

# Component Analysis and Discriminants

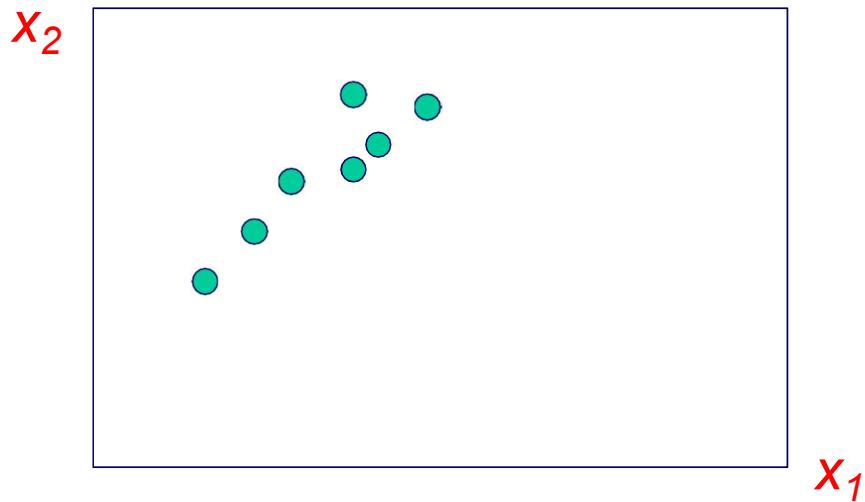
Reducing dimensionality when:

1. Classes are disregarded

Principal Component Analysis (PCA)

2. Classes are considered

Discriminant Analysis



# Component Analysis vs Discriminants

Two classical approaches for finding effective linear transformations

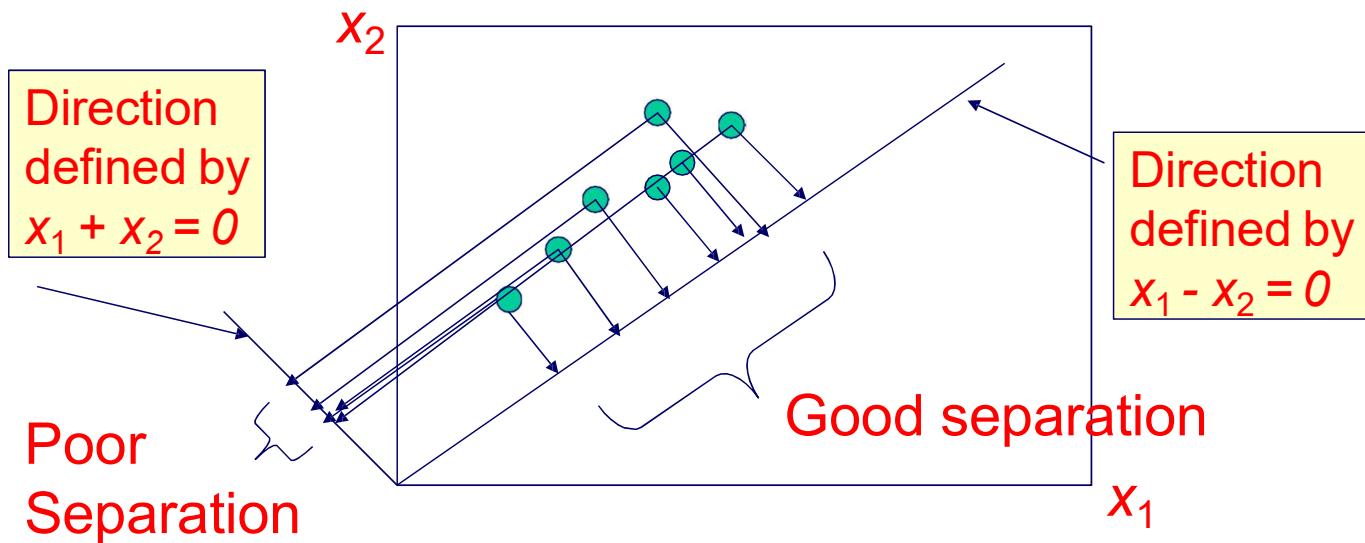
- PCA (Principal Component Analysis)  
“Projection that best **represents** data in least-square sense”
- MDA (Multiple Discriminant Analysis)  
“Projection that best **separates** the data in a least-squares sense”

# PCA: Linear Projections of Data

- Excessive dimensionality  $\mathbf{x} = [x_1, x_2, \dots x_d]$  causes
  - Computational difficulty and overfitting
  - Visualization issues
- Solution:
  - Combine features to reduce dimensionality
  - Linear combinations, e.g.,  $2x_1+3x_2+x_3$ 
    - are simple to compute and tractable
    - Project high dimensional data onto a lower dimensional space

# Projection to a lower dimensional space

- Allow computer to search for interesting directions



# Linear Projection

- Equation of a plane that passes through origin:  $x_1 + 2x_2 + 4x_3 = 0$

$$[1 \ 2 \ 4] \times \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = x_1 + 2x_2 + 4x_3 = 0$$

Can be written as  $\mathbf{w}^T \mathbf{x} = 0$  where  $\mathbf{w} = [1 \ 2 \ 4]^T$

- Projection of a point  $\mathbf{x}$  along a line with projection weights  $\mathbf{a}$  is given by:

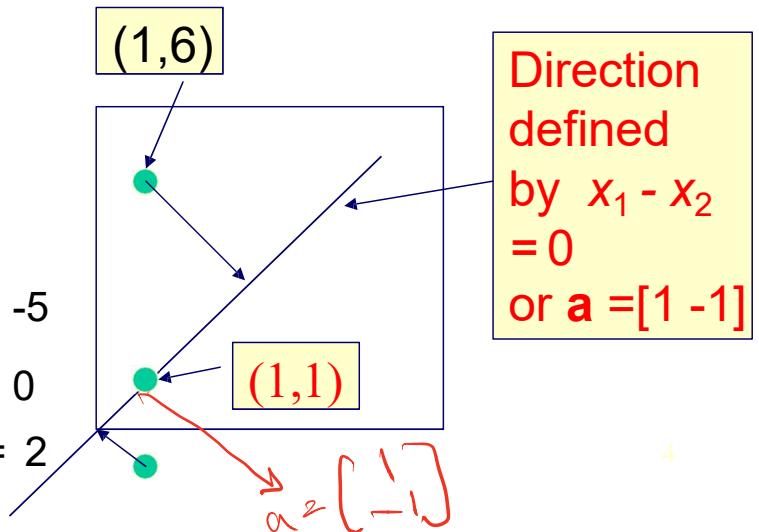
$$\mathbf{a}^T \mathbf{x} = \sum_{j=1}^d a_j x_j$$

