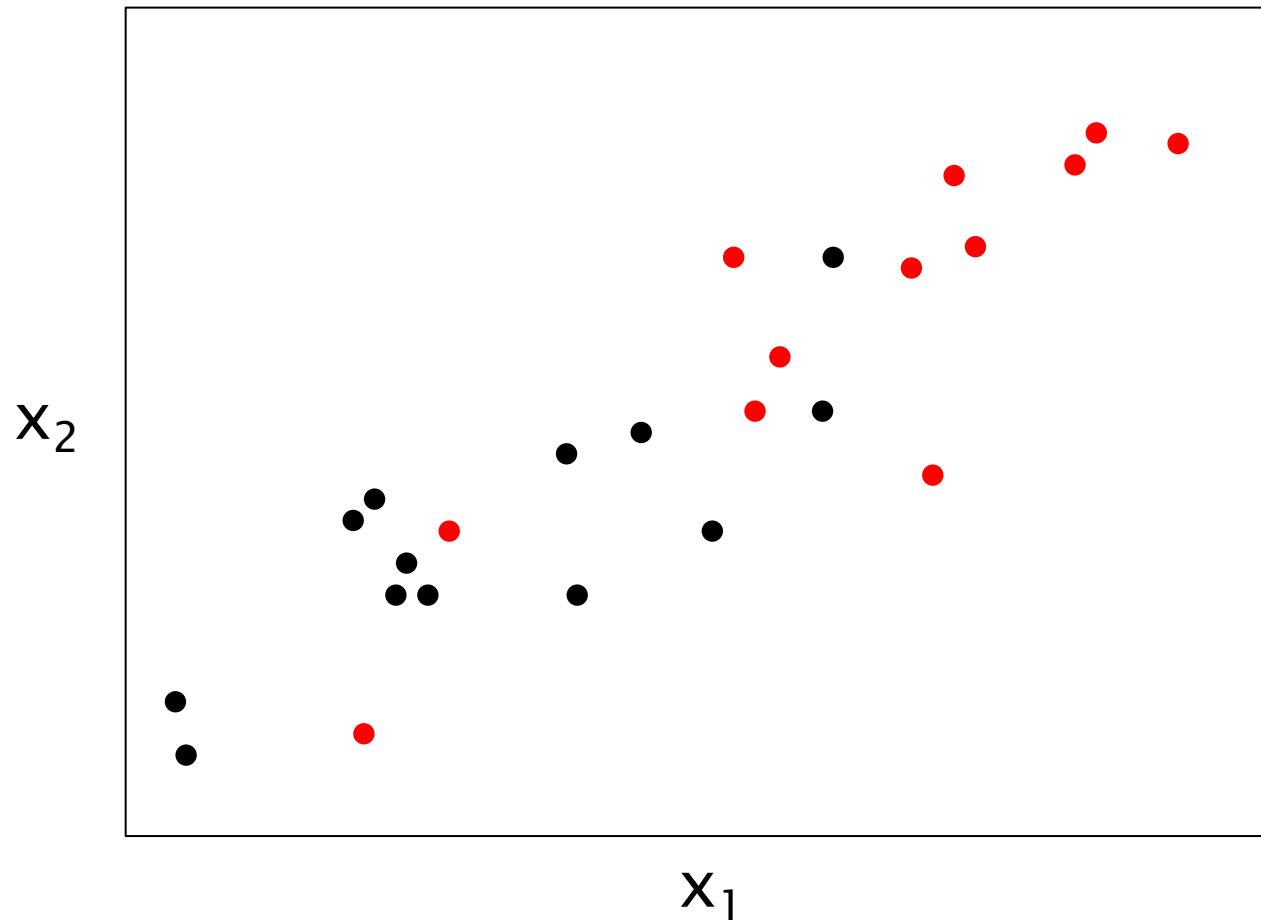


# Dimensionality Reduction by Feature Extraction

## Outline:

- Principal Component Analysis
- Independent Component Analysis
- ISOMAP

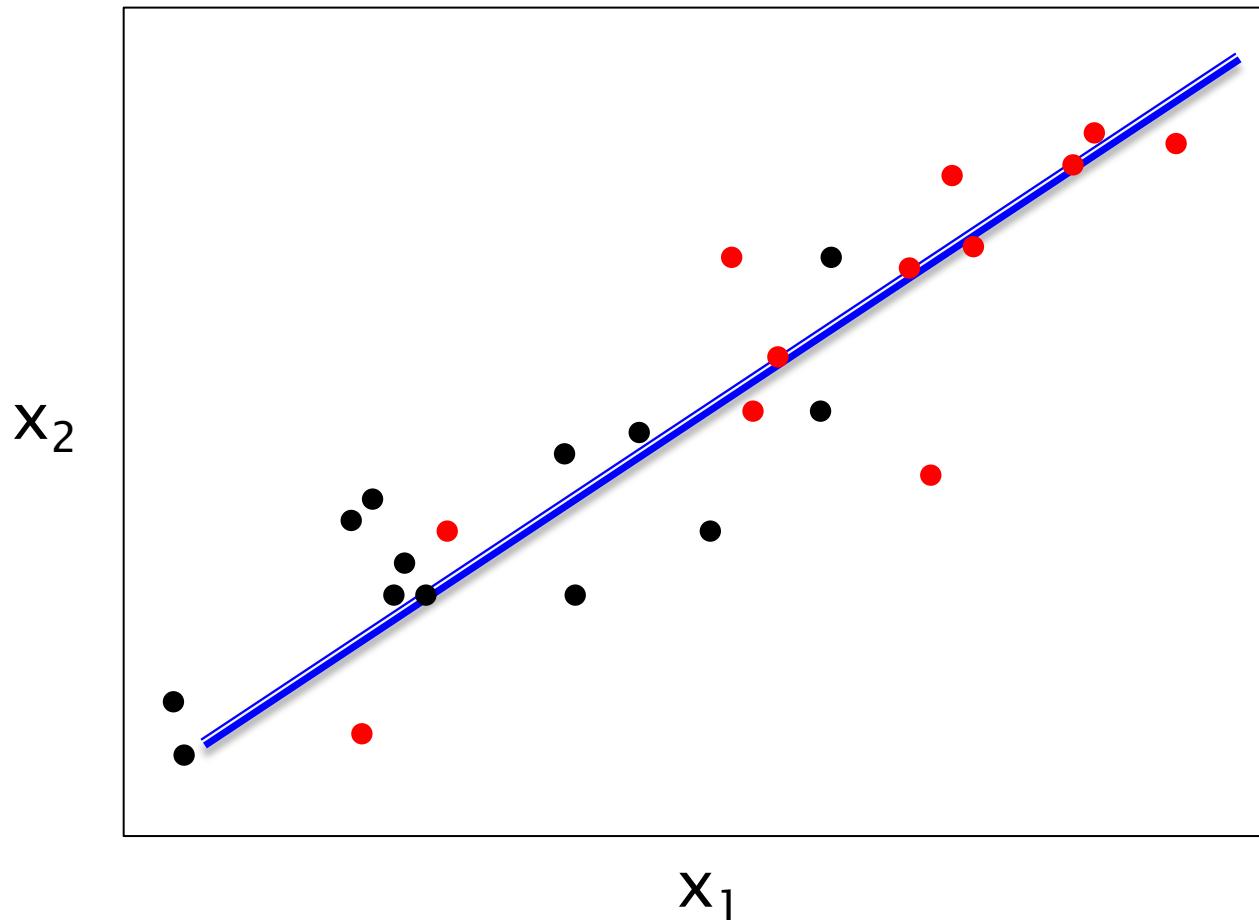
# Data



# Data

First principal component

Gives direction of  
