

# Unsupervised Learning - Dimensionality reduction.

- A very brief introduction to unsupervised learning
- Dimensionality reduction:
  - Principal Component Analysis (PCA)
  - Kernelizing PCA
  - Other non-linear dimensionality reduction techniques

## Unsupervised Learning (cont'd).

- Until now, we focused on **supervised learning** (e.g. regression and classification):
  - We have access to **labeled data**: we observe both features for each object  $\mathbf{x}_1, \dots, \mathbf{x}_m$  and the corresponding target variable  $y_1, \dots, y_m$ . The goal is to *predict the target from the features*.
- In **unsupervised learning**, we only have access to **unlabeled data**, i.e. the features  $\mathbf{x}_1, \dots, \mathbf{x}_m$ , and we want to *discover interesting things about the data* (i.e. identify the underlying structure of the data).

## Unsupervised Learning (cont'd).

- Unsupervised learning is more subjective than supervised learning (no simple goal, such as prediction in supervised learning).
- But of growing importance:
  - Easier to obtain unlabeled data than labeled data.
  - One of the next challenges of ML:
    - \* Overwhelming amount of unlabeled data available
    - \* Human don't need so many labeled examples to learn: few labeled examples and large amount of unlabeled data is often enough (semi-supervised learning)...

# Unsupervised Learning.

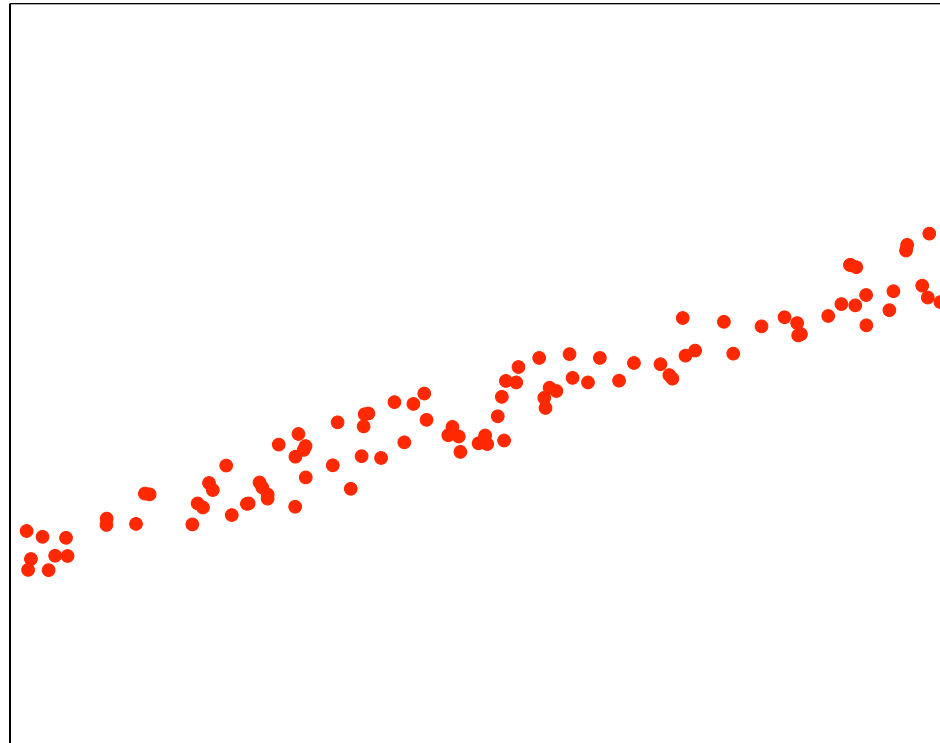
- Examples of unsupervised learning tasks:
  - **Density estimation**: estimate the distribution of the data (e.g. fit a Gaussian to the data, or a mixture of Gaussians). Closely related to generative models (e.g. GANs)...
  - **Clustering**: identify groups of similar objects (for e.g. market segmentation, data exploration). For example  $k$ -means is a very popular clustering algorithm:.
  - **Dimensionality reduction**: find a low-dimensional representation of the data (e.g. to reduce the complexity of learning algorithm, data visualization).

⇒ We focus on dimensionality reduction in this lecture.

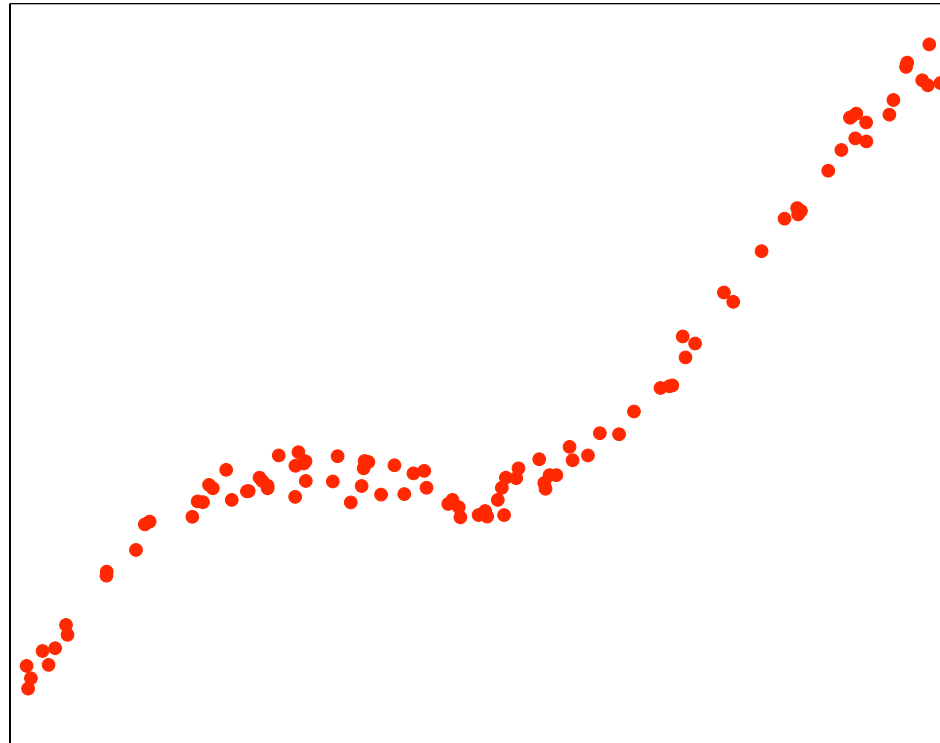
# What is dimensionality reduction?

- Dimensionality reduction (or embedding) techniques:
  - Assign instances to real-valued vectors, in a space that is much smaller-dimensional (even 2D or 3D for visualization).
  - Approximately preserve similarity/distance relationships between instances.
- Some techniques:
  - Linear: Principal components analysis
  - Non-linear
    - \* Kernel PCA
    - \* Independent components analysis
    - \* Self-organizing maps
    - \* Locally linear embeddings
    - \* Multi-dimensional scaling
    - \* Autoencoders
    - \* ...

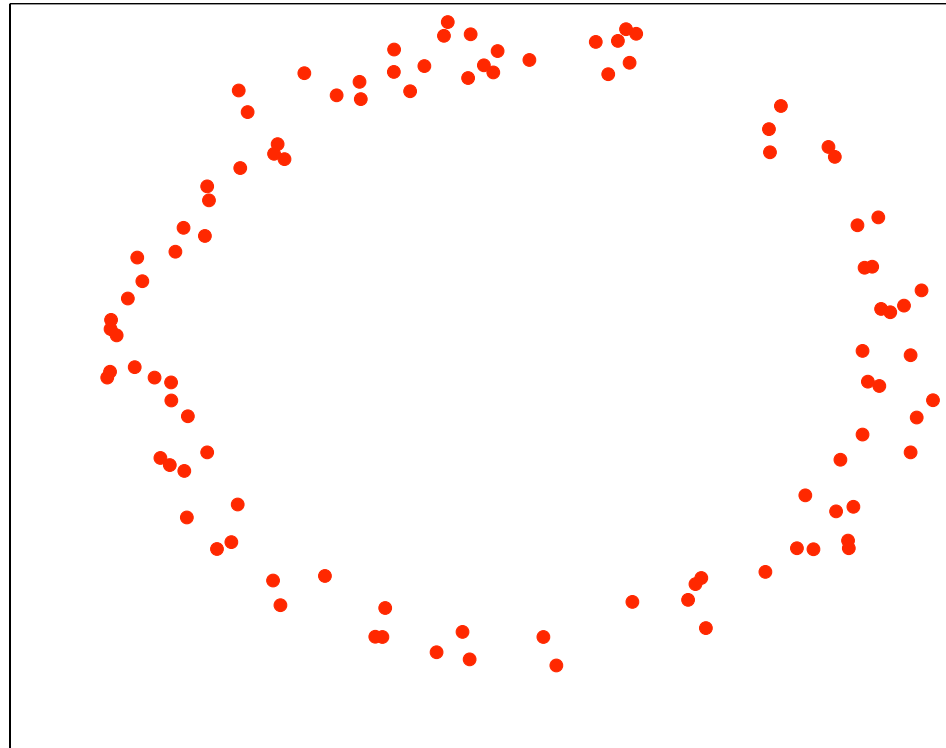
**What is the true dimensionality of this data?**



