

# Dimensionality Reduction

Slide credits : Vineeth N Balasubramanian



# ML Problems

## *Supervised Learning*

## *Unsupervised Learning*

*Discrete*

classification or  
categorization

*Continuous*

regression

clustering

dimensionality  
reduction

# What is Dimensionality Reduction

- Refers to the mapping of the original high-dimensional data onto a lower-dimensional space.
- Criterion for feature reduction can be different for different problems.
  - Unsupervised setting: minimize the information loss
  - Supervised setting: maximize the class discrimination
- Given a set of data points of  $p$  variables  $\{x_1, x_2, \dots, x_n\}$  : Compute the linear transformation (projection)

$$G \in \mathbb{R}^{p \times d} : x \in \mathbb{R}^p \rightarrow y = G^T x \in \mathbb{R}^d \quad (d \ll p)$$

$$\begin{matrix} G^T \\ \times \end{matrix} = \begin{matrix} Y \end{matrix}$$

# Dimensionality Reduction vs Feature Selection

- **Feature reduction**

- All original features are used
- The transformed features are linear combinations of the original features (in case of linear DR).

- **Feature selection**

- Only a subset of the original features are used.

# Why DR?

- Most machine learning techniques may not be effective for high-dimensional data
  - Curse of Dimensionality
  - Query accuracy and efficiency degrade rapidly as the dimension increases
  - Lower space and time complexity
  - Visualization, Data compression, Noise/irrelevant feature removal
- The intrinsic dimension may be small
  - For example, the number of genes responsible for a certain type of disease may be small

# High-dimensional data are strange

- Consider the hypersphere of radius  $r$  on a space of dimension  $d$

$$\mathcal{S} = \left\{ \mathbf{x} \mid \sum_{i=1}^d x_i^2 \leq r^2 \right\}$$

- Its volume is

$$V_d(r) = \frac{r^d \pi^{\frac{d}{2}}}{\Gamma\left(\frac{d}{2} + 1\right)}$$

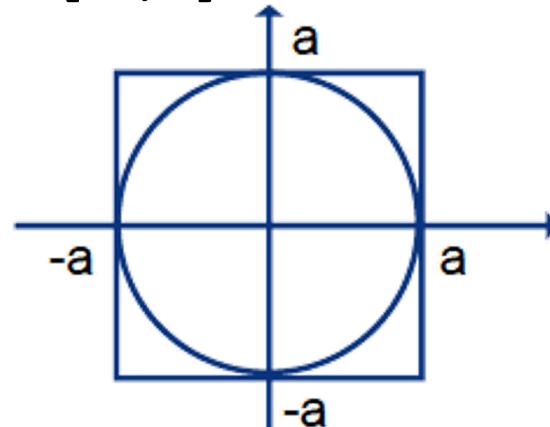
- Where  $\Gamma(n)$  is the Gamma function

$$\Gamma(n) = \int_0^\infty e^{-x} x^{n-1} dx$$

Source: N Vasconcelos

# Hypercube vs Hypersphere

- Consider the hyper-cube  $[-a,a]^d$  and the inscribed hyper-sphere, i.e.



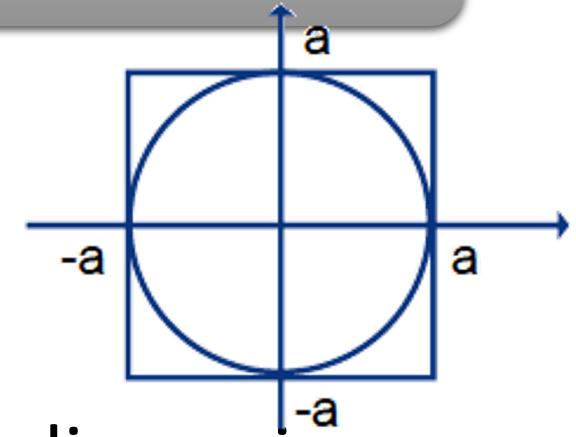
- What does your intuition tell you about the relative sizes of these two objects?
  - Volume of sphere  $\approx$  volume of cube
  - Volume of sphere  $\gg$  volume of cube
  - Volume of sphere  $\ll$  volume of cube

Source: N Vasconcelos

# Hypercube vs Hypersphere

- Let's compute the answer

$$f_d = \frac{Vol(sphere)}{Vol(cube)} = \frac{\frac{a^d \pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2}+1)}}{(2a)^d} = \frac{\pi^{\frac{d}{2}}}{2^d \Gamma(\frac{d}{2} + 1)}$$



- Sequence that does not depend on a, just on the dimension d!

d	1	2	3	4	5	6	7
$f_d$	1	.785	.524	.308	.164	.08	.037

- It goes to zero, and goes to zero fast!

Source: N Vasconcelos

# Hypercube vs Hypersphere

- Let's compute the answer

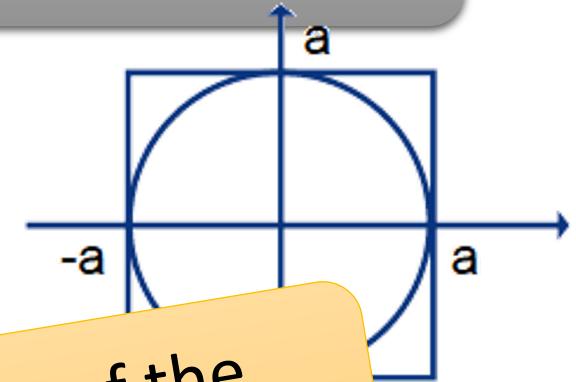
$$f_d = \frac{Vol(sphere)}{Vol(cube)} = \frac{\frac{a^d \pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2}+1)}}{(2a)^d} = \frac{\pi^{\frac{d}{2}}}{2^d \Gamma(\frac{d}{2} + 1)}$$

- Sequence that does not depend on d!

As the dimension of the space increases, the volume of the sphere is much smaller (infinitesimal) than that of the cube!

.785	.524	.308	.164	.08	.037
------	------	------	------	-----	------

- It goes to zero, and goes to zero fast!



Source: N Vasconcelos

# Hypercube vs Hypersphere

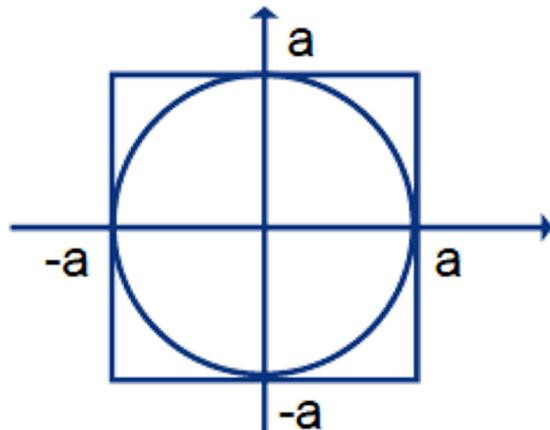
Actually not very surprising

1.  $d = 1$



Volume is the same

2.  $d = 2$

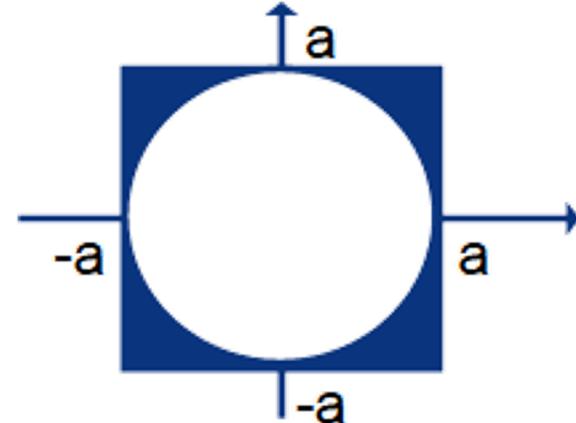


Volume of sphere is already smaller

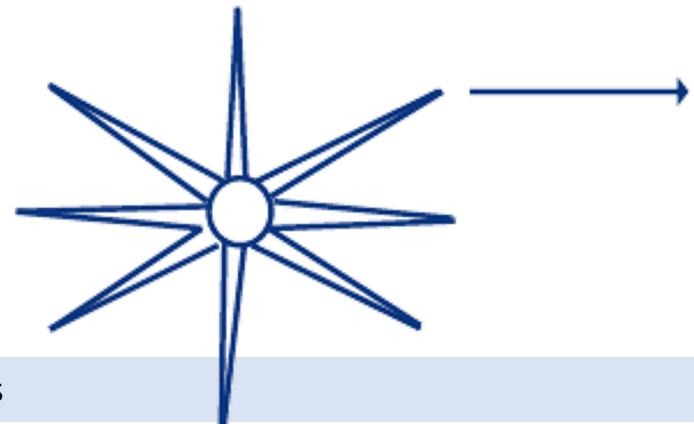
Source: N Vasconcelos

# Hypercube vs Hypersphere

- As the dimension increases the volume of the shaded corners becomes larger



- In high dimensions the picture you should have in mind is



all the volume of the cube  
is in these spikes!

