

Dimensionality Reduction

Shan-Hung Wu
shwu@cs.nthu.edu.tw

Department of Computer Science,
National Tsing Hua University, Taiwan

NetDB-ML, Spring 2015

1 Why Dimensionality Reduction?

2 Feature Selection

- Forward and Backward Selection (Supervised)

3 Feature Extraction

- Principal Component Analysis (Unsupervised, Linear)
- Linear Discriminant Analysis (Supervised, Linear)
- Kernel PCA (Unsupervised, Nonlinear)
- Isometric Feature Mapping (Unsupervised, Nonlinear)
- Locally Linear Embedding (Unsupervised, Nonlinear)

1 Why Dimensionality Reduction?

2 Feature Selection

- Forward and Backward Selection (Supervised)

3 Feature Extraction

- Principal Component Analysis (Unsupervised, Linear)
- Linear Discriminant Analysis (Supervised, Linear)
- Kernel PCA (Unsupervised, Nonlinear)
- Isometric Feature Mapping (Unsupervised, Nonlinear)
- Locally Linear Embedding (Unsupervised, Nonlinear)

Why Dimensionality Reduction?

- Given a (normalized) training dataset $\{\mathbf{x}^{(t)}, \mathbf{r}^{(t)}\}$ where $\mathbf{x}^{(t)} \in \mathbb{R}^d$, we want to reduce the input dimension from d to k , $k < d$
- Why dimensionality reduction?

Why Dimensionality Reduction?

- Given a (normalized) training dataset $\{\mathbf{x}^{(t)}, \mathbf{r}^{(t)}\}$ where $\mathbf{x}^{(t)} \in \mathbb{R}^d$, we want to reduce the input dimension from d to k , $k < d$
- Why dimensionality reduction?
- To reduce the time/space requirements in training, cross validation, testing, and prediction
- To save the cost of data collection
 - If an input is decided to be unnecessary, we can stop observing it in the future
- To make the classifier/regressor robust to small datasets
 - Fewer parameters (e.g., elements of Σ) implies lower model complexity, and lower variance due to particulars of a sample
- To extract knowledge
 - The k inputs can describe the whole data and may explain something

Feature Selection vs. Feature Extraction

- Basically, two ways to dimensionality reduction:
- **Feature selection:**
 - Keeping only k of d dimensions that is most helpful to classification/regression
- **Feature extraction:**
 - Finding a new set of k dimensions that are combinations of the original d dimensions and most helpful to classification/regression
 - The combination can be **linear** or **nonlinear**
- Can be **supervised** or **unsupervised** depending on whether or not the labels $r^{(t)}$ are used

Caution

- Is dimension reduction a preferred step in data preprocessing?

- Is dimension reduction a preferred step in data preprocessing?
Generally, yes
- But if your classifier/regressor make certain assumptions over the data, you should perform dimension reduction *only if those assumptions are still valid after the reduction*
 - Luckily, if we assume that the data are normally distributed within each class, then after linear dimension reduction, the data are still normal

1 Why Dimensionality Reduction?

2 Feature Selection

- Forward and Backward Selection (Supervised)

3 Feature Extraction

- Principal Component Analysis (Unsupervised, Linear)
- Linear Discriminant Analysis (Supervised, Linear)
- Kernel PCA (Unsupervised, Nonlinear)
- Isometric Feature Mapping (Unsupervised, Nonlinear)
- Locally Linear Embedding (Unsupervised, Nonlinear)

Sequential Forward Selection

- Denote F the set of selected inputs
- In **sequential forward selection**, we start with no input $F = \emptyset$, and add inputs one by one
- At each step, we add x_j to F if x_j decreases the error most
 - Specifically, we add x_j to F if

$$j = \arg_j \min err(F \cup x_j)$$

and

$$err(F \cup x_j) < err(F) - \varepsilon,$$

where $err(F)$ is the error of predictions (on the validation set) made by a classifier/regressor considering only the inputs in F and ε is a user-defined threshold

- We stop adding inputs if there is no decrease in error, or the decrease in error is too small

Sequential Backward Selection

- In *sequential backward selection*, we start with F containing all inputs, and remove inputs one by one
- At each step, we remove x_j from F if x_j increase the error least
 - Specifically, we remove x_j from F if

$$j = \arg_i \min \text{err}(F - x_i) \text{ and } \text{err}(F - x_j) < \text{err}(F) + \epsilon$$

- In either direction, need to perform the classification/regression $d + (d - 1) + \dots + (d - k) = O(d^2)$ times
 - But the training process in the backward direction is more costly, as F contains more inputs
 - Forward direction is preferable if k is much smaller than d

Are We Satisfied?

- The above algorithms are local search procedures and do not guarantee F to be optimal
- In some applications such as face recognition (where $\mathbf{x}^{(t)}$ are images and the inputs are pixels), feature selection is not a good method for dimensionality reduction
 - Individual pixels do not carry much discriminative information; it is the combinations of pixels that carry face identities

1 Why Dimensionality Reduction?

2 Feature Selection

- Forward and Backward Selection (Supervised)

3 Feature Extraction

- Principal Component Analysis (Unsupervised, Linear)
- Linear Discriminant Analysis (Supervised, Linear)
- Kernel PCA (Unsupervised, Nonlinear)
- Isometric Feature Mapping (Unsupervised, Nonlinear)
- Locally Linear Embedding (Unsupervised, Nonlinear)

Principal Component Analysis (1/2)

- In **Principal Component Analysis (PCA)**, we want to find k vectors $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^d$, named **principal components**, such that after projecting each instance $\mathbf{x}^{(t)}$ onto directions along $\mathbf{w}_1, \dots, \mathbf{w}_k$ and obtaining $\mathbf{z}^{(t)} = [\mathbf{w}_1 \ \dots \ \mathbf{w}_k]^\top \mathbf{x}^{(t)} \in \mathbb{R}^k$, the new data collection $\{\mathbf{z}^{(t)}\}_{t=1}^N$ helps the classifier/regressor most
- Let $k = 1$, how to find \mathbf{w}_1 ?

Principal Component Analysis (1/2)

- In *Principal Component Analysis (PCA)*, we want to find k vectors $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^d$, named *principal components*, such that after projecting each instance $\mathbf{x}^{(t)}$ onto directions along $\mathbf{w}_1, \dots, \mathbf{w}_k$ and obtaining $\mathbf{z}^{(t)} = [\mathbf{w}_1 \ \dots \ \mathbf{w}_k]^\top \mathbf{x}^{(t)} \in \mathbb{R}^k$, the new data collection $\{\mathbf{z}^{(t)}\}_{t=1}^N$ helps the classifier/regressor most
- Let $k = 1$, how to find \mathbf{w}_1 ?
 - Of course we won't pick $\mathbf{w}_1 = \mathbf{0}$, because it makes $\mathbf{z}^{(t)} = \mathbf{w}_1^\top \mathbf{x}^{(t)} = 0$ for all t and undistinguishable
 - It is plausible to pick \mathbf{w}_1 such that $\mathbf{z}^{(t)}$ is most spread out; that is, let $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$ be i.i.d. samples drawn from an (unknown) population \mathbf{x} , we can pick \mathbf{w}_1 which maximizes $\text{Var}(\mathbf{z}) = \text{Var}(\mathbf{w}_1^\top \mathbf{x}) = \mathbf{w}_1^\top \boldsymbol{\Sigma} \mathbf{w}_1$

Principal Component Analysis (2/2)

- How to obtain the covariant matrix Σ ?

