

(1)

Spring 2014

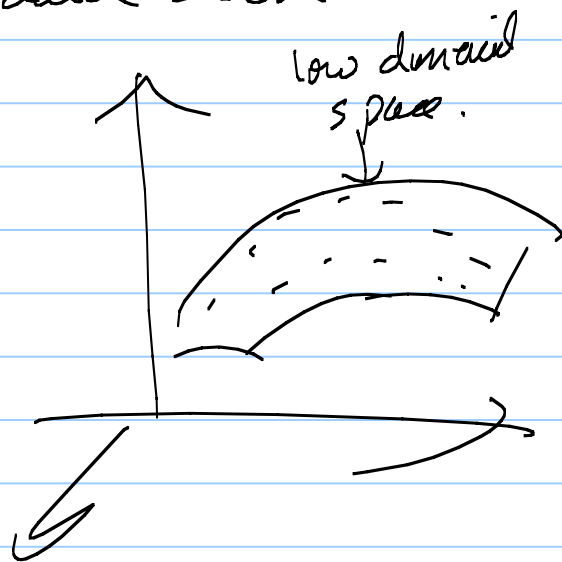
Principal Component Analysis (PCA)

Note Title

10/15/2006

One way to deal with the curse of dimensionality is to project data down onto a space of low dimensions.

There are a number of different techniques for doing this \rightarrow eg. multidimensional scaling. Too many to deal with in this course.



Now we discuss the most basic method
- Principal Component Analysis. (PCA)

(2)

CONVENTION : $\underline{\mu}^T \underline{\mu}$ is a scalar

$\underline{\mu} \underline{\mu}^T$ is a matrix

$$\mu_1^2 + \mu_2^2 + \dots + \mu_D^2$$

$$\begin{pmatrix} \mu_1^2 & \mu_1 \mu_2 & \mu_1 \mu_3 & \dots \\ & \mu_2^2 & & \\ & & & \\ & & & \dots \end{pmatrix}$$

(N.B. different convention than blackboard notes - but same as in book)

Data samples $\underline{x}_1, \dots, \underline{x}_N$

Compute the mean $\underline{\mu} = \frac{1}{N} \sum_{i=1}^N \underline{x}_i$ in D -dim space.

Compute the covariance:

$$\underline{K} = \frac{1}{N} \sum_{i=1}^N (\underline{x}_i - \underline{\mu})(\underline{x}_i - \underline{\mu})^T$$

Next compute the eigenvalues and eigenvectors of \underline{K}

$$\text{Solve } \underline{K} \underline{e} = \lambda \underline{e}$$

$$\lambda_1 > \lambda_2 > \dots > \lambda_N$$

Note: \underline{K} is symmetric - so eigenvalues are real, eigenvectors are orthogonal.

PCA reduces the dimension by

by keeping the eigenvectors \underline{e}_i with $\lambda_i > T$

Let M eigenvectors be kept.

λ threshold

Then project data \underline{x} onto the subspace spanned by the first M eigenvectors. (After subtracting out the mean).

(3)

Formally :

$$\underline{\text{Project.}} \quad \underline{x} - \underline{\mu} = \sum_{v=1}^D a_v \underline{e}_v$$

where the coefficients are given by

$$a_v = (\underline{x} - \underline{\mu}) \cdot \underline{e}_v \quad \left(\begin{array}{l} \text{orthogonality, means} \\ \underline{e}_v \cdot \underline{e}_\mu = \delta_{v\mu} \\ \text{Kronecker delta} \end{array} \right)$$

$$\underline{\text{Hence}} \quad \underline{x} = \underline{\mu} + \sum_{v=1}^D \left\{ (\underline{x} - \underline{\mu}) \cdot \underline{e}_v \right\} \underline{e}_v$$

no dimension reduction
(no compression)

$$\underline{\text{Then, approximate}} \quad \underline{x} \approx \underline{\mu} + \sum_{v=1}^M \left\{ (\underline{x} - \underline{\mu}) \cdot \underline{e}_v \right\} \underline{e}_v$$

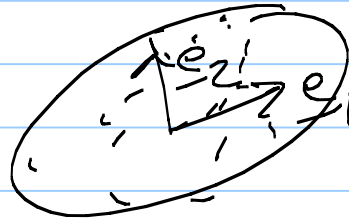
Projects the data into the M -dim subspace.