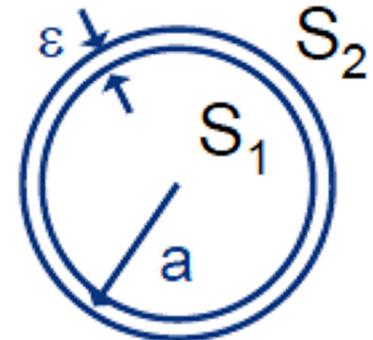


Source: N Vasconcelos

# The Curse of Dimensionality

- Consider the crust of unit hypersphere of thickness  $\epsilon$
- Let's compute the ratio of volumes

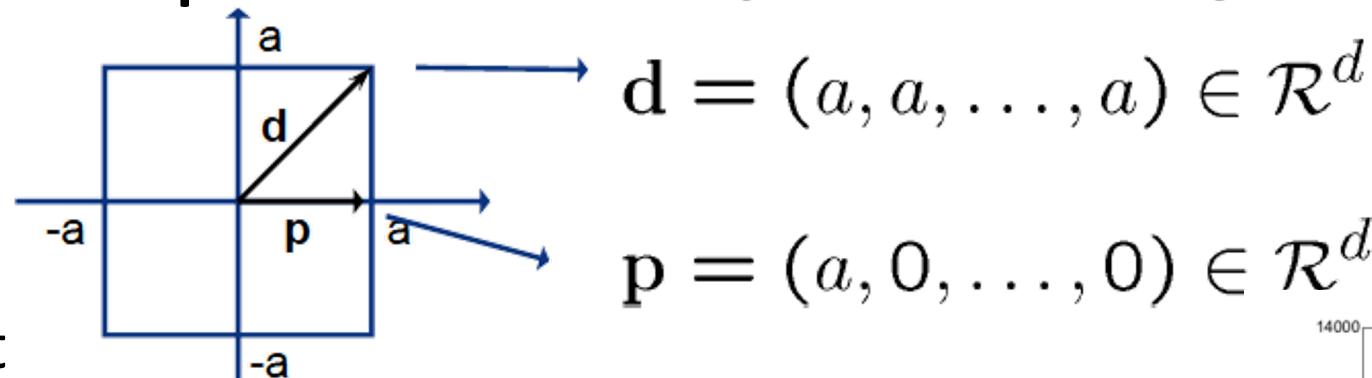
$$\frac{Vol(S_1)}{Vol(S_2)} = \frac{\frac{(a-\epsilon)^d \pi^{\frac{d}{2}}}{\Gamma\left(\frac{d}{2}+1\right)}}{\frac{a^d \pi^{\frac{d}{2}}}{\Gamma\left(\frac{d}{2}+1\right)}} = \frac{a^d \left(1 - \frac{\epsilon}{a}\right)^d}{a^d} = \left(1 - \frac{\epsilon}{a}\right)^d$$



- No matter how small  $\epsilon$  is, ratio goes to zero as  $d$  increases i. e. **“all the volume is in the crust!”**

# We can check mathematically

- Consider  $\mathbf{d}$  and  $\mathbf{p}$

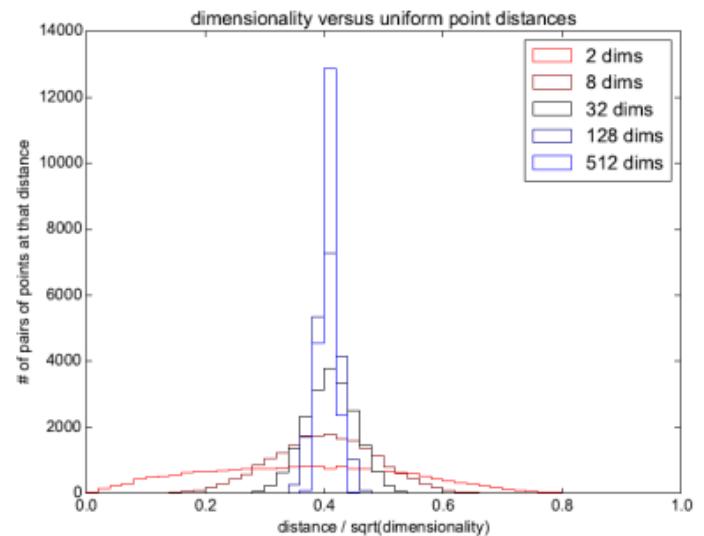


- Note that

$$\frac{\|\mathbf{d}\|^2}{\|\mathbf{p}\|^2} = \frac{da^2}{a^2} = d \rightarrow \infty$$

$$\begin{aligned} \cos\theta &= \frac{\mathbf{d}^T \mathbf{p}}{\sqrt{\|\mathbf{d}\|^2 \|\mathbf{p}\|^2}} \\ &= \frac{a^2}{\sqrt{da^2 a^2}} = \frac{1}{\sqrt{d}} \rightarrow 0 \end{aligned}$$

- $\mathbf{d}$  orthogonal to  $\mathbf{p}$  as  $\mathbf{d}$  increases and infinitely larger!!!

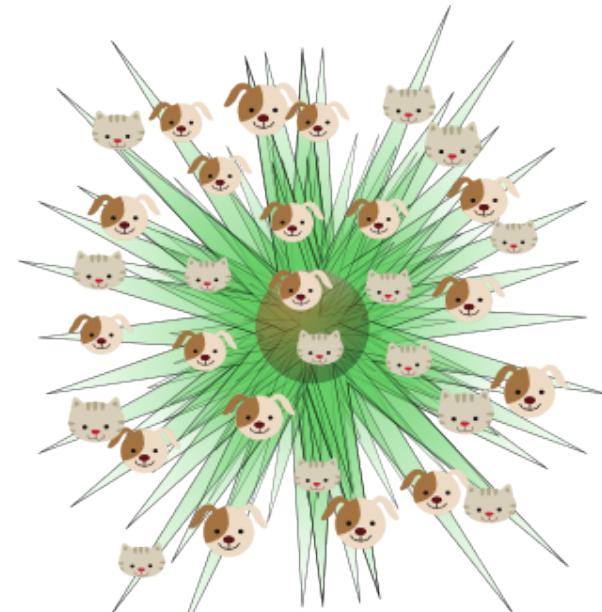
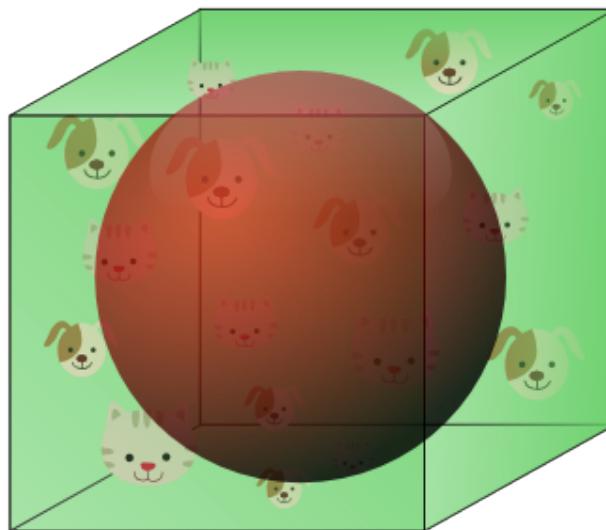
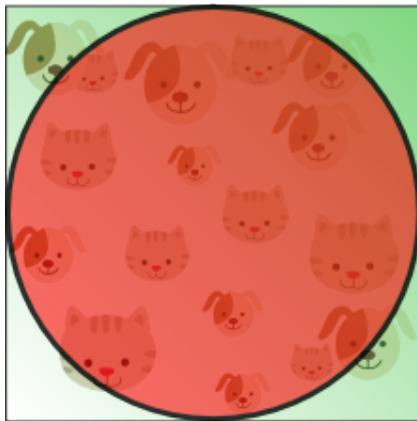
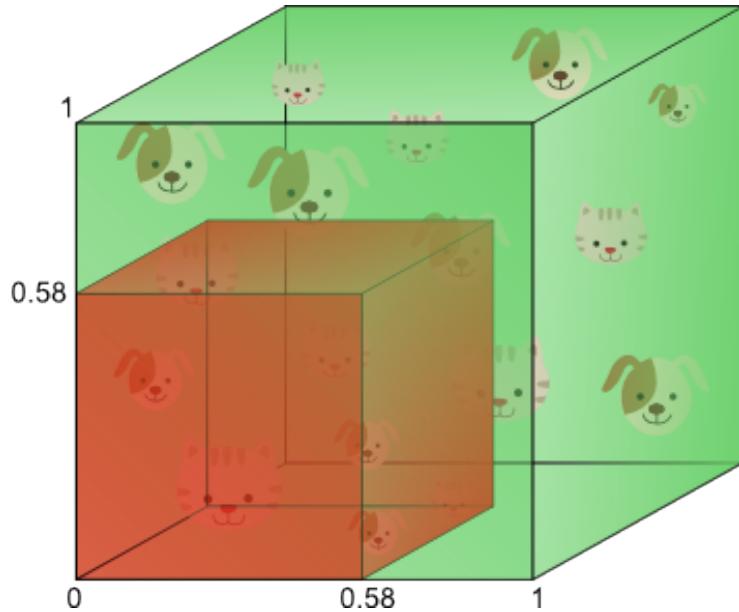
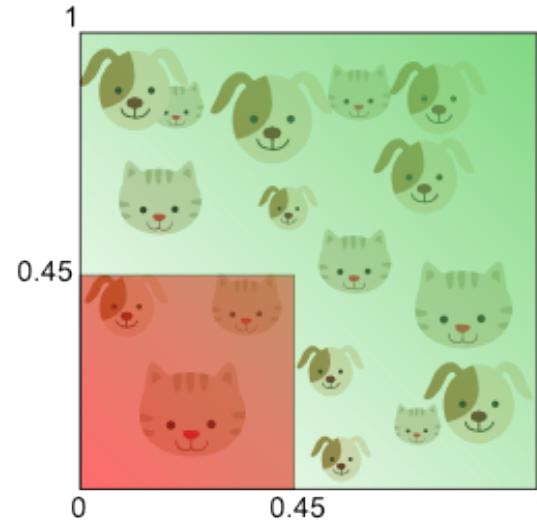


Source: N Vasconcelos

Figure 3.18: histogram of distances in  $D=2,8,32,128,512$

# Curse of dimensionality

$$\lim_{d \rightarrow \infty} \frac{\text{dist}_{\max} - \text{dist}_{\min}}{\text{dist}_{\min}} \rightarrow 0$$

