CSCI 5521: Introduction to Machine Learning

Dimension Reduction (Chpt 6)

Rui Kuang

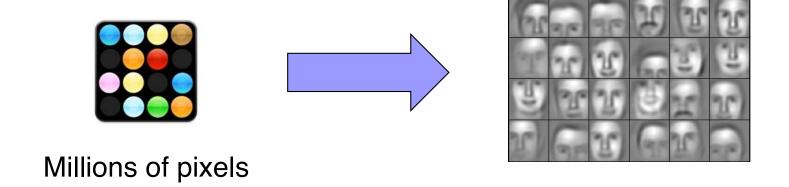
Department of Computer Science and Engineering
University of Minnesota



be.

Dimensionality Reduction: Face Recognition

 Learning a compact representation of images with millions of pixels.



Tens of eigen-faces

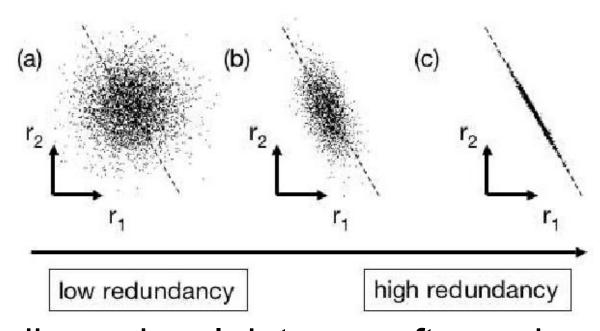


Why Reduce Dimensionality?

- Reduces time complexity: Less computation
- Reduces space complexity: Less parameters
- Saves the cost of observing the features
- Simpler models are more robust on small datasets
- More interpretable; simpler explanation
- Data visualization (structure, groups, outliers, etc) if plotted in 2 or 3 dimensions

100

Principal Components Analysis (PCA)

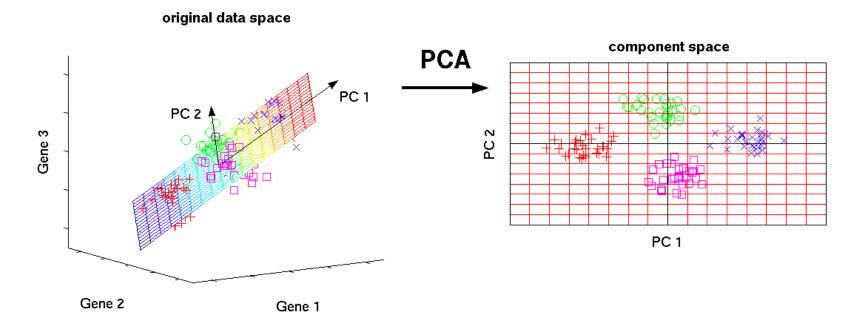


- High dimensional data are often noisy and redundant.
- We want to identify the key directions to have a compressed representation that are noise free and can be visualized
- One way is to find the principle components.



Principal Components Analysis (PCA)

- Find a low-dimensional space such that when x is projected there, information loss is minimized.
- The projection of x on the direction of w is: $z = w^T x$



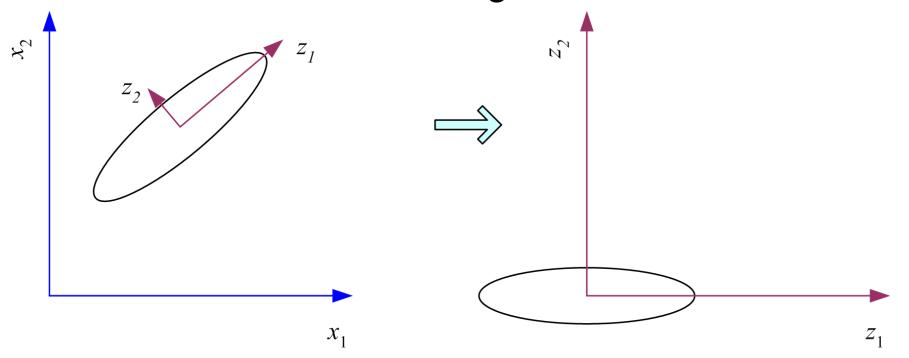


What PCA does

$$z = \mathbf{W}^{\mathsf{T}}(\mathbf{x} - \mathbf{m})$$

where the columns of **W** are the eigenvectors of \sum , and m is sample mean

Centers the data at the origin and rotates the axes



MA.

Principal Components Analysis (PCA)

Find w such that Var(z) is maximized

Var(z) = Var(
$$w^{T}x$$
) = E[($w^{T}x - w^{T}\mu$)²]
= E[($w^{T}x - w^{T}\mu$)($w^{T}x - w^{T}\mu$)]
= E[$w^{T}(x - \mu)(x - \mu)^{T}w$]
= w^{T} E[($x - \mu$)($x - \mu$)^T] $w = w^{T}$ $\sum w$

where
$$Var(\mathbf{x}) = E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] = \sum_{i=1}^{n} \mathbf{x}^T \mathbf{x}^T$$



■ Maximize Var(z) subject to ||w||=1

$$\max_{\mathbf{w}_1} \mathbf{w}_1^T \mathbf{\Sigma} \mathbf{w}_1 - \alpha (\mathbf{w}_1^T \mathbf{w}_1 - 1)$$