Principal Component Analysis (2/2)

- How to obtain the covariant matrix Σ ?
 - Recall that $\Sigma_x = E[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}_x\boldsymbol{\mu}_x^\top$
 - Assuming that \mathbf{x} is centered (e.g., z-normalized), we have an estimate
$$\mathbf{S} = \frac{1}{N} \sum_{t=1}^N \mathbf{x}^{(t)}\mathbf{x}^{(t)\top}$$
- Given arbitrary k , PCA picks $\mathbf{w}_1, \dots, \mathbf{w}_k$ such that
 - $Var(z_i) = \mathbf{w}_i^\top \mathbf{S} \mathbf{w}_i$ for $i = 1, \dots, k$ are the highest
 - $\mathbf{w}_1, \dots, \mathbf{w}_k$ are orthogonal to each other so each of them helps the classifier/regressor independently
 - $\|\mathbf{w}_i\| = 1$ for $i = 1, \dots, k$ so the direction of \mathbf{w}_i is the only factor that affects $\mathbf{w}_i^\top \mathbf{S} \mathbf{w}_i$ (it is maximized not because we pick a large \mathbf{w}_i)

Solving w_1, \dots, w_k (1/3)

- To get w_1 , we need to solve the problem $\arg_{w_1} \max w_1^T S w_1$ subject to $w_1^T w_1 = 1$
 - w_1 is a stationary point of the Lagrangian: $w_1^T S w_1 - \alpha(w_1^T w_1 - 1)$
 - Taking the partial derivatives with respect to w_1 and α we have
$$\begin{cases} 2S w_1 - 2\alpha w_1 = 0 \\ w_1^T w_1 - 1 = 0 \end{cases}$$
, implying $S w_1 = \alpha w_1$ and $w_1^T S w_1 = \alpha w_1^T w_1 = \alpha$
 - The candidates of w_1 and α are eigenvectors and eigenvalues of Σ respectively
 - Since we want to maximize $w_1^T S w_1 = \alpha$, w_1 is the eigenvector corresponding to the largest eigenvalue λ_1

Solving w_1, \dots, w_k (2/3)

- To get w_2 , we solve the problem $\arg_{w_2} \max w_2^\top S w_2$ subject to $w_2^\top w_2 = 1$ and $w_2^\top w_1 = 0$
 - w_2 is a stationary point of the Lagrangian:
 $w_2^\top S w_2 - \alpha(w_2^\top w_2 - 1) - \beta w_2^\top w_1$
 - Taking the partial derivatives with respect to w_2 , α , and β we have
$$\begin{cases} 2S w_2 - 2\alpha w_2 - \beta w_1 = 0 \\ w_2^\top w_2 - 1 = 0 \\ w_2^\top w_1 = 0 \end{cases}$$
 - This implies $S w_2 = \alpha w_2$, as
$$0 = w_1^\top 0 = 2w_1^\top S w_2 - 2\alpha w_1^\top w_2 - \beta w_1^\top w_1 = 2w_2^\top S w_1 - \beta = 2\lambda_1 w_2^\top w_1 - \beta = -\beta$$
, hence $2S w_2 - 2\alpha w_2 = 0$
 - Additionally, $w_2^\top S w_2 = \alpha$, as
$$0 = w_2^\top 0 = 2w_2^\top S w_2 - 2\alpha w_2^\top w_2 - \beta w_2^\top w_1 = 2w_2^\top S w_2 - 2\alpha$$
 - The candidates of w_2 are eigenvectors of S orthogonal to w_1
 - Since we want to maximize $w_2^\top S w_2 = \alpha$, w_2 is the eigenvector corresponding to the second largest eigenvalue λ_2
 - Recall that the eigenvectors of a symmetric matrix (S) are orthogonal

Solving w_1, \dots, w_k (3/3)

- From the above, we can see that w_i are the i th largest eigenvector of S

We can reach the same conclusion by Rayleigh's quotient [Proof]

- Define $z = W^T x$, where $W = [w_1, \dots, w_k]$
 - If x are not normalized, we can center z around the origin by letting $z = W^T (x - m)$, where m is the sample mean of x
 - In addition, we can make z_i have the unit variance by dividing itself by $\sqrt{\lambda_i}$
- $z^{(t)}$ are i.i.d. samples of z

Effects of PCA

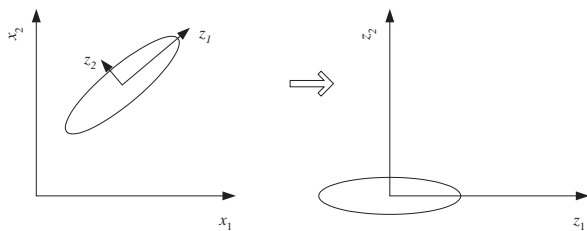
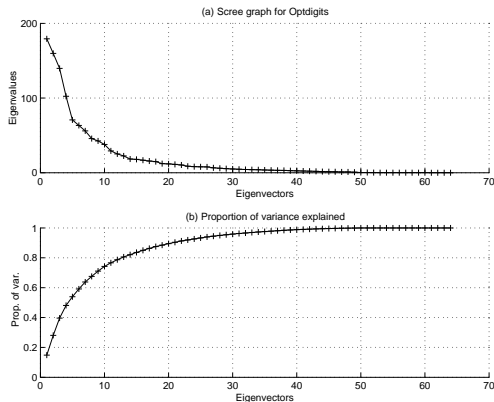


Figure : PCA centers the instances and rotates the axes to line up with the directions of the highest variance. With these new axes, the covariance matrix $\Sigma_z = \mathbf{W}^T \mathbf{S} \mathbf{W} \in \mathbb{R}^{k \times k}$ is always diagonal, making the naive Bayes' classifiers feasible.

The Scree Graph

- How do we decide a proper k ?
 - Generally, we want to pick k such that the **proportion of variance** $\frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i}$ is more than 90%
- We seek for the “elbow” in the **scree graph**:



Are We Satisfied? (1/3)

- In image processing applications (where $\mathbf{x}^{(t)}$ are images), \mathbf{w}_i themselves can also be displayed as images and be seen as templates for important features
 - In these cases, \mathbf{w}_i have a special name *eigenfaces*
- When should we use PCA?