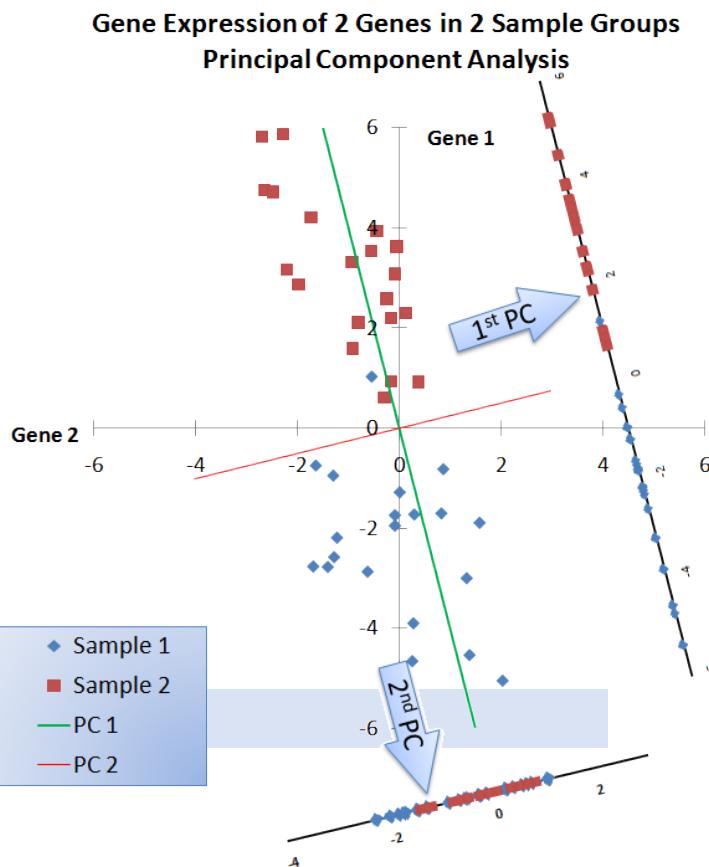
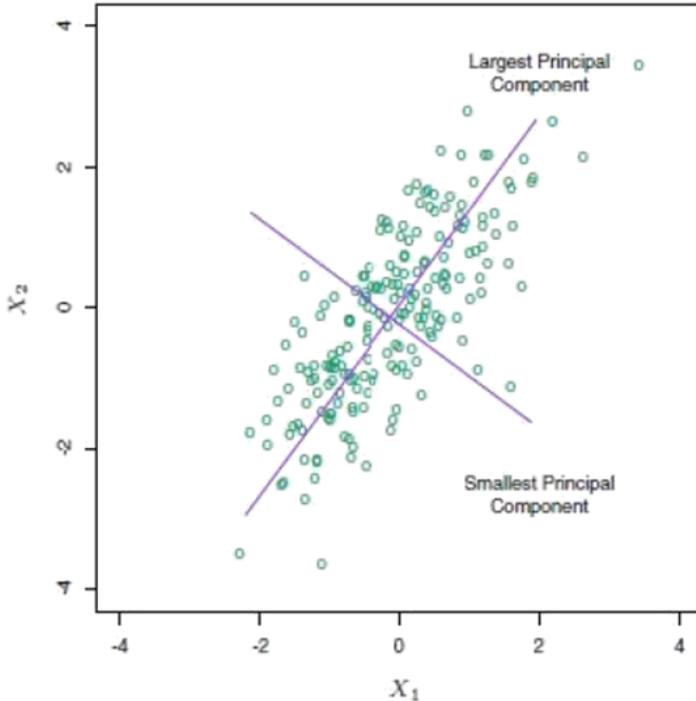


# Outline

- Linear Methods
  - Principal Component Analysis (PCA)
  - Fisher (Linear) Discriminant Analysis (LDA)

# Principal Component 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:  $z = \mathbf{w}^T \mathbf{x}$
- Find  $\mathbf{w}$  such that  $\text{Var}(z)$  is maximized,
- Our goal is to project the data onto a space having dimensionality  $M < D$  while maximizing the variance



# Principal Component 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:  $z = \mathbf{w}^T \mathbf{x}$
- Find  $\mathbf{w}$  such that  $\text{Var}(z)$  is maximized

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

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

# PCA

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

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

Using Lagrange multipliers, a la SVM

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

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  $\mathbf{w}_1$

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

How?

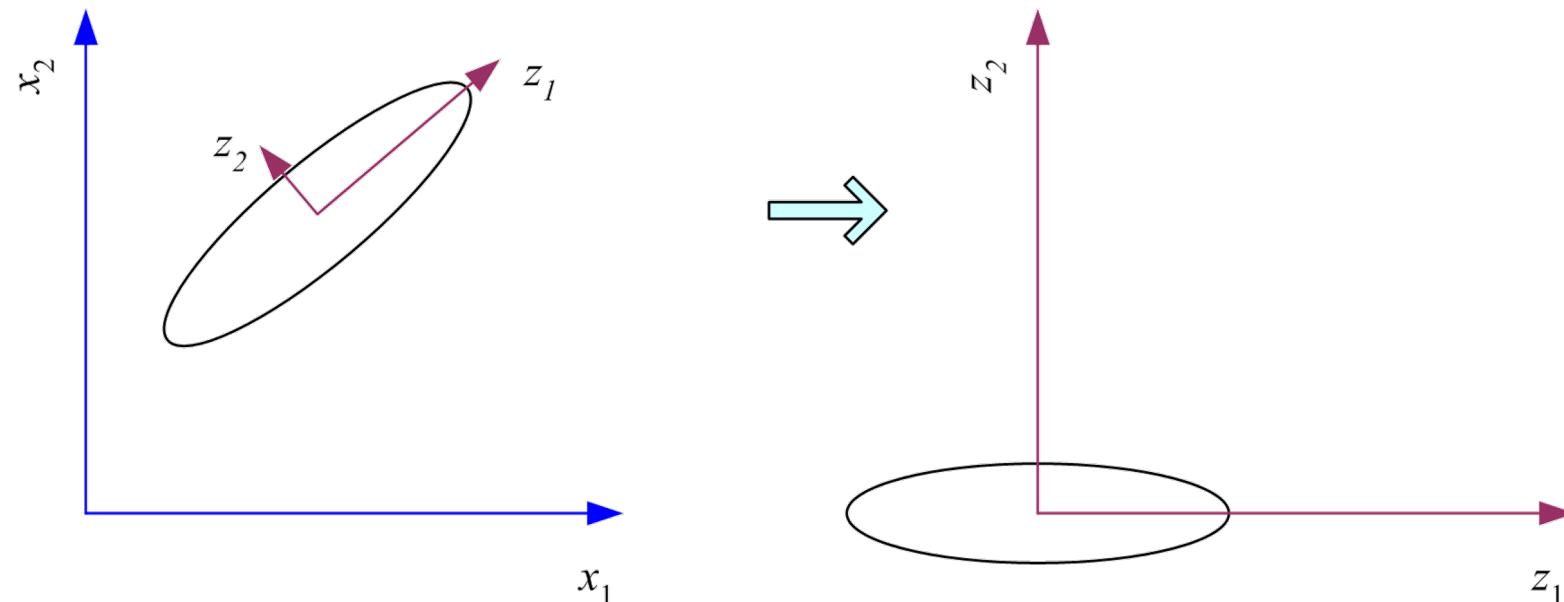
$\Sigma \mathbf{w}_2 = \alpha \mathbf{w}_2$  that is,  $\mathbf{w}_2$  is another eigenvector of  $\Sigma$

and so on.

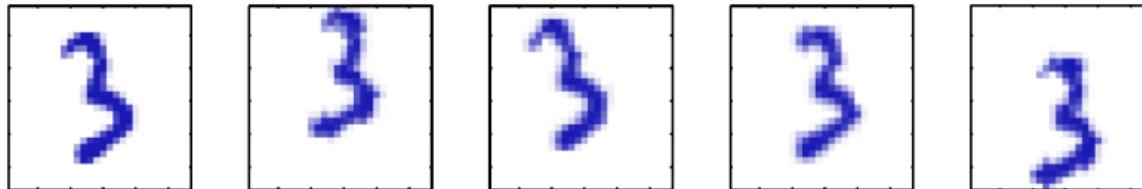
# PCA

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

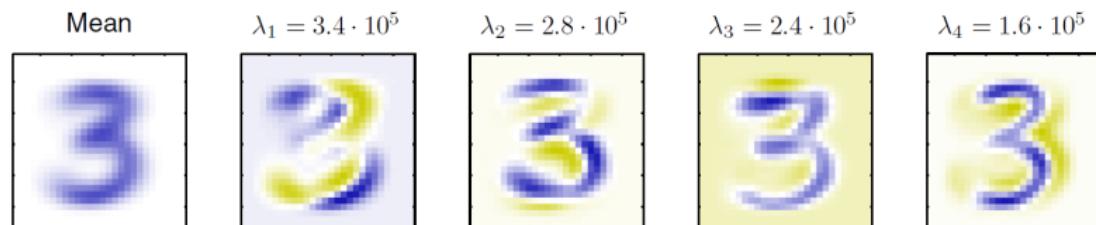
where the columns of  $\mathbf{W}$  are the eigenvectors of  $\Sigma$ , and  $\boldsymbol{\mu}$  is sample mean; Centers the data at the origin and rotates the axes



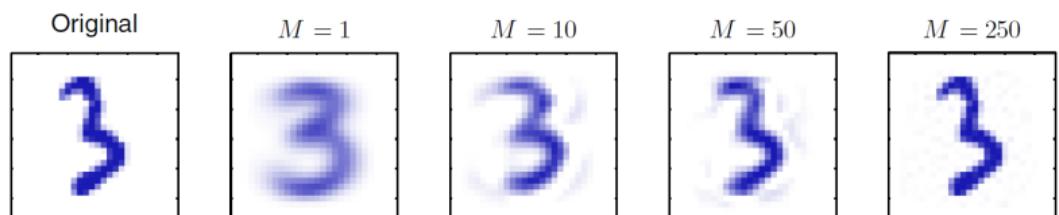
# PCA



A synthetic data set obtained by taking one of the off-line digit images and creating multiple copies in each of which the digit has undergone a random displacement and rotation within some larger image field. The resulting images each have  $100 \times 100 = 10,000$



The mean vector  $\bar{x}$  along with the first four PCA eigenvectors  $u_1, \dots, u_4$  for the off-line digits data set, together with the corresponding eigenvalues.



