

4. Dimension Reduction Techniques

Statistical Modelling & Machine Learning

Jaejik Kim

Department of Statistics
Sungkyunkwan University

STA3036

Why Dimension Reduction?

- ▶ Visualization of data (unsupervised learning).
- ▶ To improve prediction, the reduction of the number of dimensions of input space is important.
 - ▶ The curse of dimensionality.
 - ▶ Irrelevant input variables make data pattern unclear.
- ▶ Methods to reduce the # of input dimensions:
 - ▶ Dimension reduction techniques.
 - ▶ Filtering input variables (removing irrelevant inputs).
- ▶ Why dimension reduction?
 - ▶ Data visualization.
 - ▶ Reduction of computing time for predictive model building.
 - ▶ Unsupervised dimension reduction techniques do not guarantee the improvement of predictive models.

Dimension Reduction Techniques

Unsupervised learning:

- ▶ Principal component analysis (PCA): Linear projection (i.e., $\mathbf{Z} = \mathbf{XA}$, where \mathbf{A} is the $p \times p$ projection matrix).
- ▶ Principal curve / principal surface: It directly finds one or two-dimensional nonlinear principal components.
- ▶ Kernel PCA: Kernel transformation of \mathbf{X} , and then PCA.
- ▶ Non-negative matrix factorization (NMF): It works only for non-negative data (e.g., image data).
- ▶ Independent component analysis (ICA): Latent variables that are statistically independent and non-Gaussian.

Supervised learning:

- ▶ Partial least squares: Directions maximizing the covariance with Y .

Principal Component Analysis

- ▶ Let \mathbf{X} be $n \times p$ matrix with standardized variables.
- ▶ $\mathbf{Z} = \mathbf{X}\mathbf{V}$: Principal components.
 - ▶ $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_p]$, where \mathbf{v}_j is the eigenvector of $\mathbf{X}^\top \mathbf{X}$ corresponding to the j th largest eigenvalue.
- ▶ The singular value decomposition (SVD) of \mathbf{X} :

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top.$$

- ▶ \mathbf{U} is a $n \times p$ orthogonal matrix with columns which span the column space of \mathbf{X} .
- ▶ \mathbf{V} is a $p \times p$ orthogonal matrix with columns which span the row space of \mathbf{X} .
- ▶ \mathbf{D} is a $p \times p$ diagonal matrix with diagonal entries $d_1 \geq d_2 \geq \dots \geq d_q > d_{q+1} = \dots = d_p = 0$ where $q = \text{rank}(\mathbf{X})$. The entries d_1, \dots, d_p are called the *singular values* of \mathbf{X} .

Principal Component Analysis

- ▶ Eigen decomposition of $\mathbf{X}^\top \mathbf{X}$: By SVD,

$$\mathbf{X}^\top \mathbf{X} = \mathbf{V} \mathbf{D} \mathbf{U}^\top \mathbf{U} \mathbf{D} \mathbf{V}^\top = \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top.$$

- ▶ \mathbf{V} : The columns of \mathbf{V} are eigenvectors of $\mathbf{X}^\top \mathbf{X}$.
- ▶ Diagonal elements of \mathbf{D}^2 : d_1^2, \dots, d_p^2 are eigenvalues of $\mathbf{X}^\top \mathbf{X}$.
- ▶ PC: $\mathbf{Z} = \mathbf{X} \mathbf{V} = \mathbf{U} \mathbf{D} \mathbf{V}^\top \mathbf{V} = \mathbf{U} \mathbf{D}$.
- ▶ PCA using rank- q linear model: $f(\boldsymbol{\lambda}) = \mathbf{V}_q \boldsymbol{\lambda}$,
 - ▶ \mathbf{V}_q : $p \times q$ orthogonal column matrix.
 - ▶ $\boldsymbol{\lambda}$: $q \times 1$ vector of parameters.
 - ▶ Fitting $f(\boldsymbol{\lambda})$ to the data \mathbf{X} by minimizing the reconstruction error

$$\min_{\{\boldsymbol{\lambda}_i\}, \mathbf{V}_q} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{V}_q \boldsymbol{\lambda}_i)^2.$$

- ▶ For fixed \mathbf{V}_q , $\hat{\boldsymbol{\lambda}}_i = \mathbf{V}_q^\top \mathbf{x}_i$.

