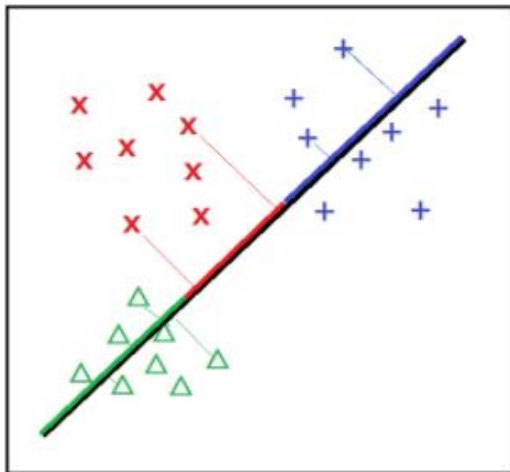# Principal Component Analysis, and

# Fisher Linear Discriminant

Dimensionality reduction / feature extraction

# Principal Component Analysis

- Originated from the work by Pearson(1901).

- Its purpose is to derive new features (variables) in the decreasing order of importance.

- Dimensionality can be reduced without loosing much information and structure present in the data.

# PRINCIPAL COMPONENT ANALYSIS [PCA]

* A method to reduce the dimensionality of the data.

* PCA seeks a projection that best represents the data in a least squares sense.

* In the new space data is so represented that the features become uncorrelated.

* In the new space distributions might become a simple one.

# PRINCIPAL COMPONENT ANALYSIS [PCA]

* A method to reduce the dimensionality of the data.

* PCA seeks a projection that best represents the data in a least squares sense.

* In the new space data is so represented that the features become uncorrelated.

* In the new space distributions might become a simpler one.

## Short comings

* Needs to find covariance matrix and its eigen vectors

* Discriminating components between classes might be lost

objective :-

Let $\mathcal{D} = \{ X_1, \ldots, X_n \}$ be the set of patterns of dimensionality $d$

We want to find $\mathcal{D}' = \{ X_1', X_2', \ldots, X_n' \}$ where each $X_i'$ is of dim. $d'$ such that $d' < d$, and

$$J = \sum_{i=1}^{n} \| X_i - X_i' \|^2$$ should be minimum possible one.

What is zero dim. projection for the data?

i.e., we want to represent the data set by just one pattern ($x_0$).

$$J = \sum_i \| X_i - X_0 \|^2$$

What is zero dim. projection for the data?

i.e., we want to represent the data set by just one pattern $(x_0)$.

$$J = \sum_i \| X_i - X_0 \|^2$$

It is easy to see that $J$ is minimized when $X_0$ is sample mean (centroid).

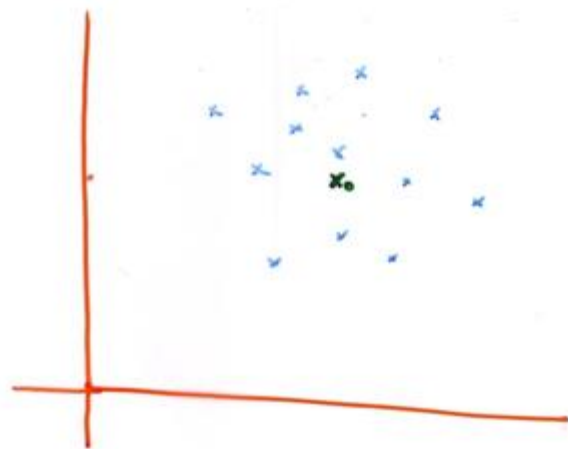i.e., $X_0 = \frac{1}{n} \sum_i X_i$

What is zero dim. projection for the data?

i.e., we want to represent the data set by just one pattern $(x_0)$.

$$J = \sum_i \| x_i - x_0 \|^2$$

It is easy to see that $J$ is minimized when $x_0$ is sample mean (centroid).

i.e., $x_0 = \frac{1}{n} \sum_i x_i$
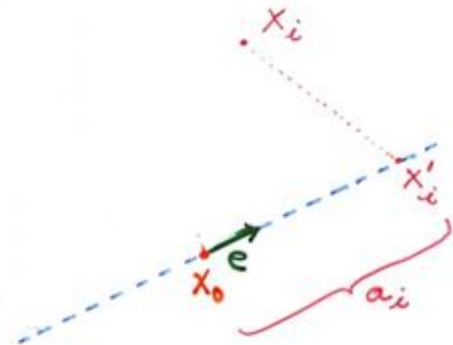
# What is good 1-Dim Representation

Let the data is projected onto a line passing through the centroid $(X_o)$

# What is good 1-Dim Representation

Let the data is projected onto a line
passing through the centroid $(X_0)$

Let $X_i' = \text{Proj}(X_i)$

$\vec{e}$ is unit vector
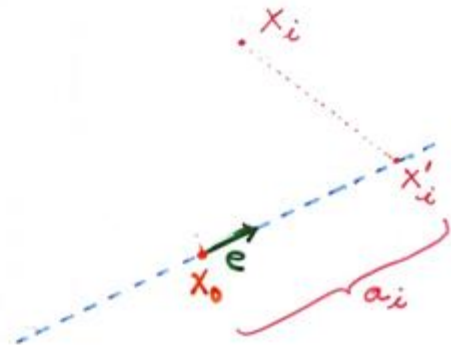in the direction of
the line.

Then

$$X_i' = X_0 + a_i\, e$$

→ scalar
dist between $X_0$ & $X_i'$

# What is good 1-Dim Representation

Let the data is projected onto a line
passing through the centroid $(X_0)$

Let $X_i' = Proj(X_i)$

$\vec{e}$ is unit vector
in the direction of
the line.