- Example:

Projection of  $(1, 6)$  along  $[1 \ -1]$ :  $1 - 6 = -5$

Projection of  $(1, 1)$  along  $[1 \ -1]$ :  $1 - 1 = 0$

Projection of  $(1, -1)$  along  $[1 \ -1]$ :  $1 + 1 = 2$

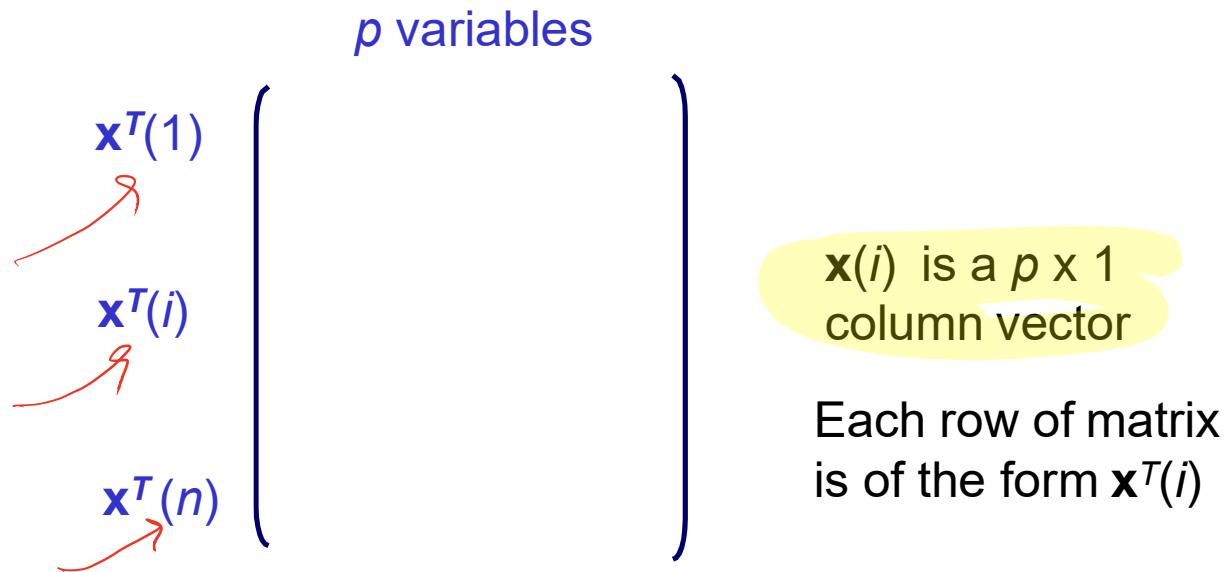


# Principal Component Analysis

- We will look at Data Matrix, Scatter Matrix and Covariance Matrix to arrive at best projection of data

# Data Matrix

Let  $\mathbf{X}$  be a  $n \times p$  *data matrix* of  $n$  samples



Assume  $\mathbf{X}$  is *centralized*, so that the mean  $m_i$  of each variable is subtracted for that variable  $\mathbf{m} = [m_1, m_2, \dots, m_p]^T$

$$C_{XX} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}(i) - \mathbf{m})(\mathbf{x}(i) - \mathbf{m})^T$$

# Scatter Matrix

$$\mathbf{S}_{p \times p} = \sum_{i=1}^n (\mathbf{x}(i) - \mathbf{m})(\mathbf{x}(i) - \mathbf{m})^T$$

- The Scatter Matrix  $\mathbf{S}$  is  $(n-1)$  times the sample Covariance Matrix
- Relationship Between Data Matrix and Scatter Matrix:
  - $\mathbf{S} = \mathbf{X}^T \mathbf{X}$  is the  $p \times p$  scatter matrix of the data since  $\mathbf{X}$  has zero mean

$$a^T x(1), \dots, a^T x(n)$$

## Projection

Let  $\mathbf{a}$  be a  $p \times 1$  column vector of projection weights that result in the largest variance when the data matrix  $\mathbf{X}$  are projected along  $\mathbf{a}$ .

The projection of any data vector  $\mathbf{x}$  is the linear combination:

$$\langle \mathbf{a}, \mathbf{x} \rangle = \mathbf{a}^T \mathbf{x} = \sum_{j=1}^p a_j x_j$$

$\begin{bmatrix} x(1) \\ \vdots \\ x(n) \end{bmatrix}_{n \times p} \cdot \begin{bmatrix} a^T \\ a^T \\ \vdots \\ a^T \end{bmatrix}_{p \times 1}$

Projected values of all data vectors in data matrix  $\mathbf{X}$  onto  $\mathbf{a}$  can be expressed as  $\mathbf{X}\mathbf{a}$  which is an  $n \times 1$  column vector.

# Variance along Projection

We define the Sample Variance along  $\mathbf{a}$  as the variance of the sample containing  $\mathbf{a}^T \mathbf{x}(i)$ 's

$$\text{Var}_{\mathbf{a}}(\mathbf{X}) \propto (\mathbf{X}\mathbf{a})^T (\mathbf{X}\mathbf{a}) = \mathbf{a}^T \underbrace{\mathbf{X}^T \mathbf{X}}_{\mathbf{S}} \mathbf{a}$$
$$= \mathbf{a}^T \mathbf{S} \mathbf{a}$$

where  $\mathbf{S}$  is the scatter matrix of  $\mathbf{S}$ .

$$\mathbf{X}\mathbf{a} = \begin{pmatrix} \mathbf{x}(1)^T \mathbf{a} \\ \mathbf{x}(2)^T \mathbf{a} \\ \vdots \\ \mathbf{x}(n)^T \mathbf{a} \end{pmatrix}$$
$$\text{Var}(\text{sample}) = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}(i)^T \mathbf{a})^2$$
$$= \frac{1}{n-1} \sum_{i=1}^n (\mathbf{X}\mathbf{a})^T (\mathbf{X}\mathbf{a})$$
$$= \frac{1}{n-1} \|\mathbf{X}\mathbf{a}\|^2$$

# Variance along Projection

Therefore, the variance is a function of both  $\mathbf{X}$  and  $\mathbf{a}$ .

Maximizing variance along  $\mathbf{a}$  is not well-defined since we can increase it without limit by increasing the size of the components of  $\mathbf{a}$ .