$$\underline{\mu} + \sum_{v=1}^M b_v \underline{e}_v \quad \cdot \parallel$$

(4)

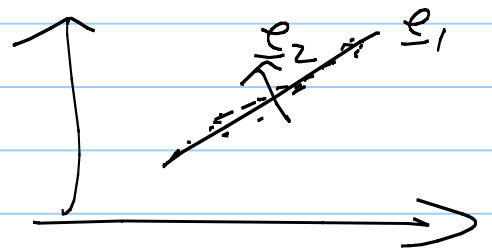
In 2-dimensions

Visually



The eigenvectors of \underline{K} correspond to the second order moments of the data.

If the data lies (almost) on a straight line, then $\lambda_1 \gg 0, \lambda_2 = 0$



PCA and Gaussian Distribution

PCA is equivalent to performing ML estimation of the parameters of a Gaussian

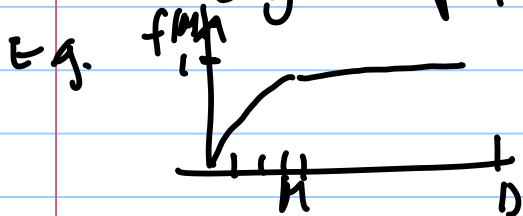
$$P(\underline{x} | \underline{\mu}, \underline{\Sigma}) = \frac{1}{\sqrt{2\pi} |\det \underline{\Sigma}|} e^{-\frac{1}{2} (\underline{x} - \underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x} - \underline{\mu})}$$

to get $\hat{\underline{\mu}}, \hat{\underline{\Sigma}}$. And then throw away the directions where the variance is small.

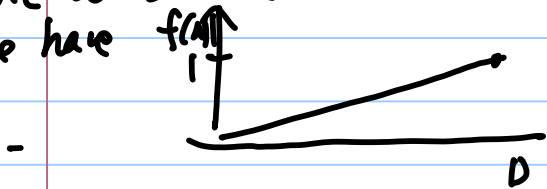
(5) When is PCA good?

It is almost always a good technique to try, because it is so simple. Obtain the eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$ and plot $f(M) = \frac{\sum_{i=1}^M \lambda_i}{\sum_{i=1}^n \lambda_i}$. $f(M)$ increases with M and takes max value 1, at $M=D$.

PCA is good if $f(M)$ asymptotes rapidly to 1. This happens if the first few eigenvalues are big and the remainder are small.



PCA is bad if all the eigenvalues are roughly equal. The worst case is when all eigenvalues are the same. Then we have



This can happen. For example if the data is a set of strings.

$$(1, 0, 0, 0, \dots) = x_1$$

$$(0, 1, 0, 0, \dots) = x_2$$

$$(0, 0, 0, 0, \dots, 0, 1) = x_n$$

(6) What is PCA doing?

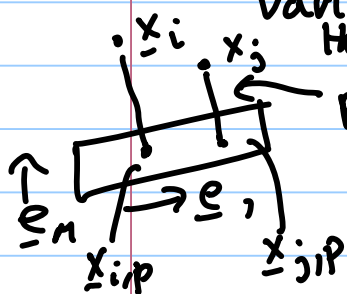
There are two equivalent ways to interpret PCA:

- (1) Minimize the projection error.
- (2) Maximize the variance of the projection.

Consider $\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$. This is the variance of the data. It is independent of the projection.

We can express $(x_i - \mu)^2 = \sum_{v=1}^M \left\{ (x_i - \mu) \cdot \underline{e}_v \right\}^2$ e_v's are the eigenvectors of the correlation

Hence $\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 = \underbrace{\frac{1}{N} \sum_{i=1}^N \sum_{v=1}^M \left\{ (x_i - \mu) \cdot \underline{e}_v \right\}^2}_{\text{variance of data within plane } \underline{e}_1, \dots, \underline{e}_n} + \underbrace{\frac{1}{N} \sum_{i=1}^N \sum_{v=n+1}^D \left\{ (x_i - \mu) \cdot \underline{e}_v \right\}^2}_{\text{projection error}}$



P-project

x_i projects to $x_{i,p} = \mu + \sum_{v=1}^M \left\{ (x_i - \mu) \cdot \underline{e}_v \right\} \underline{e}_v$
 It has projection error $\sum_{v=n+1}^D \left\{ (x_i - \mu) \cdot \underline{e}_v \right\}^2$.

