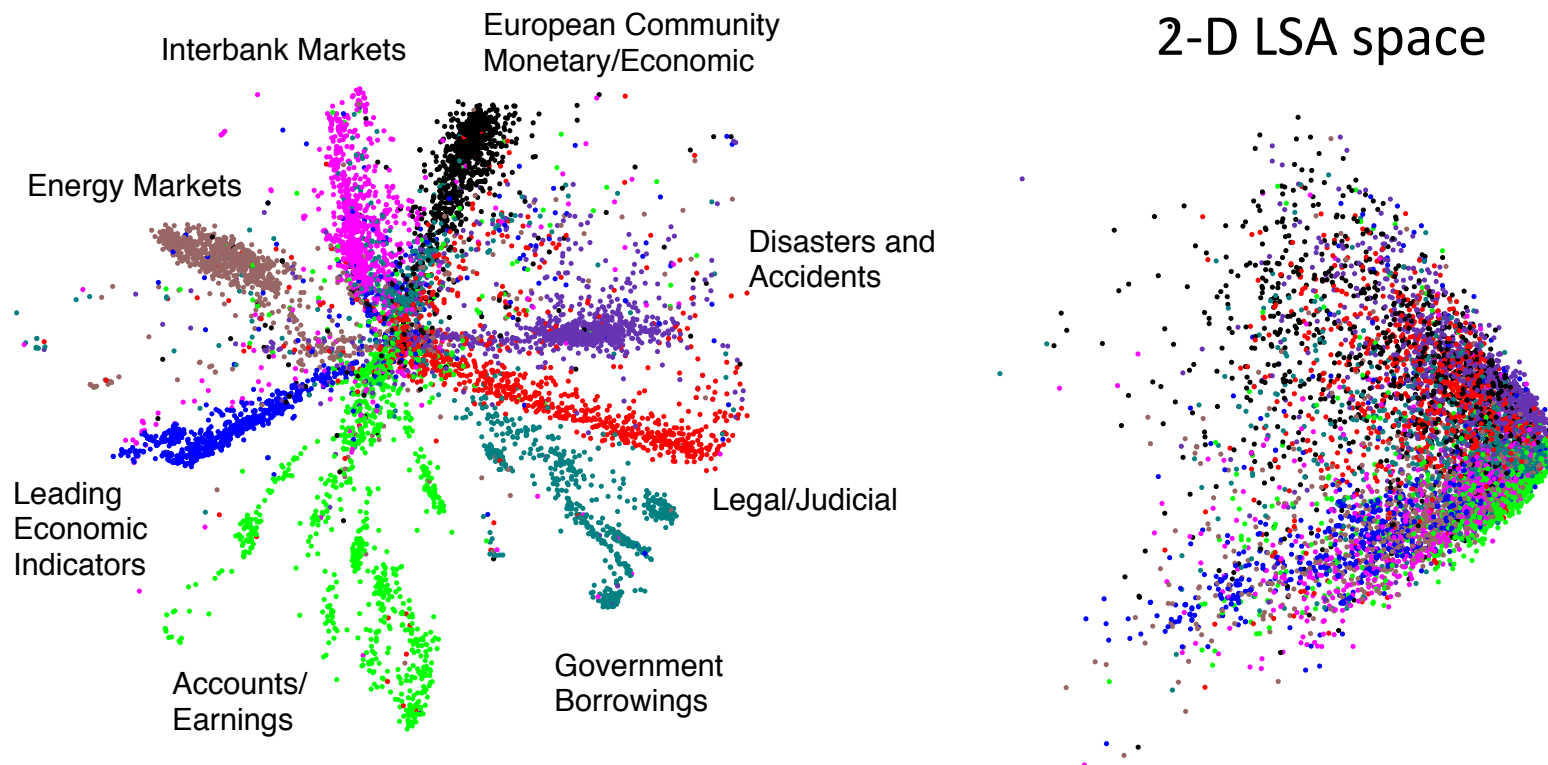
Deep Autoencoders

- We used $25 \times 25 - 2000 - 1000 - 500 - 30$ autoencoder to extract 30-D real-valued codes for Olivetti face patches.



- **Top:** Random samples from the test dataset.
- **Middle:** Reconstructions by the 30-dimensional deep autoencoder.
- **Bottom:** Reconstructions by the 30-dimensional PCA.

Information Retrieval



- The Reuters Corpus Volume II contains 804,414 newswire stories (randomly split into **402,207 training** and **402,207 test**).
- “Bag-of-words” representation: each article is represented as a vector containing the counts of the most frequently used 2000 words in the training set.

(Hinton and Salakhutdinov, Science 2006)