

# NONLINEAR DIMENSIONALITY REDUCTION

# 36

## CHAPTER CONTENTS

<b>Dimensionality Reduction with Kernel Trick</b> .....	429
Kernel PCA .....	429
Laplacian Eigenmap .....	433
<b>Supervised Dimensionality Reduction with Neural Networks</b> .....	435
<b>Unsupervised Dimensionality Reduction with Autoencoder</b> .....	436
Autoencoder .....	436
Training by Gradient Descent .....	437
Sparse Autoencoder .....	439
<b>Unsupervised Dimensionality Reduction with Restricted Boltzmann Machine</b> ....	440
Model .....	441
Training by Gradient Ascent .....	442
<b>Deep Learning</b> .....	446

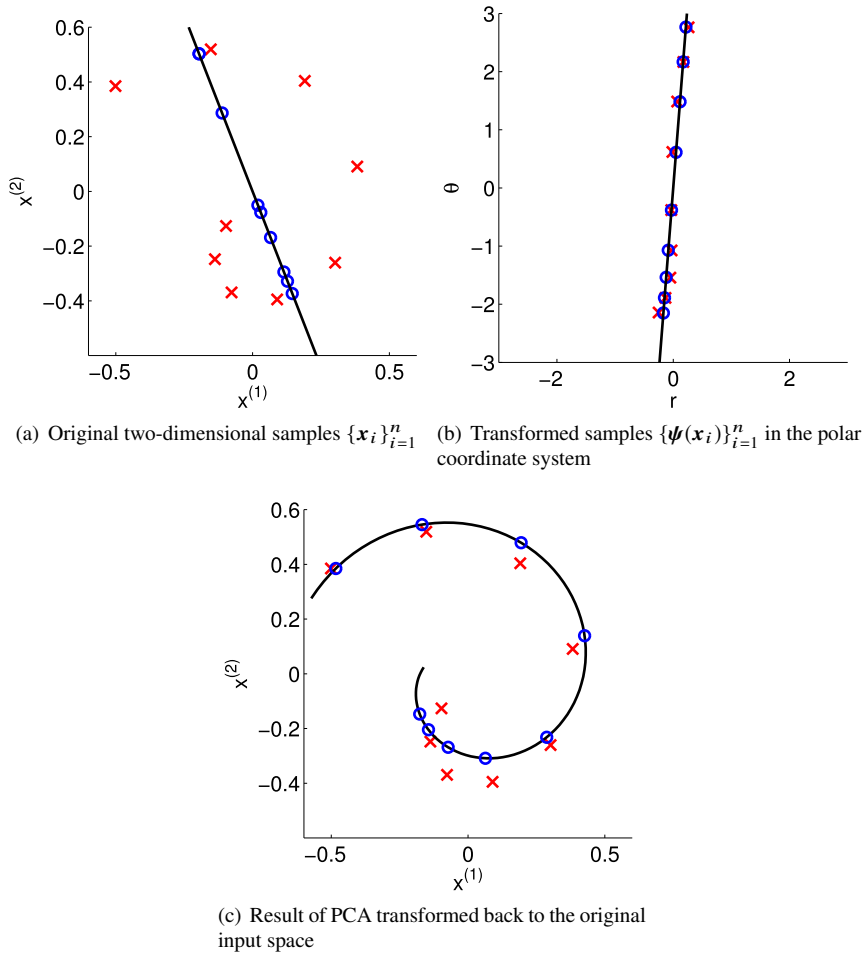
In this chapter, supervised and unsupervised methods of *nonlinear dimensionality reduction* are introduced, including approaches based on kernels and neural networks.

## 36.1 DIMENSIONALITY REDUCTION WITH KERNEL TRICK

The linear dimensionality reduction methods explained in the previous chapter handle training samples only in terms of their inner products. This means that the *kernel trick* (see Section 27.4) is applicable to obtain nonlinear methods. More specifically, training input samples  $\{\mathbf{x}_i\}_{i=1}^n$  are first transformed by a nonlinear mapping  $\psi$  and then ordinary linear dimensionality reduction algorithms are applied in the transformed space, which corresponds to nonlinear dimensionality reduction in the original input space. In this section, kernel-based nonlinear dimensionality reduction methods are introduced.

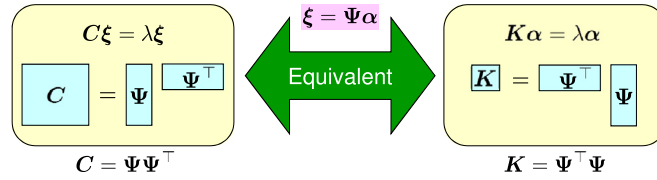
### 36.1.1 KERNEL PCA

Let us illustrate the behavior of PCA in a feature space. The original two-dimensional samples  $\mathbf{x}_i = (x_i^{(1)}, x_i^{(2)})^\top$  are plotted in Fig. 36.1(a). Due to the curved spiral shape,

**FIGURE 36.1**

Nonlinear PCA in a feature space. “x” denotes a sample, the solid line denotes the one-dimensional embedding subspace found by PCA, and “o” denotes a projected sample.

linear PCA does not work properly for this data set. If the samples are transformed into the *polar system* (i.e. radius  $r$  and angle  $\theta$ ), the spiral shape can be nicely unfolded, as plotted in Fig. 36.1(b). Then linear PCA in the transformed space can well capture the global structure of the unfolded data. Finally, transforming back the projected samples to the original input space gives nice nonlinear projection of the original spiral-shape samples, as plotted in Fig. 36.1(c).

**FIGURE 36.2**

Eigenvalue problems for PCA. Appropriately choosing the expression of eigenvalue problem depending on whether matrix  $\Psi$  is fat or skinny allows us to reduce the computational costs.

If the dimensionality of the feature space is not high, directly performing PCA as explained above is fine. That is, the eigenvalue problem for the total scatter matrix  $C$  in the feature space is solved:

$$C\xi = \lambda\xi,$$

where

$$C = \sum_{i=1}^n \psi(x_i)\psi(x_i)^\top.$$

However, if the dimensionality of the feature space is high, the computational complexity grows significantly. More extremely, if the Gaussian kernel,

$$K(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2h^2}\right), \quad (36.1)$$

is used as a nonlinear mapping, the dimensionality of the feature space is actually infinite, and thus explicitly performing PCA in the feature space is not possible.