The projected points $x_{i,p}$ have variance.

The sum of the projection error and the variance of projection are constant. So maximizing one is equivalent to minimizing the other.

Also this relationship can be expressed in terms of the eigenvalues.
 It reduces to $\sum_{v=1}^D \lambda_v = \sum_{v=1}^n \lambda_v + \sum_{v=n+1}^D \lambda_v$

To see this, $\frac{1}{N} \sum_{i=1}^N (x_i - \mu) \cdot (x_i - \mu) = \text{Trace}(\underline{C}) = \sum_{v=1}^D \lambda_v$.

The variance of the projection is $\sum_{v=1}^n \lambda_v$, by similar reasoning.

Hence projection error is $\sum_{v=n+1}^D \lambda_v$.

(7) To understand why PCA minimizes, or maximizes, these terms we must express the criterion slightly differently. Then we use Singular Value Decomposition (SVD), which is advanced material

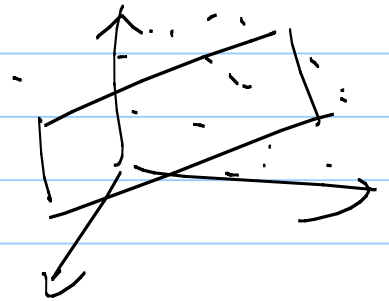
$$J(\underline{\mu}, \{a_i\}, \{e_i\}) = \sum_{k=1}^N \left\| \left(\underline{\mu} + \sum_{i=1}^M a_{ki} \underline{e}_i \right) - \underline{x}_k \right\|^2$$

Minimize J w.r.t. $\underline{\mu}, \{a_i\}, \{e_i\}$.

Data $\{ \underline{x}_k : k=1 \text{ to } N \}$

The $\{a_{ki}\}$ are projection coefficients //

Intuition: find the M -dimensional subspace s.t. the projections of the data onto this subspace have minimal error.



Minimizing J , gives the

$\{ \hat{\underline{e}}_i \}$'s to be the eigenvectors of the covariance matrix $\underline{K} = \frac{1}{N} \sum_{k=1}^N (\underline{x}_k - \underline{\mu})(\underline{x}_k - \underline{\mu})^T$

$$\underline{\mu} = \frac{1}{N} \sum_{k=1}^N \underline{x}_k$$

$\hat{a}_{ki} = (\underline{x}_k - \hat{\underline{\mu}}) \cdot \hat{\underline{e}}_i$ the projection coefficients.

(2)

To understand this fully, you must understand Singular Value Decomposition (SVD)

We can re-express the criteria as

$$J[\underline{\mu}, \{a_i, \{e_i\}\}] = \sum_{k=1}^N \sum_{b=1}^D \left((\mu_b - X_{bk}) + \sum_{i=1}^M a_{ki} e_{ib} \right)^2$$

where b denotes the vector components.

This is an example of a general class of problem.

Let $E[\Psi, e] = \sum_{a=1, k=1}^{a=D, k=N} \left(\tilde{X}_{ak} - \sum_{v=1}^M \Psi_{av} \phi_{vk} \right)^2$

Goal: minimize $E[\Psi, e]$ w.r.t. Ψ, e .

This is a bilinear problem, that can be solved by SVD.

Note: $\tilde{X}_{ak} = X_{ak} - \mu_a$
the position of the point, relative to the mean.

(9)

SVD

Note: \underline{X} is not a square matrix. So it has no eigenvalues or eigenvectors.

We can express any $N \times D$ matrix \underline{X} X_{ak} in form $\underline{X} = \underline{E} \underline{D} \underline{F}$

$$X_{ak} = \sum_{\mu=1}^M e_{a\mu} d_{\mu\nu} f_{\nu k}$$

where $\underline{D} = \{d_{\mu\nu}\}$ is a diagonal matrix ($d_{\mu\nu} = 0, \mu \neq \nu$)
 $\underline{D} = \begin{pmatrix} \sqrt{\lambda_1} & 0 \\ 0 & \sqrt{\lambda_m} \end{pmatrix}$, where the $\{\lambda_i\}$ are eigenvalues of $\underline{X} \underline{X}^T$ (equivalently of $\underline{X}^T \underline{X}$).