Are We Satisfied? (1/3)

- In image processing applications (where $\mathbf{x}^{(t)}$ are images), \mathbf{w}_i themselves can also be displayed as images and be seen as templates for important features
 - In these cases, \mathbf{w}_i have a special name *eigenfaces*
- When should we use PCA?
 - PCA is helpful only when we have a small number of large eigenvalues
 - Or, when the original inputs of \mathbf{x} are highly correlated (so the contours are stretched)
- Note the variance of each input of \mathbf{x} may be very high that affects the directions of principal components more than the correlations between inputs do
 - Generally, we perform z-normalization before applying PCA

Are We Satisfied? (2/3)

- PCA is a one-group procedure
 - Fine with the regression
 - But in the case of classification, there are multiple groups
- The Karhunen-Loeve expansion:
 - Instead of using the covariance matrix of the whole examples, we can estimate separate covariance matrices for individual classes
 - Take their average (weighted by the estimated priors) as the covariance matrix, and use its eigenvectors

1 Why Dimensionality Reduction?

2 Feature Selection

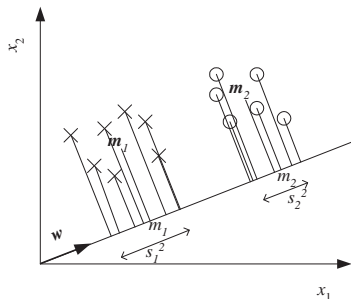
- Forward and Backward Selection (Supervised)

3 Feature Extraction

- Principal Component Analysis (Unsupervised, Linear)
- Linear Discriminant Analysis (Supervised, Linear)
- Kernel PCA (Unsupervised, Nonlinear)
- Isometric Feature Mapping (Unsupervised, Nonlinear)
- Locally Linear Embedding (Unsupervised, Nonlinear)

Linear Discriminant Analysis (1/3)

- Given a training set $\mathcal{X} = \{\mathbf{x}^{(t)}, \mathbf{r}^{(t)}\}_{t=1}^N$, in **Linear Discriminant Analysis (LDA)**, we want to find k vectors $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^d$, such that after projecting each instance $\mathbf{x}^{(t)}$ onto directions along $\mathbf{w}_1, \dots, \mathbf{w}_k$ and obtaining $\mathbf{z}^{(t)} = \mathbf{W}^\top \mathbf{x}^{(t)} \in \mathbb{R}^k$, where $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_k]$, the new data collection $\{\mathbf{z}^{(t)}\}_{t=1}^N$ has the properties:
 - Examples in different classes are as separate as possible
 - Examples of the same class are as close as possible



Linear Discriminant Analysis (2/3)

- Let $\mathbf{m}_i = \frac{1}{N_i} \sum_{t=1}^N r_i^{(t)} \mathbf{x}^{(t)}$ and $\mathbf{S}_i = \frac{1}{N_i-1} \sum_{t=1}^N r_i^{(t)} (\mathbf{x}^{(t)} - \mathbf{m}_i)(\mathbf{x}^{(t)} - \mathbf{m}_i)^\top$, where $N_i = \sum_{t=1}^N r_i^{(t)}$, be the estimated mean and variance of examples in class i respectively
- We can measure the total within-class separation by $\mathbf{S}_{within} = \sum_{i=1}^K N_i \mathbf{S}_i$
- Denote $\mathbf{m} = \frac{1}{K} \sum_{i=1}^K \mathbf{m}_i$, we can also measure the separation between classes by $\mathbf{S}_{between} = \sum_{i=1}^K N_i (\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^\top$
- For each \mathbf{w}_j , the mean and variance of the projections $z^{(t)}$ are $m_i = \frac{1}{N_i} \sum_{t=1}^N r_i^{(t)} z_j^{(t)} = \frac{1}{N_i} \sum_{t=1}^N r_i^{(t)} \mathbf{w}_j^\top \mathbf{x}^{(t)} = \mathbf{w}_j^\top \left(\frac{1}{N_i} \sum_{t=1}^N r_i^{(t)} \mathbf{x}^{(t)} \right) = \mathbf{w}_j^\top \mathbf{m}_i$ and $s_i^2 = \frac{1}{N_i-1} \sum_{t=1}^N r_i^{(t)} (z_j^{(t)} - \mathbf{w}_j^\top \mathbf{m}_i)(z_j^{(t)} - \mathbf{w}_j^\top \mathbf{m}_i)^\top = \frac{1}{N_i-1} \sum_{t=1}^N r_i^{(t)} (\mathbf{w}_j^\top \mathbf{x}^{(t)} - \mathbf{w}_j^\top \mathbf{m}_i)(\mathbf{w}_j^\top \mathbf{x}^{(t)} - \mathbf{w}_j^\top \mathbf{m}_i)^\top = \frac{1}{N_i-1} \sum_{t=1}^N r_i^{(t)} \mathbf{w}_j^\top (\mathbf{x}^{(t)} - \mathbf{m}_i)(\mathbf{x}^{(t)} - \mathbf{m}_i)^\top \mathbf{w}_j = \mathbf{w}_j^\top \mathbf{S}_i \mathbf{w}_j$ respectively
- The within- and between-class separation after projection becomes $\mathbf{w}_j^\top \mathbf{S}_{within} \mathbf{w}_j$ and $\mathbf{w}_j^\top \mathbf{S}_{between} \mathbf{w}_j$ respectively

Linear Discriminant Analysis (3/3)

- LDA picks $\mathbf{w}_1, \dots, \mathbf{w}_k$ such that
 - The **Fisher's linear discriminant** $J(\mathbf{w}_j) = \frac{\mathbf{w}_j^\top \mathbf{S}_{between} \mathbf{w}_j}{\mathbf{w}_j^\top \mathbf{S}_{within} \mathbf{w}_j}$ are the highest for $j = 1, \dots, k$
 - $\mathbf{w}_1, \dots, \mathbf{w}_k$ are orthogonal to each other
- We don't care $\|\mathbf{w}_j\|$ now since we are maximizing the ratio
- Unfortunately, the solution to the above objective is not unique since if $\mathbf{w}_1, \dots, \mathbf{w}_k$ are the solution, so do $c_1 \mathbf{w}_1, \dots, c_k \mathbf{w}_k, \forall c_j \in \mathbb{R}$
- We can instead maximize $\mathbf{w}_j^\top \mathbf{S}_{between} \mathbf{w}_j$ subject to an additional constrain $\mathbf{w}_j^\top \mathbf{S}_{within} \mathbf{w}_j = 1$

Solving $\mathbf{w}_1, \dots, \mathbf{w}_k$

- Lagrangian for \mathbf{w}_1 : $\arg_{\mathbf{w}_1} \max \mathbf{w}_1^\top \mathbf{S}_{between} \mathbf{w}_1 - \alpha(\mathbf{w}_1^\top \mathbf{S}_{within} \mathbf{w}_1 - 1)$
 - Taking the partial derivatives with respect to \mathbf{w}_1 and α and setting them to zero we have $\begin{cases} 2\mathbf{S}_{between} \mathbf{w}_1 - 2\alpha \mathbf{S}_{within} \mathbf{w}_1 = \mathbf{0} \\ \mathbf{w}_1^\top \mathbf{S}_{within} \mathbf{w}_1 - 1 = 0 \end{cases}$, leading to $(\mathbf{S}_{within})^{-1} \mathbf{S}_{between} \mathbf{w}_1 = \alpha \mathbf{w}_1$ and $\mathbf{w}_1^\top \mathbf{S}_{between} \mathbf{w}_1 = \alpha$
 - The candidates of \mathbf{w}_1 and α are eigenvectors and eigenvalues of $(\mathbf{S}_{within})^{-1} \mathbf{S}_{between}$ respectively
 - Since we want to maximize $\mathbf{w}_1^\top \mathbf{S}_{between} \mathbf{w}_1 = \alpha$, \mathbf{w}_1 is the eigenvector corresponding to the largest eigenvalue λ_1
- Similarly, \mathbf{w}_j is the j th largest eigenvector of $(\mathbf{S}_{within})^{-1} \mathbf{S}_{between}$
- We can reach the same conclusion using the Rayleigh's quotient