An original example from the off-line digits data set together with its PCA reconstructions obtained by retaining  $M$  principal components for various values of  $M$ . As  $M$  increases the reconstruction becomes more accurate and would become perfect when  $M = D = 28 \times 28 = 784$ .

# How to choose k

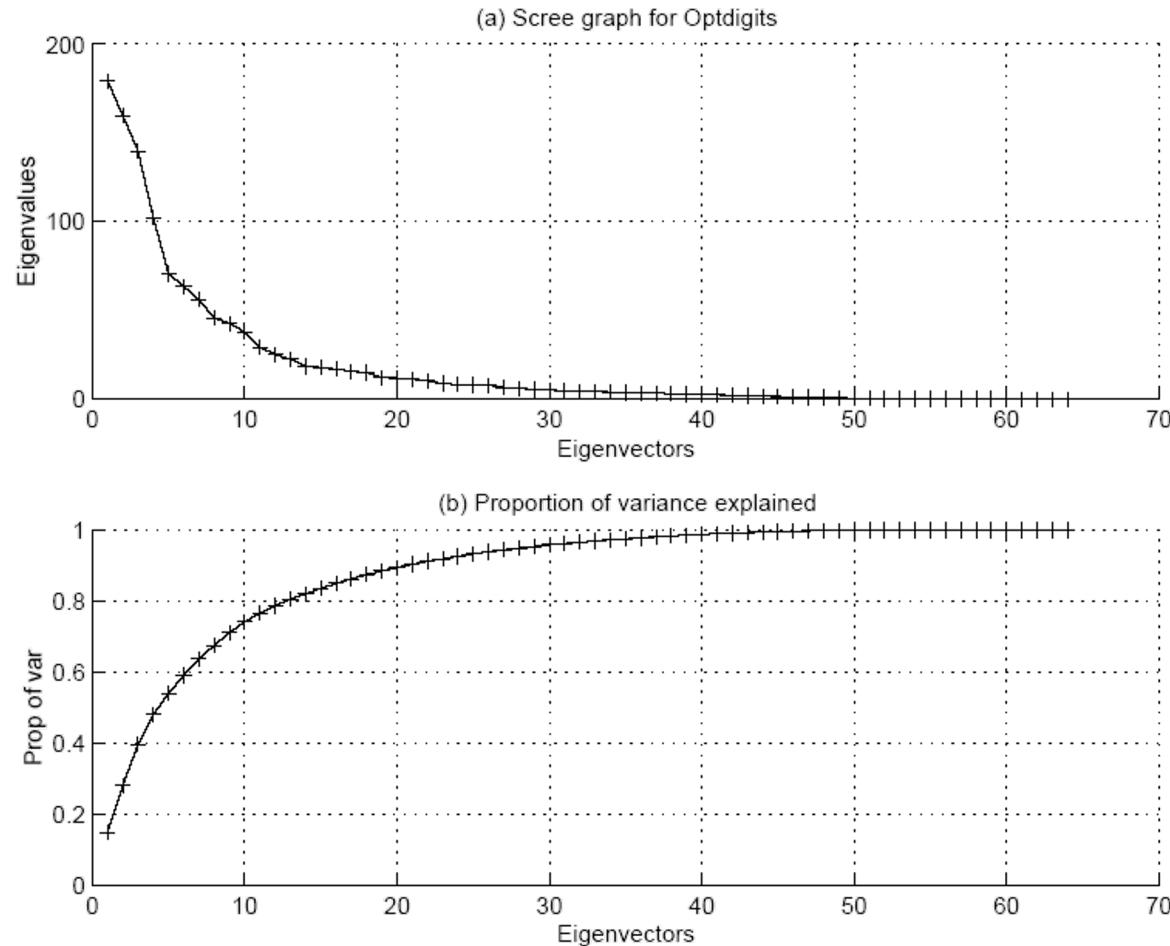
- Proportion of Variance (PoV) explained

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

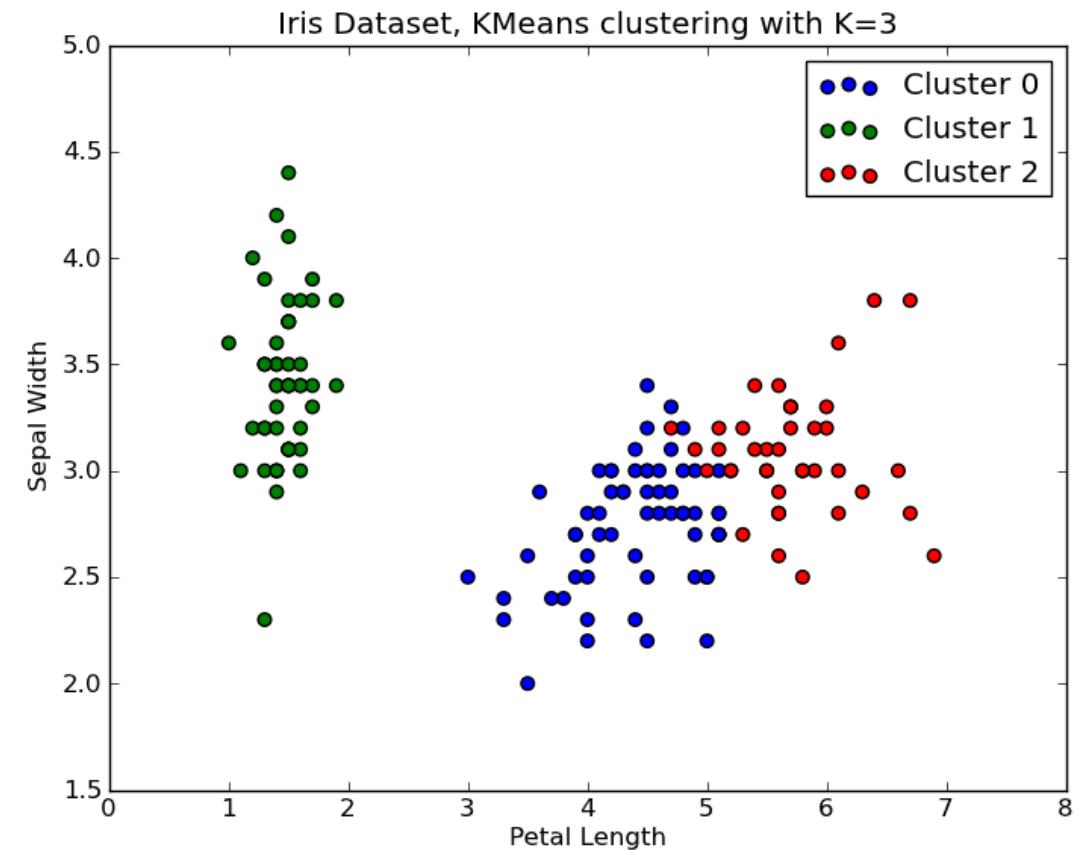
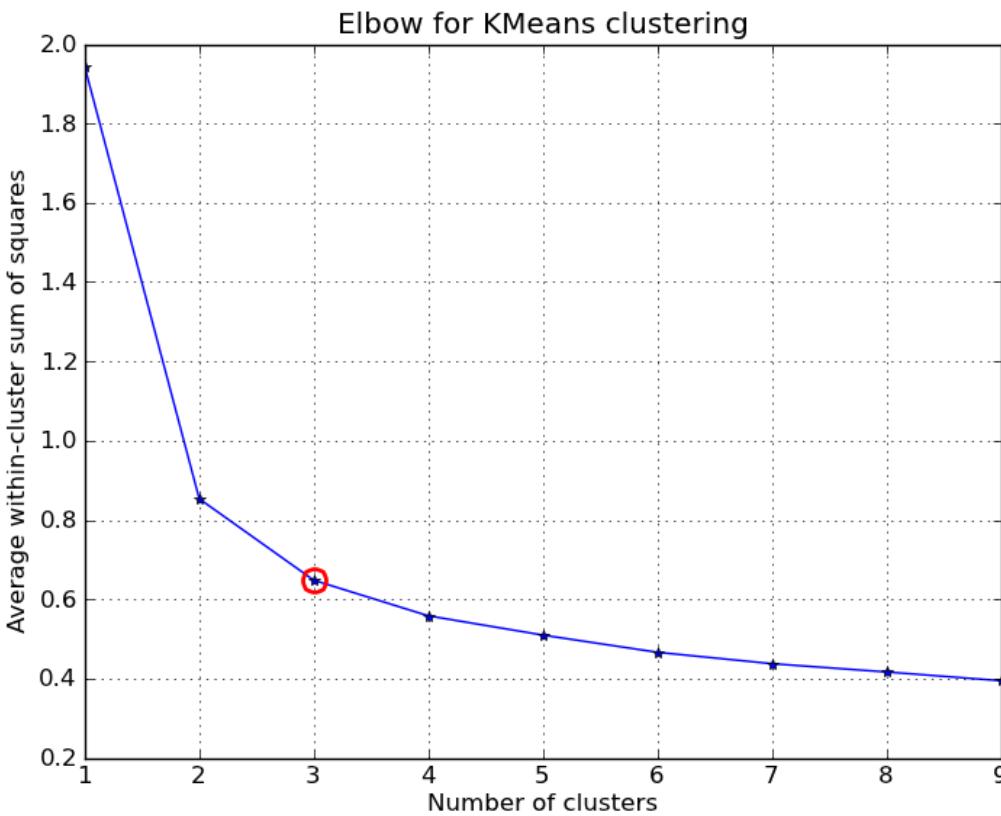
when  $\lambda_i$  are sorted in descending order