μ, ν
label the
eigenvectors.

$\underline{E} = \{e_{a\mu}\}$ are eigenvectors of $(\underline{X} \underline{X}^T)_{ab}$

$\underline{F} = \{f_{\nu k}\}$ are eigenvectors of $(\underline{X}^T \underline{X})_{kl}$

Note: For $\bar{\underline{X}}$ defined on previous page, we get that $(\bar{\underline{X}} \bar{\underline{X}}^T) = \sum_{k=1}^N (\underline{x}_{k-1}) (\underline{x}_{k-1})^T$

Note: If $(\underline{X} \underline{X}^T) \underline{e} = \lambda \underline{e}$
then $(\underline{X}^T \underline{X}) (\underline{X}^T \underline{e}) = \lambda (\underline{X}^T \underline{e})$

This relates the eigenvectors of $\underline{X} \underline{X}^T$ and of $\underline{X}^T \underline{X}$.
(Calculate the eigenvectors for the smallest matrix, then deduce those of the bigger matrix - $D < N$)

(10)

Minimize:

$$E[\psi, e] = \sum_{a=1, k=1}^{a=D, k=N} \left(\tilde{X}_{ak} - \sum_{v=1}^M \psi_{av} \phi_{vk} \right)^2$$

we set:

$$\begin{cases} \psi_{av} = \sqrt{d_{vv}} e_a^v \\ \phi_{vk} = \sqrt{d_{vv}} f_k^v \end{cases}$$

Take M biggest terms in the SVD expansion of \tilde{X} .

But there is an ambiguity.

$$\sum_{v=1}^M \psi_{av} \phi_{vk} = \psi \phi = \underline{\psi} \underline{\phi} = \underline{a}_k \quad \text{matrix multiplied.}$$

$$= \underline{\psi} \underline{A} \underline{A}^{-1} \underline{\phi} = \underline{a}_k$$

for any $M \times M$ invertible matrix \underline{A} .

$$\begin{aligned} \underline{\psi} &\rightarrow \underline{\psi} \underline{A} \\ \underline{\phi} &\rightarrow \underline{A}^{-1} \underline{\phi} \end{aligned}$$

This gets rid of the ambiguity.

For the PCA problem — we have constraints that the projection directions are orthogonal unit eigenvectors

(11)

Relate SVD to PCA (Linear Algebra)

Start with an $n \times m$ matrix \underline{X} .

$\underline{X} \underline{X}^T$ is a symmetric $n \times n$ matrix

$\underline{X}^T \underline{X}$ is a symmetric $m \times m$ matrix.

$$(\underline{X} \underline{X}^T)^T = \underline{X} \underline{X}^T$$

By standard linear algebra.

$$\underline{X} \underline{X}^T \underline{e}^\mu = \lambda^\mu \underline{e}^\mu \quad n \text{ eigenvalues } \lambda^\mu$$

eigenvectors \underline{e}^μ

eigenvectors are orthogonal

$$\underline{e}^\mu \cdot \underline{e}^\nu = \delta^{\mu\nu}$$

$$\text{Similarly } \underline{X}^T \underline{X} \underline{f}^\nu = \tau^\nu \underline{f}^\nu \quad m \text{ eigenvalues } \tau^\nu$$

eigenvectors \underline{f}^ν

$$\underline{f}^\mu \cdot \underline{f}^\nu = \delta^{\mu\nu}$$

The $\{\underline{e}^\mu\}$ and $\{\underline{f}^\mu\}$ are related

because $(\underline{X}^T \underline{X}) (\underline{X}^T \underline{e}^\mu) = \lambda^\mu (\underline{X}^T \underline{e}^\mu)$

$$(\underline{X} \underline{X}^T) (\underline{X} \underline{f}^\mu) = \tau^\mu (\underline{X} \underline{f}^\mu)$$

Hence: $\underline{X}^T \underline{e}^\mu \propto \underline{f}^\mu$, $\underline{X} \underline{f}^\mu \propto \underline{e}^\mu$ $\lambda^\mu = \tau^\mu$

If $n > m$, then there are n eigenvectors $\{\underline{e}^\mu\}$ and m eigenvectors $\{\underline{f}^\mu\}$. So some \underline{f}^μ relate to several $\{\underline{e}^\mu\}$.