Then

$$X_i' = X_0 + a_i e$$

→ scalar
dist. between $X_0$ & $X_i'$

$$J = \sum_i \| (X_0 + a_i e) - X_i \|^2$$

We need to find $a_i$ & $e$

$$J = \sum_i \left\| (x_0 + a_i e) - x_i \right\|^2$$

$$= \sum_i \left\| a_i e - (x_i - x_0) \right\|^2$$

$$= \sum a_i^2 \|e\|^2 - 2 \sum a_i e^T (x_i - x_0) + \sum \|x_i - x_0\|^2$$

$$= \sum_i a_i^2 - 2 \sum_i a_i e^T (x_i - x_0) + \sum_i \|x_i - x_0\|^2$$

Now, To find $a_j$

$$\frac{\partial J}{\partial a_j} = 2 a_j - 2 e^T (x_j - x_0) = 0$$

$$a_j = e^T (x_j - x_0)$$

To find $a_1, a_2, \ldots, a_n$

$$J = \sum_i \left\| (x_0 + a_i e) - x_i \right\|^2$$

$$= \sum_i \left\| a_i e - (x_i - x_0) \right\|^2$$

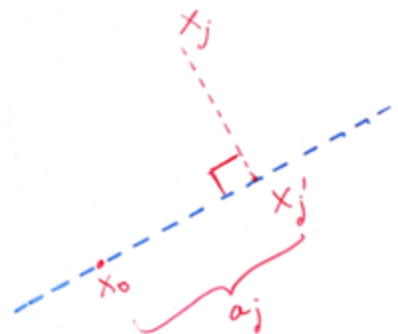$$= \sum a_i^2 \|e\|^2 - 2 \sum a_i e^T (x_i - x_0) + \sum \|x_i - x_0\|^2$$

$$= \sum_i a_i^2 - 2 \sum_i a_i e^T (x_i - x_0) + \sum_i \|x_i - x_0\|^2$$

Now, To find $a_j$

$$\frac{\partial J}{\partial a_j} = 2 a_j - 2 e^T (x_j - x_0) = 0$$

$$a_j = e^T (x_j - x_0)$$

i.e., Perpendicularly Project onto the line

<u>To Find $\vec{e}$</u>

$$J = \sum a_i^2 - 2\sum a_i \, e^T(x_i - x_0) + \sum \|x_i - x_0\|^2$$

$$\because \quad a_i = e^T(x_i - x_0) \, ,$$

$$J = \sum a_i^2 - 2\sum a_i^2 + \sum \|x_i - x_0\|^2$$

$$= -\sum a_i^2 \qquad + \sum \|x_i - x_0\|^2$$

$$= -\sum e^T(x_i - x_0)(x_i - x_0)^T e \; + \sum \|x_i - x_0\|^2$$

$$= -\left[ e^T \left( \sum (x_i - x_0)(x_i - x_0)^T \right) e \right] + \sum \|x_i - x_0\|^2$$

$$= -e^T S \, e \; + \sum \|x_i - x_0\|^2$$

## To Find $\vec{e}$

$$J = \sum a_i^2 - 2 \sum a_i \, e^T (x_i - x_0) + \sum \| x_i - x_0 \|^2$$

$$\because \quad a_i = e^T (x_i - x_0),$$

$$J = \sum a_i^2 - 2 \sum a_i^2 + \sum \| x_i - x_0 \|^2$$

$$= - \sum a_i^2 + \sum \| x_i - x_0 \|^2$$

$$= - \sum e^T (x_i - x_0)(x_i - x_0)^T e + \sum \| x_i - x_0 \|^2$$

$$= - \left[ e^T \left( \sum (x_i - x_0)(x_i - x_0)^T \right) e \right] + \sum \| x_i - x_0 \|^2$$

$$= - e^T S e + \sum \| x_i - x_0 \|^2$$

where $S =$ Scatter matrix $= \sum_i (x_i - x_0)(x_i - x_0)^T$

$S = n \cdot$ (sample Covariance Matrix)

$$J = -e^T S e + \sum \| x_i - x_0 \|^2$$

To minimize $J$, $-e^T S e$ should be minimized subject to the constraint

$$\| e \| = 1 \qquad \text{or} \qquad e^T e - 1 = 0$$

This is "Constrained Optimization" problem.

We can use the method of "Lagrange multipliers"

# Constrained Optimization

Minimize $f(v)$

subject to $g_j(v) \leq 0$, for $1 \leq j \leq n$.

Lagrangian, $\mathcal{L} = f(v) + \sum_{j=1}^{n} \alpha_j \, g_j(v)$

$\downarrow$

Lagrange Multiplier

Necessary cond. at optimal $v$ are:

(i) $\quad \nabla_v \mathcal{L} = 0$

(ii) $\quad \alpha_j \geq 0$ $\quad\Bigg\}$ for all $j = 1$ to $n$

(iii) $\alpha_j \, g_j(v) = 0$

# But, we are with equality constraint.

- So, gradient w.r.t. primal variables and gradient w.r.t. dual variables can be equated to zero.

  – Ofcourse, the Lagrange multipliers should be nonnegative.

Minimize $\quad -e^T s e$

such that $\quad e^T e - 1 = 0$

$$\mathcal{L} = (-e^T s e) + \alpha(e^T e - 1)$$

$$\nabla_e \mathcal{L} = \frac{\partial \mathcal{L}}{\partial e} = -2se + 2\alpha e = 0$$

Minimize $\quad -e^T s e$

such that $\quad e^T e - 1 = 0$

$$\mathcal{L} = (-e^T s e) + \alpha(e^T e - 1)$$

$$\nabla_e \mathcal{L} = \frac{\partial \mathcal{L}}{\partial e} = -2se + 2\alpha e = 0$$

$$se = \alpha e$$

Scatter matrix $\quad$ scalar $\quad$ vector

Minimize $\quad -e^T s e$

such that $\quad e^T e - 1 = 0$

$$\mathcal{L} = (-e^T s e) + \alpha(e^T e - 1)$$

$$\nabla_e \mathcal{L} = \frac{\partial \mathcal{L}}{\partial e} = -2Se + 2\alpha e = 0$$

$$\underset{\text{Scatter matrix}}{S} e = \underset{\text{scalar}}{\alpha}\ \underset{\text{vector}}{e}$$