[Proof: Observe that

$$\frac{\mathbf{w}_j^\top \mathbf{S}_{between} \mathbf{w}_j}{\mathbf{w}_j^\top \mathbf{S}_{within} \mathbf{w}_j} = \frac{\mathbf{w}_j^\top \mathbf{S}_{between} \mathbf{w}_j}{(\mathbf{U}^\top \mathbf{w}_j)^\top \mathbf{U}^\top \mathbf{w}_j} = \frac{(\mathbf{U}^\top \mathbf{w}_j)^\top \mathbf{U}^{-1} \mathbf{S}_{between} (\mathbf{U}^\top)^{-1} (\mathbf{U}^\top \mathbf{w}_j)}{(\mathbf{U}^\top \mathbf{w}_j)^\top \mathbf{U}^\top \mathbf{w}_j},$$

where $\mathbf{S}_{within} = \mathbf{U}\mathbf{U}^\top$]

Deciding k

- Note $\mathbf{S}_{between}$ is the sum of K matrices, namely $(\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^\top$, of rank 1
 - $(\mathbf{m}_i - \mathbf{m})(\mathbf{m}_i - \mathbf{m})^\top \stackrel{c}{\sim} [\mathbf{m}_i - \mathbf{m}, \mathbf{0}, \dots, \mathbf{0}]$
- Given K columns $\{\mathbf{m}_1 - \mathbf{m}, \dots, \mathbf{m}_K - \mathbf{m}\}$, only $K - 1$ of them are linearly independent
 - Since $\mathbf{m} = \frac{1}{K} \sum_{i=1}^K \mathbf{m}_i$, we can express $\mathbf{m}_K - \mathbf{m}$ by the linear combination of $\mathbf{m}_1 - \mathbf{m}, \dots, \mathbf{m}_{K-1} - \mathbf{m}$
- The maximum rank of $\mathbf{S}_{between}$ is $K - 1$ and we can pick at most $K - 1$ eigenvectors (i.e., $k = K - 1$)

Are We Satisfied?

- Any classifier can be used after LDA
- To be able to apply LDA, \mathbf{S}_{within} must be invertible
 - If not, we can apply PCA first to get rid of the singularity
- When should we use LDA?

Are We Satisfied?

- Any classifier can be used after LDA
- To be able to apply LDA, \mathbf{S}_{within} must be invertible
 - If not, we can apply PCA first to get rid of the singularity
- When should we use LDA?
 - LDA works best if instances in different class are distributed in groups (e.g., normal)

PCA vs. LDA

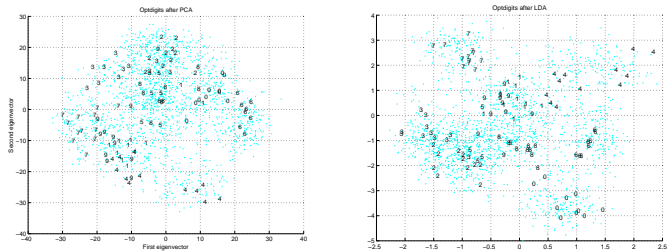


Figure : The distribution of $\mathbf{z}^{(t)}$ found by PCA (left) and LDA (right) respectively. Instances are plotted in the space of the first two attributes (out of nine). LDA, as expected, leads to a better separation between classes.

- Identify situations where PCA and LDA will find the same and totally different (orthogonal) directions [Homework]

1 Why Dimensionality Reduction?

2 Feature Selection

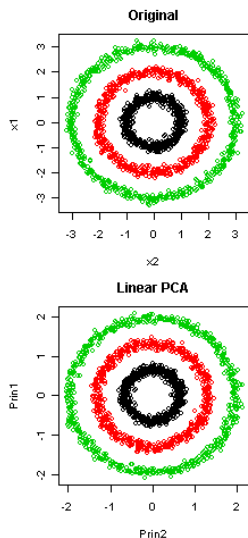
- Forward and Backward Selection (Supervised)

3 Feature Extraction

- Principal Component Analysis (Unsupervised, Linear)
- Linear Discriminant Analysis (Supervised, Linear)
- Kernel PCA (Unsupervised, Nonlinear)
- Isometric Feature Mapping (Unsupervised, Nonlinear)
- Locally Linear Embedding (Unsupervised, Nonlinear)

Why Nonlinear Mappings?

- Both PCA and LDA map examples $x^{(t)}$ to a low dimensional space (through projection W)
 - We assume that the d -dimensional input space has a linear relationship with the low dimensional space
- Sometimes, this linear mapping does not help



Kernel PCA: Basic Idea

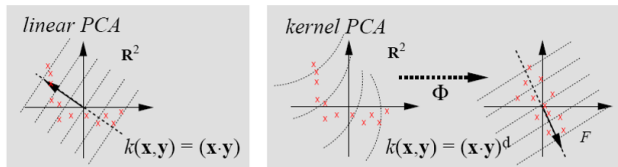


Fig. 1. Basic idea of kernel PCA: by using a nonlinear kernel function k instead of the standard dot product, we implicitly perform PCA in a possibly high-dimensional space F which is nonlinearly related to input space. The dotted lines are contour lines of constant feature value.

PCA Revisited

- Find the eigenvectors $\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(k)}$ of the matrix $\mathbf{S} = \frac{1}{N} \sum_{t=1}^N \mathbf{x}^{(t)} \mathbf{x}^{(t)\top}$ corresponding to the largest eigenvalues
- Each eigenvector $\mathbf{w}^{(i)}$ lies in the span of $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$
 - As
$$\mathbf{S} \mathbf{w}^{(i)} = \lambda^{(i)} \mathbf{w}^{(i)} \Rightarrow \mathbf{w}^{(i)} = \frac{1}{\lambda^{(i)}} \sum_{t=1}^N \mathbf{x}^{(t)} (\mathbf{x}^{(t)\top} \mathbf{w}^{(i)}) = \sum_{t=1}^N \alpha_t^{(i)} \mathbf{x}^{(t)}$$
- This holds after the lifting: $\mathbf{w}^{(i)} = \sum_{t=1}^N \alpha_t^{(i)} \Phi(\mathbf{x}^{(t)})$
 - We can instead solve N variables in $\alpha^{(i)} \in \mathbb{R}^N$ for each $\mathbf{w}^{(i)}$
- We don't need to compute $\mathbf{w}^{(i)}$ explicitly to obtain
$$\mathbf{z}^{(t)} = [\mathbf{w}^{(1)\top} \Phi(\mathbf{x}^{(t)}), \dots, \mathbf{w}^{(k)\top} \Phi(\mathbf{x}^{(t)})]^\top$$
 - $$z_i^{(t)} = \mathbf{w}^{(i)\top} \Phi(\mathbf{x}^{(t)}) = \sum_{s=1}^N \alpha_s^{(i)} k(\mathbf{x}^{(s)}, \mathbf{x}^{(t)})$$