Find a subject to

$$\|\mathbf{a}\|^2 = 1$$

# Optimization Problem

Impose a normalization constraint on  $\mathbf{a}$  i.e,

$$\mathbf{a}^T \mathbf{a} = 1$$

Optimization problem is to **maximize**

$$u(\mathbf{a}) = \mathbf{a}^T \mathbf{S} \mathbf{a}$$

Variance  
Criterion

Subject to

$$\mathbf{a}^T \mathbf{a} = 1$$

Normalization  
Criterion

# Optimization Problem

Impose a normalization constraint on  $\mathbf{a}$  i.e,

$$\mathbf{a}^T \mathbf{a} = 1$$

Optimization problem is to **maximize**

$$u(\mathbf{a}) = \mathbf{a}^T \mathbf{S} \mathbf{a} - \lambda (\mathbf{a}^T \mathbf{a} - 1)$$

Variance Criterion

Normalization Criterion

Where  $\lambda$  is a Lagrange multiplier.

Solution: Differentiating wrt  $\mathbf{a}$  yields

$$\nabla_{\mathbf{a}} u(\mathbf{a}) = 0 \quad 2 \mathbf{S} \mathbf{a} - 2\lambda \mathbf{a} = 0$$

which reduces to

Characteristic Equation of  $\mathbf{S}$ !

$$(\mathbf{S} - \lambda \mathbf{I}) \mathbf{a} = 0$$

Candidate  $\mathbf{a}$ 's : eigenvalues of  $\mathbf{S}$

# Optimization Problem

Impose a normalization constraint on  $\mathbf{a}$  i.e,

$$\mathbf{a}^T \mathbf{a} = 1$$

Optimization problem is to **maximize**

$$u(\mathbf{a}) = \mathbf{a}^T \mathbf{S} \mathbf{a} - \lambda (\mathbf{a}^T \mathbf{a} - 1)$$

Variance  
Criterion

Where  $\lambda$  is a Lagrange multiplier.

Normalization  
Criterion

Solution: Differentiating wrt  $\lambda$  yields

$$\mathbf{a}^T \mathbf{a} - 1 = 0 \Rightarrow \mathbf{a}^T \mathbf{a} = 1$$

⇒ Orthonormal eigenvalues of  $\mathbf{S}$   
are solution candidates

# Characteristic Equation

Given a  $p \times p$  matrix  $\mathbf{M}$ , a very important class of linear equations is of the form

$$\mathbf{M}_{p \times p} \mathbf{x}_{p \times 1} = \lambda \mathbf{x}_{p \times 1}$$

which can be rewritten as  $(\mathbf{M} - \lambda \mathbf{I})\mathbf{x} = 0$ .

$$x \in N(\mathbf{M} - \lambda \mathbf{I})$$

If  $\mathbf{M}$  is real and symmetric there are  $p$  possible linearly independent and **orthonormal** solution vectors (vectors  $\mathbf{x}$  that satisfy the characteristic equation) called eigenvectors,  $\mathbf{e}_1, \dots, \mathbf{e}_p$  and associated eigenvalues  $\lambda_1, \dots, \lambda_p$ , for which  $\det(\mathbf{M} - \lambda \mathbf{I}) = 0$ .

Characteristic eq. of  $M$

# Principal Components

If the matrix  $\mathbf{M}$  is the Scatter matrix  $\mathbf{S}$ , the Characteristic Equation is:

$$(\mathbf{S} - \lambda \mathbf{I})\mathbf{a} = 0$$
$$\det(\mathbf{S} - \lambda \mathbf{I})=0$$

The roots are eigenvalues of  $\mathbf{S}$ .  
Therefore, the candidate solutions to the optimization problem are eigenvectors of  $\mathbf{S}$ , and their corresponding values.

# Principal Components

If the matrix  $\mathbf{M}$  is the Scatter matrix  $\mathbf{S}$ , the Characteristic Equation is:

$$(\mathbf{S} - \lambda \mathbf{I})\mathbf{a} = 0$$
$$\det(\mathbf{S} - \lambda \mathbf{I}) = 0$$

The roots are eigenvalues of  $\mathbf{S}$ .

The corresponding eigenvectors are called the principal components.

First principal component is the eigenvector associated with the largest eigenvalue of  $\mathbf{S}$ .

# Variance along each eigenvector

The variance along each eigenvector  $\mathbf{e}_j$  is calculated as:

$$\text{Var}_{\mathbf{e}_j}(\mathbf{X}) \propto (\mathbf{X}\mathbf{e}_j)^T (\mathbf{X}\mathbf{e}_j)$$
$$= \mathbf{e}_j^T \mathbf{X}^T \mathbf{X} \mathbf{e}_j = \mathbf{e}_j^T S \underbrace{\mathbf{e}_j}_{\lambda_j \mathbf{e}_j} = \lambda_j \mathbf{e}_j^T \mathbf{e}_j$$
$$\approx \lambda_j$$

# Variance along each eigenvector

- The variance along each eigenvector  $\mathbf{e}_j$  is the eigenvalue associated with that eigenvector. Therefore, the solution to the previous optimization problem is  $\mathbf{e}_{\max}$ , the eigenvector associated with the largest eigenvalue,  $\lambda_{\max}$ .
- The maximum variance is  $\lambda_{\max}$ .

# Variance along each eigenvector

- If we want to find the direction that is *constrained to be orthogonal* to  $\mathbf{e}_{\max}$  which has the highest variance, we have  $p-1$  candidates: the other ~~recters~~ eigenvalues of  $\mathbf{S}$ .
- Obviously, it would be the direction of the eigenvector that has the second largest eigenvalue.

# Variance along each eigenvector

- If we want to find the direction that is *constrained to be orthogonal* to both previous directions, which has the highest variance, we have  $p-2$  candidates: the other eigenvalues of **S**.
- Obviously, it would be the direction of the eigenvector that has the third largest eigenvalue.

# Variance along each eigenvector

- What does the orthogonality constraint do?
- Assume that  $\mathbf{a}$  and  $\mathbf{b}$  are two vectors of  $X$  that are orthogonal, i.e.  $\mathbf{a}^T \mathbf{b} = 0$ . We show that the samples of data points projected along  $\mathbf{a}$  and  $\mathbf{b}$  are uncorrelated.

$$X_a = \begin{bmatrix} x^{(1)} a \\ \vdots \\ x^{(N)} a \end{bmatrix}$$

$$X_b = \begin{bmatrix} x^{(1)} b \\ \vdots \\ x^{(N)} b \end{bmatrix}^\top$$

$$\begin{aligned}
 & \text{Trace}(\cancel{(Xa)^T(Xb)}) = \text{Tr}(a^T X^T X b) \\
 &= \sum_i (\cancel{a(i)a})(\cancel{a(i)b}) \quad \text{not needed} \\
 &= \text{Tr}(a^T S b) = \text{Tr}(S b a^T) \\
 &= \text{Tr}(S) \text{Tr}(b a^T) = \text{Tr}(S) \text{Tr}(\cancel{a^T b}) \\
 &\quad \text{This is not true!}
 \end{aligned}$$