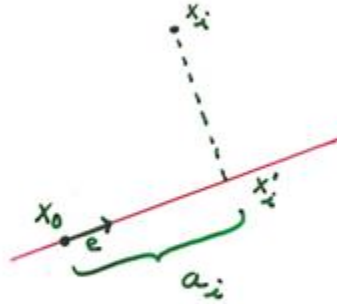This is eigen value (vector) problem

$\quad \alpha$ is eigen value $\Big\}$ for $S$
$\quad e$ is eigen vector

$\therefore \quad -e^T s e$ minimized $\Rightarrow -e^T \alpha e$ minimized

$\Rightarrow \alpha$ should be maximum.

i.e., $\quad e$ is the eigen vector for $S$ for which
$\qquad\qquad$ eigen value is maximum.
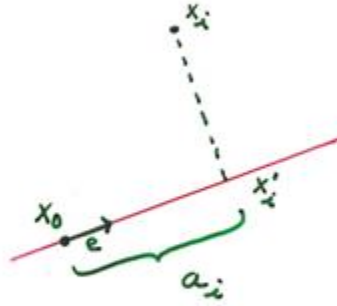
## e gives maximum Variance Direction

$$a_i = e^T (x_i - x_0)$$



$\{x_1, x_2, \ldots, x_n\}$ is represented

as $\{a_1, a_2, \ldots, a_n\}$

## e gives maximum Variance Direction

$$a_i = e^T (x_i - x_0)$$



$\{x_1, x_2, \ldots, x_n\}$ is represented

as $\{a_1, a_2, \ldots, a_n\}$

Variance of $\{a_1, \ldots, a_n\}$

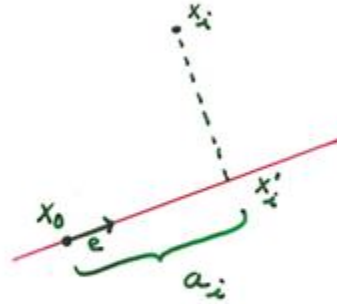$$= \frac{1}{n} \sum (a_i - a_0)^2 \qquad = \frac{1}{n} \sum a_i^2$$

$$[\because a_0 = 0]$$

$$= \frac{1}{n} e^T \left[ \sum (x_i - x_0)(x_i - x_0)^T \right] e$$

$$= \frac{1}{n} e^T S e$$

## e gives maximum Variance Direction

$$a_i = e^T (x_i - x_0)$$



$\{x_1, x_2, ...., x_n\}$ is represented

as $\{a_1, a_2, ..., a_n\}$

Variance of $\{a_1, ..., a_n\}$

$$= \frac{1}{n} \sum (a_i - a_0)^2 \quad = \quad \frac{1}{n} \sum a_i^2$$

$$[\because a_0 = 0]$$

$$= \frac{1}{n} e^T \left[ \sum (x_i - x_0)(x_i - x_0)^T \right] e$$

$$= \frac{1}{n} e^T S e$$

\* The new space is so found that $e^T S e$ is maximum possible one.

i.e., The data is Projected onto that line over which the variance is large.

## Generalization : To find $d'$ dimensions

$$x_i' = x_0 + \left[ a_{i_1} e_1 + \cdots + a_{id'} e_{d'} \right]$$

We get, $\quad J = \sum \left\| x_0 + \sum_{j=1}^{d'} a_{ij} e_j - x_i \right\|^2$

**Generalization :** To find $d'$ dimensions

$$x_i' = x_0 + \left[ a_{i1} e_1 + \cdots + a_{id'} e_{d'} \right]$$

We get, $$J = \sum \left\| x_0 + \sum_{j=1}^{d'} a_{ij} e_j - x_i \right\|^2$$

It is easy to show that

$$\left. \begin{array}{c} e_1 \\ e_2 \\ \vdots \\ e_{d'} \end{array} \right\} \text{ are eigen vectors of } S \text{ for which eigen values are maximum possible.}$$

$$S e_1 = \alpha_1 e_1$$
$$S e_2 = \alpha_2 e_2 \qquad \alpha_1 \geqslant \alpha_2 \geqslant \cdots \geqslant \alpha_{d'}$$
$$\vdots$$
$$S e_{d'} = \alpha_{d'} e_{d'}$$

Generalization : To find $d'$ dimensions

$$x'_i = x_0 + \left[ a_{i1} e_1 + \cdots + a_{id'} e_{d'} \right]$$

We get, $$J = \sum \left\| x_0 + \sum_{j=1}^{d'} a_{ij} e_j - x_i \right\|^2$$

It is easy to show that

$$\left. \begin{array}{c} e_1 \\ e_2 \\ \vdots \\ e_{d'} \end{array} \right\}$$ are eigen vectors of $S$ for which eigen values are maximum possible.

$$S e_1 = \alpha_1 e_1$$
$$S e_2 = \alpha_2 e_2$$
$$\vdots$$
$$S e_{d'} = \alpha_{d'} e_{d'}$$

$$\alpha_1 \geq \alpha_2 \geq \cdots \geq \alpha_{d'}$$

Further, $e_1, e_2, \ldots, e_{d'}$ are ortho normal

Generalization : To find $d'$ dimensions

$$x_i' = x_0 + \left[ a_{i1} e_1 + \cdots + a_{id'} e_{d'} \right]$$

We get, $$J = \sum \left\| x_0 + \sum_{j=1}^{d'} a_{ij} e_j - x_i \right\|^2$$

It is easy to show that

$$\left. \begin{array}{c} e_1 \\ e_2 \\ \vdots \\ e_{d'} \end{array} \right\} \text{ are eigen vectors of } S \text{ for which eigen values are maximum possible.}$$

$$S e_1 = \alpha_1 e_1$$
$$S e_2 = \alpha_2 e_2$$
$$\vdots$$
$$S e_{d'} = \alpha_{d'} e_{d'}$$

$$\alpha_1 \geq \alpha_2 \geq \cdots \geq \alpha_{d'}$$

Because the scatter matrix S is real and symmetric, Eigen values are real and nonnegative.