**What is the true dimensionality of this data?**

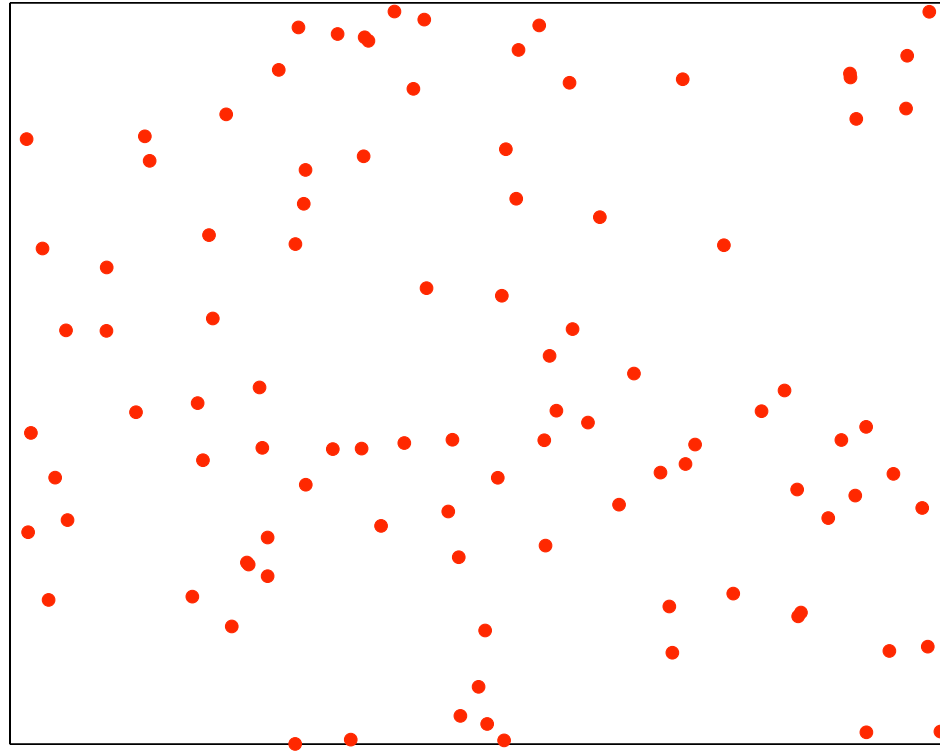


**What is the true dimensionality of this data?**





**What is the true dimensionality of this data?**

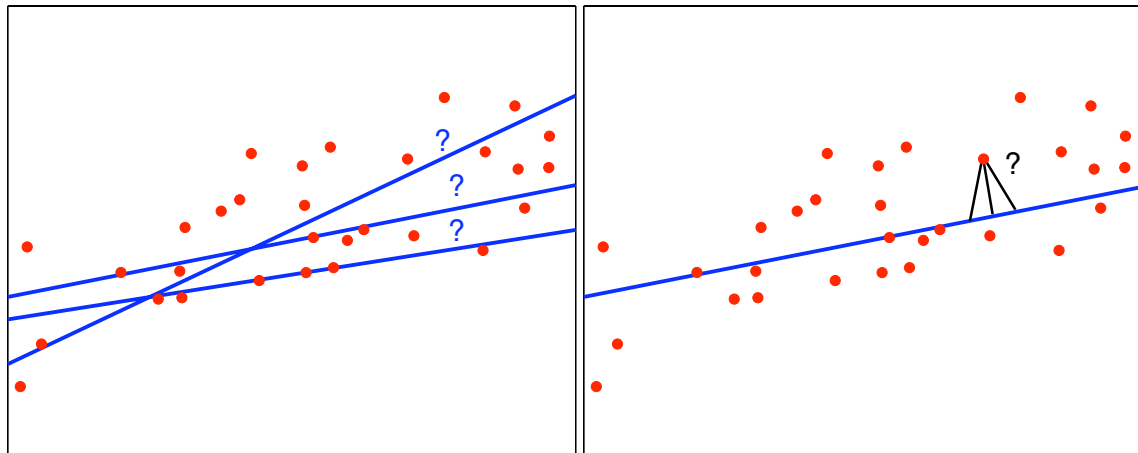


## Remarks

- All dimensionality reduction techniques are based on an implicit assumption that the data lies along some *low-dimensional manifold*
- This is the case for the first three examples, which lie along a 1-dimensional manifold despite being plotted in 2D
- In the last example, the data has been generated randomly in 2D, so no dimensionality reduction is possible without losing information
- The first three cases are in increasing order of difficulty, from the point of view of existing techniques.

# Simple Principal Component Analysis (PCA)

- Given:  $m$  instances, each being a length- $n$  real vector.
- Suppose we want a 1-dimensional representation of that data, instead of  $n$ -dimensional.
- Specifically, we will:
  - Choose a line in  $\mathbb{R}^n$  that “best represents” the data.
  - Assign each data object to a point along that line.



## Reconstruction error

- Let the line be represented as  $\mathbf{b} + \alpha \mathbf{v}$  for  $\mathbf{b}, \mathbf{v} \in \mathbb{R}^n$ ,  $\alpha \in \mathbb{R}$ .  
For convenience assume  $\|\mathbf{v}\| = 1$ .
- Each instance  $\mathbf{x}_i$  is associated with a point on the line  $\hat{\mathbf{x}}_i = \mathbf{b} + \alpha_i \mathbf{v}$ .
- We want to choose  $\mathbf{b}$ ,  $\mathbf{v}$ , and the  $\alpha_i$  to minimize the total reconstruction error over all data points, measured using Euclidean distance:

$$R = \sum_{i=1}^m \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2$$

## A constrained optimization problem!

$$\begin{array}{ll}\min & \sum_{i=1}^m \|\mathbf{x}_i - (\mathbf{b} + \alpha_i \mathbf{v})\|^2 \\ \text{w.r.t.} & \mathbf{b}, \mathbf{v}, \alpha_i, i = 1, \dots, m \\ \text{s.t.} & \|\mathbf{v}\|^2 = 1\end{array}$$

- This is a quadratic objective with quadratic constraint
- Suppose we fix a  $\mathbf{v}$  satisfying the condition, and find the best  $\mathbf{b}$  and  $\alpha_i$  given this  $\mathbf{v}$
- So, we solve:

$$\min R = \min_{\alpha, \mathbf{b}} \sum_{i=1}^m \|\mathbf{x}_i - (\mathbf{b} + \alpha_i \mathbf{v})\|^2$$

where  $R$  is the *reconstruction error*

## Solving the optimization problem (II)

- We write the gradient of  $R$  wrt to  $\alpha_i$  and set it to 0:

$$\frac{\partial R}{\partial \alpha_i} = 2\|\mathbf{v}\|^2 \alpha_i - 2\mathbf{v}\mathbf{x}_i + 2\mathbf{b}\mathbf{v} = 0 \Rightarrow \alpha_i = \mathbf{v} \cdot (\mathbf{x}_i - \mathbf{b})$$

where we used  $\|\mathbf{v}\|^2 = 1$ .

- We write the gradient of  $R$  wrt  $\mathbf{b}$  and set it to 0:

$$\nabla_{\mathbf{b}} R = 2m\mathbf{b} - 2 \sum_{i=1}^m \mathbf{x}_i + 2 \left( \sum_{i=1}^m \alpha_i \right) \mathbf{v} = 0 \Rightarrow m\mathbf{b} = \sum_{i=1}^m \mathbf{x}_i - \sum_{i=1}^m \alpha_i \mathbf{v}$$

- From above:

$$\sum_{i=1}^m \alpha_i \mathbf{v} = \left( \sum_{i=1}^m \mathbf{v}^\top (\mathbf{x}_i - \mathbf{b}) \right) \mathbf{v} = \mathbf{v}\mathbf{v}^\top \left( \sum_{i=1}^m \mathbf{x}_i - m\mathbf{b} \right)$$

## Solving the optimization problem (III)

- Combining the previous two equations we get:

$$(\mathbf{I} - \mathbf{v}\mathbf{v}^\top)m\mathbf{b} = (\mathbf{I} - \mathbf{v}\mathbf{v}^\top) \sum_{i=1}^m \mathbf{x}_i$$

- This is satisfied when:

