

UL: PRINCIPAL COMPONENT ANALYSIS

In this Chapter, we focus on methods for dimension reduction. In particular, we are interested in projective methods = find low dimensional projections that extract / summarize useful information about the data.

- Why dimension reduction?
 - Data visualization [biplots]
 - Extracting key low dimensional features
(can be used as a pre-processing step in a SL problem)
 - Curse of dimensionality / Rates of convergence .
- In this chapter we discuss Principal Component Analysis (PCA), kernel PCA (kPCA) and probabilistic PCA (pPCA) .

I - PRINCIPAL COMPONENT ANALYSIS.

I. 1. Matrix Decompositions.

We first review key concepts in linear algebra, which will be useful later on.

- The SPECTRAL DECOMPOSITION of a ($p \times p$) symmetric matrix $\underline{\underline{S}}$ is

$$\underline{\underline{S}} = \lambda_1 \underline{\underline{e}}_1 \underline{\underline{e}}_1^T + \dots + \lambda_p \underline{\underline{e}}_p \underline{\underline{e}}_p^T,$$

where → $\lambda_1, \dots, \lambda_p$ are the eigenvalues of $\underline{\underline{S}}$
→ $\underline{\underline{e}}_1, \dots, \underline{\underline{e}}_p$ are the associated normalized eigenvectors,
i.e. $\underline{\underline{e}}_i^T \underline{\underline{e}}_i = 1$ and $\underline{\underline{e}}_i^T \underline{\underline{e}}_j = 0 \quad \forall i \neq j$

Put $\underline{\underline{P}} = \begin{pmatrix} \underline{\underline{e}}_1 & \dots & \underline{\underline{e}}_p \\ (k \times k) \end{pmatrix}$ so that $\underline{\underline{P}} \underline{\underline{P}}^t = \underline{\underline{P}}^t \underline{\underline{P}} = \underline{\underline{I}}$ (2)
 $\underline{\underline{P}}$ is ORTHOGONAL

The spectral decomposition of $\underline{\underline{S}}$ can be written

$$\boxed{\underline{\underline{S}} = \underline{\underline{P}} \underline{\Lambda} \underline{\underline{P}}^t},$$

where $\underline{\Lambda}$ is a diagonal matrix, $\underline{\Lambda} = \begin{pmatrix} \lambda_1 & & 0 \\ 0 & \ddots & \lambda_p \end{pmatrix}$

→ A $(p \times p)$ symmetric matrix $\underline{\underline{S}}$ is (strictly) positive definite iff every eigenvalue is positive: $\lambda_i > 0$.

→ If $\underline{\underline{S}}$ is a (strictly) positive definite symmetric matrix, we have that

$$\underline{\underline{S}}^{-1} = \underline{\underline{P}} \underline{\Lambda}^{-1} \underline{\underline{P}}^t = \sum_{i=1}^p \frac{1}{\lambda_i} \underline{\underline{e}}_i \underline{\underline{e}}_i^t,$$

since

$$(\underline{\underline{P}} \underline{\Lambda}^{-1} \underline{\underline{P}}^t)(\underline{\underline{P}} \underline{\Lambda} \underline{\underline{P}}^t) = \underline{\underline{P}} \underline{\Lambda}^{-1} \underline{\Lambda} \underline{\underline{P}}^t = \underline{\underline{P}} \underline{\underline{P}}^t = \underline{\underline{I}}$$

→ Put $\underline{\Lambda}^{1/2} = \begin{pmatrix} \lambda_1^{1/2} & & 0 \\ 0 & \ddots & \lambda_p^{1/2} \end{pmatrix}$.

The SQUARE ROOT of a Matrix $\underline{\underline{S}}$ is
 (symm positive definite)

$$\boxed{\underline{\underline{S}}^{1/2} = \underline{\underline{P}} \underline{\Lambda}^{1/2} \underline{\underline{P}}^t}.$$

Note that . $\underline{\underline{S}}^{1/2} \underline{\underline{S}}^{1/2} = \underline{\underline{S}}$ (\rightarrow hence the name)

$$\cdot \underline{\underline{S}}^{-1/2} = (\underline{\underline{S}}^{1/2})^{-1} = \underline{\underline{P}} \underline{\Lambda}^{-1/2} \underline{\underline{P}}^t$$

$$\cdot \underline{\underline{S}}^{1/2} \underline{\underline{S}}^{-1/2} = \underline{\underline{S}}^{-1/2} \underline{\underline{S}}^{1/2} = \underline{\underline{I}}$$

$$\cdot \underline{\underline{S}}^{-1/2} \underline{\underline{S}}^{-1/2} = \underline{\underline{S}}^{-1}$$

- The SINGULAR VALUE DECOMPOSITION (SVD) of (3)
a $(n \times p)$ matrix X of rank $r \leq p$ ($p < n$) is:

$$\boxed{X = \underline{U}_r \underline{\Lambda}_r \underline{V}_r^t}$$

- where
 . $\underline{U}_r = (n \times r)$ orthonormal matrix $\underline{U}_r^t \underline{U}_r = I$
 . $\underline{V}_r = (p \times r)$ orthonormal matrix $\underline{V}_r^t \underline{V}_r = I$
 . $\underline{\Lambda}_r = (r \times r)$ diagonal matrix with entries > 0

The singular values.

→ Put $\underline{U}_r = \begin{pmatrix} | & | \\ \underline{u}_1 & \dots & \underline{u}_r \\ | & | \end{pmatrix}$, \underline{u}_i is $(n \times 1)$

$$\underline{V}_r = \begin{pmatrix} | & | \\ \underline{v}_1 & \dots & \underline{v}_r \\ | & | \end{pmatrix}, \underline{v}_i \text{ is } (p \times 1)$$

Then

$$\underline{X}^t \underline{X} = \underline{V} \underline{\Lambda} \underline{U}^t \underline{U} \underline{\Lambda} \underline{V}^t = \underline{V} \underline{\Lambda}^2 \underline{V}^t$$

$$\Rightarrow (\underline{X}^t \underline{X}) \underline{v}_i = \lambda_i^2 \underline{v}_i$$

and

$$\underline{X} \underline{X}^t = \underline{U} \underline{\Lambda} \underline{V}^t \underline{V} \underline{\Lambda} \underline{U}^t = \underline{U} \underline{\Lambda}^2 \underline{U}^t$$

$$\Rightarrow (\underline{X} \underline{X}^t) \underline{u}_i = \lambda_i^2 \underline{u}_i$$

i.e. $\underline{X} \underline{X}^t$ has eigenvalue-eigenvector pair $(\lambda_i^2, \underline{u}_i)$

$\underline{X}^t \underline{X}$ has _____ " _____ $(\lambda_i^2, \underline{v}_i)$

→ $\underline{X} \underline{X}^t$ and $\underline{X}^t \underline{X}$ have the same eigenvalues.

I.2 Derivation of Principal Components.

(4)

Given a set of observations $\mathcal{D}_n = \{x_1, \dots, x_n\}$, $x_i \in \mathbb{R}^d$, you want to find a direction $v \in \mathbb{R}^d$ along which the projections $\langle x_i, v \rangle = x_i^t v$ give a 'good' one-dimensional approximation of the original data.

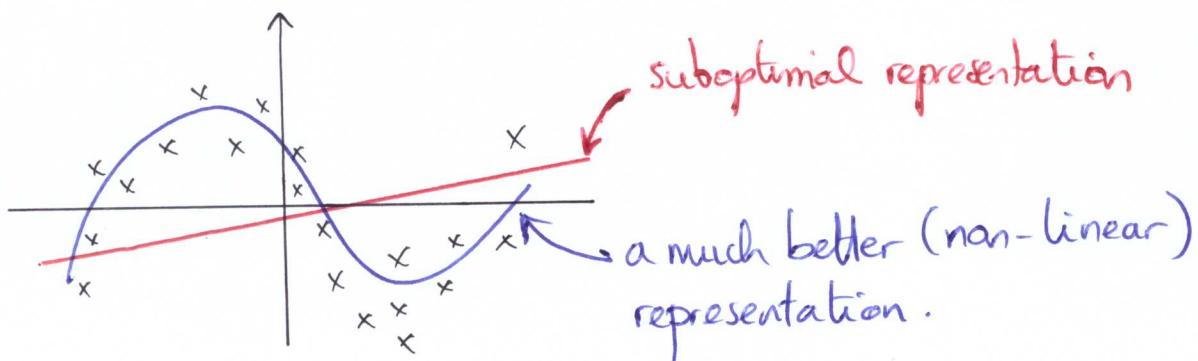
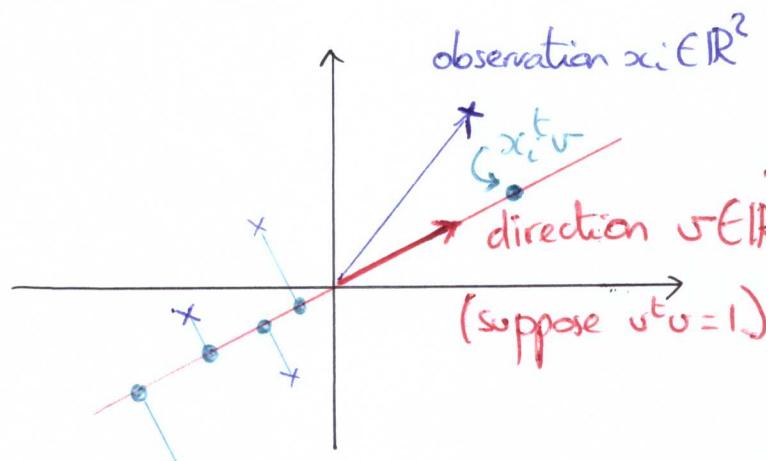
lose as little information as possible.

We consider linear combination of the covariates.

Idea: The variability observed can be accounted for by a small number of such linear combinations.

(→ the so-called principal components)

Drawback: We will only be able to detect/find linear trends, which in many cases will be suboptimal:



Notation: Let $X = (X_1, \dots, X_d)^t \in \mathbb{R}^d$ be a random vector with covariance matrix Σ .

Eigenvalues of Σ are given by $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0$.

We consider linear combinations $\begin{cases} z_1 = \underline{a}_1^t \underline{X} \\ \vdots \\ z_d = \underline{a}_d^t \underline{X} \end{cases}$ (5)
 (number of linear comb.
 considered = dimension of \underline{X})

$$\underline{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_d \end{pmatrix}$$

\underline{a}_i is a $(d \times 1)$ vector

Then $\rightarrow \text{Var } z_i = \underline{a}_i^t \sum \underline{a}_i$

$\rightarrow \text{Cov}(z_i, z_j) = \underline{a}_i^t \sum \underline{a}_j$

The principal components are uncorrelated linear combinations z_1, \dots, z_d whose variances are as large as possible.
 (variance = indicator of how much variability is in the data)

- First Principal Component = Linear combination $\underline{a}_1^t \underline{X}$ that maximizes $\text{Var}(\underline{a}_1^t \underline{X})$, subject to $\underline{a}_1^t \underline{a}_1 = 1$
 $(\text{otherwise } \text{Var}(\underline{a}_1^t \underline{X}) \text{ could be made arbitrarily large})$

- Second Principal Component = Linear combination $\underline{a}_2^t \underline{X}$ that maximizes $\text{Var}(\underline{a}_2^t \underline{X})$, subject to
 $\begin{cases} \underline{a}_2^t \underline{a}_2 = 1 \\ \text{Cov}(\underline{a}_1^t \underline{X}, \underline{a}_2^t \underline{X}) = 0 \end{cases}$

Second PC is uncorrelated with the first one.

- And So on.../...

Q: How to find directions $\underline{a}_1, \dots, \underline{a}_d$? We need the following lemma:

Lemma: Let S be a symmetric positive definite matrix with eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0$ and associated eigenvectors e_1, \dots, e_p .