To cope with this problem, let us consider the eigenvalue problem for the *kernel matrix*  $K$  [88]:

$$K\alpha = \lambda\alpha,$$

where the  $(i, i')$ th element of  $K$  is defined as

$$K_{i, i'} = K(\mathbf{x}_i, \mathbf{x}_{i'}) = \langle \psi(\mathbf{x}_i), \psi(\mathbf{x}_{i'}) \rangle.$$

The matrices  $C$  and  $K$  can be expressed by using the *design matrix*

$$\Psi = (\psi(\mathbf{x}_1), \dots, \psi(\mathbf{x}_n))$$

as

$$C = \Psi\Psi^\top \quad \text{and} \quad K = \Psi^\top\Psi.$$

This implies that the eigenvalues of  $\mathbf{C}$  and  $\mathbf{K}$  are actually common. Furthermore, eigenvector  $\alpha$  of  $\mathbf{K}$  and eigenvector  $\xi$  of  $\mathbf{C}$  are related to each other as follows (see Fig. 36.2):

$$\xi = \Psi\alpha \quad \text{and} \quad \alpha = \Psi^\top \xi.$$

While the size of covariance matrix  $\mathbf{C}$  depends on the dimensionality of the feature space, the size of kernel matrix  $\mathbf{K}$  depends only on the number of samples and is independent of the dimensionality of the feature space. Thus, if the dimensionality of the feature space is larger than the number of samples, solving the eigenvalue problem with  $\mathbf{K}$  is computationally more efficient.

As explained in Section 35.2.1, PCA requires centering of samples. However, when the eigenvalue problem with  $\mathbf{K}$  is solved, feature vectors  $\{\psi(x_i)\}_{i=1}^n$  are not explicitly handled and thus feature vectors cannot be centralized directly. In kernel PCA, the kernel matrix  $\mathbf{K}$  is *implicitly* centralized as

$$\mathbf{K} \leftarrow \mathbf{H}\mathbf{K}\mathbf{H},$$

where

$$\mathbf{H} = \mathbf{I}_n - \frac{1}{n}\mathbf{1}_{n \times n}$$

is the *centering matrix*,  $\mathbf{I}_n$  denotes the  $n$ -dimensional identity matrix, and  $\mathbf{1}_{n \times n}$  is the  $n \times n$  matrix with all ones.

Another issue to be considered is that kernel PCA requires eigenvectors to be normalized as  $\|\xi_j\| = 1$ , but solving the eigenvalue problem with  $\mathbf{K}$  usually gives eigenvectors such that  $\|\alpha_j\| = 1$ . Therefore, normalization of eigenvectors should be carried out explicitly as

$$\alpha_j \leftarrow \frac{1}{\sqrt{\lambda_j}} \alpha_j \quad \text{for } j = 1, \dots, m,$$

which comes from

$$\|\xi_j\| = \sqrt{\|\xi_j\|^2} = \sqrt{\|\Psi\alpha_j\|^2} = \sqrt{\langle \Psi^\top \Psi \alpha_j, \alpha_j \rangle} = \sqrt{\langle \mathbf{K} \alpha_j, \alpha_j \rangle} = \sqrt{\lambda_j}.$$

Summarizing the above discussions, final embedding solutions of samples  $\{x_i\}_{i=1}^n$  by kernel PCA are given by

$$(z_1, \dots, z_n) = \left( \frac{1}{\sqrt{\lambda_1}} \alpha_1, \dots, \frac{1}{\sqrt{\lambda_m}} \alpha_m \right)^\top \mathbf{H}\mathbf{K}\mathbf{H},$$

where  $\alpha_1, \dots, \alpha_m$  are eigenvectors of  $\mathbf{H}\mathbf{K}\mathbf{H}$  corresponding to the  $m$  largest eigenvalues. Similarly, an embedding solution of a new sample  $x$  by kernel PCA is given by

$$z = \left( \frac{1}{\sqrt{\lambda_1}} \alpha_1, \dots, \frac{1}{\sqrt{\lambda_m}} \alpha_m \right)^\top \mathbf{H} \left( k - \frac{1}{n} \mathbf{K} \mathbf{1}_n \right),$$

```

n=200; a=linspace(0,pi,n/2);
u=[a.*cos(a) (a+pi).*cos(a)]';
v=[a.*sin(a) (a+pi).*sin(a)]';
x=[u v]; y=[ones(1,n/2) 2*ones(1,n/2)]';

x2=sum(x.^2,2); hh=2*2^2; H=eye(n)-ones(n,n)/n;
K=exp(-(repmat(x2,1,n)+repmat(x2',n,1)-2*x*x')/hh);
G=H*K*H; [A,L]=eigs(G,2); z=(diag(diag(L).^(-1/2))*A'*G)';

figure(1); clf; hold on;
plot(z(y==1,1),z(y==1,2),'bo');
plot(z(y==2,1),z(y==2,2),'rx');

```

**FIGURE 36.3**

MATLAB code for kernel PCA with Gaussian kernels.

where  $\mathbf{1}_n$  denotes the  $n$ -dimensional vectors with all ones and

$$\mathbf{k} = (K(\mathbf{x}, \mathbf{x}_1), \dots, K(\mathbf{x}, \mathbf{x}_n))^T.$$

A MATLAB code of PCA for Gaussian kernel (36.1) is provided in Fig. 36.3, and its behavior is illustrated in Fig. 36.4. This shows that, while the original two-dimensional samples are not linearly separable, PCA for Gaussian kernel with width  $h = 2$  gives linearly separable embedding in the feature space.

