

CS376 Machine Learning

Dimensionality Reduction & PCA

김기응

Kee-Eung Kim
KAIST

kekim@kaist.ac.kr

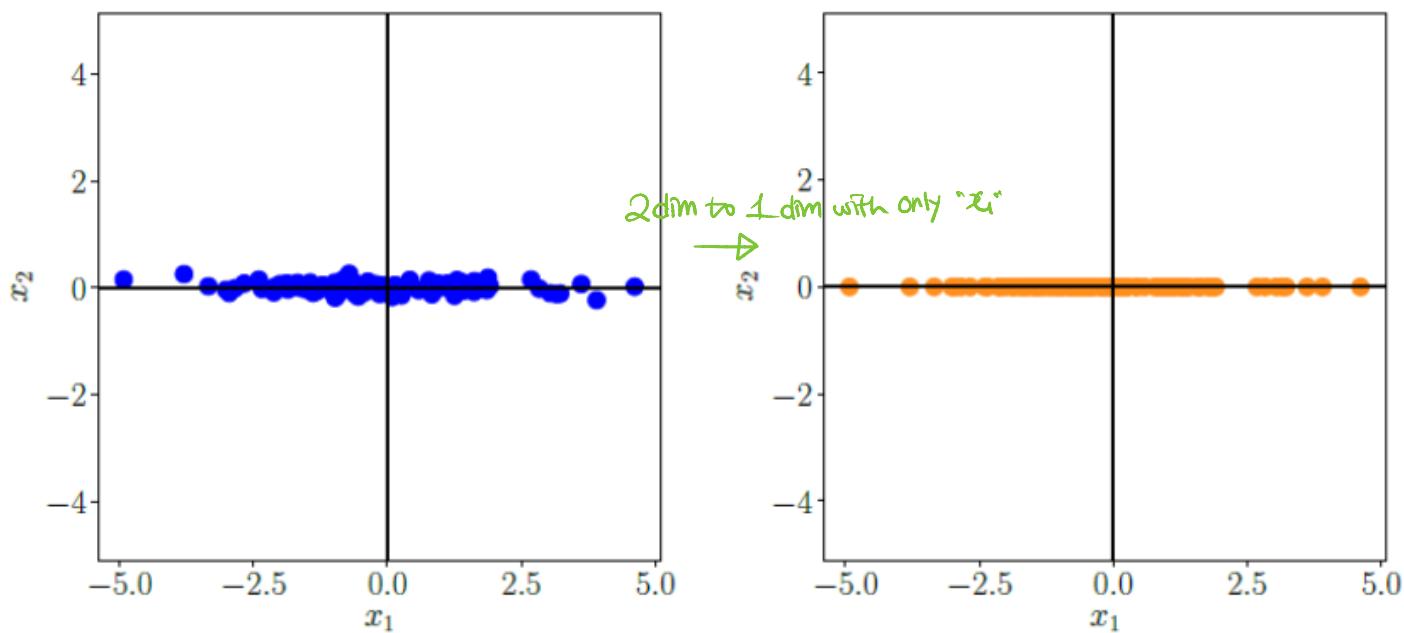


Korea Advanced Institute of Science and Technology
한국과학기술원

Motivation

Figure 10.1

Illustration:
Dimensionality
reduction. (a) The
original dataset
does not vary much
along the x_2
direction. (b) The
data from (a) can be
represented using
the x_1 -coordinate
alone with nearly no
loss.



(a) Dataset with x_1 and x_2 coordinates.

(b) Compressed dataset where only the x_1 coordinate is relevant.

Problem Setting

- Find projections \tilde{x}_n of data points x_n that are as similar as possible, but with significantly lower ^{intrinsic} dimensionality

- $X = \{x_1 \dots x_N\}$, $x_n \in \mathbb{R}^D$ with mean 0. (important)

- Data covariance matrix

$$S = \frac{1}{N} \sum_{n=1}^N x_n x_n^T$$

the columns b_1, \dots, b_M of B . Retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional code (Hotelling, 1933).

