



CRIPAC

智能感知与计算研究中心
Center for Research on Intelligent Perception and Computing



中国科学院自动化研究所
Institute of Automation
Chinese Academy of Sciences

2018

“Deep Learning Lecture”

Lecture 7: Generative Model (1)

Liang Wang

Center for Research on Intelligent Perception and Computing (CRIPAC)
National Laboratory of Pattern Recognition (NLPR)
Institute of Automation, Chinese Academy of Science (CASIA)

Outline

1 Course Review

2 Linear Factor Model

3 Autoencoder

4 DBN and RBM

Review: Regularization Strategies

- Parameter Norm Penalties
- Dataset Augmentation
- Noise Robustness
- Early Stopping
- Parameter Tying and Parameter Sharing
- Multitask Learning
- Bagging and Other Ensemble Methods
- Dropout
- Adversarial Training

Review: Optimization

- Things we have looked at last week
 - Stochastic Gradient Descent
 - Momentum Method and the Nesterov Variant
 - Adaptive Learning Methods (AdaGrad, RMSProp, Adam)
 - Batch Normalization
 - Initialization Heuristics

Outline

1 Course Review

2 Linear Factor Model

3 Autoencoder

4 DBN and RBM

Linear Factor Model

- We want to build a probabilistic model of the input $\tilde{P}(\mathbf{x})$
- Like before, we are interested in latent factors \mathbf{h} that explain \mathbf{x}
- We then care about the marginal:

$$\tilde{P}(\mathbf{x}) = \mathbb{E}_{\mathbf{h}} \tilde{P}(\mathbf{x}|\mathbf{h})$$

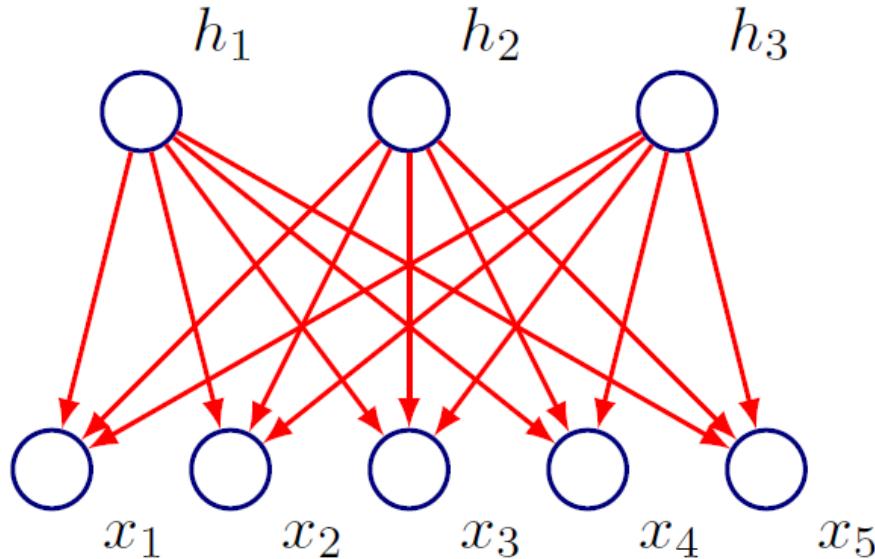
- \mathbf{h} is a representation of the data

Linear Factor Model

- The latent factor h is an encoding of the data
- Simplest decoding model: Get x after a linear transformation of x with some noise
- Formally: Suppose we sample the latent factors from a distribution $h \sim P(h)$
- Then: $x = Wh + b + \varepsilon$

Linear Factor Model

- $P(h)$ is a factorial distribution



$$\mathbf{x} = W\mathbf{h} + \mathbf{b} + \boldsymbol{\epsilon}$$

- How do learn in such a model?
- Let's look at a simple example

Probabilistic PCA

- Suppose underlying latent factor has a Gaussian distribution

$$\mathbf{h} \sim \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I})$$

- For the noise model: Assume $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$
- Then:

$$P(\mathbf{x}|\mathbf{h}) = \mathcal{N}(\mathbf{x}|W\mathbf{h} + \mathbf{b}, \sigma^2 \mathbf{I})$$

- We care about the marginal $P(\mathbf{x})$ (predictive distribution):

$$P(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{b}, WW^T + \sigma^2 \mathbf{I})$$

Probabilistic PCA

$$P(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{b}, WW^T + \sigma^2 I)$$

- How do we learn the parameters? (EM, ML Estimation)
- Let's look at the ML Estimation:
- Let $C = WW^T + \sigma^2 I$
- We want to maximize $\ell(\theta; X) = \sum_i \log P(\mathbf{x}_i|\theta)$

Probabilistic PCA: ML Estimation

$$\begin{aligned}\ell(\theta; X) &= \sum_i \log P(\mathbf{x}_i | \theta) \\ &= -\frac{N}{2} \log |C| - \frac{1}{2} \sum_i (\mathbf{x}_i - \mathbf{b}) C^{-1} (\mathbf{x}_i - \mathbf{b})^T \\ &= -\frac{N}{2} \log |C| - \frac{1}{2} \text{Tr}[(C^{-1} \sum_i (\mathbf{x}_i - \mathbf{b})(\mathbf{x}_i - \mathbf{b})^T)] \\ &= \frac{N}{2} \log |C| - \frac{1}{2} \text{Tr}[(C^{-1} S)]\end{aligned}$$

- Now fit the parameters $\theta = W, b, \sigma$ to maximize log-likelihood
- Can also use EM

Factor Analysis

- Fix the latent factor prior to be the unit Gaussian as before:

$$\mathbf{h} \sim \mathcal{N}(\mathbf{h}; \mathbf{0}, \mathbf{I})$$

- Noise is sampled from a Gaussian with a diagonal covariance:

$$\Psi = \text{diag}([\sigma_1^2, \sigma_2^2, \dots, \sigma_d^2])$$

- Still consider linear relationship between inputs and observed variables: Marginal $P(\mathbf{x}) \sim \mathcal{N}(\mathbf{x}; \mathbf{b}, \mathbf{W}\mathbf{W}^T + \Psi)$

Factor Analysis

- On deriving the posterior $P(\mathbf{h}|\mathbf{x}) = \mathcal{N}(\mathbf{h}|\boldsymbol{\mu}, \boldsymbol{\Lambda})$ we get:

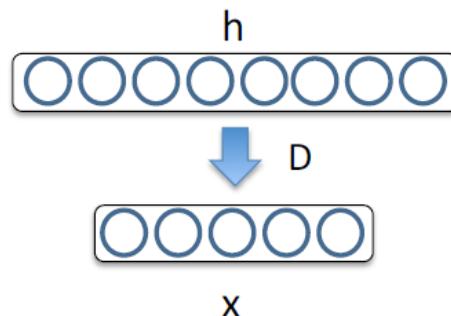
$$\boldsymbol{\mu} = \mathbf{W}^T (\mathbf{W}\mathbf{W}^T + \boldsymbol{\Psi})^{-1} (\mathbf{x} - \mathbf{b})$$

$$\boldsymbol{\Lambda} = \mathbf{I} - \mathbf{W}^T (\mathbf{W}\mathbf{W}^T + \boldsymbol{\Psi})^{-1} \mathbf{W}$$

- Parameters are coupled, makes ML estimation difficult
- Need to employ EM (or non-linear optimization)

Sparse Coding

- Sparse coding (Olshausen & Field, 1996). Originally developed to explain early visual processing in the brain (edge detection).
- For each input $x^{(t)}$ find a latent representation $h^{(t)}$ such that:
 - it is sparse: the vector $h^{(t)}$ has many zeros
 - we can good reconstruct the original input $x^{(t)}$



Sparse Coding

- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - it is sparse: the vector $\mathbf{h}^{(t)}$ has many zeros
 - we can good reconstruct the original input $\mathbf{x}^{(t)}$
- In other words:

Reconstruction: $\hat{\mathbf{x}}^{(t)}$

$$\min_{\mathbf{D}} \frac{1}{T} \sum_{t=1}^T \min_{\mathbf{h}^{(t)}} \frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2 + \lambda \|\mathbf{h}^{(t)}\|_1$$

Sparsity vs.
reconstruction control

Reconstruction error Sparsity penalty

The diagram illustrates the cost function for sparse coding. It features a red arrow pointing from the word "Reconstruction" to the squared difference term $\frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2$. Another red arrow points from the phrase "Sparsity vs. reconstruction control" to the L1 norm term $\lambda \|\mathbf{h}^{(t)}\|_1$. Below the equation, blue curly braces group the terms: the first brace groups the reconstruction error term, and the second brace groups the sparsity penalty term.

Sparse Coding

- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - it is sparse: the vector $\mathbf{h}^{(t)}$ has many zeros
 - we can good reconstruct the original input $\mathbf{x}^{(t)}$
- In other words:

$$\min_{\mathbf{D}} \frac{1}{T} \sum_{t=1}^T \min_{\mathbf{h}^{(t)}} \frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2 + \lambda \|\mathbf{h}^{(t)}\|_1$$

- we also constrain the columns of \mathbf{D} to be of norm 1
- otherwise, \mathbf{D} could grow big while \mathbf{h} becomes small to satisfy the L1 constraint

Sparse Coding

- For each $\mathbf{x}^{(t)}$ find a latent representation $\mathbf{h}^{(t)}$ such that:
 - it is sparse: the vector $\mathbf{h}^{(t)}$ has many zeros
 - we can good reconstruct the original input $\mathbf{x}^{(t)}$
- In other words:

$$\min_{\mathbf{D}} \frac{1}{T} \sum_{t=1}^T \min_{\mathbf{h}^{(t)}} \frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2 + \lambda \|\mathbf{h}^{(t)}\|_1$$

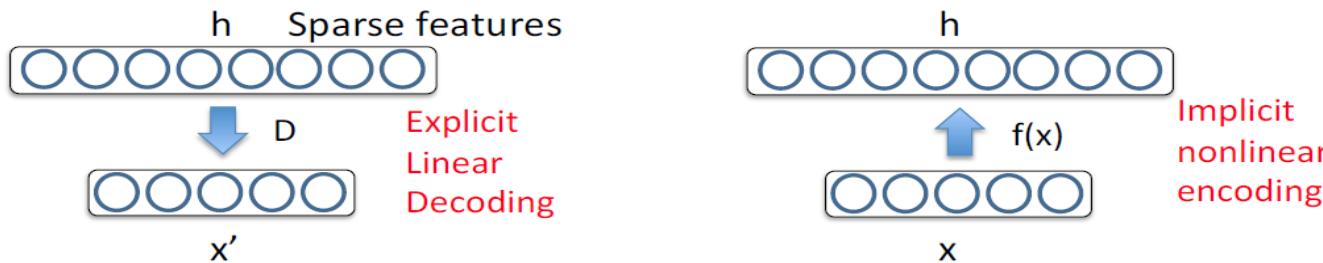
- \mathbf{D} is equivalent to the autoencoder output weight matrix
- However, $\mathbf{h}(\mathbf{x}^{(t)})$ is now a complicated function of $\mathbf{x}^{(t)}$
- Encoder is the minimization problem:

$$\mathbf{h}(\mathbf{x}^{(t)}) = \arg \min_{\mathbf{h}^{(t)}} \frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2 + \lambda \|\mathbf{h}^{(t)}\|_1$$

Interpreting Sparse Coding

Interpreting Sparse Coding

$$\min_{\mathbf{D}} \frac{1}{T} \sum_{t=1}^T \min_{\mathbf{h}^{(t)}} \frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2 + \lambda \|\mathbf{h}^{(t)}\|_1$$



- Sparse, over-complete representation h .
- Encoding $h = f(x)$ is implicit and nonlinear function of x .
- Reconstruction (or decoding) $x' = Dh$ is linear and explicit.

Sparse Coding

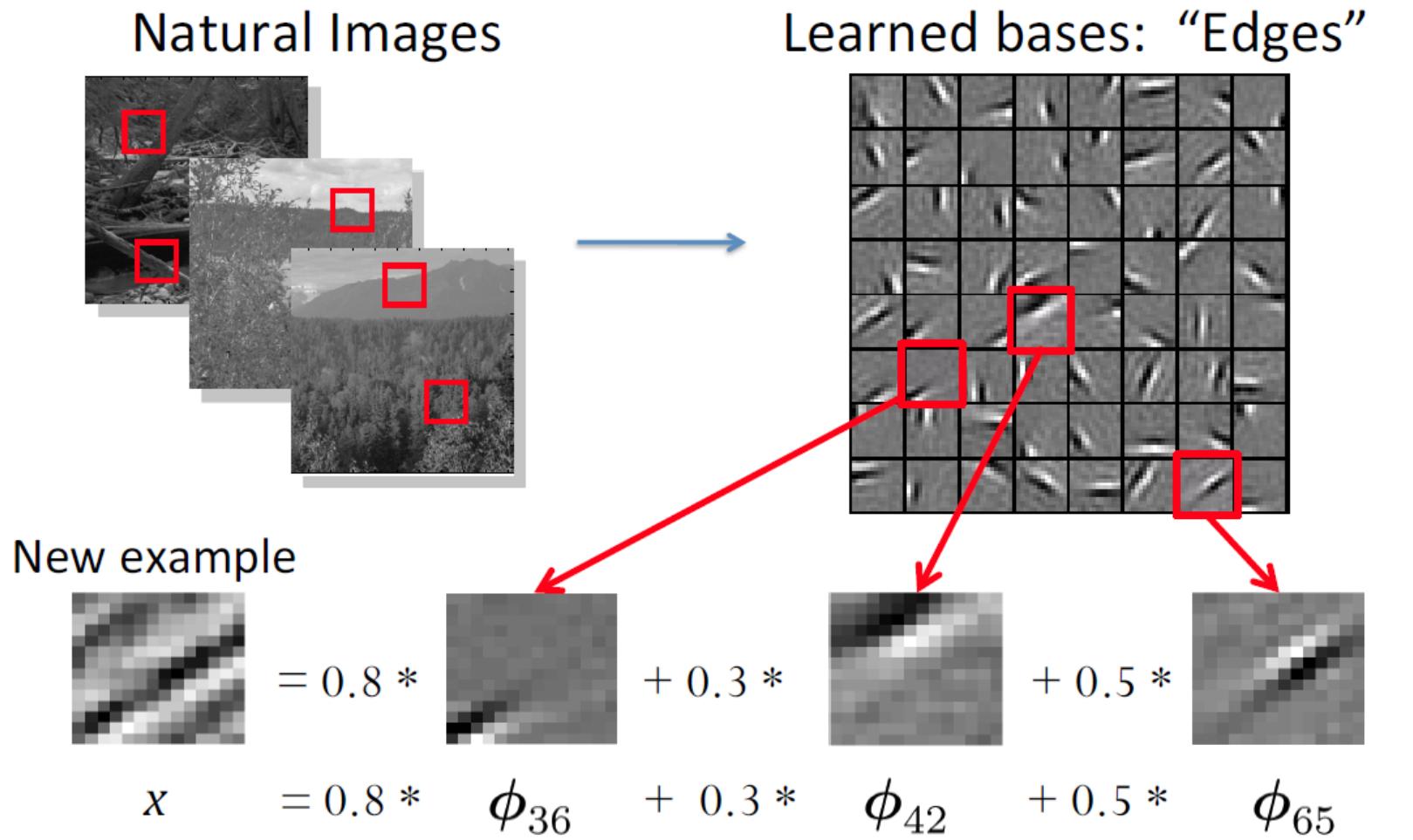
- We can also write:

$$\hat{\mathbf{x}}^{(t)} = \mathbf{D} \mathbf{h}(\mathbf{x}^{(t)}) = \sum_{\substack{k \text{ s.t.} \\ h(\mathbf{x}^{(t)})_k \neq 0}} \mathbf{D}_{\cdot, k} h(\mathbf{x}^{(t)})_k$$

$$\begin{matrix} 7 \\ = 1 \end{matrix} \begin{matrix} 9 \\ + 1 \end{matrix} \begin{matrix} 7 \\ + 1 \end{matrix} \begin{matrix} 7 \\ + 1 \end{matrix} \begin{matrix} 9 \\ + 1 \end{matrix} \begin{matrix} 7 \\ + 1 \end{matrix}$$
$$+ 1 \begin{matrix} 7 \\ + 1 \end{matrix} \begin{matrix} 7 \\ + 0.8 \end{matrix} \begin{matrix} 7 \\ + 0.8 \end{matrix} \begin{matrix} 7 \\ + 0.8 \end{matrix}$$

- D is often referred to as **Dictionary**
- In certain applications, we know what dictionary matrix to use
- In many cases, we have to learn it

Sparse Coding



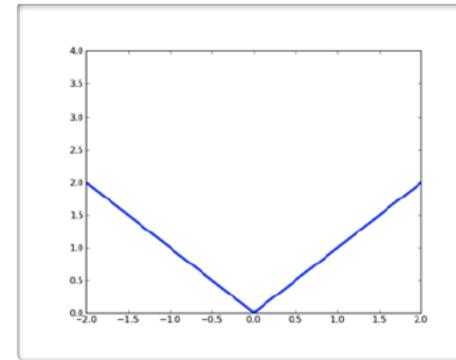
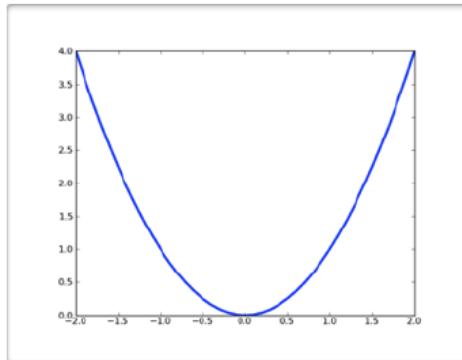
[0, 0, ... **0.8**, ..., **0.3**, ..., **0.5**, ...] = coefficients (feature representation)

Inference

- Given dictionary D , how do we compute $\mathbf{h}(\mathbf{x}^{(t)})$
 - We need to optimize:

$$l(\mathbf{x}^{(t)}) = \frac{1}{2} \|\mathbf{x}^{(t)} - \mathbf{D} \mathbf{h}^{(t)}\|_2^2 + \lambda \|\mathbf{h}^{(t)}\|_1$$

- This is Lasso.

