# Dimension Reduction

$x_i \in \mathbb{R}^P$  Our goal

$x_i \longmapsto z_i \in \mathbb{R}^q$

$q < P$

We do it by PCA

Find $q$ vectors.

$a_1, a_2, \ldots a_q \in \mathbb{R}^P$
to span $q$-dim space


Studying PCA for first time
Studying PCA for 100th time

$$\underline{A}_{P \times q} = \left( \begin{pmatrix} a_1 \end{pmatrix} \begin{pmatrix} a_2 \end{pmatrix} \cdots \begin{pmatrix} a_q \end{pmatrix} \right)$$

$$x_i \longmapsto A^T x_i \in \mathbb{R}^q$$
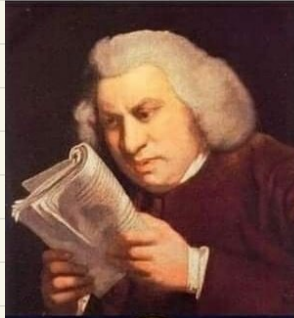
$$\|a_j\| = 1 \qquad \langle a_j, a_{j'} \rangle = 0$$

Which $q$ vectors?
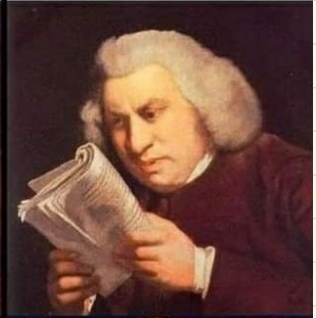
optimal $a_1, a_2, \ldots a_q$

should have least reconstruction error.

$$\frac{1}{n} \sum \|x_i - A A^T x_i\|^2$$

The problem is equivalent to $\max \ \text{Var}(a_1^T X_i)$

$\qquad s.t \ \|a_1\|=1 \qquad a_1: \text{first PC}$

$\qquad \max \ \text{Var}(a_2^T X_i)$

$\qquad St. \ \|a_2\|=1 \quad \langle a_1, a_2 \rangle =0$

$\qquad\qquad a_2: 2nd \ PC.$

# Last week:

Solve $a_1. a_2 \ldots a_q$

when $\text{Cov}(X)$ diagonal

$$\Sigma \rightarrow \begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \ddots \\ & & & \lambda_q \end{pmatrix} \qquad \lambda_1 > \lambda_2 > \cdots > \lambda_q$$

$a_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$, $a_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$, $a_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$

$A = (a_1. a_2 \cdots a_q) = \begin{pmatrix} 1 & 0 & & \\ 0 & 1 & 0 & \\ \vdots & & 1 & \\ 0 & 0 & & 1 \end{pmatrix}$ <span style="color:red">$p \times q$</span>

$\qquad\qquad \neq I_q$

Ex: $p=5, \ q=2$ $\qquad\qquad\qquad$ <span style="color:red">↙ Throw away the small variances.</span>

$z_i = A^T x_i$

$= \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 2 \\ 3 \\ 4 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

What if $\Sigma$ not diagonal?

**Theorem** (Real Spectral Thm.)

For any symmetric matrix

$\Sigma$ $(\Sigma^T = \Sigma)$, we can find another matrix $P_{p \times p}$

   s.t. $\Sigma = PDP^{-1}$, where $D$ is a $p \times p$ diagonal

      matrix.

Conclude: Every symmetric matrix can be diagonalized.

$$P^T P = I. \quad PP^T = I. \quad \Rightarrow P^{-1} = P^T$$

Let $p_1, p_2, \ldots, p_p \in \mathbb{R}^p$

$P = (p_1, p_2 \cdots p_p)$

what does $P^{-1} = P^T$ imply to $p_1, p_2, \ldots, p_p$ ?

$$P^T P = \begin{pmatrix} p_1^T \\ p_2^T \\ \vdots \\ p_p^T \end{pmatrix} (p_1, p_2, p_3 \cdots p_p) \Rightarrow \begin{pmatrix} p_1^T p_1 & & & \\ p_2^T p_1 & p_2^T p_2 & & \\ p_3^T p_1 & & p_3^T p_3 & \\ & & & \ddots \\ & & & & p_p^T p_p \end{pmatrix} = I$$