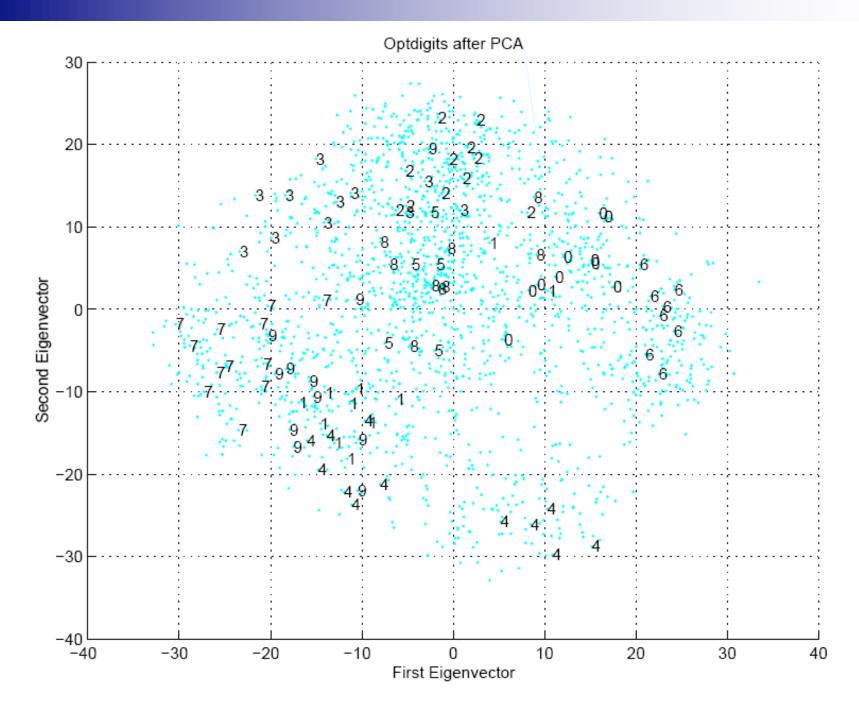
 $\sum \mathbf{w}_1 = \alpha \mathbf{w}_1$ that is, \mathbf{w}_1 is an eigenvector of \sum Choose the one with the largest eigenvalue for Var(z) to be max

■ Second principal component: Max $Var(z_2)$, s.t., $||\mathbf{w}_2||=1$ and orthogonal to \mathbf{w}_1

$$\max_{\mathbf{w}_2} \mathbf{w}_2^T \mathbf{\Sigma} \mathbf{w}_2 - \alpha (\mathbf{w}_2^T \mathbf{w}_2 - 1) - \beta (\mathbf{w}_2^T \mathbf{w}_1 - 0)$$

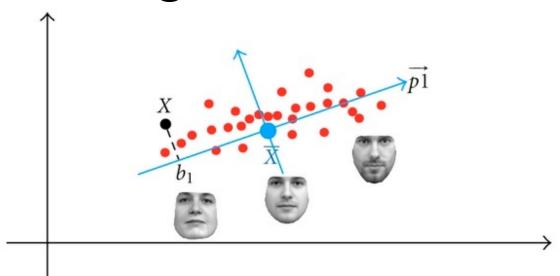
 $\sum \mathbf{w}_2 = \alpha \mathbf{w}_2$ that is, \mathbf{w}_2 is another eigenvector of \sum

and so on.

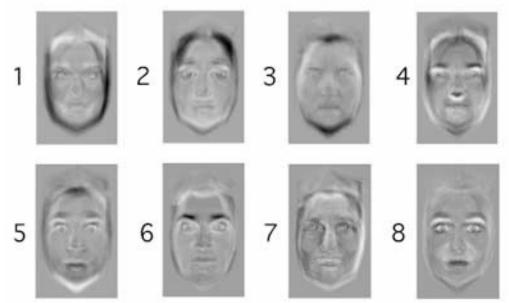


E. Alpaydin, Introduction to Machine Learning

Face Recognition



http://www.hindawi.com/journals/aans/2011/673016/



http://mathdesc.fr/documents/facerecog/PerceptionFacialExpression.htm



How to choose k?

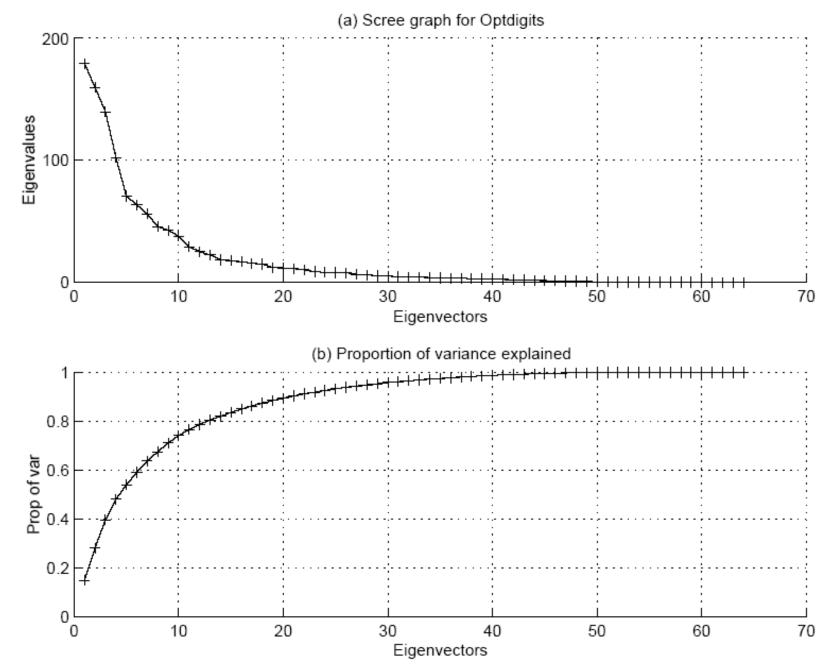
Proportion of Variance (PoV) explained

$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_k + \dots + \lambda_d}$$

when λ_i are sorted in descending order

- Typically, stop at PoV>0.9
- Scree graph plots of PoV vs k, stop at "elbow"