Solving $\alpha^{(1)}, \dots, \alpha^{(k)}$ (1/3)

- Let $\mathbf{S} = \frac{1}{N} \sum_{t=1}^N \Phi(\mathbf{x}^{(t)})\Phi(\mathbf{x}^{(t)})^\top$
- Then we can write $\mathbf{w}^{(i)}$ as $\mathbf{w}^{(i)} = \sum_{t=1}^N \alpha_t^{(i)} \Phi(\mathbf{x}^{(t)})$
- Also we have $\mathbf{S}\mathbf{w}^{(i)} = \lambda^{(i)} \mathbf{w}^{(i)}$, implying that

$$\frac{1}{N} \sum_{t=1}^N \Phi(\mathbf{x}^{(t)})\Phi(\mathbf{x}^{(t)})^\top \sum_{s=1}^N \alpha_s^{(i)} \Phi(\mathbf{x}^{(s)}) = \lambda^{(i)} \sum_{s=1}^N \alpha_s^{(i)} \Phi(\mathbf{x}^{(s)})$$

Solving $\alpha^{(1)}, \dots, \alpha^{(k)}$ (2/3)

$$\begin{aligned} & \frac{1}{N} \sum_{t=1}^N \Phi(\mathbf{x}^{(t)}) \Phi(\mathbf{x}^{(t)})^\top \sum_{s=1}^N \alpha_s^{(i)} \Phi(\mathbf{x}^{(s)}) = \lambda^{(i)} \sum_{s=1}^N \alpha_s^{(i)} \Phi(\mathbf{x}^{(s)}) \\ \Rightarrow & \sum_{t=1}^N \sum_{s=1}^N \Phi(\mathbf{x}^{(t)}) (\Phi(\mathbf{x}^{(t)})^\top \Phi(\mathbf{x}^{(i)})) \alpha_s^{(i)} = N \lambda^{(i)} \sum_{s=1}^N \alpha_s^{(i)} \Phi(\mathbf{x}^{(s)}) \\ \Rightarrow & \sum_{t=1}^N \sum_{s=1}^N \Phi(\mathbf{x}^{(l)})^\top \Phi(\mathbf{x}^{(t)}) (\Phi(\mathbf{x}^{(t)})^\top \Phi(\mathbf{x}^{(i)})) \alpha_s^{(i)} \\ & \quad = N \lambda^{(i)} \sum_{s=1}^N \alpha_s^{(i)} \Phi(\mathbf{x}^{(l)})^\top \Phi(\mathbf{x}^{(i)}), \text{ for } l = 1, \dots, N \\ \Rightarrow & \mathbf{K}^2 \boldsymbol{\alpha}^{(i)} = N \lambda^{(i)} \mathbf{K} \boldsymbol{\alpha}^{(i)} \\ \Rightarrow & \mathbf{K} \boldsymbol{\alpha}^{(i)} = N \lambda^{(i)} \boldsymbol{\alpha}^{(i)} = \tilde{\lambda}^{(i)} \boldsymbol{\alpha}^{(i)} \end{aligned}$$

- Maximizing $\lambda^{(1)}, \dots, \lambda^{(k)}$ (eigenvalues of \mathbf{S}) amounts to maximizing $\tilde{\lambda}^{(1)}, \dots, \tilde{\lambda}^{(k)}$ (eigenvalues of \mathbf{K})
- $\boldsymbol{\alpha}^{(i)}$'s are the eigenvectors corresponding the maximal eigenvalues of \mathbf{K}

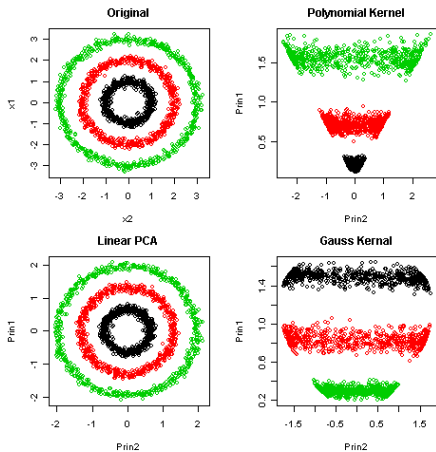
Solving $\alpha^{(1)}, \dots, \alpha^{(k)}$ (3/3)

- Note that $\mathbf{w}^{(i)}$ is normalized, i.e., $\|\mathbf{w}^{(i)}\| = 1$
 - Each $\alpha^{(i)}$ needs to be scaled properly
- $\|\mathbf{w}^{(i)}\|^2 = (\sum_{t=1}^N \alpha_t^{(i)} \Phi(\mathbf{x}^{(t)}))^\top (\sum_{t=1}^N \alpha_t^{(i)} \Phi(\mathbf{x}^{(t)})) = \alpha^{(i)\top} \mathbf{K} \alpha^{(i)} = 1 \Rightarrow \alpha^{(i)\top} \alpha^{(i)} = 1/N\lambda^{(i)} = 1/\tilde{\lambda}^{(i)}$
 - $\alpha^{(i)} \leftarrow \frac{\alpha^{(i)}}{\sqrt{\tilde{\lambda}^{(i)} \|\alpha^{(i)}\|}}$
- To project the $\mathbf{z}^{(t)} = \mathbf{w}^\top \Phi(\mathbf{x}^{(t)})$
- Similarly, $\alpha_2, \dots, \alpha_k$ are the (scaled) eigenvectors corresponding to the largest eigenvalues of \mathbf{K}

Centering $\Phi(\mathbf{x}^{(t)})$'s

- Note that by letting $\mathbf{S} = \frac{1}{N} \sum_{t=1}^N \Phi(\mathbf{x}^{(t)})\Phi(\mathbf{x}^{(t)})^\top$, we assume that $\Phi(\mathbf{x}^{(t)})$'s are centered
 - Recall that $\Sigma_{\Phi(\mathbf{x})} = E[\Phi(\mathbf{x})\Phi(\mathbf{x})^\top] - \boldsymbol{\mu}_{\Phi(\mathbf{x})}\boldsymbol{\mu}_{\Phi(\mathbf{x})}^\top$
 - In linear PCA, we can simply center $\Phi(\mathbf{x}^{(t)})$'s by a preprocessing step
- Given an arbitrary kernel function $k(\cdot)$, there is **no** guarantee that $\Phi(\mathbf{x}^{(t)})$'s will be centered in the lifted space
 - Preprocessing is infeasible
 - The model itself needs to be extended to accept uncentered instances in the lifted space. How? [Homework]

PCA vs. Kernel PCA



1 Why Dimensionality Reduction?

2 Feature Selection

- Forward and Backward Selection (Supervised)

3 Feature Extraction

- Principal Component Analysis (Unsupervised, Linear)
- Linear Discriminant Analysis (Supervised, Linear)
- Kernel PCA (Unsupervised, Nonlinear)
- Isometric Feature Mapping (Unsupervised, Nonlinear)
- Locally Linear Embedding (Unsupervised, Nonlinear)