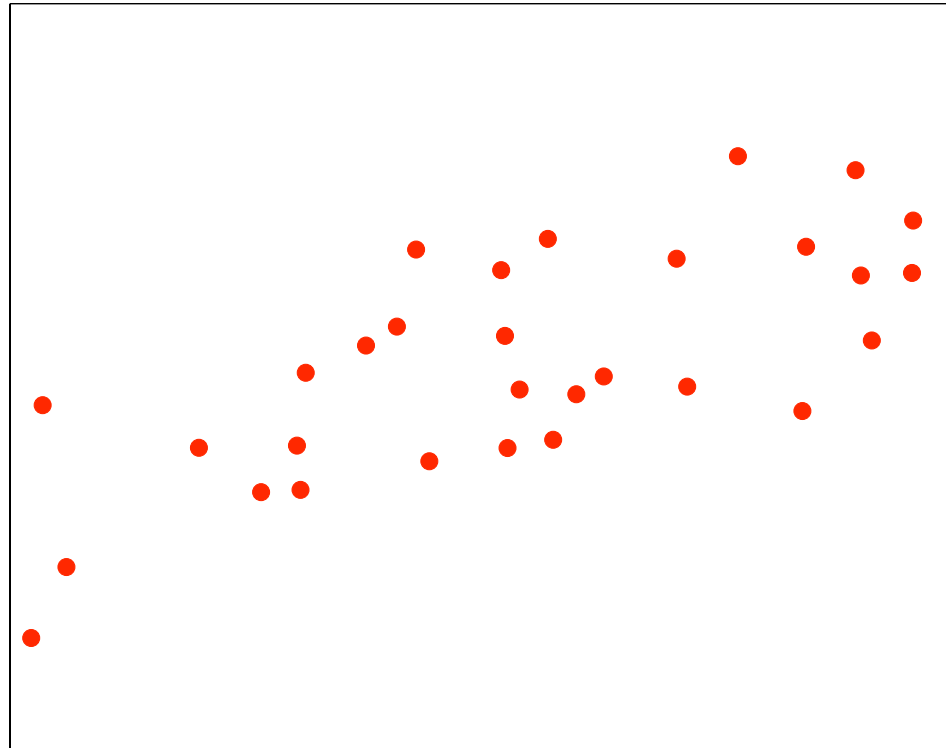
$$\mathbf{b} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$$

- This means that the line goes through the mean of the data
- By substituting  $\alpha_i = \mathbf{v}^\top(\mathbf{x}_i - \mathbf{b})$ , we get:

$$\hat{\mathbf{x}}_i = \mathbf{b} + \alpha_i \mathbf{v} = \mathbf{b} + \mathbf{v}\mathbf{v}^\top(\mathbf{x}_i - \mathbf{b})$$

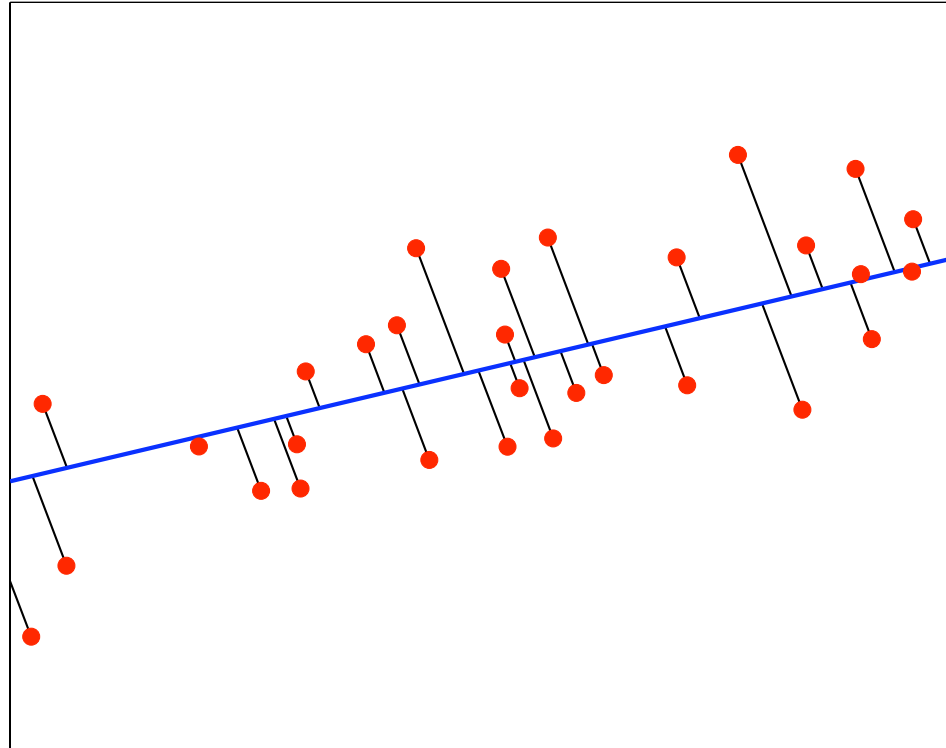
- This means that instances are projected orthogonally on the line to get the associated point.

## Example data





## Example with $\mathbf{v} \propto (1, 0.3)$



## Finding the direction of the line

- Substituting  $\hat{\mathbf{x}}_i = \mathbf{b} + \mathbf{v}\mathbf{v}^\top(\mathbf{x}_i - \mathbf{b})$ , we want to solve:

$$\min_{\mathbf{v}} \sum_{i=1}^m \|(\mathbf{I} - \mathbf{v}\mathbf{v}^\top)(\mathbf{x}_i - \mathbf{b})\|^2 \quad \text{s.t.} \quad \|\mathbf{v}\|^2 = 1$$

- Using the fact that  $\|(\mathbf{I} - \mathbf{v}\mathbf{v}^\top)(\mathbf{x}_i - \mathbf{b})\|^2 = \|\mathbf{x}_i - \mathbf{b}\|^2 - \|\mathbf{v}\mathbf{v}^\top(\mathbf{x}_i - \mathbf{b})\|^2$  (since  $\|\mathbf{v}\|^2 = 1$ ) this is equivalent to

$$\max_{\mathbf{v}} \sum_{i=1}^m \|\mathbf{v}\mathbf{v}^\top(\mathbf{x}_i - \mathbf{b})\|^2 \quad \text{s.t.} \quad \|\mathbf{v}\|^2 = 1$$

which (using  $\|\mathbf{v}\mathbf{v}^\top(\mathbf{x}_i - \mathbf{b})\|^2 = (\mathbf{v}^\top(\mathbf{x}_i - \mathbf{b}))^2$ ) can be rewritten into

$$\max_{\mathbf{v}} \sum_{i=1}^m \mathbf{v}^\top(\mathbf{x}_i - \mathbf{b})(\mathbf{x}_i - \mathbf{b})^\top \mathbf{v} \quad \text{s.t.} \quad \|\mathbf{v}\|^2 = 1$$

## Finding the direction of the line (cont'd)

- We want to solve

$$\max_{\mathbf{v}} \sum_{i=1}^m \mathbf{v}^\top (\mathbf{x}_i - \mathbf{b})(\mathbf{x}_i - \mathbf{b})^\top \mathbf{v} \quad \text{s.t.} \quad \|\mathbf{v}\|^2 = 1$$

- The Lagrangian is:

$$L(\mathbf{v}, \lambda) = \mathbf{v}^\top \left( \sum_{i=1}^m (\mathbf{x}_i - \mathbf{b})(\mathbf{x}_i - \mathbf{b})^\top \right) \mathbf{v} + \lambda - \lambda \|\mathbf{v}\|^2$$

- Let  $\mathbf{S} = \sum_{i=1}^m (\mathbf{x}_i - \mathbf{b})(\mathbf{x}_i - \mathbf{b})^\top$  be an  $n$ -by- $n$  matrix, which we will call the *scatter matrix*
- Setting  $\nabla_{\mathbf{v}} L = 0$ , the solution of the problem must satisfy

$$\mathbf{S}\mathbf{v} = \lambda\mathbf{v}$$

## Optimal choice of $\mathbf{v}$

- Recall: an *eigenvector*  $\mathbf{u}$  of a matrix  $\mathbf{A}$  satisfies  $\mathbf{A}\mathbf{u} = \lambda\mathbf{u}$ , where  $\lambda \in \mathbb{R}$  is the *eigenvalue*.
- Fact: the scatter matrix,  $\mathbf{S}$ , has  $n$  non-negative eigenvalues and  $n$  orthogonal eigenvectors.
- The equation obtained for  $\mathbf{v}$  tells us that it should be an eigenvector of  $\mathbf{S}$ .
- The  $\mathbf{v}$  that maximizes  $\mathbf{v}^\top \mathbf{S} \mathbf{v}$  is the eigenvector of  $\mathbf{S}$  with the largest eigenvalue

## What is the scatter matrix

- $\mathbf{S}$  is an  $n \times n$  matrix with

$$\mathbf{S}_{k,l} = \sum_{i=1}^m (\mathbf{x}_i(k) - \mathbf{b}(k))(\mathbf{x}_i(l) - \mathbf{b}(l))$$

- Hence,  $\mathbf{S}_{k,l}$  is proportional to the *estimated covariance* between the  $k$ th and  $l$ th dimension in the data.