E. Alpaydin, Introduction to Machine Learning



PCA Discussions

- Linearity: Linear rotation of the original space
- Gaussian assumption: mean and variance are enough to characterize the noise and redundancy.
- We choose the principal components to be orthogonal, but they don't have to be.
- Often used with clustering algorithms to cluster high-dimensional data

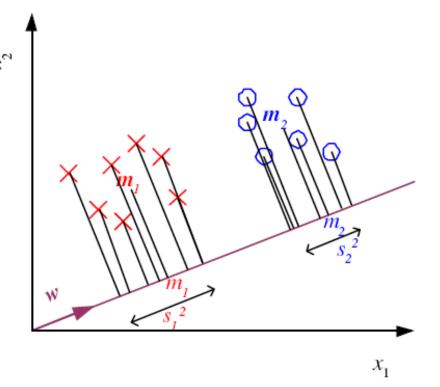


Linear Discriminant Analysis

- Find a low-dimensional space such that when **x** is projected, classes are well-separated.
- Find w that maximizes

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

$$m_1 = \frac{\sum_{t} \mathbf{w}^T \mathbf{x}^t r^t}{\sum_{t} r^t} \qquad s_1^2 = \sum_{t} (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t$$



Between-class scatter:

$$(m_1 - m_2)^2 = (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2$$

$$= \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}$$

$$= \mathbf{w}^T \mathbf{S}_B \mathbf{w} \text{ where } \mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T$$

Within-class scatter:

$$s_1^2 = \sum_t (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t$$

$$= \sum_t \mathbf{w}^T (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^T \mathbf{w} r^t = \mathbf{w}^T \mathbf{S}_1 \mathbf{w}$$
where $\mathbf{S}_1 = \sum_t (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^T r^t$

$$s_1^2 + s_2^2 = \mathbf{w}^T \mathbf{S}_W \mathbf{w} \text{ where } \mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$$



Fisher's Linear Discriminant

Find w that max

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} = \frac{\left| \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) \right|^2}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} \qquad \left(\frac{u}{v} \right)' = \frac{u'v - uv'}{v^2}$$

$$\frac{\mathbf{w}^{T}(\mathbf{m}_{1} - \mathbf{m}_{2})}{\mathbf{w}^{T}\mathbf{S}_{W}\mathbf{w}}(2(\mathbf{m}_{1} - \mathbf{m}_{2}) - \frac{\mathbf{w}^{T}(\mathbf{m}_{1} - \mathbf{m}_{2})}{\mathbf{w}^{T}\mathbf{S}_{W}\mathbf{w}}2S_{W}w) = 0$$

LDA soln:

$$\mathbf{w} = c \cdot \mathbf{S}_W^{-1} \big(\mathbf{m}_1 - \mathbf{m}_2 \big)$$



K>2 Classes

Within-class scatter:

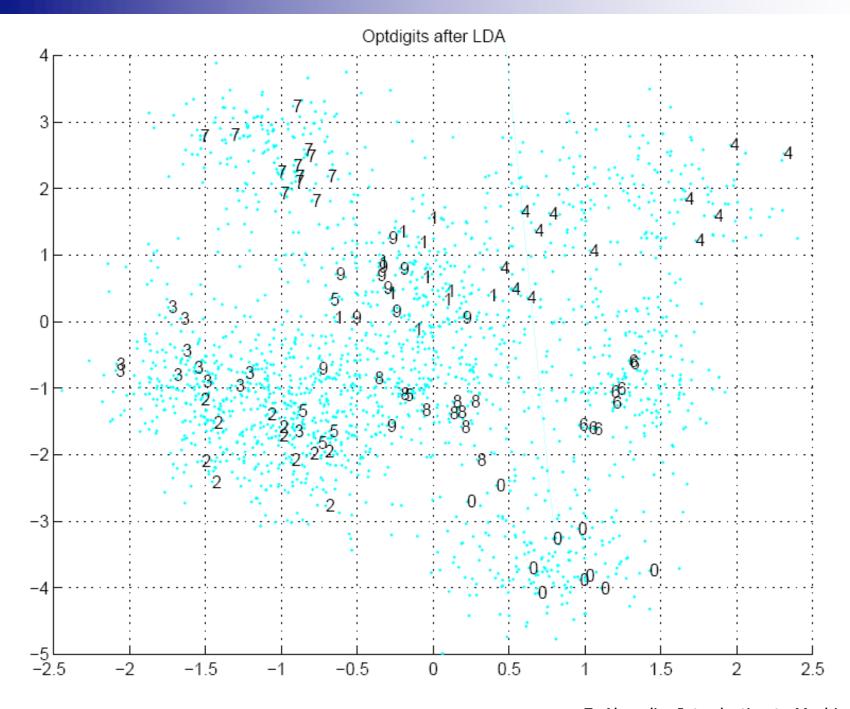
$$\mathbf{S}_{W} = \sum_{i=1}^{K} \mathbf{S}_{i} \qquad \mathbf{S}_{i} = \sum_{t} r_{i}^{t} (\mathbf{x}^{t} - \mathbf{m}_{i}) (\mathbf{x}^{t} - \mathbf{m}_{i})^{T}$$