largest variation of the data



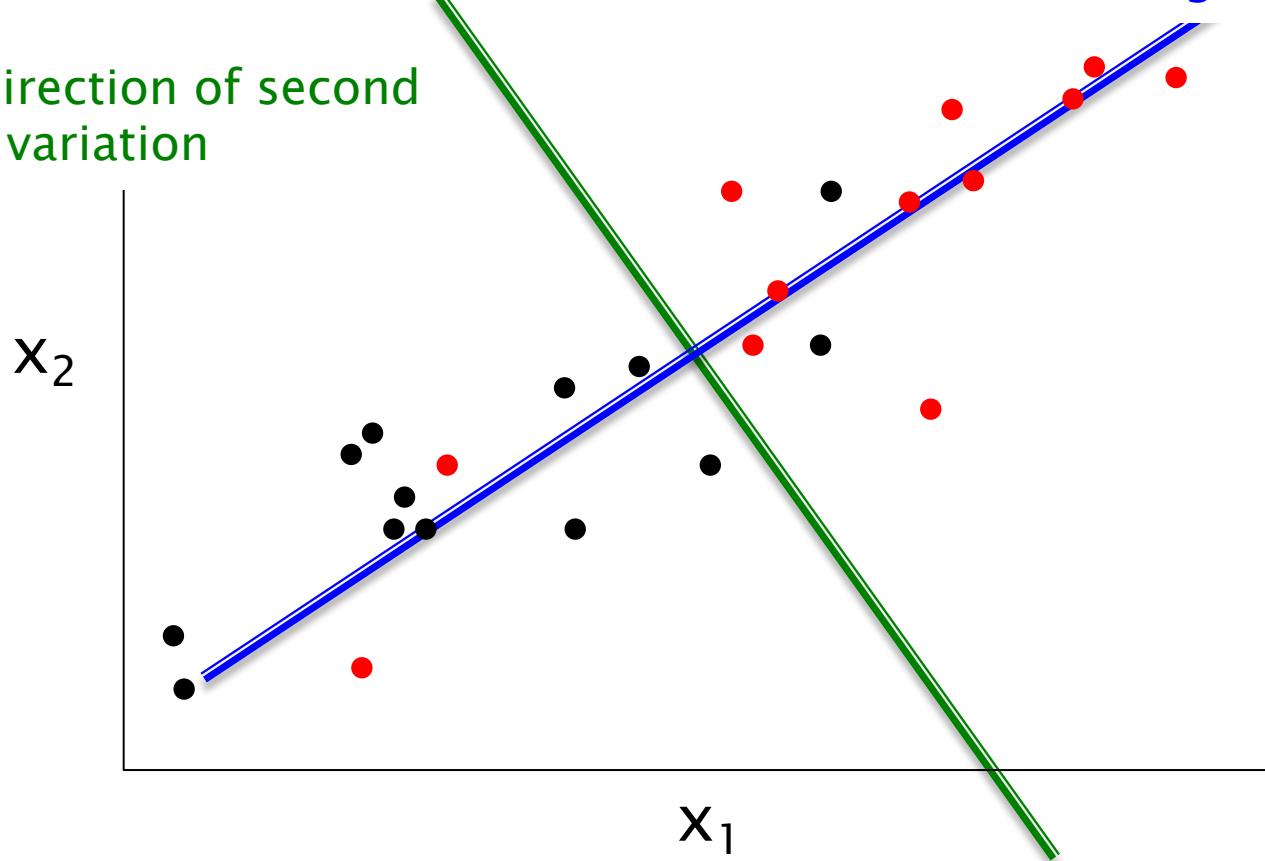
# Data

First principal component

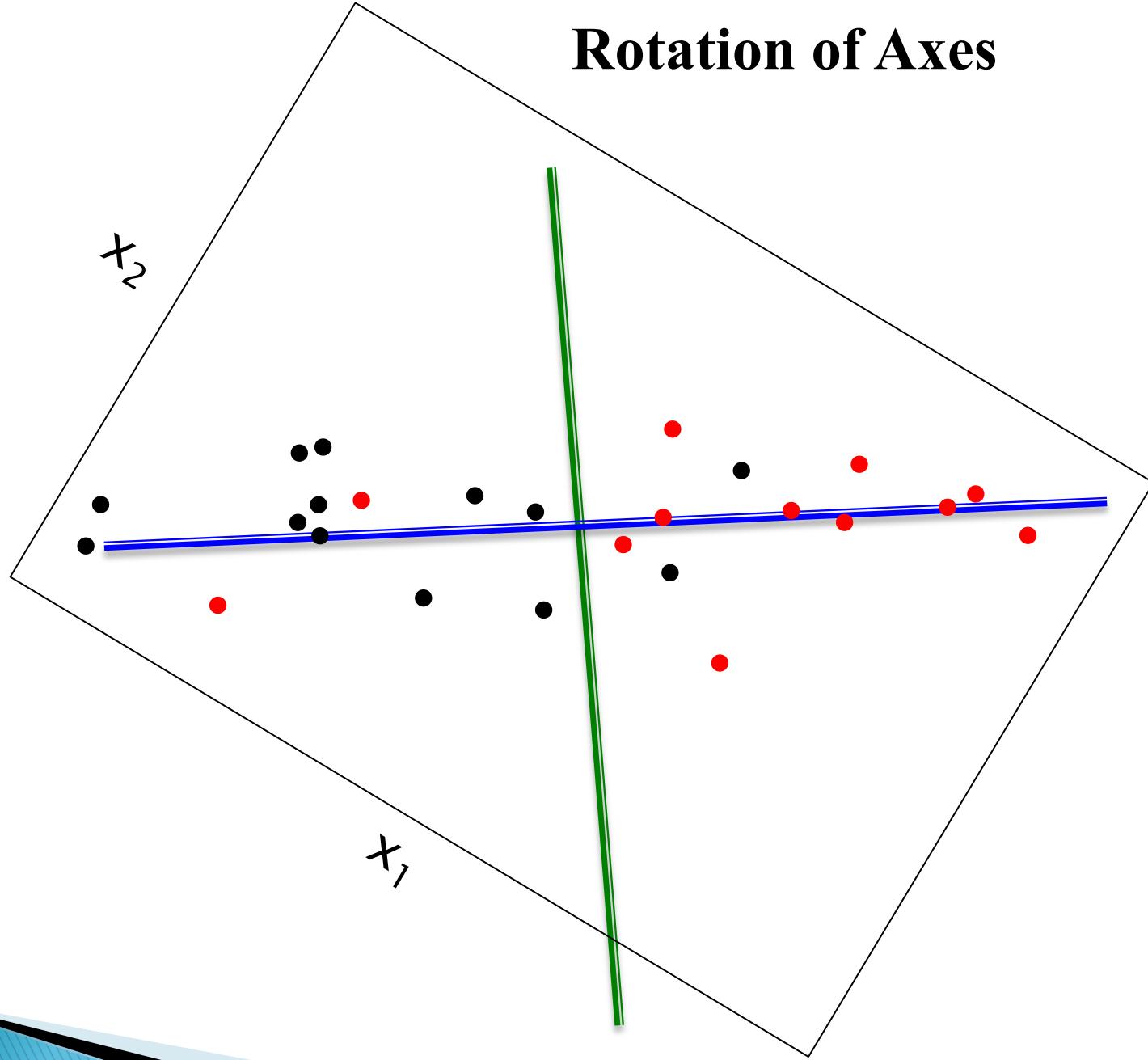
Second principal component

Gives direction of largest variation of the data

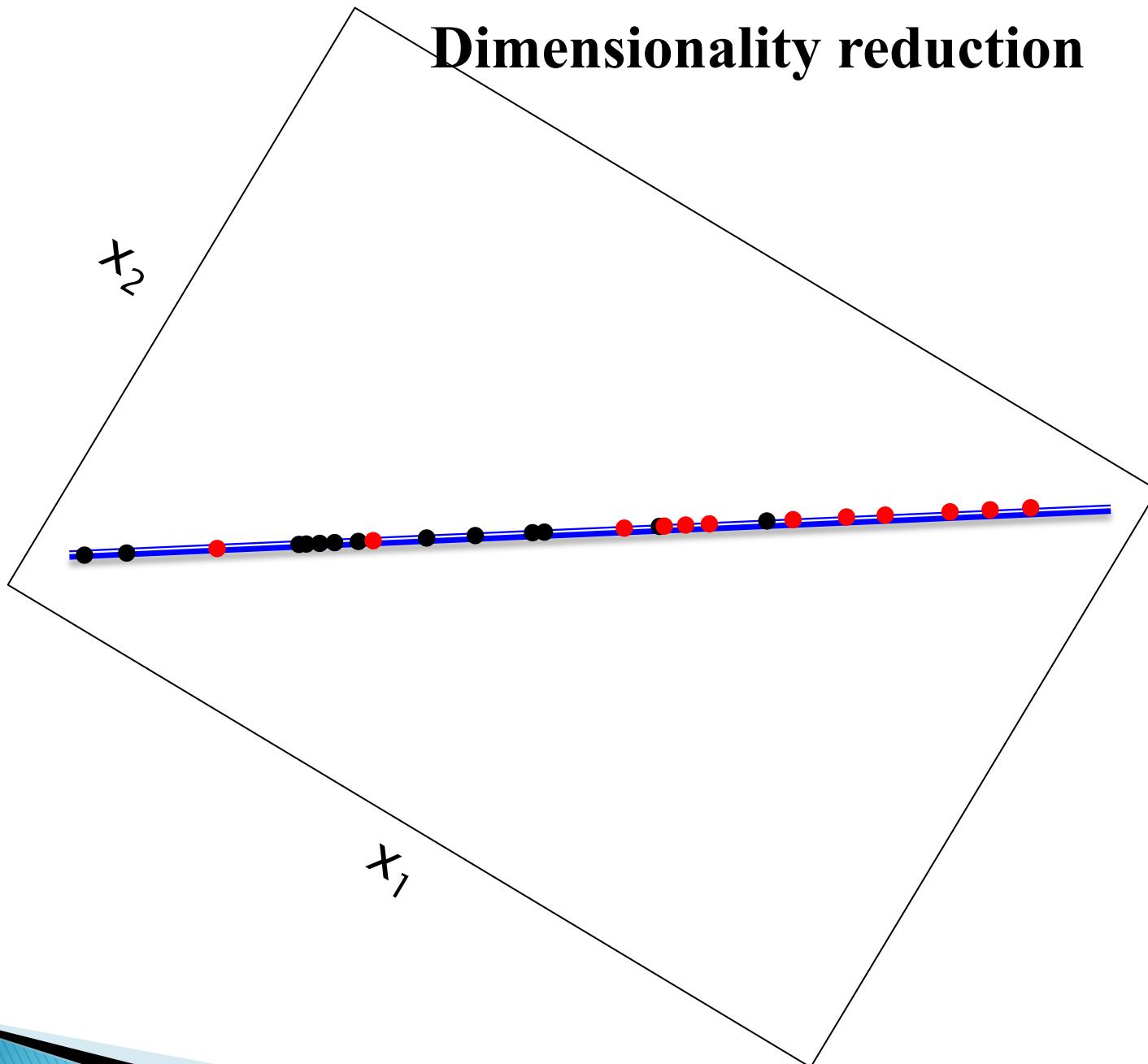
Gives direction of second largest variation



