

Inclass 20: Principal Component Analysis

[SCS4049] Machine Learning and Data Science

Seongsik Park (s.park@dgu.edu)

Department of Artificial Intelligence, Dongguk University

Dimensional reduction

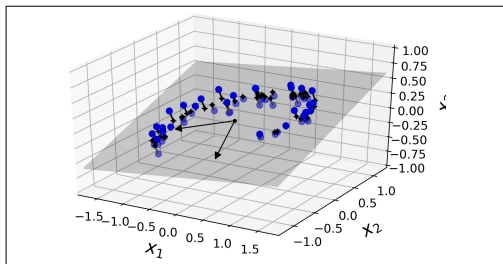


Figure 8-2. A 3D dataset lying close to a 2D subspace

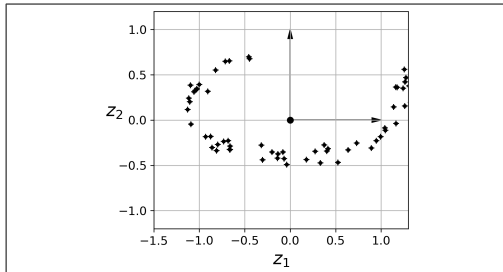


Figure 8-3. The new 2D dataset after projection

Covariance matrix

Covariance measures the strength of the linear relationship between two variables

$$\sigma_{xy} = E[(x - \mu_x)(y - \mu_y)] \quad (1)$$

Covariance matrix C for multivariate random variable X

$$C_{ij} = E[(x_i - \mu_i)(x_j - \mu_j)] \quad (2)$$

Principal component analysis (PCA)

Preserving the variance

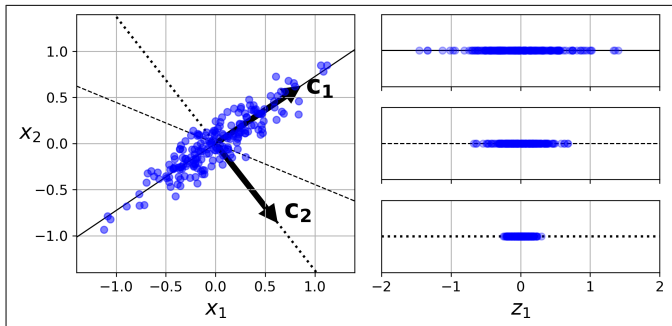


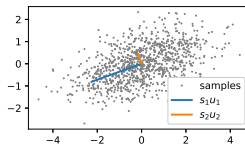
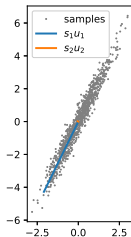
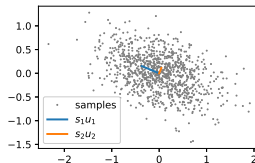
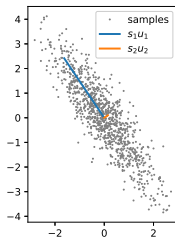
Figure 8-7. Selecting the subspace onto which to project

Principal component analysis (PCA)

For given data $x_1, x_2, \dots, x_N \in \mathbb{R}^D$

1. create a matrix $X \in \mathbb{R}^{D \times N}$ with one column vector per each sample
2. covariance matrix $\Sigma = \mathbb{E}[(X - \mathbb{E}(X))(X - \mathbb{E}(X))^T] \in \mathbb{R}^{D \times D}$
3. find singular vectors and singular values of Σ
4. principal components = largest singular values and vectors

Principal component analysis (PCA)



Principal component analysis (PCA)

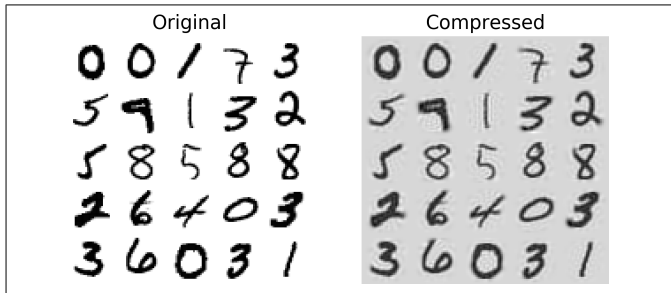


Figure 8-9. MNIST compression preserving 95% of the variance