## Recall: Covariance

- Covariance quantifies a *linear relationship* (if any) between two random variables  $X$  and  $Y$ .

$$\text{Cov}(X, Y) = E\{(X - E(X))(Y - E(Y))\}$$

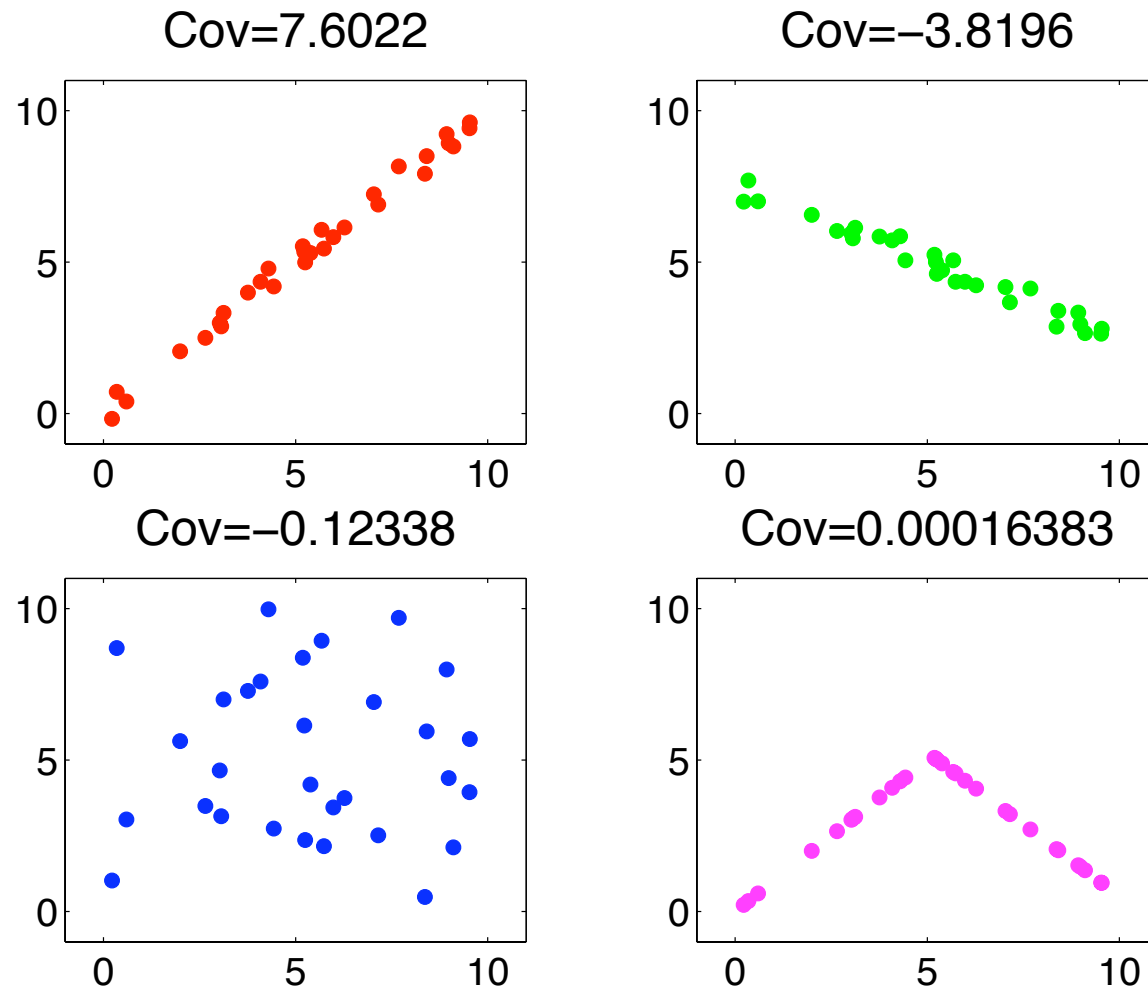
- Given  $m$  samples of  $X$  and  $Y$ , covariance can be estimated as

$$\frac{1}{m} \sum_{i=1}^m (x_i - \mu_X)(y_i - \mu_Y) ,$$

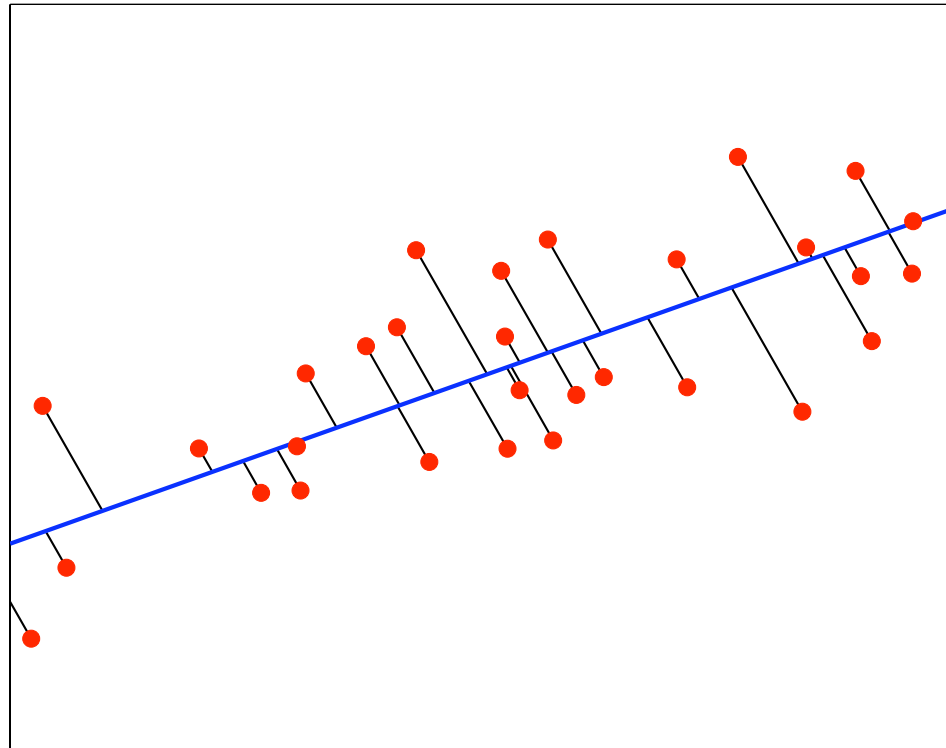
where  $\mu_X = (1/m) \sum_{i=1}^m x_i$  and  $\mu_Y = (1/m) \sum_{i=1}^m y_i$ .

- Note:  $\text{Cov}(X, X) = \text{Var}(X)$ .

## Covariance example



**Example with optimal line:  $\mathbf{b} = (0.54, 0.52)$ ,  $\mathbf{v} \propto (1, 0.45)$**





## Remarks

- The line  $\mathbf{b} + \alpha \mathbf{v}$  is the *first principal component*.
- The variance of the data along the line  $\mathbf{b} + \alpha \mathbf{v}$  is as large as along any other line.
- $\mathbf{b}$ ,  $\mathbf{v}$ , and the  $\alpha_i$  can be computed easily in polynomial time.

## Generalization to $d$ dimensions

- More generally, we can create a  $d$ -dimensional representation of our data by projecting the instances onto a hyperplane  $\mathbf{b} + \alpha^1 \mathbf{v}_1 + \dots + \alpha^d \mathbf{v}_d$ .
- If we assume the  $\mathbf{v}_j$  are of unit length and orthogonal, then the optimal choices are:
  - $\mathbf{b}$  is the mean of the data (as before)
  - The  $\mathbf{v}_j$  are orthogonal eigenvectors of  $\mathbf{S}$  corresponding to its  $d$  largest eigenvalues.
  - Each instance is projected orthogonally on the hyperplane.

## PCA: overall algorithm

1. Center the data  $\tilde{\mathbf{x}}_i = \mathbf{x}_i - \mathbf{b}$  where  $\mathbf{b} = \frac{1}{m} \sum_i \mathbf{x}_i$ .
2. (Optional step: normalize the data.)
3. Compute the top  $d$  (unit-norm) eigenvectors of  $\mathbf{S} = \sum_i \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top$ .  
(observe that  $\mathbf{S} = \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}$  where  $\tilde{\mathbf{X}} \in \mathbb{R}^{m \times n}$  is the centered data matrix)
4. Put these eigenvectors into a matrix  $\mathbf{U} \in \mathbb{R}^{n \times d}$ .
5. The PCA projection of any point  $\mathbf{x}$  is given by
  - \*  $\mathbf{U}^\top (\mathbf{x} - \mathbf{b}) \in \mathbb{R}^d$  in the latent space.
  - \*  $\mathbf{b} + \mathbf{U} \mathbf{U}^\top (\mathbf{x} - \mathbf{b}) \in \mathbb{R}^n$  in the ambient space.