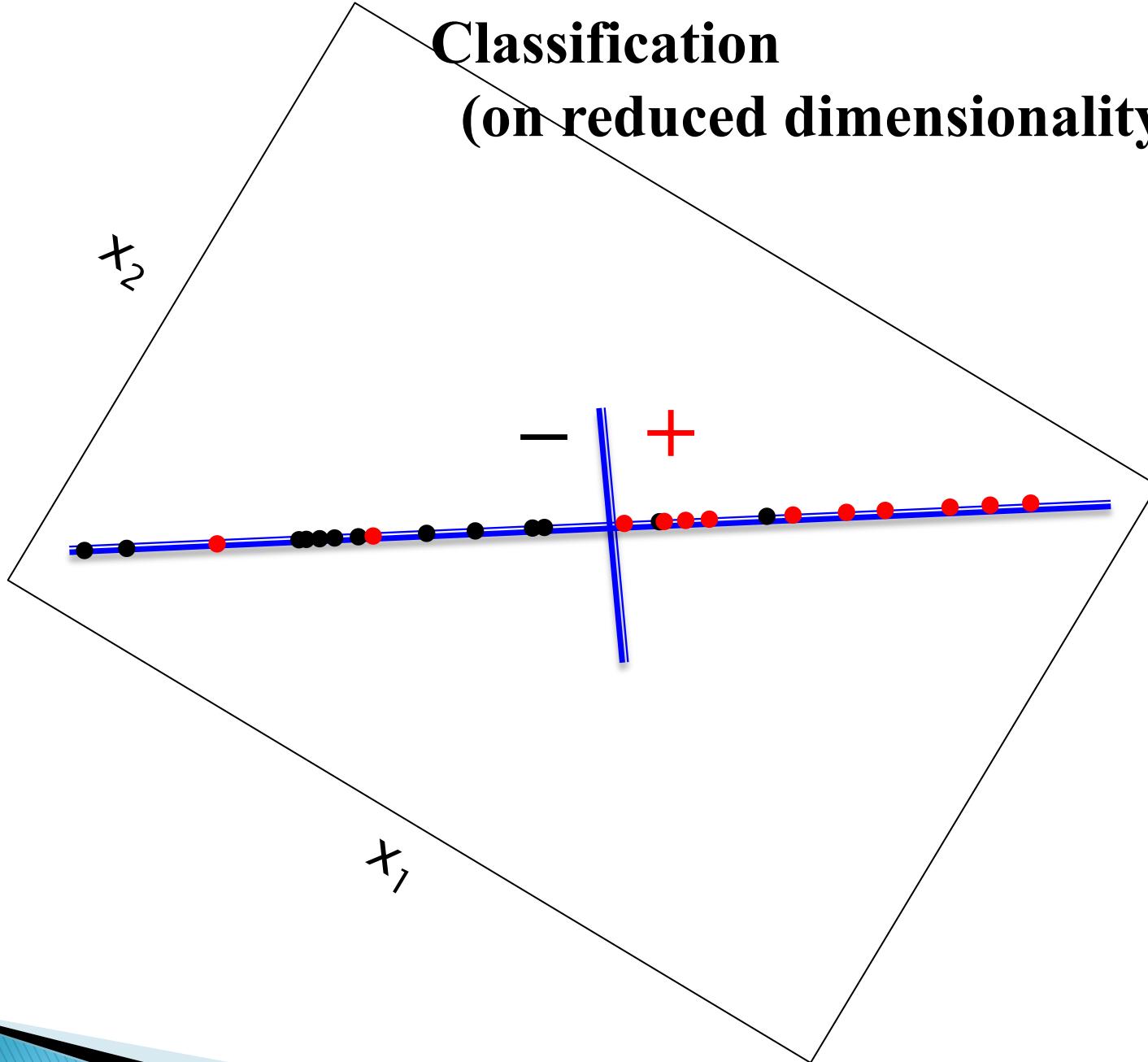
# Rotation of Axes



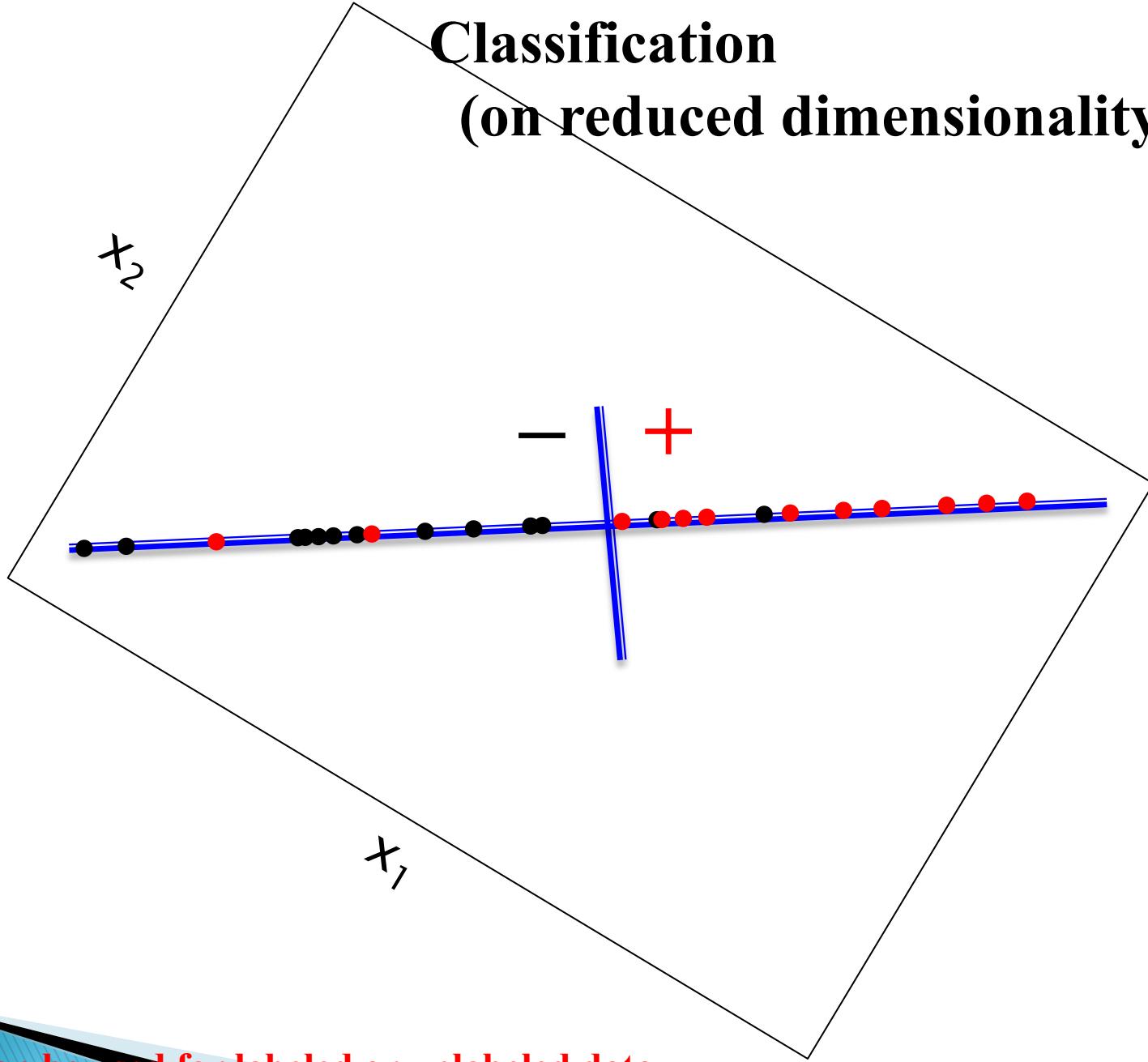
# Dimensionality reduction



# Classification (on reduced dimensionality space)



# Classification (on reduced dimensionality space)



Note: Can be used for labeled or unlabeled data.

# Principal Components Analysis (PCA)

- ▶ **Summary:** PCA finds new orthogonal axes in directions of largest variation in data.
- ▶ PCA used to create high-level features in order to improve classification and reduce dimensions of data without much loss of information.
- ▶ Used in machine learning and in signal processing and image compression (among other things).

# Background for PCA

- ▶ Suppose attributes are  $A_1$  and  $A_2$ , and we have  $n$  training examples.  $x$ 's denote values of  $A_1$  and  $y$ 's denote values of  $A_2$  over the training examples.
- ▶ Variance of an attribute:

$$\text{var}(A_1) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{(n - 1)}$$

- ▶ Covariance of two attributes:

$$\text{cov}(A_1, A_2) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1)}$$

- ▶ If covariance is positive, both dimensions increase together. If negative, as one increases, the other decreases. Zero: independent of each other.
- ▶ Covariance matrix
  - Suppose we have  $n$  attributes,  $A_1, \dots, A_n$ .
  - Covariance matrix:

$$C^{n \times n} = (c_{i,j}), \text{ where } c_{i,j} = \text{cov}(A_i, A_j)$$

	<i>Hours(H)</i>	<i>Mark(M)</i>
Data	9	39
	15	56
	25	93
	14	61
	10	50
	18	75
	0	32
	16	85
	5	42
	19	70
	16	66
	20	80
Totals	167	749
Averages	13.92	62.42

Covariance:

<i>H</i>	<i>M</i>	$(H_i - \bar{H})$	$(M_i - \bar{M})$	$(H_i - \bar{H})(M_i - \bar{M})$
9	39	-4.92	-23.42	115.23
15	56	1.08	-6.42	-6.93
25	93	11.08	30.58	338.83
14	61	0.08	-1.42	-0.11
10	50	-3.92	-12.42	48.69
18	75	4.08	12.58	51.33
0	32	-13.92	-30.42	423.45
16	85	2.08	22.58	46.97
5	42	-8.92	-20.42	182.15
19	70	5.08	7.58	38.51
16	66	2.08	3.58	7.45
20	80	6.08	17.58	106.89
Total				1149.89
Average				104.54

Table 2.2: 2-dimensional data set and covariance calculation

$$\begin{pmatrix} \text{cov}(H,H) & \text{cov}(H,M) \\ \text{cov}(M,H) & \text{cov}(M,M) \end{pmatrix}$$

$$= \begin{pmatrix} \text{var}(H) & 104.5 \\ 104.5 & \text{var}(M) \end{pmatrix}$$

$$= \begin{pmatrix} 47.7 & 104.5 \\ 104.5 & 370 \end{pmatrix}$$

Covariance matrix

# Review of Matrix Algebra

## ► Eigenvectors:

- Let  $\mathbf{M}$  be an  $n \times n$  matrix.
  - $\mathbf{v}$  is an *eigenvector* of  $\mathbf{M}$  if  $\mathbf{M} \times \mathbf{v} = \lambda \mathbf{v}$
  - $\lambda$  is called the *eigenvalue* associated with  $\mathbf{v}$
- You can always choose eigenvectors of length 1:

$$\mathbf{M} \times a\mathbf{v} = \lambda a\mathbf{v}$$

- If  $\mathbf{M}$  is symmetric with real entries, it has  $n$  eigenvectors, and they are orthogonal to one another.
- Thus, eigenvectors can be used as a new basis for a  $n$ -dimensional vector space.

# Textbook's notation

- ▶ We have original data  $\mathbf{X}$  and mean-subtracted data  $\mathbf{B}$ , and covariance matrix  $\mathbf{C} = cov(\mathbf{B})$ , where  $\mathbf{C}$  is an  $M \times N$  matrix:

$$\mathbf{C} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$$

- ▶ We find matrix  $\mathbf{V}$  such that the columns of  $\mathbf{V}$  are the  $N$  eigenvectors of  $\mathbf{C}$  and

$$\mathbf{D} = \mathbf{V}^{-1} \mathbf{C} \mathbf{V} = \begin{pmatrix} \lambda_1 & 0 & 0 & 0 \\ 0 & \lambda_2 & 0 & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \lambda_N \end{pmatrix}$$

where  $\lambda_i$  is the  $i$ -th eigenvalue of  $\mathbf{C}$ .

- ▶ Each eigenvalue in  $\mathbf{D}$  corresponds to an eigenvector in  $\mathbf{V}$ . The eigenvectors, sorted in order of decreasing eigenvalue, become the “feature vector” for PCA.

# Principal Components Analysis

---

## The Principal Components Analysis Algorithm

---

- Write  $N$  datapoints  $\mathbf{x}_i = (\mathbf{x}_{1i}, \mathbf{x}_{2i}, \dots, \mathbf{x}_{Mi})$  as row vectors
  - Put these vectors into a matrix  $\mathbf{X}$  (which will have size  $N \times M$ )
  - Centre the data by subtracting off the mean of each column, putting it into matrix  $\mathbf{B}$
  - Compute the covariance matrix  $\mathbf{C} = \frac{1}{N}\mathbf{B}^T\mathbf{B}$
  - Compute the eigenvalues and eigenvectors of  $\mathbf{C}$ , so  $\mathbf{V}^{-1}\mathbf{C}\mathbf{V} = \mathbf{D}$ , where  $\mathbf{V}$  holds the eigenvectors of  $\mathbf{C}$  and  $\mathbf{D}$  is the  $M \times M$  diagonal eigenvalue matrix
  - Sort the columns of  $\mathbf{D}$  into order of decreasing eigenvalues, and apply the same order to the columns of  $\mathbf{V}$
  - Reject those with eigenvalue less than some  $\eta$ , leaving  $L$  dimensions in the data
-

# Principal Components Analysis

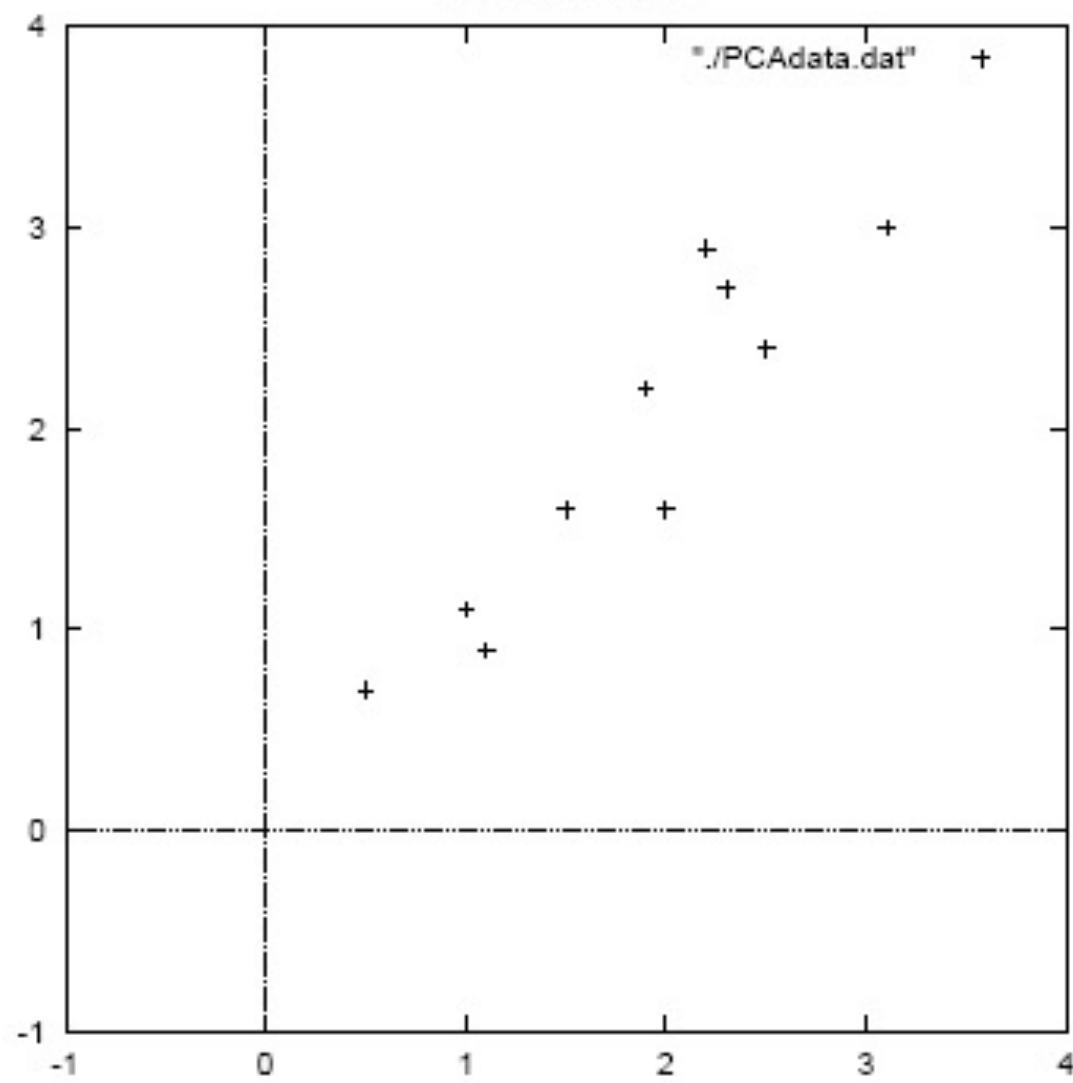
- Given original data set  $S = \{x^1, \dots, x^k\}$ , produce new set by subtracting the mean of attribute  $A_i$  from each  $x_i$ .

	$x$	$y$		$x$	$y$
Data =	2.5	2.4	DataAdjust =	.69	.49
	0.5	0.7		-1.31	-1.21
	2.2	2.9		.39	.99
	1.9	2.2		.09	.29
	3.1	3.0		1.29	1.09
	2.3	2.7		.49	.79
	2	1.6		.19	-.31
	1	1.1		-.81	-.81
	1.5	1.6		-.31	-.31
	1.1	0.9		-.71	-1.01

Mean: 1.81 1.91

Mean: 0 0

Original PCA data



2. Calculate the covariance matrix:

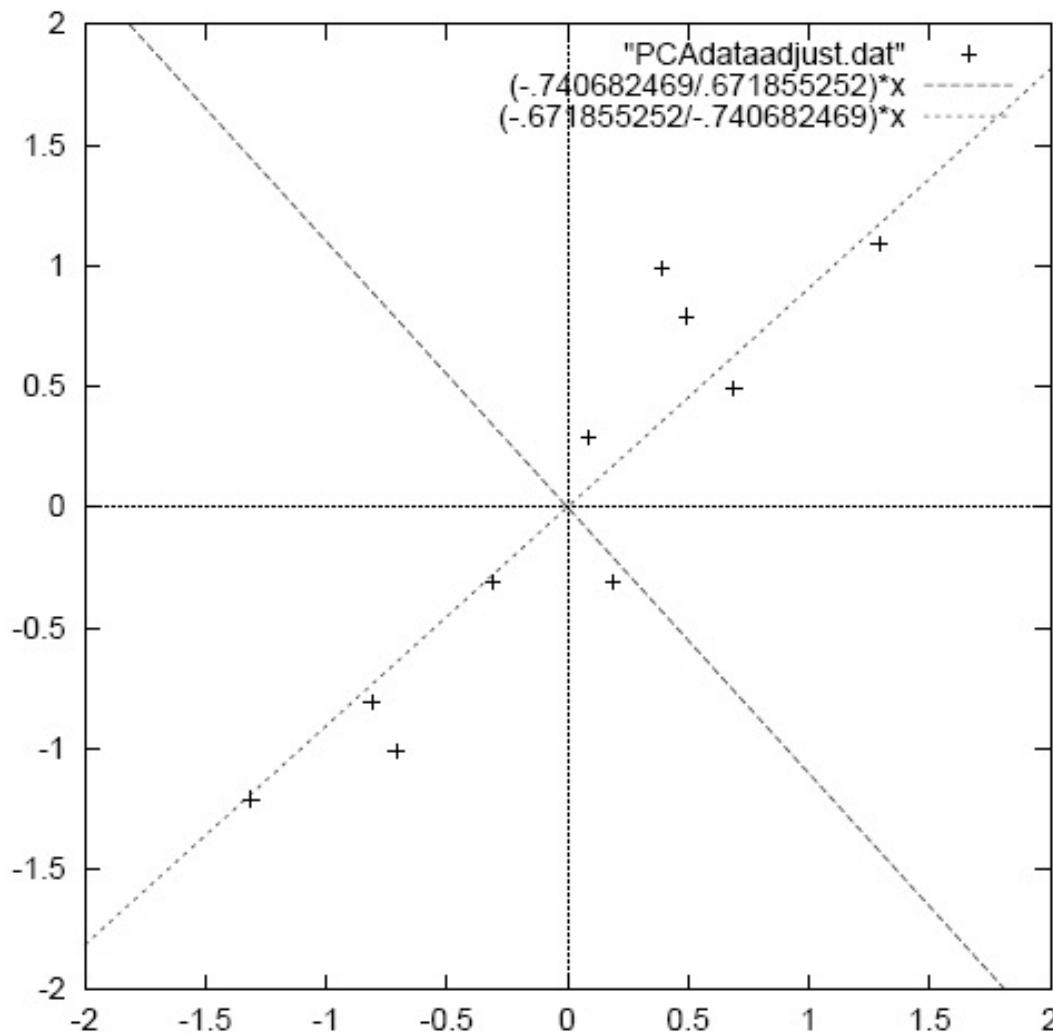
$$cov = \begin{pmatrix} & \mathbf{x} & \mathbf{y} \\ \mathbf{x} & \begin{matrix} .616555556 & .615444444 \\ .615444444 & .716555556 \end{matrix} \\ \mathbf{y} & & \end{pmatrix}$$

3. Calculate the (unit) eigenvectors and eigenvalues of the covariance matrix:

$$eigenvalues = \begin{pmatrix} .0490833989 \\ 1.28402771 \end{pmatrix}$$

$$eigenvectors = \begin{pmatrix} -.735178656 & -.677873399 \\ .677873399 & -.735178656 \end{pmatrix}$$

Mean adjusted data with eigenvectors overlayed



Eigenvector with largest eigenvalue traces linear pattern in data

Figure 3.2: A plot of the normalised data (mean subtracted) with the eigenvectors of the covariance matrix overlayed on top.

4. Order eigenvectors by eigenvalue, highest to lowest.

$$\mathbf{v}_1 = \begin{pmatrix} -.677873399 \\ -.735178956 \end{pmatrix} \quad \lambda = 1.28402771$$

$$\mathbf{v}_2 = \begin{pmatrix} -.735178956 \\ .677873399 \end{pmatrix} \quad \lambda = .0490833989$$

In general, you get  $n$  components. To reduce dimensionality to  $p$ , ignore  $n-p$  components at the bottom of the list.

Construct new “feature vector” (assuming  $\mathbf{v}_i$  is a column vector).

Feature vector =  $(\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_p)$

$$FeatureVector1 = \begin{pmatrix} \mathbf{V}_1 & \mathbf{V}_2 \\ -.677873399 & -.735178956 \\ -.735178956 & .677873399 \end{pmatrix}$$

or reduced dimension feature vector :

$$FeatureVector2 = \begin{pmatrix} -.677873399 \\ -.735178956 \end{pmatrix}$$

5. Derive the new data set.

$$TransformedData = RowFeatureVector \times RowDataAdjust$$

where  $RowDataAdjust = \text{transpose of mean-adjusted data}$

$$RowFeatureVector1 = \begin{pmatrix} -.677873399 & -.735178956 \\ -.735178956 & .677873399 \end{pmatrix}$$

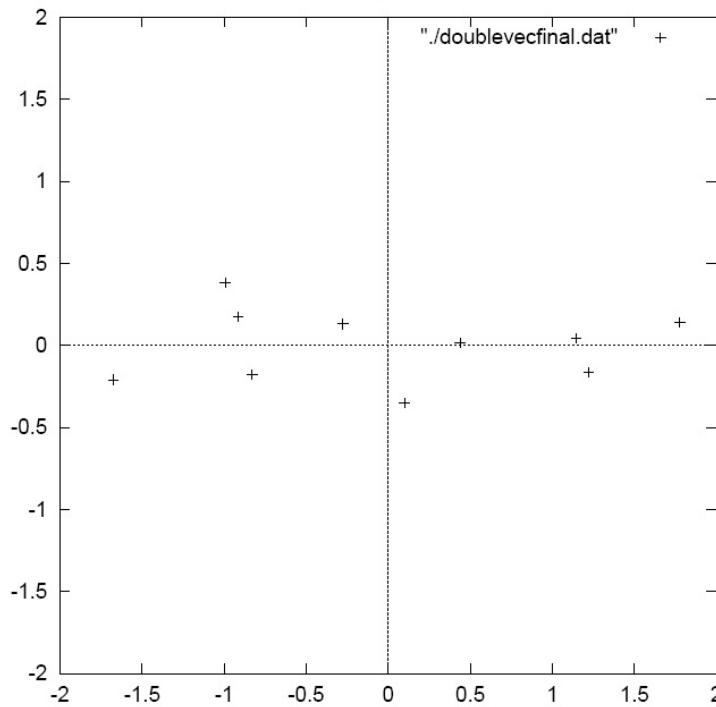
This gives original data in terms of chosen components (eigenvectors)—that is, along these axes.

$$RowFeatureVector2 = (-.677873399 \quad -.735178956)$$

$$RowDataAdjust = \begin{pmatrix} .69 & -1.31 & .39 & .09 & 1.29 & .49 & .19 & -.81 & -.31 & -.71 \\ .49 & -1.21 & .99 & .29 & 1.09 & .79 & -.31 & -.81 & -.31 & -1.01 \end{pmatrix}$$

	<i>x</i>	<i>y</i>
Transformed Data=	-.827970186	-.175115307
	1.77758033	.142857227
	-.992197494	.384374989
	-.274210416	.130417207
	-1.67580142	-.209498461
	-.912949103	.175282444
	.0991094375	-.349824698
	1.14457216	.0464172582
	.438046137	.0177646297
	1.22382056	-.162675287

Data transformed with 2 eigenvectors

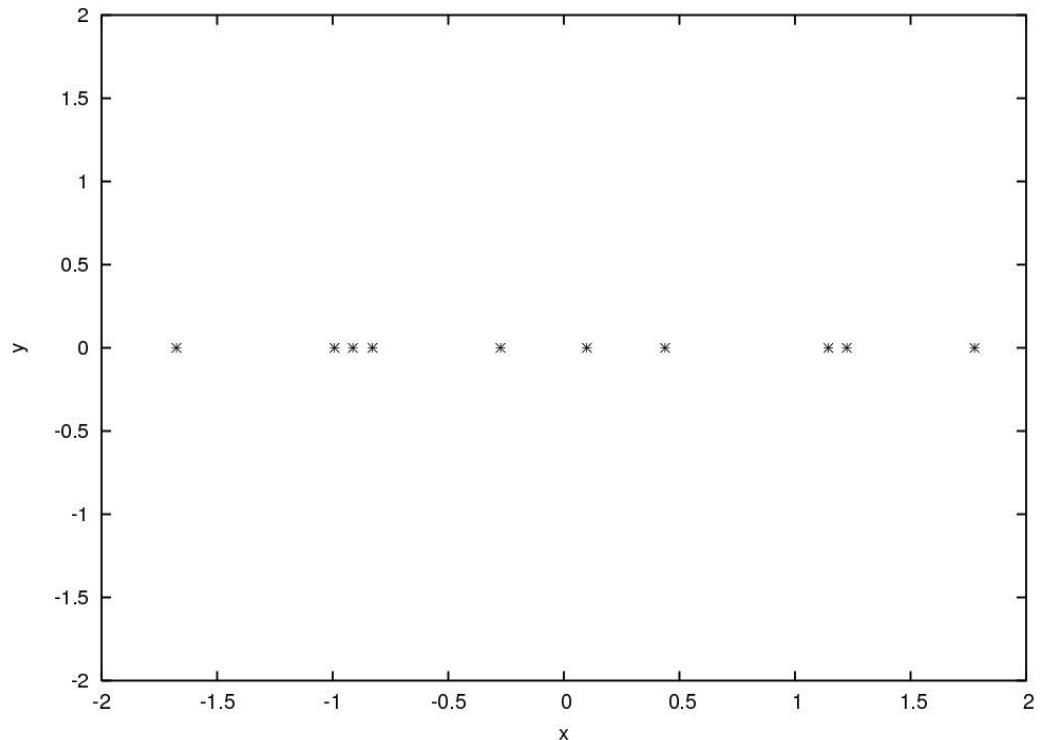


**Intuition:** We projected the data onto new axes that captures the strongest linear trends in the data set. Each transformed data point tells us how far it is above or below those trend lines.

Figure 3.3: The table of data by applying the PCA analysis using both eigenvectors, and a plot of the new data points.

Transformed Data (Single eigenvector)

<u>x</u>
-.827970186
1.77758033
-.992197494
-.274210416
-1.67580142
-.912949103
.0991094375
1.14457216
.438046137
1.22382056



# For Image Compression, How Do we Decompress Data? How to Reconstruct the original data?

We did:

$$\text{TransformedData} = \text{RowFeatureVector} \times \text{RowDataAdjust}$$

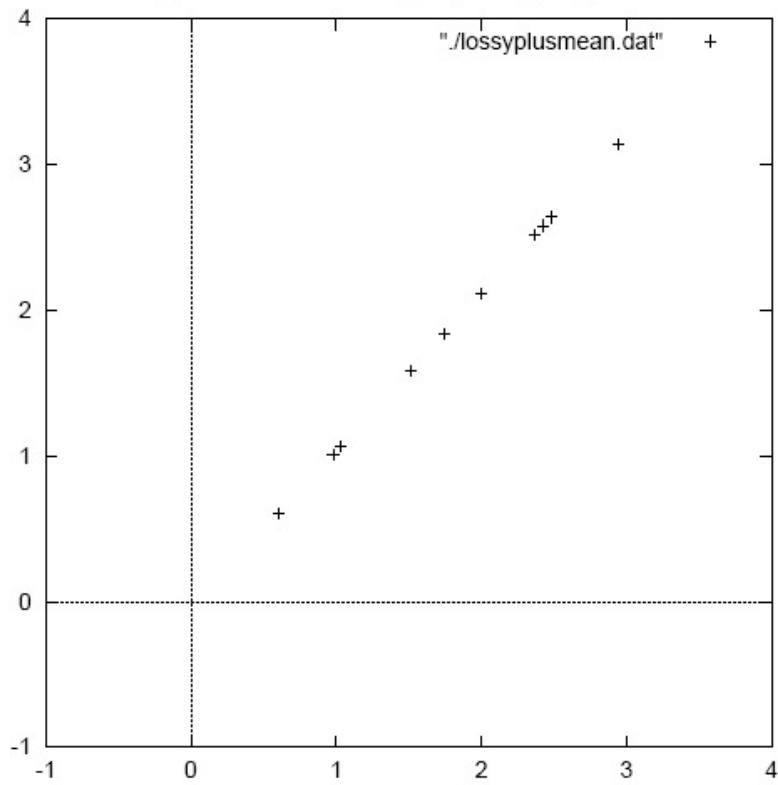
so we can do

$$\begin{aligned}\text{RowDataAdjust} &= \text{RowFeatureVector}^{-1} \times \text{TransformedData} \\ &= \text{RowFeatureVector}^T \times \text{TransformedData}\end{aligned}$$

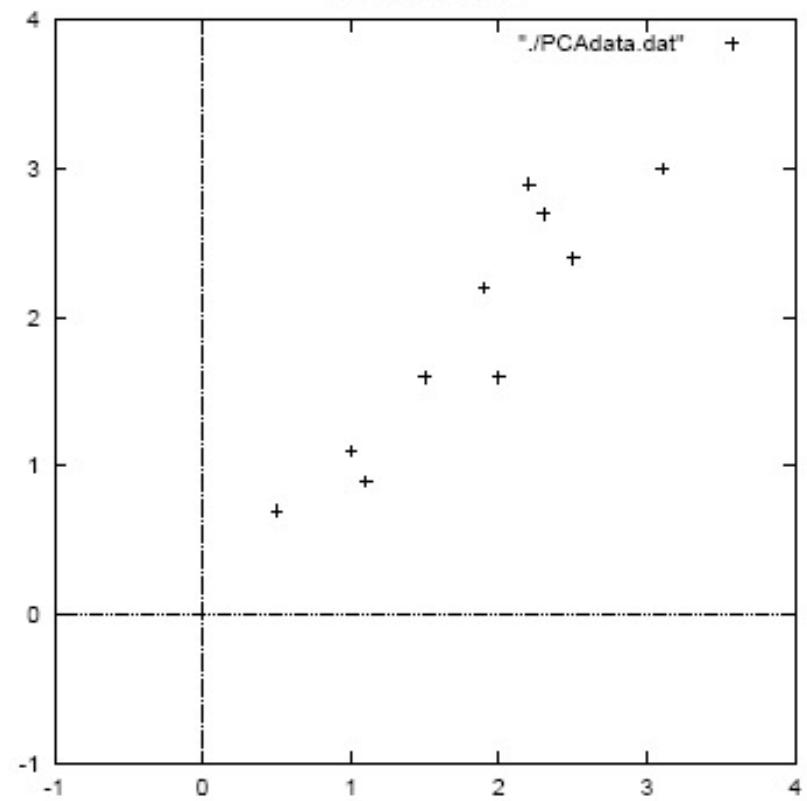
and

$$\text{RowDataOriginal} = \text{RowDataAdjust} + \text{OriginalMean}$$

Original data restored using only a single eigenvector



Original PCA data



Information is lost.  
In Image compression, this is idea of the KL transform.

- ▶ With new data, compute  $TransformedData = RowFeatureVector \times RowDataAdjust$   
where  $RowDataAdjust = transpose\ of\ mean-adjusted\ data$

# What you need to remember

- ▶ General idea of what PCA does
  - Finds new, rotated set of orthogonal axes that capture directions of largest variation
  - Allows some axes to be dropped, so data can be represented in lower-dimensional space.
  - This can improve classification performance and avoid overfitting due to large number of dimensions.
- ▶ You don't need to remember details of PCA algorithm.

# PCA WEKA on Iris data

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize Auto-WEKA Forecast CPython Scripting

Attribute Evaluator

Choose PrincipalComponents -R 0.95-A 5

Search Method

Choose Ranker -T 1.7976931348623157E308-N -1

Attribute Selection Mode

Use full training set  
 Cross-validation Folds 10 Seed 1

(Nom) class

Start Stop

Result list (right-click for options)

11:40:45 - Ranker + PrincipalComponent

Attribute selection output

```
==== Attribute Selection on all input data ===

Search Method:
    Attribute ranking.

Attribute Evaluator (unsupervised):
    Principal Components Attribute Transformer

Correlation matrix
  1   -0.11   0.87   0.82
 -0.11   1   -0.42  -0.36
  0.87  -0.42   1   0.96
  0.82  -0.36   0.96   1

eigenvalue      proportion      cumulative
  2.91082        0.7277        0.7277      -0.581petallength-0.566petalwidth-0.522sepallength+0.263sepalwidth
  0.92122        0.23031       0.95801      0.926sepalwidth+0.372sepallength+0.065petalwidth+0.021petallength

Eigenvectors
  V1      V2
-0.5224  0.3723 sepallength
  0.2634  0.9256 sepalwidth
-0.5813  0.0211 petallength
-0.5656  0.0654 petalwidth

Ranked attributes:
  0.2723  1 -0.581petallength-0.566petalwidth-0.522sepallength+0.263sepalwidth
  0.042   2  0.926sepalwidth+0.372sepallength+0.065petalwidth+0.021petallength

Selected attributes: 1,2 : 2
```

Watch: <https://www.youtube.com/watch?v=zuoMjUAPihA>

# Example: Linear discrimination using PCA for face recognition (“Eigenfaces”)

## 1. Preprocessing: “Normalize” faces

- Make images the same size
- Line up with respect to eyes
- Normalize intensities



2. Raw features are pixel intensity values (2061 features)
3. Each image is encoded as a vector  $\Gamma_i$  of these features
4. Compute “mean” face in training set:

$$\Psi = \frac{1}{M} \sum_{i=1}^M \Gamma_i$$

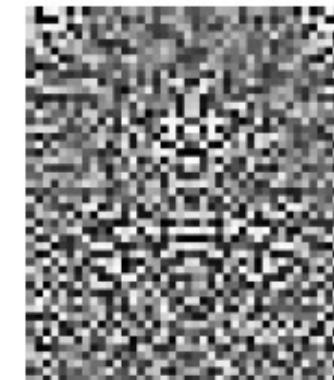
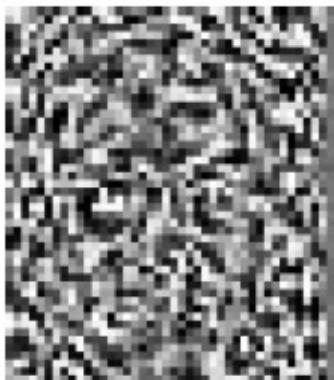


The average face and first four eigenfaces

- ▶ Subtract the mean face from each face vector
$$\Phi_i = \Gamma_i - \Psi$$
- ▶ Compute the covariance matrix  $\mathbf{C}$
- ▶ Compute the (unit) eigenvectors  $\mathbf{v}_i$  of  $\mathbf{C}$
- ▶ Keep only the first  $K$  principal components (eigenvectors)



Eigenfaces 15, 100, 200, 250, 300



Eigenfaces 400, 450, 1000, 2000

# Interpreting and Using Eigenfaces

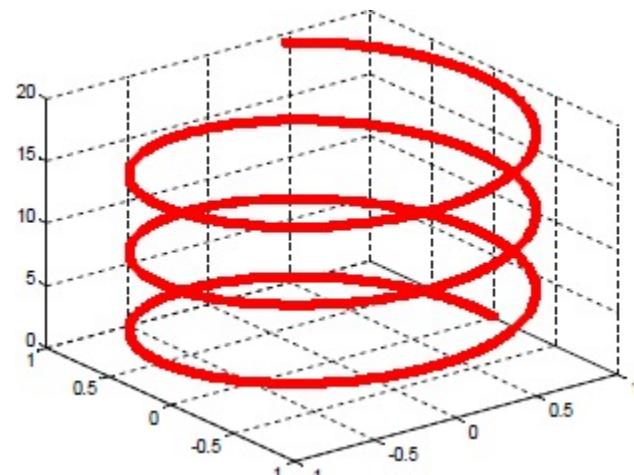
The eigenfaces encode the principal sources of variation in the dataset (e.g., absence/presence of facial hair, skin tone, glasses, etc.).

We can represent any face as a linear combination of these “basis” faces.

Use this representation for:

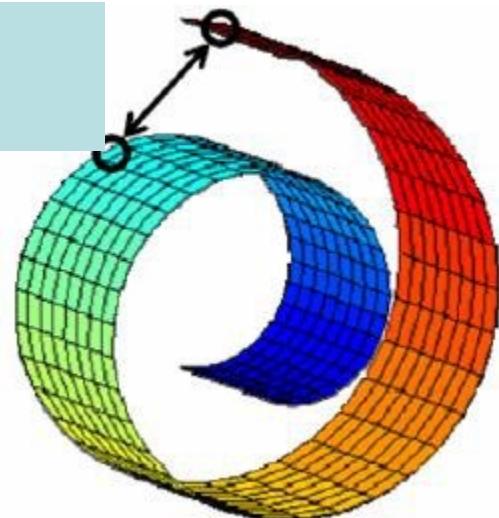
- Face recognition
  - (e.g., Euclidean distance from known faces)
- Linear discrimination
  - (e.g., “glasses” versus “no glasses”, or “male” versus “female”)