### 36.1.2 LAPLACIAN EIGENMAP

A kernelized version of locality preserving projection is called the *Laplacian eigenmap* [13]. In the generalized eigenvalue problem of locality preserving projection,

$$\mathbf{X} \mathbf{L} \mathbf{X}^T \boldsymbol{\xi} = \lambda \mathbf{X} \mathbf{D} \mathbf{X}^T \boldsymbol{\xi},$$

multiplying  $\mathbf{X}^T$  from the left-hand side and letting  $\boldsymbol{\xi} = \mathbf{X} \boldsymbol{\beta}$  yield

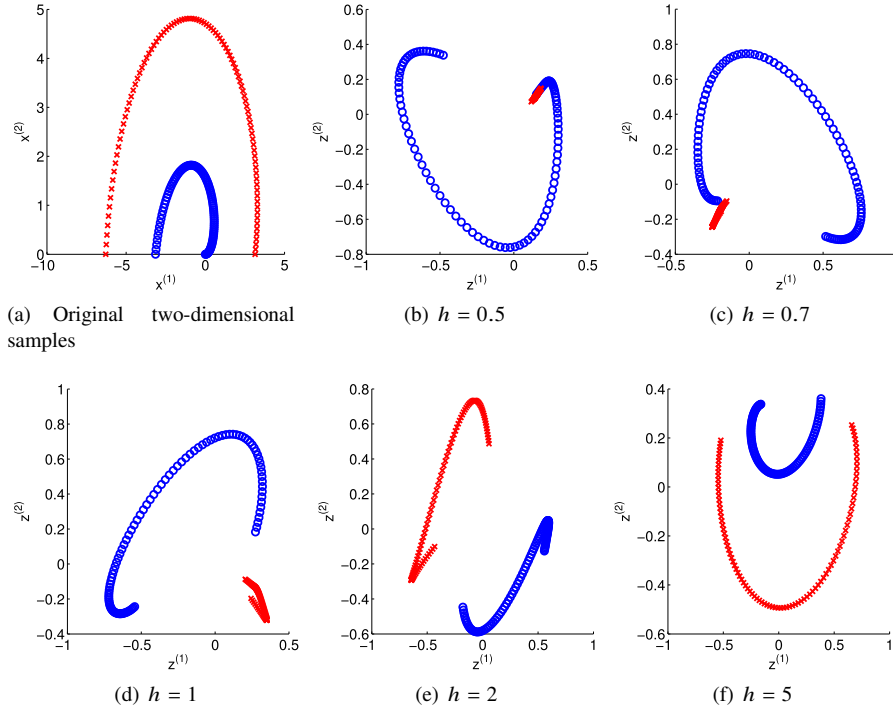
$$\mathbf{X}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \lambda \mathbf{X}^T \mathbf{X} \mathbf{D} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta}.$$

If  $\mathbf{X}^T \mathbf{X}$  is replaced with kernel matrix  $\mathbf{K}$ , a kernelized version of locality preserving projection is obtained as

$$\mathbf{K} \mathbf{L} \mathbf{K} \boldsymbol{\beta} = \lambda \mathbf{K} \mathbf{D} \mathbf{K} \boldsymbol{\beta}.$$

If  $\mathbf{K}$  is invertible, multiplying  $\mathbf{K}^{-1}$  from the left-hand side and letting  $\boldsymbol{\alpha} = \mathbf{K} \boldsymbol{\beta}$  yield

$$\mathbf{L} \boldsymbol{\alpha} = \lambda \mathbf{D} \boldsymbol{\alpha}, \tag{36.2}$$

**FIGURE 36.4**

Examples of kernel PCA with Gaussian kernels. Original two-dimensional samples are transformed to infinite-dimensional feature space by Gaussian kernels with width  $h$ , and then PCA is applied to reduce the dimensionality to two.

which is the eigenvalue problem that the Laplacian eigenmap solves. Let  $\lambda_1 \geq \dots \geq \lambda_n \geq 0$  and  $\alpha_1, \dots, \alpha_n$  be the generalized eigenvalues and generalized eigenvectors of Eq. (36.2). Then the definition of graph Laplacian matrix  $L$ , Eq. (35.2), yields  $L\mathbf{1}_n = \mathbf{0}_n$ . This means that  $\lambda_n = 0$  and  $\alpha_n = \mathbf{1}_n$  hold, which is trivial. In the Laplacian eigenmap, this trivial eigenvector is removed and the final embedding solution is obtained as

$$(\mathbf{z}_1, \dots, \mathbf{z}_n) = (\alpha_{n-1}, \alpha_{n-2}, \dots, \alpha_{n-m})^\top.$$

If the similarity matrix  $W$  is *sparse* (see Fig. 35.8),  $L = D - W$  is also sparse and thus eigenvalue problem (36.2) can be solved efficiently even if  $n$  is large.

A MATLAB code of the Laplacian eigenmap for 10-nearest neighbor similarity (see Fig. 35.8) is provided in Fig. 36.5, and its behavior is illustrated in Fig. 36.6. This shows that the “swiss roll” structure can be nicely unfolded.

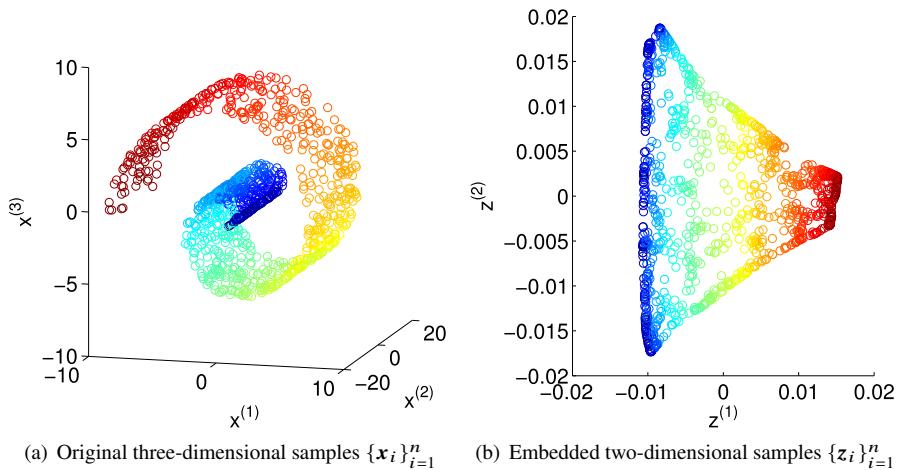
```

n=1000; k=10; a=3*pi*rand(n,1);
x=[a.*cos(a) 30*rand(n,1) a.*sin(a)];
x=x-repmat(mean(x),[n,1]); x2=sum(x.^2,2);
d=repmat(x2,1,n)+repmat(x2',n,1)-2*x*x'; [p,i]=sort(d);
W=sparse(d<=ones(n,1)*p(k+1,:)); W=(W+W'~/2);
D=diag(sum(W,2)); L=D-W; [z,v]=eigs(L,D,3,'sm');
figure(1); clf; hold on; view([15 10]);
scatter3(x(:,1),x(:,2),x(:,3),40,a,'o');
figure(2); clf; hold on; scatter(z(:,2),z(:,1),40,a,'o');

```

**FIGURE 36.5**

MATLAB code of Laplacian eigenmap for 10-nearest neighbor similarity.

