

3.5 Dimension Reduction

Many machine learning problems involve high dimensional data. In the case of classifying handwritten digits from MNIST (Mixed National Institute of Standards and Technology) database, for example, each digit is an image of size 28 by 28 (Fig. 3.13 shows 54 digits from the database), therefore the data are in 784-dimensional space. On the other hand, images of handwritten digits are known to live in low-dimensional manifolds, therefore low-dimensional representation of high-dimensional data of interest is possible in this case. This section introduces several basic tools that are found useful in materializing such low-dimensional representations. The general machinery of these tools is referred to as *dimension reduction* or *dimensionality reduction*.

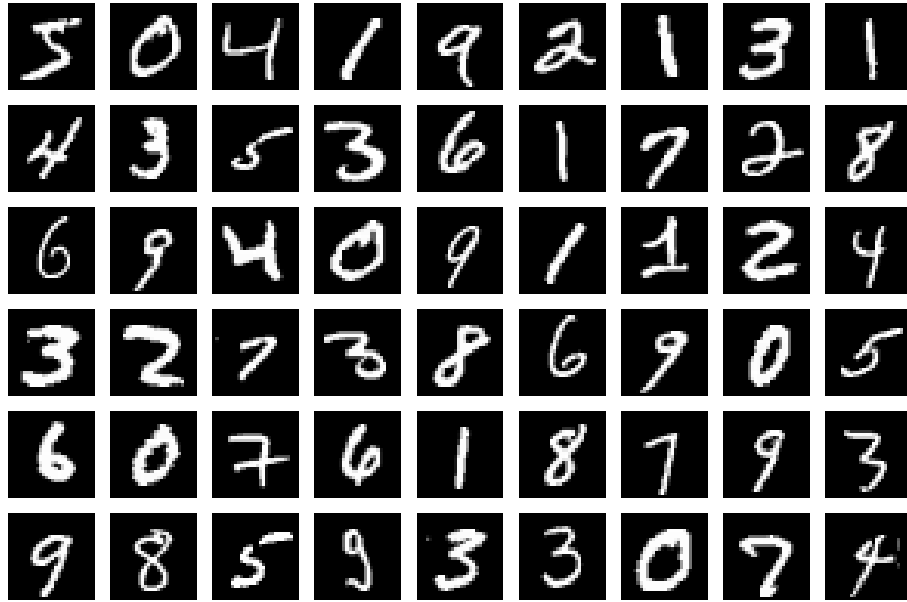


Fig. 3.13 Handwritten digits from MNIST database.

A. Singular Value Decomposition

Let X be a real-valued matrix of size m by n whose rank is r , there exist orthogonal matrices $U \in R^{m \times m}$ and $V \in R^{n \times n}$ such that

$$X = U \Sigma V^T \quad (3.33)$$

where

$$\Sigma = \begin{bmatrix} S & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$

and $S = \text{diag}\{\sigma_1, \sigma_2, \dots, \sigma_r\}$ with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$. The σ_i 's are called *singular values* of X , and the columns of U and V are called left and right singular vectors of X , respectively, and (3.33) is said to be the *singular value decomposition* (SVD) of X . The concepts of singular values and singular vectors are closely related to more popular eigenvalues and eigenvectors. The connection can be seen immediately through the products

$$\mathbf{X}\mathbf{X}^T = \mathbf{U} \begin{bmatrix} \mathbf{S}^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{U}^T, \quad \mathbf{X}^T \mathbf{X} = \mathbf{V} \begin{bmatrix} \mathbf{S}^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{V}^T$$

In words, the nonzero singular values of \mathbf{X} are the square roots of the nonzero eigenvalues of $\mathbf{X}\mathbf{X}^T$ and $\mathbf{X}^T \mathbf{X}$; the left singular vectors are the eigenvectors of $\mathbf{X}\mathbf{X}^T$; and the right singular vectors are the eigenvectors of $\mathbf{X}^T \mathbf{X}$. The MATLAB function `svd` implements (3.33) as `[U,S,V] = svd(X)`.

Following (3.33), the SVD of \mathbf{X} can be written as

$$\mathbf{X} = \mathbf{U}_r \mathbf{S}_r \mathbf{V}_r^T \quad (3.34)$$

where \mathbf{U}_r and \mathbf{V}_r are the matrices composed of the first r columns of \mathbf{U} and \mathbf{V} , respectively. Among other things, SVD is useful in finding optimal low-rank approximations of \mathbf{X} . In effect, given a positive integer $r' < r$ where r is the rank of \mathbf{X} , an optimal rank- r' approximation of \mathbf{X} , denoted by \mathbf{X}' , can be constructed via (3.34) as

$$\mathbf{X}_{r'} = \mathbf{U}_{r'} \mathbf{S}_{r'} \mathbf{V}_{r'}^T \quad (3.35)$$

where $\mathbf{U}_{r'}$ and $\mathbf{V}_{r'}$ are the matrices composed of the first r' columns of \mathbf{U}_r and \mathbf{V}_r , respectively, and $\mathbf{S}_{r'} = \text{diag}\{\sigma_1, \sigma_2, \dots, \sigma_{r'}\}$. The approximation error introduced by $\mathbf{X}_{r'}$ in Frobenius norm is controlled by the 2-norm of the remaining singular values (the residue energy), namely

$$\|\mathbf{X} - \mathbf{U}_{r'} \mathbf{S}_{r'} \mathbf{V}_{r'}^T\|_F = \left(\sum_{i=r'+1}^r \sigma_i^2 \right)^{1/2} \quad (3.36)$$

B. Dimension Reduction of Data Matrix \mathbf{X} Using SVD

One may think of matrix \mathbf{X} in (3.33) as a data array that collects a set of n data points $\{\mathbf{x}_i, i = 1, 2, \dots, n\}$ (as column vectors) in m -dimensional Euclidean space, that is, $\mathbf{X} = [\mathbf{x}_1 \mathbf{x}_2 \dots \mathbf{x}_n]$ whose rank is r . Using (3.34), we define $\mathbf{Z}_r = \mathbf{S} \mathbf{V}_r^T$ which is a matrix of size r by n , and write it as $\mathbf{Z}_r = [\mathbf{z}_1 \mathbf{z}_2 \dots \mathbf{z}_n]$ with $\mathbf{z}_i \in R^{r \times 1}$. It can be readily verified that the set $\{\mathbf{z}_i, i = 1, 2, \dots, n\}$ preserves the topological structure of the original data set $\{\mathbf{x}_i, i = 1, 2, \dots, n\}$ while reducing the dimension from m to r in the sense that the distance and “angle” between each pair of the data points as well as the length (energy) of each individual vector (signal) are preserved, namely,

$$\|\mathbf{z}_i\| = \|\mathbf{x}_i\|, \quad \|\mathbf{z}_i - \mathbf{z}_j\| = \|\mathbf{x}_i - \mathbf{x}_j\| \quad \text{and} \quad \langle \mathbf{z}_i, \mathbf{z}_j \rangle = \langle \mathbf{x}_i, \mathbf{x}_j \rangle \quad \text{for } 1 \leq i, j \leq n \quad (3.37)$$

A simple proof of (3.37) would be based on the fact that the two data sets differ only by \mathbf{U}_r as seen from (3.34), whose columns are orthonormal. This SVD-based dimension reduction is particularly attractive when $r \ll m$.

What if $r \ll m$ does not hold? In this scenario one examines the r nonzero singular values to see if the number of the dominating singular values, denoted by r' , is small. If so, based on (3.35) and (3.36) one can claim that

$$\mathbf{X} \approx \mathbf{U}_{r'} \mathbf{S}_{r'} \mathbf{V}_{r'}^T \quad (3.38)$$

Consequently, if we define $\mathbf{Z}_{r'} = \mathbf{S} \mathbf{V}_{r'}^T$ and write it as $\mathbf{Z}_{r'} = [\mathbf{z}_1 \mathbf{z}_2 \dots \mathbf{z}_n]$ with $\mathbf{z}_i \in R^{r' \times 1}$, then we

have

$$\|z_i\| \approx \|x_i\|, \|z_i - z_j\| \approx \|x_i - x_j\| \text{ and } \langle z_i, z_j \rangle \approx \langle x_i, x_j \rangle \text{ for } 1 \leq i, j \leq n \quad (3.39)$$

demonstrating the validity of using $Z_r = SV_r^T$ as a low-dimensional approximation of the original data X .

Perhaps the biggest problem with the SVD-based dimension reduction is that of computing the SVD of X , especially for large size X . A well-known solution to this problem is the Johnson-Lindenstrauss lemma [16]. Let $P \in R^{p \times m}$ be a projection matrix with $p \ll m$ that is generated by sampling p row vectors independently and uniformly on the unit sphere in the m -dimensional space which are in turn orthonormalized with each row scaled to have its 2-norm equal to $\sqrt{m/p}$. The dimension reduction of data matrix X is accomplished in two steps: First, we compute $Y = PX$ that projects the m -dimensional vectors x_i 's to p -dimensional vectors y_i 's where y_i is the i th column of Y . Next, the SVD-based dimension reduction is applied to data set Y :

$$Y = \hat{U} \hat{\Sigma} \hat{V}^T = \underbrace{\hat{U}_r}_{\substack{\text{first } r \\ \text{columns} \\ \text{of } \hat{U}}} \cdot \hat{S}_r \cdot \underbrace{\hat{V}_r^T}_{\substack{\text{first } r \\ \text{rows} \\ \text{of } \hat{V}^T}}$$

from which we obtain $\hat{Z}_r = \hat{S}_r \cdot \hat{V}_r^T$. It can be shown (Gilbert-Park-Wakin 2012) that if

$$p \geq \frac{r \log(42/\varepsilon) + \log(2/\delta)}{f(\varepsilon/\sqrt{2})}$$

where $\varepsilon \in (0, 1)$, $\delta \in (0, 1)$ and $f(\cdot)$ is quadratic, then with probability $1 - \delta$

$$(1 - \varepsilon)^{1/2} \leq \frac{\hat{\sigma}_j}{\sigma_j} \leq (1 + \varepsilon)^{1/2}$$

for $j = 1, 2, \dots, r$, and

$$\|v_j - \hat{v}_j\|_2 \leq \min \left\{ \sqrt{2}, \frac{\varepsilon \sqrt{1 + \varepsilon}}{\sqrt{1 - \varepsilon}} \max_{i \neq j} \frac{\sqrt{2} \sigma_i \sigma_j}{\min_{c \in [-1, 1]} \{ |\sigma_i^2 - \sigma_j^2 (1 + c\varepsilon)| \}} \right\}$$

C. Principal Component Analysis (PCA)

Again we consider a real-valued data matrix $X = [x_1 \ x_2 \ \dots \ x_n]$ of size m by n . Define the *average vector* and *covariance matrix* of the n component vectors x_i 's as

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad C = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T \quad (3.40)$$

respectively. Note that C is at least positive semidefinite if not positive definite, hence its SVD is identical to its eigen-decomposition

$$\mathbf{C} = \mathbf{U}\mathbf{S}\mathbf{U}^T \quad (3.41)$$

where $\mathbf{U} = [\mathbf{u}_1 \mathbf{u}_2 \cdots \mathbf{u}_m]$ is an orthogonal matrix and $\mathbf{S} = \text{diag}\{\sigma_1, \sigma_2, \dots, \sigma_m\}$ with $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_m \geq 0$, and \mathbf{u}_i is \mathbf{C} 's i th eigenvector associated with i th eigenvalue σ_i . By taking into account only the K largest eigenvalues and the associated eigenvectors, we obtain a useful approximation of variance matrix \mathbf{C} as

$$\mathbf{C} \approx \mathbf{U}_K \mathbf{S}_K \mathbf{U}_K^T \quad \text{with } \mathbf{U}_K = [\mathbf{u}_1 \mathbf{u}_2 \cdots \mathbf{u}_K] \text{ and } \mathbf{S}_K = \text{diag}\{\sigma_1, \sigma_2, \dots, \sigma_K\} \quad (3.42)$$

The usefulness of (3.42) may be understood from two different but related perspectives as follows.

(1) Since \mathbf{C} is the variance for $\{\mathbf{x}_i - \bar{\mathbf{x}}, i = 1, 2, \dots, n\}$, the K eigenvectors $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_K\}$ capture the most significant variations between the individual data points and their average, each point $\mathbf{x}_i - \bar{\mathbf{x}}$ is well represented by its projection onto the K -dimensional subspace spanned by $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_K\}$, i.e.,

$$\mathbf{z}_i = \mathbf{U}_K^T (\mathbf{x}_i - \bar{\mathbf{x}}) \quad (3.43a)$$

In effect, point \mathbf{x}_i is approximated by a point in that subspace (plus the average $\bar{\mathbf{x}}$) as

$$\mathbf{x}_i \approx \mathbf{U}_K \mathbf{z}_i + \bar{\mathbf{x}} \quad (3.43b)$$

In this way, we see how the original data set \mathbf{X} in R^m is reduced to an approximately equivalent data set $\mathbf{Z} = [\mathbf{z}_1 \mathbf{z}_2 \cdots \mathbf{z}_n]$ in space R^K (plus a single average vector $\bar{\mathbf{x}}$). This PCA-based dimensionality reduction becomes significant when $K \ll m$.

(2) Now let us consider a case in supervised multi-class classification with a training set that consists of L data classes of the form $\{(\mathbf{x}_i^{(j)}, y_j), i = 1, 2, \dots, n_j\}$ for $j = 1, 2, \dots, L$. The problem is to classify a point \mathbf{x} outside the training data to one of the L classes. To this end, PCA is applied to each data class $\{(\mathbf{x}_i^{(j)}, y_j), i = 1, 2, \dots, n_j\}$ by first computing the average vector and covariance matrix as

$$\bar{\mathbf{x}}_j = \frac{1}{n_j} \sum_{i=1}^{n_j} \mathbf{x}_i^{(j)}, \quad \mathbf{C}_j = \frac{1}{n_j - 1} \sum_{i=1}^{n_j} (\mathbf{x}_i^{(j)} - \bar{\mathbf{x}}_j)(\mathbf{x}_i^{(j)} - \bar{\mathbf{x}}_j)^T$$

then performing eigen-decomposition of \mathbf{C}_j

$$\mathbf{C}_j = \mathbf{U}_j \mathbf{S}_j \mathbf{U}_j^T$$

so as to approximate it as

$$\mathbf{C}_j \approx \mathbf{U}_K^{(j)} \mathbf{S}_K^{(j)} \mathbf{U}_K^{(j)T}$$

where $\mathbf{U}_K^{(j)} = [\mathbf{u}_1^{(j)} \mathbf{u}_2^{(j)} \cdots \mathbf{u}_K^{(j)}]$ consists of K eigenvectors of \mathbf{C}_j that are associated K largest eigenvalues and $\mathbf{S}_K^{(j)}$ is a diagonal matrix that collects the K largest eigenvalues in descent order. The classification of a testing point \mathbf{x} is carried out as follows: (i) Using each data class, compute projection $\mathbf{z}_j = \mathbf{U}_K^{(j)T}(\mathbf{x} - \bar{\mathbf{x}}_j)$ and approximate point \mathbf{x} in the j th class as $\hat{\mathbf{x}}_j = \mathbf{U}_K^{(j)} \mathbf{z}_j + \bar{\mathbf{x}}_j$; (ii) Evaluate the difference between \mathbf{x} and its approximation $\hat{\mathbf{x}}_j$ in 2-norm, i.e., $e_j = \|\mathbf{x} - \hat{\mathbf{x}}_j\|_2$ for $j = 1, 2, \dots, L$. Point \mathbf{x} is classified to class j^* if e_{j^*} reaches the minimum among $\{e_j, j = 1, 2, \dots, L\}$.

D. An Example: Classification of Handwritten Digits from MNIST Data Base Using PCA

The MNIST database provides 60,000 handwritten digits for training purposes and a separate 10,000 digits testing, where each digit is a black-and-white image of size 28 by 28. In this example the training images were partitioned into ten classes, with the images in the j th class representing digit j , for $j = 0, 1, \dots, 9$. We select from each class 500 images at random, then convert each image into a column vector of dimension 784 to construct a data matrix \mathbf{X}_j of size 784 by 500. For convenience, these ten data matrices are combined to a signal data matrix \mathbf{X}_{500} of size 784 by 5000 as a .mat file which you can download from the course website. The PCA-based classification method as described above was applied to classify the 10K testing digits. With $K = 24$ (see part C for the meaning of K), an error rate of 4.4% was achieved. The MATLAB file that implements the simulation is enclosed below.

```
% To implement the example in Sec. 3.5.D for classifying handwritten digits
% from MNIST database.
% Input:
% X: Ten classes of input data, each is of size 784 by nt (here we use nt = 500).
% n: the actual size of each data matrix is 784 by n. Hence n is limited to n <= nt.
% K: the number of "eigen-digits" to be used in PCA. Typically K << m = 784.
% Te: testing data (here we use Te28 of size 784 by 10000).
% Lte: Labels for the testing data.
% Output:
% Js: the numerical values of testing digits obtained by PCA-based classification.
% er: classification error.
% Written by W.-S. Lu, University of Victoria. Last modified: March 9, 2015.
% Example: load X500; load Te28; load Lte28;
% [Js,er] = class_pca_digits(X500,500,24,Te28,Lte28);
function [Js,er] = class_pca_digits(X,n,K,Te,Lte)
U = [];
Xb = [];
t = size(X,2);
nt = t/10;
for i = 1:10,
    Tw = X(:,((i-1)*nt+1):(i*nt));
    Ti = Tw(:,1:n);
    xbi = mean(Ti)';
    Xh = Ti - xbi*ones(1,n);
    [u,s,v] = svd(Xh*Xh');
```

```

    U = [U u(:,1:K)];
    Xb = [Xb xbj];
end
Lte = Lte(:);
M = length(Lte);
Js = zeros(M,1);
for m = 1:M,
    xw = Te(:,m);
    e = zeros(1,10);
    for j = 1:10;
        Uj = U(:,(j-1)*K+1):(j*K));
        xbj = Xb(:,j);
        zj = Uj'*(xw-xbj);
        xjh = Uj*zj + xbj;
        e(j) = norm(xw-xjh);
    end
    [emin,ind] = min(e);
    Js(m) = ind - 1;
end
E = (Lte ~= Js);
er = sum(E)/M

```

Problems

3.14 Prove (3.37). (Hint: use Eq. (3.34))

3.15 The Frobenius norm of a real-valued matrix $A = \{a_{i,j}\}$ is equal to the square root of the sum of all $a_{i,j}^2$.

(i) Prove that $\|A\|_F^2 = \text{trace}(AA^T)$.

(ii) Use the result from part (i) to show Eq. (3.36).

3.16 The computational complexity of performing the eigen-decomposition of C in (3.41) becomes intractably high when m is very large. In the case of $n \ll m$, the computational burden may be greatly reduced and this exercise problem introduces such a technique.

(i) Define $\hat{X} = \frac{1}{\sqrt{n-1}} [x_1 - \bar{x} \quad x_2 - \bar{x} \quad \cdots \quad x_n - \bar{x}]$. Verify that $C = \hat{X} \hat{X}^T$.

(ii) Show that the nonzero eigenvalues of C are the same as the nonzero eigenvalues of $\Phi = \hat{X}^T \hat{X}$.

(iii) Let v be an eigenvector of Φ , show that $\hat{X} v$ is an eigenvector of C .

(iv) Based on the results from parts (i) to (iii), develop an efficient method to compute U and S in (3.41) without computing the eigenvalues and eigenvectors of C .