Between-class scatter (among means):

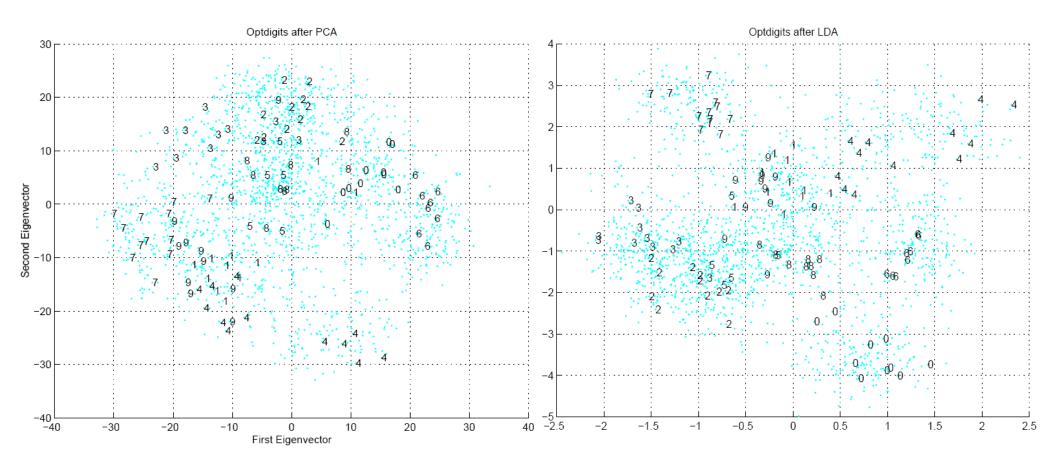
$$\mathbf{S}_{B} = \sum_{i=1}^{K} N_{i} (\mathbf{m}_{i} - \mathbf{m}) (\mathbf{m}_{i} - \mathbf{m})^{T} \qquad \mathbf{m} = \frac{1}{K} \sum_{i=1}^{K} \mathbf{m}_{i}$$

Find W that max

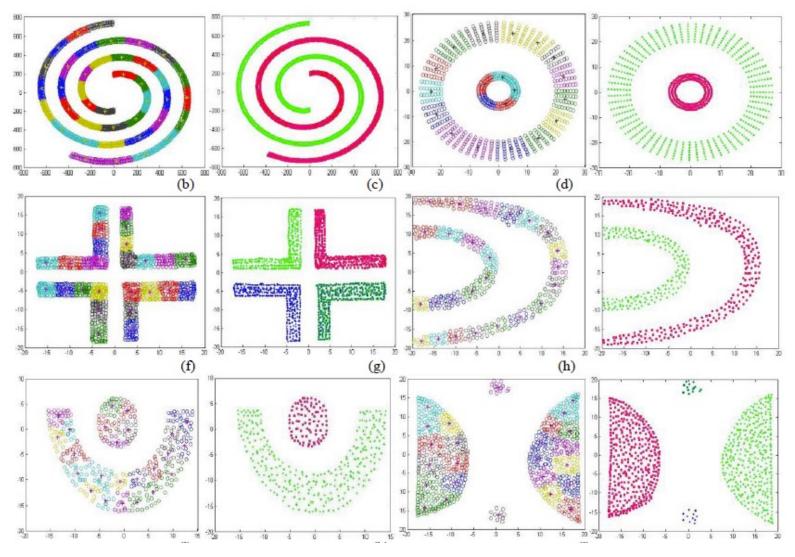
$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|}$$
 The largest eigenvectors of $\mathbf{S}_W^{-1} \mathbf{S}_B$ Maximum rank of K -1