Then

$$\max_{x \neq 0} \frac{x^T S x}{x^T x} = \lambda_1, \text{ attained at } x = e_1,$$

$$\min_{x \neq 0} \frac{x^T S x}{x^T x} = \lambda_p, \text{ attained at } x = e_p$$

$$\max_{x \perp e_1, \dots, e_l} \frac{x^T S x}{x^T x} = \lambda_{l+1}, \text{ attained at } x = e_{l+1}$$

(l=1, ..., k-1)

proof = Let • $P = \begin{pmatrix} | & | \\ e_1 & \dots & e_p \\ | & | \end{pmatrix}$ • $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$

Put • $S^{1/2} = P \Lambda^{1/2} P^T$

• $z = P^T x$

(note that $x \neq 0 \Rightarrow P^T x = z \neq 0$)

since orthogonal $\Rightarrow 0 \neq x = P z$

* We have

$$\frac{x^T S x}{x^T x} = \frac{x^T S^{1/2} S^{1/2} x}{x^T x} = \frac{x^T (P \Lambda^{1/2} P^T) (P \Lambda^{1/2} P^T) x}{x^T (P P^T) x} \stackrel{I}{=} \frac{z^T \Lambda z}{z^T z}$$

$$\Rightarrow \frac{x^T S x}{x^T x} = \frac{\sum_i \lambda_i z_i^2}{\sum z_i^2} \leq \frac{\sum_i (\max \lambda_j) z_i^2}{\sum z_i^2} = \lambda_1$$

Moreover, setting $x = e_1$ gives $z = P^T e_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$ since

$$e_j^T e_i = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{o/w} \end{cases}$$

For this choice of x we have $\frac{z^T \Lambda z}{z^T z} = \lambda_1 = \text{the max value.}$

$$\Rightarrow \max_{x \neq 0} \frac{x^T S x}{x^T x} = \frac{e_1^T \Lambda e_1}{e_1^T e_1} = \lambda_1$$

* A similar argument replacing max by min gives (7)

$$\min_{x \neq 0} \frac{x^t S x}{x^t x} = \frac{e_p^t \Sigma e_p}{e_p^t e_p} = \lambda_p.$$

* Now, $x = Pz$

$$= z_1 e_1 + \dots + z_p e_p$$

$$\text{So } x \perp e_1, \dots, e_l \Rightarrow 0 = e_i^t x = z_1 e_i^t e_1 + \dots + z_p e_i^t e_p \quad (i \leq l) \\ = z_i$$

Thus,

$$\frac{x^t S x}{x^t x} = \frac{\sum_{i=l+1}^k \lambda_i z_i^2}{\sum_{l+1}^k z_i^2} \leq \lambda_{l+1}$$

Taking $z_{l+1} = 1$ and $z_{l+2} = \dots = z_p$ gives the max ■

We have the following result:

The i -th principal component of \underline{X} is given by

$$z_i = \underline{e}_i^t \underline{X} \quad (= \text{coordinate of } \underline{X} \text{ in the } \underline{e}_i \text{ direction}),$$

where $\underline{e}_1, \dots, \underline{e}_d$ are the eigenvectors of $\underline{\Sigma}$ (cov matrix of \underline{X}) associated with the eigenvalues $\lambda_1 \geq \dots \geq \lambda_d > 0$.

$$\text{Moreover, } \begin{cases} \text{var } z_i = \underline{e}_i^t \underline{\Sigma} \underline{e}_i = \lambda_i \\ \text{cov}(z_i, z_j) = \underline{e}_i^t \underline{\Sigma} \underline{e}_j = 0 \end{cases}$$

Remark: In practice, we use the eigenvalues and eigenvectors of the sample covariance matrix S , calculated from a random sample $\mathcal{D}_n = \{\underline{X}_1, \dots, \underline{X}_n\}$, $\underline{X}_i \in \mathbb{R}^d$:

$$S = \frac{1}{n} \underline{X}^t \underline{X}, \text{ where } \begin{matrix} (d \times d) & (d \times n) & (n \times d) \end{matrix}$$

$$\underline{X} = \begin{pmatrix} X_{11} & \dots & X_{1d} \\ \vdots & & \vdots \\ X_{n1} & \dots & X_{nd} \end{pmatrix} = \begin{pmatrix} \underline{X}_1^t \\ \vdots \\ \underline{X}_n^t \end{pmatrix}$$

Using the SVD decomposition of $\underline{X} = \underline{U} \Lambda \underline{V}^t$,

$\underline{X}^t \underline{X} = \underline{V} \Lambda^2 \underline{X}^t$ so that $\underline{X}^t \underline{X}$ has eigenvalue-eigenvector pair $(\lambda_i^2, \underline{v}_i)$ (see page 3)

Thus \underline{S} has e-e pair $(\frac{\lambda_i^2}{n}, \underline{v}_i)$ [Assume matrix \underline{X} centered: $\underline{S} = n^{-1} \underline{X}^t \underline{X}$]

- The i -th principal component is thus $\underline{z}_i = \underline{X} \underline{v}_i$,
 $(n \times 1) \quad (n \times d) (d \times 1)$

and its variance is $\frac{\lambda_i^2}{n}$.

[n realizations of the i -th PC
of \underline{X}]

Putting all d PC into a matrix =

$$\underline{z} = \underline{X} \underline{V} = \underline{U} \Lambda \underline{V}^t \underline{V} = \underline{U} \Lambda \underline{U}^t \Rightarrow \underline{z}_i = \lambda_i \underline{u}_i$$

$(n \times d) \quad (n \times d) \quad (d \times d) \quad (n \times 1) \quad (n \times 1)$

$$\begin{matrix} & d \\ \begin{matrix} n & | & | & | & | \\ & | & - & - & | \\ & | & & & | \end{matrix} & = \underline{z} \end{matrix}$$

n realizations
of the
first PC n realizations of the
second PC, and so on ... / ...

Proof = To find the first principal component, we need to

solve

$$\begin{aligned} & \text{maximize } \underline{a}^t \underline{\Sigma} \underline{a} \\ & \underline{a} \neq 0 \\ & \text{subject to } \underline{a}^t \underline{a} = 1 \end{aligned}$$

\Leftrightarrow

$$\begin{aligned} & \text{maximize } \frac{\underline{a}^t \underline{\Sigma} \underline{a}}{\underline{a}^t \underline{a}} \\ & \underline{a} \neq 0 \\ & \text{subject to } \underline{a}^t \underline{a} = 1 \end{aligned}$$

Recall that $\max_{\underline{a} \neq 0} \frac{\underline{a}^t \underline{\Sigma} \underline{a}}{\underline{a}^t \underline{a}} = \frac{\underline{e}_1^t \Lambda \underline{e}_1}{\underline{e}_1^t \underline{e}_1} = \lambda_1$, with $\underline{e}_1^t \underline{e}_1 = 1$
 $\quad \quad \quad$ (Lemma p.5) $\underline{e}_1^t \underline{\Sigma} \underline{e}_1 = \text{Var } \underline{z}_1$

And similarly for the subsequent principal components ■

I.3. Eckart-Young Theorem

(9)

Let \underline{X} be an $(n \times d)$ matrix with $n \geq d$. A common problem (with diverse applications) is to find approximations of \underline{X} that are of rank r lower than rank \underline{X} .

"[↑]
low-rank
approximation problem"

Finding a 'best' low-rank approximation requires to specify an objective function to minimize, in terms of the matrix norm of the difference between the original matrix \underline{X} and its approximate. When the norm is the FROBENIUS NORM, the Eckart-Young (EY) theorem provides an answer.

EY Theorem.

- Let \underline{X} be an $(n \times d)$ matrix of observations of rank $d \leq n$. (centered)
- $\underline{X} = U \Lambda V^t$ its SVD decomposition.
- $(n \times d)(n \times d)(d \times d)(d \times d)$

Then the solution to the optimization problem

$$\underline{A}^* = \underset{\substack{A \in \mathbb{R}^{n \times d} \\ \text{rank } \underline{A} = r \leq d}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \sum_{j=1}^d (x_{ij} - a_{ij})^2 \right\}$$

\underline{A} has entries a_{ij}

Frobenius norm of $\underline{X} - \underline{A}$.

is given by

$$\underline{A}^* = U \Lambda^* V^t$$

$$\text{where } \Lambda^* = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_r & 0 & \dots & 0 \end{pmatrix}$$

$$\text{Moreover, } \sum_{i=1}^n \sum_{j=1}^d (x_{ij} - a_{ij})^2 = \sum_{k=r+1}^d \lambda_k^2$$

Sum of the variance of the remaining PC.

(10)

• Interpretation.

- Recall that the i -th PC is $\underline{z}_i = \underline{X} \underline{v}_i = \lambda_i \underline{u}_i$.

Put $\underline{z}_i = \begin{pmatrix} z_{1i} \\ \vdots \\ z_{ni} \end{pmatrix} \in \mathbb{R}^n$, $i = 1, \dots, d$, $\underline{z} = \begin{pmatrix} z_{11} & z_{12} & \dots & z_{1d} \\ | & | & & | \\ z_{n1} & z_{n2} & \dots & z_{nd} \end{pmatrix}$

- The j -th row of \underline{X} (corresponding to the j -th observation) is thus approximated by the j -th row of $\underline{A}^* = \underline{U} \underline{\Lambda}^* \underline{V}^t$, i.e.

our rank r approx

$$\underline{\hat{x}}_j^t = \underline{a}_j^t = \begin{pmatrix} z_{j1} \\ \vdots \\ z_{jr} \end{pmatrix}_{(1 \times d)}^t + \begin{pmatrix} z_{j2} \\ \vdots \\ z_{jr} \end{pmatrix}_{(1 \times d)}^t + \dots + \begin{pmatrix} z_{jr} \\ \vdots \\ z_{jr} \end{pmatrix}_{(1 \times d)}^t, \quad r \leq d, \quad j = 1, \dots, n$$

Principal components

Principal directions
= eigenvectors of
the sample cov matrix

$S = n^{-1} \underline{X}^t \underline{X}$, \underline{X} assumed
centered

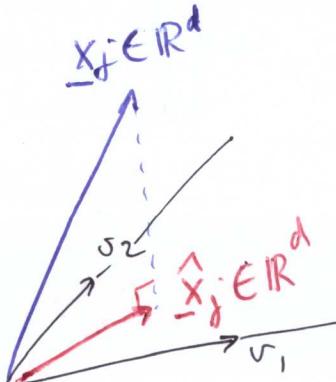
$$\underline{A}^* = \begin{pmatrix} -\frac{a_1^t}{\lambda_1} \\ \vdots \\ -\frac{a_n^t}{\lambda_n} \end{pmatrix}$$

- Note that $z_{ji} = \underline{x}_j^t \underline{v}_i$ since

$$\begin{bmatrix} \underline{z}_i \\ \vdots \\ \underline{z}_i \end{bmatrix} = \begin{bmatrix} -x_1^t \\ -x_j^t \\ -x_n^t \end{bmatrix} \begin{bmatrix} z_{1i} \\ \vdots \\ z_{ni} \end{bmatrix}$$

$$\Rightarrow \underline{\hat{x}}_j^t = \langle \underline{x}_j, \underline{v}_1 \rangle \underline{v}_1^t + \dots + \langle \underline{x}_j, \underline{v}_r \rangle \underline{v}_r^t$$

$\underline{\hat{x}}_j^t$ = \perp projection of \underline{x}_j onto the column space of the first r columns of \underline{V} . $(d \times d)$



Take Away Message

The approximation $\underline{\hat{x}}_j^t$ is an orthogonal projection of \underline{x}_j ; minimizing the squared error between \underline{x}_j and any approximating vector

Proof of the EY theorem.

(10a)

Step I

Let $\underline{e}_1, \dots, \underline{e}_r$ be r orthonormal vectors, and

$$\underline{\underline{E}} = \begin{pmatrix} \underline{e}_1^T & \dots & \underline{e}_r^T \end{pmatrix} ; \quad \underline{\underline{E}}^T \underline{\underline{E}} = \underline{\underline{I}}_{(r \times r)}$$

↑ our goal will be to find the entries of $\underline{\underline{E}}$.