Further, $e_1, e_2, \ldots, e_{d'}$ are ortho normal

# Representation in New Space

$$X_i' = X_0 + [a_{i_1} e_1 + \cdots + a_{id'} e_{d'}]$$

$\because X_0, \ e_1, .., e_{d'}$ are fixed, So

$X_i'$ can be represented as $\begin{bmatrix} a_{i_1} \\ a_{i_2} \\ \vdots \\ a_{id'} \end{bmatrix} = Y_i$

# Representation in New Space

$$X_i' = X_0 + [a_{i_1} e_1 + \cdots + a_{id'} e_{d'}]$$

$\because X_0, \ e_1, \ldots, e_{d'}$ are fixed, So

$X_i'$ can be represented as $\begin{bmatrix} a_{i_1} \\ a_{i_2} \\ \vdots \\ a_{id'} \end{bmatrix} = Y_i$

Since $\quad a_{i_1} = e_1^T(X_i - X_0)$

$$a_{i_2} = e_2^T(X_i - X_0)$$
$$\vdots$$
$$a_{id'} = e_{d'}^T(X_i - X_0)$$

$$Y_i = \begin{bmatrix} e_1^T \\ e_2^T \\ \vdots \\ e_{d'}^T \end{bmatrix}(X_i - X_0)$$

Transformation matrix

Let us call this, $P = \begin{bmatrix} - e_1 - \\ - e_2 - \\ \vdots \\ - e_{d'} - \end{bmatrix}_{d' \times d}$

# The projection matrix P

- $P^t P = I$, But $P$ need not be square.

  If $P$ is square, i.e., use all eigen vectors, then $P$ is orthogonal.

  Not only $P$ is orthogonal, $P$ is a rotation matrix.

# The projection matrix P

- $P^t P = I$, But $P$ need not be square.

  If $P$ is square, i.e., use all eigen vectors, then $P$ is orthogonal.

  Not only $P$ is orthogonal, $P$ is a rotation matrix.

- PCA basically, translates (so that origin becomes centroid) and does rotation (so that features are uncorrelated).

Data

| $x$ | $y$ |
| --- | --- |
| 2.5 | 2.4 |
| 0.5 | 0.7 |
| 2.2 | 2.9 |
| 1.9 | 2.2 |
| 3.1 | 3.0 |
| 2.3 | 2.7 |
| 2 | 1.6 |
| 1 | 1.1 |
| 1.5 | 1.6 |
| 1.1 | 0.9 |

Ref: http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf

## Data

| x | y |
|---|---|
| 2.5 | 2.4 |
| 0.5 | 0.7 |
| 2.2 | 2.9 |
| 1.9 | 2.2 |
| 3.1 | 3.0 |
| 2.3 | 2.7 |
| 2 | 1.6 |
| 1 | 1.1 |
| 1.5 | 1.6 |
| 1.1 | 0.9 |

## Mean subtracted Data

| x | y |
|---|---|
| .69 | .49 |
| -1.31 | -1.21 |
| .39 | .99 |
| .09 | .29 |
| 1.29 | 1.09 |
| .49 | .79 |
| .19 | -.31 |
| -.81 | -.81 |
| -.31 | -.31 |
| -.71 | -1.01 |

Ref: http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf

Data

| $x$ | $y$ |
|-----|-----|
| 2.5 | 2.4 |
| 0.5 | 0.7 |
| 2.2 | 2.9 |
| 1.9 | 2.2 |
| 3.1 | 3.0 |
| 2.3 | 2.7 |
| 2 | 1.6 |
| 1 | 1.1 |
| 1.5 | 1.6 |
| 1.1 | 0.9 |

Mean subtracted Data

| $x$ | $y$ |
|------|-------|
| .69 | .49 |
| -1.31 | -1.21 |
| .39 | .99 |
| .09 | .29 |
| 1.29 | 1.09 |
| .49 | .79 |
| .19 | -.31 |
| -.81 | -.81 |
| -.31 | -.31 |
| -.71 | -1.01 |

$$cov = \begin{pmatrix} .61 & .61 \\ .61 & .71 \end{pmatrix}$$

Ref: http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf

| Data | |
|------|-----|
| $x$ | $y$ |
| 2.5 | 2.4 |
| 0.5 | 0.7 |
| 2.2 | 2.9 |
| 1.9 | 2.2 |
| 3.1 | 3.0 |
| 2.3 | 2.7 |
| 2 | 1.6 |
| 1 | 1.1 |
| 1.5 | 1.6 |
| 1.1 | 0.9 |

Mean subtracted Data

| $x$ | $y$ |
|------|-----|
| .69 | .49 |
| -1.31 | -1.21 |
| .39 | .99 |
| .09 | .29 |
| 1.29 | 1.09 |
| .49 | .79 |
| .19 | -.31 |
| -.81 | -.81 |
| -.31 | -.31 |
| -.71 | -1.01 |

$$cov = \begin{pmatrix} .61 & .61 \\ .61 & .71 \end{pmatrix}$$

Transformed Data (Single eigenvector)

| $x$ |
|-----|
| -.827970186 |
| 1.77758033 |
| -.992197494 |
| -.274210416 |
| -1.67580142 |
| -.912949103 |
| .0991094375 |
| 1.14457216 |
| .438046137 |
| 1.22382056 |

The data after transforming using only
the most significant eigenvector

**Data**

| x | y |
|---|---|
| 2.5 | 2.4 |
| 0.5 | 0.7 |
| 2.2 | 2.9 |
| 1.9 | 2.2 |
| 3.1 | 3.0 |
| 2.3 | 2.7 |
| 2 | 1.6 |
| 1 | 1.1 |
| 1.5 | 1.6 |
| 1.1 | 0.9 |