(12)

Claim: we can express

$$\underline{X} = \sum_{\mu} \alpha^{\mu} \underline{e}^{\mu} \underline{f}^{\mu T} \quad \text{for some } \alpha^{\mu}$$
$$\underline{X}^T = \sum_{\mu} \alpha^{\mu} \underline{f}^{\mu} \underline{e}^{\mu T} \quad \left(\text{we will solve for } \alpha^{\mu} \text{ later.} \right)$$

Verify the claim

$$\underline{X} \underline{f}^{\nu} = \sum_{\mu} \alpha^{\mu} \underline{e}^{\mu} \underline{f}^{\mu T} \underline{f}^{\nu}$$
$$= \sum_{\mu} \alpha^{\mu} \delta_{\mu\nu} \underline{e}^{\mu} = \alpha^{\nu} \underline{e}^{\nu} \quad \parallel$$

$$\underline{X} \underline{X}^T = \sum_{\mu, \nu} \alpha^{\nu} \underline{e}^{\nu} \underline{f}^{\nu T} \alpha^{\mu} \underline{f}^{\mu} \underline{e}^{\mu T}$$
$$= \sum_{\mu, \nu} \alpha^{\nu} \alpha^{\mu} \underline{e}^{\nu} \delta_{\mu\nu} \underline{e}^{\mu T} = \sum_{\mu} (\alpha^{\mu})^2 \underline{e}^{\mu} \underline{e}^{\mu T}.$$

Similarly $\underline{X}^T \underline{X} = \sum_{\mu} (\alpha^{\mu})^2 \underline{f}^{\mu} \underline{f}^{\mu T}$, so $(\alpha^{\mu})^2 = \lambda^{\mu}$ ||
(Because we can express a symmetric matrix in form $\sum_{\mu} \lambda_{\mu} \underline{e}^{\mu} \underline{e}^{\mu T}$ ||)

$\underline{X} = \sum_{\mu} \alpha^{\mu} \underline{e}^{\mu} \underline{f}^{\mu T}$ is the SVD of \underline{X}

In coordinates: $X_{ai} = \sum_{\mu} \alpha^{\mu} e_a^{\mu} f_i^{\mu}$

$$X_{ai} = \sum_{\mu, \nu} e_a^{\mu} \alpha^{\mu} \delta_{\mu\nu} f_i^{\nu}$$

$$\underline{X} = \underline{E} \underline{D} \underline{F} \quad E_{a\mu} = e_a^{\mu}, \quad D_{\mu\nu} = \alpha^{\mu} \delta_{\mu\nu}$$
$$F_{\nu i} = f_i^{\nu}$$

(13)

Fisher's Linear Discriminant

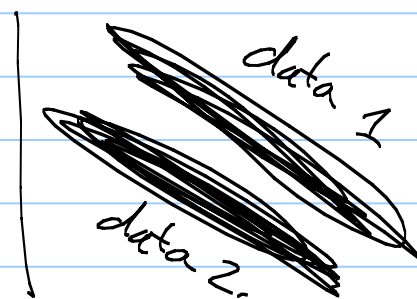
Note Title

10/15/2006

PCA may not be the best way to reduce the dimension if the goal is discrimination.

Suppose you want to discriminate between two classes of data 1 & 2.

If you put both sets of data into PCA, you will get this

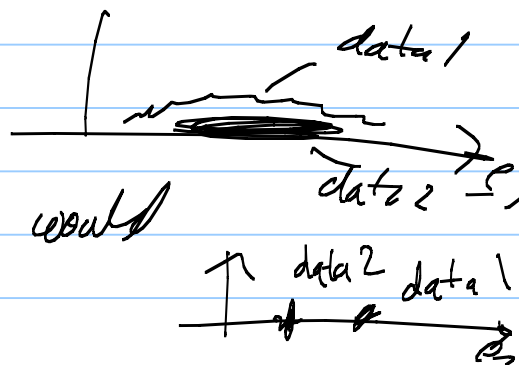


The best axis, according to PCA is in the worst direction for segmentation.