Principal Component Analysis

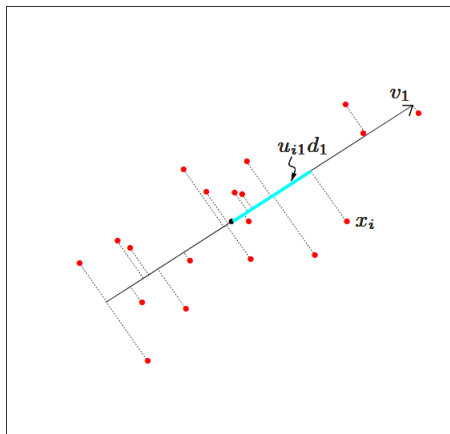
- ▶ Minimization of the reconstruction error becomes

$$\min_{\mathbf{V}_q} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{V}_q \mathbf{V}_q^\top \mathbf{x}_i)^2.$$

- ▶ The solution of \mathbf{V}_q : The matrix with the first q columns of \mathbf{V} from SVD.
- ▶ The first q PC's: $\mathbf{Z}_q = \mathbf{XV}_q = \mathbf{UDV}^\top \mathbf{V}_q = \mathbf{U}_q \mathbf{D}_q$.

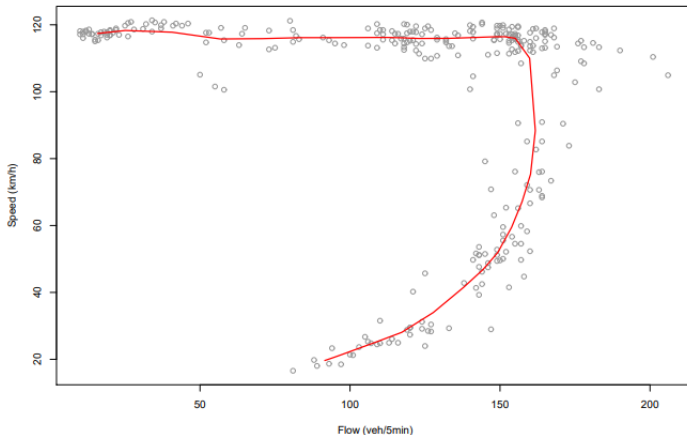
Principal Component Analysis

- ▶ The first PC line in $\mathbf{X} = (X_1, X_2)$ space:
 - ▶ The first PC direction $\mathbf{v}_1 = \mathbf{V}_1$ ($q = 1$).
 - ▶ $\hat{\lambda}_i = u_{i1}d_1$: Distance along the PC line from the origin.



Principal Curves / Surfaces

Example: Speed-Flow data from a Californian 'freeway'.



Principal Curves

- ▶ Principal curves:
 - ▶ Smooth curves passing through the middle of the data cloud.
 - ▶ One-dimensional curved approximation to a set of data points in p -dimension.
 - ▶ X_1, \dots, X_p have nonlinear relationships.
 - ▶ Visualization of data.
- ▶ The principal curve for random variables $\mathbf{X} \in \mathbb{R}^p$:
 $\mathbf{f}(\lambda) = [f_1(\lambda), \dots, f_p(\lambda)]$
 - ▶ A smooth curve in \mathbb{R}^p .
 - ▶ A vector function with p coordinates indexed by λ .
 - ▶ Each function $f_j(\lambda)$ is a coordinate function of a single parameter λ .
 - ▶ λ : The arc-length along the curve from some fixed origin.

Principal Curves

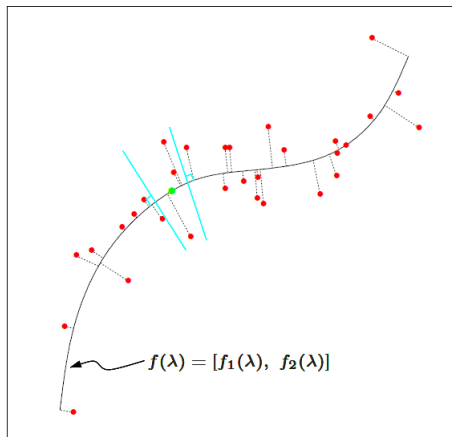


Figure: Each point on the curve is the average of all data points that project there.

Principal Curves

- ▶ Principal curve:

$$\mathbf{f}(\lambda) = E(\mathbf{X} | \lambda_f(\mathbf{X}) = \lambda),$$

- ▶ $\lambda_f(\mathbf{X})$: The closest point on curve to a data point \mathbf{x} (Euclidean).
- ▶ $\mathbf{f}(\lambda)$: The average of all data points that project to $\mathbf{f}(\lambda)$
 \Rightarrow 'self-consistency' property.
- ▶ In practice, for continuous multivariate distributed \mathbf{X} , principal curves are not unique.