Our goal is to minimize the quantity

$$\sum_{i=1}^n \sum_{j=1}^d (x_{ij} - a_{ij})^2 = \sum_{i=1}^n (\underline{x}_i - \underline{a}_i)^T (\underline{x}_i - \underline{a}_i)$$

over the set of matrices $\underline{\underline{A}} \in \mathbb{R}^{n \times d}$, with rank $\underline{\underline{A}} = r \leq d$.

$$\underline{\underline{A}} = \begin{pmatrix} \underline{a}_1^T \\ \vdots \\ \underline{a}_n^T \end{pmatrix}$$

$$\underline{\underline{X}} = \begin{pmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{pmatrix}$$

This guy is our \underline{a}_i .
We will first derive the optimal
 \underline{b}_i minimizing the sq error, and
then the columns of $\underline{\underline{E}}$

First, let's compute the square norm of $\underline{x}_i - \underline{\underline{E}} \underline{b}_i$,
for an arbitrary vector $\underline{b}_i \in \mathbb{R}^r$.

$$\begin{aligned} \underline{x}_i - \underline{\underline{E}} \underline{b}_i &= \underline{x}_i - \underline{\underline{E}} \underline{\underline{E}}^T \underline{x}_i + \underline{\underline{E}} \underline{\underline{E}}^T \underline{x}_i - \underline{\underline{E}} \underline{b}_i \\ &= (\underline{\underline{I}} - \underline{\underline{E}} \underline{\underline{E}}^T) \underline{x}_i + \underline{\underline{E}} (\underline{\underline{E}}^T \underline{x}_i - \underline{b}_i) . \end{aligned}$$

$$\bullet (x_i - Eb_i)^t (x_i - Eb_i) = x_i^t (I - EE^t) x_i \quad \text{independent} \\ + (E^t x_i - b_i)^t E^t E (E^t x_i - b_i)$$

Since the cross products vanish : $(I - EE^t)E^t = E - E = 0$
 $E^t(I - EE^t) = E^t - E^t = 0$

$$\Rightarrow (x_i - Eb_i)^t (x_i - Eb_i) = x_i^t (I - EE^t) x_i \\ + \underbrace{(E^t x_i - b_i)^t (E^t x_i - b_i)}$$

This term is positive unless b_i is chosen such that $b_i = E^t x_i$
 $\Rightarrow Eb_i = EE^t x_i$

Now, we are looking
for the columns of E

\downarrow
 \Rightarrow With $a_i = \underbrace{EE^t x_i}_{(dx_1)(dx_r)(rxl)(dx_1)}$, we get

$$\sum_{i=1}^n (x_i - a_i)^t (x_i - a_i) = \sum_{i=1}^n x_i^t (I - \underbrace{EE^t}_{\text{independent}}) x_i \\ = \underbrace{\sum_{i=1}^n x_i^t x_i}_{\text{minimizing this}} - \underbrace{\sum_{i=1}^n x_i^t EE^t x_i}_{\text{maximizing this one}}$$

minimizing this = maximizing this one.

OK, so it remains to find a matrix E which will make $\sum_{i=1}^n x_i^t EE^t x_i$ as large as possible.

(\triangle Recall that E has l vectors of unit length.)

Well,

$$\sum_{i=1}^n x_i^t EE^t x_i = \sum_{i=1}^n \text{Tr}(x_i^t \underbrace{EE^t x_i}_{\text{independent}}) = \sum_{i=1}^n \text{Tr}(EE^t x_i x_i^t)$$

and introduce $S = \frac{1}{n} \sum_{i=1}^n x_i x_i^t$ = sample cov matrix
 $(d \times d)$ $= n^{-1} X^t X$ (obs are centered) (12)

$$\Rightarrow \sum_{i=1}^n x_i^t E E^t x_i = n \text{Tr}(E E^t S)$$

$$= n \text{Tr}(E^t S E).$$

In other words, the best choice for E maximizes the sum of the diagonal elements of $E^t S E$.

↳ Selecting e_1 to maximize $e_1^t S e_1$, the first diagonal element of $E^t S E$, subject to $e_1^t e_1 = 1$, gives $e_1 = v_1$ (= eigenvector of the sample covariance matrix see page 8)

↳ For $e_2 \perp v_1$, $e_2^t S e_2$ is maximized by $e_2 = v_2$.
 $e_2^t e_2 = 1$

↳ And so on.

⇒ We find that the matrix E maximizing $\sum_{i=1}^n x_i^t E E^t x_i$, denoted E^* , is thus $E^* = \begin{pmatrix} v_1 & \dots & v_r \end{pmatrix}$.
 (dxr) $(dx1)$ $(dx1)$

⇒ The best approximation matrix $A^* = \begin{pmatrix} a_1^{*t} \\ \vdots \\ a_n^{*t} \end{pmatrix}$ is such that

- $(A^*)^t = \begin{pmatrix} a_1^* & \dots & a_n^* \end{pmatrix}$
- $= \begin{pmatrix} E^* E^{*t} x_1 & \dots & E^* E^{*t} x_n \end{pmatrix}$
- $$(n \times d) A^* = \begin{pmatrix} -x_1^t E^* E^{*t} \\ -x_n^t E^* E^{*t} \end{pmatrix}$$

(13)

$$\Rightarrow A^* = X E^* E^{*t}$$

$$= U \Lambda V^t E^* E^{*t}$$

$\underbrace{V^t E^*}_{(d \times r)}$

$$V^t E^* = \begin{pmatrix} -v_1^t & & \\ & \ddots & \\ & & -v_d^t \end{pmatrix} \begin{pmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{pmatrix}$$

$$= \begin{pmatrix} 1 & 0 & & \\ 0 & 1 & & \\ & & \ddots & \\ & & & 0 \end{pmatrix} \xrightarrow[r]{\downarrow} \xrightarrow[d]{\downarrow}$$

so that $V^t E^* E^{*t} = \begin{pmatrix} 1 & & & \\ & \ddots & & \\ & & 1 & \\ & & & 0 \end{pmatrix} \begin{pmatrix} -v_1^t & & \\ & \ddots & \\ & & -v_r^t \end{pmatrix}$

$$= \begin{pmatrix} -v_1^t & & & \\ & \ddots & & \\ & & -v_r^t & \\ & & & 0 \end{pmatrix} \xrightarrow[d]{\downarrow}$$

Complete the missing rows with v_{r+1}^t, v_d^t , and set the associated singular values $\lambda_{r+1}, \dots, \lambda_d$ to zero.

It's the same!

$$\Rightarrow A^* = U \Lambda^* V^t, \text{ where } \Lambda^* = (\lambda_1 - \lambda_0, \dots, \lambda_d - \lambda_0),$$

as asserted.

Step II . Residual error.

With this choice of A , it remains to compute the residual term

$$\sum_{i=1}^n (x_i - a_i^*)^t (x_i - a_i^*) = \underbrace{\sum_{i=1}^n x_i^t x_i}_{= n \text{Tr } S} - \underbrace{\sum_{i=1}^n x_i^t E^* E^{*t} x_i}_{= n \text{Tr}(E^{*t} S E^*)}$$

$$= n \left(\frac{\lambda_1^2}{n} + \dots + \frac{\lambda_d^2}{n} \right)$$

$$= \lambda_1^2 + \dots + \lambda_d^2$$

Term $\text{Tr}(\underbrace{\mathbf{E}^{*t} \mathbf{S} \mathbf{E}^*}_{(r \times r)}) \rightarrow$ the i -th diagonal element of $\mathbf{E}^{*t} \mathbf{S} \mathbf{E}^*$ is

$$\mathbf{v}_i^t \mathbf{S} \mathbf{v}_i = \mathbf{v}_i^t \underbrace{\frac{\lambda_i^2}{n} \mathbf{v}_i}_\text{page 8} = \frac{\lambda_i^2}{n}$$

$$\Rightarrow \text{Tr}(\mathbf{E}^{*t} \mathbf{S} \mathbf{E}^*) = \frac{\lambda_1^2}{n} + \dots + \frac{\lambda_r^2}{n},$$

and

$$\sum_{i=1}^n (\mathbf{x}_i - \mathbf{a}_i^*)^t (\mathbf{x}_i - \mathbf{a}_i^*) = \left(\lambda_1^2 + \dots + \lambda_d^2 \right) - \left(\lambda_1^2 + \dots + \lambda_r^2 \right)$$

$$= \sum_{k=r+1}^d \lambda_k^2$$

Application = Visualisation of data with more than two predictors (aka BIPLOTS)

→ With only two variables, a scatterplot provides a visual information on both the sampling unit (observation) and the predictor.

→ With more than two predictors, one may consider a matrix array of scatterplots. Alternatively, we may plot the first two principal components:

Recall that the j -th observation (j -th row of $\underline{\mathbf{X}}_{(n \times d)}$) is approximated by

est rank r
approxim. in the
Frobenius norm

$\rightarrow \underline{\mathbf{x}}_j^t \approx z_{j1} \mathbf{v}_1^t + z_{j2} \mathbf{v}_2^t + \dots + z_{jr} \mathbf{v}_r^t, \quad j=1, \dots, n$
corresponding to the j -th row of $\underline{\mathbf{A}}^*$, where $\text{rank } \underline{\mathbf{A}}^* = r \leq d$
(top of page 10)

The rank-2 approximation is thus

$$\underline{\mathbf{x}}_j^t \approx z_{j1} \mathbf{v}_1^t + z_{j2} \mathbf{v}_2^t \quad j=1, \dots, n$$

(1x1) (1x1) (1x1)

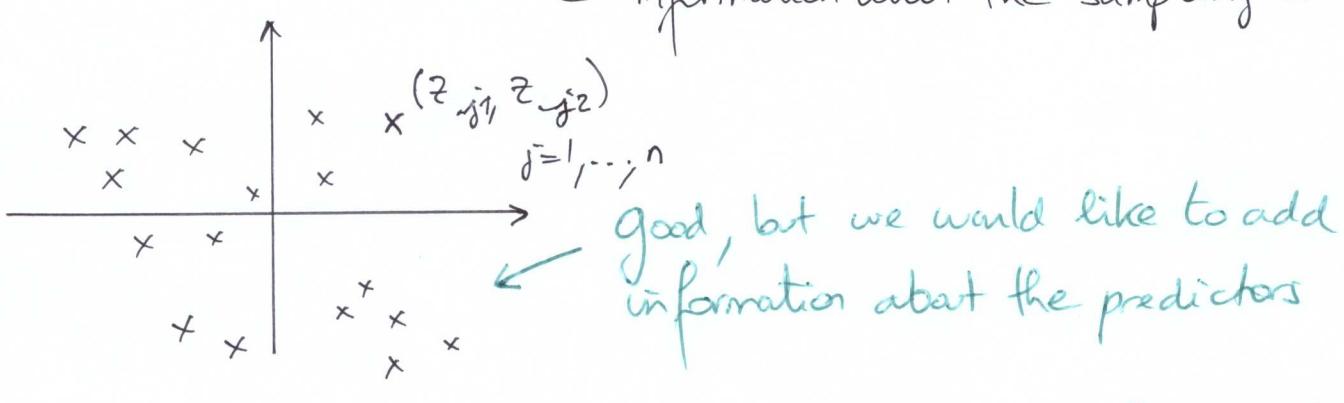
In other words:

$$\begin{pmatrix} \hat{x}_{j1} \\ \vdots \\ \hat{x}_{jd} \end{pmatrix} = z_{j1} \begin{pmatrix} v_{11} \\ \vdots \\ v_{d1} \end{pmatrix} + z_{j2} \begin{pmatrix} v_{12} \\ \vdots \\ v_{d2} \end{pmatrix}$$

↑ ↑
weights weights

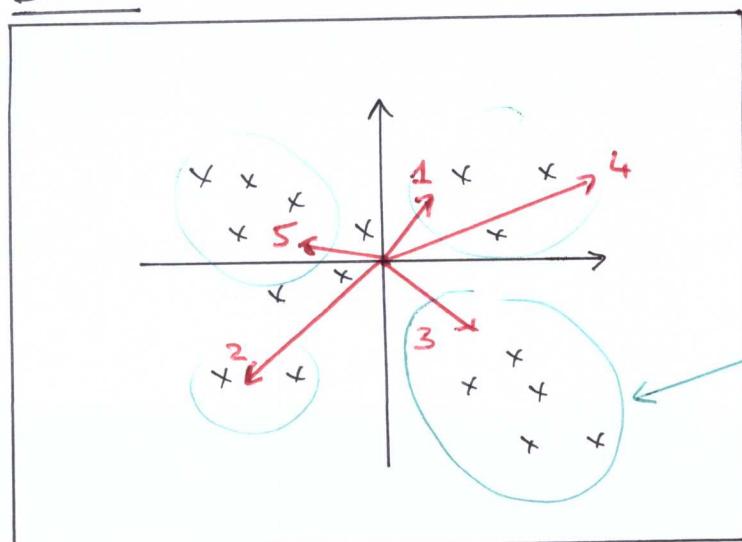
(z_{j1}, z_{j2}) ← 1st predict
 (z_{j1}, z_{j2}) ← dth predict

(z_{j1}, z_{j2}) = coordinate of \underline{x}_j in the plane defined by the first two eigenvectors.
= information about the sampling unit j .