Manifold-Preserved Nonlinear Mapping

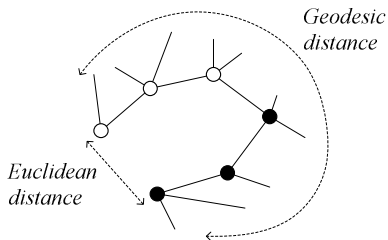
- Kernel PCA does not assumed a particular nonlinear mapping
 - Performance largely depends on the selection of kernel function
- In some applications, assuming some particular form of nonlinear mapping will be more helpful
- Consider an example where $\mathbf{x}^{(t)}$ are face photos of size 100×100 pixels
 - A series of one's photos taken from different angles forms a trajectory in the 10000-dimensional space
 - The collection of people's photos defines a *manifold* in the 10000-dimensional space
 - If we map the 10000-dimensional space linearly to a low dimensional space, the structure of the manifold may not be preserved

The Basic Idea

- Nonlinear methods rebuild the manifold in a low dimensional space
 - ① Treat the manifold as a space
 - ② Measure the relations/constraints between $\mathbf{x}^{(t)}$ in that space (which reflect the structure of the manifold)
 - ③ Find instances $\mathbf{z}^{(t)}$ in a low dimensional space that obey the constraints most (so preserve the structure)
- The final $\mathbf{z}^{(t)}$ may not relate to $\mathbf{x}^{(t)}$ linearly
- Two popular algorithms:
 - Isometric feature mapping (Isomap)
 - Locally linear embedding (LLE)
- It might be a good idea to review the topology now

Isomap

- The *Isometric feature mapping (Isomap)* measures the *geodesic distances* between examples in step 2; and then in step 3, find instances $z^{(t)}$ in a low dimensional space with mutual distances as close as the geodesic distances as possible
 - Geodesic distance is the distance along the manifold that the data lies in, as opposed to the Euclidean distance in the input space



Measuring the Geodesic Distances

- Recall that a manifold resembles a Euclidean space at a small enough scale
 - One topological property of the Euclidean space is that points are connected
- A point $\mathbf{x}^{(r)}$ and another $\mathbf{x}^{(s)}$ are directly connected if $\mathbf{x}^{(s)}$ lies in the same neighborhood with $\mathbf{x}^{(r)}$, and their geodesic distance can simply be the Euclidean distance
 - Any other choice?

Measuring the Geodesic Distances

- Recall that a manifold resembles a Euclidean space at a small enough scale
 - One topological property of the Euclidean space is that points are connected
- A point $\mathbf{x}^{(r)}$ and another $\mathbf{x}^{(s)}$ are directly connected if $\mathbf{x}^{(s)}$ lies in the same neighborhood with $\mathbf{x}^{(r)}$, and their geodesic distance can simply be the Euclidean distance
 - Any other choice? Mahalanobis distance, at the cost of computing the covariance matrix Σ
- How about points not in the same neighborhood?

Measuring the Geodesic Distances

- Recall that a manifold resembles a Euclidean space at a small enough scale
 - One topological property of the Euclidean space is that points are connected
- A point $\mathbf{x}^{(r)}$ and another $\mathbf{x}^{(s)}$ are directly connected if $\mathbf{x}^{(s)}$ lies in the same neighborhood with $\mathbf{x}^{(r)}$, and their geodesic distance can simply be the Euclidean distance
 - Any other choice? Mahalanobis distance, at the cost of computing the covariance matrix Σ
- How about points not in the same neighborhood?
 - In an atlas, neighborhoods are overlapped
 - The geodesic distance between two points that are not in the same neighborhood can be calculated by the length of their shortest path

Identifying the Neighborhood

- How identify to the neighboring points of $x^{(r)}$?

Identifying the Neighborhood

- How identify to the neighboring points of $\mathbf{x}^{(r)}$?
- We assume that $\mathbf{x}^{(r)}$ and $\mathbf{x}^{(s)}$ lie in the same neighborhood if:
 - $\|\mathbf{x}^{(r)} - \mathbf{x}^{(s)}\| < \varepsilon$; or
 - $\mathbf{x}^{(s)}$ is one of the n -nearest neighbors of $\mathbf{x}^{(r)}$
- User-specific parameters ε and n are usually small, but must be chosen carefully to make sure that the network is still connected

Finding $\{z^{(t)}\}_{t=1}^N$ (1/3)

- Let $d_{r,s}$ be the geodesic distance between $x^{(r)}$ and $x^{(s)}$, $1 \leq r, s \leq N$, now we want to find $z^{(r)}$ and $z^{(s)}$ in a k -dimensional space such that their mutual Euclidean distances are as close to $d_{r,s}$ as possible
- Exact solution (without error) may exist only when k is **larger** than d
 - There may be no way to “straighten” the geodesic distances between points in a space of dimension d (or lower)
- The solution $\{z^{(t)}\}_t$ is not unique, as we can shift all $z^{(t)}$ together to get another solution
 - To constrain the solution, we assume that $\sum_{t=1}^N z_i^{(t)} = 0$ for $i = 1, \dots, k$

Finding $\{z^{(t)}\}_{t=1}^N$ (2/3)

- Let $\mathbf{Z} = \begin{bmatrix} z^{(1)} \\ \vdots \\ z^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times k}$, we want to find the relationship between \mathbf{Z} and $d_{r,s}$
- For any $z^{(r)}$ and $z^{(s)}$, we have $d_{r,s}^2 = \|z^{(r)} - z^{(s)}\|^2 = b_{r,r} + b_{s,s} - 2b_{r,s}$, where $b_{r,s} = \sum_{i=1}^k z_i^{(r)} z_i^{(s)} = (z^{(r)})^\top z^{(s)}$
- Therefore, $\mathbf{B} = \begin{bmatrix} b_{1,1} & \cdots & b_{1,N} \\ \vdots & \ddots & \vdots \\ b_{N,1} & \cdots & b_{N,N} \end{bmatrix} = \mathbf{Z}\mathbf{Z}^\top$

Finding $\{z^{(t)}\}_{t=1}^N$ (3/3)

- On the other hand, let $T = \sum_{t=1}^N b_{t,t} = \sum_t \sum_i (z_i^{(t)})^2$, we get
 - $\sum_r d_{r,s}^2 = T + Nb_{s,s} - 2 \sum_r \sum_i z_i^{(r)} z_i^{(s)} =$
 $T + Nb_{s,s} - 2 \sum_i \left(\sum_r z_i^{(r)} \right) z_i^{(s)} = T + Nb_{s,s}$
 - $\sum_s d_{r,s}^2 = Nb_{r,r} + T$
 - $\sum_r \sum_s d_{r,s}^2 = 2NT$
- Each element in \mathbf{B} can be expressed by the geodesic distance by $b_{i,j} = \frac{1}{2} (b_{i,i} + b_{j,j} - d_{i,j}^2) = \frac{1}{2} \left(\frac{1}{N} \sum_s d_{i,s}^2 + \frac{1}{N} \sum_r d_{r,j}^2 - \frac{1}{N^2} \sum_r \sum_s d_{r,s}^2 - d_{i,j}^2 \right)$
- Note \mathbf{B} is symmetric and can be written as $\mathbf{B} = \mathbf{U}\mathbf{D}\mathbf{U}^\top$, where the columns of \mathbf{U} are eigenvectors
- We find $Z = \mathbf{U}\mathbf{D}^{1/2}$

Deciding k

- Given $\mathbf{Z} \in \mathbb{R}^{N \times k}$, the rank of $\mathbf{Z}\mathbf{Z}^T$ is at most k
- In the case that $\mathbf{B} \in \mathbb{R}^{N \times N}$ has full rank, we need $k = N$ to obtain the exact solution
- However, for dimension reduction, we want $k < d$ (and N)
- This leads to a **low rank approximation problem**: given a small k , $k < N$, we want to find a matrix $\tilde{\mathbf{B}}$ such that $\|\mathbf{B} - \tilde{\mathbf{B}}\|_F$ is minimized, subject to $\text{rank}(\tilde{\mathbf{B}}) = k$
 - Why Frobenius norm?