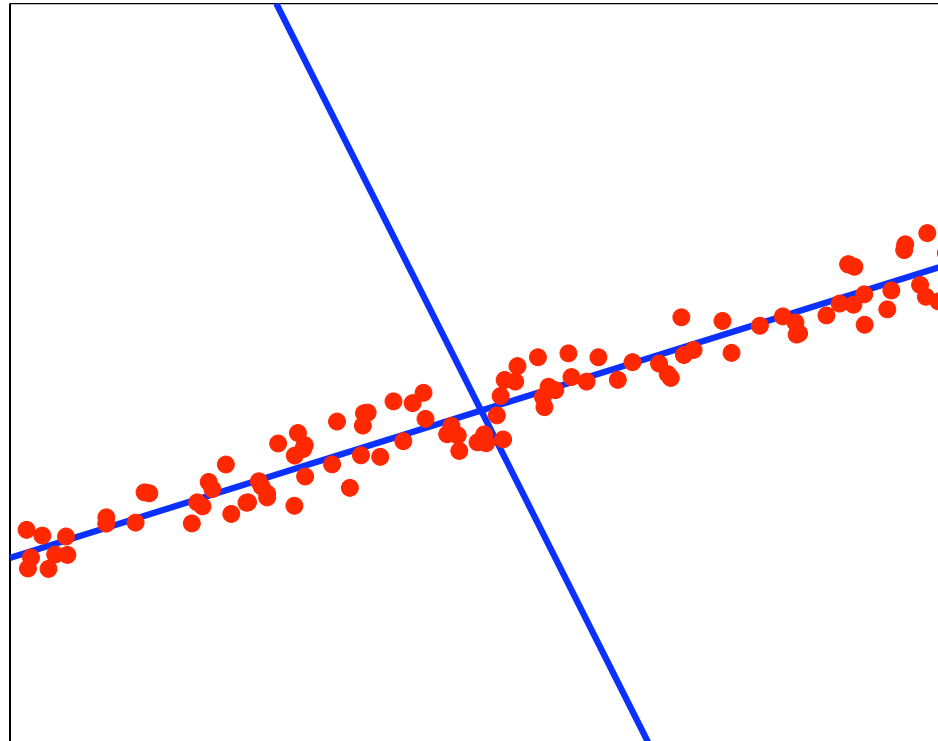
## Remarks

- $\mathbf{b}$ , the eigenvalues, the  $\mathbf{v}_j$ , and the projections of the instances can all be computed in polynomial time.
- The magnitude of the  $j^{th}$ -largest eigenvalue,  $\lambda_j$ , tells you how much variability in the data is captured by the  $j^{th}$  principal component
- So you have feedback on how to choose  $d$ !
- When the eigenvalues are sorted in decreasing order, the proportion of the variance captured by the first  $d$  components is:

$$\frac{\lambda_1 + \cdots + \lambda_d}{\lambda_1 + \cdots + \lambda_d + \lambda_{d+1} + \cdots + \lambda_n}$$

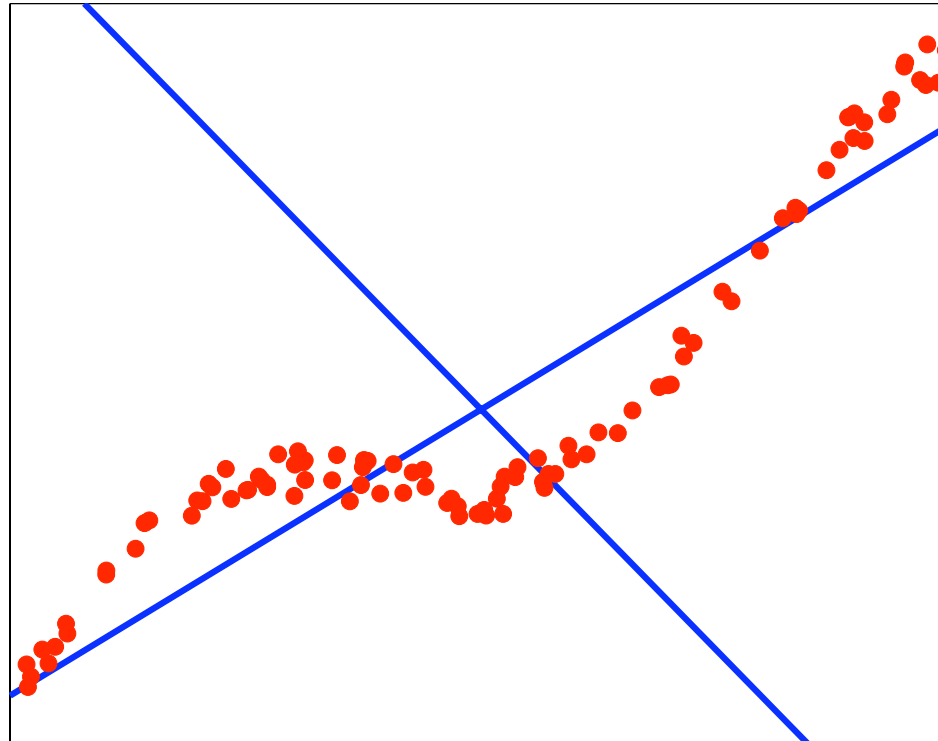
- So if a “big” drop occurs in the eigenvalues at some point, that suggests a good dimension cutoff

**Example:**  $\lambda_1 = 0.0938$ ,  $\lambda_2 = 0.0007$



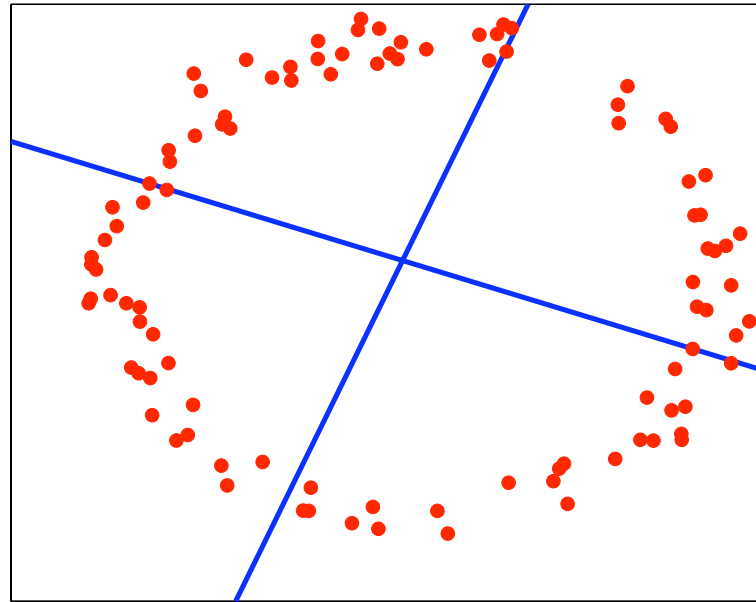
The first eigenvalue accounts for most variance, so the dimensionality is 1

**Example:**  $\lambda_1 = 0.1260$ ,  $\lambda_2 = 0.0054$



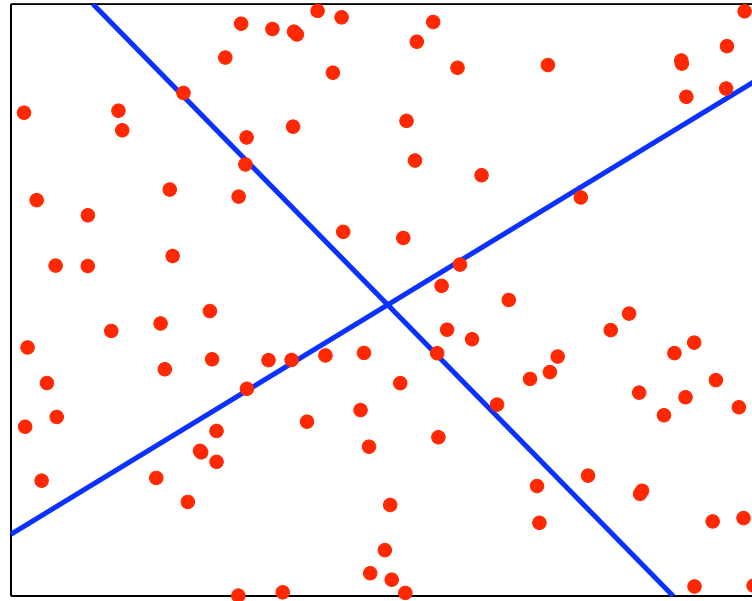
The first eigenvalue accounts for most variance, so the dimensionality is 1  
(despite some non-linear structure in the data)

**Example:**  $\lambda_1 = 0.0884, \lambda_2 = 0.0725$



- Each eigenvalue accounts for about half the variance, so the PCA-suggested dimension is 2
- Note that this is the *linear* dimension
- The true “non-linear” dimension of the data is 1 (using polar coordinates)

**Example:**  $\lambda_1 = 0.0881, \lambda_2 = 0.0769$



- Each eigenvalue accounts for about half the variance, so the PCA-suggested dimension is 2
- In this case, the non-linear dimension is also 2 (data is fully random)
- Note that *PCA cannot distinguish non-linear structure from no structure*
- This case and the previous one yield a very similar PCA analysis