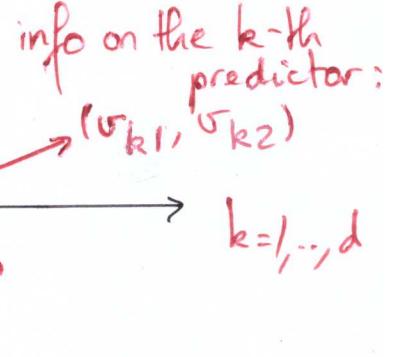


Putting these two in a single plot yields:

BIPLOT



The third variable is responsible for grouping these points here, which all have a high value associated with the 3rd predictor



II - PROBABILISTIC PCA

(16)

Before introducing probabilistic PCA (pPCA), let's summarize the conventional PCA results we obtained:

(i) Conventional PCA is a non-probabilistic approach: no probability distribution is attached to $\underline{\underline{X}} = \begin{pmatrix} \underline{x_1^t} \\ \vdots \\ \underline{x_n^t} \end{pmatrix}$

(ii) The principal directions of $\underline{\underline{X}}$ are given by the eigenvectors of the sample covariance matrix → page 8

$$\underline{\underline{S}} = \frac{1}{n} \sum_{i=1}^n (\underline{x_i} - \bar{\underline{x}})(\underline{x_i} - \bar{\underline{x}})^t = \frac{1}{n} \underline{\underline{X}}^t \underline{\underline{X}}$$

if $\underline{\underline{X}}$ is centered

Corresponding to the columns of the matrix $\underline{\underline{V}}$ in the SVD decomposition of $\underline{\underline{X}} = \underline{\underline{U}} \Lambda \underline{\underline{V}}^t$ → page 3.

(iii) The best rank r approximation (in the square error sense) of $\underline{\underline{x_j}}$ is

$$\begin{aligned} \hat{\underline{x_j}} &= z_{j1} \underline{v}_1 + \dots + z_{jr} \underline{v}_r + \bar{\underline{x}} \\ &= \langle \underline{x_j}, \underline{v}_1 \rangle \underline{v}_1 + \dots + \langle \underline{x_j}, \underline{v}_r \rangle \underline{v}_r + \bar{\underline{x}} \\ &= \left[\begin{array}{c|c} \frac{1}{\sqrt{r}} & \frac{1}{\sqrt{r}} \\ \hline \underline{v}_1 & \vdots \\ \hline \vdots & \vdots \\ \hline \underline{v}_r & \frac{1}{\sqrt{r}} \end{array} \right] \begin{bmatrix} z_{j1} \\ \vdots \\ z_{jr} \end{bmatrix} + \bar{\underline{x}} \end{aligned}$$

$(d \times 1)$

(1) → page 10

$$\boxed{\hat{\underline{x_j}} = \underline{\underline{V}}_r \underline{z}_r + \bar{\underline{x}}} \quad \begin{array}{l} \text{"reconstruction from the PC"} \\ \text{vector } \mathbb{R}^r \text{ whose entries correspond to the first } r \text{ PC of the } j\text{-th observation} \end{array}$$

The reduced-dimensionality transformation of the j -th observation is

(2) → page 10

$$\boxed{\underline{z}_r = \underline{\underline{V}}_r^t (\underline{x_j} - \bar{\underline{x}})} \quad \text{"extraction of the PC".}$$

→ p PCA challenges (i) by deriving the principal components (17) using a probabilistic model for the \underline{x} s. The model is motivated by relation (1), which links the principal components (soon to be the hidden / latent variables) to the original observations.

→ In view of (1), consider the model:

$$\text{For } i=1, \dots, n, \quad \underline{x}_i = \underline{W} \underline{\zeta}_i + \underline{\mu} + \underline{\xi}_i \quad (r \leq d)$$

(1) $\underline{\zeta}_i$ Principal Component

where • $\underline{\zeta}_i \in \mathbb{R}^r$ is an r-dimensional vector of latent variables.

- generally, $r < d$ so that the variable $\underline{\zeta}$ offers a more parsimonious description of the data.
- $\underline{\zeta}_i$ are assumed to be independent, with distribution $\underline{\zeta}_i \sim \mathcal{N}(0, \underline{\Sigma}_r)$

- $\underline{W} \in \mathbb{R}^{d \times r}$ = matrix of parameters to be estimated.
- $\underline{\xi}_i \sim \mathcal{N}(0, \underline{\Psi})$, with $\underline{\Psi}$ diagonal.

The choice of a diagonal covariance matrix for $\underline{\xi}_i$ ensures that conditionally on $\{\underline{\zeta}_i\}$ observations $\{\underline{x}_i\}$ are independent \Rightarrow the dependency structure is explained by a small number of latent variable. What is left unexplained appears in the noise term $\underline{\xi}_i$, which differs for each observation.

Given model (1), observations \underline{x}_i are also normally distributed:

$$\underline{x}_i \sim \mathcal{N}(\underline{\mu}, \underline{W}\underline{W}^t + \underline{\Psi})$$

invariant under post-multiplication

of \underline{W} by an orthogonal (\equiv rotation) matrix \underline{Q} :

$$(\underline{W}\underline{Q})(\underline{Q}^t \underline{W}^t) = \underline{W}(\underline{Q}\underline{Q}^t)\underline{W}^t = \underline{W}\underline{W}^t.$$

FACTOR ANALYSIS (FA)
(not p PCA yet)

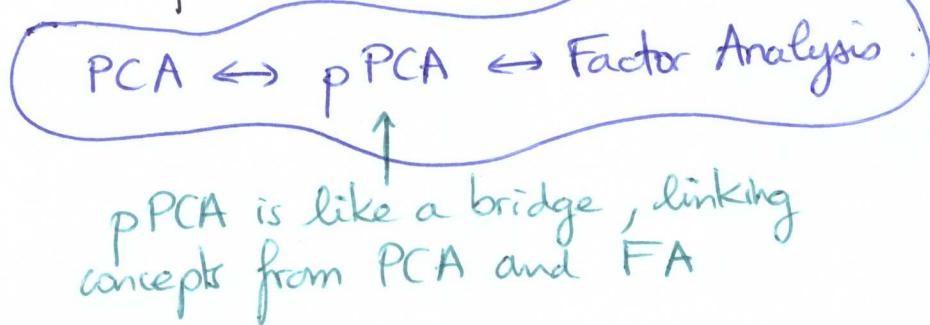
The probabilistic formulation allows ML estimation of the parameters \underline{W} and $\underline{\Psi}$, using an iterative procedure (EM algorithm), since no closed analytical expression of the MLE of \underline{W} and $\underline{\Psi}$ are available in general. (18)

Q: Is there a clear correspondance between \underline{W} and $\underline{V_r}$?
 = For example, do the columns of \underline{W} and $\underline{V_r}$ span the same space?

In general, the answer is no.

However, links between the two exist in the special case of an ISOTROPIC noise model, ie when $\underline{\Psi} = \sigma^2 \underline{I_d}$. And this is exactly the model considered by Tipping & Bishop. They named this model pPCA.

Informally,



pPCA model

$$\underline{x_i} = \underline{W} \underline{z_i} + \underline{\mu} + \underline{\xi_i}, \quad i=1, \dots, n$$

$$(dx_i) \quad (d \times r)(r \times 1) \quad (dx_1) \quad (dx_1)$$

$$\underline{z_i} \sim \mathcal{N}(0, \underline{I_r})$$

$$\underline{\xi_i} \sim \mathcal{N}(0, \sigma^2 \underline{I_d}), \quad r \leq d$$

$$\Rightarrow \underline{x_i} \sim \mathcal{N}(\underline{\mu}, \underline{W} \underline{W}^T + \sigma^2 \underline{I_d})$$

$$\text{Put } \underline{C} = \underline{W} \underline{W}^T + \sigma^2 \underline{I_d}$$

$$(d \times d) \quad (d \times r)(r \times d)$$

Goal: Estimation of \underline{W} and σ^2 + Extraction of the hidden variable and/or reconstruction from the hidden variables.

How: Maximum likelihood, of course.

(19)

$$\begin{aligned} \mathcal{L} &= \ln \left[\prod_{i=1}^n \frac{1}{(2\pi)^{d/2} |C|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu})^t C^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \right\} \right] \\ &= -\frac{n}{2} \left\{ d \ln(2\pi) + \ln |C| \right\} - \frac{1}{2} \underbrace{\sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})^t C^{-1} (\mathbf{x}_i - \boldsymbol{\mu})}_{\sum_{i=1}^n \text{Tr} \{ (\mathbf{x}_i - \boldsymbol{\mu})^t C^{-1} (\mathbf{x}_i - \boldsymbol{\mu}) \}} \\ &= \sum_{i=1}^n \text{Tr} \{ C^{-1} (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^t \} \\ &= n \text{Tr} (C^{-1} S), \end{aligned}$$

where we defined the (sample) covariance matrix as

$$\underline{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^t.$$

$$\boxed{\mathcal{L} = -\frac{n}{2} \left\{ d \ln(2\pi) + \ln |C| + \text{Tr} (C^{-1} S) \right\}}$$

→ The MLE of $\underline{\boldsymbol{\mu}}$ is the mean of the observations:

$$\hat{\boldsymbol{\mu}}_{\text{MLE}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i. \quad \text{You may again pre-process the data by centering } \underline{\mathbf{x}}.$$

→ The MLE of $\begin{cases} \underline{W} \\ \sigma^2 \end{cases}$ → Explicit expression by direct maximization of \mathcal{L} . Gives no insight about the link between $\hat{\underline{W}}_{\text{MLE}}$ and \underline{Y}_r .

Using EM algorithm.
(why bother? It allows us to deal with missing components in the \mathbf{x}_i 's.)

Remark: $\underline{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_{\text{MLE}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_{\text{MLE}})^t$, so centering the data yields $\underline{S} = \frac{1}{n} \underline{\mathbf{X}}^t \underline{\mathbf{X}}$. Substitute this expression back to \mathcal{L}

(A) Direct Maximization of \mathcal{L} .

(20)

The derivative of \mathcal{L} with respect to W is found to be:

$$\frac{\partial \mathcal{L}}{\partial W} = n(C^{-1}S C^{-1}W - C^{-1}w), \text{ where } C = WW^t + \sigma^2 I \quad (\text{d} \times \text{d})$$

\Rightarrow Stationary points satisfy the relation $SC^{-1}W = w$.

Assume $W \neq 0$ [Solution $W=0$ can be shown to be a minimum of \mathcal{L}]

& $C \neq S$ [Unrealistic Scenario of exact covariance model]

- Consider the SVD decomposition of W :

$$W = U L V^t, \text{ where } \begin{cases} U = \begin{pmatrix} | & | \\ u_1 & \dots & u_r \\ | & | \end{pmatrix}, \text{ columns are orthonormal} \\ L = \begin{pmatrix} l_1 & & 0 \\ 0 & \ddots & \\ & & l_r \end{pmatrix} \text{ some of these might be zero.} \\ V = \text{orthogonal matrix.} \end{cases}$$

Our goal: identify U , L and V .

- First, we compute the product $C^t W$:

$$\begin{aligned} C^t W &= (WW^t + \sigma^2 I)^{-1} W \\ &= W(\sigma^2 I + W^t W)^{-1} \\ &= W L V^t (\sigma^2 I + V L U^t U L V^t)^{-1} \\ &= W L V^t (\sigma^2 I + V L^2 V^t)^{-1} \\ &= W L (\sigma^2 V + V L^2)^{-1} \\ &= W L (V(\sigma^2 I + L^2))^{-1} \\ &= W L (L^2 + \sigma^2 I) V^t \end{aligned}$$

Useful: $\forall P, Q,$
 $\sigma^2 P + P Q P = P(\sigma^2 I + QP)$
 $(\sigma^2 I + PQ)P \Downarrow$
 $P(\sigma^2 I + QP)^{-1} = (\sigma^2 I + PQ)^{-1}P$

$V^t V = I$
and thus $(V^t)^{-1} = V^{-1}$

$$\text{Now, } S C^{-1} W = S \underbrace{U L (L^2 + \sigma^2 I)^{-1} V^t}_{\begin{array}{c} \parallel \\ W \\ \parallel \\ U L V^t \end{array}} \quad (21)$$

Multiplying by V on the right-hand side and making use of $V^t V = I$, we get

$$S U L (L^2 + \sigma^2 I)^{-1} = U L$$

$$S U L = U L \underbrace{(L^2 + \sigma^2 I)}_{\substack{\text{diagonal matrices.} \\ \text{They commute}}}$$

$$(Δ) \boxed{S U L = U (L^2 + \sigma^2 I) L}$$

We now need to consider two cases; depending on whether $l_j = 0$ or not.