**Mean subtracted Data**

| x | y |
|---|---|
| .69 | .49 |
| -1.31 | -1.21 |
| .39 | .99 |
| .09 | .29 |
| 1.29 | 1.09 |
| .49 | .79 |
| .19 | -.31 |
| -.81 | -.81 |
| -.31 | -.31 |
| -.71 | -1.01 |

$$cov = \begin{pmatrix} .61 & .61 \\ .61 & .71 \end{pmatrix}$$

Eigenvalues are 1.28, 0.049

Corresponding Eigenvectors are $\begin{pmatrix} -.67 \\ -.73 \end{pmatrix}$ and $\begin{pmatrix} -.73 \\ .67 \end{pmatrix}$
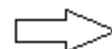
**Transformed Data (Single eigenvector)**

| x |
|---|
| -.827970186 |
| 1.77758033 |
| -.992197494 |
| -.274210416 |
| -1.67580142 |
| -.912949103 |
| .0991094375 |
| 1.14457216 |
| .438046137 |
| 1.22382056 |

The data after transforming using only the most significant eigenvector

Ref: http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf

## Data

| $x$ | $y$ |
|-----|-----|
| 2.5 | 2.4 |
| 0.5 | 0.7 |
| 2.2 | 2.9 |
| 1.9 | 2.2 |
| 3.1 | 3.0 |
| 2.3 | 2.7 |
| 2 | 1.6 |
| 1 | 1.1 |
| 1.5 | 1.6 |
| 1.1 | 0.9 |

## Mean subtracted Data

| $x$ | $y$ |
|------|------|
| .69 | .49 |
| -1.31 | -1.21 |
| .39 | .99 |
| .09 | .29 |
| 1.29 | 1.09 |
| .49 | .79 |
| .19 | -.31 |
| -.81 | -.81 |
| -.31 | -.31 |
| -.71 | -1.01 |

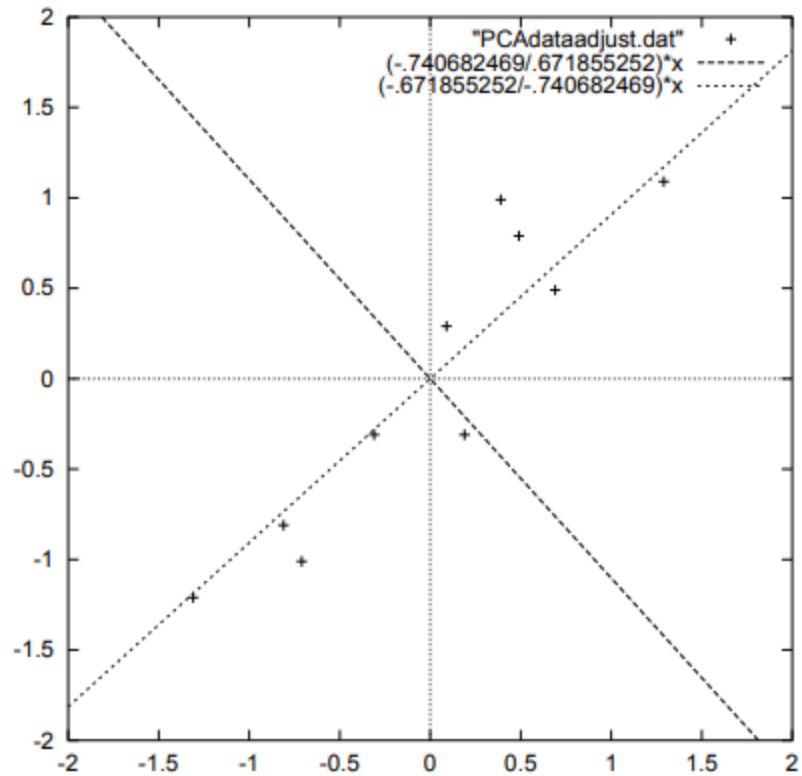$$cov = \begin{pmatrix} .61 & .61 \\ .61 & .71 \end{pmatrix}$$

