

Dimension Reduction (Chpt 6)

Rui Kuang

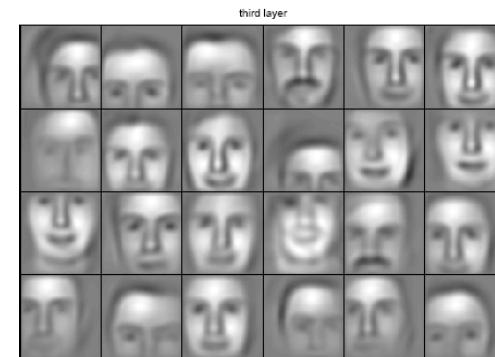
Department of Computer Science and Engineering
University of Minnesota

Dimensionality Reduction: Face Recognition

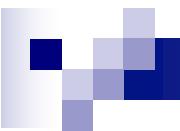
- Learning a compact representation of images with millions of pixels.



Millions of pixels



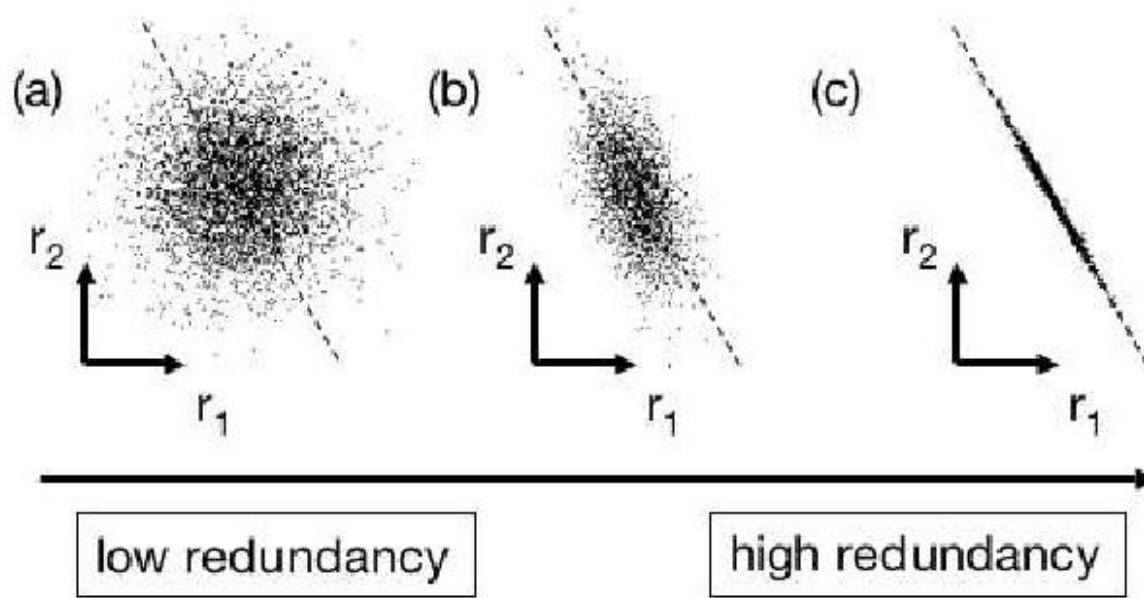
Tens of eigen-faces



Why Reduce Dimensionality?

- Reduces time complexity: Less computation
- Reduces space complexity: Less parameters
- Saves the cost of observing the feature
- 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

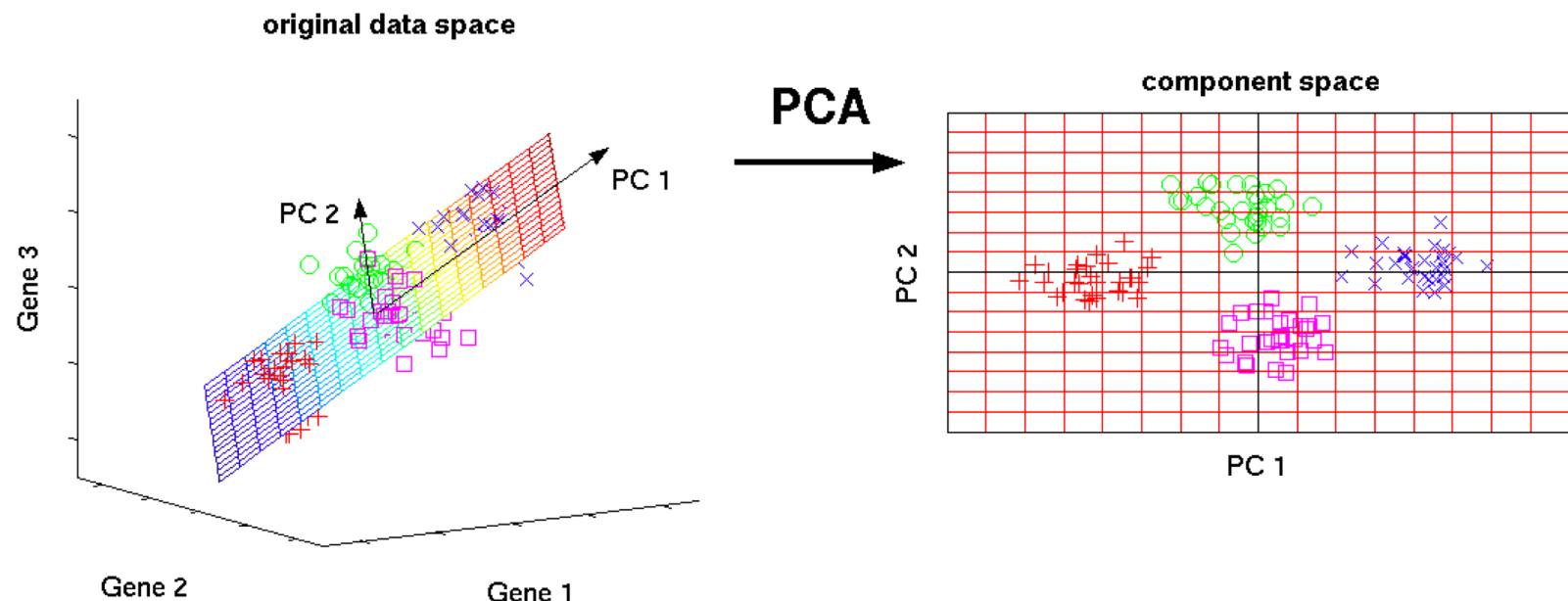
Principal Components Analysis (PCA)



- 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 \mathbf{x} is projected there, information loss is minimized.
- The projection of \mathbf{x} on the direction of \mathbf{w} is: $z = \mathbf{w}^T \mathbf{x}$