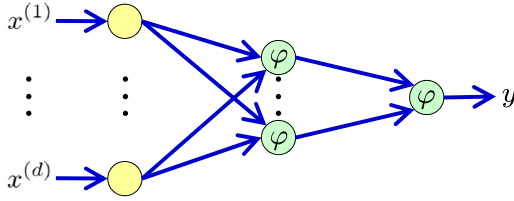
**FIGURE 36.6**

Example of Laplacian eigenmap for 10-nearest neighbor similarity.

## 36.2 SUPERVISED DIMENSIONALITY REDUCTION WITH NEURAL NETWORKS

A *neural network* is a nonlinear model having hierarchical structure. If the number of layers is three (i.e. input, hidden, and output layers), its function is expressed as follows (see Section 21.3):

$$f_{\theta}(x) = \sum_{j=1}^b \alpha_j \phi(x; \beta_j),$$

**FIGURE 36.7**

Dimensionality reduction by neural network. The number of hidden nodes is smaller than the number of input (and output) nodes.

where  $\phi(\mathbf{x}; \boldsymbol{\beta})$  is a basis function parameterized by  $\boldsymbol{\beta}$ . As explained in Section 22.5, LS training of neural networks based on input-output paired training samples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  can be carried out by the gradient method called the *error back-propagation* algorithm.

If the number of hidden units,  $b$ , is set to be less than the dimensionality of input  $\mathbf{x}$ , such a neural network can be used for supervised dimensionality reduction. More specifically, the output of the hidden layer,

$$\mathbf{z} = (\phi(\mathbf{x}; \boldsymbol{\beta}_1), \dots, \phi(\mathbf{x}; \boldsymbol{\beta}_b))^T,$$

can be regarded as a dimensionality reduction solution (Fig. 36.7).

### 36.3 UNSUPERVISED DIMENSIONALITY REDUCTION WITH AUTOENCODER

In this section, the back-propagation approach is applied to unsupervised dimensionality reduction.

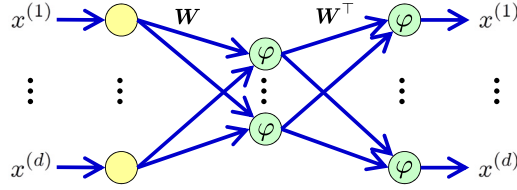
#### 36.3.1 AUTOENCODER

The *autoencoder* [14, 116] is a three-layer neural network given by

$$\mathbf{x} = \boldsymbol{\varphi}(\mathbf{W}^\top \boldsymbol{\varphi}(\mathbf{W}\mathbf{x} + \mathbf{c}) + \mathbf{b}),$$

where  $\mathbf{W}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  are parameters. As illustrated in Fig. 36.8, the autoencoder has a bottleneck layer, and the connection weights between the first and second layers and the connection weights between the second and third layers are shared, i.e.  $\mathbf{W}$  and  $\mathbf{W}^\top$  are used. The connection weights are learned so that output is as close to input as possible, by which a compressed expression of input can be obtained in the bottleneck layer. Let us use the *sigmoidal activation function* (see Fig. 21.6) for the



**FIGURE 36.8**

Autoencoder. Input and output are the same and the number of hidden nodes is smaller than the number of input nodes.

second and third layers:

$$\varphi(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x}^\top \mathbf{w} - a)}.$$

### 36.3.2 TRAINING BY GRADIENT DESCENT

For  $d$ -dimensional input  $\mathbf{x} = (x^{(1)}, \dots, x^{(d)})^\top$  and  $d$ -dimensional output

$$\mathbf{y} = (y^{(1)}, \dots, y^{(d)})^\top = \boldsymbol{\varphi}(\mathbf{W}^\top \boldsymbol{\varphi}(\mathbf{W}\mathbf{x} + \mathbf{c}) + \mathbf{b}),$$

the parameters  $\mathbf{W}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  are learned by the *stochastic gradient* algorithm introduced in Section 15.3:

$$\begin{aligned} \mathbf{W} &\leftarrow \mathbf{W} - \varepsilon \nabla_{\mathbf{W}} J(\mathbf{W}, \mathbf{b}, \mathbf{c}), \\ \mathbf{b} &\leftarrow \mathbf{b} - \varepsilon \nabla_{\mathbf{b}} J(\mathbf{W}, \mathbf{b}, \mathbf{c}), \\ \mathbf{c} &\leftarrow \mathbf{c} - \varepsilon \nabla_{\mathbf{c}} J(\mathbf{W}, \mathbf{b}, \mathbf{c}), \end{aligned}$$

where the *squared loss*,

$$J(\mathbf{W}, \mathbf{b}, \mathbf{c}) = \frac{1}{2} \sum_{k=1}^d (x^{(k)} - y^{(k)})^2, \quad (36.3)$$

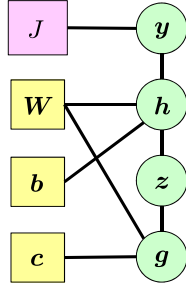
may be used for regression (this corresponds to the log-likelihood of Gaussian distributions), and the *cross entropy loss*,

$$J(\mathbf{W}, \mathbf{b}, \mathbf{c}) = - \sum_{k=1}^d (x^{(k)} \log y^{(k)} + (1 - x^{(k)}) \log(1 - y^{(k)})), \quad (36.4)$$

may be used for binary classification (this corresponds to the log-likelihood of binomial distributions).