- $l_j \neq 0$, then relation (Δ) becomes

$$\underline{S} \begin{pmatrix} l_1 & | & \underline{u}_1 \\ | & \dots & | \\ l_r & | & \underline{u}_r \end{pmatrix} = \begin{pmatrix} l_1 (l_1^2 + \sigma^2) & | & \underline{u}_1 \\ | & \dots & | \\ l_r (l_r^2 + \sigma^2) & | & \underline{u}_r \end{pmatrix}$$

So that the j -th column satisfies

$$\underline{S} \cancel{l_j} \underline{u}_j = \cancel{l_j} (l_j^2 + \sigma^2) \underline{u}_j$$

$$\underline{S} \underline{u}_j = (l_j^2 + \sigma^2) \underline{u}_j$$

$\Rightarrow \underline{u}_j$ is an eigenvector of \underline{S} , with eigenvalue $(l_j^2 + \sigma^2)$.

Let λ_j^s the j-th eigenvalue of \underline{S} , so that

$$l_j = (\lambda_j^s - \sigma^2)^{1/2}$$

We identified
one column
of \underline{U}

... and also
one term in the
diagonal of \underline{S} .

Related to the singular value λ_j of \underline{X}
by $\lambda_j^s = n^{-1} \lambda_j^2$ (if \underline{X} is centered)

- $\ell_j = 0$, then u_j is arbitrary.

(22)

All stationary points can thus be written:

$$(o) \quad \underline{W} = \underline{U}_r (\underline{K}_r - \sigma^2 \underline{I}_r)^{1/2} \underline{R}$$

where

- \underline{U}_r is a $(d \times r)$ matrix comprising r column eigenvectors of \underline{S} (which ones?)
- \underline{K}_r is a $(r \times r)$ diagonal matrix with elements:
 $k_j = \begin{cases} \lambda_j^s & = \text{the corresponding eigenvalue to } u_j \\ \sigma^2 & \end{cases}$
 gives $\ell_j = 0$ so that the associated column in \underline{U}_r can indeed be chosen arbitrarily.
- \underline{R} = an arbitrary orthogonal matrix = rotation
 (indeed, the matrix \underline{V} in the SVD decomposition of \underline{W} is left unspecified.)

Now the question is = which column eigenvectors of \underline{S} shall we select to get to the global maximum?

→ Substitute expression (o) for \underline{W} back into the likelihood \mathcal{L} , whose expression is given on page 19.

We need to compute $\ln |\mathbf{C}|$ and $\text{Tr}(\mathbf{C}^{-1} \mathbf{S})$, $\mathbf{C} = \sigma^2 \mathbf{I} + \mathbf{W} \mathbf{W}^t$.

Useful facts: for an $(n \times n)$ matrix A ,

$$|\mathbf{A}| = \prod_{i=1}^n \lambda_i \quad , \quad \lambda_i = \text{eigenvalues of } A$$

$$\text{Tr } A = \sum_{i=1}^n \lambda_i$$

→ First, $\ln |C| \leq 1$.

(23)

The determinant of $C = \sigma^2 I + \underline{W} \underline{W}^t$ is equal to the product of its eigenvalues.

$$C = \sigma^2 I + U_r (K_r - \sigma^2 I) U_r^t$$

Clearly, u_j are eigenvectors of C :

$$C u_j = \sigma^2 u_j + U_r \begin{pmatrix} \cdots & \lambda_k^s - \sigma^2 & \cdots \\ & 0 & \cdots \\ \uparrow & & \end{pmatrix} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \leftarrow j\text{-th position}$$

Terms on the diagonal are either $\lambda_k^s - \sigma^2$ or 0.

Let q' = number of non-zero terms.

$$= \begin{cases} \lambda_j^s u_j \\ \sigma^2 u_j \end{cases} \text{, depending on whether the } j\text{-th diagonal term is zero or not.}$$

⇒ Denote by $\lambda_1^s, \dots, \lambda_{q'}^s$ the eigenvalues corresponding to the eigenvectors 'retained' in W . It follows that

$$\ln |C| = \sum_{j=1}^{q'} \ln \lambda_j^s + (d - q') \ln \sigma^2. \quad 0 \leq q' \leq r$$

→ Next, $\text{Tr}(C^{-1}S)$.

The trace of $C^{-1}S$ is equal to the sum of its eigenvalues:

$$C^{-1} S e_j = (\sigma^2 I + U_r (K_r - \sigma^2 I) U_r^t)^{-1} S e_j = p e_j$$

pair(eigenvalue, eigenvector)
we are looking for.

$$\Rightarrow S e_j = p C e_j.$$

Take $e_j = \underline{u}_j$ = eigenvector of S .

(24)

$$\begin{aligned} \text{Then } C \underline{u}_j &= (\sigma^2 I + U_r (K_r - \sigma^2 I) U_r^t) \underline{u}_j \\ &= \sigma^2 \underline{u}_j + \begin{cases} (\lambda_j^S - \sigma^2) \underline{u}_j \\ 0 \end{cases} \\ &= \begin{cases} \lambda_j^S \underline{u}_j \\ \sigma^2 \underline{u}_j \end{cases}, \text{ depending on whether the } j\text{-th} \\ &\text{diagonal term of } K_r - \sigma^2 \text{ is} \\ &\text{zero or not.} \end{aligned}$$

↙ Eigenvalues of $C^{-1}S$ satisfy: $\lambda_j^S = \mu \times \begin{cases} \lambda_j^S \\ \sigma^2 \end{cases}$
 since $S \underline{u}_j = \lambda_j^S \underline{u}_j$

$\Rightarrow C^{-1}S$ has q' eigenvalues equal to 1 and $(d-q')$ eigenvalues equal to λ_j^S / σ^2 , where $\lambda_{q+1}^S, \dots, \lambda_d^S$ are the eigenvalues corresponding to the eigenvectors not 'retained' in W .

$$\boxed{\text{Tr}(C^{-1}S) = \frac{1}{\sigma^2} \sum_{j=q+1}^d \lambda_j^S + q'}$$

The likelihood can now be written:

$$\boxed{\mathcal{L} = -\frac{n}{2} \left\{ d \ln 2\pi + \sum_{j=1}^{q'} \ln \lambda_j^S + (d-q') \ln \sigma^2 + \frac{1}{\sigma^2} \sum_{j=q+1}^d \lambda_j^S + q' \right\}}$$

Maximizing \mathcal{L} with respect to σ^2 gives:

$$\frac{\partial \mathcal{L}}{\partial \sigma^2} = -\frac{n}{2} \left\{ \frac{d-q'}{\sigma^2} - \frac{1}{\sigma^4} \sum_{j=q+1}^d \lambda_j^S \right\} = 0$$

$$\Rightarrow \boxed{\sigma^2 = \frac{1}{d-q'} \sum_{j=q'+1}^d \lambda_j^s}$$

(25)

Note that we do not know yet at this stage which λ_j^s are selected.

Substituting this expression back into \mathcal{L} gives:

$$\mathcal{L} = -\frac{n}{2} \left\{ d \ln 2\pi + \sum_{j=1}^{q'} \ln \lambda_j^s + (d-q') \ln \left(\frac{1}{d-q'} \sum_{j=q'+1}^d \lambda_j^s \right) + (d-q') + q' \right\}$$

$$\mathcal{L} = -\frac{n}{2} \left\{ d \ln 2\pi + d + \sum_{j=1}^{q'} \ln \lambda_j^s + (d-q') \ln \left(\sum_{j=q'+1}^d \lambda_j^s / (d-q') \right) \right\}.$$

Since $\sum_{j=1}^d \lambda_j^s = \text{Tr } S$, $\sum_{j=1}^{q'} \ln \lambda_j^s = \text{Constant} - \sum_{j=q'+1}^d \ln \lambda_j^s$.

\Rightarrow Maximizing \mathcal{L} is equivalent to minimizing the quantity

$$\ln \left(\frac{1}{d-q'} \sum_{j=q'+1}^d \lambda_j^s \right) + \frac{1}{d-q'} \sum_{j=1}^{q'} \ln \lambda_j^s$$

\Updownarrow

for a fixed q' .

Minimizing

$$\ln \left(\frac{1}{d-q'} \sum_{j=q'+1}^d \lambda_j^s \right) - \frac{1}{d-q'} \sum_{j=q'+1}^d \ln \lambda_j^s =: E$$

Advantage: this quantity to minimize depends only on the discounted values -

Remark: Jensen's inequality $\ln \left(\frac{1}{n} \sum x_i \right) \geq \frac{1}{n} \sum \ln x_i$ ensures $E \geq 0$.

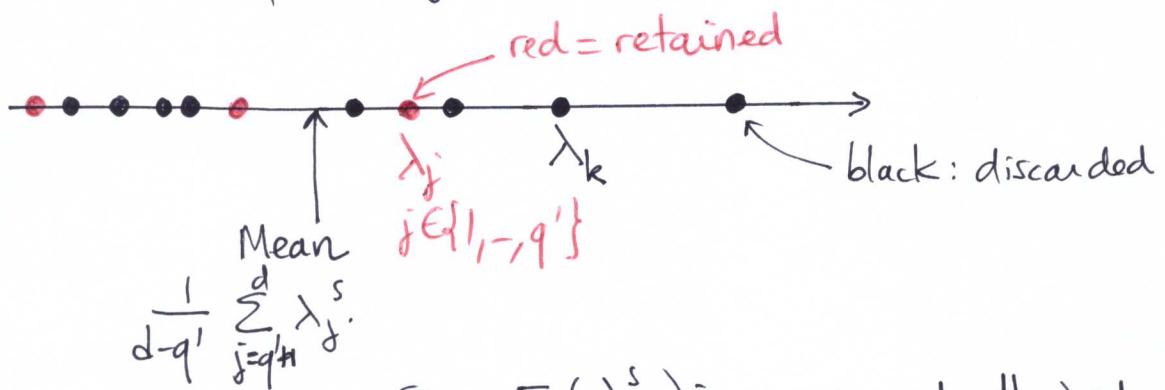
We consider the minimization of E by selecting appropriately the eigenvalues of S . (26)

Suppose that the $(d-q')$ discarded eigenvalues have been chosen arbitrarily, and consider how a single value $\lambda_k^s \in \{\lambda_{q'+1}^s, \dots, \lambda_d^s\}$ affects the value of E :

$$\frac{\partial E}{\partial \lambda_k^s} = \frac{1}{\sum_{j=q'+1}^d \lambda_j^s} - \frac{1}{(d-q') \lambda_k^s}$$

$\Rightarrow E(\lambda_k^s)$ as a function of λ_k^s has a unique minimum at $\frac{1}{d-q'} \sum_{j=q'+1}^d \lambda_j^s$ = mean of the discarded eigenvalues.

Consider the following scenario:



Since $E(\lambda_k^s)$ is convex and attained its minimum at the mean of the discarded eigenvalues, if there is a $\lambda_j^s, j \in \{1, \dots, q'\}$ between λ_k^s and the mean, the swapping λ_k^s with λ_j^s will necessarily decrease E .

\Rightarrow Unless the discarded eigenvalues $\lambda_{q'+1}^s, \dots, \lambda_d^s$ are chosen to be adjacent amongst the ordered eigenvalues of S , we can always swap two eigenvalues to decrease the value of E .