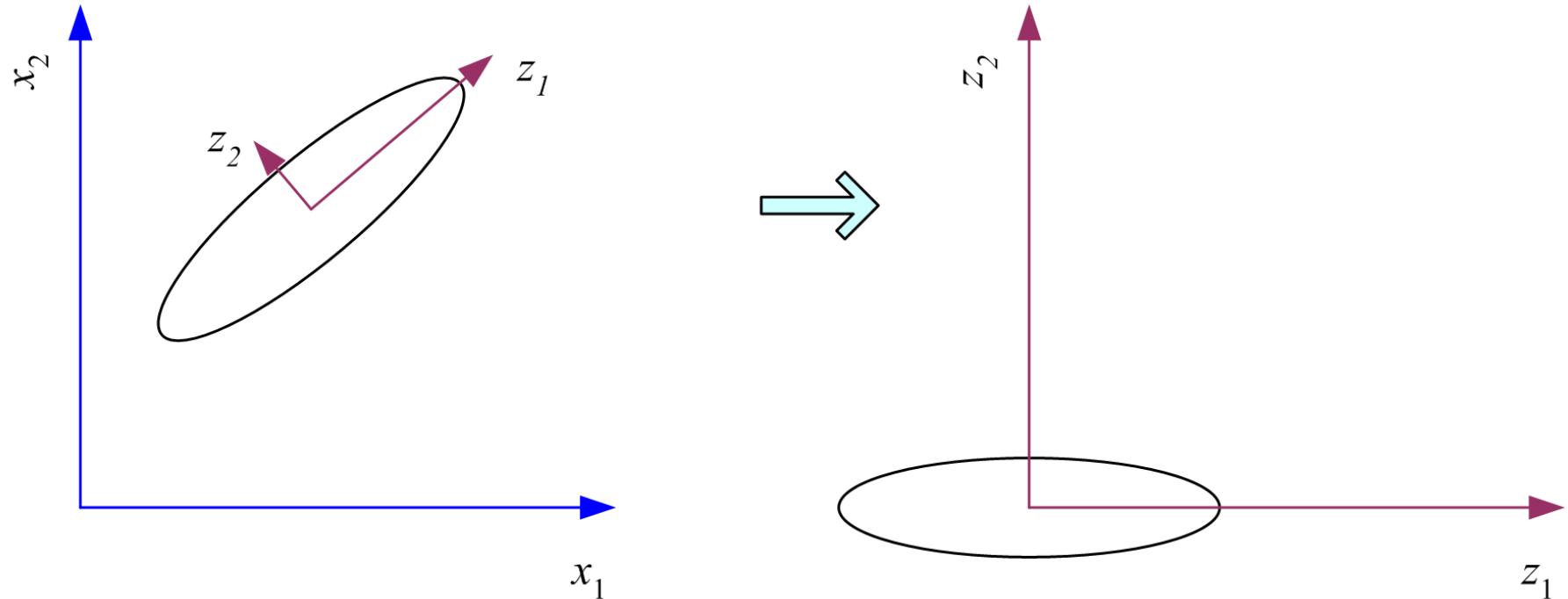


What PCA does

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

where the columns of \mathbf{W} are the eigenvectors of Σ , and \mathbf{m} is sample mean

Centers the data at the origin and rotates the axes



Principal Components Analysis (PCA)

- Find w such that $\text{Var}(z)$ is maximized

$$\begin{aligned}\text{Var}(z) &= \text{Var}(w^T x) = E[(w^T x - w^T \mu)^2] \\ &= E[(w^T x - w^T \mu)(w^T x - w^T \mu)] \\ &= E[w^T(x - \mu)(x - \mu)^T w] \\ &= w^T E[(x - \mu)(x - \mu)^T] w = w^T \Sigma w\end{aligned}$$

where $\text{Var}(x) = E[(x - \mu)(x - \mu)^T] = \Sigma$

- Maximize $\text{Var}(z)$ subject to $\|\mathbf{w}\|=1$

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

$\Sigma \mathbf{w}_1 = \alpha \mathbf{w}_1$ that is, \mathbf{w}_1 is an eigenvector of Σ

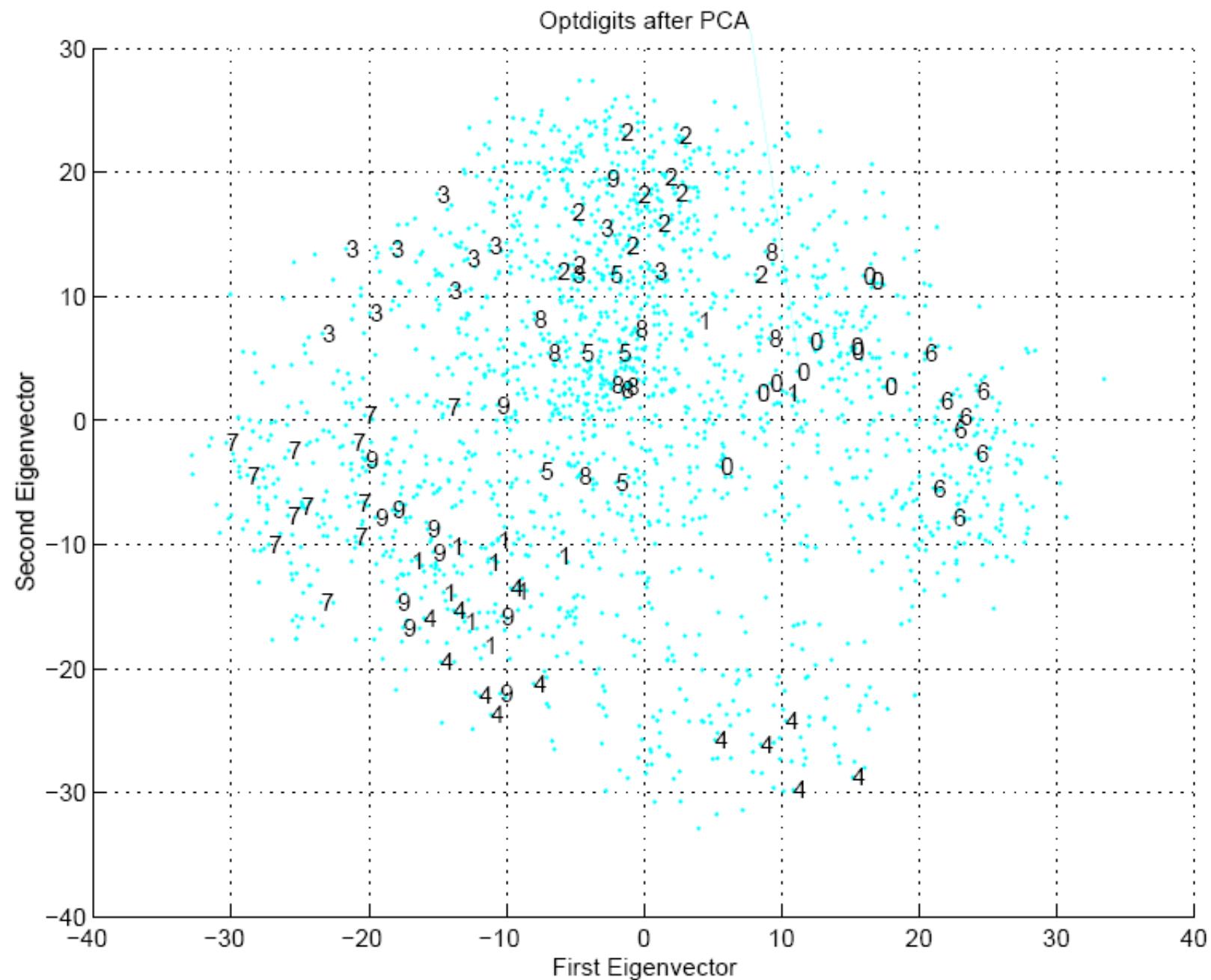
Choose the one with the largest eigenvalue for $\text{Var}(z)$ to be max