For

$$\mathbf{y} = \boldsymbol{\varphi}(\mathbf{h}), \quad \mathbf{h} = \mathbf{W}^\top \mathbf{z} + \mathbf{b}, \quad \mathbf{z} = \boldsymbol{\varphi}(\mathbf{g}), \quad \text{and} \quad \mathbf{g} = \mathbf{W}\mathbf{x} + \mathbf{c}, \quad (36.5)$$

**FIGURE 36.9**

Chain rule for autoencoder.

the *chain rule* of the derivative (Fig. 36.9) yields

$$\begin{aligned}\frac{\partial J}{\partial W_{j,k}} &= \frac{\partial J}{\partial y^{(k)}} \frac{\partial y^{(k)}}{\partial h_k} \frac{\partial h_k}{\partial W_{j,k}} + \frac{\partial J}{\partial y^{(k)}} \frac{\partial y^{(k)}}{\partial h_k} \frac{\partial h_k}{\partial z^{(j)}} \frac{\partial z^{(j)}}{\partial g_j} \frac{\partial g_j}{\partial W_{j,k}}, \\ \frac{\partial J}{\partial b_k} &= \frac{\partial J}{\partial y^{(k)}} \frac{\partial y^{(k)}}{\partial h_k} \frac{\partial h_k}{\partial b_k}, \\ \frac{\partial J}{\partial c_j} &= \sum_{k=1}^d \frac{\partial J}{\partial y^{(k)}} \frac{\partial y^{(k)}}{\partial h_k} \frac{\partial h_k}{\partial z^{(j)}} \frac{\partial z^{(j)}}{\partial g_j} \frac{\partial g_j}{\partial c_j}.\end{aligned}$$

Then the gradients are given by

$$\begin{aligned}\nabla_W J &= z(\nabla_b J)^\top + (\nabla_c J)\mathbf{x}^\top, \\ \nabla_b J &= \begin{cases} (y - \mathbf{x}) * \mathbf{y} * (\mathbf{1} - \mathbf{y}) & \text{(squared loss),} \\ \mathbf{y} - \mathbf{x} & \text{(cross entropy loss),} \end{cases} \\ \nabla_c J &= W \nabla_b J * \mathbf{z} * (\mathbf{1} - \mathbf{z}),\end{aligned}$$

where “\*” denotes the elementwise product.

A MATLAB code for dimensionality reduction by the autoencoder is provided in Fig. 36.10, and its behavior is illustrated in Fig. 36.11. Here, the autoencoder is trained with 100 noiseless images of hand-written digit “2” (see Section 12.5 for the details of the hand-written digit data set). Then a test image illustrated in Fig. 36.11(b), which is not included in the training data set, is input to the trained autoencoder. Since the autoencoder is trained only with noiseless images, noise components may be eliminated in the second layer. Then a denoised image illustrated in Fig. 36.11(c) is obtained as output, which looks much better than the output of linear PCA illustrated in Fig. 36.11(d).

```

load digit.mat; x=X(:,1:100,2); [d,n]=size(x); g=min(x,[],2);
x=(x-repmat(g,1,n))./repmat(max(x,[],2)-g,1,n); x=(x>0.5);
m=10; e=0.01/n; W0=randn(m,d); b0=randn(d,1); c0=rand(m,1);
for o=1:100000
    Z=1./(1+exp(-W0*x-repmat(c0,[1 n])));
    Y=1./(1+exp(-W0'*Z-repmat(b0,[1 n])));
    nb=Y-x; % Cross entropy loss
    % nb=(Y-x).*Y.*(1-Y); % Squared loss
    nc=(W0*nb).*Z.*(1-Z); W=W0-e*(Z*nb'+nc*x');
    b=b0-e*sum(nb,2); c=c0-e*sum(nc,2);
    if norm(W-W0)+norm(b-b0)+norm(c-c0)<0.003, break, end
    W0=W; b0=b; c0=c;
end
t=T(:,1,2)>0.5; u=t; u(rand(d,1)>0.9)=1; u(rand(d,1)>0.9)=0;
z=1./(1+exp(-W*t-c)); y=1./(1+exp(-W'*z-b));

figure(1); clf; hold on; colormap gray
subplot(1,3,1); imagesc(reshape(t,[16 16]))'
subplot(1,3,2); imagesc(reshape(u,[16 16]))'
subplot(1,3,3); imagesc(reshape(y>0.5,[16 16]))'

```

**FIGURE 36.10**

MATLAB code for denoising autoencoder. See Section 12.5 for details of hand-written digit data set “digit.mat.”

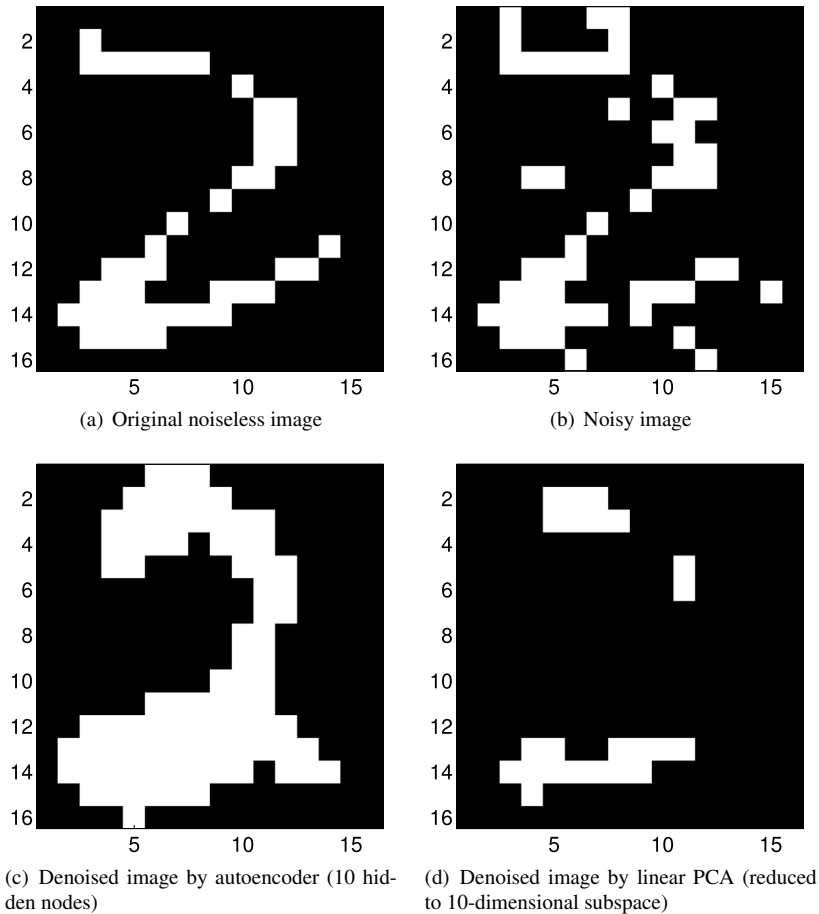
### 36.3.3 SPARSE AUTOENCODER

Let us consider a nonbottleneck neural network, i.e. the number of hidden nodes,  $m$ , is larger than the number of input (and output) nodes,  $d$ . If this architecture is used as an autoencoder, just the identity mapping may be learned as a trivial solution. However, if the activations in the hidden layer,  $\mathbf{z} = (z^{(1)}, \dots, z^{(m)})$  (see Eq. (36.5)), are enforced to be *sparse*, meaningful features may be extracted.