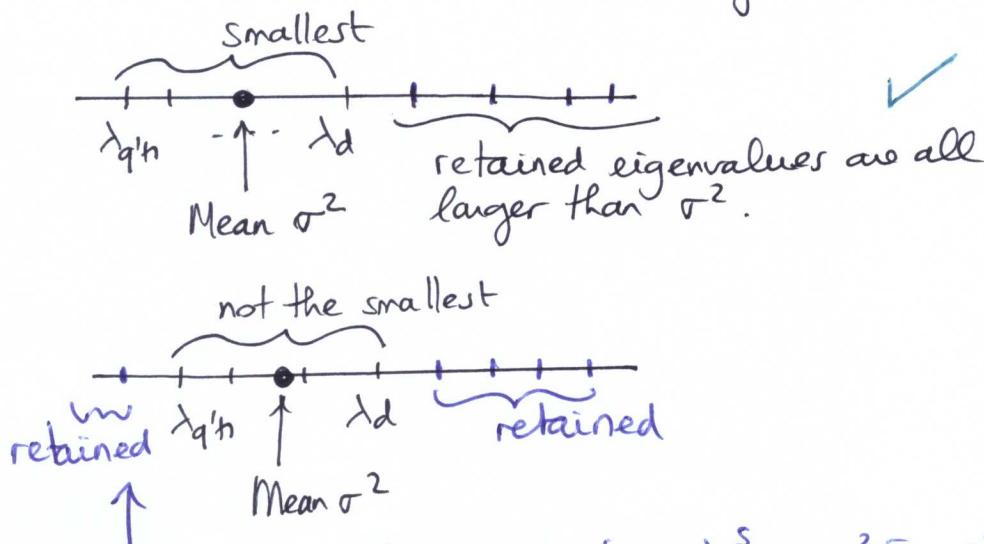
\Rightarrow The discarded eigenvalues must be adjacent.

We now conclude that they need to be the smallest:

Recall the expression of the MLE of σ^2 (top of page 25)

$$(A) \sigma^2 = \frac{1}{d-q'} \sum_{j=q'+1}^d \lambda_j^s.$$

The retained eigenvectors must have eigenvalues $\lambda_j^s > \sigma^2$ (corresponding to $\lambda_j \neq 0$). In view of (A), this condition is violated unless the discarded eigenvalues are the smallest ones:



Conclusion: E must be minimized when $\lambda_{q'+1}^s, \dots, \lambda_d^s$ are the smallest ($d-q'$) eigenvalues.

$\Rightarrow L$ is maximized when $\lambda_1^s, \dots, \lambda_{q'}^s$ are the principal eigenvalues of $\underline{\underline{S}}$.

Also, E is further minimized as a function of q' when there are the fewest terms in the sum, which corresponds to $q' = r$; i.e. when no λ_j is zero.

Rx: L is minimized when $q' = 0$, ie when $\underline{\underline{W}} = 0$.

— Summary: Consider the pPCA model ————— (28)

$$\begin{cases} \underline{x} = \underline{W}\underline{z} + \underline{\mu} + \underline{\varepsilon} \\ \underline{z} \sim \mathcal{N}(0, \mathbb{I}_r) \\ \underline{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_d) \end{cases} \quad \begin{cases} \underline{x} \in \mathbb{R}^d \\ \underline{z} \in \mathbb{R}^r \\ \underline{W} \in \mathbb{R}^{d \times r} \end{cases}$$

Let $\mathcal{D} = \{\underline{x}_1, \dots, \underline{x}_n\}$ be a sample of size n , and $\underline{S} = \frac{1}{n} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})^t$ be the sample covariance matrix, with eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0$.

The MLE of \underline{W} and σ^2 are given by

$$\begin{cases} \hat{\underline{W}}_{ML} = V_r (\Lambda_r - \sigma^2 \mathbb{I})^{1/2} R \\ \hat{\sigma}_{ML}^2 = \frac{1}{d-r} \sum_{j=r+1}^d \lambda_j, \end{cases}$$

where

- $\Lambda_r = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \ddots & \lambda_r \end{pmatrix}$ = diagonal matrix of the r largest eigenvalues of S .
- V_r = matrix of the associated eigenvectors of S .
($d \times r$)
- R = an arbitrary orthogonal rotation matrix.
($r \times r$)

Remarks.

(i) $\hat{\sigma}_{ML}^2$ = average variance 'lost' per discarded dimension.

(ii) Reduced-dimensionality transformation & reconstruction.

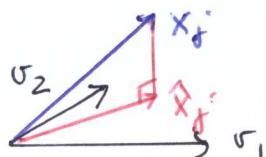
PCA

$$\underline{z}_r^j = V_r^t (\underline{x}_j - \bar{\underline{x}})$$

= reduced dim. trans.

\uparrow \uparrow
in \mathbb{R}^r in \mathbb{R}^d

$\hat{\underline{x}}_j = V_r \underline{z}_r^j + \bar{\underline{x}}$
= reconstruction
= orthogonal projection of \underline{x}_j onto the space spanned by $\{V_1, \dots, V_r\}$



p PCA

Posterior distribution is
 $z_j | x_j \sim \mathcal{N}(M^{-1}W^T(x_j - \mu), \sigma^2 M^{-1})$
 where
 $M = W^T W + \sigma^2 I$
 $(r \times r)$

\Rightarrow An observation $x_j \in \mathbb{R}^d$ in the original space is summarized by its posterior mean in the latent space:

$$E(z_j | x_j) = M^{-1} \hat{W}_{ML}^T (x_j - \bar{x})$$

$(r \times 1)$ = reduced dim. trans.

$$\hat{x}_j = \hat{W}_{ML} E(z_j | x_j) + \bar{x} \quad (29)$$

= reconstruction from the latent variable.

Q: Is \hat{x}_j = orthogonal projection of x_j onto the subspace spanned by $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$?

A: In general, NO.
 However:



\rightarrow As $\sigma^2 \rightarrow 0$, the posterior mean

$$E(z_j | x_j) \rightarrow (\hat{W}_{ML}^T \hat{W}_{ML})^{-1} \hat{W}_{ML}^T (x_j - \bar{x})$$

so that in the limit,

$$\hat{x}_j = \underbrace{\hat{W}_{ML} (\hat{W}_{ML}^T \hat{W}_{ML})^{-1} \hat{W}_{ML}^T (x_j - \bar{x}) + \bar{x}}$$

= orthogonal projection of x_j onto the column space of \hat{W}_{ML} = space spanned by $\{\mathbf{v}_1, \dots, \mathbf{v}_r\}$. We recover conventional PCA.

Note that as $\sigma^2 \rightarrow 0$ the posterior covariance is zero and the density becomes singular.

(B) EM algorithm for p PCA.

(30)

Alternatively, we can use the EM algorithm to maximize numerically the likelihood (potential advantage in the EM approach will be discussed later).

- The complete data log-likelihood is :

$$\mathcal{L}_c = \sum_{i=1}^n \ln f(\underline{x}_i, \underline{z}_i),$$

where

$$f(\underline{x}_i, \underline{z}_i) = f(\underline{x}_i | \underline{z}_i) f(\underline{z}_i)$$

$$= \frac{1}{(2\pi\sigma^2)^{d/2}} \exp \left\{ -\frac{1}{2\sigma^2} \|\underline{x}_i - W\underline{z}_i - \boldsymbol{\mu}\|^2 \right\} \frac{1}{(2\pi)^{r/2}} \exp \left\{ -\frac{1}{2} \underline{z}_i^T \underline{z}_i \right\}$$

$$\Rightarrow \mathcal{L}_c = \sum_{i=1}^n \left\{ -\frac{d}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\underline{x}_i - W\underline{z}_i - \boldsymbol{\mu}\|^2 - \frac{r}{2} \ln 2\pi - \frac{1}{2} \underline{z}_i^T \underline{z}_i \right\}$$

Omitting terms independent of the model parameters :

$$\begin{aligned} \mathcal{L}_c &= - \sum_{i=1}^n \left\{ \frac{d}{2} \ln \sigma^2 + \frac{1}{2\sigma^2} \|\underline{x}_i - W\underline{z}_i - \boldsymbol{\mu}\|^2 + \frac{1}{2} \underline{z}_i^T \underline{z}_i \right\} \\ &= - \sum_{i=1}^n \left\{ \frac{d}{2} \ln \sigma^2 + \frac{1}{2\sigma^2} \|\underline{x}_i - \boldsymbol{\mu}\|^2 + \frac{1}{2\sigma^2} \|W\underline{z}_i\|^2 \right. \\ &\quad \left. - \frac{1}{\sigma^2} \langle W\underline{z}_i, \underline{x}_i - \boldsymbol{\mu} \rangle + \frac{1}{2} \underline{z}_i^T \underline{z}_i \right\}, \end{aligned}$$

$$\begin{aligned} \text{where } \|W\underline{z}_i\|^2 &= \langle W\underline{z}_i, W\underline{z}_i \rangle = \underline{z}_i^T W^T W \underline{z}_i \\ &= \text{Tr} (\underline{z}_i^T W^T W \underline{z}_i) \\ &= \text{Tr} (W^T W \underline{z}_i \underline{z}_i^T) \end{aligned}$$

$$\cdot \underline{z}_i^T \underline{z}_i = \text{Tr} (\underline{z}_i^T \underline{z}_i) = \text{Tr} (\underline{z}_i \underline{z}_i^T)$$

\Rightarrow

$$\begin{aligned} \mathcal{L}_c &= - \sum_{i=1}^n \left\{ \frac{d}{2} \ln \sigma^2 + \frac{1}{2\sigma^2} \|\underline{x}_i - \boldsymbol{\mu}\|^2 + \frac{1}{2\sigma^2} \text{Tr} (W^T W \underline{z}_i \underline{z}_i^T) \right. \\ &\quad \left. - \frac{1}{\sigma^2} \underline{z}_i^T W^T (\underline{x}_i - \boldsymbol{\mu}) + \frac{1}{2} \text{Tr} (\underline{z}_i \underline{z}_i^T) \right\} \end{aligned}$$

- E-step.

We need to compute the expected value of the complete log-likelihood, keeping the observed variables fixed, with respect to the conditional distribution $p(\underline{z} | \underline{x})$.

$$\text{Here } p(\underline{z} | \underline{x}) \sim \mathcal{N}(M^{-1}W^t(\underline{x} - \mu), \sigma^2 M^{-1})$$

where

$$M = W^t W + \sigma^2 I.$$

(rxr)

Thus

$$\rightarrow \langle z_i \rangle \equiv E_{p(\underline{z} | \underline{x}_i)} z_i = M^{-1} W^t (\underline{x}_i - \mu)$$

Computed with the
'old' values of
the parameters

$$\rightarrow \langle z_i z_i^t \rangle \equiv E_{p(\underline{z} | \underline{x}_i)} (z_i z_i^t) = \underbrace{\sigma^2 M^{-1}}_{\text{Covariance matrix}} + \langle z_i \rangle \langle z_i \rangle^t$$

=>

$$\begin{aligned} E_{p(\underline{z} | \underline{x})} \mathcal{L}_c &= - \sum_{i=1}^n \left\{ \frac{d}{2} \ln \sigma^2 + \frac{1}{2\sigma^2} \|\underline{x}_i - \mu\|^2 + \frac{1}{2\sigma^2} \text{Tr} (W^t W \langle z_i z_i^t \rangle) \right. \\ &\quad \left. - \frac{1}{\sigma^2} \langle z_i^t \rangle W^t (\underline{x}_i - \mu) + \frac{1}{2} \text{Tr} (\langle z_i z_i^t \rangle) \right\} \end{aligned}$$

- M-step. $\langle \mathcal{L}_c \rangle$ is maximized with respect to W and σ^2 .

$$\frac{\partial \langle \mathcal{L}_c \rangle}{\partial W} = - \sum_{i=1}^n \left\{ \frac{1}{2\sigma^2} \frac{\partial}{\partial W} \text{Tr} (W^t W \langle z_i z_i^t \rangle) - \frac{1}{\sigma^2} \frac{\partial}{\partial W} \langle z_i \rangle^t W^t (\underline{x}_i - \mu) \right\}$$

Useful: $\frac{\partial}{\partial W} \text{Tr} (W^t W A) = W A^t + W A$
 $(r \times d)(d \times r)(r \times r)$

$$\frac{\partial}{\partial W} (v^t W^t v) = v v^t$$
 $(1 \times r)(r \times d)(d \times 1) (d \times 1) (1 \times r)$

$WE \mathbb{R}^{d \times r}$

$$\Rightarrow \frac{\partial \langle \mathcal{L}_c \rangle}{\partial W} = - \sum_{i=1}^n \left\{ \frac{2}{2\sigma^2} W \langle z_i z_i^t \rangle - \frac{1}{\sigma^2} (x_i - \mu) \langle z_i \rangle^t \right\} = 0 \quad (32)$$

$$\Rightarrow \sum_{i=1}^n \tilde{W} \langle z_i z_i^t \rangle = \sum_{i=1}^n (x_i - \mu) \langle z_i \rangle^t$$

$$\tilde{W} = \left(\sum_{i=1}^n (x_i - \mu) \langle z_i \rangle^t \right) \left(\sum_{i=1}^n \langle z_i z_i^t \rangle \right)^{-1}$$

Next,

$$\frac{\partial \langle \mathcal{L}_c \rangle}{\partial \sigma^2} = - \sum_{i=1}^n \left\{ \frac{d}{2\sigma^2} - \frac{1}{2\sigma^4} \|x_i - \mu\|^2 - \frac{1}{2\sigma^4} \text{Tr}(W^t W \langle z_i z_i^t \rangle) + \frac{1}{\sigma^4} \langle z_i \rangle^t W^t (x_i - \mu) \right\} = 0$$

$$\tilde{\sigma}^2 \underset{\text{nd}}{=} \sum_{i=1}^n \left\{ \|x_i - \mu\|^2 + \text{Tr}(\tilde{W}^t \tilde{W} \langle z_i z_i^t \rangle) - 2 \langle z_i \rangle^t \tilde{W}^t (x_i - \mu) \right\}.$$

$$\tilde{\sigma}^2 = \frac{1}{nd} \sum_{i=1}^n \left\{ \|x_i - \mu\|^2 + \text{Tr}(\langle z_i z_i^t \rangle \tilde{W}^t \tilde{W}) - 2 \langle z_i \rangle^t \tilde{W}^t (x_i - \mu) \right\}.$$

EM algorithm Version #1.