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

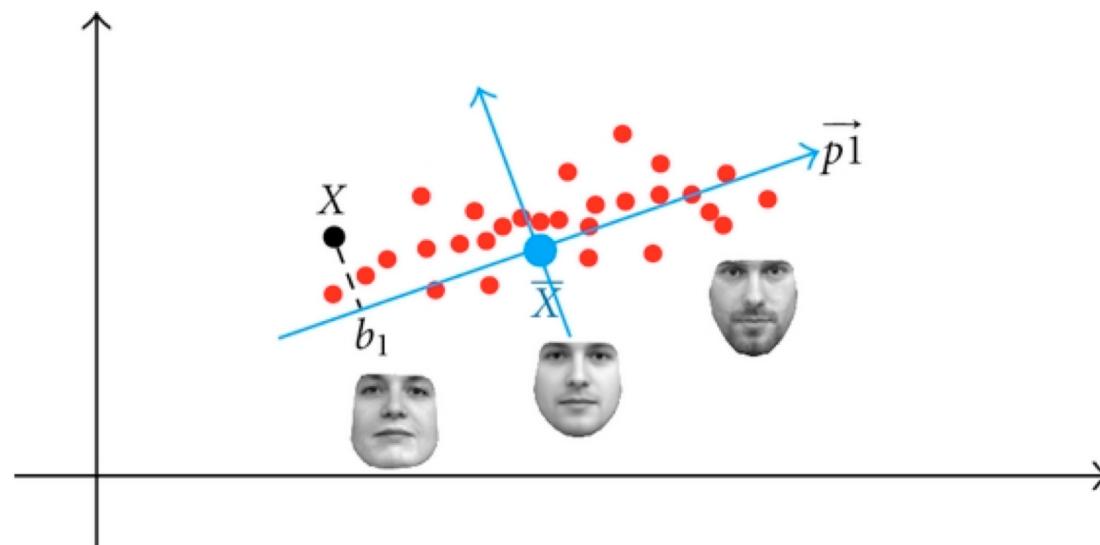
$$\max_{\mathbf{w}_2} \mathbf{w}_2^T \Sigma \mathbf{w}_2 - \alpha(\mathbf{w}_2^T \mathbf{w}_2 - 1) - \beta(\mathbf{w}_2^T \mathbf{w}_1 - 0)$$

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

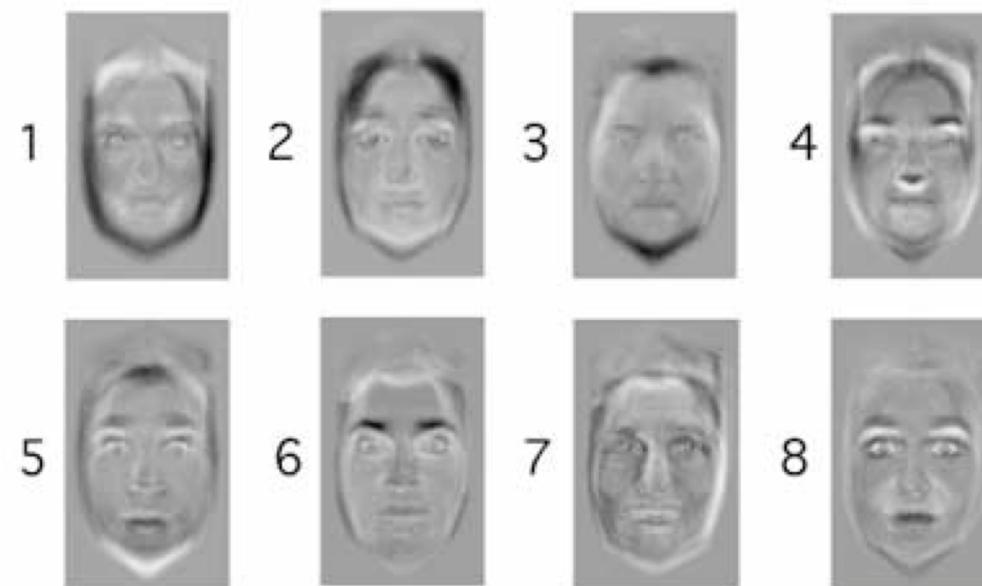
and so on.



Face Recognition



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



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

How to choose k ?

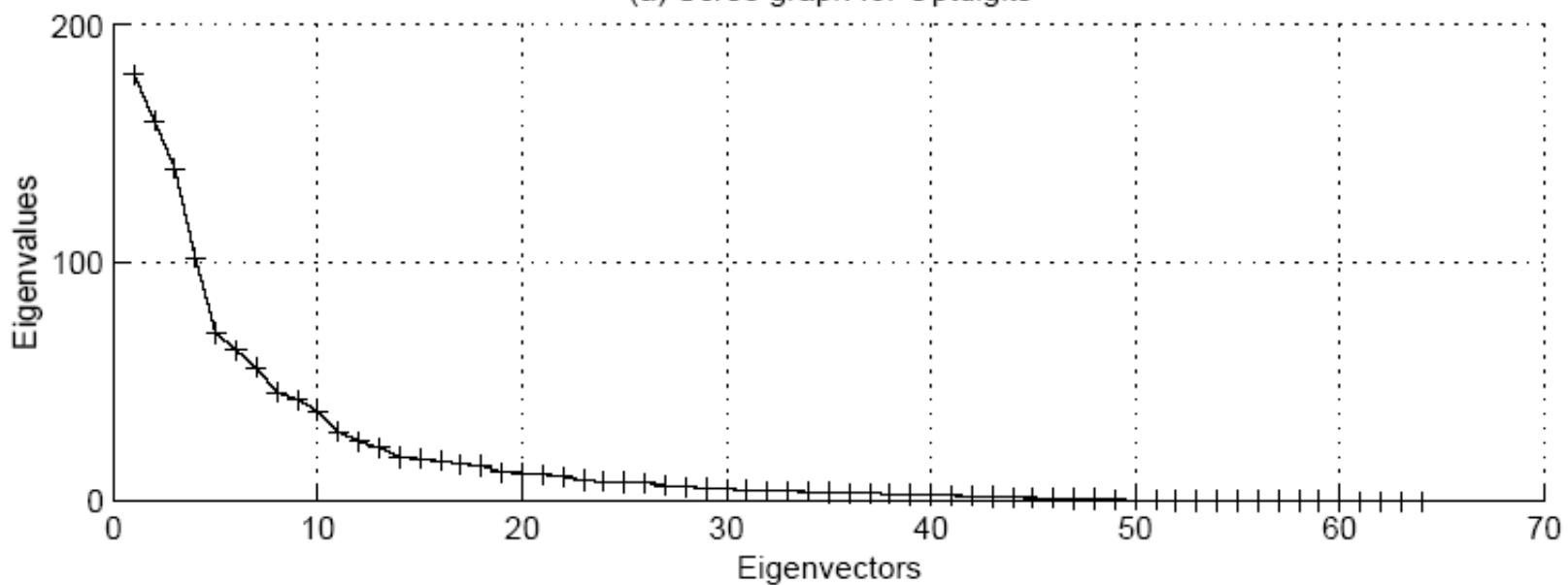
- Proportion of Variance (PoV) explained