Principal Curves

- ▶ To find a principal curve, we need to iteratively find $[f_1(\lambda), \dots, f_p(\lambda)]$.
- ▶ **Algorithm:** For given data points $\mathbf{x}_1, \dots, \mathbf{x}_n$,
 - (1) Initialize $f_1(\lambda), \dots, f_p(\lambda)$ using the first principal component direction.
 - (2) Initialize $\lambda_1, \dots, \lambda_n$ using the first principal component values.
 - (3) Update $\hat{f}_j(\lambda) \leftarrow \hat{E}[X_j | \hat{\lambda}_f(\mathbf{X}) = \lambda]$ for $j = 1, \dots, p$ using a smoother.
 - (4) Update $\hat{\lambda}_f(\mathbf{X}) \leftarrow \arg \max_{\lambda} \| \mathbf{X} - \hat{\mathbf{f}}(\lambda) \|$ by the projection on the curve.
 - (5) Iterate (3) and (4) until convergence.
- ▶ Tuning parameter: df of smoother in step (3).

Principal Surfaces

- ▶ Let the surface be defined in $q < p$ dimensions.
- ▶ The most commonly used is 2D surface. When $q = 2$,
 $\mathbf{f}(\lambda_1, \lambda_2) = [\mathbf{f}_1(\lambda_1, \lambda_2), \dots, \mathbf{f}_p(\lambda_1, \lambda_2)]$.
- ▶ Initial using plane defined by the first two principal component directions.
- ▶ Use 2D scatterplot smoother in algorithm to estimate $f_j(\lambda_1, \lambda_2)$, $j = 1, \dots, p$.
- ▶ Data points can be plotted by estimated nonlinear coordinates $(\hat{\lambda}_1(\mathbf{x}_i), \hat{\lambda}_2(\mathbf{x}_i))$.

Principal Surfaces

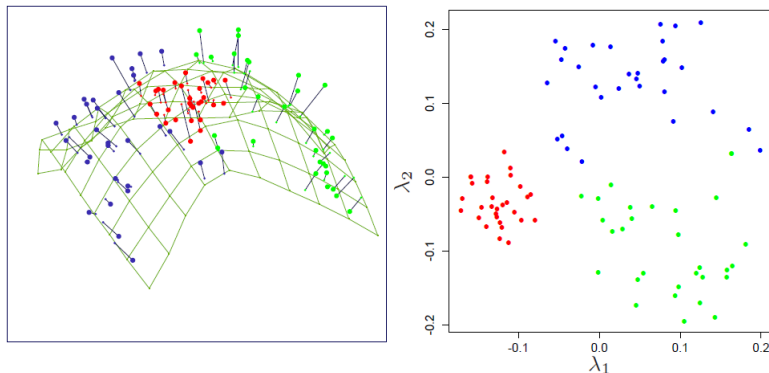
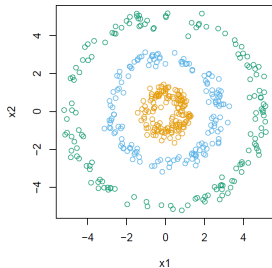


Figure: Left panel: Fitted two-dimensional surface; Right panel: Projections of data points onto the surface, resulting in coordinates $(\hat{\lambda}_1, \hat{\lambda}_2)$.

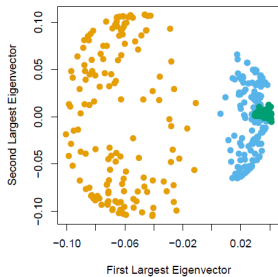
- ▶ Kernel PCA: Nonlinear version of PCA.
- ▶ Nonlinear transformation of inputs \Rightarrow PCA for the transformed feature space.
- ▶ PCA:
 - ▶ Let $\mathbf{K} = \mathbf{X}\mathbf{X}^\top$, where \mathbf{X} is the standardized input matrix, and let $\mathbf{Z} = \mathbf{U}\mathbf{D}$, where \mathbf{Z} is principal component matrix.
 - ▶ $\mathbf{K} = \{\langle \mathbf{x}_i, \mathbf{x}_{i'} \rangle\}$: It consists of inner-product terms.
 - ▶ By SVD, $\mathbf{K} = \mathbf{U}\mathbf{D}\mathbf{V}^\top \mathbf{V}\mathbf{D}\mathbf{U}^\top = \mathbf{U}\mathbf{D}^2\mathbf{U}^\top \Rightarrow$ Eigen decomposition of \mathbf{K} .
 - ▶ This means that the principal component matrix can be obtained by the eigen decomposition of \mathbf{K} .
- ▶ Kernel PCA simply mimics this procedure of PCA.