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

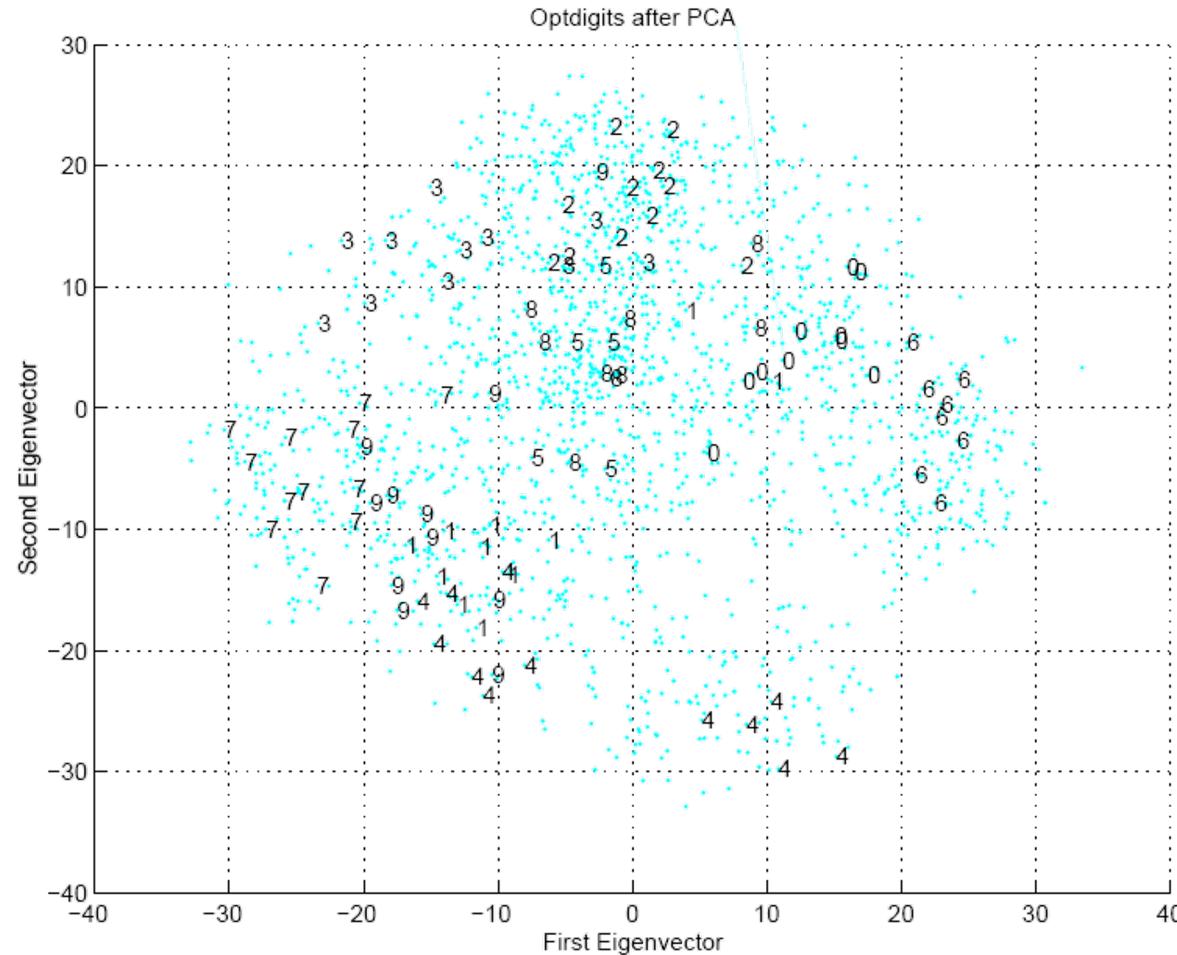
# Illustration



# Elbow plot in Clustering



# Example



# PCA for high dimensional data

$N \ll D$

$$\frac{1}{N} \mathbf{X}^T \mathbf{X} \mathbf{u}_i = \lambda_i \mathbf{u}_i.$$

$O(N^3)$  instead of  $O(D^3)$

$$\frac{1}{N} \mathbf{X} \mathbf{X}^T (\mathbf{X} \mathbf{u}_i) = \lambda_i (\mathbf{X} \mathbf{u}_i).$$

$$\frac{1}{N} \mathbf{X} \mathbf{X}^T \mathbf{v}_i = \lambda_i \mathbf{v}_i$$

$$\left( \frac{1}{N} \mathbf{X}^T \mathbf{X} \right) (\mathbf{X}^T \mathbf{v}_i) = \lambda_i (\mathbf{X}^T \mathbf{v}_i)$$

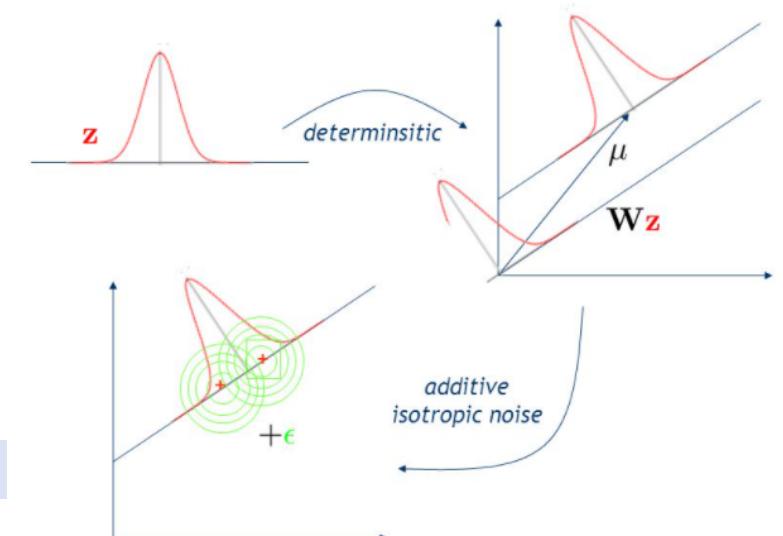
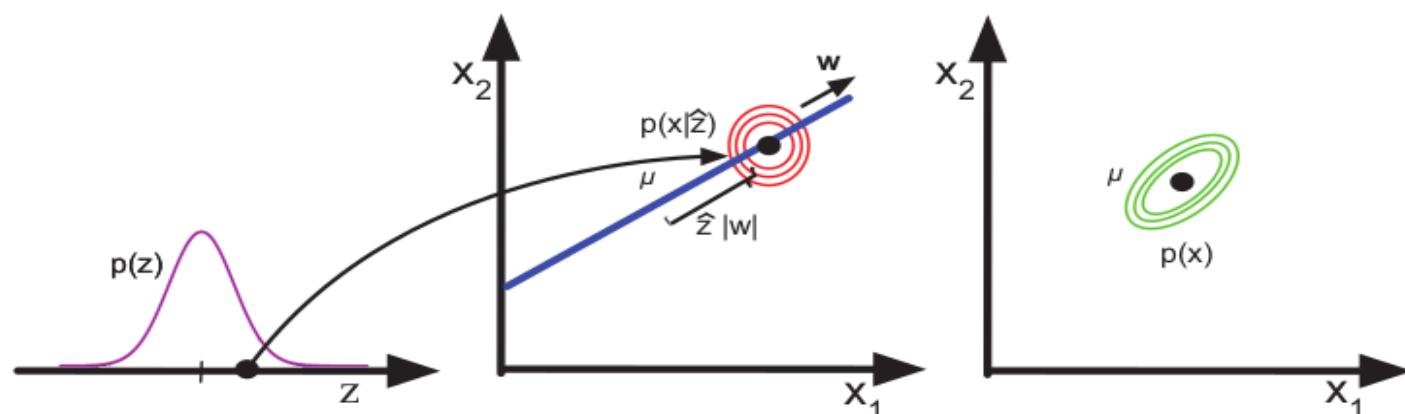
# Probabilistic PCA

- Maximum likelihood solution of a probabilistic latent variable model.
- latent variable  $\mathbf{z}$  corresponding to the principal-component subspace.
- Gaussian prior distribution  $p(\mathbf{z})$  over the latent variable, together with a Gaussian conditional distribution  $p(\mathbf{x}|\mathbf{z})$  for the observed variable  $\mathbf{x}$  conditioned on the value of the latent variable.

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}).$$

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$$

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}$$



**Figure 12.1** Illustration of the PPCA generative process, where we have  $L = 1$  latent dimension generating  $D = 2$  observed dimensions. Based on Figure 12.9 of (Bishop 2006b).

# Probabilistic PCA

- Maximum likelihood solution of a probabilistic latent variable model.
- latent variable  $\mathbf{z}$  corresponding to the principal-component subspace.
- Gaussian prior distribution  $p(\mathbf{z})$  over the latent variable, together with a Gaussian conditional distribution  $p(\mathbf{x}|\mathbf{z})$  for the observed variable  $\mathbf{x}$  conditioned on the value of the latent variable.

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}). \quad p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2\mathbf{I}) \quad \mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}$$

- Parameter Estimation : marginal likelihood  $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) d\mathbf{z}.$

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C})$$

- $$\begin{aligned}\mathbb{E}[\mathbf{x}] &= \mathbb{E}[\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}] = \boldsymbol{\mu} \\ \text{cov}[\mathbf{x}] &= \mathbb{E}[(\mathbf{W}\mathbf{z} + \boldsymbol{\epsilon})(\mathbf{W}\mathbf{z} + \boldsymbol{\epsilon})^T] \\ \mathbf{C} &= \mathbb{E}[\mathbf{W}\mathbf{z}\mathbf{z}^T\mathbf{W}^T] + \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T] = \mathbf{W}\mathbf{W}^T + \sigma^2\mathbf{I}\end{aligned}$$

# Probabilistic PCA

- Maximum likelihood solution of a probabilistic latent variable model.
- latent variable  $\mathbf{z}$  corresponding to the principal-component subspace.
- Gaussian prior distribution  $p(\mathbf{z})$  over the latent variable, together with a Gaussian conditional distribution  $p(\mathbf{x}|\mathbf{z})$  for the observed variable  $\mathbf{x}$  conditioned on the value of the latent variable.