More specifically, in addition to the loss function  $J(\mathbf{W}, \mathbf{b}, \mathbf{c})$  (see Eqs (36.3) and (36.4)), a sparsity-inducing regularization term is minimized at the same time:

$$\min_{\mathbf{W}, \mathbf{b}, \mathbf{c}} J(\mathbf{W}, \mathbf{b}, \mathbf{c}) + \lambda \sum_{j=1}^m R(z^{(j)}),$$

where  $\lambda > 0$  is the regularization parameter and  $R(z)$  is the regularization functional, e.g. the  $\ell_1$ -norm  $R(z) = |z|$  (see Chapter 24), or the KL divergence (see Section 14.2)

**FIGURE 36.11**

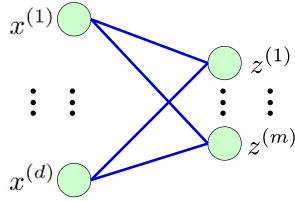
Example of denoising autoencoder.

from a constant  $\rho \in (0, 1)$ :

$$R(z) = \rho \log \frac{\rho}{z/m} + (1 - \rho) \log \frac{(1 - \rho)}{1 - z/m}.$$

## 36.4 UNSUPERVISED DIMENSIONALITY REDUCTION WITH RESTRICTED BOLTZMANN MACHINE

Although the autoencoder is an unsupervised neural network, it is trained in a supervised way by regarding output as input. In this section, another unsupervised

**FIGURE 36.12**

Restricted Boltzmann machine.

neural network called the *restricted Boltzmann machine* [55, 96] is introduced, which is trained in an unsupervised way.

### 36.4.1 MODEL

The restricted Boltzmann machine is a two-layered neural network illustrated in Fig. 36.12, where the left layer for  $d$ -dimensional input  $\mathbf{x} = (x^{(1)}, \dots, x^{(d)})^\top$  is called the *visible layer* and the right layer for  $m$ -dimensional hidden variable

$$\mathbf{z} = (z^{(1)}, \dots, z^{(m)})^\top \in \{0, 1\}^m$$

is called the *hidden layer*. Fig. 36.12 shows that any unit in the input layer and any unit in the hidden layer are completely connected. On the other hand, units in the input layer are not connected to each other and units in the hidden layer are not connected to each other. For the moment, let us assume that input is binary:  $\mathbf{x} \in \{0, 1\}^d$ .

The restricted Boltzmann machine is a model of joint probability  $q(\mathbf{x}, \mathbf{z})$  given by

$$q(\mathbf{x}, \mathbf{z}) = \frac{e^{-E(\mathbf{x}, \mathbf{z})}}{\sum_{\mathbf{x}' \in \{0, 1\}^d} \sum_{\mathbf{z}' \in \{0, 1\}^m} e^{-E(\mathbf{x}', \mathbf{z}')}}.$$

Here,  $E(\mathbf{x}, \mathbf{z})$  is called the *energy function* defined as

$$E(\mathbf{x}, \mathbf{z}) = -\mathbf{b}^\top \mathbf{x} - \mathbf{c}^\top \mathbf{z} - \mathbf{z}^\top \mathbf{W} \mathbf{x},$$

where  $\mathbf{W}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  are parameters. Because there is no connection between the input and hidden units, the conditional probabilities can be factorized as

$$q(\mathbf{z}|\mathbf{x}) = \prod_{j=1}^m q(z^{(j)}|\mathbf{x}) \quad \text{and} \quad q(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^d q(x^{(k)}|\mathbf{z}), \quad (36.6)$$

where each conditional probability is given by the *sigmoidal function*:

$$q(z^{(j)} = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\sum_{k=1}^d W_{j,k} x^{(k)} - c_j)},$$

$$q(x^{(k)} = 1|\mathbf{z}) = \frac{1}{1 + \exp(-\sum_{j=1}^m W_{j,k} z^{(j)} - b_k)}.$$

### 36.4.2 TRAINING BY GRADIENT ASCENT

For training the restricted Boltzmann machine, let us employ MLE (see [Chapter 12](#)). However, since hidden variable  $\mathbf{z}$  is not observable, MLE cannot be performed directly. Here, let us consider the *marginal* model,

$$q(\mathbf{x}) = \sum_{\mathbf{z} \in \{0,1\}^m} q(\mathbf{x}, \mathbf{z}),$$

and perform MLE for this model:

$$\max_{\mathbf{W}, \mathbf{b}, \mathbf{c}} L(\mathbf{W}, \mathbf{b}, \mathbf{c}),$$

where

$$L(\mathbf{W}, \mathbf{b}, \mathbf{c}) = \frac{1}{n} \sum_{i=1}^n \log \left( \sum_{\mathbf{z} \in \{0,1\}^m} q(\mathbf{x}_i, \mathbf{z}) \right).$$

The gradient of the log-likelihood is given by

$$\frac{\partial L}{\partial \theta} = -\frac{1}{n} \sum_{i=1}^n \sum_{\mathbf{z} \in \{0,1\}^m} \frac{\partial E}{\partial \theta} q(\mathbf{z}|\mathbf{x} = \mathbf{x}_i) + \sum_{\mathbf{x} \in \{0,1\}^d} \sum_{\mathbf{z} \in \{0,1\}^m} \frac{\partial E}{\partial \theta} q(\mathbf{x}, \mathbf{z}), \quad (36.7)$$

where  $\theta$  represents either  $\mathbf{W}$ ,  $\mathbf{b}$ , or  $\mathbf{c}$ , and  $\frac{\partial E}{\partial \theta}$  is given by

$$\frac{\partial E}{\partial \mathbf{W}} = -\mathbf{z}^\top \mathbf{x}, \quad \frac{\partial E}{\partial \mathbf{b}} = -\mathbf{x}, \quad \text{and} \quad \frac{\partial E}{\partial \mathbf{c}} = -\mathbf{z}.$$