We must show that

$$\sum_i (\cancel{a(i)a}) (\cancel{a(i)b}) = \sum_i (\cancel{a(i)} a) (\cancel{a(i)} b)$$

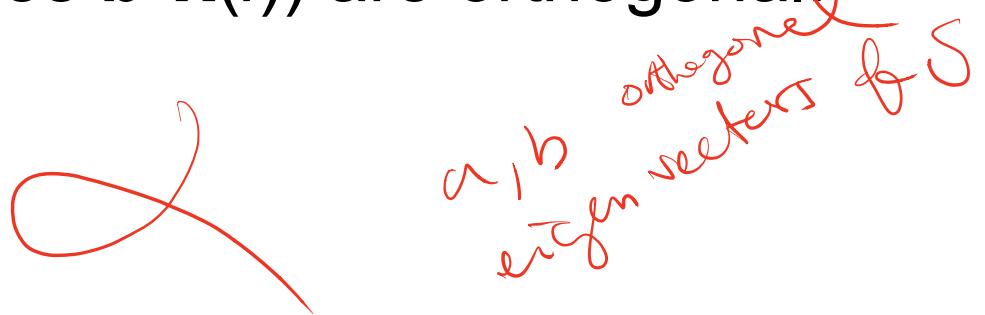
$(Xa)^T(Xb) = 0$  which is equivalent to  $\sum_i (\cancel{a(i)} a) (\cancel{a(i)} b)$

$\Rightarrow (Xa)^T(Xb) = a^T X^T X b$

but  $a^T X^T X b = a^T S b = a^T \cancel{b} b = \cancel{a^T b} = 0$

# Variance along each eigenvector

- This means  $\mathbf{Xa}$  (vector of the samples  $\mathbf{a}^T \mathbf{x}(i)$ ) and  $\mathbf{Xb}$  (vector of the samples  $\mathbf{b}^T \mathbf{x}(i)$ ) are orthogonal.



# Variance along each eigenvector

- We have  $p$  linearly independent and orthogonal eigenvectors, so we can continue doing this and find  $p$  directions that have the highest variance subject to being orthogonal to all other directions.

# Variance along each eigenvector

- This means the samples obtained from projection of the data points along each direction are **uncorrelated** from one another!

# Principal Directions

- Therefore, we are looking for directions that have the highest possible variance, which give us projected  $\mathbf{x}(i)$ 's that are uncorrelated from one another.
- Those directions are called **principal directions or principal component loading vectors.**

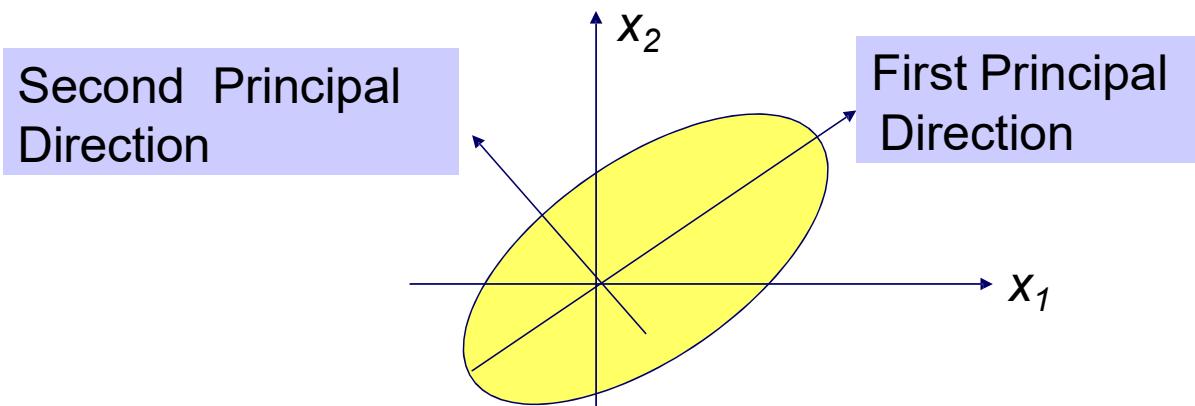
# First Principal Direction

~~direction~~

- First principal component is the eigenvector associated with the largest eigenvalue of  $\mathbf{S}$ .
- It is the direction along which  $\mathbf{X}$  has the largest variation.
- Second Principal direction is:
  - In a direction orthogonal to first.  
Therefore, it is uncorrelated with the first principal direction.
  - Has second largest eigenvalue.

# Other Principal Directions

- All principal directions are orthogonal to one another, so they are uncorrelated.
- The variation along each principal direction depends on how large its eigenvalue is.



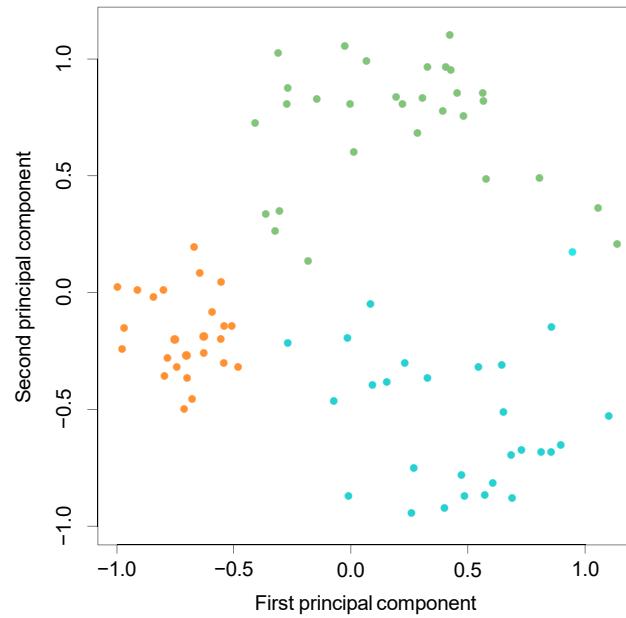
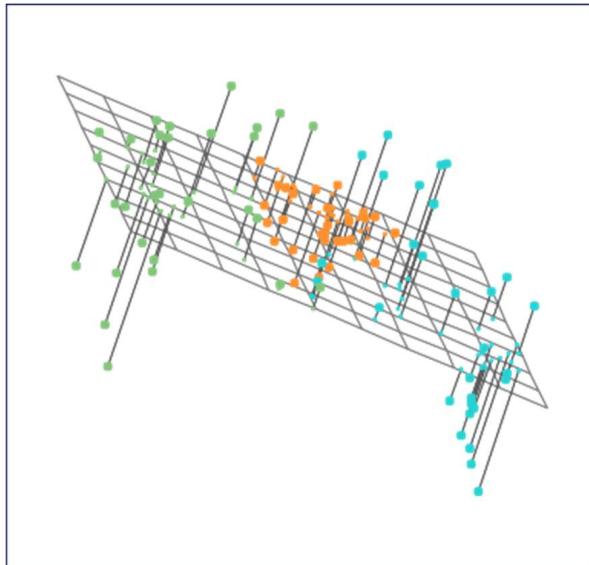
# Principal Components

- The component of each data point along each principal direction is called a **principal components** or **principal component score** of that data point.
- We call the component of each data point along the  $i^{\text{th}}$  principal direction, its  $i^{\text{th}}$  principal component.
- Principal components are found by projecting data points onto principal directions.