- ▶ Kernel matrix $\mathbf{K} = \{K(\mathbf{x}_i, \mathbf{x}_{i'})\}$: An inner-product matrix of the implicit features $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_{i'}) \rangle$.
- ▶ From the eigen decomposition of \mathbf{K} , \mathbf{U} and \mathbf{D} can be obtained.
- ▶ The m th kernel PC: \mathbf{z}_m = the m th column of $\mathbf{Z} = \mathbf{U}\mathbf{D}$.
- ▶ $z_{im} = \sum_{j=1}^n \alpha_{jm} K(\mathbf{x}_i, \mathbf{x}_j)$, where $\alpha_{jm} = u_{jm} / d_m$.
- ▶ We do NOT have to specify transformed feature $\phi(\mathbf{x})$. To construct \mathbf{K} , a kernel function should be chosen.
- ▶ Kernel functions:
 - ▶ Radial: $K(\mathbf{x}_i, \mathbf{x}_{i'}) = \exp(-\|\mathbf{x}_i - \mathbf{x}_{i'}\|^2 / c)$.
 - ▶ Polynomial: $K(\mathbf{x}_i, \mathbf{x}_{i'}) = (1 + \mathbf{x}_i^\top \mathbf{x}_{i'})^d$.
 - ▶ Neural network: $K(\mathbf{x}_i, \mathbf{x}_{i'}) = \tanh(\mathbf{x}_i^\top \mathbf{x}_{i'} + \delta)$.

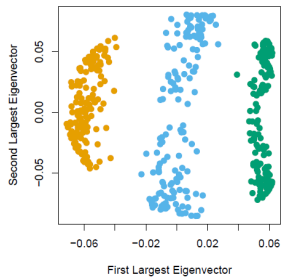
Kernel PCA



Radial Kernel ($c=2$)



Radial Kernel ($c=10$)



Non-negative Matrix Factorization (NMF)

- ▶ Alternative approach to PCA.
- ▶ Data and components are assumed to be non-negative (e.g., image data).
- ▶ The $n \times p$ input matrix \mathbf{X} is approximated by

$$\mathbf{X} \approx \mathbf{WH},$$

- ▶ \mathbf{W} : $n \times r$; \mathbf{H} : $r \times p$; $r \leq \max(n, p)$.
 - ▶ $x_{ij}, w_{ik}, h_{kj} \geq 0$.
- ▶ \mathbf{W} and \mathbf{H} can be found by maximizing

$$L(\mathbf{W}, \mathbf{H}) = \sum_{i=1}^n \sum_{j=1}^p [x_{ij} \log(\mathbf{WH})_{ij} - (\mathbf{WH})_{ij}].$$

- ▶ This is the log-likelihood from a model in which $x_{ij} \sim \text{Poisson}$ with mean $(\mathbf{WH})_{ij}$.

- ▶ Algorithm converging to a local maximum of $L(\mathbf{W}, \mathbf{H})$:

- (1) Initialize \mathbf{W} & \mathbf{H} .

- (2) Update w_{ik} by

$$w_{ik} \leftarrow w_{ik} \frac{\sum_{j=1}^P h_{kj} x_{ij} / (\mathbf{WH})_{ij}}{\sum_{j=1}^P h_{kj}}.$$

- (3) Update h_{kj} by

$$h_{kj} \leftarrow h_{kj} \frac{\sum_{i=1}^P w_{ik} x_{ij} / (\mathbf{WH})_{ij}}{\sum_{i=1}^P w_{ik}}.$$

- (4) Iterate steps (2) & (3) until convergence.

- ▶ This algorithm is related to the iterative proportional scaling algorithm for log-linear models.

- ▶ NMF may not be unique. This means that the optimum value of the algorithm depends on initial values.
- ▶ The non-uniqueness could hamper the interpretability of the factorization.
- ▶ Approximation to the columns of $\mathbf{X} \Rightarrow$ Columns of \mathbf{W} (primary non-negative components).
- ▶ Approximation to the rows of $\mathbf{X} \Rightarrow$ Rows of \mathbf{H} (archetypal data points).