E. Alpaydin, Introduction to Machine Learning



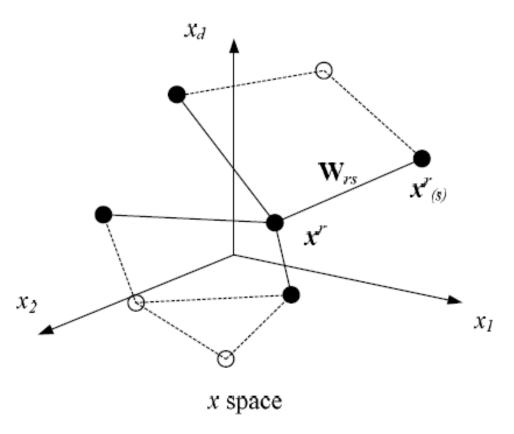


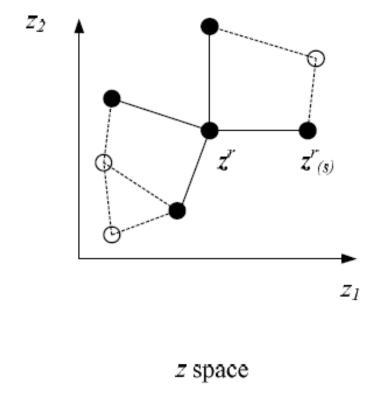




Locally Linear Embedding

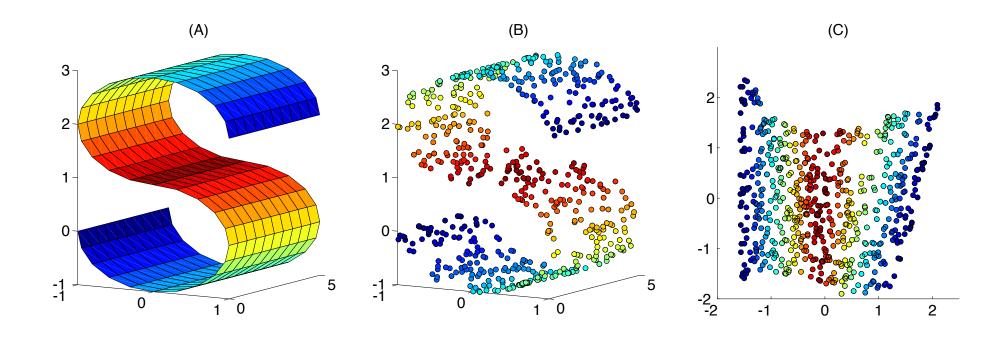
- Mapping to a new space allowing linear embedding.
- Use overlapping hyperplanes to approximate the nonlinear surface.







LLE on S Manifold



NA.

Locally Linear Embedding

- 1. Given \mathbf{x}^r find its neighbors $\mathbf{x}^s_{(r)}$
- 2. Find W_{rs} that minimize

$$E(\mathbf{W} \mid X) = \sum_{r} \left\| \mathbf{x}^{r} - \sum_{s} \mathbf{W}_{rs} \mathbf{x}_{(r)}^{s} \right\|^{2}, \text{subj: } \sum_{s} W_{rs} = 1$$

3. Find the new coordinates z^r that minimize

$$E(\mathbf{z} \mid \mathbf{W}) = \sum_{r} \left\| z^r - \sum_{s} \mathbf{W}_{rs} z_{(r)}^s \right\|^2, \text{subj: } \sum_{r} z^r = 0 \text{ and } \sum_{r} (z^r)^T z^r = 1$$

re.

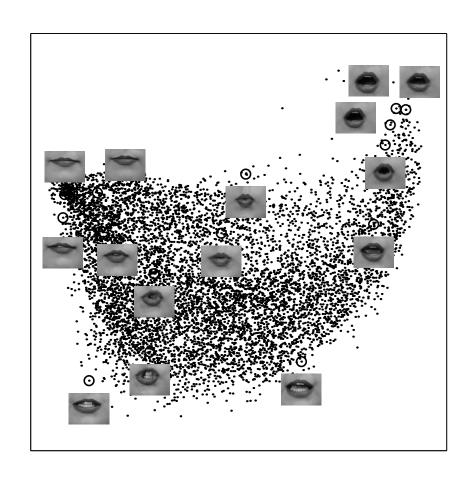
Locally Linear Embedding

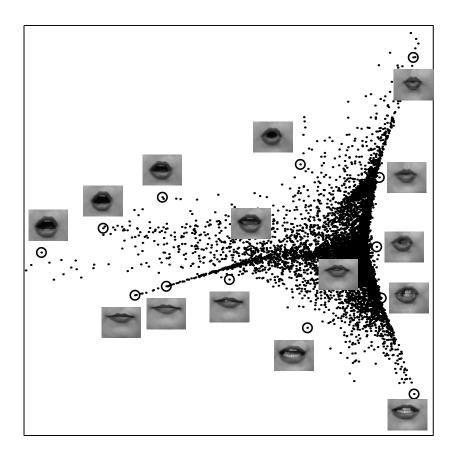
$$E(\mathbf{z} \mid \mathbf{W}) = \sum_{r} \left\| z^{r} - \sum_{s} \mathbf{W}_{rs} z_{(r)}^{s} \right\|^{2}, \text{subj: } \sum_{r} z^{r} = 0 \text{ and } \sum_{r} (z^{r})^{T} z^{r} = 1$$
$$= \sum_{r,s} (\delta_{rs} - W_{rs} - W_{sr} + \sum_{i} W_{is} W_{ir}) (z^{r})^{T} z^{s}$$

- With n neighbors and d dimensions, d≤n-1. d is often a bit smaller.
- Take k+1 lowest eigenvectors and discard the first one.