Thanks to factorization (36.6), the first term in Eq. (36.7) can be computed naively as

$$-\frac{1}{n} \sum_{i=1}^n \sum_{\mathbf{z} \in \{0,1\}^m} \frac{\partial E}{\partial \theta} q(\mathbf{z}|\mathbf{x} = \mathbf{x}_i) = -\frac{1}{n} \sum_{i=1}^n \frac{\partial E}{\partial \theta} \prod_{j=1}^m \sum_{z^{(j)} \in \{0,1\}} q(z^{(j)}|\mathbf{x} = \mathbf{x}_i).$$

On the other hand, the second term in Eq. (36.7) can be approximated by the empirical distribution of samples  $\{\mathbf{x}_i\}_{i=1}^n$  as

$$\begin{aligned} \sum_{\mathbf{x} \in \{0,1\}^d} \sum_{\mathbf{z} \in \{0,1\}^m} \frac{\partial E}{\partial \theta} q(\mathbf{x}, \mathbf{z}) &= \sum_{\mathbf{x} \in \{0,1\}^d} \sum_{\mathbf{z} \in \{0,1\}^m} \frac{\partial E}{\partial \theta} q(\mathbf{z}|\mathbf{x}) p(\mathbf{x}) \\ &\approx \frac{1}{n} \sum_{i=1}^n \sum_{\mathbf{z} \in \{0,1\}^m} \frac{\partial E}{\partial \theta} q(\mathbf{z}|\mathbf{x}_i). \end{aligned}$$

1. Let  $\hat{\mathbf{x}}_i = \mathbf{x}_i$  for  $i = 1, \dots, n$ .
2. Generate  $\{\hat{\mathbf{z}}_i\}_{i=1}^n$  from  $\{\hat{\mathbf{x}}_i\}_{i=1}^n$  following  $q(\mathbf{z}|\mathbf{x} = \hat{\mathbf{x}}_i)$ .
3. Generate  $\{\hat{\mathbf{x}}_i\}_{i=1}^n$  from  $\{\hat{\mathbf{z}}_i\}_{i=1}^n$  following  $q(\mathbf{x}|\mathbf{z} = \hat{\mathbf{z}}_i)$ .
4. Iterate 2–3 until convergence.

**FIGURE 36.13**

Contrastive divergence algorithm for restricted Boltzmann machine. Note that  $q(\mathbf{z}|\mathbf{x} = \hat{\mathbf{x}}_i)$  and  $q(\mathbf{x}|\mathbf{z} = \hat{\mathbf{z}}_i)$  can be factorized as Eq. (36.6), which allows efficient computation.

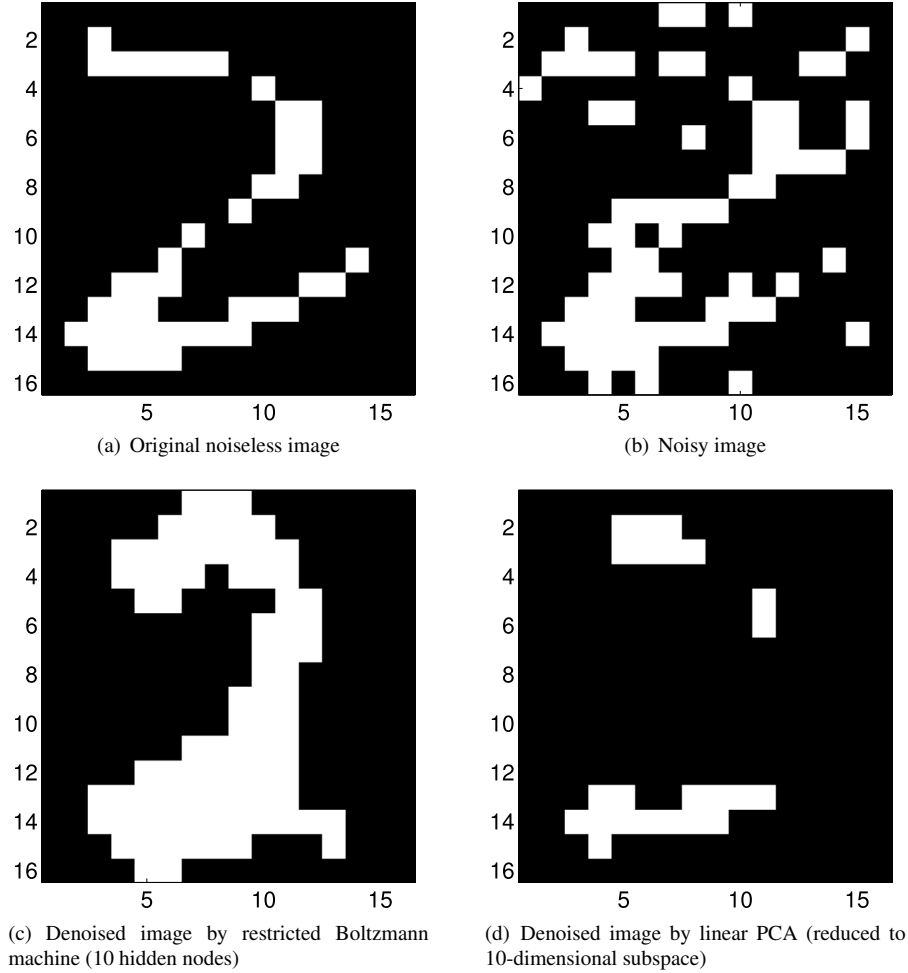
```
load digit.mat; x=X(:,1:100,2); [d,n]=size(x); g=min(x,[],2);
x=(x-repmat(g,1,n))./repmat(max(x,[],2)-g,1,n); x=(x>0.5);
m=10; e=0.01/n; W0=randn(m,d); b0=randn(d,1); c0=rand(m,1);
for o=1:100000
    pZ=1./(1+exp(-W0*x-repmat(c0,1,n))); Z=1*(rand(m,n)<pZ);
    pY=1./(1+exp(-W0'*Z-repmat(b0,1,n))); Y=1*(rand(d,n)<pY);
    pX=1./(1+exp(-W0*Y-repmat(c0,1,n))); W=W0+e*(pZ*x'-pX*Y');
    b=b0+e*(x*prod(pZ)'-Y*prod(pX)');
    c=c0+e*(sum(pZ,2)-sum(pX,2));
    if norm(W-W0)+norm(b-b0)+norm(c-c0)<0.007, break, end
    W0=W; b0=b; c0=c;
end
t=T(:,1,2)>0.5; u=t; u(rand(d,1)>0.9)=1; u(rand(d,1)>0.9)=0;
z=1./(1+exp(-W*t-c)); y=1./(1+exp(-W'*z-b));

figure(1); clf; hold on; colormap gray
subplot(1,3,1); imagesc(reshape(t,[16 16]))'
subplot(1,3,2); imagesc(reshape(u,[16 16]))'
subplot(1,3,3); imagesc(reshape(y>0.5,[16 16]))'
```