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}).$$

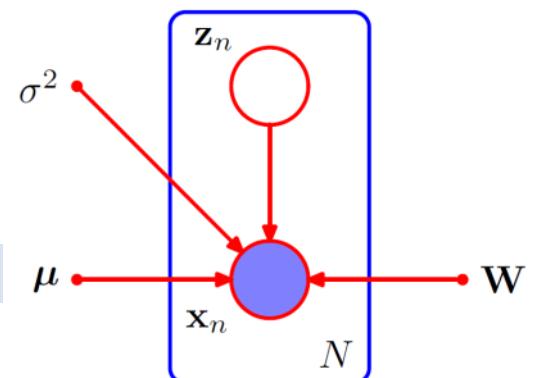
$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$$

$$\mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}$$

- Parameter Estimation : marginal likelihood  $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) d\mathbf{z}.$
- Predictive distribution : finding  $\mathbf{z}$  for an  $\mathbf{x}$

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C})$$

$$p(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\mathbf{M}^{-1}\mathbf{W}^T(\mathbf{x} - \boldsymbol{\mu}), \sigma^{-2}\mathbf{M}).$$



# Probabilistic PCA

- Gaussian prior distribution  $p(\mathbf{z})$  over the latent variable, together with a Gaussian conditional distribution  $p(\mathbf{x}|\mathbf{z})$  for the observed variable  $\mathbf{x}$  conditioned on the value of the latent variable.

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}). \quad p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I}) \quad \mathbf{x} = \mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}$$

- Parameter Estimation : marginal likelihood  $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) d\mathbf{z}$ .  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C})$

$$\begin{aligned} \ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{W}, \sigma^2) &= \sum_{n=1}^N \ln p(\mathbf{x}_n|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) \\ &= -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{C}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{C}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}). \end{aligned}$$

$$\ln p(\mathbf{X}|\mathbf{W}, \boldsymbol{\mu}, \sigma^2) = -\frac{N}{2} \left\{ D \ln(2\pi) + \ln |\mathbf{C}| + \text{Tr}(\mathbf{C}^{-1} \mathbf{S}) \right\}$$

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T.$$

$$\mathbf{W}_{\text{ML}} = \mathbf{U}_M (\mathbf{L}_M - \sigma^2 \mathbf{I})^{1/2} \mathbf{R}$$

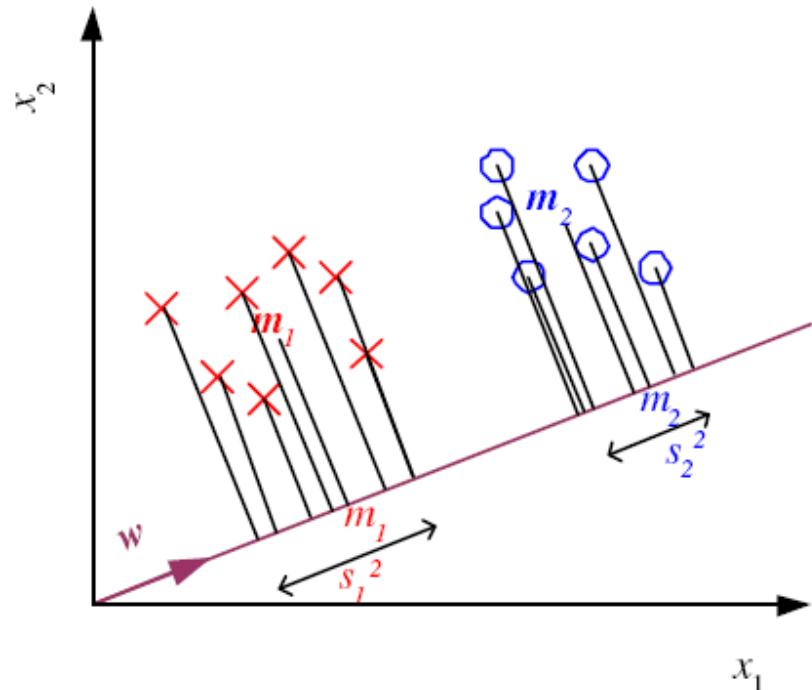
where  $\mathbf{U}_M$  is a  $D \times M$  matrix whose columns are given by any subset (of size  $M$ ) of the eigenvectors of the data covariance matrix  $\mathbf{S}$ , the  $M \times M$  diagonal matrix  $\mathbf{L}_M$  has elements given by the corresponding eigenvalues  $\lambda_i$ , and  $\mathbf{R}$  is an arbitrary  $M \times M$  orthogonal matrix.

# Outline

- Linear Methods
  - Principal Component Analysis (PCA)
  - Fisher (Linear) Discriminant Analysis (LDA)

# Fishers Linear Discriminant

- Supervised linear DR method
- Find a low-dimensional space such that when  $\mathbf{x}$  is projected, classes are well-separated.
- Simplest measure of the separation of the classes, when projected onto  $\mathbf{w}$ , is the separation of the projected class means

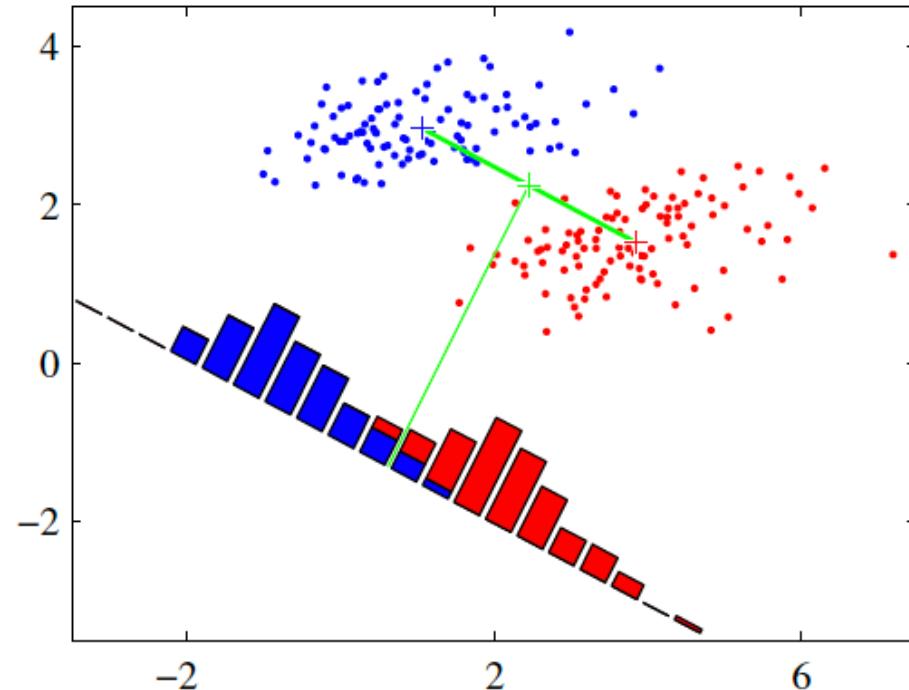


$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} \mathbf{x}_n. \quad \text{Maximize} \quad \mathbf{m}_2 - \mathbf{m}_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)$$

Subj. to

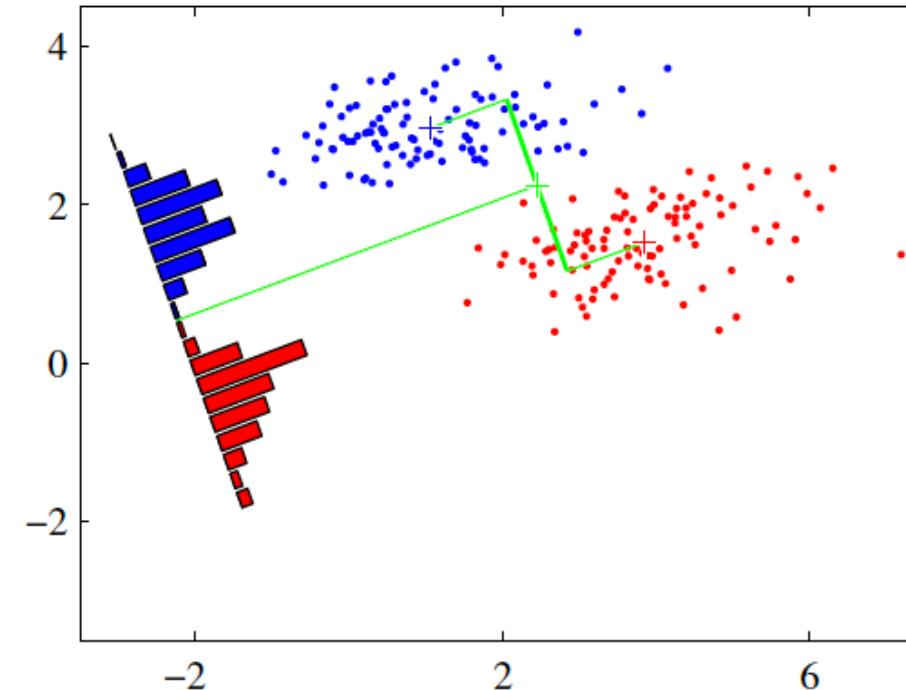
$$\sum_i w_i^2 = 1.$$

# Fishers Linear Discriminant



$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in \mathcal{C}_1} \mathbf{x}_n,$$

$$\mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in \mathcal{C}_2} \mathbf{x}_n.$$