Deciding k

- Given $\mathbf{Z} \in \mathbb{R}^{N \times k}$, the rank of $\mathbf{Z}\mathbf{Z}^T$ is at most k
- In the case that $\mathbf{B} \in \mathbb{R}^{N \times N}$ has full rank, we need $k = N$ to obtain the exact solution
- However, for dimension reduction, we want $k < d$ (and N)
- This leads to a **low rank approximation problem**: given a small k , $k < N$, we want to find a matrix $\tilde{\mathbf{B}}$ such that $\|\mathbf{B} - \tilde{\mathbf{B}}\|_F$ is minimized, subject to $\text{rank}(\tilde{\mathbf{B}}) = k$
 - Why Frobenius norm? To minimize the error of the approximated distances
 - Recall that the best approximation is given by $\tilde{\mathbf{B}} = \tilde{\mathbf{U}}\tilde{\mathbf{D}}\tilde{\mathbf{U}}^T$, where $\tilde{\mathbf{D}} \in \mathbb{R}^{k \times k}$ contains only the k largest eigenvalues and $\tilde{\mathbf{U}} = [\mathbf{u}_1, \dots, \mathbf{u}_k] \in \mathbb{R}^{N \times k}$ contains the corresponding eigenvectors
- Finally, we let $\mathbf{Z} = \tilde{\mathbf{U}}\tilde{\mathbf{D}}^{1/2}$

Are We Satisfied? (1/2)

- Isomap is also a one-group process
 - This is less severe since the “structure” between examples is preserved
 - We can also let $d'_{r,s} = (1 - \alpha)d_{r,s} + \alpha c_{r,s}$, where $c_{r,s}$ is the distance between classes $\mathbf{x}^{(r)}$ and $\mathbf{x}^{(s)}$ belong to, and the parameter α can be tuned using the cross validation
- The major problem of Isomap is that it does not learn a mapping between $\mathbf{x}^{(t)}$ and $\mathbf{z}^{(t)}$
 - $Z = \tilde{U}\tilde{D}^{1/2}$ implies that $z_j^{(t)} = \lambda_j u_j^{(t)}$, where $u_j^{(t)}$ is the t th component of the eigenvector $\mathbf{u}_j \in \mathbb{R}^{N \times 1}$ in \tilde{U}
 - Given a new instance \mathbf{x}' , we need to rerun the whole algorithm using the $N+1$ instances to get \mathbf{z}'

Are We Satisfied? (2/2)

- Solution?

Are We Satisfied? (2/2)

- Solution?
 - Taking the advantage that a manifold is locally linear, we can identify examples $\mathbf{x}^{(s)}$ in the same neighborhood as \mathbf{x} and calculate \mathbf{w} such that $\|\mathbf{x} - \sum_s w_s \mathbf{x}^{(s)}\|^2$ is minimized
 - Calculate \mathbf{z} by interpolation: $\mathbf{z} = \sum_s w_s \mathbf{z}^{(s)}$
- The cost is that we need to store the whole set of $\{\mathbf{x}^{(t)}, \mathbf{z}^{(t)}\}_{t=1}^N$

1 Why Dimensionality Reduction?

2 Feature Selection

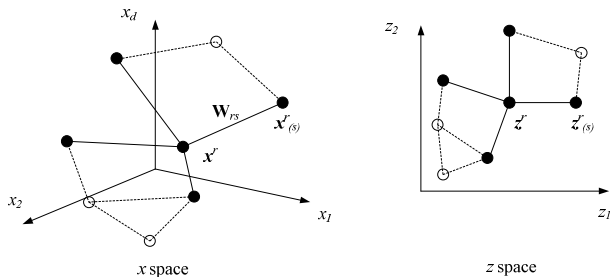
- Forward and Backward Selection (Supervised)

3 Feature Extraction

- Principal Component Analysis (Unsupervised, Linear)
- Linear Discriminant Analysis (Supervised, Linear)
- Kernel PCA (Unsupervised, Nonlinear)
- Isometric Feature Mapping (Unsupervised, Nonlinear)
- Locally Linear Embedding (Unsupervised, Nonlinear)

Locally Linear Embedding (1/2)

- Based on the similar idea above, the **Locally Linear Embedding (LLE)** represents each example as a linear combination of its nearby points in step 2; and then in step 3, find instances $z^{(t)}$ in a low dimensional space which preserve the combinations most



Locally Linear Embedding (2/2)

- Step 2: For each example $\mathbf{x}^{(r)}$ and all its nearby points $\mathbf{x}^{(s)}$, $s \neq r$, find weights $w_{r,s}$ such that $\|\mathbf{x}^{(r)} - \sum_s w_{r,s} \mathbf{x}^{(s)}\|^2$ is minimized, subject to $\sum_s w_{r,s} = 1$
 - The constraint ensures that after translating all the points together by some vector \mathbf{c} , the combination is still valid; i.e.,
$$\|(\mathbf{x}^{(r)} + \mathbf{c}) - \sum_s w_{r,s} (\mathbf{x}^{(s)} + \mathbf{c})\|^2 = \|\mathbf{x}^{(r)} - \sum_s w_{r,s} \mathbf{x}^{(s)}\|^2$$
- Step 3: Find $\{\mathbf{z}^{(t)}\}_{t=1}^N$ such that $\sum_r \|\mathbf{z}^{(r)} - \sum_s w_{r,s} \mathbf{z}^{(s)}\|^2$ is minimized, subject to $\frac{1}{N} \sum_t \mathbf{z}^{(t)} = \mathbf{0}$ ($E[\mathbf{z}] = \mathbf{0}$) and $\frac{1}{N-1} \sum_t (\mathbf{z}^{(t)} - \mathbf{0})(\mathbf{z}^{(t)} - \mathbf{0})^\top = \mathbf{I}$ ($Cov(\mathbf{z}) = \mathbf{I}$)
 - The first constraint is similar to that of Isomap and ensures an unique solution
 - The second guarantees that attributes of \mathbf{z} a) are uncorrelated; **b)** have the same (unit) variance (this is stronger)