- Initialise W, σ^2 .

- Repeat until convergence

$$\text{E-step} \quad \langle z_i \rangle = M^{-1} W^t (x_i - \mu)$$

$$\langle z_i z_i^t \rangle = \sigma^2 M^{-1} + \langle z_i \rangle \langle z_i \rangle^t$$

$$\tilde{W} = \left(\sum_{i=1}^n (x_i - \mu) \langle z_i \rangle^t \right) \left(\sum_{i=1}^n \langle z_i z_i^t \rangle \right)^{-1}$$

$$\tilde{\sigma}^2 = \frac{1}{nd} \sum_{i=1}^n \left\{ \|x_i - \mu\|^2 + \text{Tr}(\langle z_i z_i^t \rangle \tilde{W}^t \tilde{W}) - 2 \langle z_i \rangle^t \tilde{W}^t (x_i - \mu) \right\}$$

$$W \leftarrow \tilde{W}$$

$$\tilde{\sigma} \leftarrow \tilde{\sigma}^2$$

Simplifications:

(33)

$$\begin{aligned}
 \tilde{W} &= \left(\sum_{i=1}^n (x_i - \mu) \langle z_i \rangle^t \right) \left(\sum_{i=1}^n \langle z_i z_i^t \rangle \right)^{-1} \\
 &= \left(\sum_{i=1}^n (x_i - \mu) \underbrace{(x_i - \mu)^t W M^{-1}}_{= \langle z_i \rangle^t} \right) \left(\sum_{i=1}^n \left\{ \sigma^2 M^{-1} + M^{-1} W^t (x_i - \mu) (x_i - \mu)^t W M^{-1} \right\} \right)^{-1} \\
 &= \frac{1}{n} \left\{ \sum_{i=1}^n (x_i - \mu) (x_i - \mu)^t \right\} W M^{-1} \left(\sigma^2 M^{-1} + M^{-1} W^t \left\{ \frac{1}{n} \sum (x_i - \mu) (x_i - \mu)^t \right\} W M^{-1} \right)^{-1} \\
 &= S W M^{-1} \left(\sigma^2 M^{-1} + M^{-1} W^t S W M^{-1} \right)^{-1} \\
 \tilde{W} &= S W \left(\sigma^2 I + M^{-1} W^t S W \right)^{-1} \\
 &\quad \equiv \text{Fusion of E and M steps for } \tilde{W}.
 \end{aligned}$$

$$\begin{aligned}
 \tilde{\sigma}^2 &= \frac{1}{nd} \sum_{i=1}^n \left\{ \text{Tr} (x_i - \mu) (x_i - \mu)^t + \text{Tr} (\langle z_i z_i^t \rangle \tilde{W}^t \tilde{W}) - 2 \langle z_i \rangle^t \tilde{W}^t (x_i - \mu) \right\} \\
 &\quad \uparrow \uparrow \\
 &\quad \text{These are computed with the 'dd' versions of } W \text{ and } \sigma^2. \\
 &= \frac{1}{d} \left\{ \text{Tr} S + \frac{1}{n} \sum_{i=1}^n \text{Tr} \left([\sigma^2 M^{-1} + M^{-1} W^t (x_i - \mu) (x_i - \mu)^t W M^{-1}] \tilde{W}^t \tilde{W} \right) \right. \\
 &\quad \left. - \frac{2}{n} \sum_{i=1}^n \underbrace{(x_i - \mu)^t W M^{-1} \tilde{W}^t (x_i - \mu)}_{\text{Tr} ((x_i - \mu)^t W M^{-1} \tilde{W}^t (x_i - \mu))} \right\} \\
 &= \text{Tr} (W M^{-1} \tilde{W}^t (x_i - \mu) (x_i - \mu)^t) \\
 &\Rightarrow \frac{2}{n} \sum_{i=1}^n (\dots) = 2 \text{Tr} (W M^{-1} \tilde{W}^t S) \\
 &= 2 \text{Tr} (S W M^{-1} \tilde{W}^t) \\
 &= \text{Tr} \sigma^2 M^{-1} \tilde{W}^t \tilde{W} \\
 &+ \text{Tr} M^{-1} W^t S W M^{-1} \tilde{W}^t \tilde{W} \\
 &= \text{Tr} \{ \tilde{W} (\sigma^2 M^{-1} + M^{-1} W^t S W M^{-1}) \tilde{W}^t \} \\
 &= \text{Tr} \{ \cancel{S W (\sigma^2 I + M^{-1} W^t S W)^{-1}} (\sigma^2 I + M^{-1} W^t S W) M^{-1} \tilde{W}^t \} \\
 &= \text{Tr} \{ S W M^{-1} \tilde{W}^t \}
 \end{aligned}$$

Putting everything together, the M-step for $\tilde{\sigma}^2$ simplifies to: (34)

$$\tilde{\sigma}^2 = \frac{1}{d} \left[\text{Tr} S + \text{Tr} \{ SWM^{-1}\tilde{w}^t\} - 2 \text{Tr} \{ SWM^{-1}\tilde{w}^t\} \right]$$

$$\tilde{\sigma}^2 = \frac{1}{d} \text{Tr} (S - SWM^{-1}\tilde{w}^t)$$

\equiv fusion of E and M steps for $\tilde{\sigma}^2$.

EM algorithm Version #2.

- Initialise w, σ^2
- Repeat until convergence:

$$\tilde{w} = SW(\sigma^2 I + M^{-1}w^t S w)^{-1}$$

$$\tilde{\sigma}^2 = \frac{1}{d} \text{Tr} (S - SWM^{-1}\tilde{w}^t)$$

$$w \leftarrow \tilde{w}$$

$$\sigma \leftarrow \tilde{\sigma}$$

pPCA

Remarks.

(i) Data enters into the EM algorithm only through the covariance matrix \underline{S} (compare with the analytical solution).

(ii) $\underline{S} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^t$ \rightarrow computed in $O(nd^2)$ operations.
 $\underbrace{(dx1) \quad (1xd)}_{O(d^2) \text{ operations}}$

Computing \underline{S} may be expensive & so time may be saved by not computing \underline{S} directly, as suggested by the EM algorithm.

$\hookrightarrow S$ appears through $\text{Tr} S$ and SW only.

(35)

$$\begin{aligned} \text{Tr } S &= \frac{1}{n} \sum_{i=1}^n \text{Tr } (x_i - \mu)(x_i - \mu)^t \\ &= \frac{1}{n} \sum_{i=1}^n \underbrace{(x_i - \mu)^t}_{(d \times d)} \underbrace{(x_i - \mu)}_{(d \times 1)} \rightarrow O(nd) \text{ operations.} \\ &\quad \text{d operations, repeated n times} \end{aligned}$$

$$\begin{aligned} \bullet SW &= \frac{1}{n} \left\{ \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^t \right\} w \\ &\quad \underbrace{\qquad\qquad\qquad}_{(d \times d)} \underbrace{\qquad\qquad\qquad}_{(d \times r)}, \\ &\quad O(d^2r) \end{aligned}$$

$$\Rightarrow O(d^2r) + \underbrace{O(nd^2)}_{\text{computation of } S}, \text{ to compute } SW \text{ this way.}$$

However,

$$\begin{aligned} SW &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu) \left[\underbrace{(x_i - \mu)^t w}_{(1 \times d)(d \times r)} \right] \\ &\quad \underbrace{\qquad\qquad\qquad}_{O(dr)} \\ &\quad \underbrace{\qquad\qquad\qquad}_{O(dr) + O(dr)} \\ &\quad O(ndr) \end{aligned}$$

Total of $O(ndr)$ operations \Rightarrow If $r \ll d$, then computing SW without explicitly computing S is more efficient.

→ Analytical solution requires the computation of $\underline{\underline{S}}$.

→ We might get away in the EM algorithm without direct computation of $\underline{\underline{S}}$.

\Rightarrow Iterative approach $\underset{\uparrow}{\text{MAY}}$ be more efficient.

Why 'may': because $\underline{\underline{S}}$ is computed only once, while the efficient way of calculating SW requires its evaluation at each iteration of the algorithm \Rightarrow potential efficiency gain depends on the number of iterations needed to converge.

III - KERNEL PCA

(36)

- Conventional PCA is looking for direction in the data with maximum variance. The directions are given by the eigenvectors of the covariance matrix Σ . Assuming the data is centered, the eigenvectors can be computed from the sample covariance matrix

$$S = \frac{1}{n} X^t X = \frac{1}{n} \sum_{j=1}^n x_j x_j^t.$$

(dxd) (dxn) (nxd)

Eigenvalue-eigenvector pair (λ, v) of S satisfy $Sv = \lambda v$;

$$\frac{1}{n} \sum_{i=1}^n x_i x_i^t v = \lambda v$$

$\underbrace{\quad}_{= \langle x_i, v \rangle}$

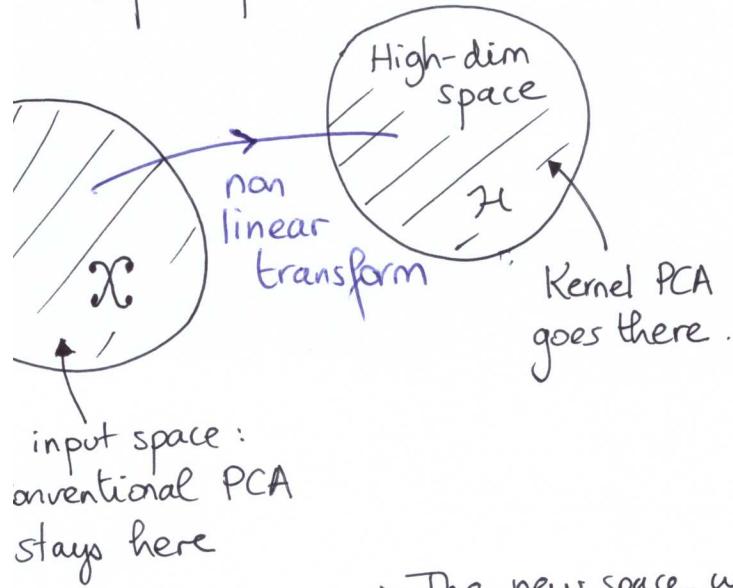
$= \text{a dot product!}$

$$\Rightarrow v = \frac{1}{n\lambda} \sum_{i=1}^n \langle x_i, v \rangle x_i$$

$$= \sum_{i=1}^n \alpha_i x_i$$

i.e. all solutions lie in the span of x_1, \dots, x_n .

- Kernel PCA is a non-linear generalization of conventional PCA: PCA is performed in a high-dimensional space \mathcal{H} (possibly of infinite dimension), which is non-linearly related to the original input space X .



The technique used is that of the 'kernel trick', that we encountered when we discussed kernel SVM: the kernel substitution allows us to take an algorithm expressed in terms of dot products $x^t x$ and generalize it by replacing the dot products with a kernel.