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

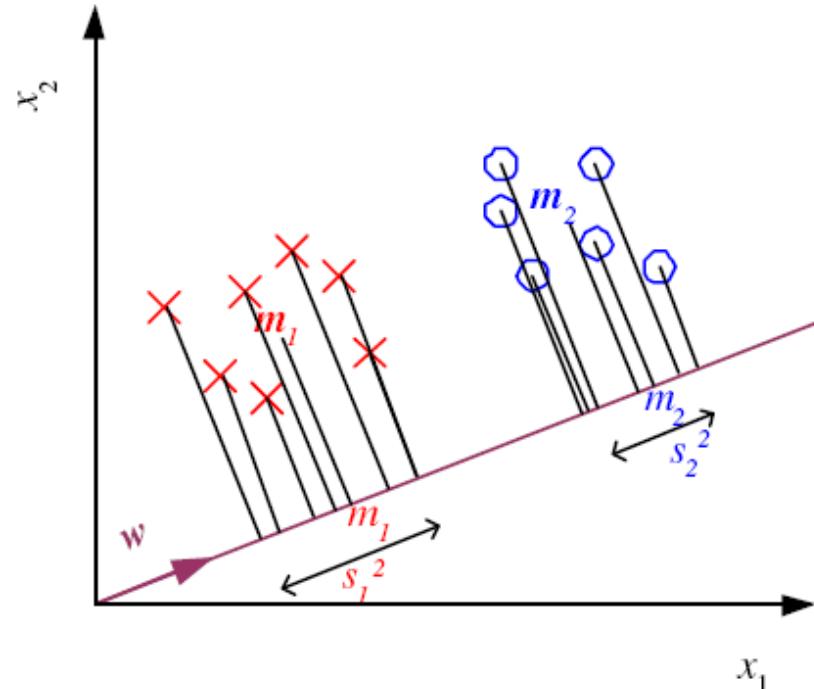
$$\text{Subj. to } \sum_i w_i^2 = 1.$$

# Fishers Linear Discriminant

- Idea proposed by Fisher is to maximize a function that will give a large separation between the projected class means while also giving a small variance within each class, thereby minimizing the class overlap.
- Find  $\mathbf{w}$  that maximizes

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

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



# Fishers Linear Discriminant

- 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:

$$\begin{aligned}s_1^2 &= \sum_t (\mathbf{w}^T \mathbf{x}^t - m_1)^2 r^t \\&= \sum_t \mathbf{w}^T (\mathbf{x}^t - \mathbf{m}_1)(\mathbf{x}^t - \mathbf{m}_1)^T \mathbf{w} r^t = \mathbf{w}^T \mathbf{S}_1 \mathbf{w} \\&\text{where } \mathbf{S}_1 = \sum_t (\mathbf{x}^t - \mathbf{m}_1)(\mathbf{x}^t - \mathbf{m}_1)^T r^t \\s_1^2 + s_2^2 &= \mathbf{w}^T \mathbf{S}_W \mathbf{w} \text{ where } \mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2\end{aligned}$$

# Fisher's Linear Discriminant

- Find  $\mathbf{w}$  that max

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

- FLD soln:

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

- Alternative parametric soln:

$$\mathbf{w} = \Sigma^{-1} (\mu_1 - \mu_2)$$

when  $p(\mathbf{x} | C_i) \sim \mathcal{N}(\mu_i, \Sigma)$

# Fisher's Linear Discriminant

- Within-class scatter:

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

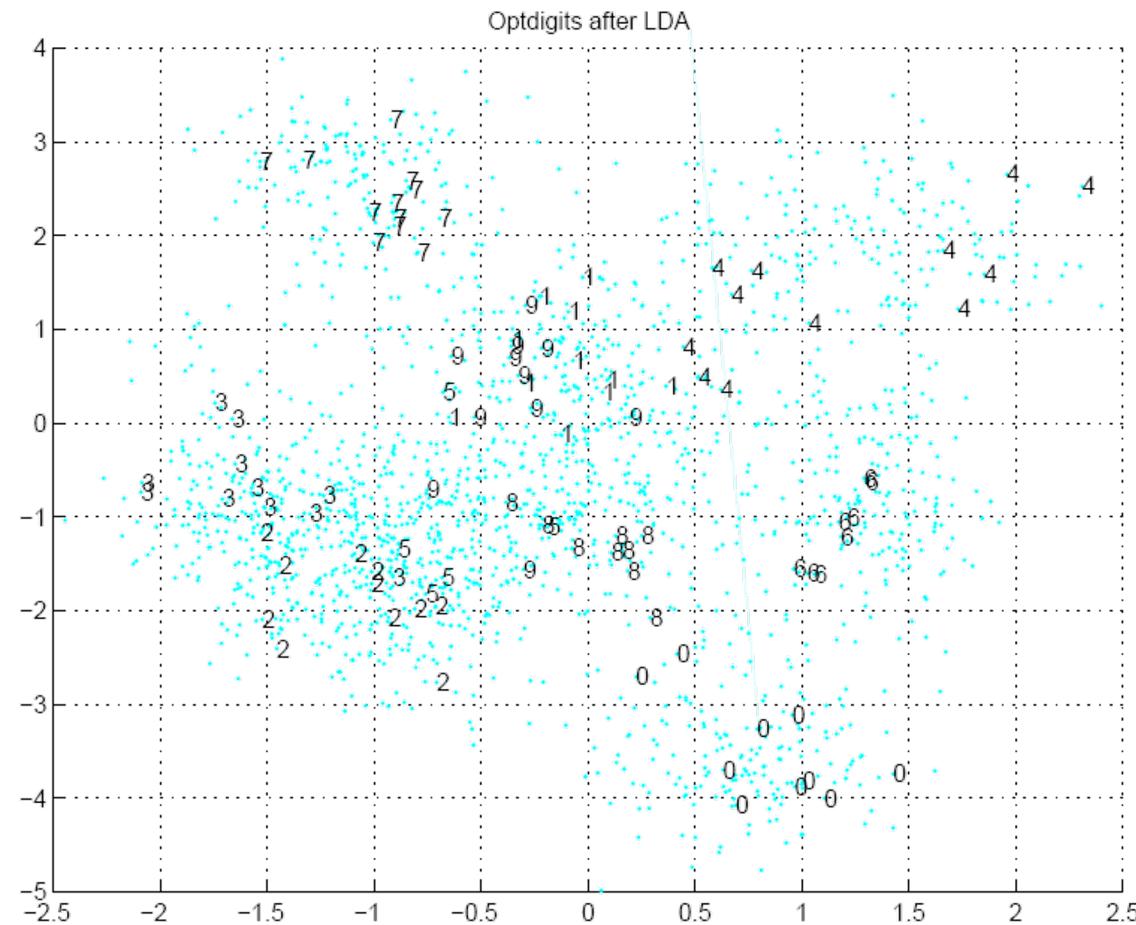
- Between-class scatter:

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

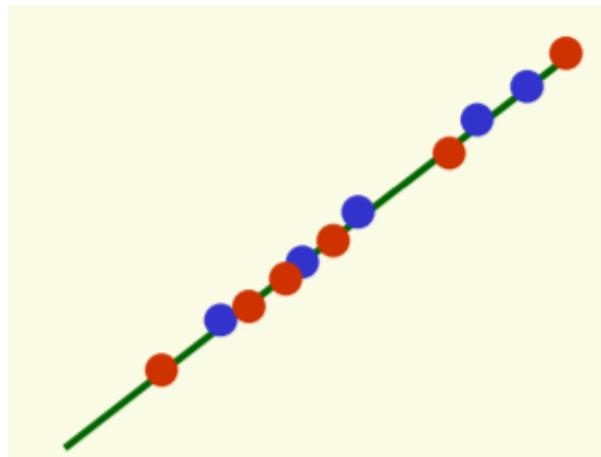
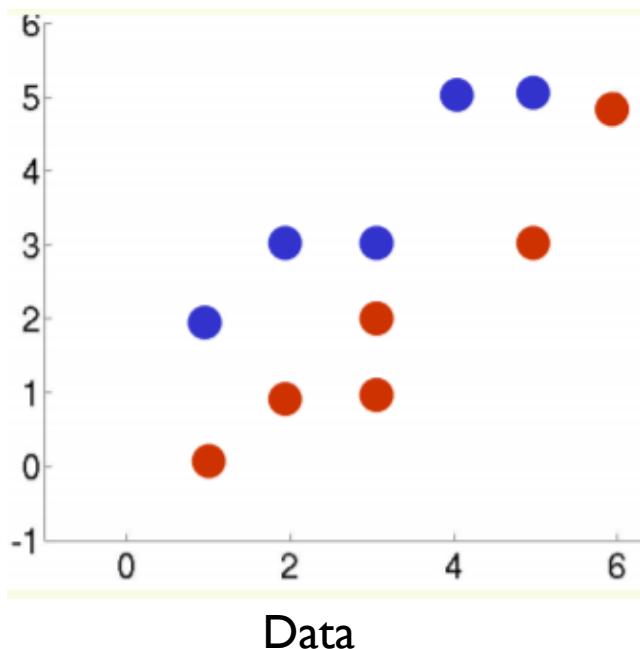
- Find  $\mathbf{W}$  that max

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_w \mathbf{W}|} \quad \text{The largest eigenvectors of } \mathbf{S}_w^{-1} \mathbf{S}_B$$