# Alternative Derivation of PCA

- DHS Text gives another way to derive the fact that eigenvectors of scatter matrix are principal components
- By setting up the squared error criterion function of the data and the vector  $\mathbf{e}$  that minimizes it satisfies the characteristic equation
- Can be generalized to from one-dimensional projection to a dimension  $p'$  which are the eigenvectors of  $\mathbf{S}$  forming the principal components

# Another Interpretation of Principal Components



The first three principal components of a data set span the three-dimensional hyperplane that is closest to the  $n$  observations,

# PCA find the hyperplane closest to the observations

- The first principal component **loading vector** (direction) has a very special property: it defines the line in  $p$ -dimensional space that is **closest** to the  $n$  observations (using average squared Euclidean distance as a measure of closeness)

# PCA find the hyperplane closest to the observations

- The notion of principal directions as the dimensions that are closest to the  $n$  observations extends beyond just the first principal direction.
- For instance, the first two principal directions of a data set span the plane that is closest to the  $n$  observations, in terms of average squared Euclidean distance.

# Projection into $k$ Eigenvectors

Variance of data projected into first  $k$  eigenvectors is proportional to

$$\sum_{j=1}^k \lambda_j$$

That is because:

$$\text{Var}_{\mathbf{e}_j}(\mathbf{X}) \propto (\mathbf{X}\mathbf{e}_j)^T (\mathbf{X}\mathbf{e}_j)$$

$$= \lambda_j$$

$$\mathbf{e}_j^T \mathbf{X}^T \mathbf{X} \mathbf{e}_j = \mathbf{e}_j^T S \mathbf{e}_j$$

# Projection into $k$ Eigenvectors

and:

$$\text{Var}_{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_k}(\mathbf{X}) \propto$$

Projections along eigenvectors:

$$T = [x_{e_1} \quad x_{e_2} \quad \dots \quad x_{e_k}]$$
$$= X \mathcal{V}_k$$
$$\mathcal{V}_k (\mathbf{e}_1 - \dots - \mathbf{e}_k)$$

$$\begin{aligned}
 &= \text{Var of Samples} \times \sum_{i,j} T_{ij}^2 \\
 &= \|T\|_F^2 = (\sqrt{\text{Tr}(T^T T)})^2 = \text{Tr}(T^T T)
 \end{aligned}$$

$$\begin{aligned}
 &\approx \text{Tr}[(XV)^T(XV)] \\
 &\approx \text{Tr}(V^T X^T X V) = \text{Tr}(V_k^T S V_k)
 \end{aligned}$$

$$\begin{aligned}
 &\text{Tr}([v_1 \dots v_k]^T [sv_1 \dots sv_k]) \\
 &\approx \text{Tr}([v_1 \dots v_k]^T [\lambda_1 v_1 \dots \lambda_k v_k])
 \end{aligned}$$

$$\approx \text{Tr}\left(\left[\begin{array}{c} v_1^T \\ \vdots \\ v_k^T \end{array}\right] [\lambda_1 v_1 \dots \lambda_k v_k]\right)$$

$$= \text{Tr} \left( \begin{bmatrix} \underbrace{\lambda_1 u_1 u_1^T}_{\lambda_1} & \underbrace{\lambda_2 u_2 u_2^T}_{\lambda_2} & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \ddots & \lambda_k u_k u_k^T \end{bmatrix} \right)$$

$$= \sum_{i=1}^k \lambda_i$$

If we use all eigenvalues

Var  $\underset{\text{e-ep}}{\longrightarrow}$   $\sum_{k=1}^p \lambda_k$

On the other hand

Var  $\underset{\text{e-ep}}{\longrightarrow}$   $\text{Tr}(V^T S V)$

$$\geq \text{Tr}(V V^T S) = \text{Tr}(S)$$

$V \in \mathbb{R}^{P \times P}$

$$\text{Tr}(X^T X) = \sum_{ij} X_{ij}^2 = \sum_{ij} X_{ij}^2$$

$\text{Var}(\text{all } g_i)$  =  $\sum_{j=1}^p \text{Var}(X_j)$

remember that Variance  
has a normalization factor  $n-1$

Variance of the whole data

= Variance of the principal  
components of  $\sum_{i=1}^p \lambda_i$

# Projection into $k$ Eigenvectors

- Variance of data projected into first  $k$  eigenvectors is proportional to  
$$\sum_{j=1}^k \lambda_j$$

Variance of First PC  
The whole Variance  
Variance of  $k$  PC  
Variance of PPC
- The proportion of variance lost in approximating true data matrix  $\mathbf{X}$  using only first  $k$  eigenvectors is  
$$\frac{\sum_{j=k+1}^p \lambda_j}{\sum_{l=1}^p \lambda_l}$$

Usually 5 to 10 principal components capture 90% of variance in the data

$$\frac{\sum_{j=1}^k \lambda_j}{\sum_{j=1}^{n-1} \lambda_j}$$

# How Many Principal Components Should One Use?

When Using PCA for classification, one can consider the number of components used as a hyper parameter and use cross-validation to select the number of components.