## Remarks

- Outliers have a big effect on the covariance matrix, so they can affect the eigenvectors quite a bit
- A simple examination of the pairwise distances between instances can help discard points that are very far away (for the purpose of PCA)
- If the variances in the original dimensions vary considerably, they can “muddle” the true correlations. There are two solutions:
  - Work with the correlation (covariance rescaled to  $(-1, 1)$ ) of the original data, instead of covariance matrix (which provides one type of normalization)
  - Normalize the input dimensions individually (possibly based on domain knowledge) before PCA
- PCA is most often performed using Singular Value Decomposition (SVD)
- In certain cases, the eigenvectors are meaningful; e.g. in vision, they can be displayed as images (“eigenfaces”)

## Eigenfaces example

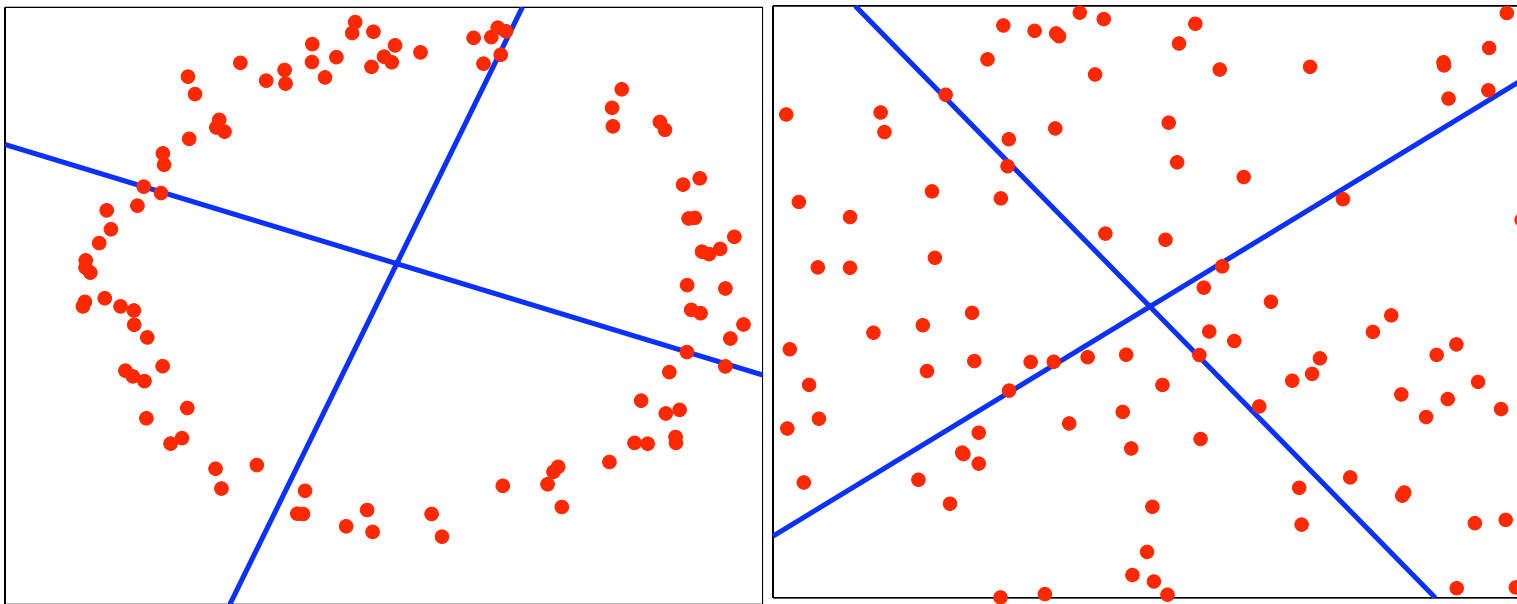


- A set of faces on the left and the corresponding eigenfaces (principal components) on the right
- Note that faces have to be centred and scaled ahead of time
- The components are in the same space as the instances (images) and can be used to reconstruct the images

## Uses of PCA

- Pre-processing for a supervised learning algorithm, e.g. for image data, robotic sensor data
- Used with great success in image and speech processing
- Visualization
- Exploratory data analysis
- Removing the linear component of a signal (before fancier non-linear models are applied)

## Difficult example



- PCA will make no difference between these examples, because the structure on the left is not linear
- Are there ways to find non-linear, low-dimensional manifolds?

## Making PCA non-linear

- Suppose that instead of using the points  $\mathbf{x}_i$  as is, we wanted to go to some different *feature space*  $\phi(\mathbf{x}_i) \in \mathbb{R}^N$
- E.g. using polar coordinates instead of cartesian coordinates would help us deal with the circle
- In the higher dimensional space, we can then do PCA
- The result will be non-linear in the original data space!
- Similar idea to support vector machines

## PCA in feature space (I)

- Suppose for now that the data is centered in feature space, i.e.  
 $\sum_{i=1}^m \phi(\mathbf{x}_i) = \mathbf{0}$

- The scatter matrix is:

$$\mathbf{S} = \sum_{i=1}^m \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top = \mathbf{\Phi}^\top \mathbf{\Phi} \in \mathbb{R}^{N \times N} \quad \text{where } \mathbf{\Phi}_{i,:} = \phi(\mathbf{x}_i)^\top$$

- The eigenvectors are:

$$\mathbf{S} \mathbf{v}_j = \lambda_j \mathbf{v}_j, \quad j = 1, \dots, N \quad (N \text{ is the dim. of the feature space})$$

- We want to avoid explicitly going to feature space - instead we want to work with *kernels and the Gram matrix*  $\mathbf{K} = \mathbf{\Phi} \mathbf{\Phi}^\top \in \mathbb{R}^{m \times m}$ :

$$\mathbf{K}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_j)^\top \phi(\mathbf{x}_i)$$

## PCA in feature space (II)

- Let  $\mathbf{v} \in \mathbb{R}^N$  be any eigenvector of the scatter matrix. We have

$$\lambda \mathbf{v} = \mathbf{S} \mathbf{v} = \Phi^\top \Phi \mathbf{v} = \sum_{i=1}^m \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \mathbf{v}$$

$\Rightarrow$  The eigenvectors can be written as a linear combinations of features:

$$\mathbf{v} = \sum_{i=1}^m \frac{1}{\lambda} (\phi(\mathbf{x}_i)^\top \mathbf{v}) \phi(\mathbf{x}_i) = \sum_{i=1}^m (\mathbf{a})_i \phi(\mathbf{x}_i) = \Phi^\top \mathbf{a}$$

- Finding an eigenvector  $\mathbf{v}$  of the scatter matrix is equivalent to finding the vector of coefficients  $\mathbf{a} \in \mathbb{R}^m$  (since  $\mathbf{v} = \Phi^\top \mathbf{a}$ )!

## PCA in feature space (III)

- By substituting  $\mathbf{v} = \Phi^\top \mathbf{a}$  back into the eigenvector equation we get:

$$\mathbf{S}\mathbf{v} = \lambda\mathbf{v} \Rightarrow \Phi^\top \Phi \mathbf{v} = \Phi^\top \Phi \Phi^\top \mathbf{a} = \lambda \Phi^\top \mathbf{a} \Rightarrow \Phi^\top \mathbf{K} \mathbf{a} = \lambda \Phi^\top \mathbf{a}$$

- A small trick: multiplying by  $\Phi$  to the left gives us

$$\Phi \Phi^\top \mathbf{K} \mathbf{a} = \lambda \Phi \Phi^\top \mathbf{a} \Rightarrow \mathbf{K}^2 \mathbf{a} = \lambda \mathbf{K} \mathbf{a}$$

- We can remove a factor of  $\mathbf{K}$  from both sides of the matrix (this will only affect eigenvectors with eigenvalues 0, which will not be principle components anyway):

$$\mathbf{K} \mathbf{a} = \lambda \mathbf{a}$$

$\Rightarrow$  For any eigenvector  $\mathbf{v}$  of the scatter matrix (in feature space), the corresponding vector of coefficients  $\mathbf{a}$  is an eigenvector of the Gram matrix (with the same eigenvalue)!



## PCA in feature space (IV)

- We know that  $\mathbf{a}$  is an eigenvector of  $\mathbf{K}$  but we don't know its norm yet...
- Remember that the eigenvector  $\mathbf{v}$  of  $\mathbf{S}$  must be of norm 1, this implies a dual normalization condition for the vector  $\mathbf{a}$ :

$$\|\mathbf{v}\|^2 = \mathbf{v}^\top \mathbf{v} = 1 \Rightarrow \mathbf{a}^\top \Phi \Phi^\top \mathbf{a} = \mathbf{a}^\top \mathbf{K} \mathbf{a} = 1$$

- Plugging this into  $\mathbf{K} \mathbf{a} = \lambda \mathbf{a}$  we get  $\|\mathbf{a}\|^2 = \frac{1}{\lambda}$ .  
→ We can rescale a unit-norm eigenvector  $\mathbf{z}$  of  $\mathbf{K}$  to obtain  $\mathbf{a} = \frac{1}{\sqrt{\lambda}} \mathbf{z}$ .
- As before, for a new point  $\mathbf{x}$ , let  $\mathbf{k}_x \in \mathbb{R}^m$  be defined by  $(\mathbf{k}_x)_i = K(\mathbf{x}, \mathbf{x}_i)$ .  
The projection of  $\mathbf{x}$  onto the  $j$ th principal components is:

$$\phi(\mathbf{x})^\top \mathbf{v}_j = \phi(\mathbf{x})^\top \Phi^\top \mathbf{a}_j = \mathbf{k}_x^\top \mathbf{a}_j$$

where  $\mathbf{v}_j$  is the  $j$ th eigenvector of  $\mathbf{S}$  and  $\mathbf{a}_j$  is the scaled  $j$ th eigenvector of  $\mathbf{K}$ !