LLE on Lip Images

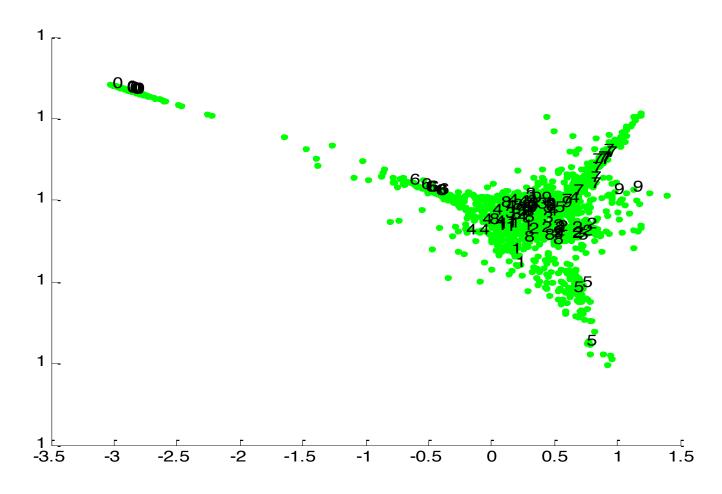




PCA LLE

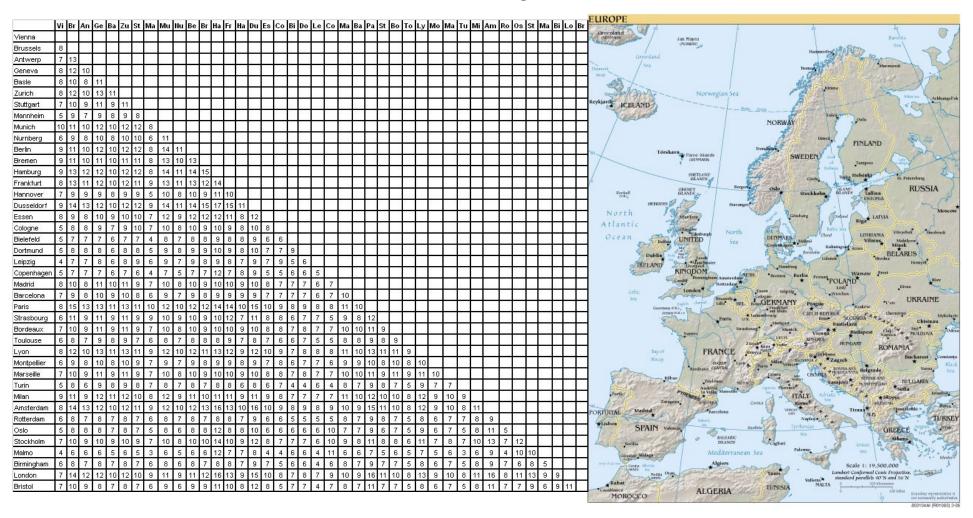


LLE on Optdigits



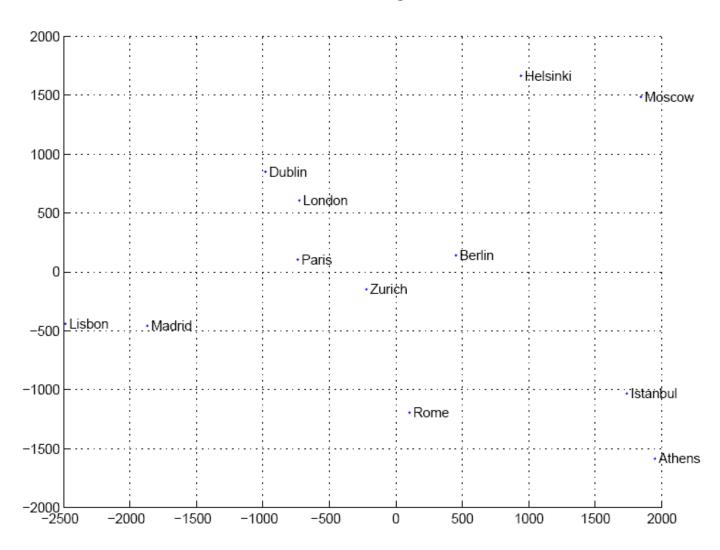
Matlab source from http://www.cs.toronto.edu/~roweis/lle/code.html

Construct Map By Distances





Map of Europe by MDS



M

Multidimensional Scaling

Given pairwise distances between N points,

$$d_{ij}, i,j = 1,...,N$$

place on a low-dim map s.t. distances are preserved.

 $z = g(x \mid \theta)$ Find θ that min Sammon stress

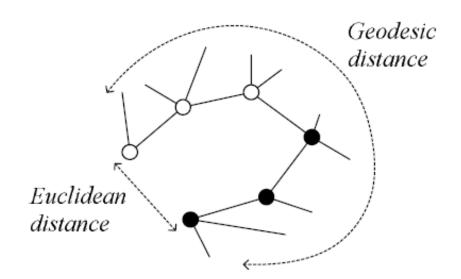
$$E(\theta \mid \mathcal{X}) = \sum_{r,s} \frac{\left(\left\|\mathbf{z}^r - \mathbf{z}^s\right\| - \left\|\mathbf{x}^r - \mathbf{x}^s\right\|\right)^2}{\left\|\mathbf{x}^r - \mathbf{x}^s\right\|^2}$$

$$= \sum_{r,s} \frac{\left(\left\|\mathbf{g}(\mathbf{x}^r \mid \theta) - \mathbf{g}(\mathbf{x}^s \mid \theta)\right\| - \left\|\mathbf{x}^r - \mathbf{x}^s\right\|^2}{\left\|\mathbf{x}^r - \mathbf{x}^s\right\|^2}$$



Isomap

 Geodesic distance is the distance along the manifold that the data lies in, as opposed to the Euclidean distance in the input space

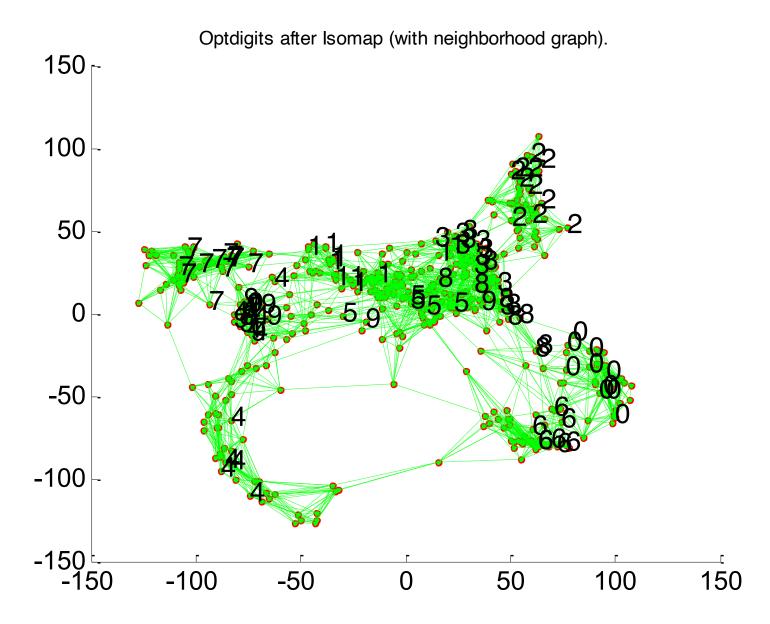




Isomap

- Instances r and s are connected in the graph if $||x^r-x^s|| < \varepsilon$ or if x^s is one of the k neighbors of x^r . The edge length is $||x^r-x^s||$
- For two nodes r and s not connected, the distance is equal to the shortest path between them (force faraway data points to go through their neighbors - unfolding).
- Once the NxN distance matrix is thus formed, use MDS to find a lower-dimensional mapping
- No general mapping function; rerun for each new point





Matlab source from http://web.mit.edu/cocosci/isomap/isomap.html



Summary

- Linear projection:
 - □ **PCA:** Project **x** to **z** to maximize the variance
 - LDA: Projection maximize within-class scatter and minimize between class scatter
- Non-linear embedding
 - □ Isomap: Use geodesic distance along the manifold instead of Euclidean distance
 - □ LLE: Linear patch assumption (linear combinations of neighbors)



Subset Selection vs Extraction

- Feature selection (subset selection):
 Choosing k<d important features, ignoring the remaining d k
 Subset selection algorithms
- Feature extraction: Project the original x_i, i =1,...,d dimensions to new k<d dimensions, z_j, j =1,...,k
- Principal components analysis (PCA), linear discriminant analysis (LDA), factor analysis (FA)



Subset Selection

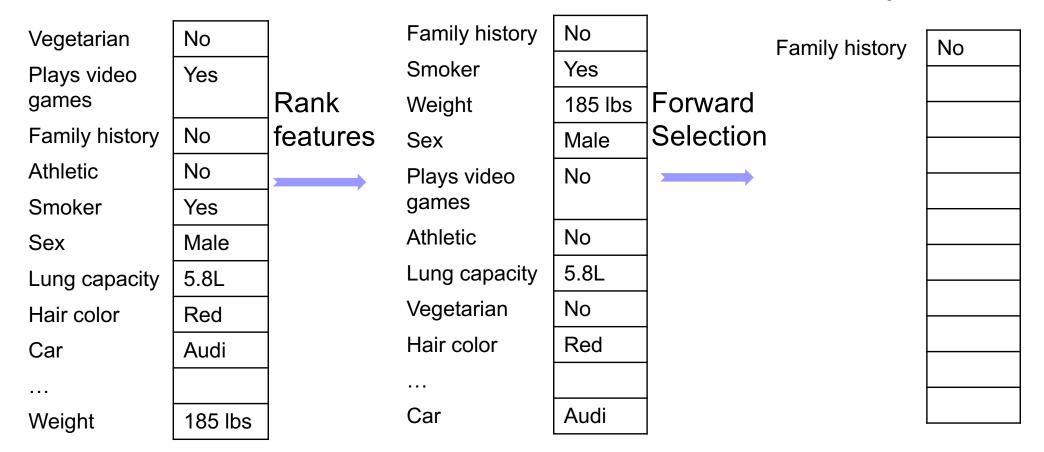
- There are 2^d subsets of d features
- Forward search: Add the best feature at each step
 - \square Set of features F initially \emptyset .
 - □ At each iteration, find the best new feature $j = \operatorname{argmin}_i E(F \cup x_i)$
 - \square Add x_j to F if $E(F \cup x_j) < E(F)$
- Backward search: Start with all features and remove one at a time, if possible.
- Hill-climbing $O(d^2)$ algorithm
- Floating search (Add k, remove l)



Forward Selection

- 1. Add the highest ranked feature
- 2. Check classification performance

Classification Accuracy: 75%





Forward Selection

- 1. Add the highest ranked feature
- 2. Check classification performance
- 3. Add the next highest ranked feature