# How Many Principal Components Should One Use?

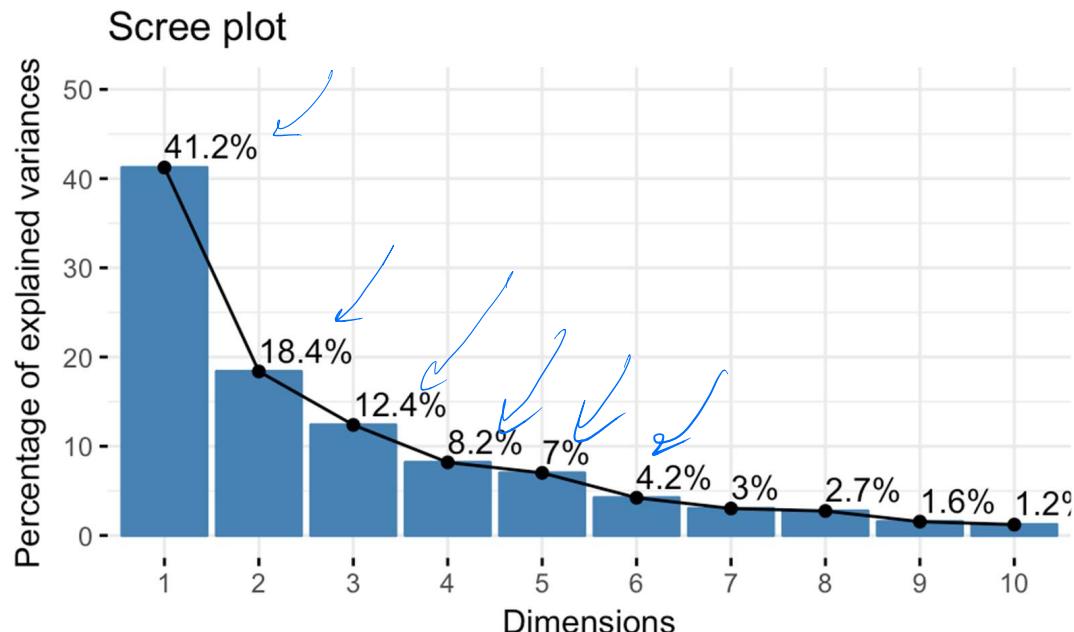
An alternative method to determine the number of principal components merely from unlabeled  $\mathbf{X}$  is to look at a **Scree Plot**, which is the plot of eigenvalues ordered from largest to the smallest.

# How Many Principal Components Should One Use?

- The number of component is determined at the point, beyond which the remaining eigenvalues are all relatively small and of comparable size (Jollife 2002, Peres-Neto, Jackson, and Somers (2005)).

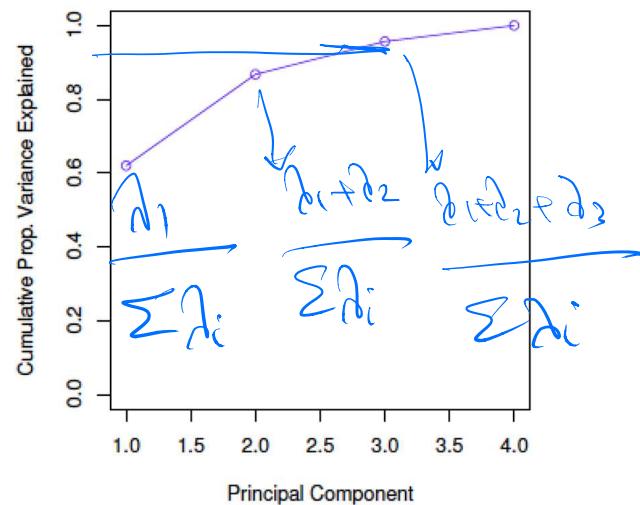
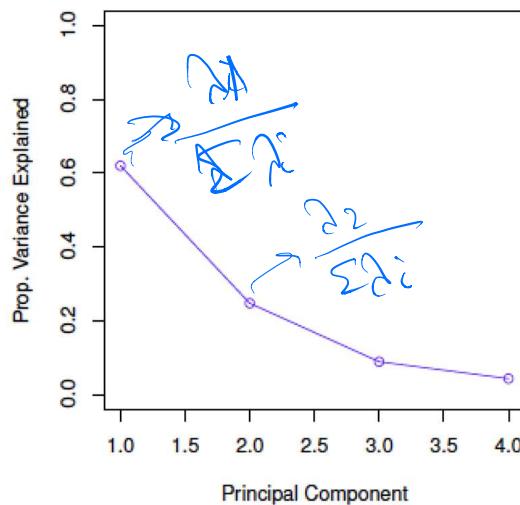
# How Many Principal Components Should One Use?

## Scree plot



# Proportion Variance Explained: continued

- The PVEs sum to one. We sometimes display the cumulative PVEs.



# PCA and SVD

- Recall that *any real matrix*  $\mathbf{X}_{n \times p}$  can be decomposed into the following form:

$$\mathbf{U}\Sigma\mathbf{V}^T$$

$n \times n$   $n \times p$   $p \times p$

which is called the Singular Value Decomposition (SVD) of  $\mathbf{X}$ .

$\mathbf{U}_{n \times n}$  and  $\mathbf{V}_{p \times p}$  are orthogonal matrices.

$$\mathbf{U}^T \mathbf{U} = \mathbf{I} \quad \mathbf{V}^T \mathbf{V} = \mathbf{I}$$

# PCA and SVD

- The matrix  $\Sigma$  is of the form

$$\Sigma$$

where  $\sigma_1, \dots, \sigma_r$  are called the singular values of  $\mathbf{X}$  and  $r = \text{rank}(\mathbf{X})$

Obviously,  $r \leq \min(n, p)$

# PCA and SVD

- Defining the matrix  $\Sigma'$  as:

$$\Sigma' = \text{diag}(\sigma_1, \dots, \sigma_r)$$

We can rewrite the SVD in its  
*compact form*

$$V = [v_1 \ v_2]$$

$$[U_1 \ U_2] \begin{bmatrix} \Sigma' & 0 \\ 0 & 0 \end{bmatrix} V^T$$

$$\begin{bmatrix} v_1^T \\ v_2^T \end{bmatrix} = U_1 \Sigma' V^T$$

# PCA and SVD

- On the other hand, we know the principal components are the eigenvalues of  $\mathbf{X}^T \mathbf{X}$ . Plugging in  $\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^T$  and  $\mathbf{X} = \mathbf{U}_1 \Sigma' \mathbf{V}_1^T$ :

$$\mathbf{X}^T \mathbf{X} = (\mathbf{U} \Sigma \mathbf{V}^T)^T (\mathbf{U} \Sigma \mathbf{V}^T)$$

$$\Rightarrow \mathbf{V} \Sigma^T \underbrace{\mathbf{U}^T \mathbf{U}}_{\mathbf{I}} \Sigma \mathbf{V}^T = \mathbf{V} \Sigma^T \Sigma \mathbf{V}^T$$
$$= \mathbf{V} \left[ \frac{\Sigma'^2}{\sigma_1^2} + \dots + \frac{\sigma_n^2}{\sigma_1^2} \right] \mathbf{V}^T$$
$$= \mathbf{V} \Sigma'^2 \mathbf{V}^T$$

$\lambda_{\text{eig}}$  are eigenvalues of  $X^T X$

but we showed that they are

also  $\sigma_i^{-2}$ .