Eigenvalues are 1.28, 0.049

Corresponding Eigenvectors are $\begin{pmatrix} -.67 \\ -.73 \end{pmatrix}$ and $\begin{pmatrix} -.73 \\ .67 \end{pmatrix}$
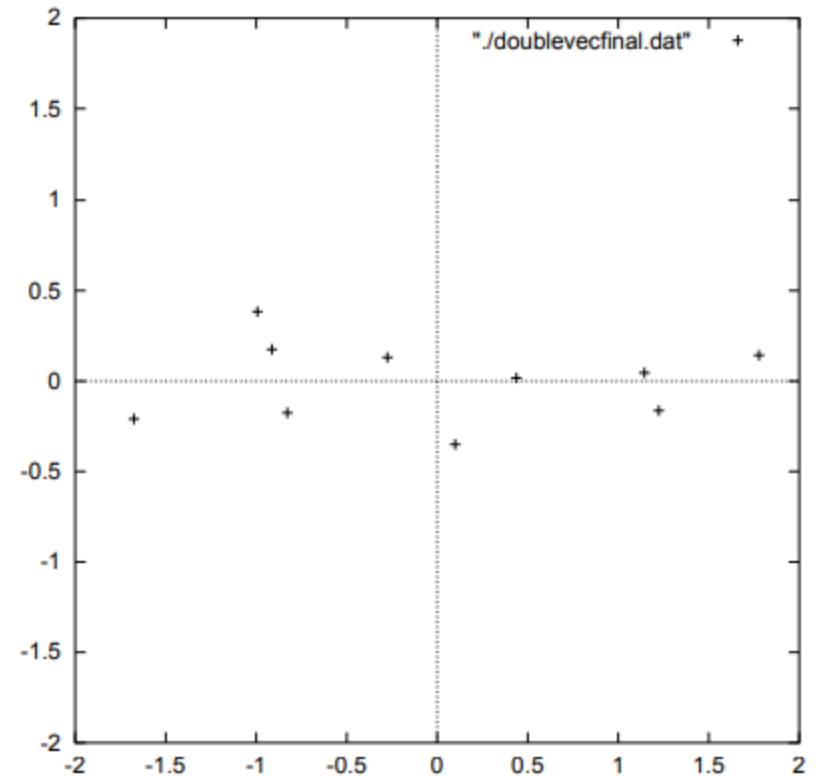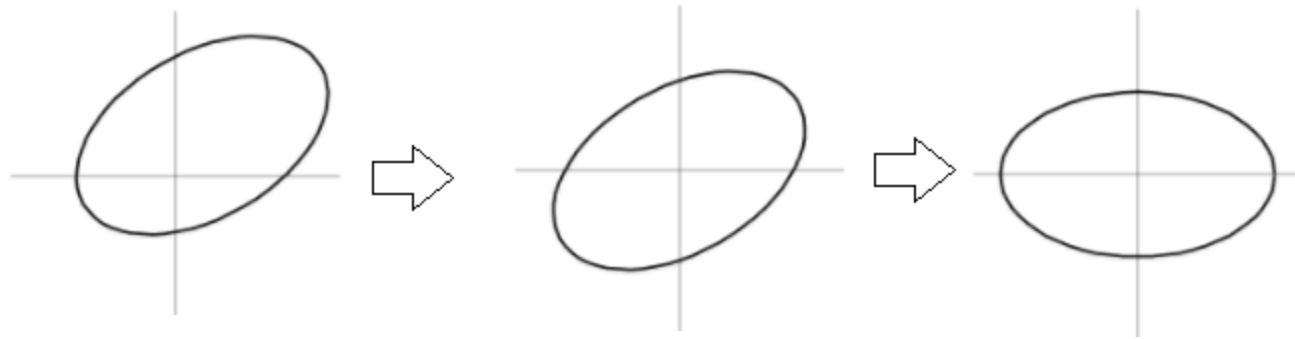
What is the projection matrix?

## Transformed Data (Single eigenvector)

| $x$ |
|-----|
| -.827970186 |
| 1.77758033 |
| -.992197494 |
| -.274210416 |
| -1.67580142 |
| -.912949103 |
| .0991094375 |
| 1.14457216 |
| .438046137 |
| 1.22382056 |

The data after transforming using only the most significant eigenvector

Ref: http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf

**Data**

| x | y |
|---|---|
| 2.5 | 2.4 |
| 0.5 | 0.7 |
| 2.2 | 2.9 |
| 1.9 | 2.2 |
| 3.1 | 3.0 |
| 2.3 | 2.7 |
| 2 | 1.6 |
| 1 | 1.1 |
| 1.5 | 1.6 |
| 1.1 | 0.9 |

**Mean subtracted Data**

| x | y |
|---|---|
| .69 | .49 |
| -1.31 | -1.21 |
| .39 | .99 |
| .09 | .29 |
| 1.29 | 1.09 |
| .49 | .79 |
| .19 | -.31 |
| -.81 | -.81 |
| -.31 | -.31 |
| -.71 | -1.01 |

$$cov = \begin{pmatrix} .61 & .61 \\ .61 & .71 \end{pmatrix}$$

Eigenvalues are 1.28, 0.049

Corresponding Eigenvectors are $\begin{pmatrix} -.67 \\ -.73 \end{pmatrix}$ and $\begin{pmatrix} -.73 \\ .67 \end{pmatrix}$

What is the projection matrix?

$$P = \begin{pmatrix} -.67 \\ -.73 \end{pmatrix}^{t}$$

**Transformed Data (Single eigenvector)**

| x |
|---|
| -.827970186 |
| 1.77758033 |
| -.992197494 |
| -.274210416 |
| -1.67580142 |
| -.912949103 |
| .0991094375 |
| 1.14457216 |
| .438046137 |
| 1.22382056 |

The data after transforming using only the most significant eigenvector

Ref: http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf

Mean adjusted data with eigenvectors overlayed

Data transformed with 2 eigenvectors

Ref: http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf

# HOW IS THAT FEATURES ARE UNCORRELATED IN THE NEW SPACE?

## In the new space features are un-correlated !

$$Y_i = P(X_i - X_0) = PX_i - PX_0$$

Mean in the new space, $Y_0 = P(X_0 - X_0) = \vec{0}$

Scatter matrix in new space $= S'$

$$S' = \sum_i (Y_i - Y_0)(Y_i - Y_0)^T = \sum_i Y_i Y_i^T$$

$$= \sum_i (PX_i - PX_0)(PX_i - PX_0)^T$$

$$= \sum_i P(X_i - X_0)(X_i - X_0)^T P^T$$

$$= P S P^T \qquad \text{S is the scatter matrix of original space.}$$

$$= \begin{bmatrix} -e_1- \\ -e_2- \\ \vdots \\ -e_{d'}- \end{bmatrix} S \cdot \begin{bmatrix} | & | & & | \\ e_1 & e_2 & \cdots & e_{d'} \\ | & | & & | \end{bmatrix}$$

$$= \begin{bmatrix} -e_1- \\ -e_2- \\ \vdots \\ -e_{d'}- \end{bmatrix} \begin{bmatrix} \alpha_1 e_1 & \cdots & \alpha_{d'} e_{d'} \end{bmatrix}$$

$$= \begin{bmatrix} \alpha_1 & & & 0 \\ & \alpha_2 & & \\ & & \ddots & \\ 0 & & & \alpha_{d'} \end{bmatrix}$$

# A Transformation; Covariance Matrix = I

The following transformation matrix can give $\Sigma = I$ in the new space

$$\begin{pmatrix} (\alpha_1)^{-1/2} & & & & \\ & (\alpha_2)^{-1/2} & & & \\ & & \cdot & & \\ & & & \cdot & \\ & & & & (\alpha_d)^{-1/2} \end{pmatrix} P$$