- The PCA can not discover the structure of a data set in a spiral form

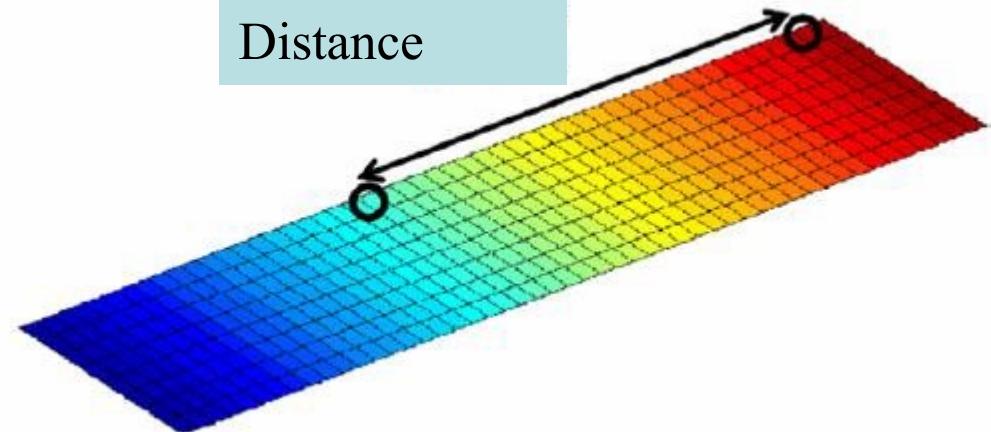


# Euclidian Distance vs. Geodesic Distance

Euclidian  
Distance



Geodesic  
Distance



# Possible solutions

- Kernel PCA
- ISOMAP

# Kernel PCA

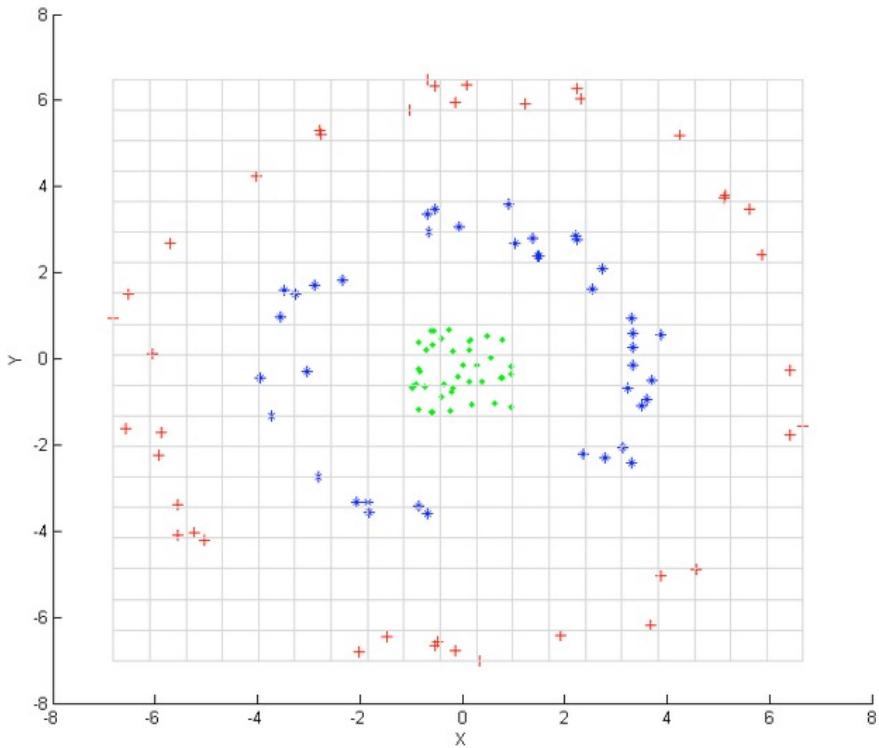
---

## The Kernel PCA Algorithm

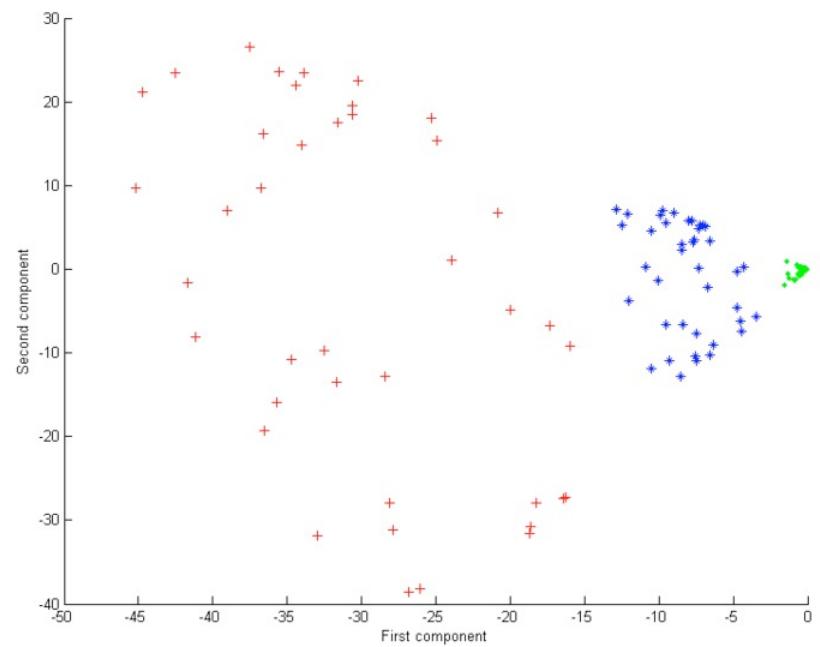
---

- Choose a kernel and apply it to all pairs of points to get matrix  $\mathbf{K}$  of distances between the pairs of points in the transformed space
  - Compute the eigenvalues and eigenvectors of  $\mathbf{K}$
  - Normalise the eigenvectors by the square root of the eigenvalues
  - Retain the eigenvectors corresponding to the largest eigenvalues
-

Original data



Data after kernel PCA



# Kernel Eigenfaces

(Yang et al., Face Recognition Using Kernel Eigenfaces, 2000)

Training data:  $\sim 400$  images, 40 subjects

Original features: 644 pixel gray-scale values.

Transform data using kernel PCA, reduce dimensionality to number of components giving lowest error.

Test: new photo of one of the subjects

Recognition done using nearest neighbor classification

Table 3: Experimental results on AT&T database

Method	Reduced space	Error Rate (%)
Eigenface	30	2.75
Kernel PCA, d=2	50	2.50
Kernel PCA, d=3	50	2.00
Kernel PCA, d=4	60	2.25
Kernel PCA, d=10	80	2.25

# Independent Component Analysis (ICA)

- ▶ Independent component analysis attempts to decompose a multivariate signal into independent non-Gaussian signals.
- ▶ As an example, sound is usually a signal that is composed of the numerical addition, at each time  $t$ , of signals from several sources.
- ▶ ICA finds the independent components by maximizing the statistical independence of the estimated components.

# Isometric Feature Mapping (ISOMAP)

- The goal of ISOMAP is to find a non-linear manifold containing the data
- We use the fact that for close points, the Euclidean distance is a good approximation to the geodesic distance on the manifold
- We build a graph connecting each point to its  $k$  nearest neighbors

# ISOMAP

- The lengths of the geodesics are then estimated by searching the length of the shortest path between two points in the graph
- Thereafter, we compute distances to determine a position of points in a space of reduced dimensions

# ISOMAP

- ISOMAP [Tenebaum *et al.*, 2000]
  - For neighboring samples, Euclidian distance provides a good approximation to geodesic distance
  - For distant points, geodesic distance can be approximated with a sequence of steps between clusters of neighboring points

# ISOMAP

ISOMAP operates in three steps:

1. Build the neighborhood graph  $G$
2. For each pair of points of  $G$ , compute the shortest path (the geodesic distance)
3. Using the classical MDS (Multidimensional Scaling) on geodesic distances

Euclidian Distance is replaced by Geodesic Distance

# ISOMAP Algorithm

- Step 1
  - Build the neighborhood graph based on the distances  $d_X(i,j)$  in the input space X.
  - This can be performed in two different ways:
    - Connect each point to all points within a fixed radius  $\varepsilon$
    - Connect each point to all of its  $k$  nearest neighbors
  - A weighted neighborhood graph  $G$  is obtained, where  $d_X(i,j)$  is the weight of each edge between neighboring points

# ISOMAP Algorithm

- Step 2
  - Estimate the geodesic distances  $d_M(i,j)$  between all pair of points on the manifold  $M$  by computing their shortest path distances  $d_G(i,j)$  in the graph  $G$ .
  - It can be done using Dijkstra's or Floyd's algorithm.

# ISOMAP Algorithm

- Step 3
  - Apply classical MDS to the matrix of graph distances D.

*Multidimensional scaling (MDS) is a means of visualizing the level of similarity of individual vectors of a dataset. MDS is used to translate "information about the pairwise distances into a configuration of points mapped into a Cartesian space.*

# Complexity of ISOMAP

- For large datasets ISOMAP can be quite slow :
  - Step 1: Complexity of k-nearest neighbors  $O(n^2 D)$
  - Step 2 : Complexity of the Djikstra algorithm  $O(n^2 \log n + n^2 k)$
  - Step 3 : MDS complexity  $O(n^2 d)$