---

# PCA and SVD

- Meaning that right singular vectors  $\mathbf{V}$  are principal directions and that squared singular values are proportional to the variance of data along each principal component:

Columns of  $\mathbf{V}$  are the eigenvectors of  $\mathbf{X}^T \mathbf{X}$ , which means they are principal directions.

Projections of  $x(i)$ 's on  
the principal directions are

principal components

$$XV = U\Sigma V^T V = U\Sigma \underbrace{\Sigma}_{n \times n} \underbrace{V^T V}_{n \times p} \quad n \times r$$

$$\cancel{XV_i} = U_i \underbrace{\Sigma}_{n \times r} \underbrace{V_i^T V_i}_{r \times r} = U_i \underbrace{\Sigma'}_{n \times r} \quad r \times r$$

$V$   $r \times p$  Principal components

# PCA and SVD

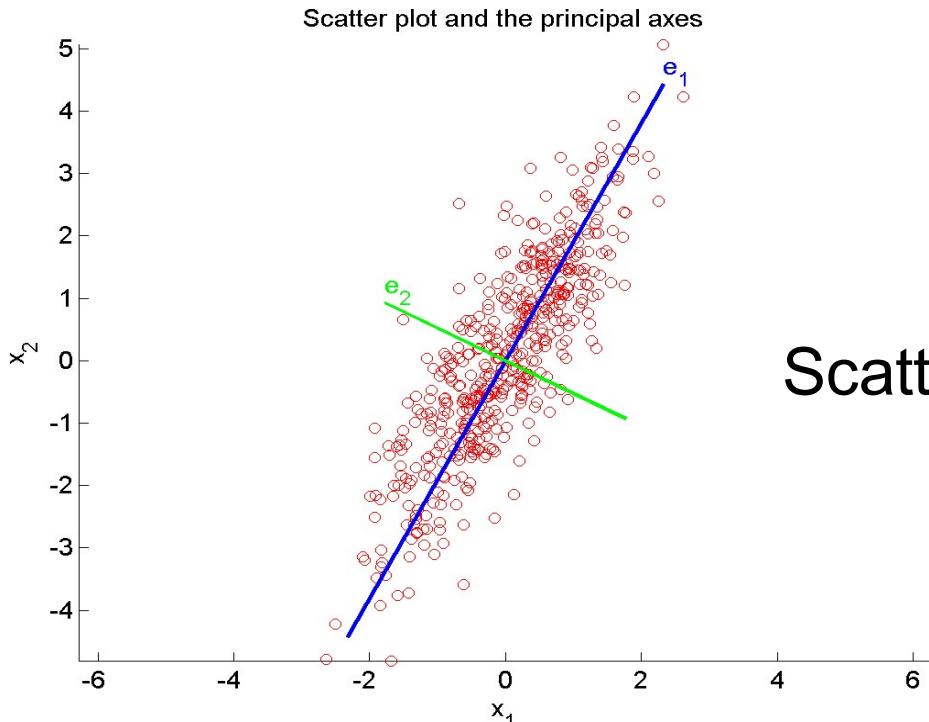
- The projections of data points onto principal directions, i.e. principal components/scores are obtained from  $\mathbf{X}\mathbf{V}_1$  (or  $\mathbf{X}\mathbf{V}$ ):

$$\mathbf{U}, \Sigma \quad \xleftarrow{\quad} \quad \mathbf{U}, \Sigma$$

## PCA and SVD

- Note that the number of meaningful principal components (with nonzero  $\lambda_i$ ) depends on the rank of the matrix  $\mathbf{X}$ .
- One can use the full SVD and find  $p$  directions, but if the rank of the data matrix is less than  $p$  [ $r \leq \min(p,n)$ ] then  $n-r$  principal directions will have a zero eigenvalues. How do you interpret that?

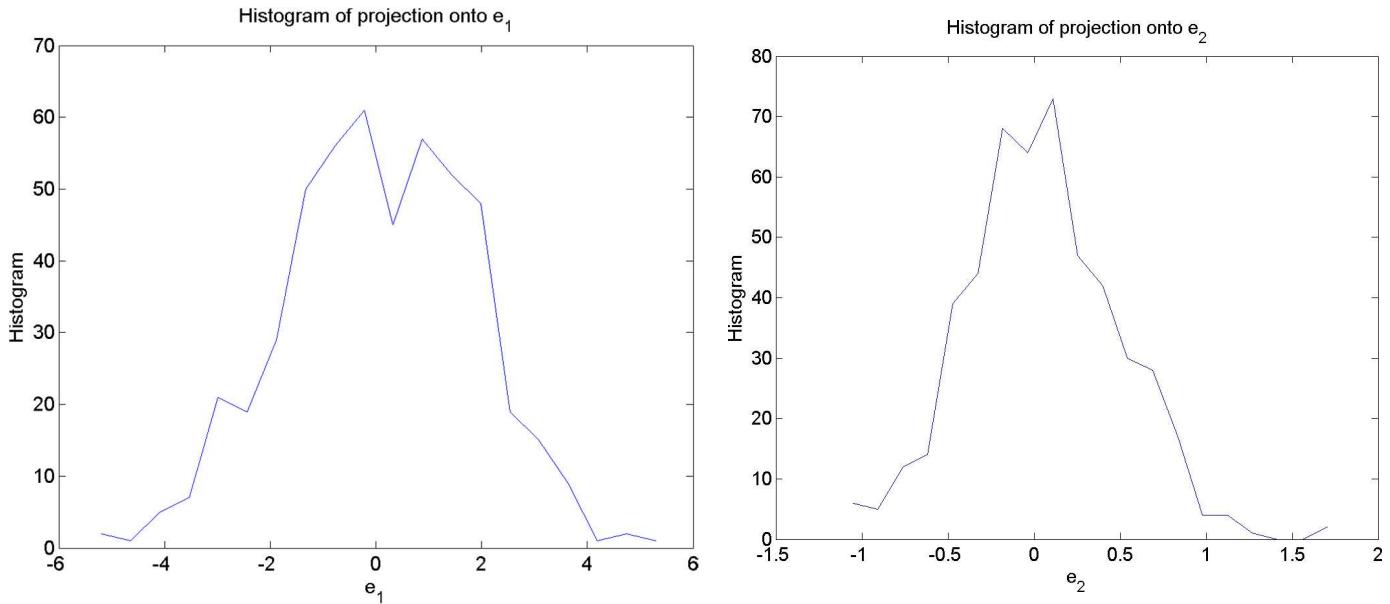
# Examples



Scatter plot.