\Rightarrow The new space we live in is an RKHS, with reproducing kernel $K(\cdot, \cdot)$. ↑ our \mathcal{H}

\Rightarrow Original features x_i are now transformed using the feature map $\Phi: \Phi(x_i)$ ↑ not unique

$1 \leq i \leq n$

However, the feature map is not unique, and we want to express (37) the principal components in \mathcal{H} in terms of $K \Rightarrow$ we need to formulate the problem in terms of dot products in \mathcal{H} . (recall that $K(x, y) = \Phi^t(x) \Phi(y)$; the feature map can be extracted using Mercer theorem for example)

For the moment, assume that the new features $\Phi(x_i)$ are centered:

$$\sum_{i=1}^n \Phi(x_i) = 0 \quad \text{We return to this later}$$

The sample covariance matrix S in \mathcal{H} is $S = \frac{1}{n} \sum_{i=1}^n \Phi(x_i) \Phi(x_i)^t$.

Its eigenvector expansion is $S v_k = \lambda_k v_k$ $\uparrow k=1, \dots, M$.

As mentioned before, we are looking for solution to this eigenvalue-eigenvector problem without working explicitly in the feature space; i.e. working with kernels

From the definition of S , we get

$$\frac{1}{n} \sum_{i=1}^n \Phi(x_i) \boxed{\Phi(x_i)^t v_k} = \lambda_k v_k \quad (*)$$

a scalar

$\Rightarrow v_k$ can be expressed as a linear combination of the features:

$$v_k = \sum_{j=1}^n \alpha_{kj} \Phi(x_j)$$

Substituting this expression back into (*) yields

$$\frac{1}{n} \sum_{i,j} \alpha_{ki} \Phi(x_i) \underbrace{\Phi(x_i)^t \Phi(x_j)}_{K(x_i, x_j)} = \lambda_k \sum_{j=1}^n \alpha_{kj} \Phi(x_j)$$

$K(x_i, x_j)$ \leftarrow We are on the right track!

$$\sum_{i,j} \alpha_{ki} \underbrace{\Phi^t(x_m) \Phi(x_i)}_{K(x_i, x_j)} = n \lambda_k \sum_{j=1}^n \alpha_{kj} \underbrace{\Phi^t(x_k) \Phi(x_j)}$$

$\Phi^t(x_m)$
for some m

$$\Rightarrow \sum_{i,j} \alpha_{ki} K(x_m, x_i) K(x_i, x_j) = n \lambda_k \sum_j \alpha_{kj} K(x_k, x_j) \quad (38)$$

We can rewrite this equation in matrix notation, noticing that the m -th row of K^2 is the j -th column.

$$K^2_{mj} = \sum_{i=1}^n K(x_m, x_i) K(x_j, x_i)$$

We obtain $\boxed{(*)} K^2 \alpha_k = n \lambda_k K \alpha_k$, where $\alpha_k = \begin{pmatrix} \alpha_{k1} \\ \vdots \\ \alpha_{kn} \end{pmatrix} \in \mathbb{R}^n$

Solution to this problem can be found by solving the following eigenvalue problem

$$\boxed{(**)} K \alpha_k = n \lambda_k \alpha_k$$

\leftarrow We have removed a factor K on both sides.

Clearly, solutions to this problem are also solutions to $K^2 \alpha_k = n \lambda_k K \alpha_k$.

Moreover, we will show shortly that for any solution to $(**)$ having eigenvalue λ_k , we can add any multiple of an eigenvector of K having 0 eigenvalue, and obtain a solution to $(*)$ that also has eigenvalue λ_k . We also show that these modifications do not affect the principal components, that we shortly derive. Good.

As for conventional PCA, we normalize the solution so that $v_k^t v_k = 1$

$$\text{i.e. } \left(\sum_{i=1}^n \alpha_{ki} \Phi(x_i) \right)^t \left(\sum_{i=1}^n \alpha_{ki} \Phi(x_i) \right) = 1$$

$$\sum_{i,j} \alpha_{ki} \alpha_{kj} \Phi(x_i)^t \Phi(x_j)$$

$$\sum_{i,j} \alpha_{ki} \alpha_{kj} K(x_i, x_j) = \alpha_k^t K \alpha_k = 1$$

From **(**)** we get that $K\alpha_k = n \lambda_k \alpha_k$, so that

$$n \lambda_k \alpha_k^T \alpha_k = 1$$

(3g)

- Summary:
- (i) Solve the eigenvalue problem for $K \rightarrow (\mu_k, \alpha_k)$
 - (ii) Eigenvalues of S satisfy $\mu_k = n \lambda_k \quad k=1, \dots, D$
 - (iii) Normalize eigenvectors s.t. $n \lambda_k \alpha_k^T \alpha_k = 1$.
 - (iv) Compute v_k using $v_k = \sum_{j=1}^n \alpha_{kj} \Phi(x_j)$

Once the eigenvalue problem is solved, the resulting PC can also be expressed in term of K :

\Rightarrow we never do step (iv) in practice since

→ projection of a point x onto the eigenvector v_k is:

$$\Phi(x)^T v_k = \sum_{j=1}^n \alpha_{kj} \Phi(x)^T \Phi(x_j)$$

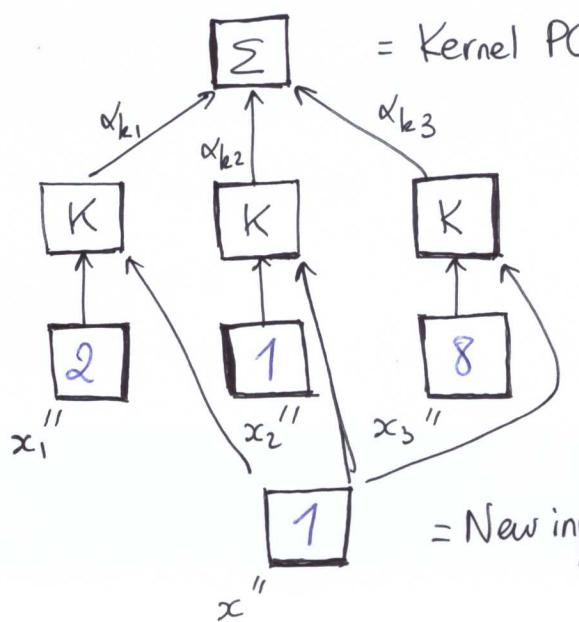
$$z_k(x) := \sum_{j=1}^n \alpha_{kj} K(x, x_j). \quad \text{k-th KERNEL PC.}$$

$z_k(x)$ = k-th principal component of a new input vector x .

\Rightarrow Replace step (iv) in the summary above by:

$$(iv) \text{ Compute the } k\text{-th PC by } z_k(x) = \sum_{j=1}^n \alpha_{kj} K(x, x_j)$$

A weighted linear combination of comparisons of the new input vector x with training features x_1, \dots, x_n .



= Comparison using kernel K

= Training samples x_1, x_2, x_3

Once the features are extracted, you can use your favorite classifier on the kernel PC

Remarks: (i) As mentioned before, the solutions of $(*)$ and $(**)$ differ only by eigenvectors of K having 0 eigenvalues that do not affect the kernel PC:

$$(*) \quad K^2 \tilde{\alpha}_k = n \lambda_k K \tilde{\alpha}_k \rightarrow \text{let } \tilde{\alpha}_k \text{ be a solution with eigenvalue } \lambda_k$$

$$(**) \quad K \alpha_k = n \lambda_k \alpha_k \rightarrow \text{let } \alpha_k \text{ be a solution with the same } \lambda_k.$$

Write $\tilde{\alpha}_k = \alpha_k + b_k$. Then, with $K b_k = 0$,

$$K \alpha_k = n \lambda_k \alpha_k \Rightarrow K \tilde{\alpha}_k = n \lambda_k \tilde{\alpha}_k$$

$$\Rightarrow K^2 \tilde{\alpha}_k = n \lambda_k K \tilde{\alpha}_k$$

$$K b_k = 0 \Rightarrow \sum_{j=1}^n b_{kj} K(x, x_j) = 0 \quad \forall x \text{ &}$$

$$\Phi^t(x) \tilde{\alpha}_k = \sum_j \tilde{\alpha}_{kj} \Phi^t(x) \Phi(x_j) = \sum_j \alpha_{kj} K(x, x_j) = z_k$$

(ii) The data is usually not centered in \mathbb{H} . We do this manually:

$$\tilde{\Phi}(x_i) = \Phi(x_i) - \frac{1}{n} \sum_{\ell=1}^n \Phi(x_\ell).$$

$$\text{Put } \tilde{K}_{ij} = \tilde{\Phi}^t(x_i) \tilde{\Phi}(x_j)$$

The eigenvalue problem becomes $\tilde{K} \tilde{\alpha}_k = n \lambda_k \tilde{\alpha}_k$.

$$\tilde{K}_{ij} = \left(\tilde{\Phi}(x_i) - \frac{1}{n} \sum_{\ell=1}^n \tilde{\Phi}(x_\ell) \right)^t \left(\tilde{\Phi}(x_j) - \frac{1}{n} \sum_{\ell=1}^n \tilde{\Phi}(x_\ell) \right)$$

$$= K_{ij} - \frac{1}{n} \sum_{\ell} \Phi(x_i)^t \Phi(x_\ell)$$

$$- \frac{1}{n} \sum_{\ell} \tilde{\Phi}(x_i)^t \tilde{\Phi}(x_\ell)$$

$$+ \frac{1}{n^2} \sum_{\ell, m} \tilde{\Phi}(x_\ell)^t \tilde{\Phi}(x_m)$$

$$\Rightarrow \tilde{K}_{ij} = K_{ij} - \frac{1}{n} \sum_{\ell} K_{i\ell} - \frac{1}{n} \sum_{\ell} K_{j\ell} + \frac{1}{n^2} \sum_{\ell, m} K_{\ell, m}$$

Symmetric!

In matrix form,

$$\tilde{K} = K - \frac{1}{n} K - K \frac{1}{n} + \frac{1}{n} K \frac{1}{n},$$

where $\frac{1}{n} = \begin{pmatrix} \frac{1}{n} & \frac{1}{n} & \dots \\ \vdots & \ddots & \ddots \end{pmatrix}$

\Rightarrow In practice, pick a kernel K , compute \tilde{K} , solve the eigenvalue problem $\tilde{K} \tilde{\alpha}_k = n \lambda_k \tilde{\alpha}_k$, and compute the k -th kernel PC.

(iii) The eigenvalue problem $K \alpha_k = n \lambda_k \alpha_k$ reduces to the conventional PCA eigenvalue problem $S v_k = \lambda_k v_k$, for the linear kernel

Feature map $\Phi(x) = x$ \rightarrow $K(x, y) = x^t y$. Indeed,

$$(K \alpha_k)_i = \sum_{j=1}^n \alpha_{kj} K(x_i, x_j) = \sum_{j=1}^n \alpha_{kj} x_i^t x_j.$$

i -th row of $K \alpha_k$

Plugging this expression back into $K \alpha_k = n \lambda_k \alpha_k$ yields:

$$\begin{aligned} (x_i x_i) &\left(\sum_{j=1}^n \alpha_{kj} x_i^t x_j \right) = n \lambda_k \alpha_{ki} \\ &\left(\sum_{j=1}^n \alpha_{kj} x_i x_i^t x_j \right) = n \lambda_k \alpha_{ki} x_i \\ \text{sum over } i &\left(\sum_{j=1}^n \left(\sum_{i=1}^n x_i x_i^t \right) \alpha_{kj} x_j \right) = n \lambda_k \left(\sum_{i=1}^n \alpha_{ki} x_i \right) \\ n S \underbrace{\left(\sum_{j=1}^n \alpha_{kj} x_j \right)}_{:= v_k} &= n \lambda_k \underbrace{\left(\sum_{i=1}^n \alpha_{ki} x_i \right)}_{:= v_k} \end{aligned}$$

$$\begin{aligned} z_k(x) &= \overline{\Phi(x)^t v_k} \\ &= x^t v_k \\ &= k\text{-th PC.} \end{aligned} \quad \rightarrow \quad \Rightarrow \quad S v_k = \lambda_k v_k, \text{ as required.}$$

(iv) A disadvantage of kernel PCA is that it involves finding the eigenvectors of a $(D \times D)$ matrix, rather than an $(d \times d)$ matrix for conventional PCA. (42)

- (v) In conventional PCA, the number of PC cannot exceed d .

In kernel PCA, the dimension of the feature space is usually much larger than d .

\Rightarrow we can find a number of non-linear PC which exceeds the dimension of \underline{x} .

- (vi) Schölkopf, Smola & Müller used kernel PCA to train an SVM classifier on hand-written digit data, and showed that they can improve the classification performance compared with linear PCA.