Projecting datasets onto e_1 given.



The second direction (e_2) would be far better. How to get this



(14)

Fisher's Linear Discriminant gives a way to find a better projection direction.

n_1 samples \underline{x}_i from class \mathcal{X}_1

n_2 samples \underline{x}_i from class \mathcal{X}_2

Goal: find a vector \underline{w} , project data onto this axis (i.e. $\underline{x}_i \cdot \underline{w}$) so that the data is well separated.

Define the sample means

$$\underline{m}_i = \frac{1}{N_i} \sum_{\underline{x} \in \mathcal{X}_i} \underline{x}$$

Define scatter matrices

$$\underline{S}_i = \sum_{\underline{x} \in \mathcal{X}_i} (\underline{x} - \underline{m}_i) (\underline{x} - \underline{m}_i)^T$$

Define the between-class scatter.

$$\underline{S}_B = (\underline{m}_1 - \underline{m}_2) (\underline{m}_1 - \underline{m}_2)^T$$

within-class scatter

$$\underline{S}_W = \underline{S}_1 + \underline{S}_2.$$

(15)

Now project onto an (unknown) direction $\underline{\omega}$.

$$\hat{m}_i = \frac{1}{n_i} \sum_{x \in X_i} \underline{\omega} \cdot \underline{x} = \underline{\omega} \cdot \underline{m}_i$$

// The means of the projections are the projections of the means.

The scatter of the projected points is

$$\begin{aligned} S_i^2 &= \sum_{x \in X_i} (\underline{\omega} \cdot \underline{x} - \underline{\omega} \cdot \underline{m}_i)^2 \\ &= \underline{\omega}^T \underline{S}_i \underline{\omega} \end{aligned}$$

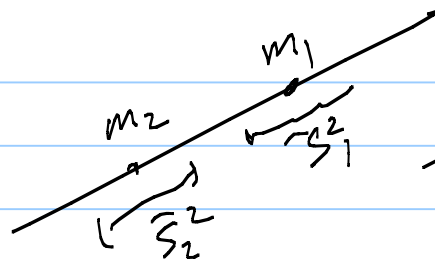
Fisher's Criterion:

Choose the projection direction $\underline{\omega}$ to maximize:

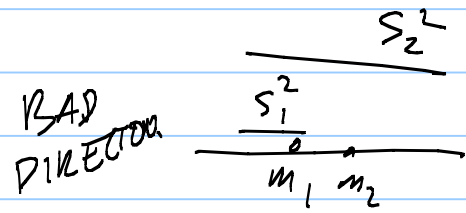
$$J(\underline{\omega}) = \frac{|\hat{m}_1 - \hat{m}_2|^2}{S_1^2 + S_2^2}$$

maximizes the ratio of the between-class distance to the within-class scatter.

(16)



This is a good projection direction.



Result: The projection direction that maximizes $J(\underline{w})$ is $\underline{w} = S^{-1}(\underline{m}_1 - \underline{m}_2)$.

Proof.
maximize $\underline{w}^T S \underline{w} - \lambda (\underline{w}^T S \underline{w} - \tau)$
 $\frac{\partial}{\partial \underline{w}}$ Lagrange multiplier const.

$$\rightarrow S \underline{w} - \lambda S \underline{w} = 0$$

Hence $S^{-1} S \underline{w} = \lambda \underline{w}$.

But $S = (\underline{m}_1 - \underline{m}_2)(\underline{m}_1 - \underline{m}_2)^T$
 $S \underline{w} = \rho (\underline{m}_1 - \underline{m}_2)$ for some ρ .

Hence $S \hat{\underline{w}} \propto (\underline{m}_1 - \underline{m}_2)$.

(17)

Fisher's Linear Discriminant

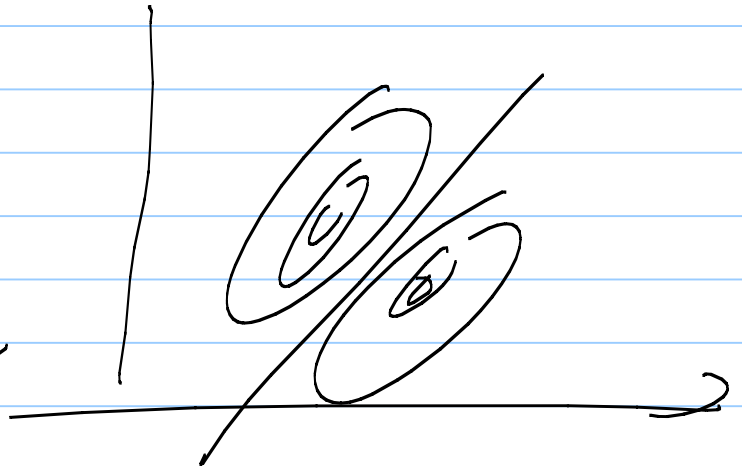
If the data comes from two Gaussians with same covariance $\underline{\Sigma}$ and means $\underline{\mu}_1, \underline{\mu}_2$.

Then the Bayes classifier is a straight line whose normal is the direction \underline{w}

$$\underline{w} \cdot \underline{x} + w_0 = 0, \quad \underline{w} = \underline{\Sigma}^{-1} (\underline{\mu}_1 - \underline{\mu}_2).$$

But. if the data comes from two Gaussians with different covariances,

then Bayes classifier is a quadratic curve, so it differs from Fisher's linear discriminant.



(18)

Multiple Classes:

For c classes, compute $c-1$ discriminants
Project D -dimensional features into $c-1$ space.

Within-class

$$S_w = S_1 + \dots + S_{c-1}$$

Between-class

$$S_b = S_{\text{total}} - S_w$$

S_{total} is the scatter
matrix for all the classes.

$$= \sum_{i=1}^c n_i (\underline{m}_i - \underline{m})(\underline{m}_i - \underline{m})^T$$

Multiple Discriminant Analysis:

Seek vectors $\omega_i: i=1, \dots, c-1$,
Project samples to $c-1$ dim space:

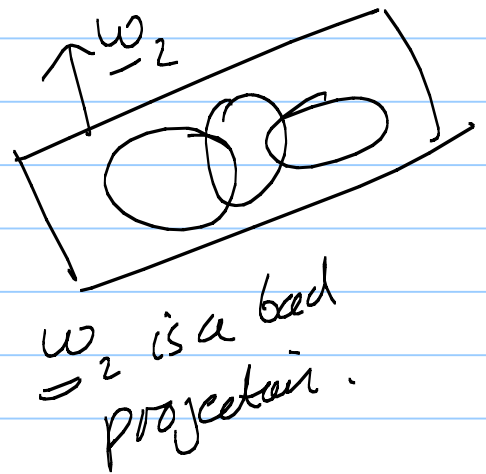
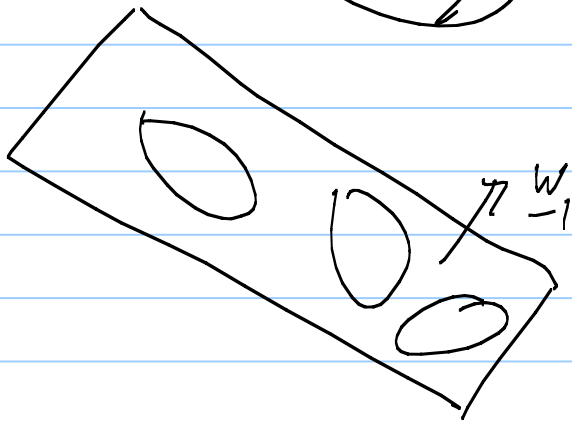
$$(\omega_i \cdot X, \dots, \omega_{c-1} \cdot X) = \underline{\omega}^T X$$

$$\text{Criterion is } J(\omega) = \frac{|\underline{\omega}^T S_b \underline{\omega}|}{|\underline{\omega}^T S_w \underline{\omega}|} \quad \text{I.e. is determinant}$$

(19)

The solution is given by the
eigenvectors, whose eigenvalues are the
C-1 largest in $\sum_B w = \lambda \sum_w w$.

w_1
is a
good
projection
of the data.



w_2 is a bad
projection.