Scatter plot (red dots) and the principal axes for a bivariate sample. The blue line shows the axis  $e_1$  with the greatest variance and the green line shows the axis  $e_2$  with the smallest variance. Features are now uncorrelated.

# Examples



Projection onto  $e_1$ .

(c) Projection onto  $e_2$ .

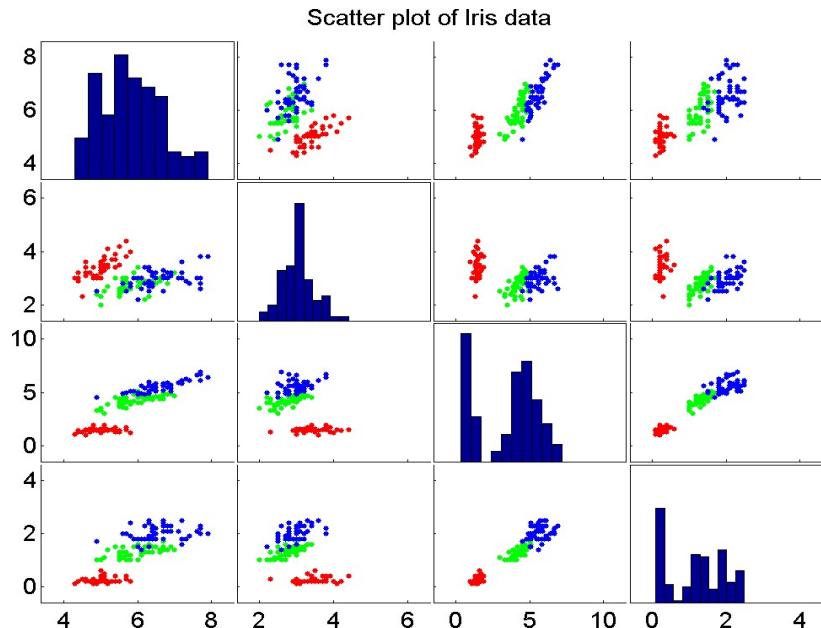
4 variables

3 species

50 samples/class

# Examples

- Setosa
- Versicolor
- Virginica



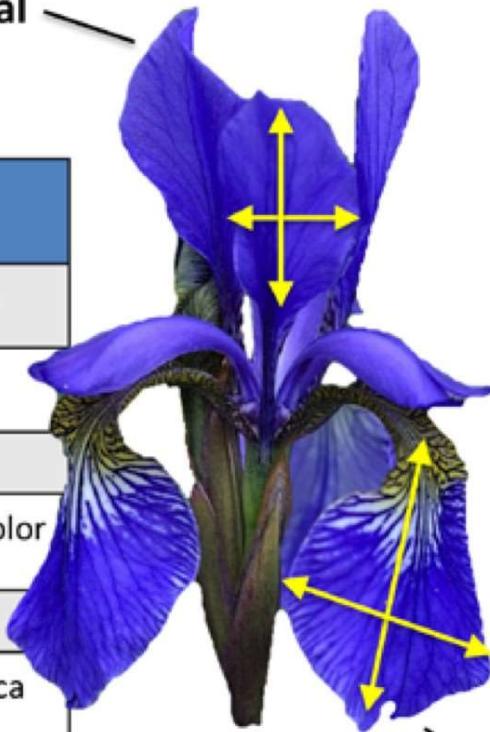
Scatter plot of the iris data. Diagonal cells show the histogram for each feature. Other cells show scatters of pairs of features  $x_1, x_2, x_3, x_4$  in top-down and left-right order. Red, green and blue points represent samples for the setosa, versicolor and virginica classes, respectively.

## Samples

(instances, observations)

	Sepal length	Sepal width	Petal length	Petal width	Class label
1	5.1	3.5	1.4	0.2	Setosa
2	4.9	3.0	1.4	0.2	Setosa
...					
50	6.4	3.5	4.5	1.2	Versicolor
...					
150	5.9	3.0	5.0	1.8	Virginica

Petal



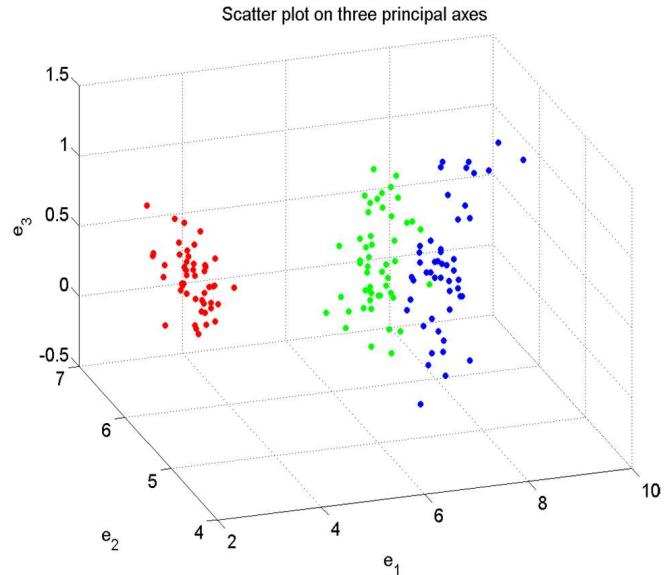
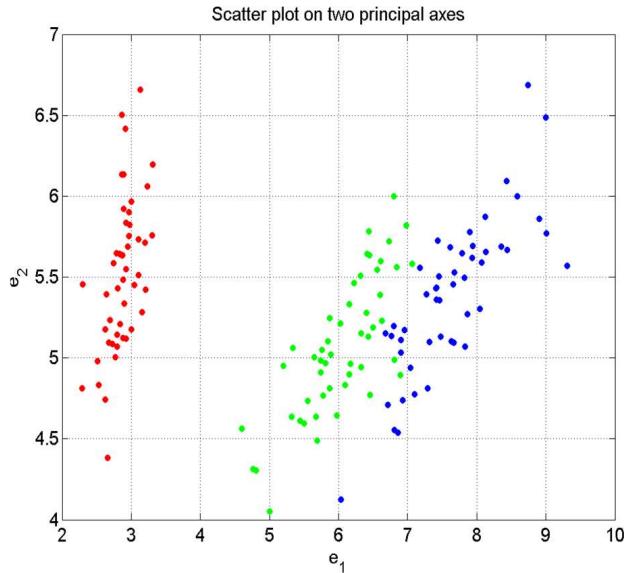
Sepal

Features

(attributes, measurements, dimensions)

Class labels  
(targets)

# Examples



Scatter plot of the projection of the iris data onto the first two and the first three principal axes. Red, green and blue points represent samples for the setosa, versicolor and virginica classes, respectively.