**FIGURE 36.14**

MATLAB code for denoising restricted Boltzmann machine. See Section 12.5 for details of handwritten digit data set “digit.mat.”

However, since this is exactly the same as the negative of the first term, the gradient approximately computed in this way always vanishes and thus is useless in practice. To cope with this problem, let us consider another set of samples  $\{\hat{\mathbf{x}}_i\}_{i=1}^n$  generated by the *contrastive divergence algorithm* [54], described in Fig. 36.13, which is based on *Gibbs sampling* (see Section 19.3.3).

**FIGURE 36.15**

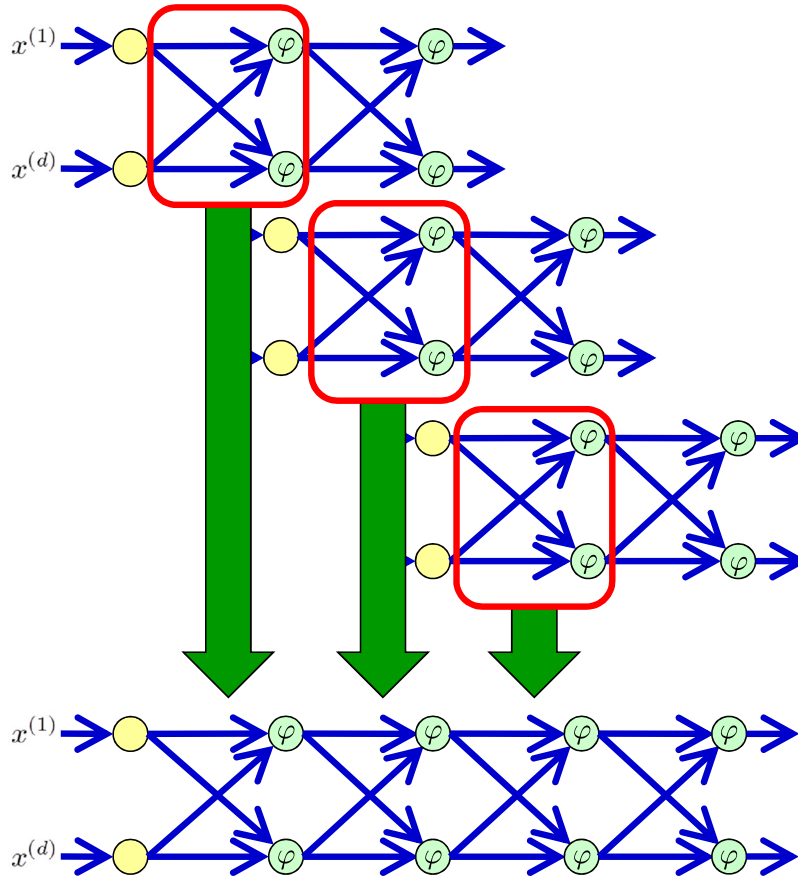
Example of denoising restricted Boltzmann machine.

Finally, the gradient of the log-likelihood is approximated as

$$\frac{\partial L}{\partial \theta} \approx \frac{1}{n} \sum_{i=1}^n \frac{\partial E}{\partial \theta} \prod_{j=1}^m \sum_{z^{(j)} \in \{0,1\}} \left( -q(z^{(j)} | \mathbf{x} = \mathbf{x}_i) + q(z^{(j)} | \mathbf{x} = \hat{\mathbf{x}}_i) \right).$$

Based on this approximated gradient, the parameters  $\mathbf{W}$ ,  $\mathbf{b}$ , and  $\mathbf{c}$  can be learned by *gradient ascent* (see Section 15.3).



**FIGURE 36.16**

Construction of deep neural network by stacking.

When input  $x$  is continuous, the energy function may be replaced with

$$E(x, z) = \frac{1}{2} x^T x - b^T x - c^T z - z^T W x,$$

which results in the Gaussian conditional probability:

$$q(x^{(k)} = 1 | z) = \frac{1}{\sqrt{2\pi}} \exp \left( -\frac{1}{2} \left( x^{(k)} - b_k - \sum_{j=1}^m z^{(j)} W_{j,k} \right)^2 \right).$$

A MATLAB code for dimensionality reduction by the restricted Boltzmann machine is provided in [Fig. 36.14](#), and its behavior is illustrated in [Fig. 36.15](#). This is

the same denoising experiment as Fig. 36.11, and the restricted Boltzmann machine also works much better than linear PCA.

## 36.5 DEEP LEARNING

In Section 22.5, the *error back-propagation* algorithm for supervised training of neural networks was introduced. As pointed out there, due to the hierarchical structure of neural networks, gradient-based training is prone to be trapped by a local optimal solution.

It was experimentally demonstrated that training autoencoders or restricted Boltzmann machines layer by layer by unsupervised learning and stacking them (see Fig. 36.16) can produce good initialization of connection weights [55, 86]. In other words, such deep autoencoders and deep Boltzmann machines can be regarded as good feature extractors for succeeding supervised learning tasks.

Extensive research on deep learning is on going and latest information can be found, e.g. from “<http://deeplearning.net/>.”