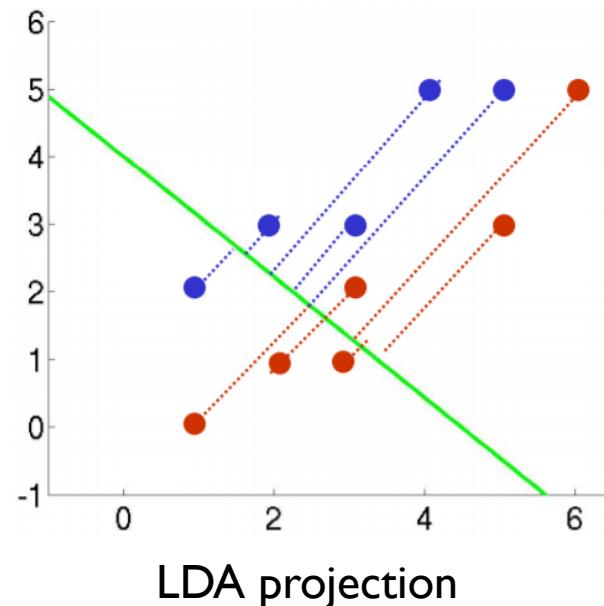
# FLD: Illustration



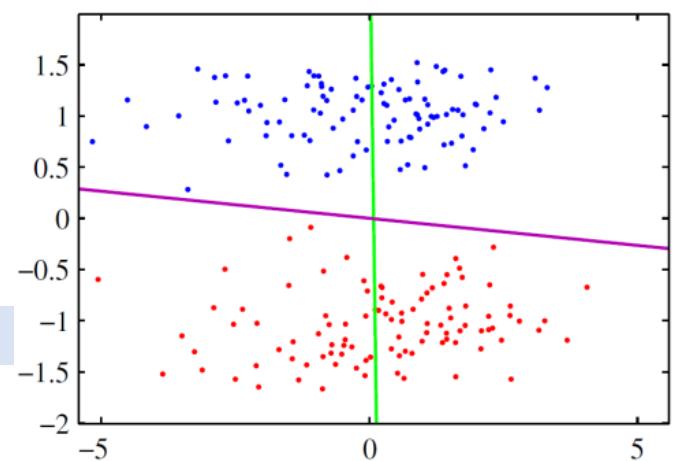
# LDA: Illustration



PCA projection



LDA projection



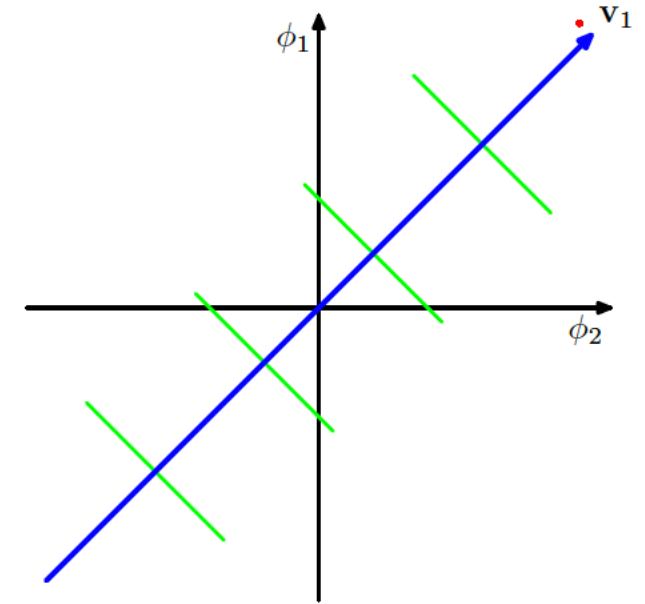
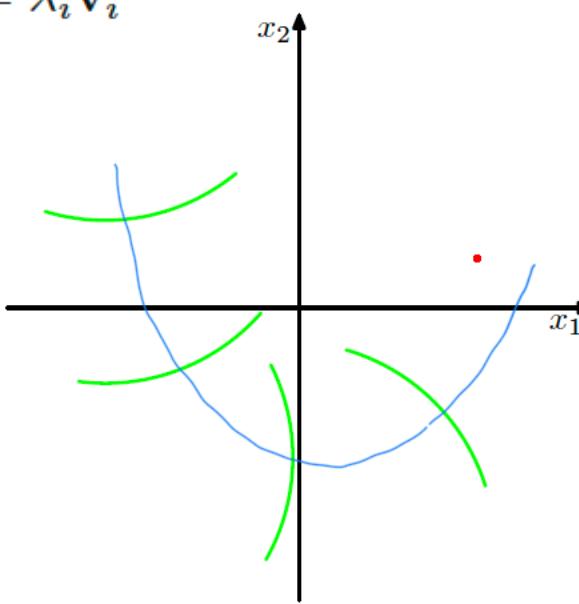
# Kernel PCA

- PCA
  - nonlinear transformation  $\phi(x)$  into an M-dimensional feature space, so that each data point  $x_n$  is thereby projected onto a point  $\phi(x_n)$ .

$$\mathbf{C} = \frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \quad \mathbf{C}\mathbf{v}_i = \lambda_i \mathbf{v}_i$$

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \{ \phi(\mathbf{x}_n)^T \mathbf{v}_i \} = \lambda_i \mathbf{v}_i$$

$$\mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n).$$



# Kernel PCA

- Standard PCA

- Nonlinear transformation  $\phi(x)$  into an M-dimensional feature space, so that each data point  $x_n$  is thereby projected onto a point  $\phi(x_n)$ . Aim is to find non-linear principal components using kernels without requiring to explicitly project to the M-dimensional feature space

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

$$\mathbf{C}\mathbf{v}_i = \lambda_i \mathbf{v}_i$$

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \{ \phi(\mathbf{x}_n)^T \mathbf{v}_i \} = \lambda_i \mathbf{v}_i$$

$$\mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n).$$

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T, \quad \mathbf{S}\mathbf{u}_i = \lambda_i \mathbf{u}_i$$

$$\frac{1}{N} \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \sum_{m=1}^N a_{im} \phi(\mathbf{x}_m) = \lambda_i \sum_{n=1}^N a_{in} \phi(\mathbf{x}_n).$$

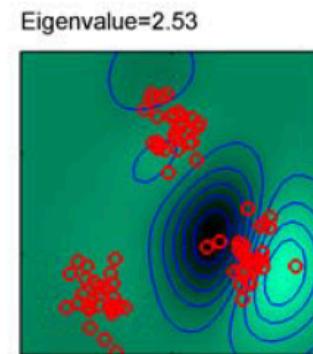
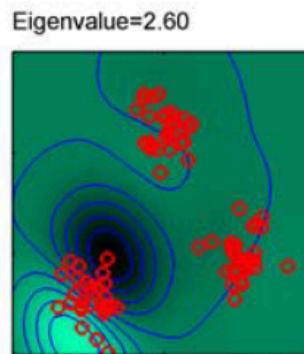
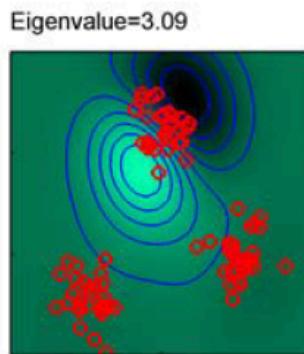
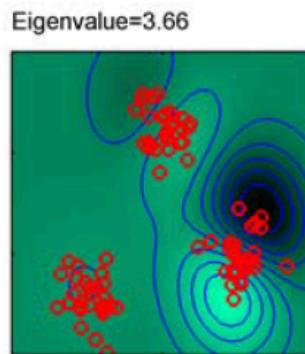
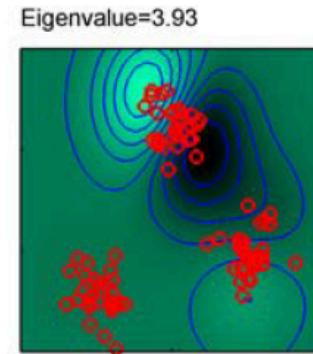
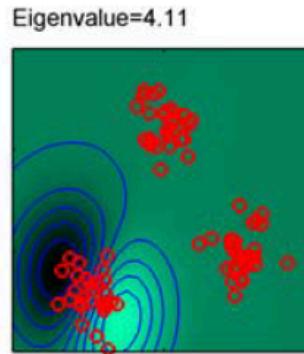
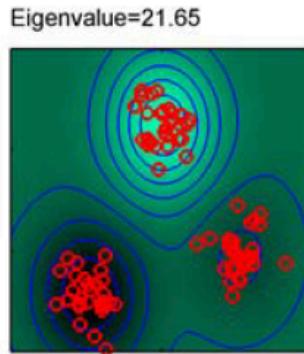
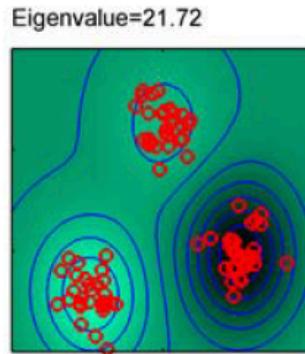
Assuming  $k(\mathbf{x}_n, \mathbf{x}_m) = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$

$$\frac{1}{N} \sum_{n=1}^N k(\mathbf{x}_l, \mathbf{x}_n) \sum_{m=1}^m a_{im} k(\mathbf{x}_n, \mathbf{x}_m) = \lambda_i \sum_{n=1}^N a_{in} k(\mathbf{x}_l, \mathbf{x}_n).$$

$$\mathbf{K}^2 \mathbf{a}_i = \lambda_i N \mathbf{K} \mathbf{a}_i \quad \mathbf{K} \mathbf{a}_i = \lambda_i N \mathbf{a}_i$$

$$y_i(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{v}_i = \sum_{n=1}^N a_{in} \phi(\mathbf{x})^T \phi(\mathbf{x}_n) = \sum_{n=1}^N a_{in} k(\mathbf{x}, \mathbf{x}_n)$$

# Kernel PCA



# Kernel Regression

- Many linear models for regression and classification can be reformulated in terms of a dual representation in which the kernel function arises

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\} \phi(\mathbf{x}_n) = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a} \quad a_n = -\frac{1}{\lambda} \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}.$$

# Kernel Regression

- Many linear models for regression and classification can be reformulated in terms of a dual representation in which the kernel function arises

$$J(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

$$\mathbf{w} = -\frac{1}{\lambda} \sum_{n=1}^N \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\} \phi(\mathbf{x}_n) = \sum_{n=1}^N a_n \phi(\mathbf{x}_n) = \Phi^T \mathbf{a}$$

$$a_n = -\frac{1}{\lambda} \{\mathbf{w}^T \phi(\mathbf{x}_n) - t_n\}.$$

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \Phi \Phi^T \Phi \Phi^T \mathbf{a} - \mathbf{a}^T \Phi \Phi^T \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \Phi \Phi^T \mathbf{a}$$

Gram matrix  $\mathbf{K} = \Phi \Phi^T$ ,  $K_{nm} = \phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m) = k(\mathbf{x}_n, \mathbf{x}_m)$

$$J(\mathbf{a}) = \frac{1}{2} \mathbf{a}^T \mathbf{K} \mathbf{K} \mathbf{a} - \mathbf{a}^T \mathbf{K} \mathbf{t} + \frac{1}{2} \mathbf{t}^T \mathbf{t} + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}. \quad \mathbf{a} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}.$$

$$y(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \mathbf{a}^T \Phi \phi(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t}$$

we determine the parameter vector  $\mathbf{a}$  by inverting an  $N * N$  matrix, whereas in the original parameter space formulation we had to invert an  $M * M$  matrix in order to determine  $\mathbf{w}$ .