# Drawback of PCA

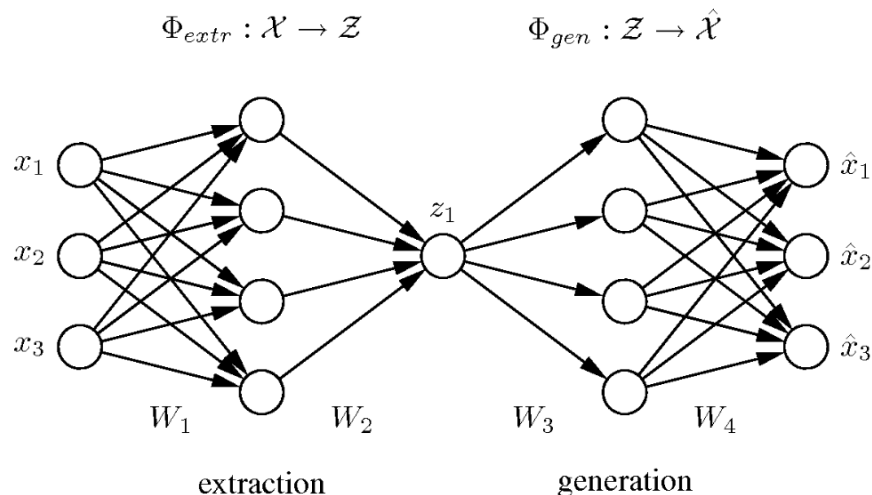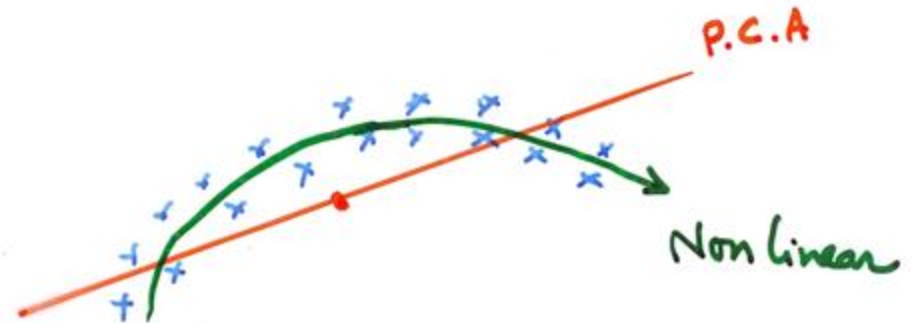* PCA seeks directions that are efficient for representation

  But doesn't take class-labels into account

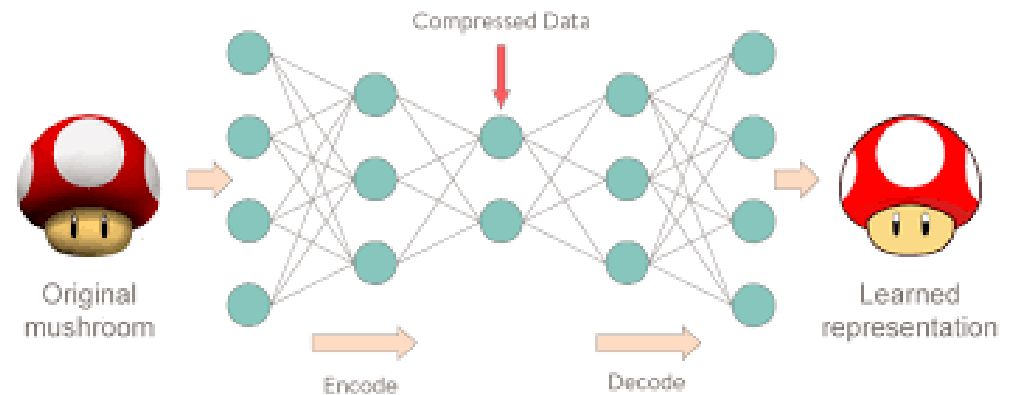  So, the new space may not be good for classification problem



* 2nd P.C. is better than 1st P.C.

Non-linear Transformation can give better results



$\Phi_{extr} : \mathcal{X} \to \mathcal{Z}$

$\Phi_{gen} : \mathcal{Z} \to \hat{\mathcal{X}}$

$x_1$

$x_2$

$x_3$

$z_1$

$\hat{x}_1$

$\hat{x}_2$

$\hat{x}_3$

$W_1$    $W_2$    $W_3$    $W_4$

extraction

generation

**Autoencoder networks and Kernel PCA can find these kind of non-linear mappings..**

Compressed Data

Original mushroom

Encode

Decode

Learned representation

**Autoencoder networks and Kernel PCA can find these kind of non-linear mappings..**

Have you realized that PCA is unsupervised?

# FISHER LINEAR DISCRIMINANT (FLD)

# PCA is unsupervised

- With a labeled (training) data set, for the classification problem, principal components based dimensionality reduction may be bad.

# PCA is unsupervised

- With a labeled (training) data set, for the classification problem, principal components based dimensionality reduction may be bad.



* 2nd P.C. is better than 1st P.C.

# Fisher Linear Discriminant

- The objective is to find linear projections of the patterns which is good for classification.

- Class-labels are taken into account.