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

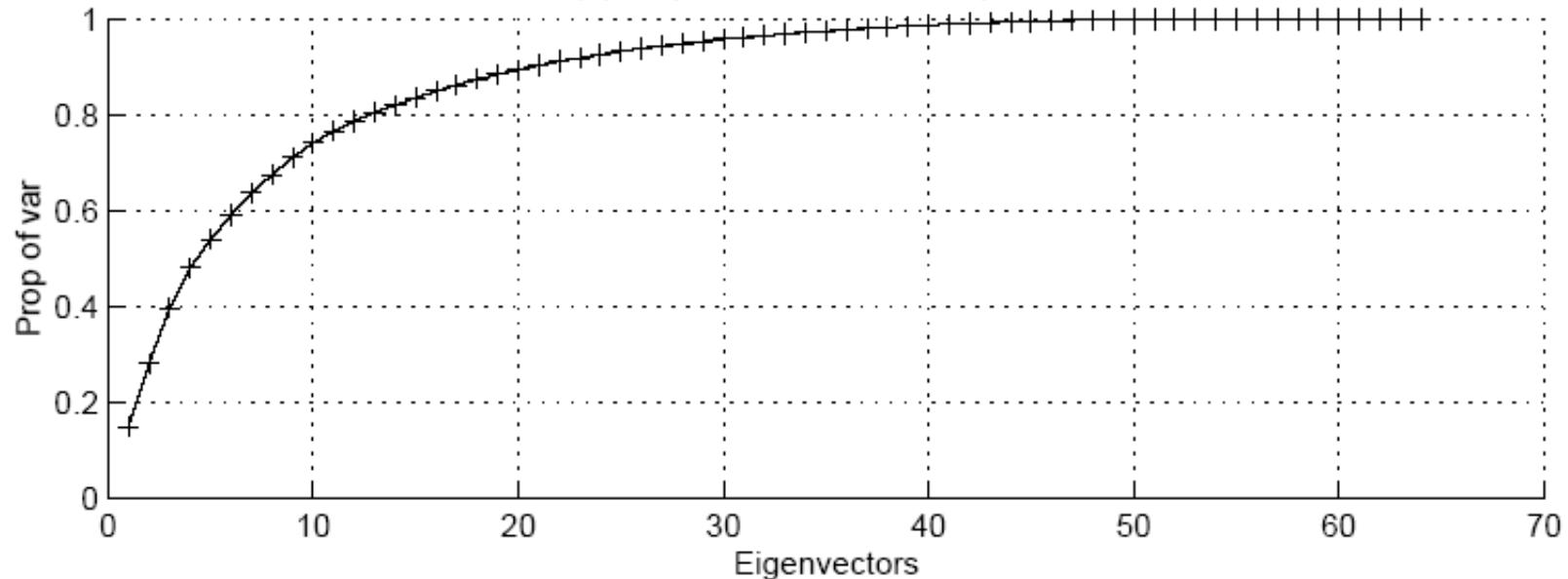
when λ_i are sorted in descending order

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

(a) Scree graph for Optdigits



(b) Proportion of variance explained



PCA Discussions

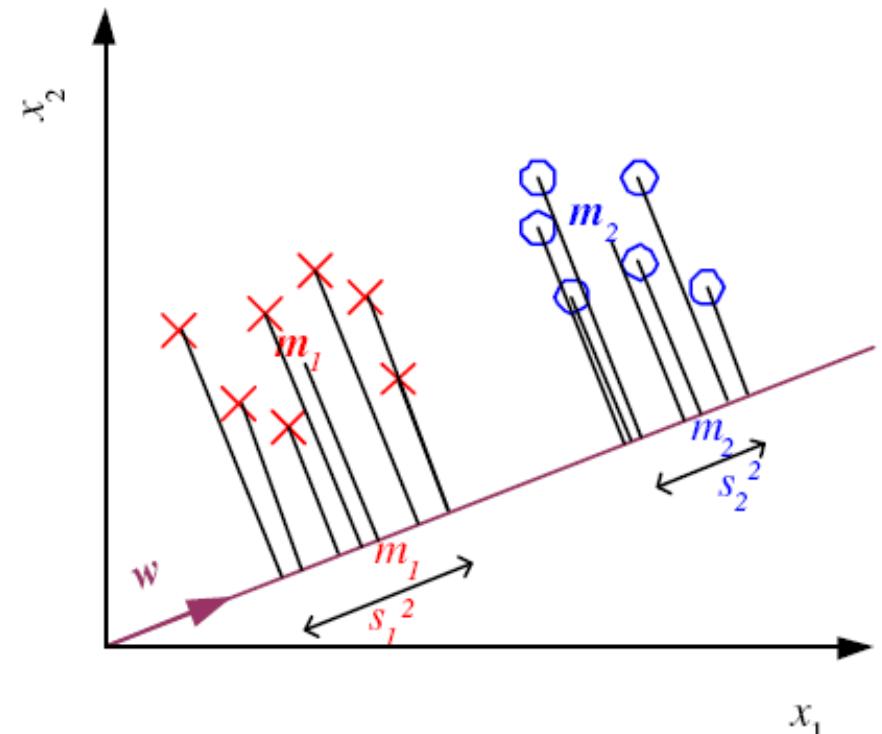
- Linearity: Linear rotation of the original space
- Gaussian assumption: mean and variance are enough to characterize the noise and redundancy.
- We choose the principal components to be orthogonal, but they don't have to be.
- Often used with clustering algorithms to cluster high-dimensional data

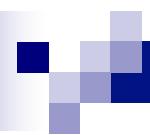
Linear Discriminant Analysis

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

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

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





■ Between-class scatter:

$$\begin{aligned} (m_1 - m_2)^2 &= (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 \\ &= \mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w} \\ &= \mathbf{w}^T \mathbf{S}_B \mathbf{w} \text{ where } \mathbf{S}_B = (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \end{aligned}$$

■ Within-class scatter:

$$\begin{aligned} s_1^2 &= \sum_t (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t \\ &= \sum_t \mathbf{w}^T (\mathbf{x}^t - \mathbf{m}_1) (\mathbf{x}^t - \mathbf{m}_1)^T \mathbf{w} r^t = \mathbf{w}^T \mathbf{S}_1 \mathbf{w} \end{aligned}$$

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

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

Fisher's Linear Discriminant

- Find \mathbf{w} that max

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} = \frac{|\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)|^2}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

$$\left(\frac{u}{v} \right)' = \frac{u'v - uv'}{v^2}$$

$$\frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} (2(\mathbf{m}_1 - \mathbf{m}_2) - \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}} 2\mathbf{S}_W \mathbf{w}) = 0$$

- LDA soln:

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

K>2 Classes

- Within-class scatter:

$$\mathbf{S}_W = \sum_{i=1}^K \mathbf{S}_i \quad \mathbf{S}_i = \sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i) (\mathbf{x}^t - \mathbf{m}_i)^T$$

- Between-class scatter (among means):

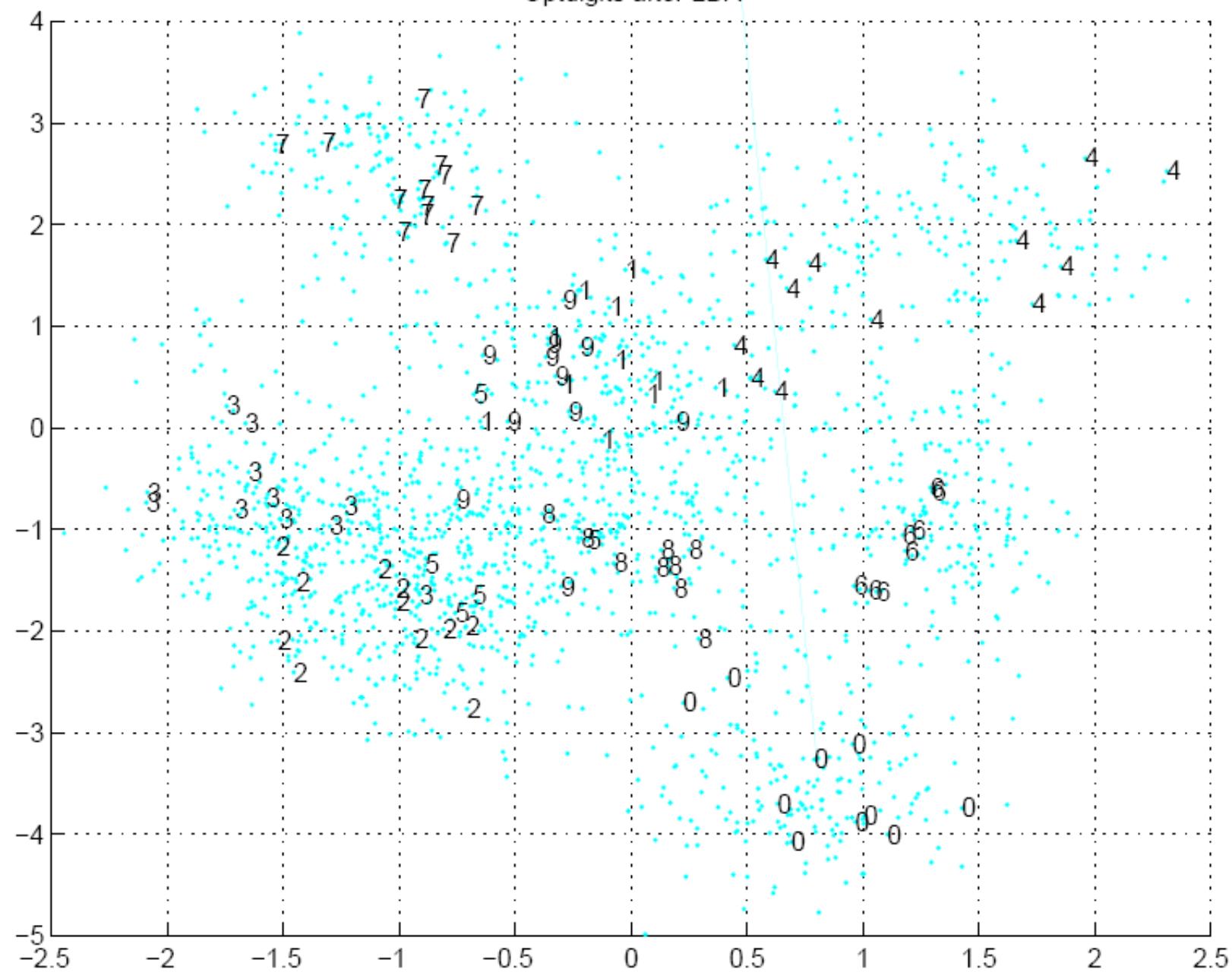
$$\mathbf{S}_B = \sum_{i=1}^K N_i (\mathbf{m}_i - \mathbf{m}) (\mathbf{m}_i - \mathbf{m})^T \quad \mathbf{m} = \frac{1}{K} \sum_{i=1}^K \mathbf{m}_i$$

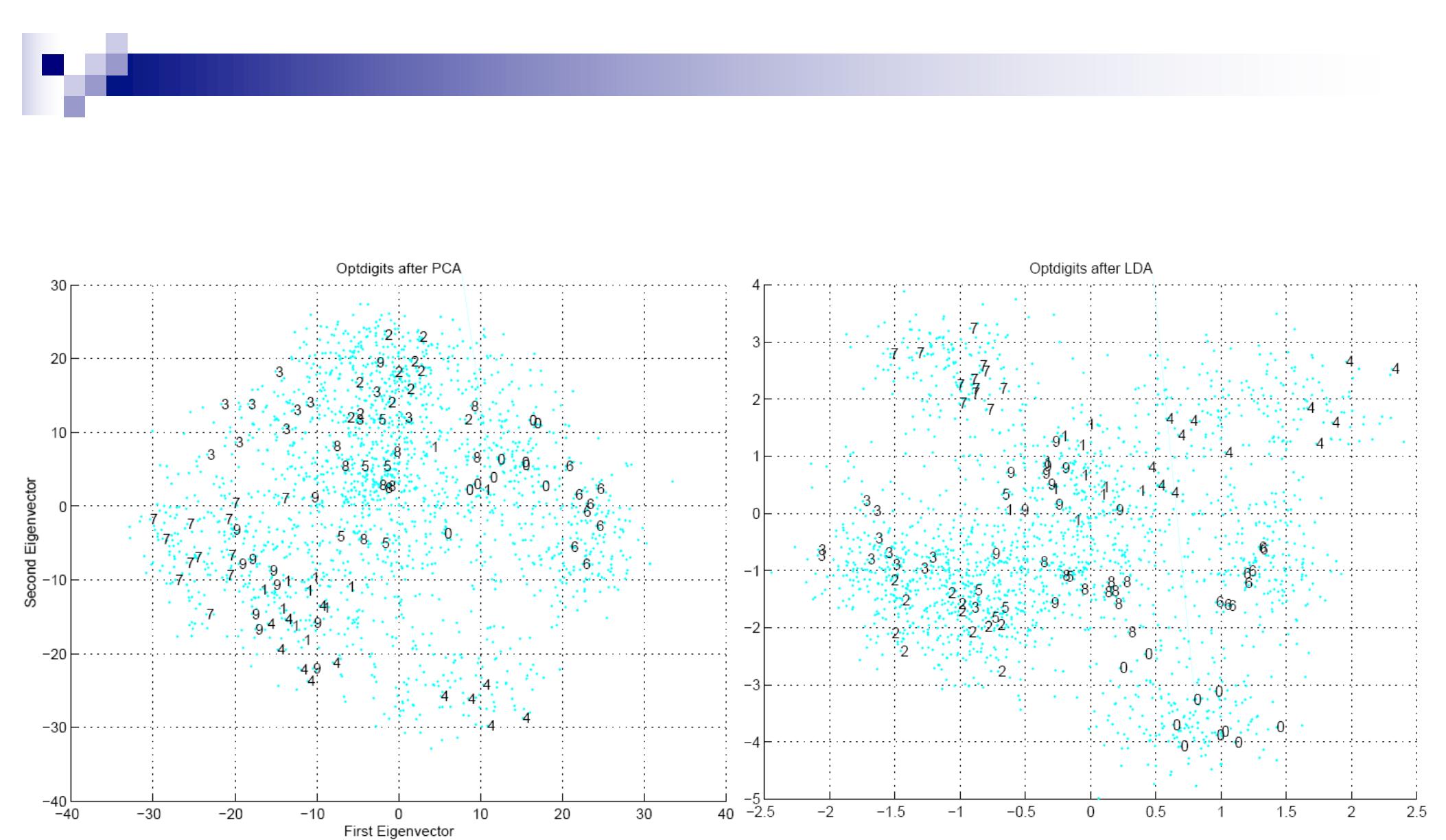
- Find \mathbf{W} that max

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|}$$

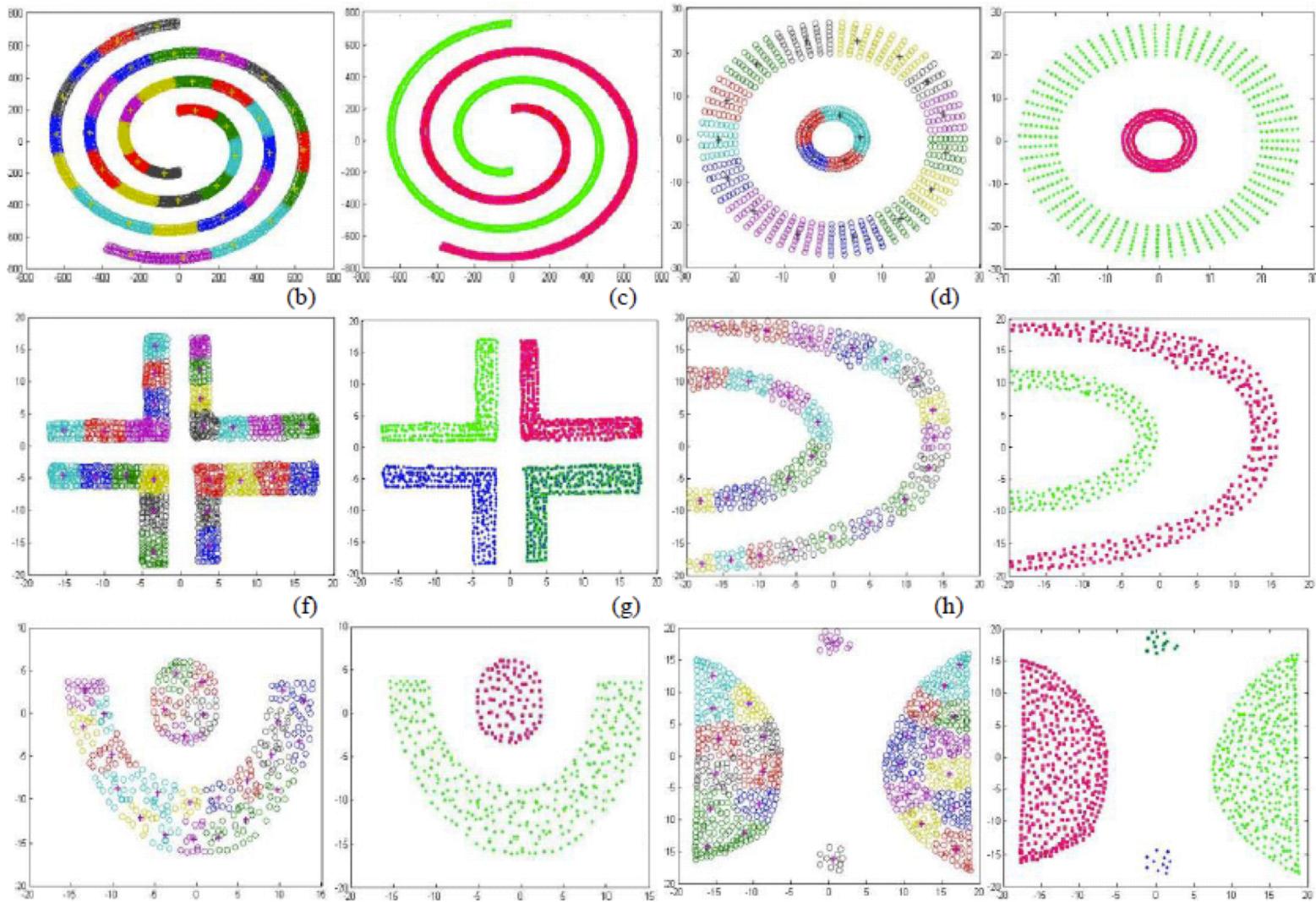
The largest eigenvectors of $\mathbf{S}_W^{-1} \mathbf{S}_B$
Maximum rank of $K-1$

Optdigits after LDA



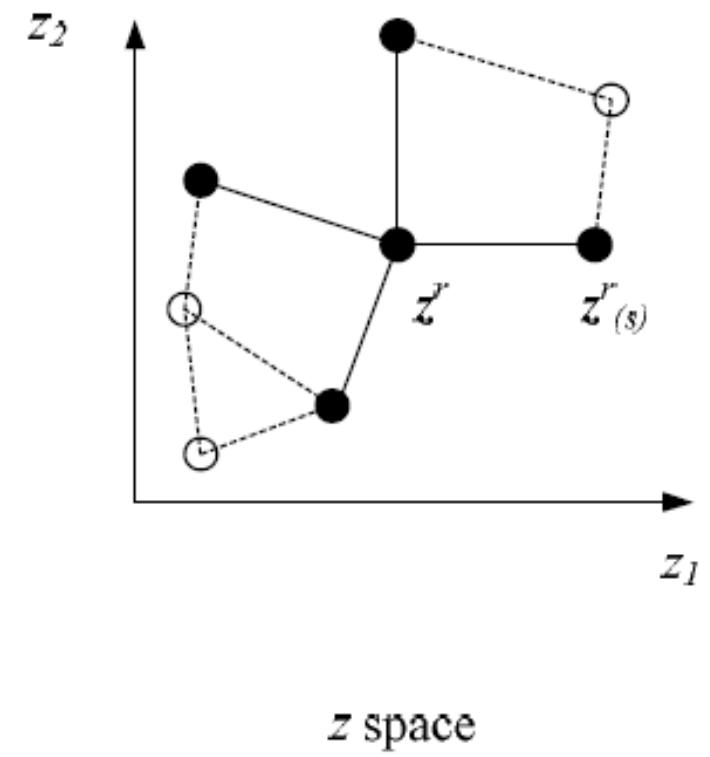
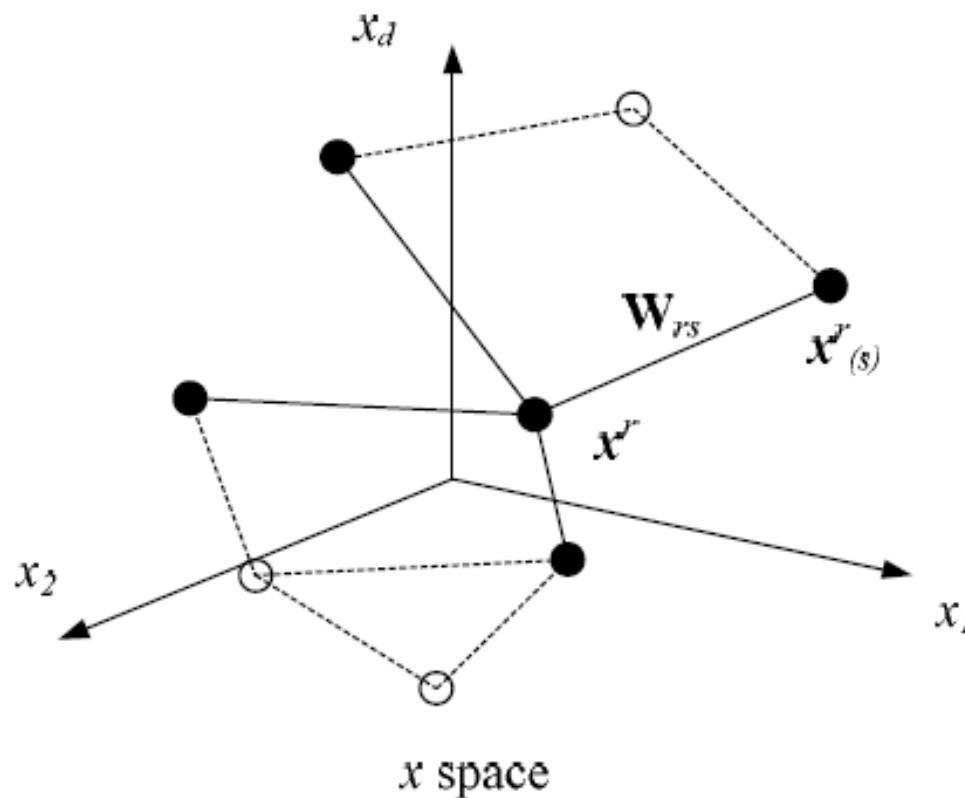


Projection of Non-linear Data

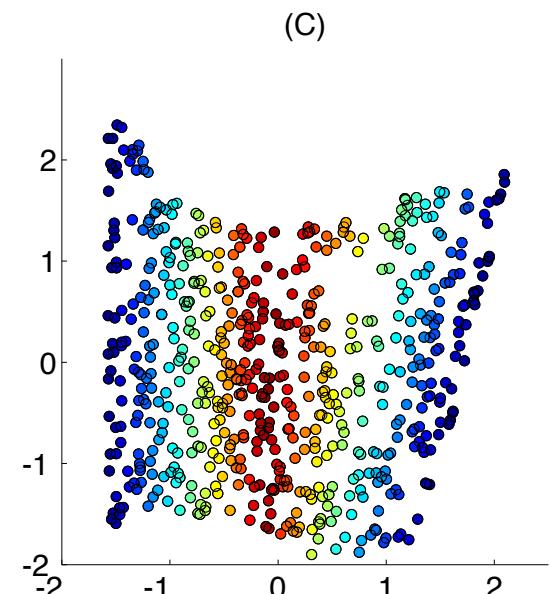
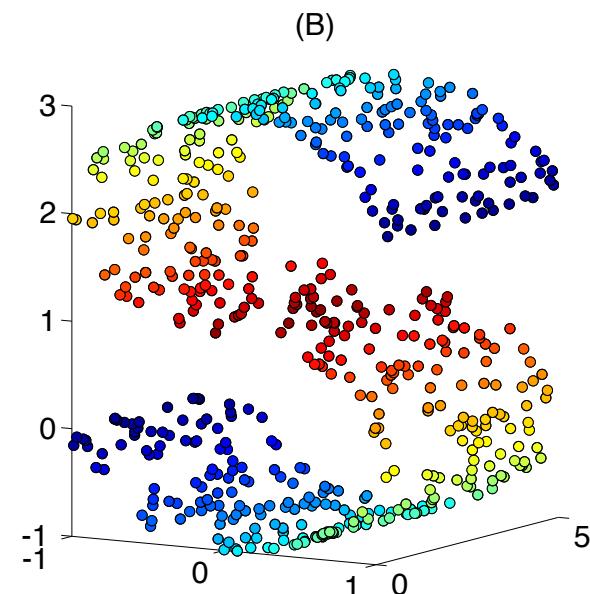
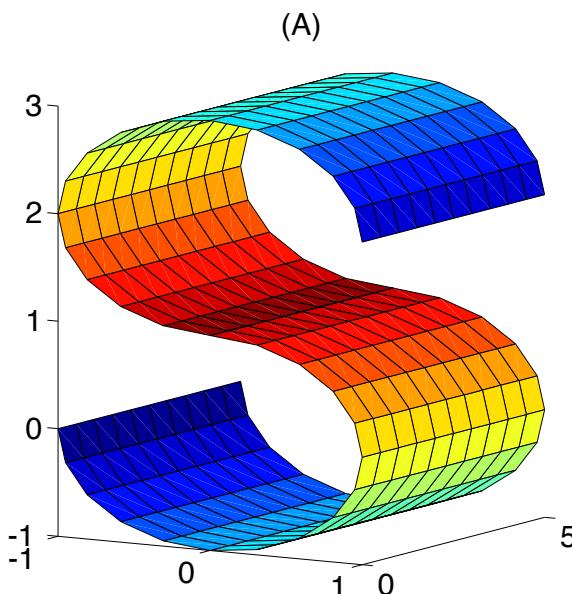