Classification Accuracy: 75%→95%

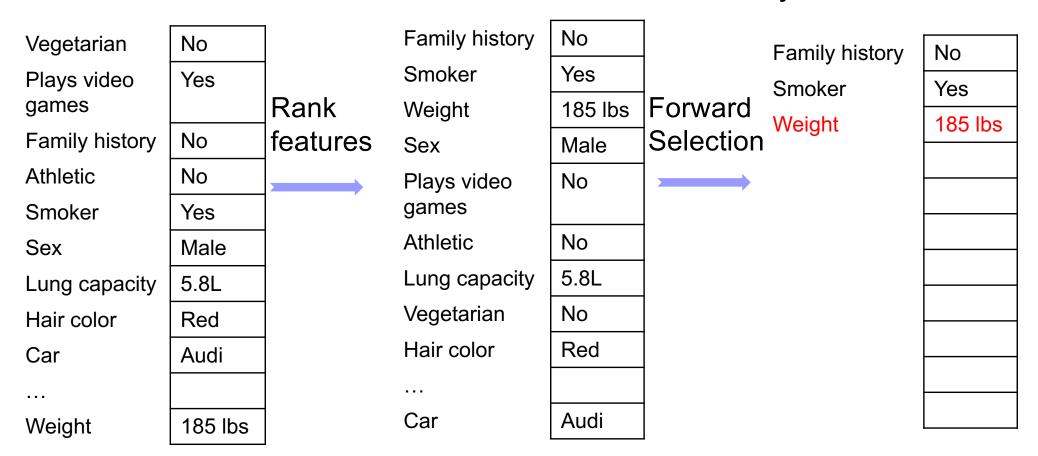
\/a = ata = a	Nia	1	Family history	No]		
Vegetarian	No		,		1	Family history	No
Plays video	Yes		Smoker	Yes		Smoker	Yes
games		Rank	Weight	185 lbs	Forward	Officker	163
Family history	No	features	Sex	Male	Selection		
Athletic	No		Plays video	No	——		
Smoker	Yes	,	games				
Sex	Male		Athletic	No			
Lung capacity	5.8L		Lung capacity	5.8L			
Hair color	Red		Vegetarian	No			
Car	Audi		Hair color	Red			
Weight	185 lbs	-	Car	Audi			



Forward Selection

- 1. Add the highest ranked feature
- 2. Check classification performance
- 3. Add the next highest ranked feature

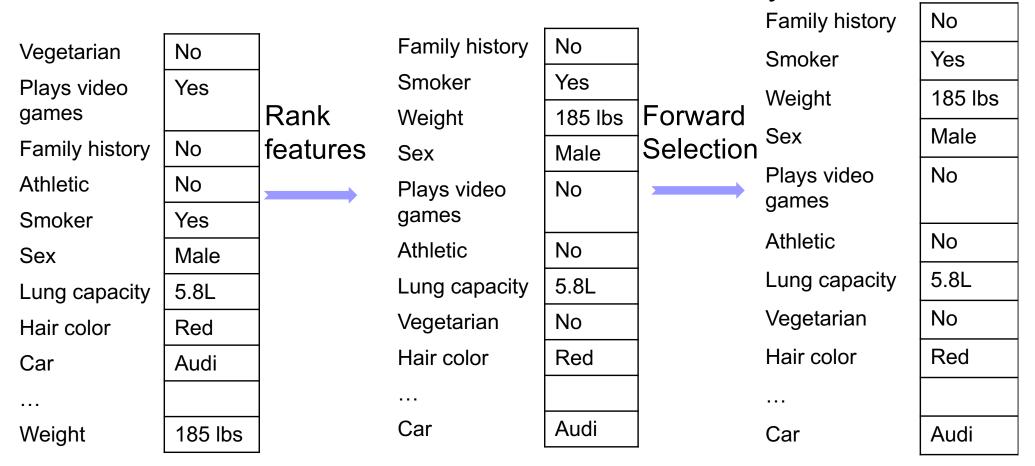
Classification Accuracy: 95%→80%





- 1. Remove the lowest ranked feature
- 2. Check classification performance

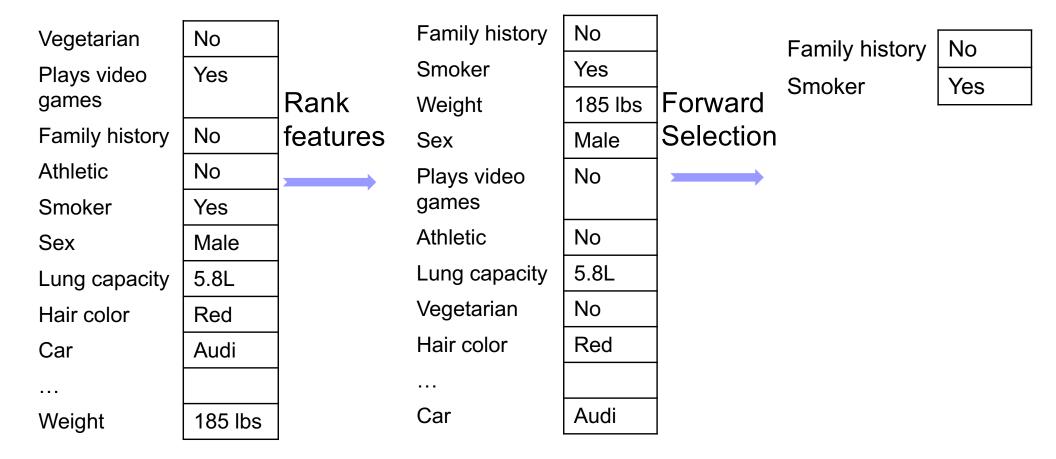
Classification Accuracy: 60%→75%





- 1. Remove the lowest ranked feature
- 2. Check classification performance
- 3. Remove the next lowest ranked feature until performance worse

Classification Accuracy: 95%





Feature Selection

- NP-hard to search through all the combinations
 - Need heuristic solutions
- The assumption is based on the maximum classification performance.
 - There might be more than one subset of features that can give the optimal classification performance.
- Easy to ignore the relation among the features
 - Combinations of several non-informative features might be meaningful