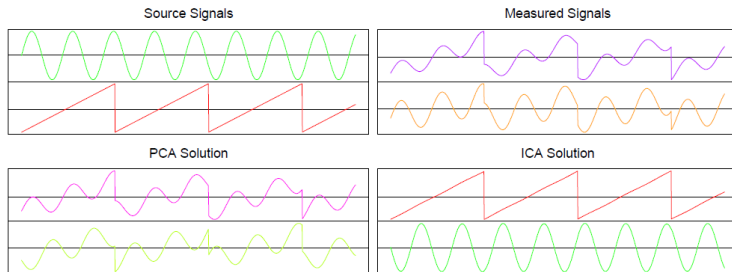
Independent Component Analysis (ICA)

- ▶ PC's are uncorrelated linear combinations of $X = (X_1, \dots, X_p)$.
- ▶ Statistical independent components (linear combinations of X)
 \Rightarrow Independent component analysis.
- ▶ ICA problem:

$$X = \mathbf{A}S,$$

- ▶ $X = (X_1, \dots, X_p)^\top$: A $p \times 1$ random vector representing multivariate input measurements.
- ▶ $S = (S_1, \dots, S_p)^\top$: A $p \times 1$ latent source vector. S_1, \dots, S_p are independent random variables.
- ▶ \mathbf{A} : $p \times p$ mixing matrix.

- ▶ The goal of ICA:
 - ▶ Estimation of \mathbf{A} .
 - ▶ Estimation the source distributions $S_j \sim f_{S_j}$, $j = 1, \dots, p$.
- ▶ Cocktail party problem:
 - ▶ p independent sources of sound in a room.
 - ▶ p microphones placed around the room hear different mixtures.
 - ▶ Recovery of the sources.



- ▶ Assume that $E(S) = \mathbf{0}$ and $\text{Cov}(S) = \mathbf{I}$.

$$\text{Cov}(X) = \text{Cov}(\mathbf{A}S) = \mathbf{A}\text{Cov}(S)\mathbf{A}^\top = \mathbf{A}\mathbf{A}^\top.$$

- ▶ $X = \mathbf{A}S$: Each correlated X_j is represented as a linear expansion in the uncorrelated, unit variance variables S_j .
- ▶ For any given orthogonal $p \times p$ matrix \mathbf{R} ,

$$X = \mathbf{A}S = \mathbf{A}\mathbf{R}^\top \mathbf{R}S = \mathbf{A}^*S^*,$$

- ▶ $\text{Cov}(S^*) = \mathbf{R}\text{Cov}(S)\mathbf{R}^\top = \mathbf{I}$.
- ▶ Non-uniqueness of \mathbf{A} and S under uncorrelated and Gaussian assumptions.
- ▶ Impossible to identify any particular latent source variables as unique underlying sources.

- ▶ To recover \mathbf{A} , ICA assumes statistical independence and non-Gaussianity of S .
- ▶ Under the uncorrelatedness and Gaussian assumptions, the second order moment (covariance) of random variables is required.
- ▶ Independence and non-Gaussianity require higher order moments.
- ▶ Like PC, Independent components are linear combinations of X . However, the linear combinations of X should be chosen to be as independent as possible.



Figure: Source $S = (S_1, S_2)^T$ is two independent uniform random variables. Data are generated by $X = AS$.

- ▶ There are a lot of methods to estimate \mathbf{A} & S in ICA.
- ▶ The most methods start with this:
 - ▶ Let $\mathbf{\Sigma} = \text{Cov}(X)$ and $X = \mathbf{A}S$. Then, $S = \mathbf{A}^{-1}X$.
 - ▶ Let $\mathbf{A}^{-1} = \mathbf{W}\mathbf{\Sigma}^{-\frac{1}{2}}$ for some non-singular \mathbf{W} . Then,
 $S = \mathbf{A}^{-1}X = \mathbf{W}\mathbf{\Sigma}^{-\frac{1}{2}}X$ with $\text{Cov}(S) = \mathbf{I}$ and orthonormal \mathbf{W} .
- ▶ To find S , transform the data $\tilde{X} = \mathbf{\Sigma}^{-\frac{1}{2}}X$, and then seek an orthogonal matrix \mathbf{W} so that the components $S = \mathbf{W}\tilde{X}$ are independent.

Methods to find the orthonormal matrix \mathbf{W} :

- ▶ Mutual information and entropy, maximizing non-Gaussianity.
 - ▶ FastICA, Infomax methods.
 - ▶ Kullback-Leibler divergence: A measure of dependence between two random vectors.
 - ▶ Using the K-L divergence, it looks for low-entropy or non-Gaussian projections \mathbf{W} .
- ▶ Likelihood methods:
 - ▶ ProdDenICA.
 - ▶ It directly fits the independent density model $f_S(s) = \prod_{j=1}^P f_j(s_j)$ to find \mathbf{W} .
- ▶ etc.