## Normalizing the feature space

- In general, the features  $\phi(\mathbf{x}_i)$  may not have mean 0
- We want to work with  $\tilde{\phi}(\mathbf{x}_i) = \phi(\mathbf{x}_i) - \frac{1}{m} \sum_{k=1}^m \phi(\mathbf{x}_k)$
- The corresponding kernel matrix entries are given by:

$$\tilde{\mathbf{K}}_{i,j} = \tilde{K}(\mathbf{x}_i, \mathbf{x}_j) = \tilde{\phi}(\mathbf{x}_i)^\top \tilde{\phi}(\mathbf{x}_j)$$

- After some algebra, we get:

$$\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{O}_{1/m} \mathbf{K} - \mathbf{K} \mathbf{O}_{1/m} + \mathbf{O}_{1/m} \mathbf{K} \mathbf{O}_{1/m}$$

and

$$\tilde{\mathbf{k}}_{\mathbf{x}} = \mathbf{k}_{\mathbf{x}} - \mathbf{O}_{1/m} \mathbf{k}_{\mathbf{x}} - \mathbf{K} \mathbf{1}_{1/m} + \mathbf{O}_{1/m} \mathbf{K} \mathbf{1}_{1/m}$$

where  $\mathbf{O}_{1/m}$  (resp.  $\mathbf{1}_{1/m}$ ) is the matrix (resp. vector) with all elements equal to  $1/m$ .

## Kernel PCA: overall algorithm

1. Pick a kernel and build the Gram matrix  $\mathbf{K} \in \mathbb{R}^{m \times m}$ .
2. Compute the Gram matrix of the centered the data in the feature space:

$$\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{O}_{1/m}\mathbf{K} - \mathbf{K}\mathbf{O}_{1/m} + \mathbf{O}_{1/m}\mathbf{K}\mathbf{O}_{1/m}$$

3. Compute the top  $d$  eigenvalues and (unit-norm) eigenvectors of  $\tilde{\mathbf{K}}$ .
4. Put the eigenvectors into a matrix  $\mathbf{U} \in \mathbb{R}^{n \times d}$  and the corresponding eigenvalues in a diagonal matrix  $\mathbf{D} \in \mathbb{R}^{d \times d}$ .
5. The PCA projection of any point  $\mathbf{x}$  is given by

$$\hat{\mathbf{x}} = \mathbf{D}^{-1/2}\mathbf{U}^\top \tilde{\mathbf{k}}_{\mathbf{x}}$$

where  $\tilde{\mathbf{k}}_{\mathbf{x}}$  is defined as in the previous slide.

(Note that multiplying by  $\mathbf{D}^{-1/2}$  corresponds to rescaling the unit-norm eigenvectors of  $\tilde{\mathbf{K}}$  to get the vectors of coefficients  $\mathbf{a}_j$ ).

## Representation obtained by kernel PCA

- Each  $y_j = \phi(\mathbf{x})^\top \mathbf{v}_j = \mathbf{a}_j^\top \mathbf{k}_x$  is the coordinate of  $\phi(\mathbf{x})$  along one of the feature space axis  $\mathbf{v}_j$
- Since the  $\mathbf{v}_j$ 's are orthogonal, the projection of  $\phi(\mathbf{x})$  onto the space spanned by the top  $d$  eigenvectors is:

$$\Pi\phi(\mathbf{x}) = \sum_{j=1}^d y_j \mathbf{v}_j = \sum_{j=1}^d (\mathbf{a}_j^\top \mathbf{k}_x) \Phi^\top \mathbf{a}_j$$

- The *reconstruction error in feature space* can be evaluated as:

$$\|\phi(\mathbf{x}) - \Pi\phi(\mathbf{x})\|^2$$

This can be re-written by expanding the norm; we obtain dot-products which can all be replaced by kernels

- Note that the error will be 0 on the training data if enough  $\mathbf{v}_j$  are retained

## Alternative reconstruction error measures

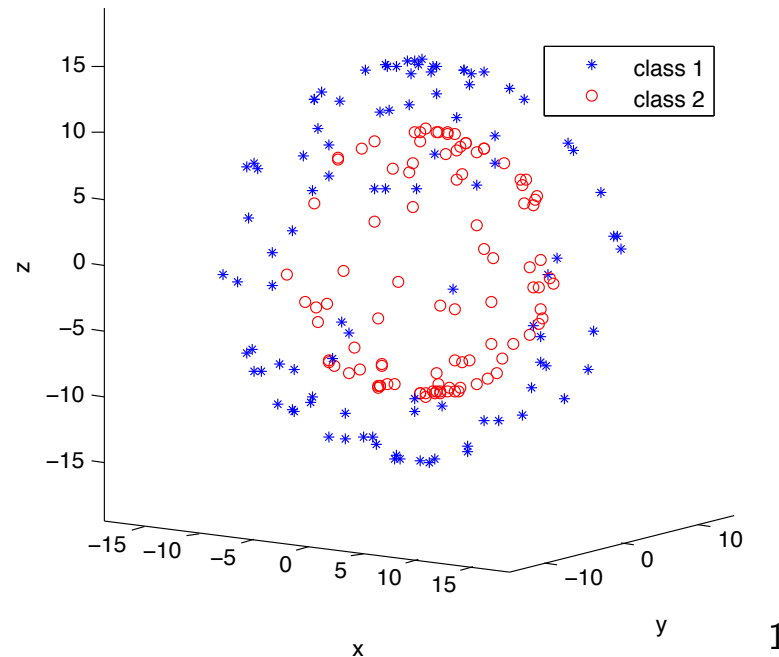
- An alternative way of measuring performance is by looking at how well kernel PCA preserves distances between data points
- In this case, the Euclidian distance in kernel space between points  $\phi(\mathbf{x}_i)$  and  $\phi(\mathbf{x}_j)$ ,  $d_{ij}$ , is:

$$\|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 = K(\mathbf{x}_i, \mathbf{x}_i) + K(\mathbf{x}_j, \mathbf{x}_j) - 2K(\mathbf{x}_i, \mathbf{x}_j)$$

- The distance  $\hat{d}_{ij}$  between the projected points in kernel space is defined as above, but with  $\phi(\mathbf{x}_i)$  replaced by  $\Pi\phi(\mathbf{x}_i)$ .
- The average of  $d_{ij} - \hat{d}_{ij}$  over all pairs of points is a measure of reconstruction error
- Note that reconstruction error in the original space of the  $\mathbf{x}_i$  is very difficult to compute, because it requires taking  $\Pi\phi(\mathbf{x})$  and finding its pre-image in the original feature space, which is not always feasible (though approximations exist)

# Example: Two concentric spheres

two concentric spheres data

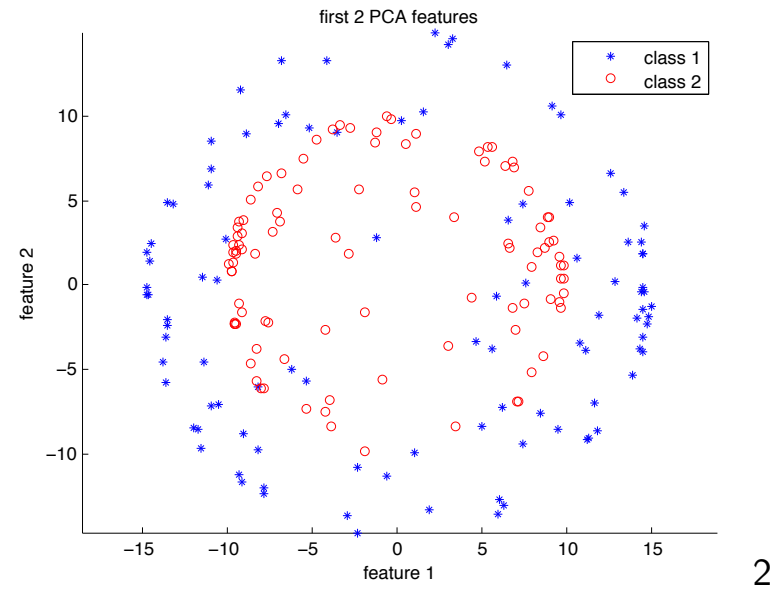


- Colours are used for clarity in the picture, but the data is presented unlabelled
- We want to project from 3D to 2D