- $w_{r,s}$ can be solved by considering one example $\mathbf{x}^{(r)}$ by one
- Let $\mathbf{w}^{(r)} = [w_{r,1}, \dots, w_{r,n}]^\top$, we can rewrite the objective in step 2 as
$$\|\mathbf{x}^{(r)} - \sum_s w_{r,s} \mathbf{x}^{(s)}\|^2 = \|\sum_s w_{r,s} (\mathbf{x}^{(r)} - \mathbf{x}^{(s)})\|^2 = \mathbf{w}^{(r)\top} \mathbf{G}^{(r)\top} \mathbf{G}^{(r)} \mathbf{w}^{(r)}$$
 - $\mathbf{G}^{(r)} = [\mathbf{x}^{(r)} - \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(r)} - \mathbf{x}^{(n)}] \in \mathbb{R}^{d \times n}$
 - Subject to $\mathbf{1}^\top \mathbf{w}^{(r)} = 1$
- Taking the partial derivatives of the Lagrangian $\mathbf{w}^{(r)\top} \mathbf{G}^{(r)\top} \mathbf{G}^{(r)} \mathbf{w}^{(r)} - \alpha(\mathbf{1}^\top \mathbf{w}^{(r)} - 1)$ with respect to $\mathbf{w}^{(r)}$ and α and setting them to zero, we have
$$\begin{cases} 2\mathbf{G}^{(r)\top} \mathbf{G}^{(r)} \mathbf{w}^{(r)} - \alpha \mathbf{1} = 0 \\ \mathbf{1}^\top \mathbf{w}^{(r)} = 1 \end{cases}$$
 - $\mathbf{w}^{(r)} = \frac{(\mathbf{G}^{(r)\top} \mathbf{G}^{(r)})^{-1} \mathbf{1}}{\mathbf{1}^\top (\mathbf{G}^{(r)\top} \mathbf{G}^{(r)})^{-1} \mathbf{1}}$ can be solve analytically [Proof]
 - Note for $\mathbf{G}^{(r)\top} \mathbf{G}^{(r)}$ to be invertible, we need make sure that $n \leq d$

Solving $\{z^{(t)}\}_{t=1}^N$ (1/3)

- Now we are given $\mathbf{W} \in \mathbb{R}^{N \times N}$, a matrix with nonzero elements corresponding to the weights $w_{r,s}$ found in step 2

- $\mathbf{W}\mathbf{1} = \mathbf{1}$

- Let $\mathbf{Z} = \begin{bmatrix} z^{(1)} \\ \vdots \\ z^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times k}$. Since the attributes of z are uncorrelated,

the columns $\mathbf{c}_1, \dots, \mathbf{c}_k \in \mathbb{R}^{N \times 1}$ of \mathbf{Z} are orthogonal to each other

- To get \mathbf{c}_1 , we consider only the first attribute of $z^{(t)}$ and rewrite the objective in step 3 as

$$\begin{aligned} \sum_r \left\| z_1^{(r)} - \sum_s w_{r,s} z_1^{(s)} \right\|^2 &= \sum_r z_1^{(r)2} - \sum_r z_1^{(r)} \left(\sum_s w_{r,s} z_1^{(s)} \right) - \\ &\sum_r \left(\sum_s w_{r,s} z_1^{(s)} \right) z_1^{(r)} + \sum_r \left(\sum_s w_{r,s} z_1^{(s)} \right)^2 = \\ &\mathbf{c}_1^\top \mathbf{c}_1 - \mathbf{c}_1^\top (\mathbf{W}\mathbf{c}_1) - (\mathbf{W}\mathbf{c}_1)^\top \mathbf{c}_1 + (\mathbf{W}\mathbf{c}_1)^\top (\mathbf{W}\mathbf{c}_1) = \\ &((\mathbf{I} - \mathbf{W})\mathbf{c}_1)^\top ((\mathbf{I} - \mathbf{W})\mathbf{c}_1) = \mathbf{c}_1^\top \mathbf{M}\mathbf{c}_1, \text{ where } \mathbf{M} = (\mathbf{I} - \mathbf{W})^\top (\mathbf{I} - \mathbf{W}) \end{aligned}$$

- Subject to $\frac{1}{N-1} \mathbf{c}_1^\top \mathbf{c}_1 = 1$

Solving $\{z^{(t)}\}_{t=1}^N$ (2/3)

- Taking the partial derivatives of the Lagrangian $\mathbf{c}_1^\top \mathbf{M} \mathbf{c}_1 - \alpha(\mathbf{c}_1^\top \mathbf{c}_1 - N + 1)$ with respect to \mathbf{c}_1 and α and setting them to zero, we have $\begin{cases} 2\mathbf{M}\mathbf{c}_1 - 2\alpha\mathbf{c}_1 = \mathbf{0} \\ \mathbf{c}_1^\top \mathbf{c}_1 - N + 1 = 0 \end{cases}$, implying $\mathbf{M}\mathbf{c}_1 = \alpha\mathbf{c}_1$ and $\mathbf{c}_1^\top \mathbf{M} \mathbf{c}_1 = (N-1)\alpha$
 - \mathbf{c}_1 is the eigenvector of \mathbf{M} corresponding to the smallest eigenvalue
- Similarly, \mathbf{c}_j which is orthogonal to $\mathbf{c}_1, \dots, \mathbf{c}_{j-1}$ is the eigenvector corresponding to the j th smallest eigenvalue
 - \mathbf{M} is symmetric and has orthogonal eigenvectors

Solving $\{z^{(t)}\}_{t=1}^N$ (3/3)

- Note we did not enforce the constraint $\frac{1}{N} \sum_t z^{(t)} = \frac{1}{N} \mathbf{Z}^\top \mathbf{1} = \mathbf{0}$ (or $\frac{1}{N} \mathbf{c}_j^\top \mathbf{1} = 0$ for all $1 \leq j \leq k$)
- Notice that the first eigenvector is always $\mathbf{1}$
 - Since $\sum_s w_{r,s} = 1$, we have
$$\mathbf{M}\mathbf{1} = (\mathbf{I} - \mathbf{W})^\top (\mathbf{I} - \mathbf{W})\mathbf{1} = (\mathbf{I} - \mathbf{W})^\top (\mathbf{1} - \mathbf{1}) = \mathbf{0}$$
 - \mathbf{M} is positive semidefinite and 0 must be the smallest eigenvalue [Proof]
- To be orthogonal to $\mathbf{1}$, all other eigenvectors must have components summed to 0, by virtue of orthogonality
- We can simply discard $\mathbf{1}$ and let \mathbf{c}_j correspond to the $(j+1)$ th eigenvector to enforce $\mathbf{c}_j^\top \mathbf{1} = 0$
- Finally, $z_j^{(t)}$ equals to the t th component of the eigenvector $\mathbf{c}_j \in \mathbb{R}^{N \times 1}$ of \mathbf{M}

Summary of Nonlinear Methods

- LLE has a similar problem with Isomap in that there is no mapping between $\mathbf{x}^{(t)}$ and $\mathbf{z}^{(t)}$
 - The interpolation technique applies to LLE
- It can be shown that LLE is equivalent to kernel PCA with the “LLE kernel”
- Both Isomap and LLE reconstruct the manifold in a k -dimensional space by patching the overlapping neighborhoods
 - In Isomap, geodesic distances are calculated hop by hop and preserved in the low dimensional space
 - In LLE, weights of combination are preserved hop by hop in the low dimensional space