Locally Linear Embedding

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



LLE on S Manifold



Locally Linear Embedding

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

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

3. Find the new coordinates \mathbf{z}^r that minimize

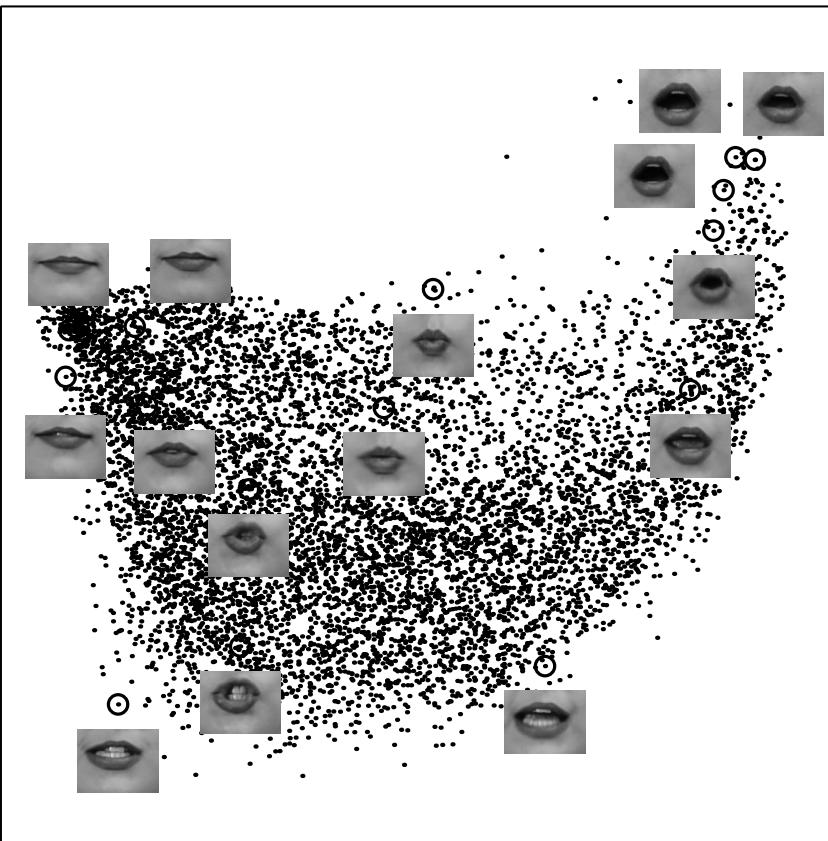
$$E(\mathbf{z} | \mathbf{W}) = \sum_r \left\| \mathbf{z}^r - \sum_s \mathbf{W}_{rs} \mathbf{z}^s_{(r)} \right\|^2, \text{ subj: } \sum_r \mathbf{z}^r = \mathbf{0} \text{ and } \sum_r (\mathbf{z}^r)^T \mathbf{z}^r = 1$$

Locally Linear Embedding

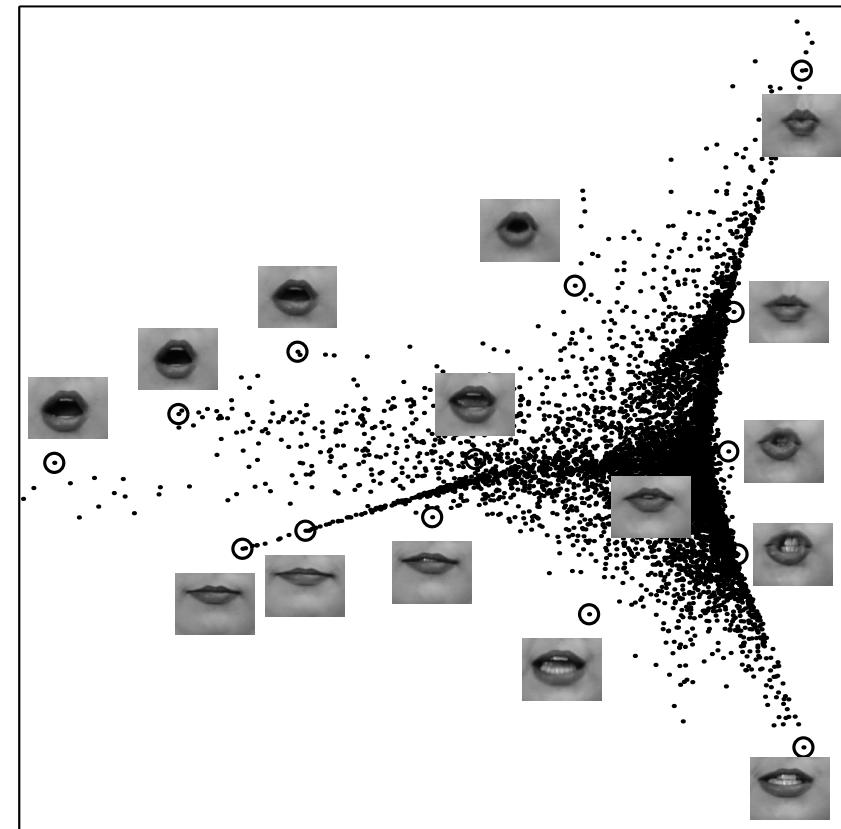
$$\begin{aligned} E(\mathbf{z} | \mathbf{W}) &= \sum_r \left\| z^r - \sum_s \mathbf{W}_{rs} z_{(r)}^s \right\|^2, \text{ subj: } \sum_r z^r = 0 \text{ and } \sum_r (z^r)^T z^r = 1 \\ &= \sum_{r,s} (\delta_{rs} - W_{rs} - W_{sr} + \sum_i W_{is} W_{ir}) (z^r)^T z^s \end{aligned}$$