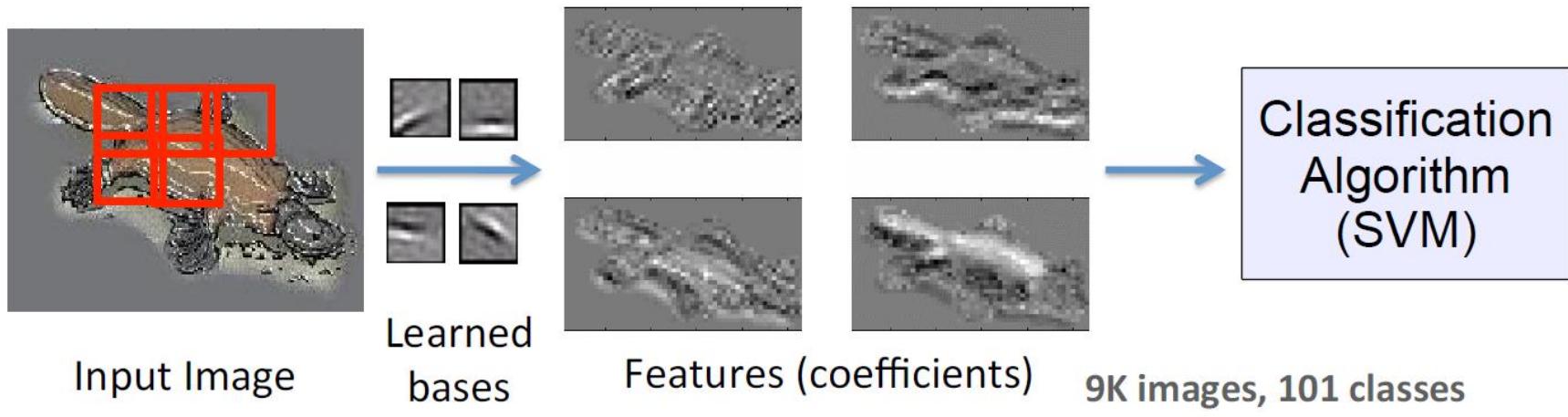


- We could use a gradient descent method:

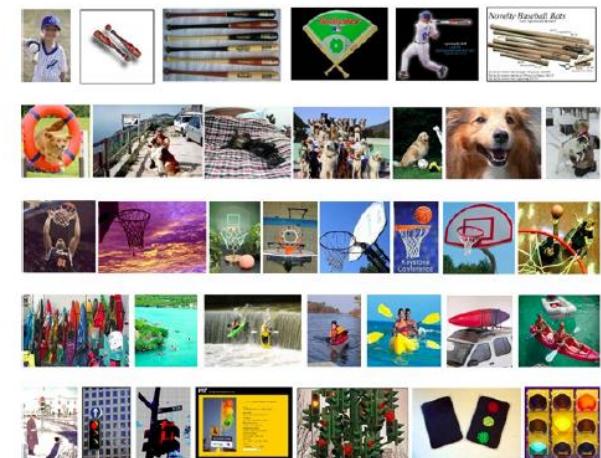
$$\nabla_{\mathbf{h}^{(t)}} l(\mathbf{x}^{(t)}) = \mathbf{D}^\top (\mathbf{D} \mathbf{h}^{(t)} - \mathbf{x}^{(t)}) + \lambda \text{sign}(\mathbf{h}^{(t)})$$

Image Classification

- Evaluated on Caltech101 object category dataset



Algorithm	Accuracy
Baseline (Fei-Fei et al., 2004)	16%
PCA	37%
Sparse Coding	47%



Lee et al., NIPS 2006

Outline

1 Course Review

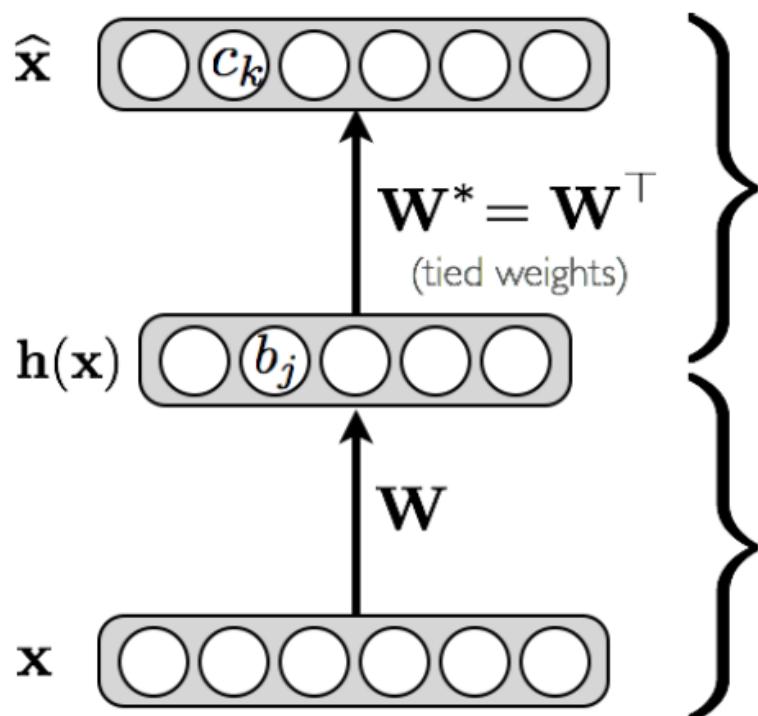
2 Linear Factor Model

3 Autoencoder

4 DBN and RBM

Autoencoder

- Feed-forward neural network trained to reproduce its input at the output layer



Decoder

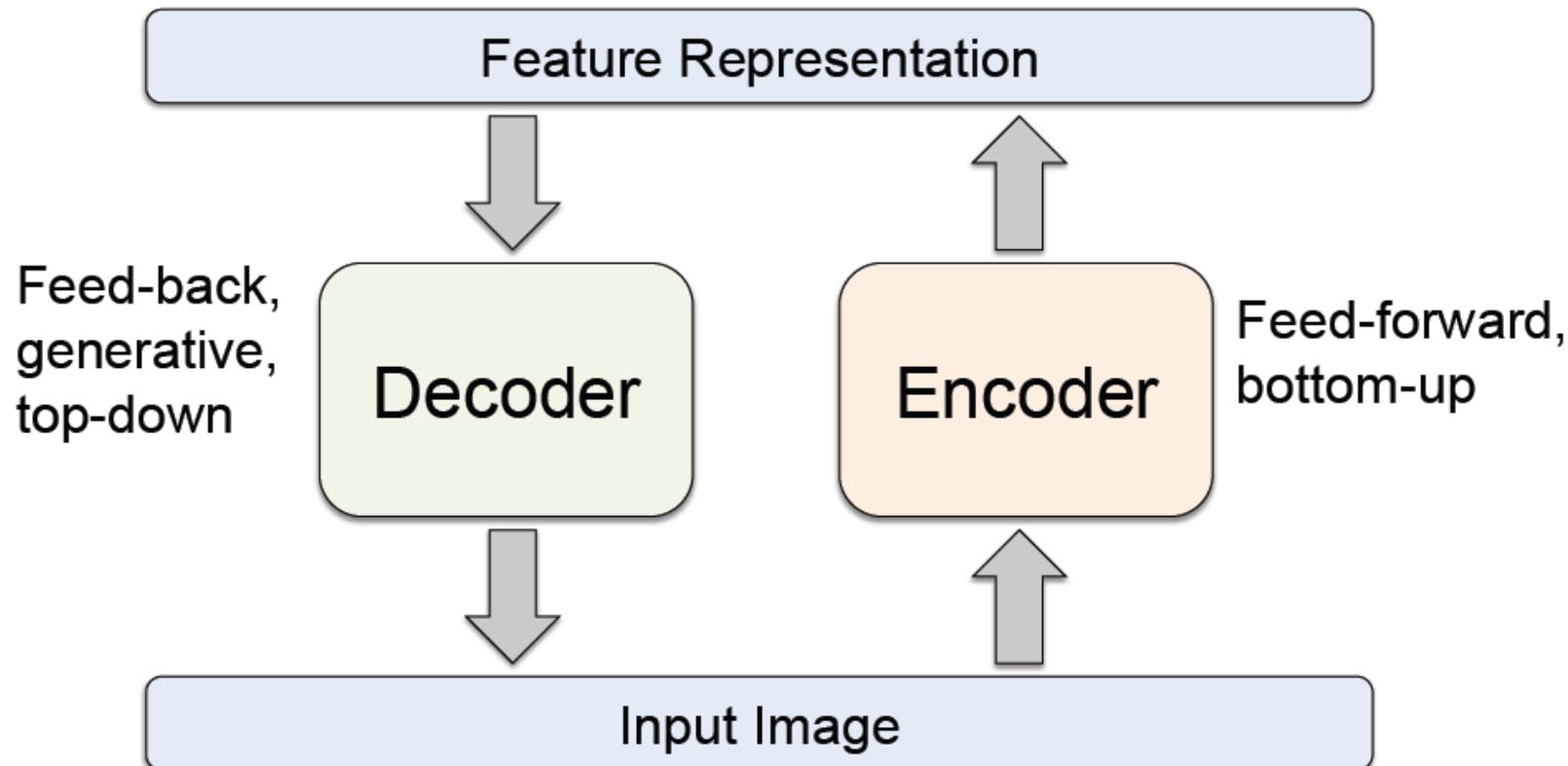
$$\begin{aligned}\hat{x} &= o(\hat{\mathbf{a}}(x)) \\ &= \text{sigm}(\mathbf{c} + \mathbf{W}^* \mathbf{h}(x))\end{aligned}$$

For binary units
|

Encoder

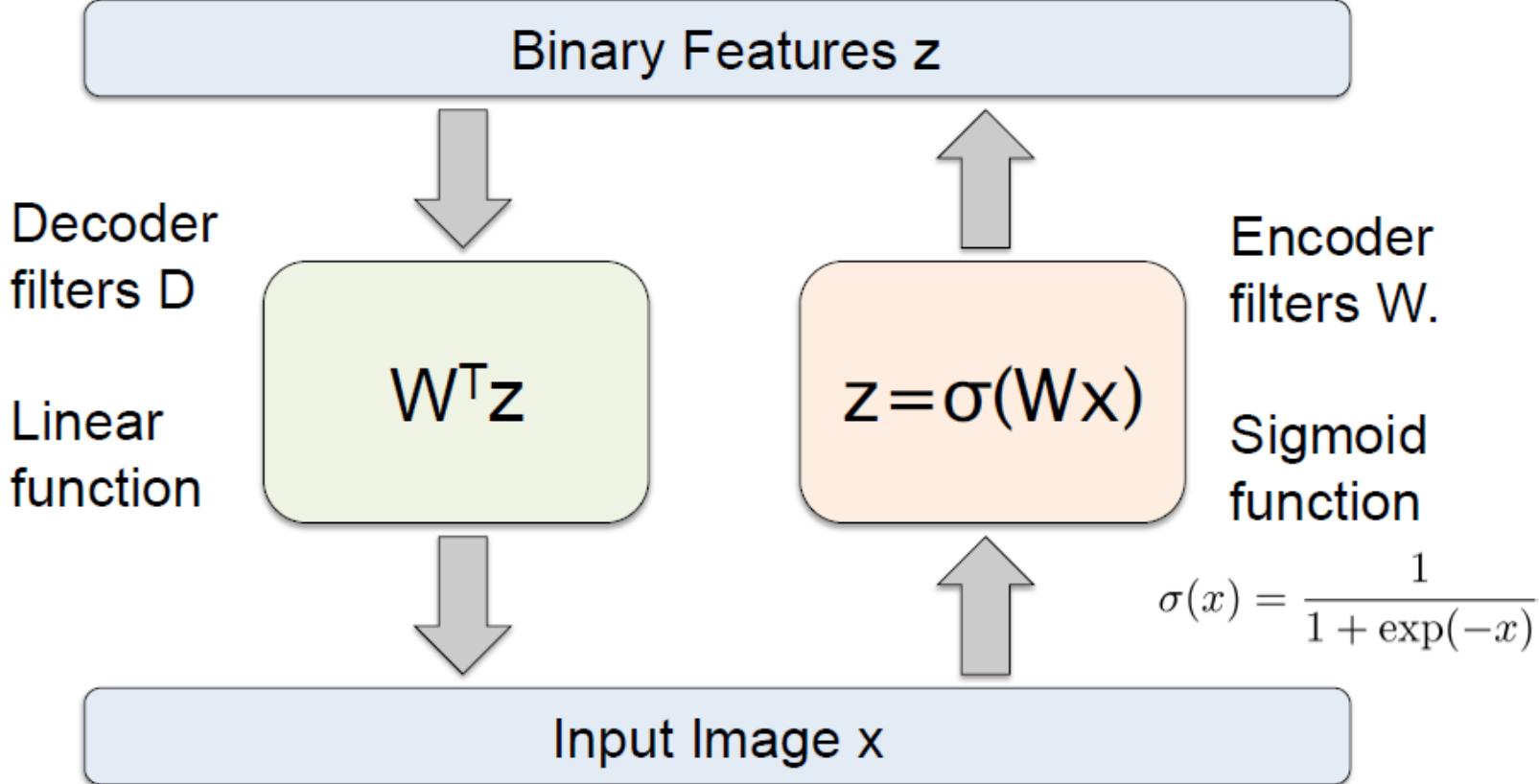
$$\begin{aligned}\mathbf{h}(x) &= g(\mathbf{a}(x)) \\ &= \text{sigm}(\mathbf{b} + \mathbf{Wx})\end{aligned}$$

Autoencoder

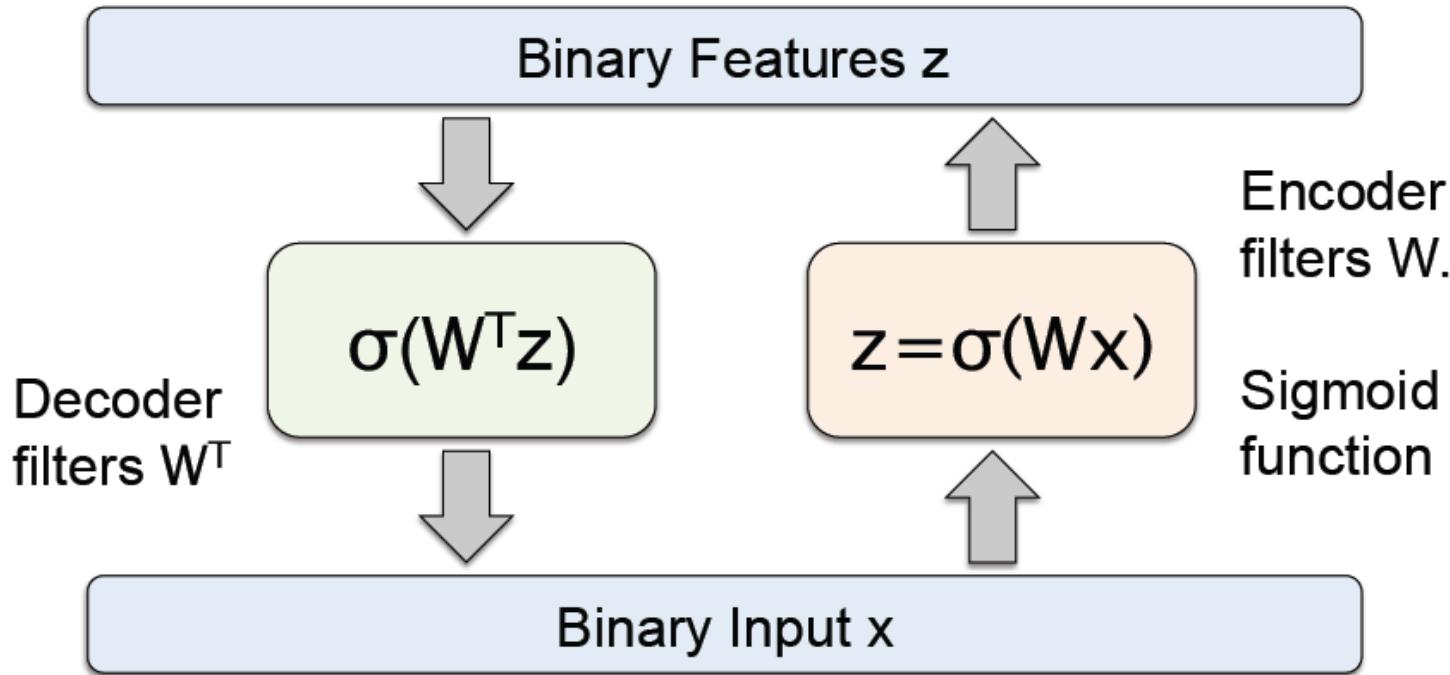


- Details of what goes inside the encoder and decoder matter!
- Need constraints to avoid learning an identity.

Autoencoder



Another Autoencoder Model



- Need additional constraints to avoid learning an identity.
- Relates to Restricted Boltzmann Machines
- Encoder and Decoder filters can be different.

Loss Function

- Loss function for binary inputs

$$l(f(\mathbf{x})) = - \sum_k (x_k \log(\hat{x}_k) + (1 - x_k) \log(1 - \hat{x}_k))$$

- Cross-entropy error function (reconstruction loss) $f(\mathbf{x}) \equiv \hat{\mathbf{x}}$
- Loss function for real-valued inputs

$$l(f(\mathbf{x})) = \frac{1}{2} \sum_k (\hat{x}_k - x_k)^2$$

- sum of squared differences (reconstruction loss)
- we use a linear activation function at the output

Loss Function

- For both cases, the gradient $\nabla_{\hat{\mathbf{a}}(\mathbf{x}^{(t)})} l(f(\mathbf{x}^{(t)}))$ has a very simple form:

$$\nabla_{\hat{\mathbf{a}}(\mathbf{x}^{(t)})} l(f(\mathbf{x}^{(t)})) = \hat{\mathbf{x}}^{(t)} - \mathbf{x}^{(t)} \quad f(\mathbf{x}) \equiv \hat{\mathbf{x}}$$

- Parameter gradients are obtained by backpropagating the gradient $\nabla_{\hat{\mathbf{a}}(\mathbf{x}^{(t)})} l(f(\mathbf{x}^{(t)}))$ like in a regular network
 - important: when using tied weights ($\mathbf{W}^* = \mathbf{W}^\top$), $\nabla_{\mathbf{W}} l(f(\mathbf{x}^{(t)}))$ is the sum of two gradients
 - this is because \mathbf{W} is present in the encoder and in the decoder

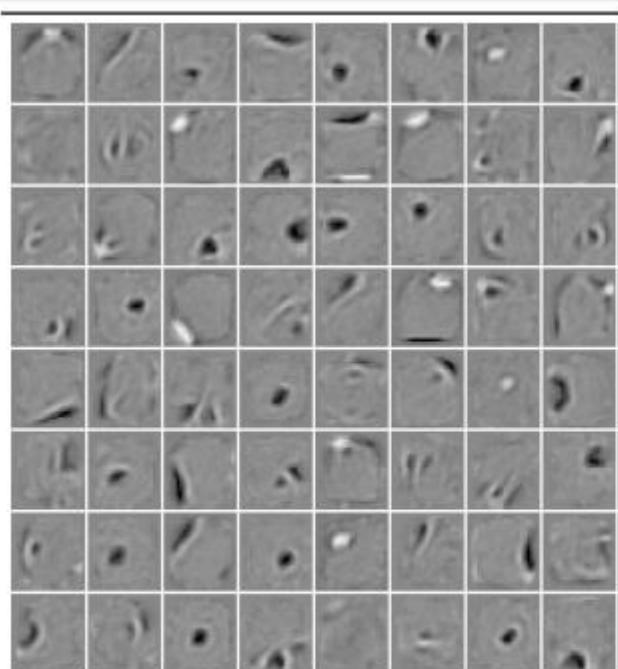
Example: MNIST

- MNIST dataset:

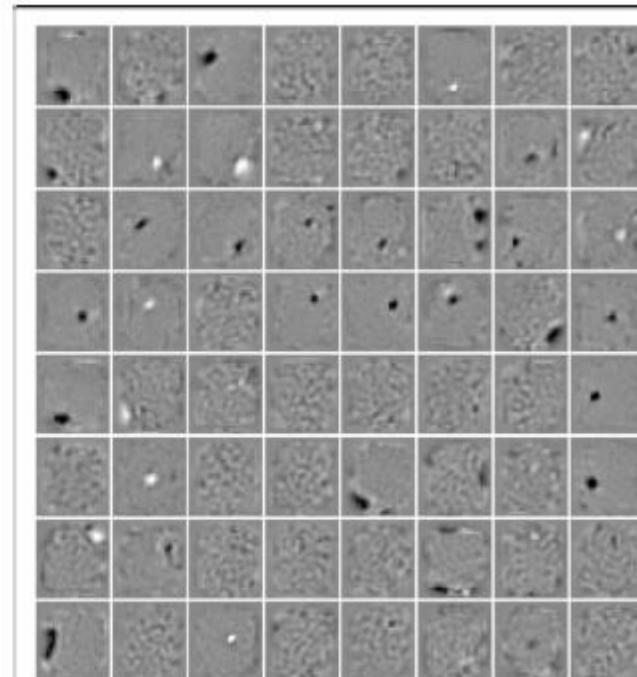
3	8	6	9	6	4	5	3	8	4	5	2	3	8	4	8
1	5	0	5	9	7	4	1	0	3	0	6	2	9	9	4
1	3	6	8	0	7	1	6	8	9	0	3	8	3	7	7
8	4	4	1	2	9	8	1	1	0	6	6	5	0	1	1
7	2	7	3	1	4	0	5	0	6	8	7	6	8	9	9
4	0	6	1	9	2	1	3	9	4	4	5	6	6	1	7
2	8	6	9	7	0	9	1	6	2	8	3	6	4	9	5
8	6	8	7	8	8	6	9	1	7	6	0	9	6	7	0

Learned Features

- MNIST dataset:

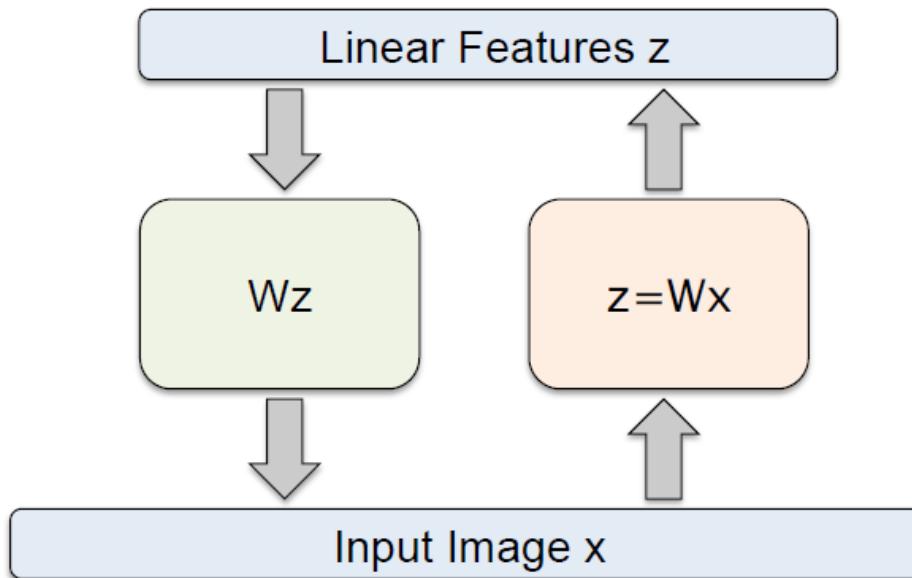


RBM



Autoencoder

Optimality of the Linear Autoencoder



- If the hidden and output layers are linear, it will learn hidden units that are a linear function of the data and minimize the squared error.
- The K hidden units will span the same space as the first k principal components. The weight vectors may not be orthogonal.

- With nonlinear hidden units, we have a nonlinear generalization of PCA.

Optimality of the Linear Autoencoder

- Let us consider the following theorem:

- let \mathbf{A} be any matrix, with singular value decomposition $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^\top$
 - Σ is a diagonal matrix
 - \mathbf{V}, \mathbf{U} are orthonormal matrices (columns/rows are orthonormal vectors)
- let $\mathbf{U}_{\cdot, \leq k} \Sigma_{\leq k, \leq k} \mathbf{V}_{\cdot, \leq k}^\top$ be the decomposition where we keep only the k largest singular values
- then, the matrix \mathbf{B} of rank k that is closest to \mathbf{A} :

$$\mathbf{B}^* = \underset{\mathbf{B} \text{ s.t. } \text{rank}(\mathbf{B})=k}{\arg \min} \|\mathbf{A} - \mathbf{B}\|_F$$

is $\mathbf{B}^* = \mathbf{U}_{\cdot, \leq k} \Sigma_{\leq k, \leq k} \mathbf{V}_{\cdot, \leq k}^\top$

Optimality of the Linear Autoencoder

$$\min_{\theta} \sum_t \frac{1}{2} \sum_i (x_i^{(t)} - \underbrace{\hat{x}_i^{(t)}}_{\text{based on linear encoder}})^2 \geq \min_{\mathbf{W}^*, \mathbf{h}(\mathbf{X})} \frac{1}{2} \|\overbrace{\mathbf{X} - \mathbf{W}^* \mathbf{h}(\mathbf{X})}^{\text{matrix where columns are } \mathbf{x}^{(t)}}\|_F^2$$

matrix of all hidden layers
(could be any encoder)

$$\arg \min_{\mathbf{W}^*, \mathbf{h}(\mathbf{X})} \frac{1}{2} \|\mathbf{X} - \mathbf{W}^* \mathbf{h}(\mathbf{X})\|_F^2 = (\mathbf{W}^* \leftarrow \mathbf{U}_{\cdot, \leq k} \Sigma_{\leq k, \leq k}, \mathbf{h}(\mathbf{X}) \leftarrow \mathbf{V}_{\cdot, \leq k}^\top)$$

based on previous theorem, where $\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^\top$
and k is the hidden layer size

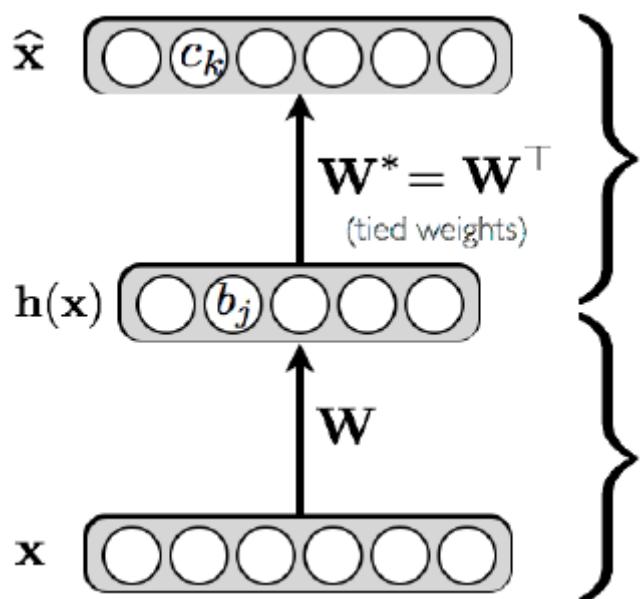
Let's show $\mathbf{h}(\mathbf{X})$ is a linear encoder:

$$\begin{aligned}
 \mathbf{h}(\mathbf{X}) &= \mathbf{V}_{\cdot, \leq k}^\top \\
 &= \mathbf{V}_{\cdot, \leq k}^\top (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{X}) && \leftarrow \text{multiplying by identity} \\
 &= \mathbf{V}_{\cdot, \leq k}^\top (\mathbf{V} \Sigma^\top \mathbf{U}^\top \mathbf{U} \Sigma \mathbf{V}^\top)^{-1} (\mathbf{V} \Sigma^\top \mathbf{U}^\top \mathbf{X}) && \leftarrow \text{replace with SVD} \\
 \\
 &= \mathbf{V}_{\cdot, \leq k}^\top \mathbf{V} (\Sigma^\top \Sigma)^{-1} \mathbf{V}^\top \mathbf{V} \Sigma^\top \mathbf{U}^\top \mathbf{X} && \leftarrow \mathbf{V}(\Sigma^\top \Sigma)^{-1} \mathbf{V}^\top \mathbf{V} \Sigma^\top \Sigma \mathbf{V}^\top = \mathbf{I} \\
 &= \mathbf{V}_{\cdot, \leq k}^\top \mathbf{V} (\Sigma^\top \Sigma)^{-1} \Sigma^\top \mathbf{U}^\top \mathbf{X} && \leftarrow \mathbf{V}^\top \mathbf{V} = \mathbf{I} \text{ (orthonormal)} \\
 &= \mathbf{I}_{\leq k, \cdot} (\Sigma^\top \Sigma)^{-1} \Sigma^\top \mathbf{U}^\top \mathbf{X} && \leftarrow \text{idem} \\
 &= \mathbf{I}_{\leq k, \cdot} \Sigma^{-1} (\Sigma^\top)^{-1} \Sigma^\top \mathbf{U}^\top \mathbf{X} && \leftarrow (\Sigma^\top \Sigma)^{-1} = \Sigma^{-1} (\Sigma^\top)^{-1} \\
 &= \mathbf{I}_{\leq k, \cdot} \Sigma^{-1} \mathbf{U}^\top \mathbf{X} \\
 &= \underbrace{\Sigma_{\leq k, \leq k}^{-1} (\mathbf{U}_{\cdot, \leq k})^\top}_{\text{this is a linear encoder}} \mathbf{X} && \leftarrow \text{multiplying by } \mathbf{I}_{\leq k, \cdot} \text{ selects the } k \text{ first rows}
 \end{aligned}$$

Optimality of the Linear Autoencoder

- So an optimal pair of encoder and decoder is

$$\mathbf{h}(\mathbf{x}) = \underbrace{\left(\Sigma_{\leq k, \leq k}^{-1} (\mathbf{U}_{\cdot, \leq k})^\top \right) \mathbf{x}}_{\mathbf{W}}$$
$$\hat{\mathbf{x}} = \underbrace{(\mathbf{U}_{\cdot, \leq k} \Sigma_{\leq k, \leq k}) \mathbf{h}(\mathbf{x})}_{\mathbf{W}^*}$$

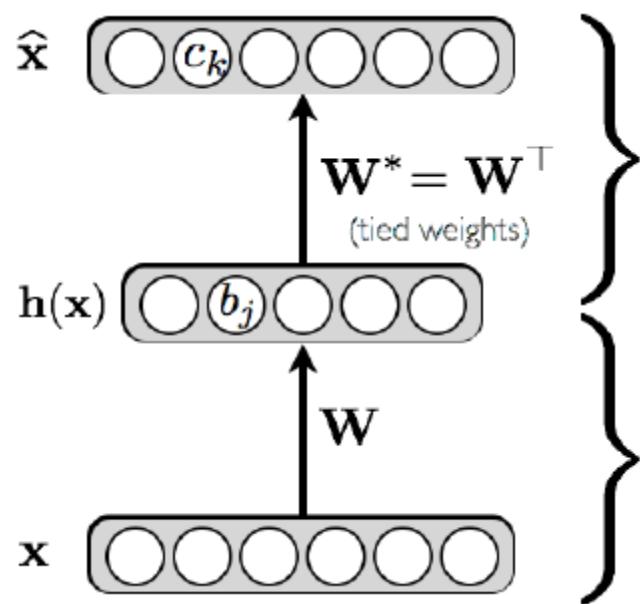


- for the sum of squared difference error ?
- for an autoencoder with a linear decoder
- where optimality means “has the lowest training reconstruction error”

Optimality of the Linear Autoencoder

- So an optimal pair of encoder and decoder is

$$\mathbf{h}(\mathbf{x}) = \underbrace{\left(\Sigma_{\leq k, \leq k}^{-1} (\mathbf{U}_{\cdot, \leq k})^\top \right) \mathbf{x}}_{\mathbf{W}} \quad \widehat{\mathbf{x}} = \underbrace{(\mathbf{U}_{\cdot, \leq k} \Sigma_{\leq k, \leq k}) \mathbf{h}(\mathbf{x})}_{\mathbf{W}^*}$$



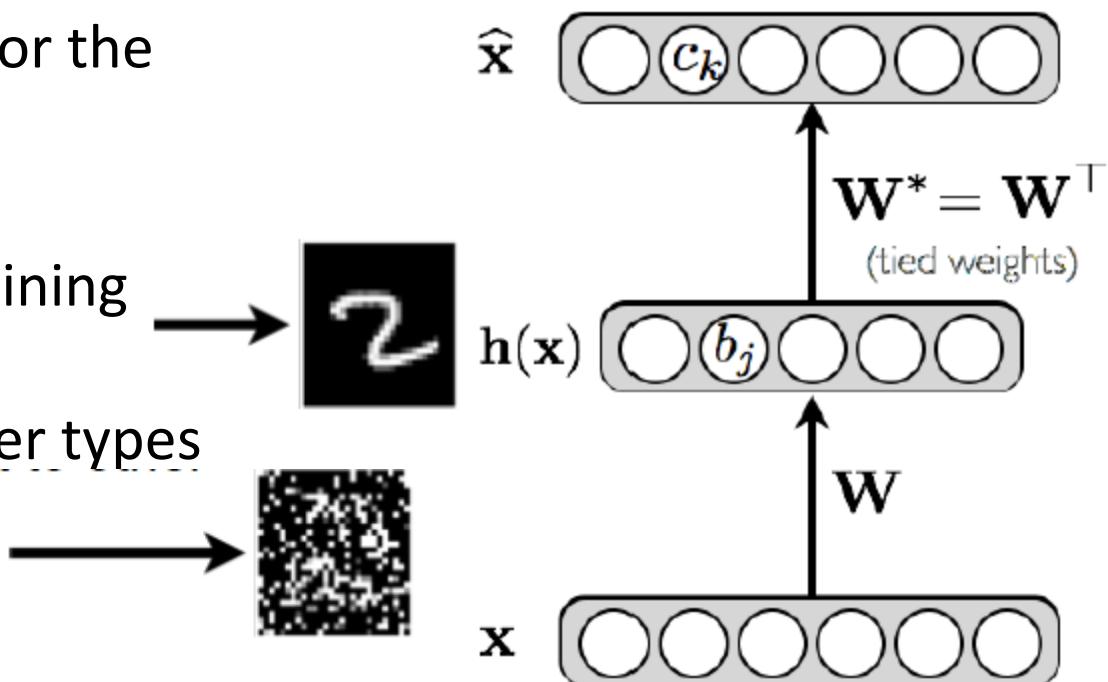
- If inputs are normalized as follows:

$$\mathbf{x}^{(t)} \leftarrow \frac{1}{\sqrt{T}} \left(\mathbf{x}^{(t)} - \frac{1}{T} \sum_{t'=1}^T \mathbf{x}^{(t')} \right)$$

- encoder corresponds to Principal Component Analysis (PCA)
- singular values and (left) vectors = the eigenvalues/vectors of covariance matrix

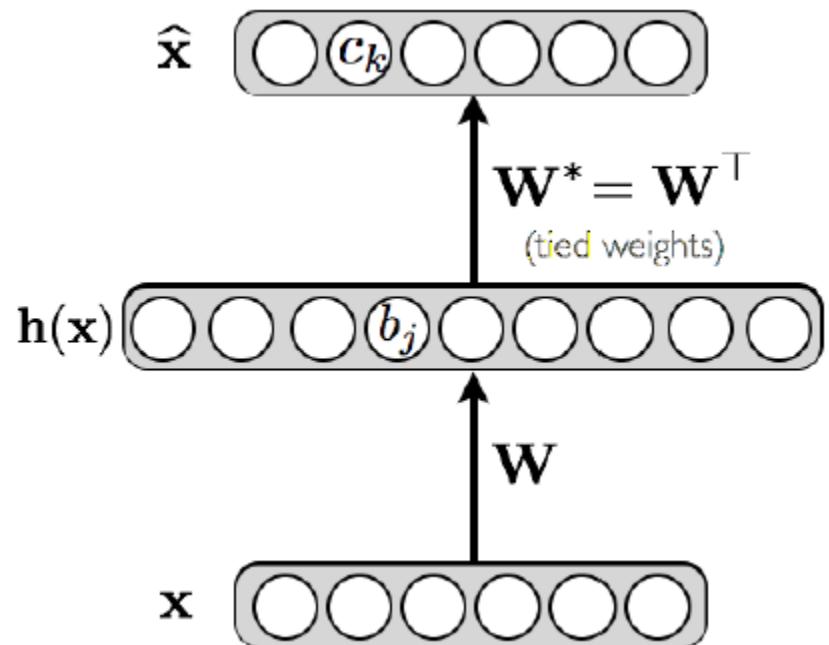
Undercomplete Representation

- Hidden layer is undercomplete if smaller than the input layer (bottleneck layer, e.g. dimensionality reduction):
 - hidden layer “compresses” the input
 - will compress well only for the training distribution
- Hidden units will be
 - good features for the training distribution
 - will not be robust to other types of input



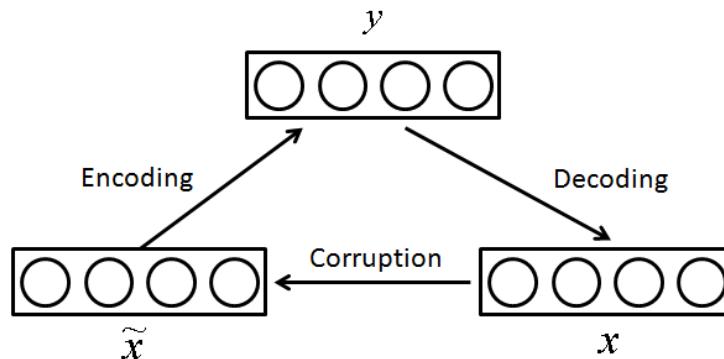
Overcomplete Representation

- Hidden layer is overcomplete if greater than the input layer
 - no compression in hidden layer
 - each hidden unit could copy a different input component
- No guarantee that the hidden units will extract meaningful structure



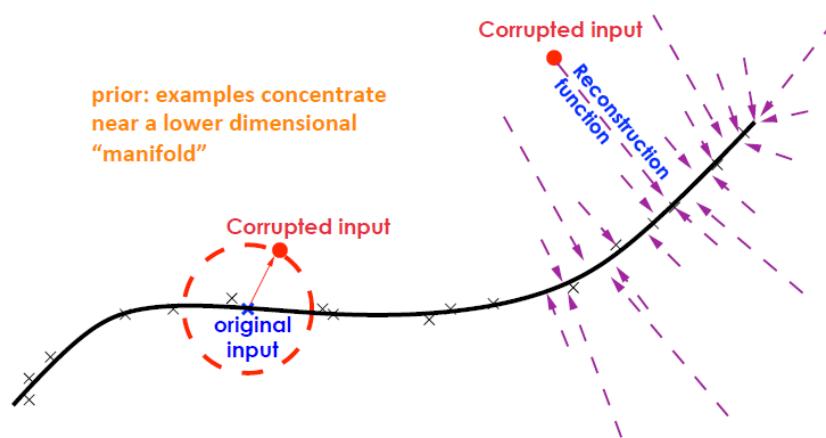
Denoising Autoencoder

- In order to force the hidden layer to discover more robust features and prevent it from simply learning the identity function, train the autoencoder to reconstruct the input from a corrupted version of it
 - Encode the input (preserve the information about the input)
 - Undo the effect of a corruption process stochastically applied to the input of the auto-encoder
- To convert the autoencoder to a denoising autoencoder, all we need to do is to add a stochastic corruption step operating on the input
 - Randomly sets some of the inputs (as many as half of them) to zero. Hence the denoising auto-encoder is trying to predict the corrupted (i.e. missing) values from the uncorrupted (i.e., non-missing) values, for randomly selected subsets of missing patterns.
 - The input can be corrupted in other ways



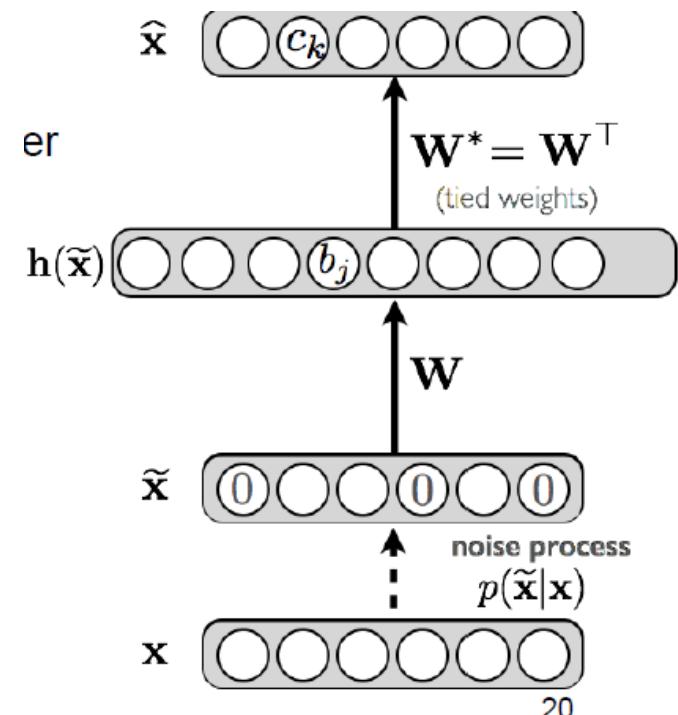
Denoising Autoencoder

- The learner must capture the structure of the input distribution in order to optimally undo the effect of the corruption process, with the reconstruction essentially being a nearby but higher density point than the corrupted input
- The denoising autoencoder is learning a reconstruction function that corresponds to a vector field pointing towards high-density regions (the manifold where examples concentrate)
- Denosing autoencoder basically learns in $r(\tilde{\mathbf{x}}) - \tilde{\mathbf{x}}$ a vector pointing in the direction $\frac{\partial \log P(\tilde{\mathbf{x}})}{\partial \tilde{\mathbf{x}}}$

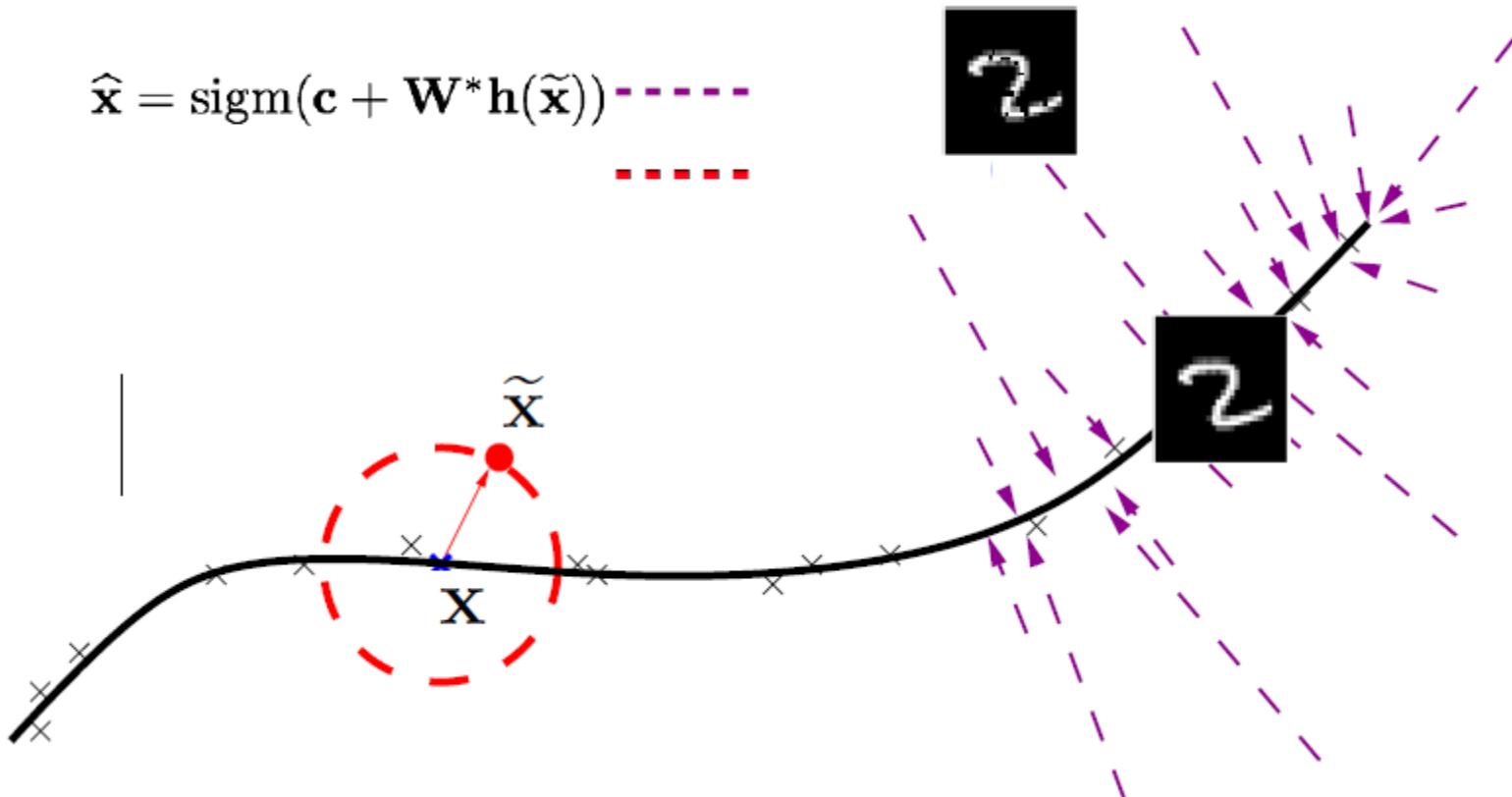


Denoising Autoencoder

- Idea: representation should be robust to introduction of noise:
 - random assignment of subset of inputs to 0, with probability
 - Similar to dropouts on the input layer
 - Gaussian additive noise
- Reconstruction $\hat{\mathbf{x}}$ computed from the corrupted input $\tilde{\mathbf{x}}$
- Loss function compares $\hat{\mathbf{x}}$ reconstruction with the noiseless input \mathbf{x}

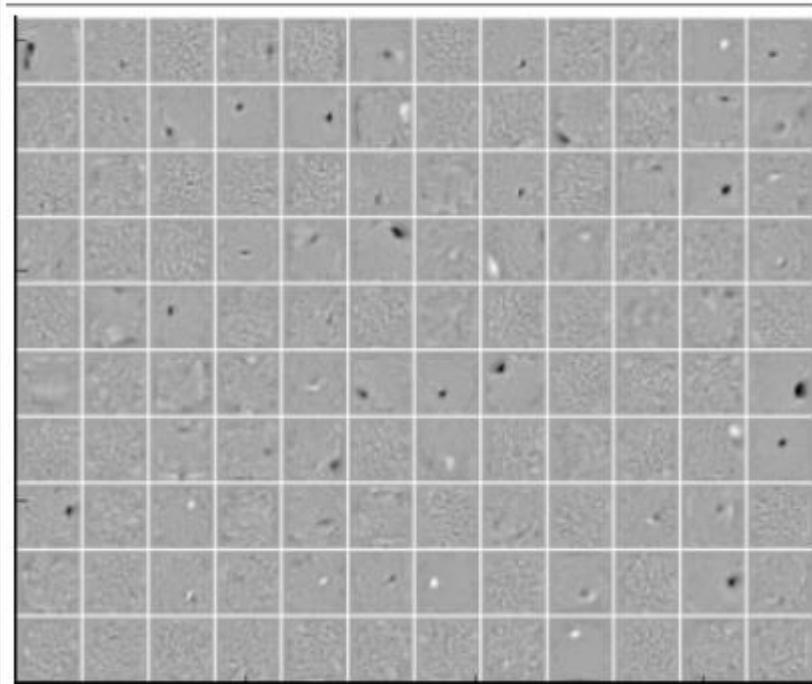


Denoising Autoencoder

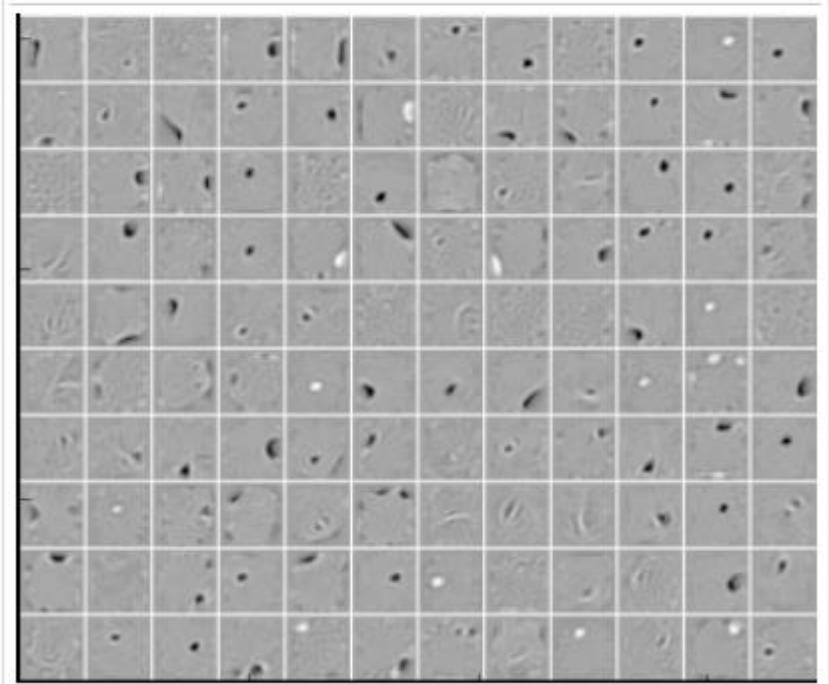


Learned Filters

Non-corrupted

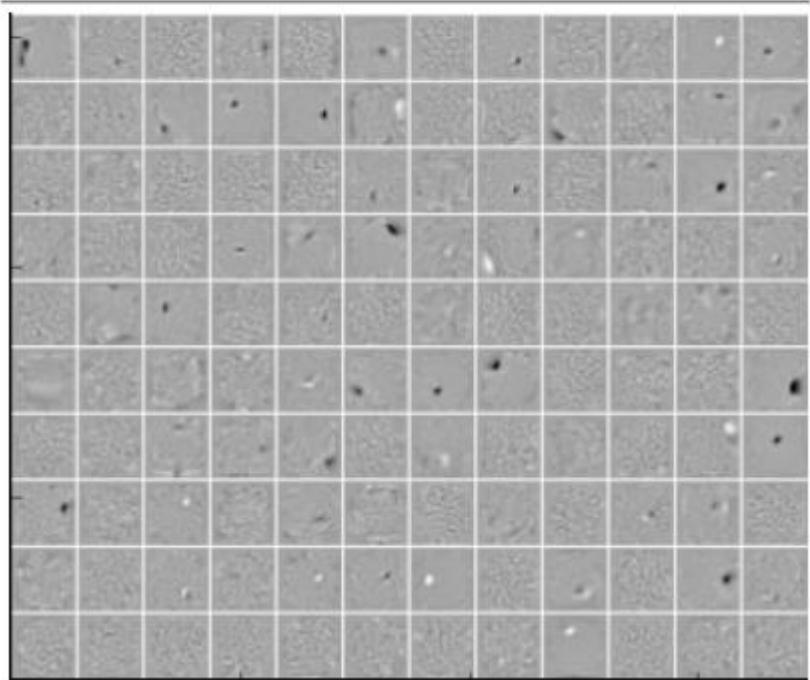


25% corrupted input

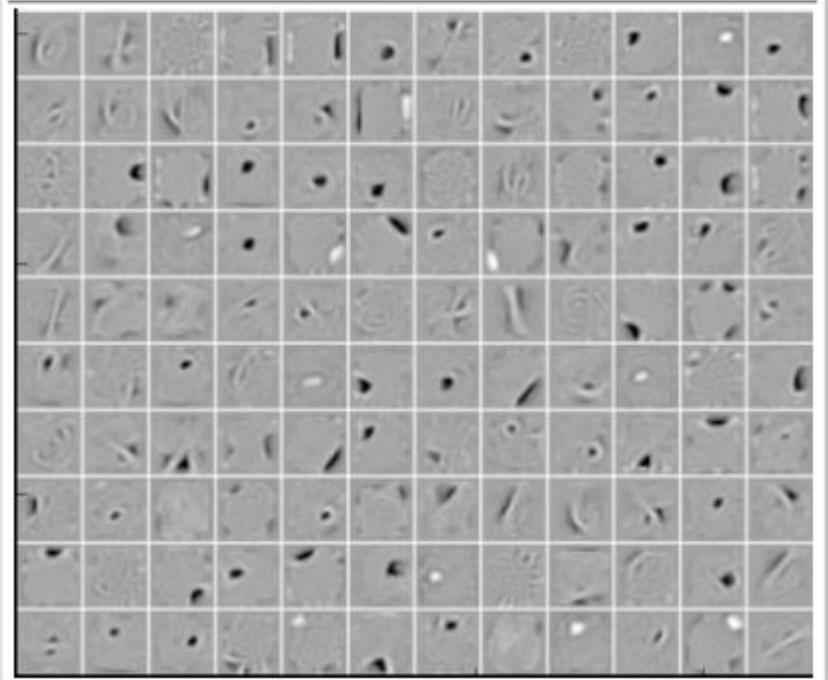


Learned Filters

Non-corrupted

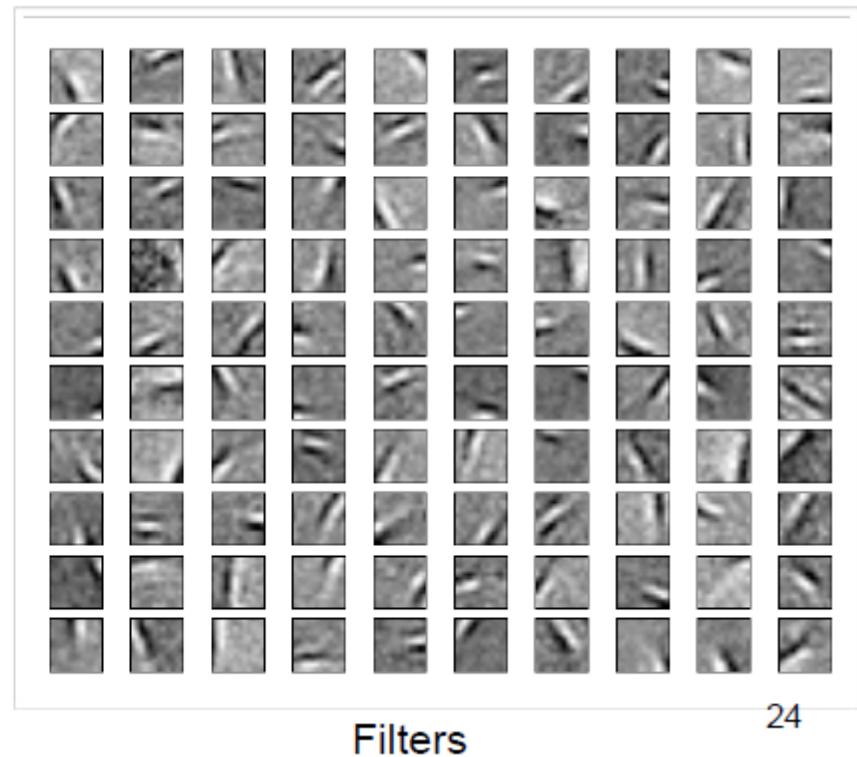
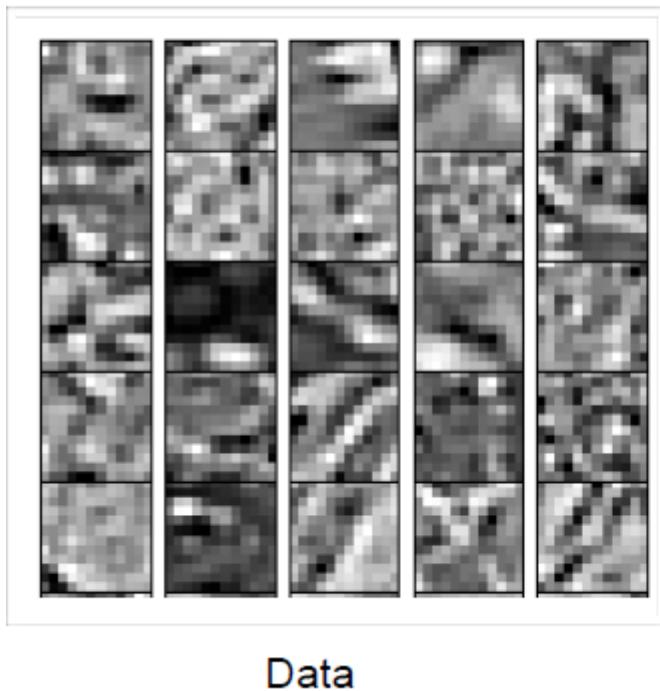


50% corrupted input



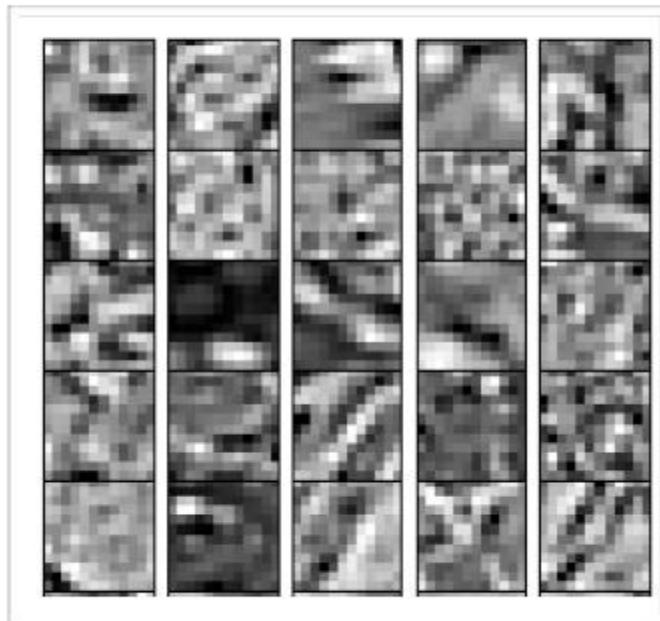
Squared Error Loss

- Training on natural image patches, with squared loss
 - PCA may not be the best solution

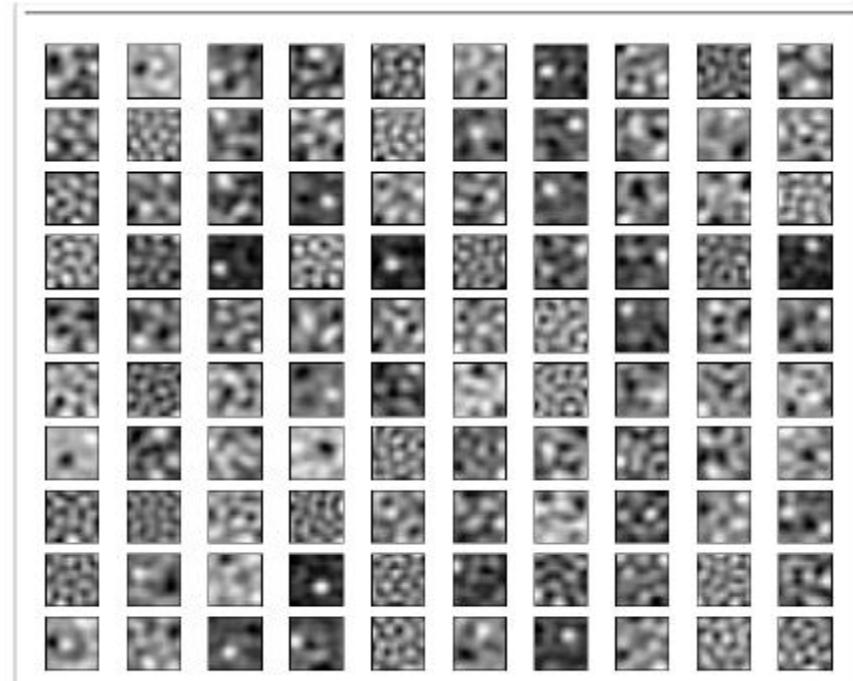


Squared Error Loss

- Training on natural image patches, with squared loss
 - Not equivalent to weight decay



Data

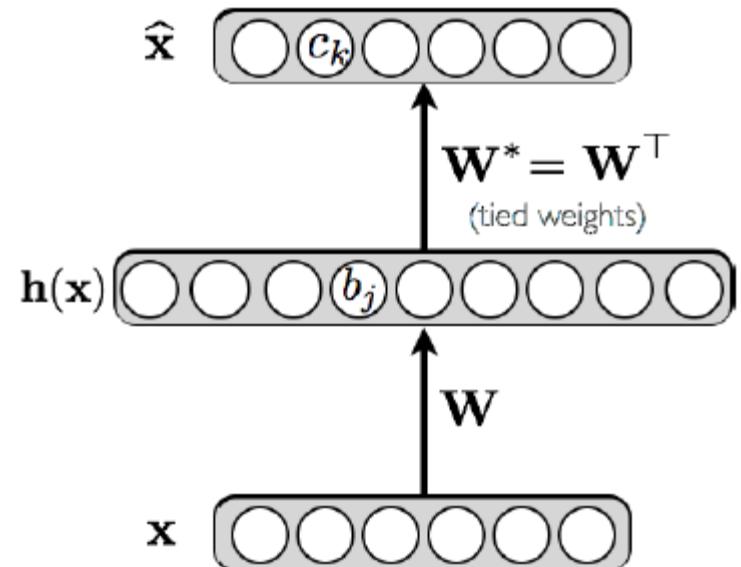


Filters

25

Contractive Autoencoder

- Alternative approach to avoid uninteresting solutions
 - add an explicit term in the loss that penalizes that solution
- We wish to extract features that only reflect variations observed in the training set
 - we'd like to be invariant to the other variations



Contractive Autoencoder

- Consider the following loss function:

$$l(f(\mathbf{x}^{(t)})) + \lambda \|\nabla_{\mathbf{x}^{(t)}} \mathbf{h}(\mathbf{x}^{(t)})\|_F^2$$

 Reconstruction Loss  Jacobian of Encoder

- For the **binary observations**:

$$l(f(\mathbf{x}^{(t)})) = - \sum_k \left(x_k^{(t)} \log(\hat{x}_k^{(t)}) + (1 - x_k^{(t)}) \log(1 - \hat{x}_k^{(t)}) \right)$$

$$\|\nabla_{\mathbf{x}^{(t)}} \mathbf{h}(\mathbf{x}^{(t)})\|_F^2 = \sum_j \sum_k \left(\frac{\partial h(\mathbf{x}^{(t)})_j}{\partial x_k^{(t)}} \right)^2$$

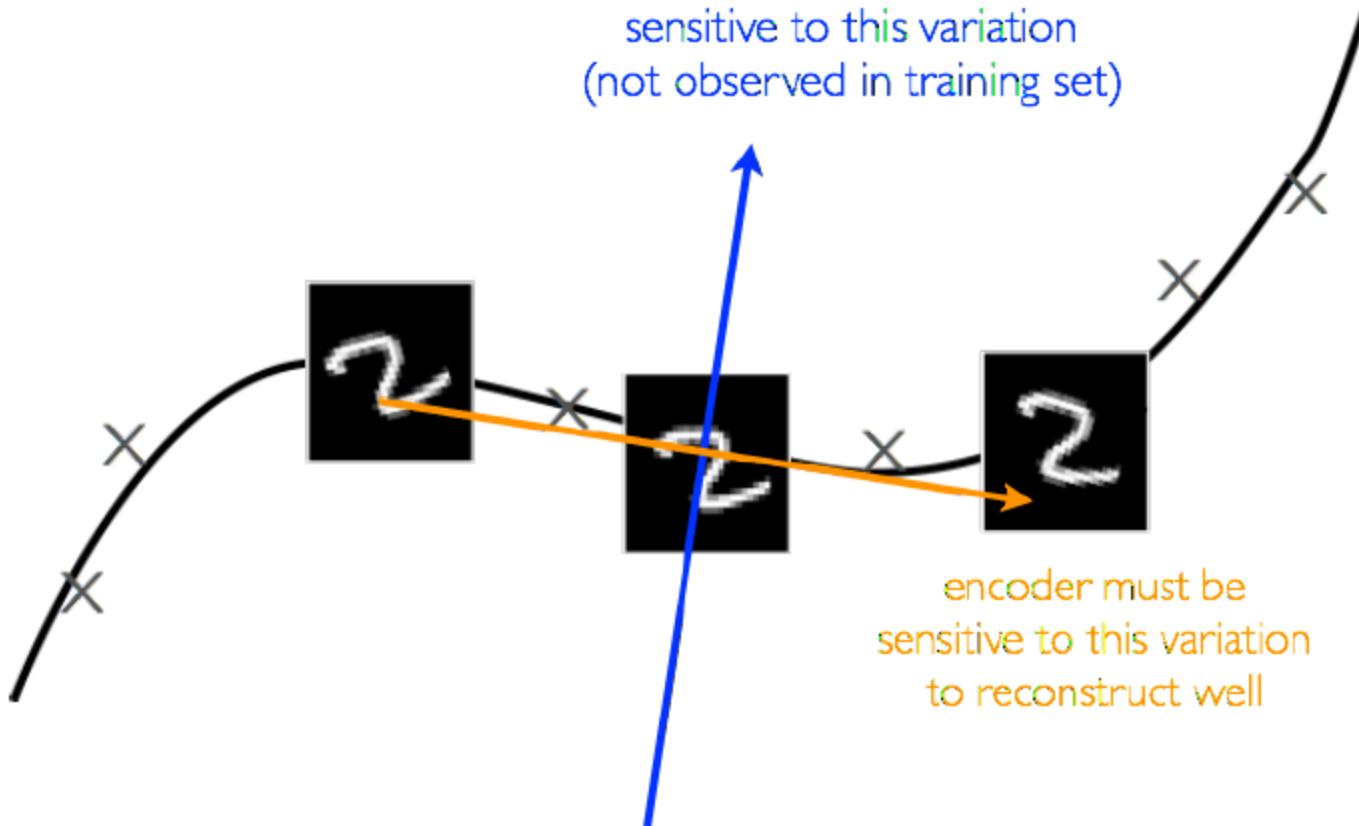
Encoder throws
away all information


Autoencoder attempts to
preserve all information

Contractive Autoencoder

- Illustration:

encoder doesn't need to be
sensitive to this variation
(not observed in training set)

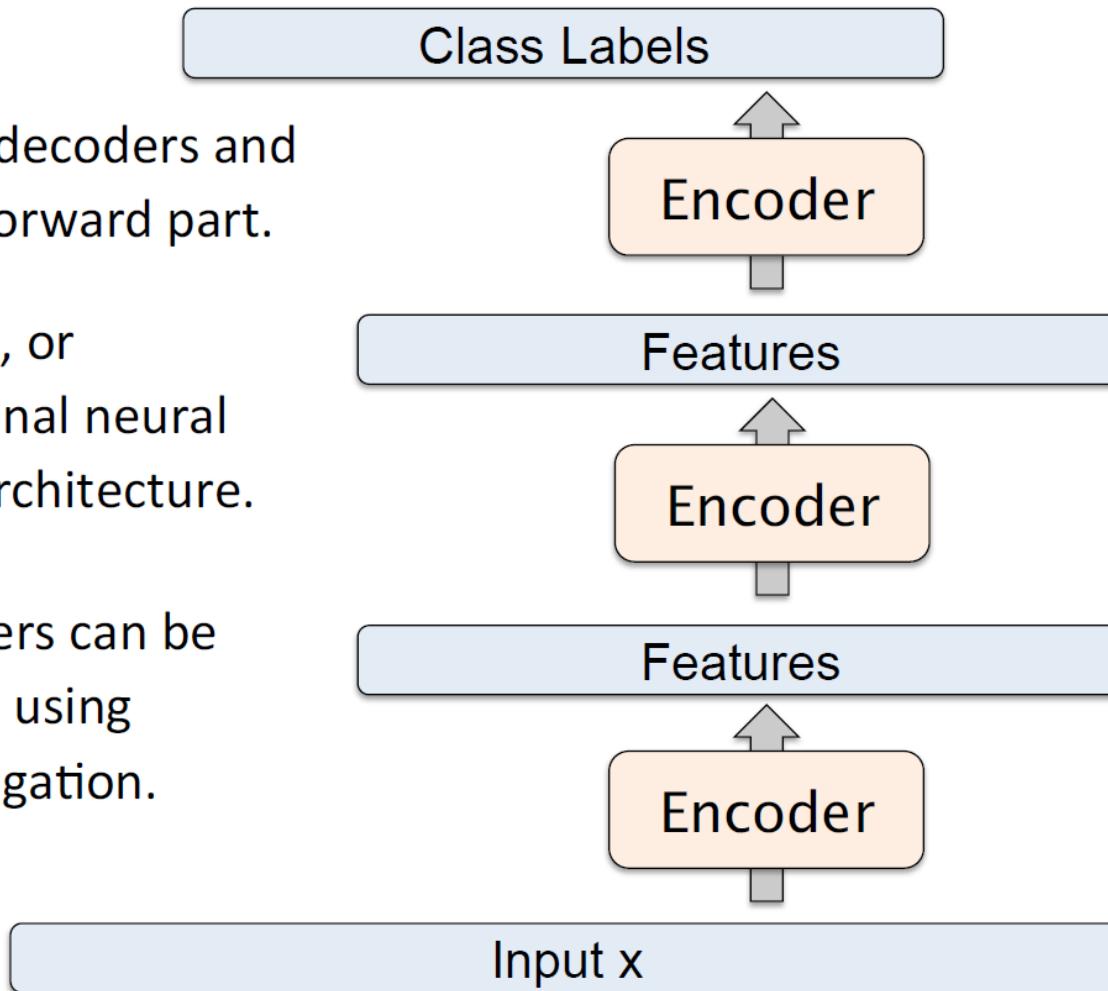


Pros and Cons

- Advantage of denoising autoencoder: simpler to implement
 - requires adding one or two lines of code to regular autoencoder
 - no need to compute Jacobian of hidden layer
- Advantage of contractive autoencoder: gradient is deterministic
 - can use second order optimizers (conjugate gradient, LBFGS, etc.)
 - might be more stable than denoising autoencoder, which uses a sampled gradient

Stacked Autoencoder

- Remove decoders and use feed-forward part.
- Standard, or convolutional neural network architecture.
- Parameters can be fine-tuned using backpropagation.



Stacked Autoencoder

- Remove decoders and use feed-forward part.
- Standard, or convolutional neural network architecture.
- Parameter fine-tuning backprop

Class Labels

Encoder

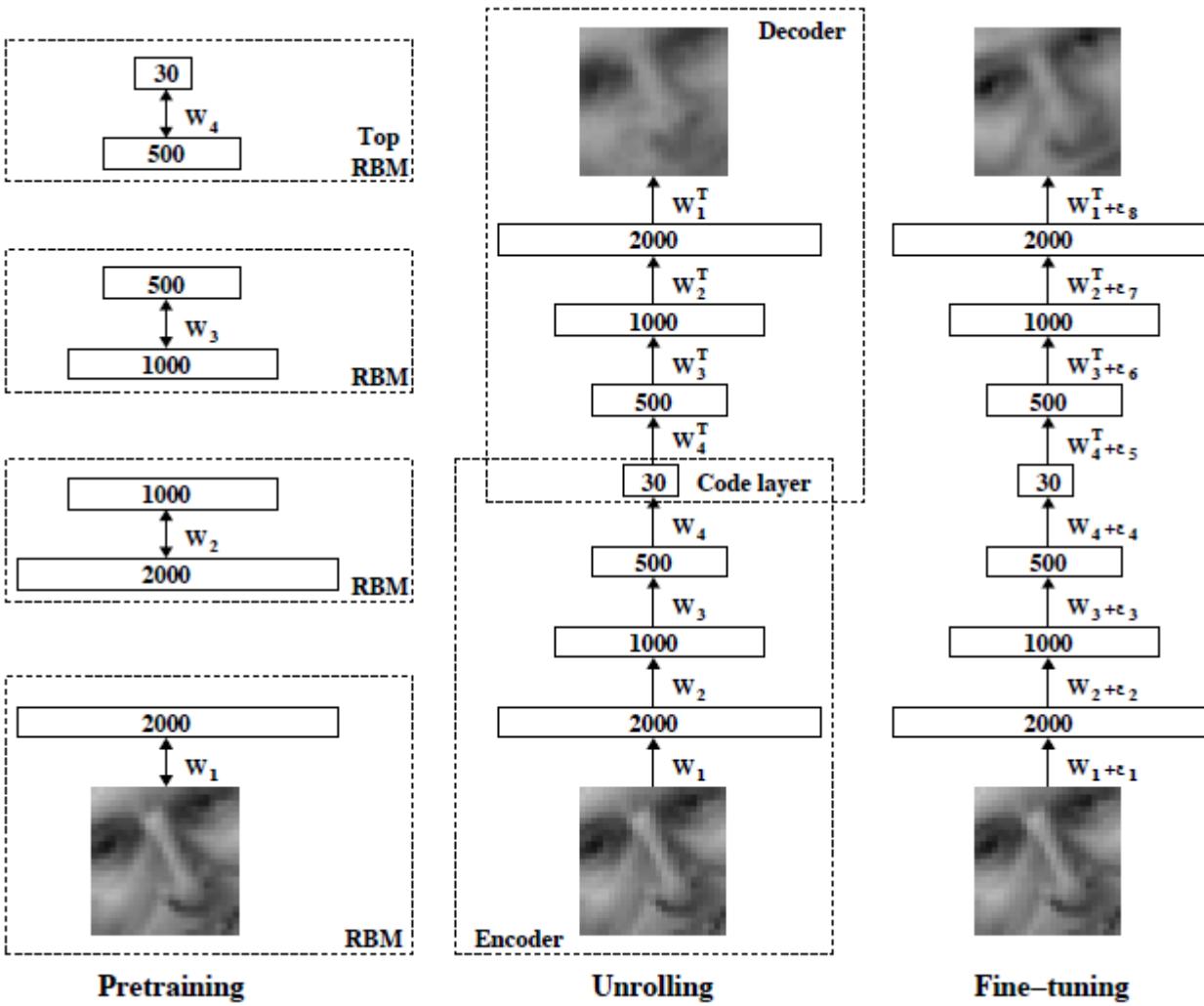
Features

Encoder

Top-down vs. bottom up?
Is there a more rigorous
mathematical formulation?

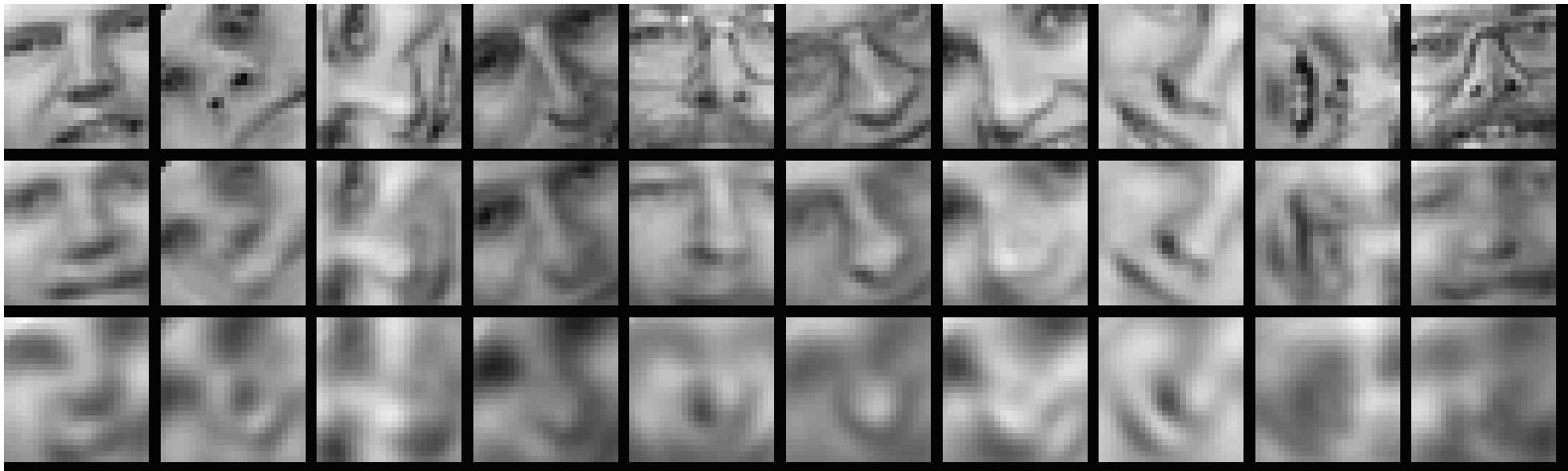
Input x

Deep Autoencoder



Deep Autoencoder

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

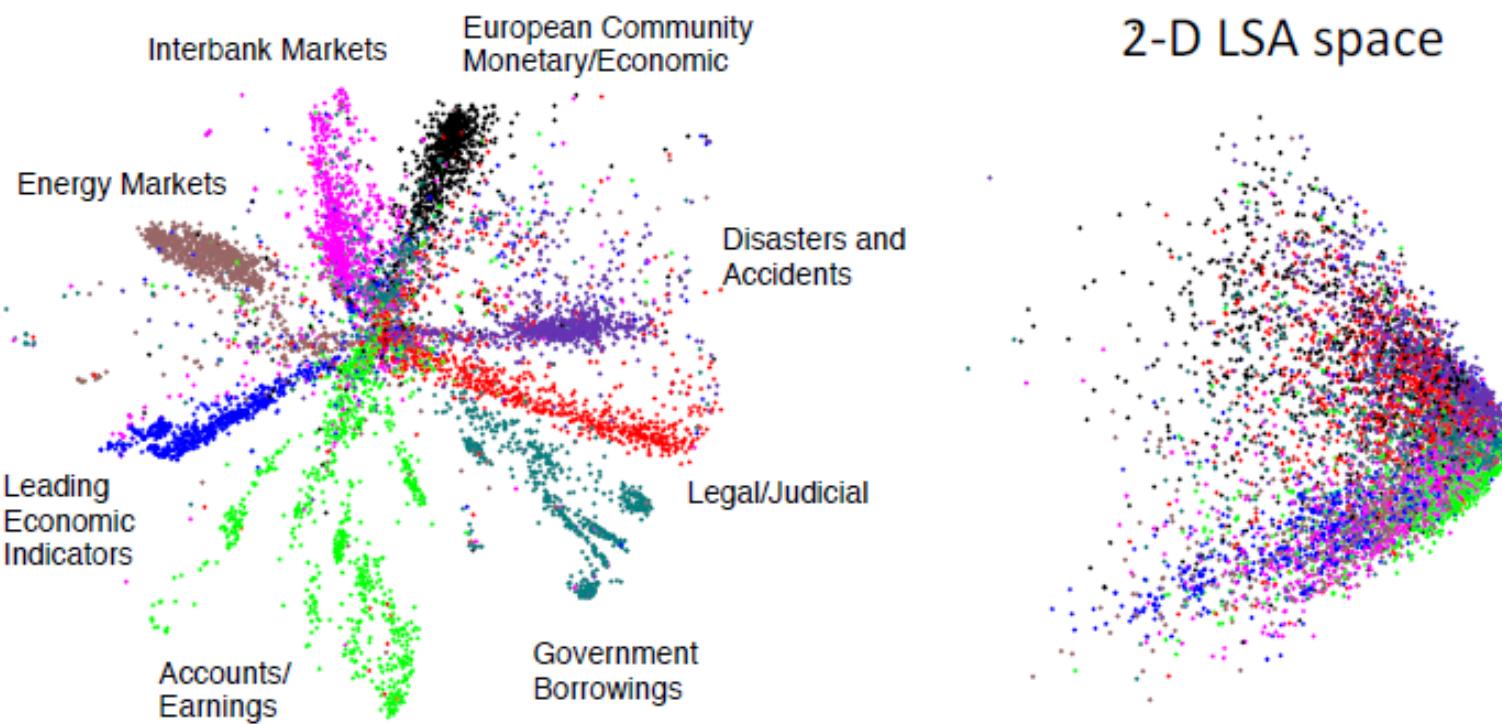


- **Top:** Random samples from the test dataset.
- **Middle:** Reconstructions by the 30-dimensional deep autoencoder.
- **Bottom:** Reconstructions by the 30-dimensional PCA

Deep Autoencoder

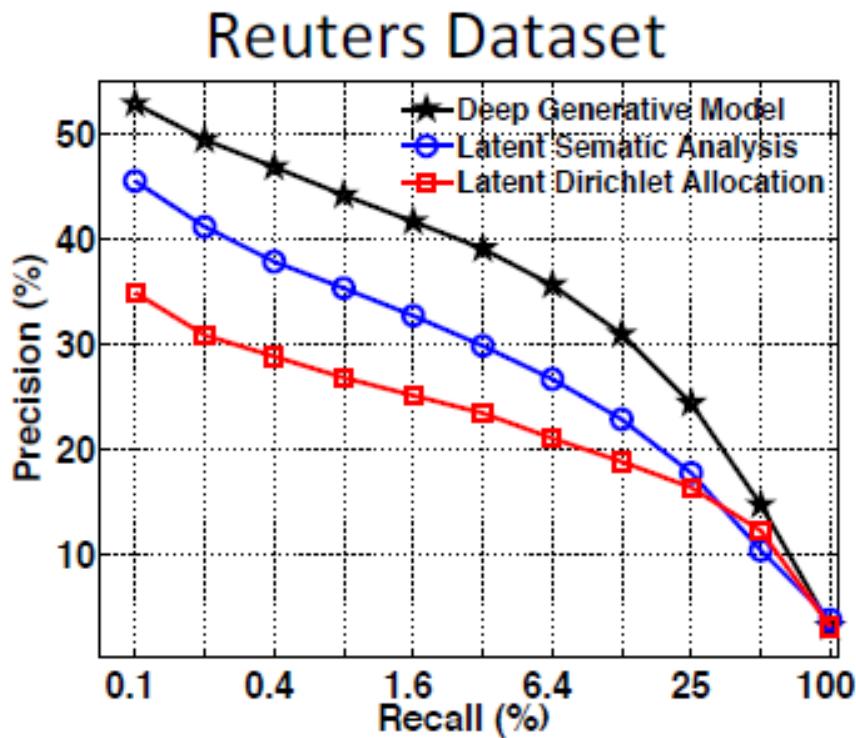
- Very difficult to optimize deep autoencoders using backpropagation
- Pre-training + fine-tuning
 - First train a stack of RBMs
 - Then “unroll” them
 - Then fine-tune with backpropagation

Information Retrieval



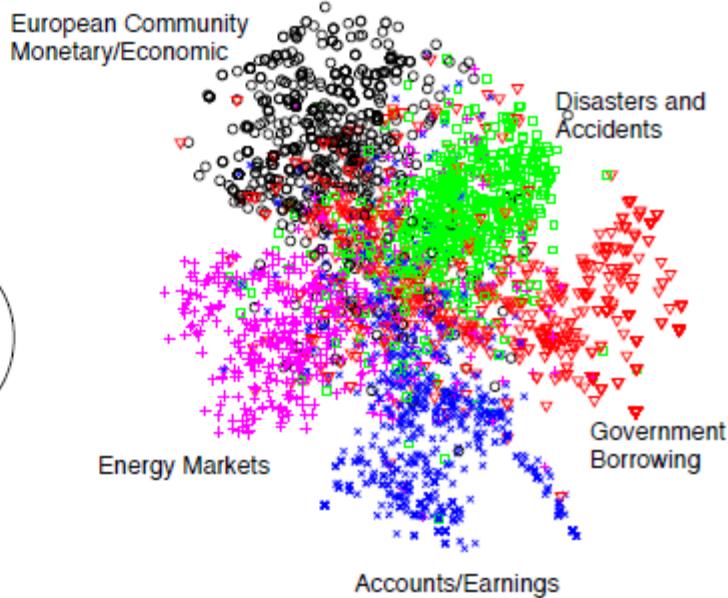
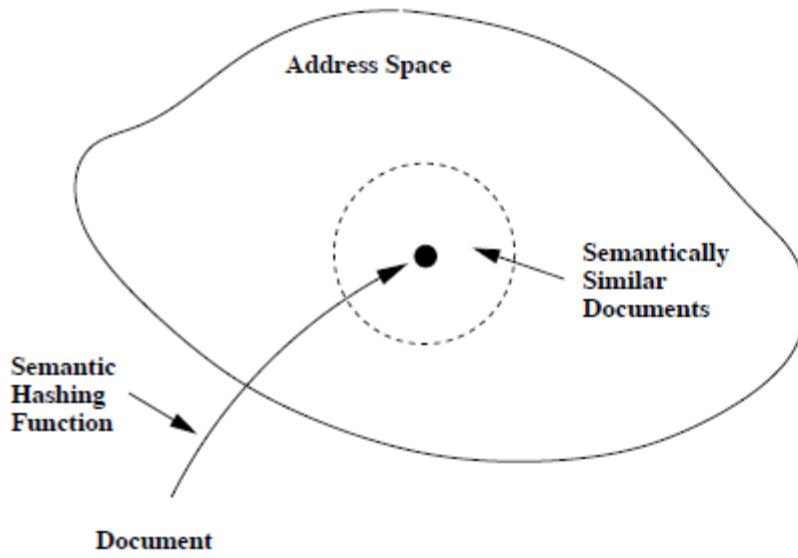
- The Reuters Corpus Volume II contains 804,414 newswire stories (randomly split into **402,207 training and 402,207 test**).
- “Bag-of-words” representation: each article is represented as a vector containing the counts of the most frequently used 2000 words in the training set.

Information Retrieval



- Reuters dataset: 804,414 newswire stories.
- Deep generative model significantly outperforms LSA and LDA topic models

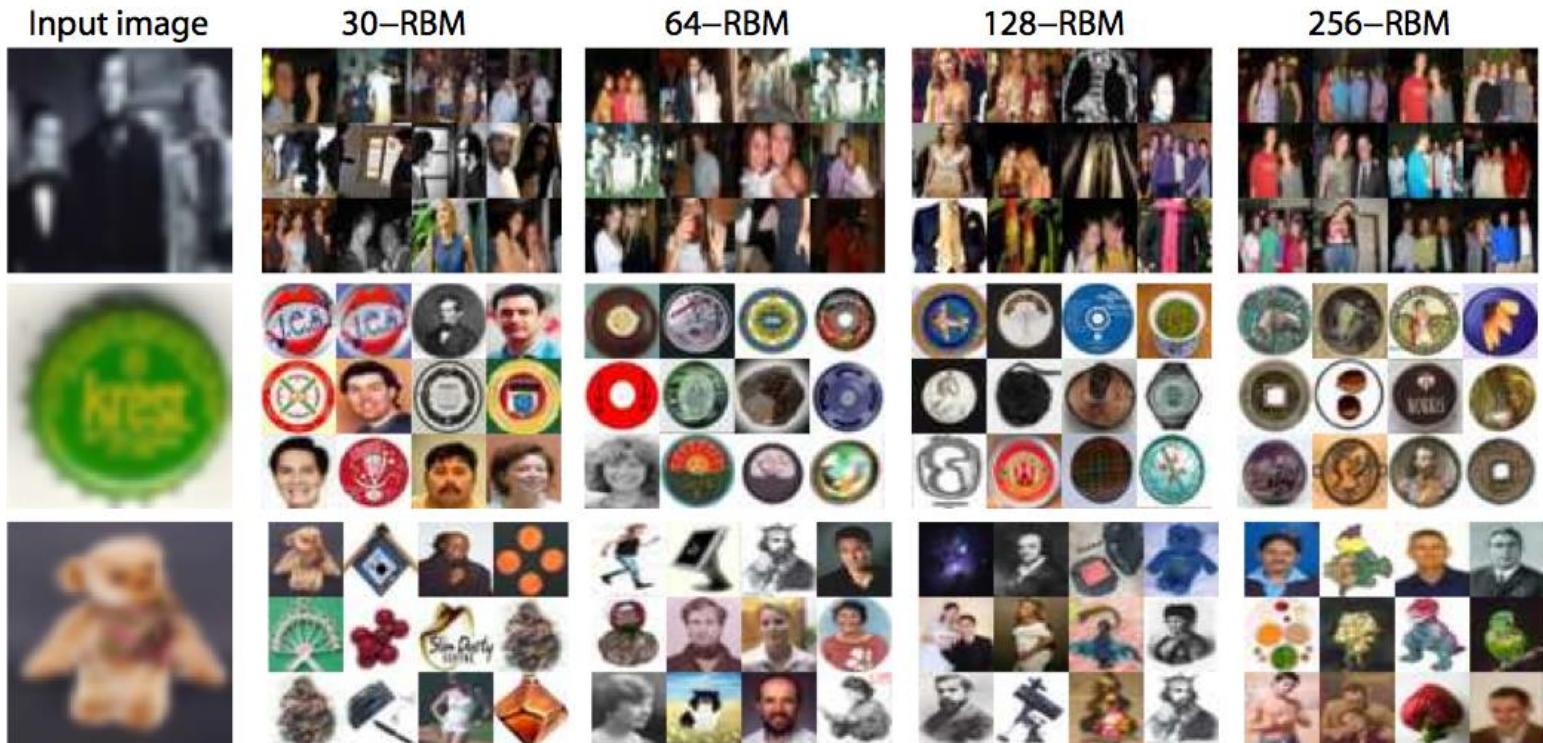
Semantic Hashing



- Learn to map documents into **semantic 20-D binary codes**.
- Retrieve similar documents stored at the nearby addresses **with no search at all**.

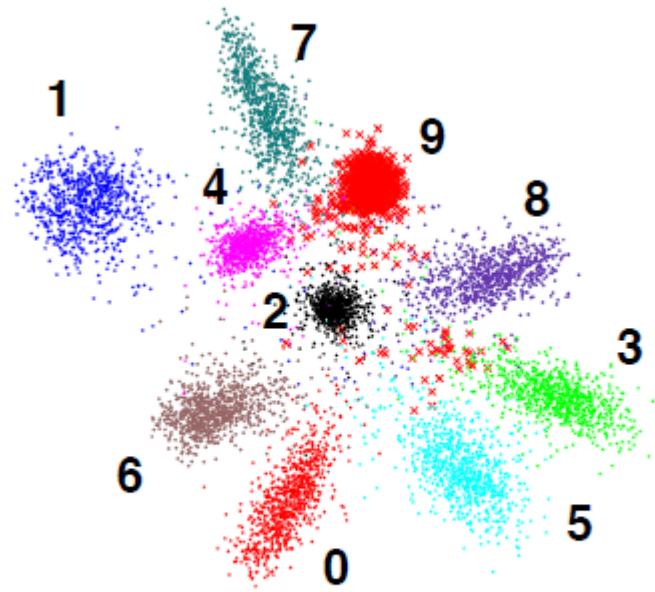
Searching Image Database using Binary Codes

- Map images into binary codes for fast retrieval



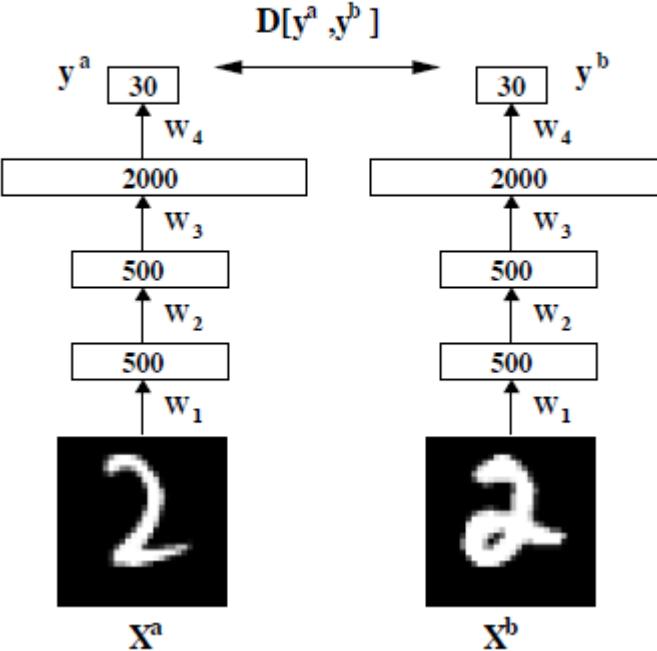
- Small Codes, Torralba, Fergus, Weiss, CVPR 2008
- Spectral Hashing, Y. Weiss, A. Torralba, R. Fergus, NIPS 2008
- Kulis and Darrell, NIPS 2009, Gong and Lazebnik, CVPR 2011
- Norouzi and Fleet, ICML 2011,

Learning Similarity Measures



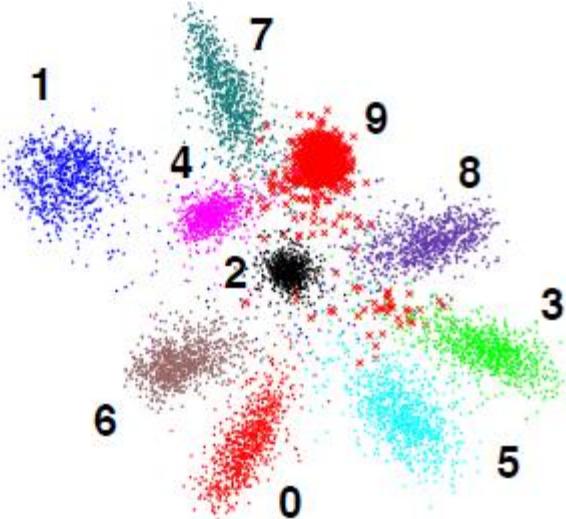
Related to Siamese Networks of LeCun.

Maximize the Agreement



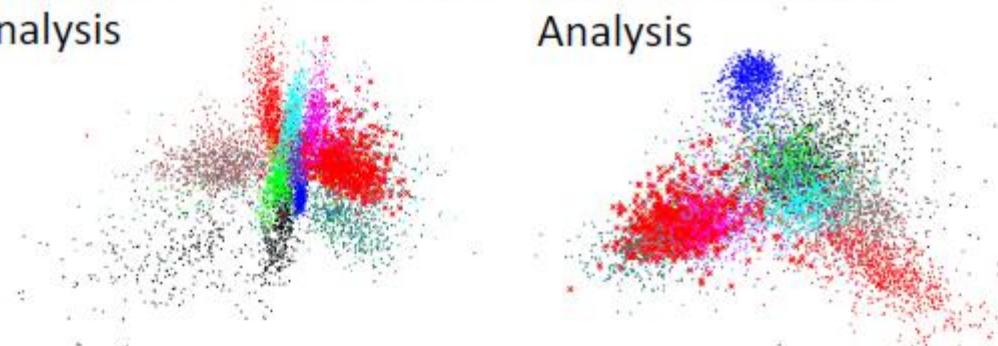
- Learn a nonlinear transformation of the input space.
- Optimize to make KNN perform well in the low-dimensional feature space

Learning Similarity Measures

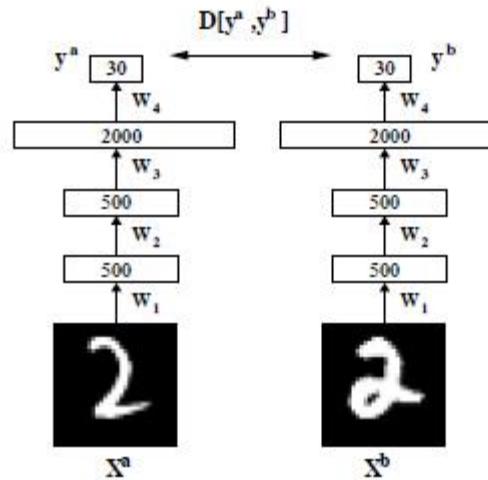


Neighborhood Component Analysis

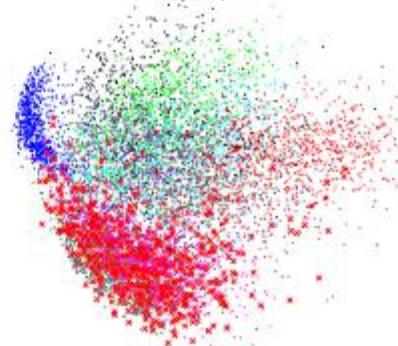
Linear discriminant Analysis



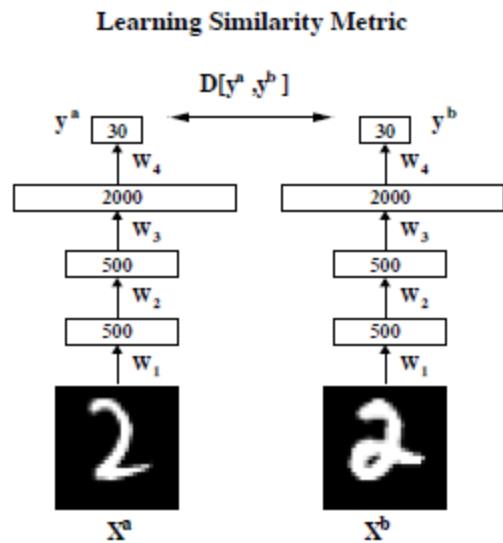
Learning Similarity Metric



PCA



Learning Similarity Measures



- As we change unit 25 in the code layer, ``3'' image turns into ``5'' image
- As we change unit 42 in the code layer, thick ``3'' image turns into skinny ``3''.

Outline

1 Course Review

2 Linear Factor Model

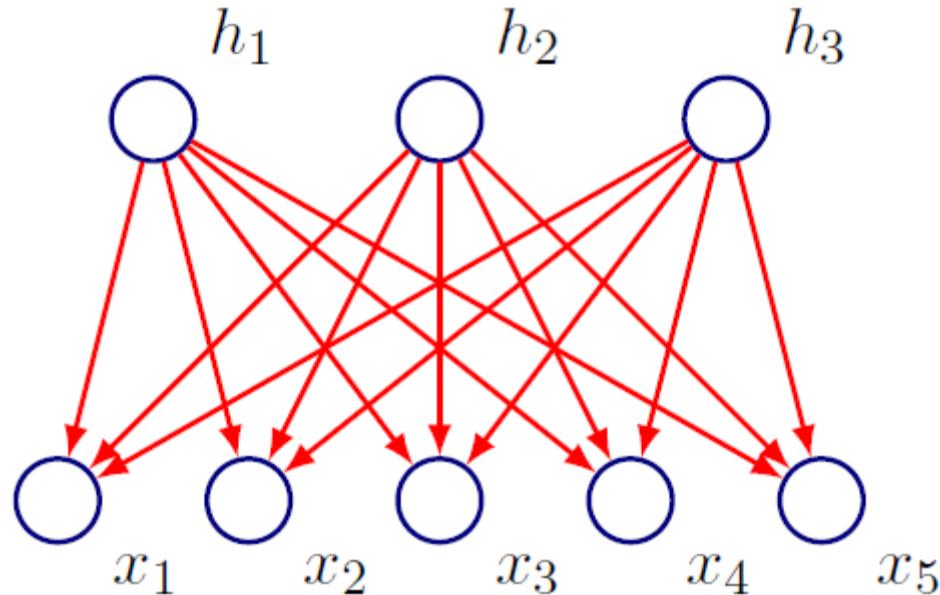
3 Autoencoder

4 DBN and RBM

More General Models

- Suppose $P(h)$ can not be assumed to have a nice Gaussian form
- The decoding of the input from the latent states can be a complicated non-linear function
- Estimation and inference can get complicated!

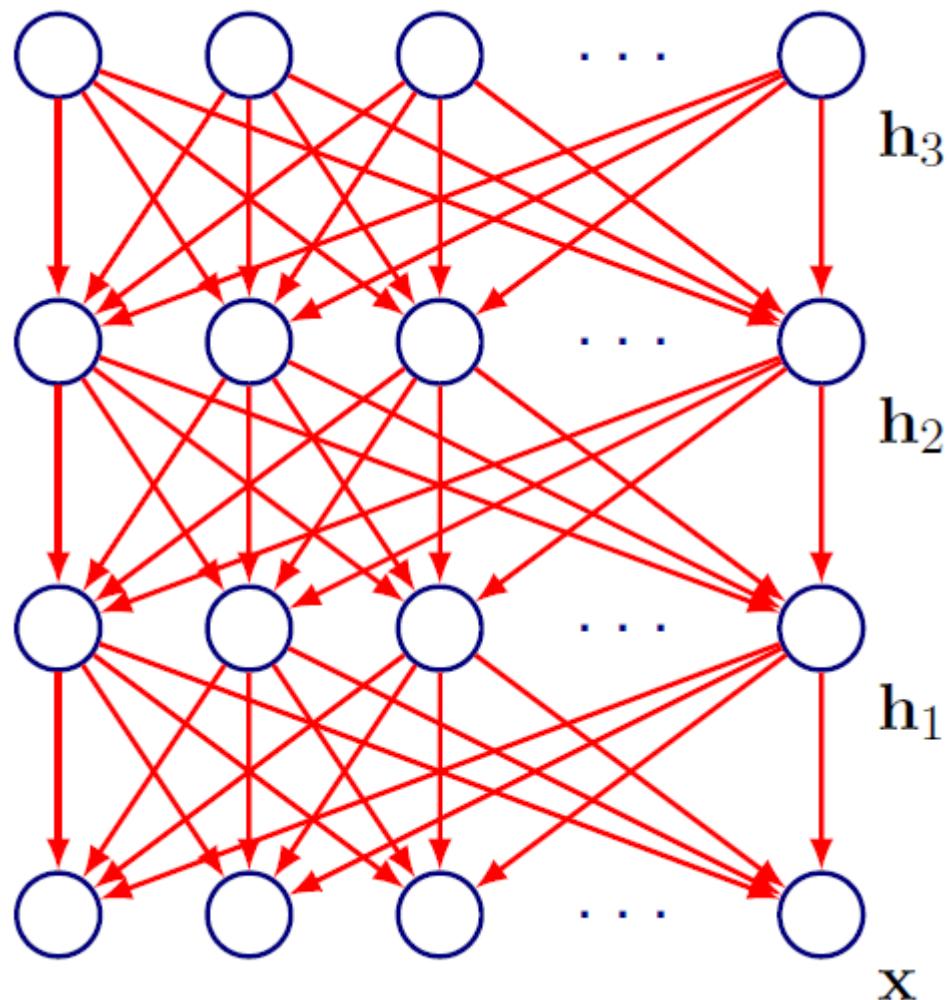
Earlier we had:



Quick Review

- Generative models can be modeled as directed graphical models
- The nodes represent random variables and arcs indicate dependency
- Some of the random variables are observed, others are hidden

Sigmoid Belief Networks



- Just like a feedforward network, but with arrows reversed.

Sigmoid Belief Networks

- Let $\mathbf{x} = \mathbf{h}^0$. Consider binary activations, then:

$$P(\mathbf{h}_i^k = 1 | \mathbf{h}^{k+1}) = \text{sigm}(b_i^k + \sum_j W_{i,j}^{k+1} \mathbf{h}_j^{k+1})$$

- The joint probability factorizes as:

$$P(\mathbf{x}, \mathbf{h}^1, \dots, \mathbf{h}^l) = P(\mathbf{h}^l) \left(\prod_{k=1}^{l-1} P(\mathbf{h}^k | \mathbf{h}^{k+1}) \right) P(\mathbf{x} | \mathbf{h}^1)$$

- Marginalization yields $P(\mathbf{x})$, intractable in practice except for very small models

Sigmoid Belief Networks

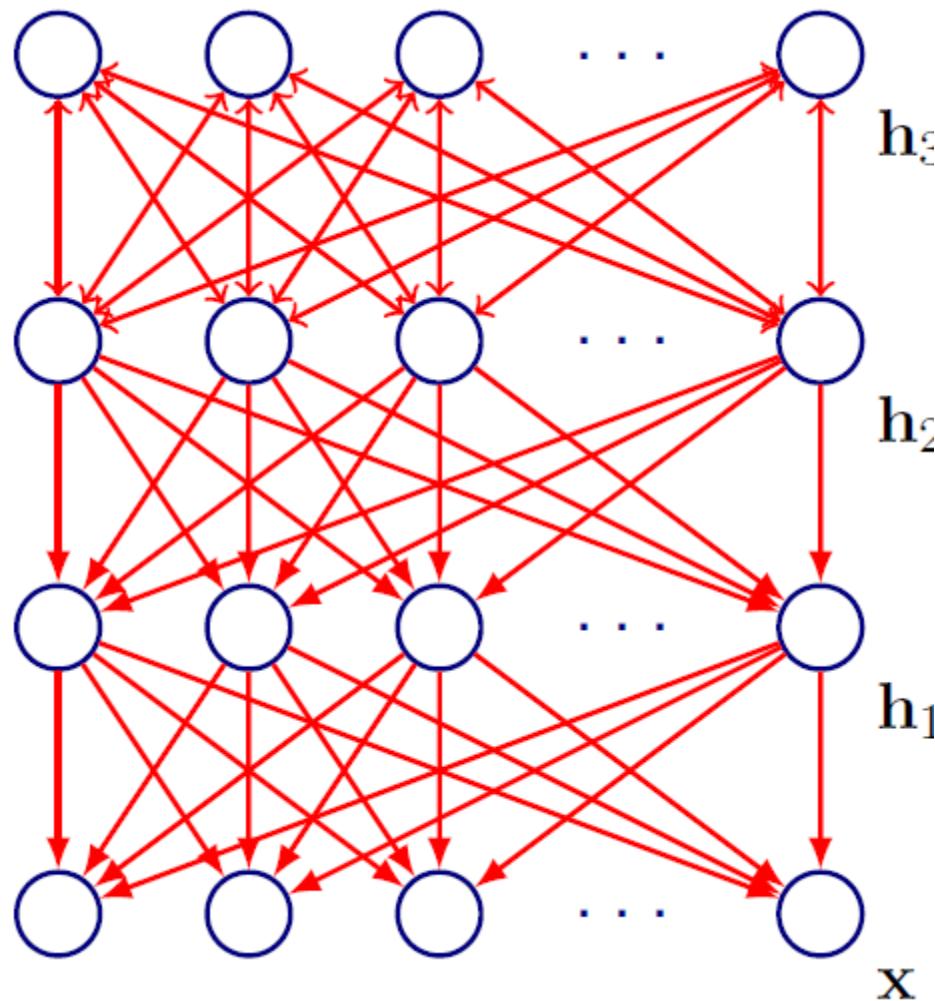
$$P(\mathbf{x}, \mathbf{h}^1, \dots, \mathbf{h}^l) = P(\mathbf{h}^l) \left(\prod_{k=1}^{l-1} P(\mathbf{h}^k | \mathbf{h}^{k+1}) \right) P(\mathbf{x} | \mathbf{h}^1)$$

- The top level prior is chosen as factorizable:
 $P(\mathbf{h}^l) = \prod_i P(\mathbf{h}_i^l)$
- A single (Bernoulli) parameter is needed for each \mathbf{h}_i in case of binary units
- Deep Belief Networks are like Sigmoid Belief Networks except for the top two layers

Sigmoid Belief Networks

- General case models are called Helmholtz Machines
- Two key references:
 - G. E. Hinton, P. Dayan, B. J. Frey, R. M. Neal: The Wake-Sleep Algorithm for Unsupervised Neural Networks, In Science, 1995
 - R. M. Neal: Connectionist Learning of Belief Networks, In Artificial Intelligence, 1992

Deep Belief Networks



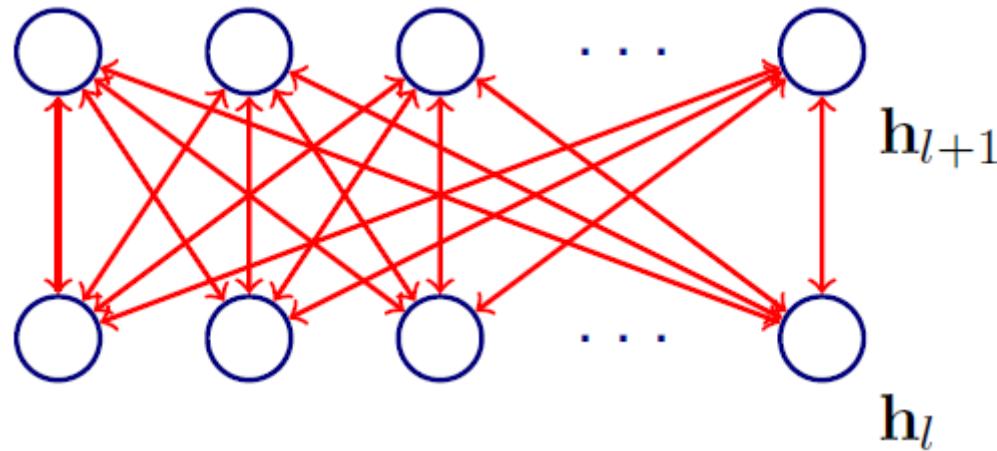
- The top two layers now have undirected edges

Deep Belief Networks

- The joint probability changes as:

$$P(\mathbf{x}, \mathbf{h}^1, \dots, \mathbf{h}^l) = P(\mathbf{h}^l, \mathbf{h}^{l-1}) \left(\prod_{k=1}^{l-2} P(\mathbf{h}^k | \mathbf{h}^{k+1}) \right) P(\mathbf{x} | \mathbf{h}^1)$$

Deep Belief Networks



- The top two layers are a Restricted Boltzmann Machine
- A RBM has the joint distribution:

$$P(\mathbf{h}^{l+1}, \mathbf{h}^l) \propto \exp(\mathbf{b}' \mathbf{h}^{l-1} + \mathbf{c}' \mathbf{h}^l + \mathbf{h}^l W \mathbf{h}^{l-1})$$

- We will return to RBMs and training procedures in a while, but first we look at the mathematical machinery that will make our task easier

Greedy Layer-wise Training of DBNs

- Reference: G. E. Hinton, S. Osindero and Y-W Teh: A Fast Learning Algorithm for Deep Belief Networks, In Neural Computation, 2006.
- First Step: Construct a RBM with input x and a hidden layer h , train the RBM
- Stack another layer on top of the RBM to form a new RBM. Fix W^1 , sample from $P(h^1|x)$, train W^2 as RBM
- Continue till k layers
- Implicitly defines $P(x)$ and $P(h)$ (variational bound justifies layerwise training)
- Can then be discriminatively fine-tuned using backpropagation

Energy Based Models

- Energy-Based Models assign a scalar energy with every configuration of variables under consideration
- Learning: Change the energy function so that its final shape has some desirable properties
- We can define a probability distribution through an energy:

$$P(\mathbf{x}) = \frac{\exp^{-(\text{Energy}(\mathbf{x}))}}{Z}$$

- Energies are in the log-probability domain:

$$\text{Energy}(\mathbf{x}) = \log \frac{1}{(ZP(\mathbf{x}))}$$

Energy Based Models

$$P(\mathbf{x}) = \frac{\exp^{-\text{Energy}(\mathbf{x})}}{Z}$$

- Z is a normalizing factor called the Partition Function

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

- How do we specify the energy function?

Product of Experts Formulation

- In this formulation, the energy function is:

$$\text{Energy}(\mathbf{x}) = \sum_i f_i(\mathbf{x})$$

- Therefore:

$$P(\mathbf{x}) = \frac{\exp^{-(\sum_i f_i(\mathbf{x}))}}{Z}$$

- We have the product of experts:

$$P(\mathbf{x}) \propto \prod_i P_i(\mathbf{x}) \propto \prod_i \exp^{(-f_i(\mathbf{x}))}$$

Product of Experts Formulation

$$P(\mathbf{x}) \propto \prod_i P_i(\mathbf{x}) \propto \prod_i \exp(-f_i(\mathbf{x}))$$

- Every expert f_i can be seen as enforcing a constraint on \mathbf{x}
- If f_i is large $\Rightarrow P_i(\mathbf{x})$ is small i.e. the expert thinks \mathbf{x} is implausible (constraint violated)
- If f_i is small $\Rightarrow P_i(\mathbf{x})$ is large i.e. the expert thinks \mathbf{x} is plausible (constraint satisfied)
- Contrast this with mixture models

Latent Variables

- x is observed, let's say h are **hidden factors** that explain x
- The probability then becomes:

$$P(x, h) = \frac{\exp^{-(\text{Energy}(x, h))}}{Z}$$

- We only care about the marginal:

$$P(x) = \sum_h \frac{\exp^{-(\text{Energy}(x, h))}}{Z}$$

Latent Variables

$$P(\mathbf{x}) = \sum_{\mathbf{h}} \frac{\exp^{-(\text{Energy}(\mathbf{x}, \mathbf{h}))}}{Z}$$

- We introduce another term in analogy from statistical physics:
free energy:

$$P(\mathbf{x}) = \frac{\exp^{-(\text{FreeEnergy}(\mathbf{x}))}}{Z}$$

- Free Energy is just a marginalization of energies in the log-domain:

$$\text{FreeEnergy}(\mathbf{x}) = - \log \sum_{\mathbf{h}} \exp^{-(\text{Energy}(\mathbf{x}, \mathbf{h}))}$$

Latent Variables

$$P(\mathbf{x}) = \frac{\exp^{-\text{FreeEnergy}(\mathbf{x})}}{Z}$$

- Likewise, the partition function:

$$Z = \sum_{\mathbf{x}} \exp^{-\text{FreeEnergy}(\mathbf{x})}$$

- We have an expression for $P(\mathbf{x})$ (and hence for the data log-likelihood). Let us see how the gradient looks like

Data Log-Likelihood Gradient

$$P(\mathbf{x}) = \frac{\exp^{-\text{FreeEnergy}(\mathbf{x})}}{Z}$$

- The gradient is simply working from the above:

$$\begin{aligned}\frac{\partial \log P(\mathbf{x})}{\partial \theta} &= -\frac{\partial \text{FreeEnergy}(\mathbf{x})}{\partial \theta} \\ &\quad + \frac{1}{Z} \sum_{\tilde{\mathbf{x}}} \exp^{-\text{FreeEnergy}(\tilde{\mathbf{x}})} \frac{\partial \text{FreeEnergy}(\tilde{\mathbf{x}})}{\partial \theta}\end{aligned}$$

- Note that $P(\tilde{\mathbf{x}}) = \exp^{-\text{FreeEnergy}(\tilde{\mathbf{x}})}$

Data Log-Likelihood Gradient

- The expected log-likelihood gradient over the training set has the following form:

$$\mathbb{E}_{\tilde{P}} \left[\frac{\partial \log P(\mathbf{x})}{\partial \theta} \right] = \mathbb{E}_{\tilde{P}} \left[\frac{\partial \text{FreeEnergy}(\mathbf{x})}{\partial \theta} \right] + \mathbb{E}_P \left[\frac{\partial \text{FreeEnergy}(\mathbf{x})}{\partial \theta} \right]$$

Data Log-Likelihood Gradient

$$\mathbb{E}_{\tilde{P}} \left[\frac{\partial \log P(\mathbf{x})}{\partial \theta} \right] = \mathbb{E}_{\tilde{P}} \left[\frac{\partial \text{FreeEnergy}(\mathbf{x})}{\partial \theta} \right] + \mathbb{E}_P \left[\frac{\partial \text{FreeEnergy}(\mathbf{x})}{\partial \theta} \right]$$

- \tilde{P} is the empirical training distribution
- Easy to compute!

Data Log-Likelihood Gradient

$$\mathbb{E}_{\tilde{P}} \left[\frac{\partial \log P(\mathbf{x})}{\partial \theta} \right] = \mathbb{E}_{\tilde{P}} \left[\frac{\partial \text{FreeEnergy}(\mathbf{x})}{\partial \theta} \right] + \mathbb{E}_P \left[\frac{\partial \text{FreeEnergy}(\mathbf{x})}{\partial \theta} \right]$$

- P is the model distribution (exponentially many configurations!)
- Usually very hard to compute!
- Resort to Markov Chain Monte Carlo to get a stochastic estimator of the gradient

A Special Case

- Suppose the energy has the following form:

$$\text{Energy}(\mathbf{x}, \mathbf{h}) = -\beta(\mathbf{x}) + \sum_i \gamma_i(\mathbf{x}, \mathbf{h}_i)$$

- The free energy, and numerator of log likelihood can be computed tractably!
- What is $P(\mathbf{x})$?
- What is the FreeEnergy(\mathbf{x})?

A Special Case

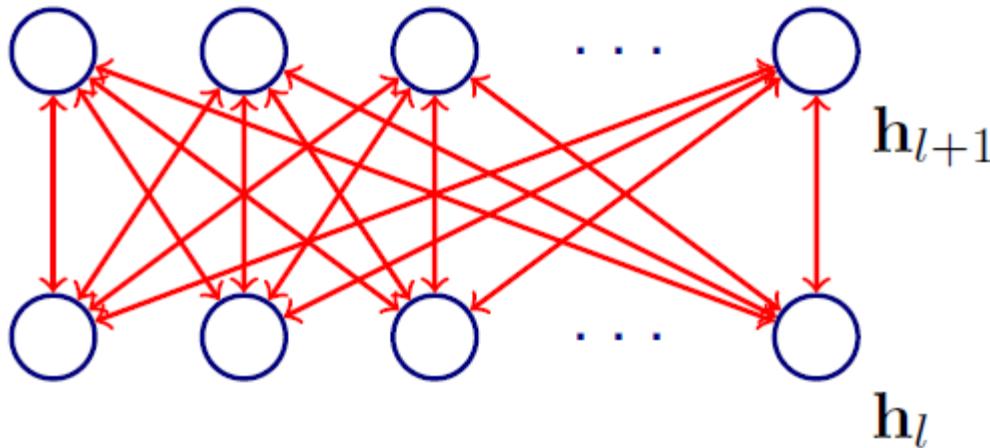
- The likelihood term:

$$P(\mathbf{x}) = \frac{\exp^{\beta(\mathbf{x})}}{Z} \prod_i \sum_{\mathbf{h}_i} \exp^{-\gamma_i(\mathbf{x}, \mathbf{h}_i)}$$

- The Free Energy term:

$$\begin{aligned}\text{FreeEnergy}(\mathbf{x}) &= -\log P(\mathbf{x}) - \log Z \\ &= -\beta - \sum_i \log \sum_{\mathbf{h}_i} \exp^{-\gamma_i(\mathbf{x}, \mathbf{h}_i)}\end{aligned}$$

Restricted Boltzmann Machines

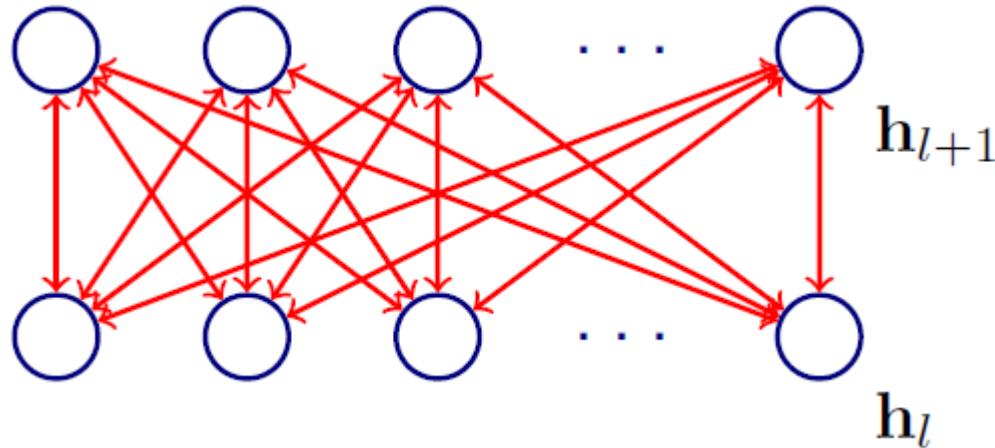


- Recall the form of energy:

$$\text{Energy}(\mathbf{x}, \mathbf{h}) = -\mathbf{b}^T \mathbf{x} - \mathbf{c}^T \mathbf{h} - \mathbf{h}^T W \mathbf{x}$$

- Takes the earlier nice form with $(\mathbf{x}) = \mathbf{b}^T \mathbf{x}$ and
 $\gamma_i(\mathbf{x}, \mathbf{h}_i) = \mathbf{h}_i(\mathbf{c}_i + W_i \mathbf{x})$
- Originally proposed by Smolensky (1987) who called them Harmoniums as a special case of Boltzmann Machines

Restricted Boltzmann Machines



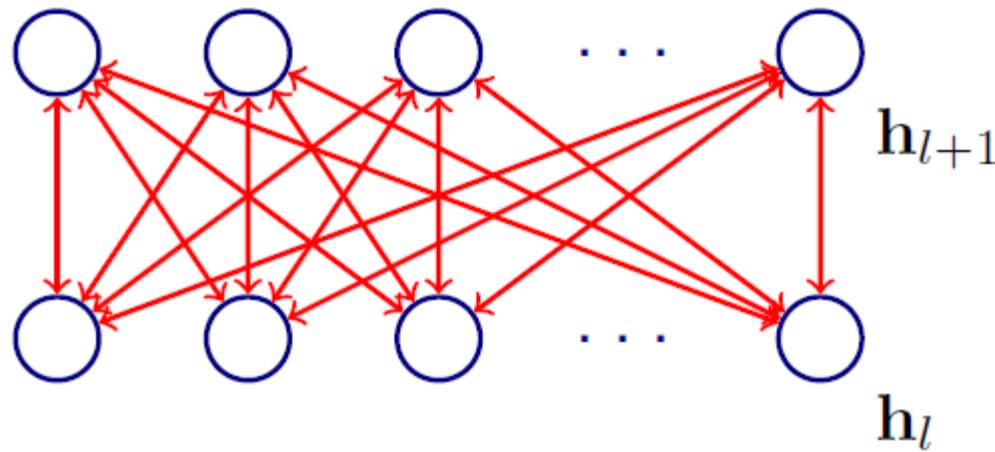
- As seen before, the Free Energy can be computed efficiently:

$$\text{FreeEnergy}(\mathbf{x}) = -\mathbf{b}^T \mathbf{x} - \sum_i \log \sum_{\mathbf{h}_i} \exp^{\mathbf{h}_i(\mathbf{c}_i + \mathbf{W}_i \mathbf{x})}$$

- The conditional probability:

$$P(\mathbf{h}|\mathbf{x}) = \frac{\exp (\mathbf{b}^T \mathbf{x} + \mathbf{c}^T \mathbf{h} + \mathbf{h}^T \mathbf{W} \mathbf{x})}{\sum_{\tilde{\mathbf{h}}} \exp (\mathbf{b}^T \mathbf{x} + \mathbf{c}^T \tilde{\mathbf{h}} + \tilde{\mathbf{h}}^T \mathbf{W} \mathbf{x})} = \prod_i P(\mathbf{h}_i|\mathbf{x})$$

Restricted Boltzmann Machines



- x and h play symmetric roles:

$$P(\mathbf{x}|\mathbf{h}) = \prod_i P(x_i|\mathbf{h})$$

- The common transfer (for the binary case):

$$P(h_i = 1|\mathbf{x}) = \sigma(\mathbf{c}_i + W_i \mathbf{x})$$

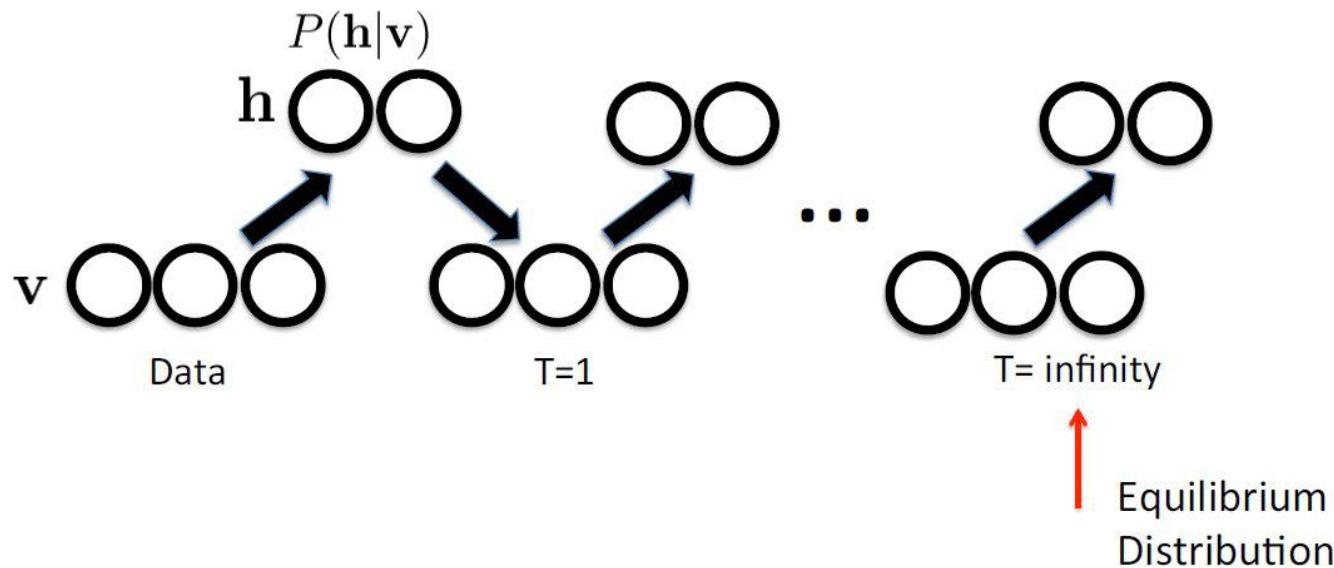
$$P(x_j = 1|h) = \sigma(\mathbf{b}_j + W_{:,j}^T \mathbf{h})$$

Approximate Learning and Gibbs Sampling

$$\mathbb{E}_{\tilde{P}} \left[\frac{\partial \log P(\mathbf{x})}{\partial \theta} \right] = \mathbb{E}_{\tilde{P}} \left[\frac{\partial \text{FreeEnergy}(\mathbf{x})}{\partial \theta} \right] + \mathbb{E}_P \left[\frac{\partial \text{FreeEnergy}(\mathbf{x})}{\partial \theta} \right]$$

- We saw the expression for Free Energy for a RBM. But the second term was intractable. How do learn in this case?
- Replace the average over all possible input configurations by samples
- Run Markov Chain Monte Carlo (Gibbs Sampling):
- First sample $\mathbf{x}_1 \sim P(\mathbf{x})$, then $\mathbf{h}_1 \sim P(\mathbf{h}|\mathbf{x}_1)$, then $\mathbf{X}_2 \sim P(\mathbf{x}|\mathbf{h}_1)$, then $\mathbf{h}_2 \sim P(\mathbf{h}|\mathbf{x}_2)$ till \mathbf{x}_{k+1}

Approximate Learning, Alternating Gibbs Sampling



- We have already seen: $P(x|h) = \prod_i P(x_i|h)$
$$P(h|x) = \prod_i P(h_i|x)$$
- With: $P(h_i = 1|x) = \sigma(c_i + W_i x)$ and
$$P(x_j = 1|h) = \sigma(b_j + W_{:,j}^T h)$$

Training RBM: Contrastive Divergence Algorithm

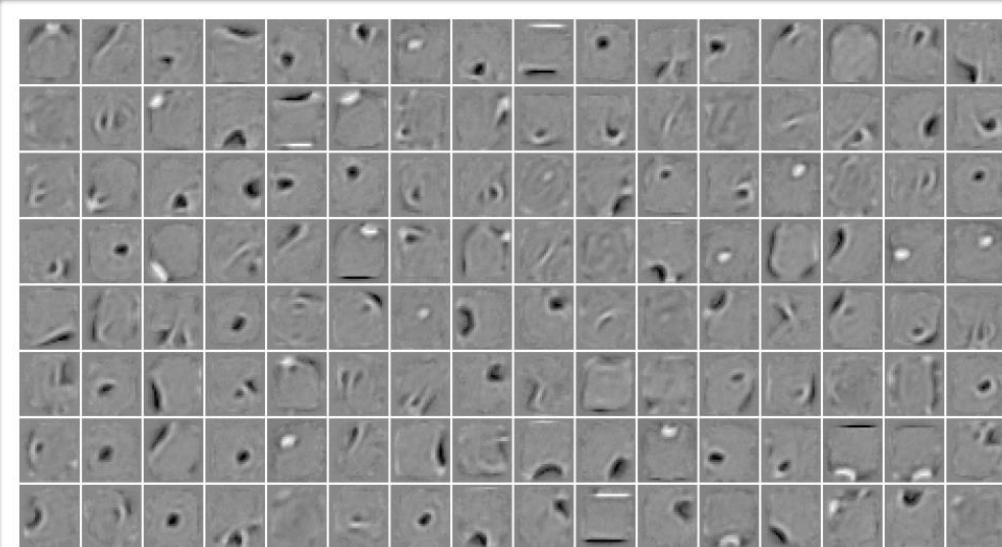
- Start with a training example on the visible units
- Update all the hidden units in parallel
- Update all the visible units in parallel to obtain a reconstruction
- Update all the hidden units again
- Update model parameters
- Aside: Easy to extend RBM (and contrastive divergence) to the continuous case

Example: MNIST

Original images:

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

Learned features:



(Larochelle et al., JMLR 2009)

Acknowledgement

Some of the materials in these slides are drawn inspiration from:

- Shubhendu Trivedi and Risi Kondor, University of Chicago, Deep Learning Course
- Hung-yi Lee, National Taiwan University, Machine Learning and having it Deep and Structured course
- Xiaogang Wang, The Chinese University of Hong Kong, Deep Learning Course
- Fei-Fei Li, Standord University, CS231n Convolutional Neural Networks for Visual Recognition course

Next time

- Model Computation and Inference

Questions?

Thank You !



WeChat Group for Deep Learning