---

<sup>1</sup>Wang, 2012

## Example: Two concentric spheres - PCA



2

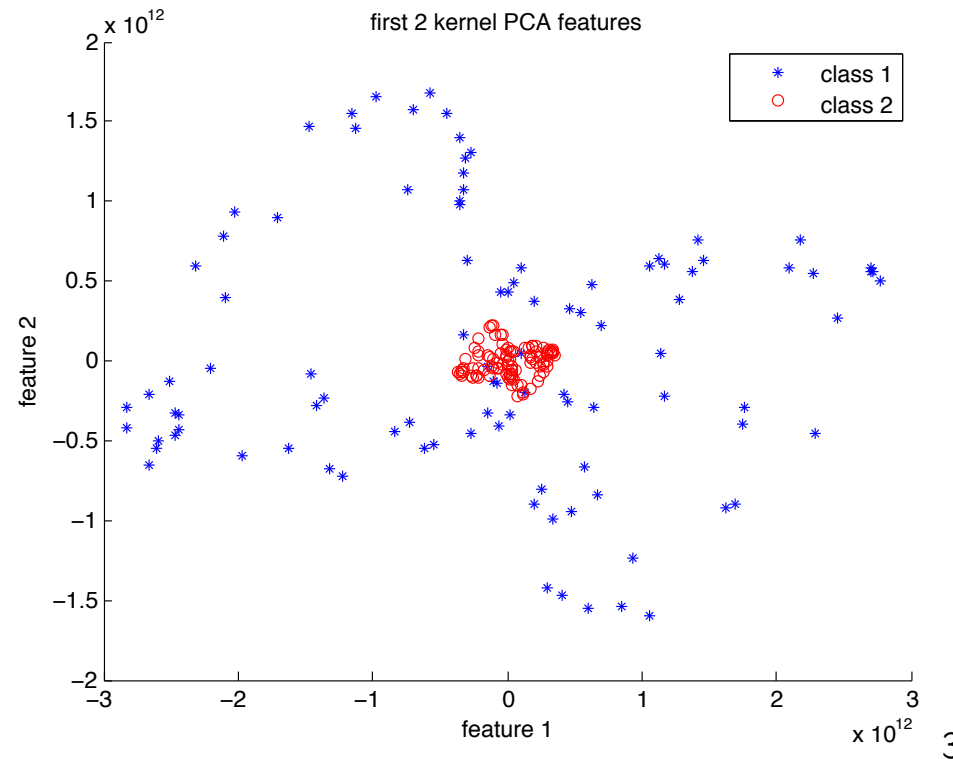
Note that PCA is unable to separate the points from the two spheres

---

<sup>2</sup>Wang, 2012



## Example: Kernel PCA with Polynomial Kernel ( $d = 5$ )

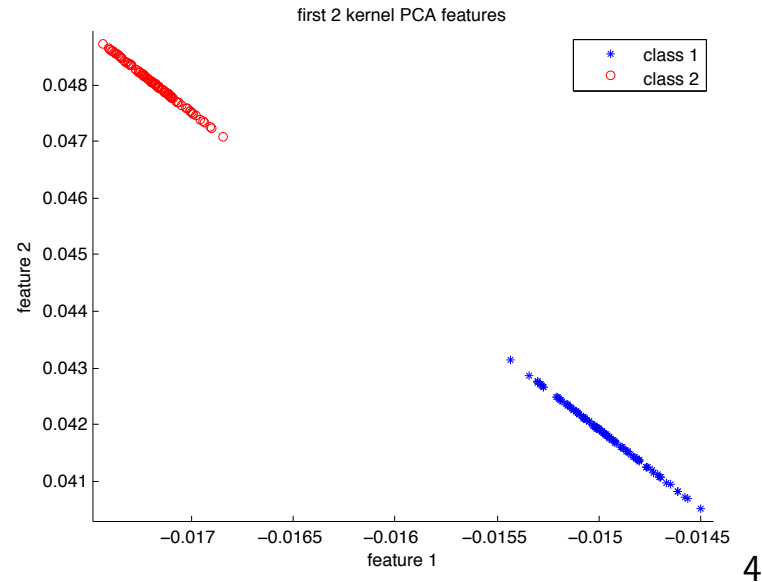


- Points from one sphere are much closer together, the others are scattered
- The projected data is not linearly separable

---

<sup>3</sup>Wang, 2012

## Example: Kernel PCA with Gaussian Kernel ( $\sigma = 20$ )



- Points from the two spheres are really well separated
- Note that the choice of parameter for the kernel matters!
- Validation can be used to determine good kernel parameter values

---

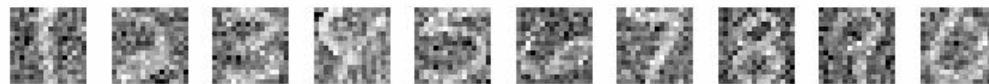
<sup>4</sup>Wang, 2012

## Example: De-noising images

Original data



Data corrupted with Gaussian noise



Result after linear PCA



Result after kernel PCA, Gaussian kernel



## PCA vs Kernel PCA

- Kernel PCA can give a good re-encoding of the data when it lies along a non-linear manifold
- The kernel matrix is  $m \times m$ , so kernel PCA will have difficulties if we have lots of data points
- In this case, we may need to use dictionary methods to pick a subset of the data
- For general kernels, we may not be able to easily visualize the image of a point in the input space, though visualization still works for simple kernels

## Locally Linear Embedding

- $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$  lies on a  $k$ -dimensional manifold.
- $\Rightarrow$  Each point and its neighbors lie close to a *locally linear* patch of the manifold.
- We try to reconstruct each point from its neighbors:

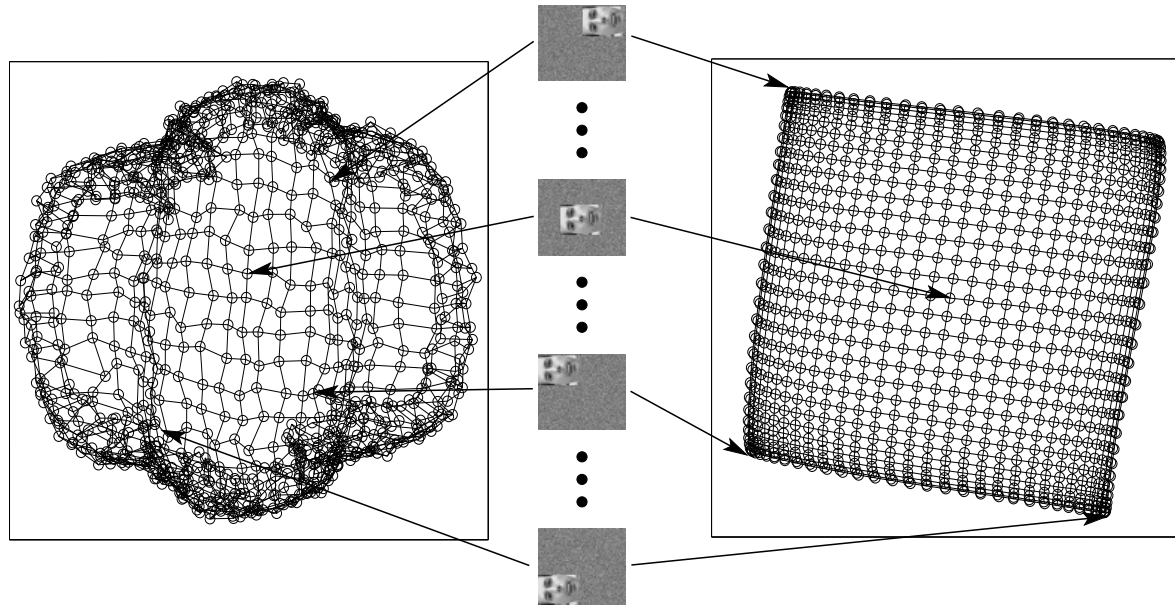
$$\min_{\mathbf{W}} \sum_i \|\mathbf{x}_i - \sum_j \mathbf{W}_{i,j} \mathbf{x}_j\|^2$$

s.t.  $\mathbf{W}\mathbf{1} = \mathbf{1}$  and  $\mathbf{W}_{i,j} = 0$  if  $\mathbf{x}_j \notin \text{neighbors}(\mathbf{x}_i)$

- $\Rightarrow$  For each point the weights are invariant to rotation, scaling and translations: the weights  $\mathbf{W}_{i,j}$  capture intrinsic geometric properties of each neighborhood.
- These local properties of each neighborhood should be preserved by the embedding:

$$\min_{\mathbf{z}_1, \dots, \mathbf{z}_m \in \mathbb{R}^k} \sum_i \|\mathbf{z}_i - \sum_j \mathbf{W}_{i,j} \mathbf{z}_j\|^2$$

# PCA vs Locally Linear Embedding



[Saul, L. K., & Roweis, S. T. (2000). An introduction to locally linear embedding.]