# Fisher Linear Discriminant

$\Omega = \{\omega_1, \omega_2\}$

$D_1$ = Set of patterns belonging to $\omega_1$

$D_2$ = " " $\omega_2$

Let $\vec{W}$ denotes direction of a line

$$Y_1 = \{ y_i = W^T x_i \mid x_i \in D_1 \}$$

$$Y_2 = \{ y_j = W^T x_j \mid x_j \in D_2 \}$$

Objective : find $W$ s.t. $Y_1$ & $Y_2$ are well seperated

**FIGURE 3.5.** Projection of the same set of samples onto two different lines in the directions marked **w**. The figure on the right shows greater separation between the red and black projected points. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

**bad line to project to, classes are mixed up**

**good line to project to, classes are well separated**

$$J(w) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2} \qquad \text{maximize}$$

where $\quad \tilde{m}_1 = $ Mean of samples in $Y_1$

$\tilde{m}_2 = \qquad$ " $\qquad$ " $Y_2$

$\tilde{s}_1^2 = $ Within class scatter for $Y_1$

$\qquad = \sum_{y \in Y_1} (y - \tilde{m}_1)^2$

$\tilde{s}_2^2 = $ Within class scatter for $Y_2$

$$(\tilde{m}_1 - \tilde{m}_2)^2 = \left[ W^T (m_1 - m_2) \right]^2$$

mean of patterns in $D_2$

" " " $D_1$

$$= W^T (m_1 - m_2)(m_1 - m_2)^T W$$

$$= W^T S_B W$$

Between class Scatter

$$\left(\tilde{m}_1 - \tilde{m}_2\right)^2 = \left[W^T(m_1 - m_2)\right]^2$$

mean of patterns in $D_2$

" " " $D_1$

$$= W^T(m_1 - m_2)(m_1 - m_2)^T W$$

$$= W^T S_B W$$

Between class Scatter

$$\tilde{s}_i^2 = \sum_{y \in Y_i} (y - \tilde{m}_i)^2$$

$$= \sum_{x \in D_i} \left(W^T x - W^T m_i\right)^2$$

$$= \sum W^T(x - m_i)(x - m_i)^T W$$

$$= W^T S_i W$$

Within class scatter for $D_i$

$$\left(\tilde{m}_1 - \tilde{m}_2\right)^2 = \left[W^T(m_1 - m_2)\right]^2$$

mean of patterns in $D_2$

" " " $D_1$

$$= W^T(m_1 - m_2)(m_1 - m_2)^T W$$

$$= W^T S_B W$$

→ Between Class Scatter

$$\tilde{s}_i^2 = \sum_{y \in Y_i} (y - \tilde{m}_i)^2$$

$$= \sum_{x \in D_i} \left(W^T x - W^T m_i\right)^2$$

$$= \sum W^T(x - m_i)(x - m_i)^T W$$

$$= W^T S_i W$$

→ Within class scatter for $D_i$

$$\tilde{s}_1^2 + \tilde{s}_2^2 = W^T S_1 W + W^T S_2 W$$

$$= W^T(S_1 + S_2)W$$

$$= W^T S_W W$$

→ Total "within class scatter".

$$J(W) = \frac{W^T S_B W}{W^T S_W W}$$

W that maximizes J must satisfy

$$S_B W = \lambda S_W W \longrightarrow \textcircled{1}$$

$$S_W^{-1} S_B W = \lambda W$$

W is eigen vector for $S_W^{-1} S_B$

$\lambda$ is eigen value

$$J(w) = \frac{w^T S_B w}{w^T S_w w}$$

W that maximizes J must satisfy

$$S_B W = \lambda S_w W \longrightarrow \textcircled{1}$$

$$S_w^{-1} S_B W = \lambda W$$

W is eigen vector for $S_w^{-1} S_B$
$\lambda$ is eigen value

Reason is given in the supplementary slides towards the end.

$$J(W) = \frac{W^T S_B W}{W^T S_W W}$$

$W$ that maximizes $J$ must satisfy

$$S_B W = \lambda S_W W \qquad \longrightarrow \text{①}$$

$$S_W^{-1} S_B W = \lambda W$$

$W$ is eigen vector for $S_W^{-1} S_B$
$\lambda$ is eigen value

But there is no need to solve the eigen value problem.

$$S_B W = (m_1 - m_2)\underbrace{(m_1 - m_2)^T W}_{\text{scalar}}$$

$$= K(m_1 - m_2)$$

From ①,

$$K(m_1 - m_2) = \lambda S_W W$$

$$W = \frac{K}{\lambda} S_W^{-1}(m_1 - m_2)$$

Since only direction of $W$ is important

$$W = S_W^{-1}(m_1 - m_2)$$

# Extension to get more than 1D data

- Multiple discriminant analysis. Popularly known as Linear Discriminant Analysis (LDA).

- Similar to Kernel PCA, kernel Fisher discriminant is there.

Some supplementary material

# SUPPLEMENTARY

# Fisher Linear Discriminant Derivation

- Thus our objective function can be written:

$$J(v) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{s}_1^2 + \tilde{s}_2^2} = \frac{v^t S_B v}{v^t S_W v}$$

- Minimize $J(v)$ by taking the derivative w.r.t. $v$ and setting it to 0

$$\frac{d}{dv} J(v) = \frac{\left(\frac{d}{dv} v^t S_B v\right) v^t S_W v - \left(\frac{d}{dv} v^t S_W v\right) v^t S_B v}{\left(v^t S_W v\right)^2}$$

$$= \frac{(2 S_B v) v^t S_W v - (2 S_W v) v^t S_B v}{\left(v^t S_W v\right)^2} = 0$$

Instead of *W, V* is used.
Ref: http://www.csd.uwo.ca/~olga/Courses/CS434a_541a/Lecture8.pdf

## Fisher Linear Discriminant Derivation

- Need to solve $v^t S_W v (S_B v) - v^t S_B v (S_W v) = 0$

$$\Rightarrow \frac{v^t S_W v (S_B v)}{v^t S_W v} - \frac{v^t S_B v (S_W v)}{v^t S_W v} = 0$$

$$\Rightarrow S_B v - \frac{v^t S_B v (S_W v)}{v^t S_W v} = 0$$
$$= \lambda$$

$$\Rightarrow S_B v = \lambda S_W v$$

generalized eigenvalue problem

Instead of *W, V* is used.
Ref: http://www.csd.uwo.ca/~olga/Courses/CS434a_541a/Lecture8.pdf

- [http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf](http://www.cs.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf)