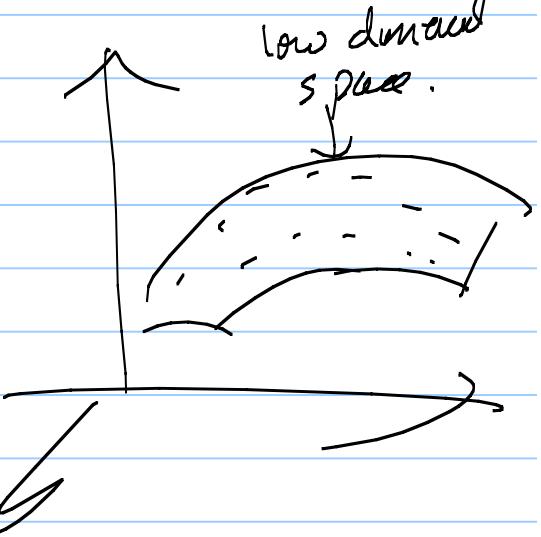


# Principal Component Analysis (PCA)

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 → e.g. multidimensional scaling. Too many to deal with in this course.



Now we discuss the most basic method

- Principal Component Analysis. (PCA)

CONVENTION:  $\underline{\mu}^T \underline{\mu}$  is a scalar  $\mu_1^2 + \mu_2^2 + \dots + \mu_D^2$   
 $\underline{\mu} \underline{\mu}^T$  is a matrix  $(\mu_1^2 \ \mu_1 \mu_2 \ \mu_1 \mu_3 \dots)$ .  
 $\mu_2^2 \dots$

(2) (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:

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

Next compute the eigenvalues and  
eigenvectors of  $\underline{\Sigma}$

Solve  $\underline{\Sigma} \underline{e} = \lambda \underline{e}$

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

Note:  $\underline{\Sigma}$  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.

$\chi$  hold  
threshold

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

(3)

Formally :

Proj.

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

where the coefficients are given by

$$a_v = (\underline{x} - \underline{\mu}) \cdot \underline{e}_v \quad \begin{array}{l} \text{(Orthogonality, meas} \\ \underline{e}_v \cdot \underline{e}_\mu = \delta_{\mu v} \end{array}$$

Here

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

Kronecker delta

no dimension reduction  
(no compression)

Then, approximate

$$\underline{x} \approx \underline{\mu} + \sum_{v=1}^m ((\underline{x} - \underline{\mu}) \cdot \underline{e}_v) \underline{e}_v$$

Projects the data into the  $m$ -dim subspace.

$$\underline{\mu} + \sum_{v=1}^m b_v \underline{e}_v \quad //$$

(4)

In 2-dimensions

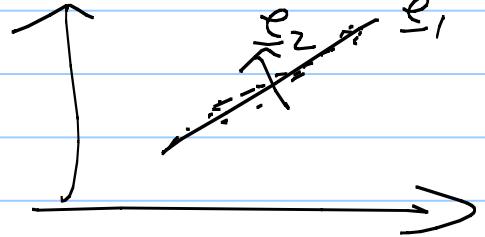
Visually



The eigenvectors of

$\Sigma$  correspond to the second order moments of the data.

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



PCA and Gaussian Distributions

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

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

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

(5)

## Cost Function for PCA

$$J(\underline{\mu}, \{a_i\}, \{e\}) = \sum_{k=1}^n \|(\underline{\mu} + \sum_{i=1}^m a_{ki} \underline{e}_i) - \underline{x}_k\|^2$$

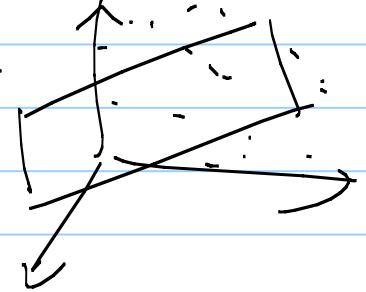
Minimizing  $J$  w.r.t.  $\underline{\mu}, \{a_i\}, \{e\}$

Data  $\{x_k : k=1 \text{ to } n\}$ ,

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

$\{\hat{e}_i\}$ 's to be the eigenvectors of

the covariance matrix  $K = \frac{1}{n} \sum_{k=1}^n (\underline{x}_k - \underline{\mu})(\underline{x}_k - \underline{\mu})^T$

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

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

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

We can re-express the criteria as

$$J[\mu, \{a\}, \{e\}] = \sum_{b=1}^N \sum_{i=1}^D ((\mu_b - x_{bi}) + \sum_{i=1}^M a_{ki} e_{ib})^2$$

where  $b$  denotes the vector components.

This is an example of a general class of problem.

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

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

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

Note:  $\tilde{x}_{ak} = x_{ak} - m_a$   
 the position of the point, relative to  
 the mean.

## (7) 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

$$\underline{X} = \underline{E} \underline{D} \underline{F}$$

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

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  $\underline{X}^T \underline{X}$ ).

$\underline{e}_{a\mu}$  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  $(\underline{X} \bar{\underline{X}}^T) = \sum_{k=1}^N (\underline{X}_k - \underline{M})(\underline{X}_k - \underline{M})^T$

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

$$\text{then } (\underline{X}^T \underline{X}) (\underline{X}^T \underline{e}) = \lambda (\underline{X}^T \underline{e})$$

This relates the eigenvectors of  $\underline{X} \bar{\underline{X}}^T$  and of  $\underline{X} \underline{X}^T$ .

(Calculate the eigenvectors for the smallest matrix, then deduce those of the bigger matrix -  $D < N$ )

(8)

Maximizing:

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

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  $\underline{X}$ .

But there is an ambiguity.

$$\sum_{v=1}^M \psi_{av} \phi_{vk} = \underline{\underline{\psi}} \underline{\underline{\phi}}_{ak} \quad \text{matrix multiplication.}$$
$$= \underline{\underline{\psi}} \underline{\underline{A}} \underline{\underline{A}}^{-1} \underline{\underline{\phi}}_{ak}$$

for any  $M \times M$  invertible matrix  $\underline{\underline{A}}$  gets rid of

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

This gets the ambiguity.

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

(9) Relate SVD to PCA (Linear Algebra)

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

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

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

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

By standard linear algebra.

$\underline{\underline{X}} \underline{\underline{X}}^T \underline{\underline{e}}^m = \lambda^m \underline{\underline{e}}^m$   $n$  eigenvalues  $\lambda^m$   
 eigenvectors  $\underline{\underline{e}}^m$   
 eigenvectors are orthogonal  $\underline{\underline{e}}^m \cdot \underline{\underline{e}}^v = \delta_{mv}$ .

Similarly  $\underline{\underline{X}} \underline{\underline{X}}^T \underline{\underline{f}}^n = \tau^n \underline{\underline{f}}^n$   $m$  eigenvalues  $\tau^n$   
 eigenvectors  $\underline{\underline{f}}^n$   
 $\underline{\underline{f}}^n \cdot \underline{\underline{f}}^v = \delta^{nv}$ .

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

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

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

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

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

(10)

Claim: we can express

$$\underline{\underline{X}} = \sum_{\mu} \underline{\lambda^{\mu}} \underline{\underline{e^{\mu}}} \underline{\underline{f^{\mu T}}} \quad \text{for some } \underline{\lambda^{\mu}}$$

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

Verify the claim

$$\underline{\underline{X}} \underline{\underline{f^{\nu}}} = \sum_{\mu} \underline{\lambda^{\mu}} \underline{\underline{e^{\mu}}} \underline{\underline{f^{\mu}}} \underline{\underline{f^{\nu}}} \\ = \sum_{\mu} \underline{\lambda^{\mu}} \delta_{\mu\nu} \underline{\underline{e^{\mu}}} = \underline{\lambda^{\nu}} \underline{\underline{e^{\nu}}} \quad //$$

$$\begin{aligned} \underline{\underline{X}} \underline{\underline{X}}^T &= \sum_{\mu, \nu} \underline{\lambda^{\nu}} \underline{\underline{e^{\nu}}} \underline{\underline{f^{\nu T}}} \underline{\lambda^{\mu}} \underline{\underline{f^{\mu}}} \underline{\underline{e^{\mu T}}} \\ &= \sum_{\mu, \nu} \underline{\lambda^{\nu}} \underline{\lambda^{\mu}} \underline{\underline{e^{\nu}}} \delta_{\mu\nu} \underline{\underline{e^{\mu T}}} = \sum_{\mu} (\underline{\lambda^{\mu}})^2 \underline{\underline{e^{\mu}}} \underline{\underline{e^{\mu T}}}. \end{aligned}$$

Similarly

$$\underline{\underline{X}}^T \underline{\underline{X}} = \sum_{\mu} (\underline{\lambda^{\mu}})^2 \underline{\underline{f^{\mu}}} \underline{\underline{f^{\mu T}}}, \text{ so } (\underline{\lambda^{\mu}})^2 = \underline{\lambda^{\mu}} //$$

(Because we can express a symmetric matrix in form  $\sum_{\mu} \underline{\lambda^{\mu}} \underline{\underline{e^{\mu}}} \underline{\underline{e^{\mu T}}}$  //

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

In coordinates:  $X_{ai} = \sum_{\mu} \underline{\lambda^{\mu}} \underline{\underline{e^{\mu}_a}} \underline{\underline{f^{\mu}_i}}$

$$X_{ai} = \sum_{\mu, \nu} \underline{\underline{e^{\mu}_a}} \underline{\lambda^{\mu}} \delta_{\mu\nu} \underline{\underline{f^{\nu}_i}}$$

$$\underline{\underline{X}} = \underline{\underline{E}} \underline{\underline{D}} \underline{\underline{F}} \quad E_{ai} = \underline{\underline{e^{\mu}_a}}, D_{\mu\nu} = \underline{\lambda^{\mu}} \delta_{\mu\nu}, F_{\nu i} = \underline{\underline{f^{\nu}_i}}$$

(11)

## Effectiveness of PCA.

In practice, PCA is often effective at unimany reduction of data dimension

But it will not be effective for  
some problems -

For example, if the data is a set  
of strings

$$(1, 0, 0, 0, \dots) = x_1$$
$$(0, 1, 0, 0, \dots) = \bar{x}_2$$

$$(6, 0, 0, 0, \dots, 0) = x_n$$

then the eigenvalues do not fall off  
as PCA requires.