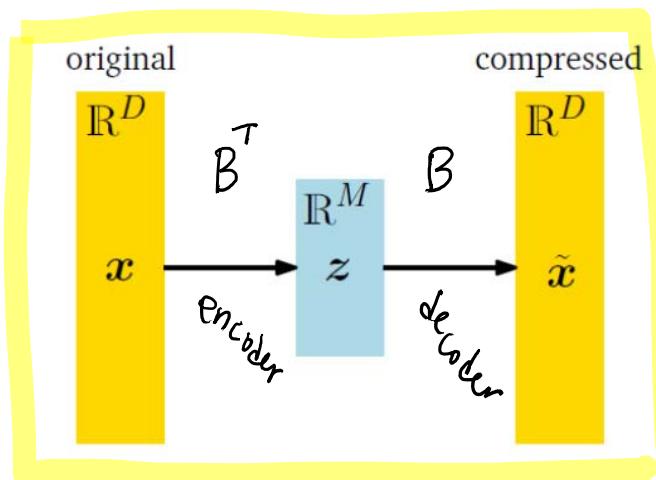
- Find $z_n = B^T x_n \in \mathbb{R}^M$

$$B \triangleq [b_1 \dots b_M] \in \mathbb{R}^{D \times M}$$

$\{b_1 \dots b_M\}$ are orthonormal.

- $B^{-1} = B^T$, so $\tilde{x}_n = B z_n \in \mathbb{R}^D$

Figure 10.2
Graphical illustration of PCA.
In PCA, we find a compressed version \tilde{x} of original data x that has an intrinsic lower-dimensional representation z .



$$x \xrightarrow[\mathbb{E}IR^D]{} z \xrightarrow[\mathbb{E}IR^M]{} \tilde{x} \xrightarrow[\mathbb{E}IR^B]{} \hat{x}$$

10.2 Maximum Variance Perspective.

- Retain most informations after data compression \Leftrightarrow ~~구현하는~~ Capture the largest amount of variance in the low-dimensional code.

10.2.1 Direction with Maximal Variance

* single vector $b_1 \in \mathbb{R}^D$ that maximize the variance of projected data \Rightarrow .

* maximize the variance of first coordinate z_1 of $z \in \mathbb{R}^M$

$$\Rightarrow V_1 := \mathbb{V}[z_1] = \frac{1}{N} \sum_{n=1}^N z_n^2 \quad (z_n \in \mathbb{R}^M, z \in \mathbb{R}^D)$$

z_1 의 variance \rightarrow ~~하는~~ b_1 .

$$z_n = b_1^T x_n \rightarrow z = B^T X$$

data covariance matrix

$$\therefore V_1 = \frac{1}{N} \sum_{n=1}^N (b_1^T x_n)^2 = \frac{1}{N} \sum_{n=1}^N b_1^T x_n x_n^T b_1 = b_1^T \left(\frac{1}{N} \sum_{n=1}^N x_n x_n^T \right) b_1 = b_1^T S b_1 \quad \therefore \text{orth } \|b_1\|^2 = 1$$

\Rightarrow Optimization problem

$$\max_{b_1} b_1^T S b_1 \quad \text{s.t. } \|b_1\|^2 = 1$$

$$\mathcal{L}(b_1, \lambda) = b_1^T S b_1 + \lambda_1 (1 - b_1^T b_1)$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = 2b_1^T S - 2\lambda_1 b_1^T \quad \frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - b_1^T b_1$$

$\therefore Sb_1 = \lambda_1 b_1 \quad \& \quad b_1^T b_1 = 1 \Leftrightarrow \lambda_1: \text{eigenvalue}, b_1: \text{eigenvector}$

$$\Rightarrow V_1 = b_1^T S b_1 = \lambda_1 b_1^T b_1 = \lambda_1 \Leftrightarrow \text{variance } \tilde{x}_1 = \text{eigenvalue } \tilde{x}_1$$

\therefore variance 값이 클려면 eigenvalue 값도 커야해.

이런식으로 $z \in \mathbb{R}^M$ 中有 M-개의 큰 eigenvalue 값

뽑고, 해당 eigenvector \tilde{x}_1 을 구성하기

$$x \xrightarrow[\mathbb{E}R^0]{} z \xrightarrow[\mathbb{E}R^m]{} \tilde{x} \xrightarrow[\mathbb{E}R^0]{} \hat{x}$$

10.2.2 M-dimensional subspace with Maximal Variance

* Symmetric Σ matrix \Leftrightarrow eigenvalue $\in \text{REAL}$

eigenvector \in orthonormal.

* M 번째까지는 큰 eigenvalue 3 \Rightarrow 고, M 번째는 남는 모든 이유는 쓰고 싶어

\rightarrow subtract the effect of the first $m-1$ principal components from the data,
thereby trying to find principal components that compress the
remaining information.

$$\hat{x} := x - \sum_{i=1}^{m-1} b_i b_i^\top x = x - B_{m-1} x$$

$$\Rightarrow V_m = V[Z_m] = \frac{1}{N} \sum_{n=1}^N Z_m^2 = \frac{1}{N} \sum_{n=1}^N (b_m^\top \hat{z}_n)^2 = b_m^\top \hat{\Sigma} b_m.$$

\Rightarrow 똑같이 optimize 하면 " b_m 은 eigenvector of $\hat{\Sigma}$ with largest eigenvalue for $\hat{\Sigma}$ "

그런데 사실 $S, \hat{\Sigma}$ 은 같은 sets of eigenvectors 가지고 있음.

\hookrightarrow D 개의 distinct Σ eigenvector 를 지님.

(1) 만약 Σ 의 $m+1$ 번째 eigenvector b_1, b_2, \dots, b_m $\neq 0$. b_i 는 Σ 의 eigenvector 라면 $Sb_i = \lambda b_i$

$$\hat{\Sigma} b_i = \frac{1}{N} \hat{X} \hat{X}^\top b_i = \frac{1}{N} (x - B_{m-1} x)(x - B_{m-1} x)^\top b_i = (\Sigma - S B_{m-1} - B_{m-1}^\top S + B_{m-1}^\top S B_{m-1}) b_i$$

① $i \geq m$ $B_{m-1} b_i = 0$

② $i < m$. $B_{m-1} b_i = b_i$

$$\hat{\Sigma} b_i = S b_i = \lambda b_i$$

$$\hat{\Sigma} b_i = 0 b_i = 0 \Rightarrow \underline{\text{eigenvalue}} = 0.$$

b_1, b_2, \dots, b_m : Span Σ ULT SPACE of $\hat{\Sigma}$

\therefore every eigenvector of $\Sigma = \hat{\Sigma}$

Σ 의 eigenvector of $(m+1)$ dim principal subspace \Rightarrow eigenvalue of $\hat{\Sigma} = 0$.

* M dimensional subspace on project Σ \Rightarrow variance of data equals sum of the $m+1$ eigenvalues with high values : $V_m = \sum_{i=1}^m \lambda_i$

Maximum Variance Perspective

- First, we aim to maximize variance of the first coord z_1 ,

$$V_1 \triangleq V[z_1] = \frac{1}{N} \sum_{n=1}^N z_{1,n}^2$$

$$\underline{x} \xrightarrow{B^T} \underline{z} \xrightarrow{B} \tilde{\underline{x}}$$

where $z_{1,n} = b_1^T x_n$ (From $\underline{z} = B^T \underline{x}$)

- $V_1 = \frac{1}{N} \sum_{n=1}^N (b_1^T x_n)^2 = \frac{1}{N} \sum_{n=1}^N b_1^T x_n x_n^T b_1 = b_1^T \left(\frac{1}{N} \sum_{n=1}^N x_n x_n^T \right) b_1 = b_1^T S b_1$
- Thus, we need solve constrained optimization problem "data covariance matrix"

$$\max_{b_1} b_1^T S b_1 \text{ subject to } \|b_1\|^2 = 1$$

$$\Leftrightarrow \mathcal{L}(b_1, \lambda) = b_1^T S b_1 + \lambda_1 (1 - b_1^T b_1)$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = 2b_1^T S - 2\lambda_1 b_1^T = 0 \quad \frac{\partial \mathcal{L}}{\partial \lambda} = 1 - b_1^T b_1 = 0$$

$$\Leftrightarrow \begin{cases} S b_1 = \lambda_1 b_1 \\ b_1^T b_1 = 1 \end{cases} \Rightarrow b_1 \text{ is an eigenvector of } S, \quad \lambda_1 \text{ is the corresponding eigenvalue.}$$

$$V_1 = b_1^T S b_1 = \lambda_1 b_1^T b_1 = \lambda_1 \Rightarrow \text{choose eigenvector with}$$

"First principal component" largest eigenvalue! ⇒ largest variance! ✨

Maximum Variance Perspective

- Assume we have found $m-1$ principal component. m-1 개까지는 이전에
찾은 방식으로 eigenvector
를 만들었음.
- That is, we have found b_1, \dots, b_{m-1} , and find remaining PCs that compress remaining information → 마지막 m개는 남아있는 정보 가장 잘
이유를 수 있는 애로 고르기.

$$\hat{X} \triangleq X - \sum_{i=1}^{m-1} b_i b_i^T X \quad \text{where } X = [x_1, \dots, x_N] \in \mathbb{R}^{D \times N} (!)$$

원래 집약하는 애

$$\max V_m = V[\hat{z}_m] = \frac{1}{N} \sum_{n=1}^N \hat{z}_m^2 = \frac{1}{N} \sum_{n=1}^N (b_m^T \hat{X}_n)^2 = b_m^T \hat{S} b_m$$

Subject to $\|b_m\|=1$ and $b_i^T b_m = 0, i=1 \dots m-1$
 $\Rightarrow b_m$ is the eigenvector with largest eigenvalue of \hat{S}

- This is also eigenvector of S

$$\begin{aligned} \hat{S} &= \frac{1}{N} \sum_{n=1}^N \hat{x}_n \hat{x}_n^T = \frac{1}{N} \sum_{n=1}^N \left(x_n - \sum_{i=1}^{m-1} b_i b_i^T x_n \right) \left(x_n - \sum_{i=1}^{m-1} b_i b_i^T x_n \right)^T \\ &= \frac{1}{N} \sum_{n=1}^N \left[x_n x_n^T + \left(\sum_{i=1}^{m-1} b_i b_i^T \right) x_n x_n^T \left(\sum_{i=1}^{m-1} b_i b_i^T \right) \right. \\ &\quad \left. - x_n x_n^T \left(\sum_{i=1}^{m-1} b_i b_i^T \right) - \left(\sum_{i=1}^{m-1} b_i b_i^T \right) x_n x_n^T \right] \end{aligned}$$

Maximum Variance Perspective

This is also eigenvector of S

$$\hat{S} = \frac{1}{N} \sum_{n=1}^N \hat{x}_n \hat{x}_n^T = \frac{1}{N} \sum_{n=1}^N \left(x_n - \sum_{i=1}^{m-1} b_i b_i^T x_n \right) \left(x_n - \sum_{i=1}^{m-1} b_i b_i^T x_n \right)^T$$

$$\begin{aligned} &= \frac{1}{N} \sum_{n=1}^N \left[x_n x_n^T + \left(\sum_{i=1}^{m-1} b_i b_i^T \right) x_n x_n^T \left(\sum_{i=1}^{m-1} b_i b_i^T \right) \right. \\ &\quad \left. - x_n x_n^T \left(\sum_{i=1}^{m-1} b_i b_i^T \right) - \left(\sum_{i=1}^{m-1} b_i b_i^T \right) x_n x_n^T \right] \end{aligned}$$

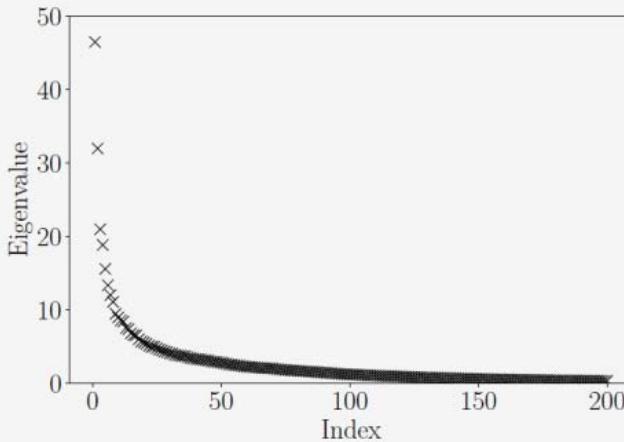
$$\hat{S} b_m = \frac{1}{N} \sum_{n=1}^N \left[x_n x_n^T b_m + 0 \right] - 0 = S b_m$$

last 0: $\frac{1}{N} \sum_{n=1}^N \sum_{i=1}^{m-1} b_i b_i^T x_n x_n^T b_m = \sum_{i=1}^{m-1} b_i b_i^T \underbrace{\left(\frac{1}{N} \sum_{n=1}^N x_n x_n^T \right)}_S b_m$

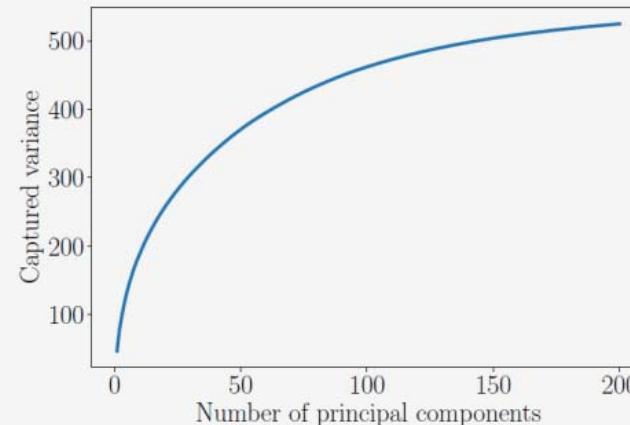
$$= \sum_{i=1}^{m-1} b_i b_i^T b_m = 0$$

$$\therefore V_m = b_m^T \hat{S} b_m = b_m^T S b_m$$

Maximum Variance Perspective



(a) Eigenvalues (sorted in descending order) of the data covariance matrix of all digits '8' in the MNIST training set.



(b) Variance captured by the principal components.

- The amount of variance captured by PCA using M PCs

$$V_M = \sum_{m=1}^M \lambda_m$$

The amount of variance lost by PCA

$$\Sigma_M = \sum_{j=M+1}^D \lambda_j = V_D - V_M$$

relative measure: $\frac{V_M}{V_D}$, $1 - \frac{V_M}{V_D}$

Projection Perspective

- PCA as an algorithm that directly minimizes approximation error ($\|x_n - \tilde{x}_n\|$)
 - Assume ordered orthonormal basis $B = (b_1, \dots, b_D)$,
ie $b_i^T b_j = 1$ iff $i=j$ and 0 otherwise
 - Any $x \in \mathbb{R}^D$, $x = \sum_{d=1}^D \zeta_d b_d = \sum_{m=1}^M \underline{\zeta_m} b_m + \sum_{j=M+1}^D \underline{\zeta_j} b_j$
 $\tilde{x} = \sum_{m=1}^M \underline{\zeta_m} b_m \in U \subseteq \mathbb{R}^D$ with $\dim(U) = M$
- Rem $\underline{\zeta_m}$ and $\underline{\zeta_m}$ are not identical? identical?

Projection Perspective

a $\min_{\substack{z_1 \dots z_n \\ b_1 \dots b_m}} J_M \triangleq \frac{1}{N} \| x_n - \tilde{x}_n \|_2^2$

$B^T Z$

①

Finding
 $z_1 \dots z_n$

$$\frac{\partial J_M}{\partial z_{in}} = \frac{\partial J_M}{\partial \tilde{x}_n} \frac{\partial \tilde{x}_n}{\partial z_{in}}, \quad \left\{ \begin{array}{l} \frac{\partial J_M}{\partial \tilde{x}_n} = -\frac{2}{N} (x_n - \tilde{x}_n)^T \in \mathbb{R}^{1 \times D} \\ \frac{\partial \tilde{x}_n}{\partial z_{in}} = \frac{\partial}{\partial z_{in}} \left(\sum_{m=1}^M z_{mn} b_m \right) = b_i \end{array} \right.$$

$$= -\frac{2}{N} (x_n - \tilde{x}_n)^T b_i = -\frac{2}{N} (x_n - \sum_{m=1}^M z_{mn} b_m)^T b_i$$

단위 길이 orthogonality 때문

$$= -\frac{2}{N} (x_n^T b_i - z_{in} b_i^T b_i) = -\frac{2}{N} (x_n^T b_i - z_{in}) = 0$$

$$\therefore z_{in} = x_n^T b_i = b_i^T x_n \quad (z_{in} \text{ is the coordinate of ortho. projection of } x_n \text{ onto } b_i)$$

Note ① z_m and z_{in} are identical for $m=1 \dots M$

Low-Rank Approximation Perspective

② $z_n = \beta^T x_n$, $\tilde{z}_n = \underbrace{\beta}_{D \times M} \underbrace{x_n}_{M}$ $B = [b_1 \cdots b_M] \in R^{D \times M}$

$D \rightarrow M \rightarrow D$

Finding
 $b_1 \cdots b_M$

$$\tilde{x}_n = \sum_{m=1}^M z_m b_m = \sum_{m=1}^M (x_n^T b_m) b_m = \left(\sum_{m=1}^M b_m b_m^T \right) x_n$$

$$x_n = \sum_{d=1}^D z_d b_d = \sum_{d=1}^D (x_n^T b_d) b_d = \left(\sum_{d=1}^D b_d b_d^T \right) x_n$$

$$= \left(\sum_{m=1}^M b_m b_m^T \right) x_n + \left(\sum_{j=M+1}^D b_j b_j^T \right) x_n$$

$$x - \tilde{x}_n = \left(\sum_{j=M+1}^D b_j b_j^T \right) x_n = \sum_{j=M+1}^D (x_n^T b_j)^T b_j$$

low rank \leftarrow Approx. View

$$\tilde{x}_n = \left(\sum_{m=1}^M b_m b_m^T \right) x_n = BB^T x_n \quad (BB^T \text{ is of rank } \leq M)$$

$$J_M = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \|x_n - BB^T x_n\|^2 = \frac{1}{N} \sum_{n=1}^N \|(I - BB^T)x_n\|^2$$

(Try to find the best rank- M approximation of BB^T to I)

Low-Rank Approximation Perspective

$$\begin{aligned} J_M &= \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \left\| \sum_{j=M+1}^D (b_j^T x_n) b_j \right\|^2 = \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D (b_j^T x_n)^2 \\ &= \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D b_j^T x_n x_n^T b_j = \sum_{j=M+1}^D b_j^T \left(\frac{1}{N} \sum_{n=1}^N x_n x_n^T \right) b_j = \sum_{j=M+1}^D b_j^T S b_j \\ &= \sum_{j=M+1}^D \text{tr}(b_j^T S b_j) = \sum_{j=M+1}^D \text{tr}(b_j b_j^T S) = \text{tr}\left(\left(\sum_{j=M+1}^D b_j b_j^T\right) S\right) \\ &= \sum_{j=M+1}^D \lambda_j \quad (\text{minimize the variance of data when projected on to the subspace we ignore}) \\ &= (\text{Select } (D-M) \text{ eigenvectors with smallest eig values}) \\ &= (\text{Select } D \text{ eigenvectors with largest eig values}) \end{aligned}$$

PCA Illustration

Example 10.3 (MNIST Digits Embedding)

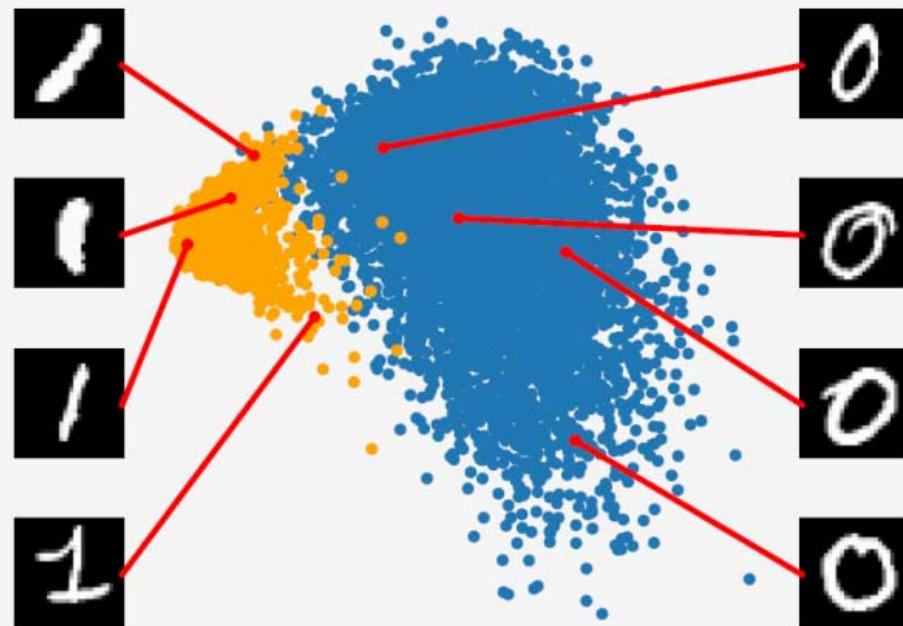


Figure 10.10
Embedding of
MNIST digits 0
(blue) and 1
(orange) in a
two-dimensional
principal subspace
using PCA. Four
embeddings of the
digits '0' and '1'
in the principal
subspace are
highlighted in red
with their
corresponding
original digit.