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

LLE on Lip Images

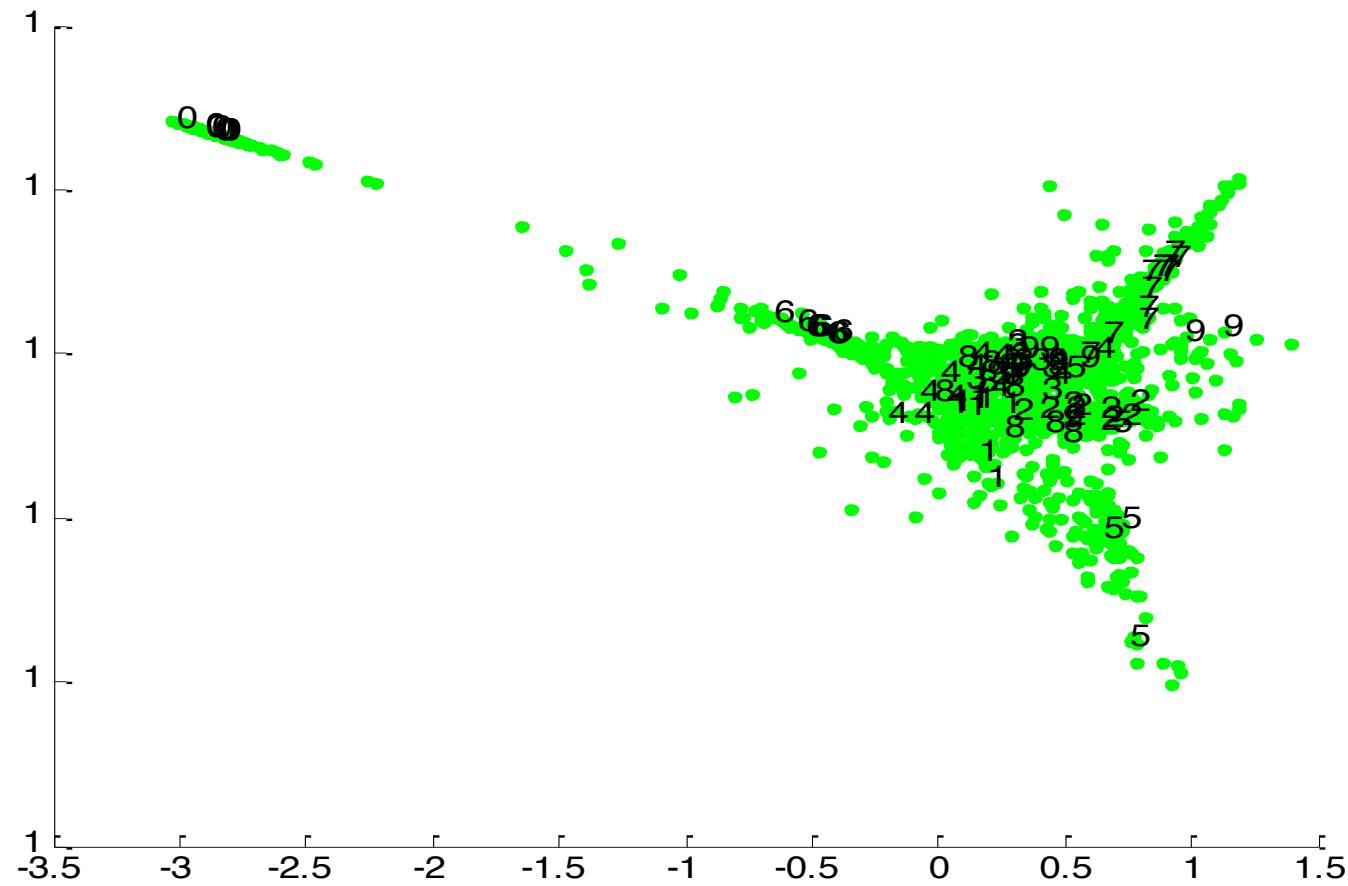


PCA



LLE

LLE on Optdigits



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

Summary

■ Linear projection:

- **PCA**: Project \mathbf{x} to \mathbf{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

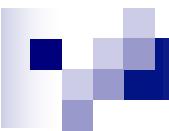
- Feature selection (subset selection):
Choosing $k < d$ important features, ignoring
the remaining $d - k$

Subset selection algorithms

- Feature extraction: Project the
original x_i , $i = 1, \dots, d$ dimensions to
new $k < d$ dimensions, z_j , $j = 1, \dots, k$
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
 - Set of features F initially \emptyset .
 - At each iteration, find the best new feature
$$j = \operatorname{argmin}_i E(F \cup x_i)$$
 - 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%

The diagram illustrates the process of feature selection, specifically Forward Selection, starting from a list of ranked features.

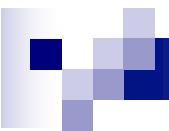
Rank features:

Vegetarian	No
Plays video games	Yes
Family history	No
Athletic	No
Smoker	Yes
Sex	Male
Lung capacity	5.8L
Hair color	Red
Car	Audi
...	
Weight	185 lbs

Forward Selection:

The process starts with the feature "Family history". It then moves to the next highest ranked feature, "Smoker", and adds it to the selected set. This continues until all features have been evaluated.

Family history	No
Smoker	Yes
Weight	185 lbs
Sex	Male
Plays video games	No
Athletic	No
Lung capacity	5.8L
Vegetarian	No
Hair color	Red
...	
Car	Audi



Forward Selection

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

Classification

Accuracy: 75%→95%

The diagram illustrates the forward selection process for feature selection. It starts with a list of features and their values, which are then ranked. The top-ranked features are selected for the model.

Rank features

Vegetarian	No
Plays video games	Yes
Family history	No
Athletic	No
Smoker	Yes
Sex	Male
Lung capacity	5.8L
Hair color	Red
Car	Audi
...	
Weight	185 lbs

Forward Selection

Family history	No
Smoker	Yes
Weight	185 lbs
Sex	Male
Plays video games	No
Athletic	No
Lung capacity	5.8L
Vegetarian	No
Hair color	Red
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%

The diagram illustrates the process of feature selection, specifically forward selection, starting from a full set of features and ranking them based on their importance.

Rank features:

Vegetarian	No
Plays video games	Yes
Family history	No
Athletic	No
Smoker	Yes
Sex	Male
Lung capacity	5.8L
Hair color	Red
Car	Audi
...	
Weight	185 lbs

Forward Selection:

Family history	No
Smoker	Yes
Weight	185 lbs
Sex	Male
Plays video games	No
Athletic	No
Lung capacity	5.8L
Vegetarian	No
Hair color	Red
Car	Audi

The process starts with all features listed on the left. A blue arrow labeled "Rank features" points to the list of all features. Another blue arrow labeled "Forward Selection" points to the subset of features on the right. The feature "Weight" is highlighted in red in the "Forward Selection" list, indicating it is the next feature selected.

Backward Elimination

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

Vegetarian	No
Plays video games	Yes
Family history	No
Athletic	No
Smoker	Yes
Sex	Male
Lung capacity	5.8L
Hair color	Red
Car	Audi
...	
Weight	185 lbs

Rank features



Family history	No
Smoker	Yes
Weight	185 lbs
Sex	Male
Plays video games	No
Athletic	No
Lung capacity	5.8L
Vegetarian	No
Hair color	Red
...	
Car	Audi

Classification
Accuracy: 60% → 75%

Forward Selection



Family history	No
Smoker	Yes
Weight	185 lbs
Sex	Male
Plays video games	No
Athletic	No
Lung capacity	5.8L
Vegetarian	No
Hair color	Red
...	
Car	Audi

Backward Elimination

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

Classification Accuracy: 95%

Vegetarian	No
Plays video games	Yes
Family history	No
Athletic	No
Smoker	Yes
Sex	Male
Lung capacity	5.8L
Hair color	Red
Car	Audi
...	
Weight	185 lbs

No
Yes
No
No
Yes
Male
5.8L
Red
Audi
185 lbs

Rank features



Family history	No
Smoker	Yes
Weight	185 lbs
Sex	Male
Plays video games	No
Athletic	No
Lung capacity	5.8L
Vegetarian	No
Hair color	Red
...	
Car	Audi

No
Yes
185 lbs
Male
No

Forward Selection



Family history	No
Smoker	Yes

No
Yes



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