"0"

# Summary

$\mathbb{P}^{-1} = \mathbb{P}^T$  $\mathbb{P} = (p_1, p_2, p_3)$ $\rightleftarrows$ $\|p_{\vec{i}}\| = 1$

$\langle p_{\vec{i}}, p_{\vec{j}} \rangle = 0$   We call such matrix  orthonormal

breakfast + lunch = brunch; orthogonal + normal $''$

eg.

$$\Sigma = \begin{pmatrix} 34 & 12 \\ 12 & 41 \end{pmatrix}$$

We can verify:            $P$      $\times$    $D$     $\times$    $P^{-1}$

$$\begin{pmatrix} 34 & 12 \\ 12 & 41 \end{pmatrix} = \begin{pmatrix} \frac{3}{5} & \frac{-4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{pmatrix} \begin{pmatrix} 50 & 0 \\ 0 & 25 \end{pmatrix} \begin{pmatrix} \frac{3}{5} & \frac{4}{5} \\ \frac{-4}{5} & \frac{3}{5} \end{pmatrix}$$

$\downarrow$                                                   $P^T = P^{-1}$

column length = 1

We can apply R.S.T. to solve PCA

non diagonal $\Sigma$, $\Sigma = PDP^{-1}$, $D$ = diagonal.

e.g.

$\max\limits_{a_1 \in \mathbb{R}^p}$ Var $(a_1^T X)$        Note: Var $(a_1^T X)$

$\qquad \|a_1\| = 1$                 $= a_1^T \Sigma a_1$

$\qquad\qquad\qquad\qquad\qquad = a_1^T P D P^T a_1$

$$\downarrow$$

$$= \underbrace{(P^T a_1)^T}_{b_1^T} D \underbrace{P^T a_1}_{b_1}$$

Define $b_1 = P^T a_1$.

objective: max $b_1^T D b_1$

constraint

$\|a\| = 1 \iff a^T a = 1$

$a^T a = a_1^T P P^{-1} a_1$

$\qquad = (P^T a_1)^T P^T a_1$

$\qquad = b^T b_1 = 1$

$\qquad \iff \|b_1\| = 1$

max $b_1^T D b_1$

s.t. $\|b_1\| = 1$

$b_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

$\nearrow (P^{-1} = P^T)$

$b_1 = P^T a_1 \implies a_1 = P b_1$

Quiz:

$$\Sigma = \begin{pmatrix} 34 & 12 \\ 12 & 41 \end{pmatrix}$$

$$\Sigma = PDP^T$$

$$P = \begin{pmatrix} \frac{3}{5} & \frac{-4}{5} \\ \frac{4}{5} & \frac{3}{5} \end{pmatrix}, \quad D = \begin{pmatrix} 50 & 0 \\ 0 & 25 \end{pmatrix}$$   What is $a_1$? (first PC)

                                     $b_1$?

$\Rightarrow \Sigma = PDP^T$ ← Do it yourself because I went to
                                the washroom.
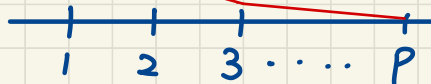
## How to choose $q$?

It depends.    $$\Sigma = \begin{pmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_p \end{pmatrix}$$   $\lambda_1 > \lambda_2 > \cdots > \lambda_p$

$$\hat{\Sigma} = \begin{pmatrix} \lambda_1 & \lambda_2 & 0 \\ \hline 0 & & \lambda_3 \end{pmatrix}$$   $q = 2.$, $\lambda_3$ will be your information loss.

$y_1 = \dfrac{\lambda_1}{\sum\limits_{j=1}^{p} \lambda_j}$ (percentage of your $\lambda$ explained)

↙ elbow point.

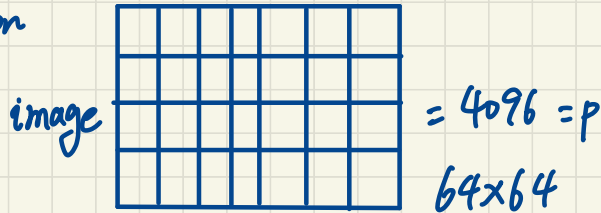$$\Sigma = \begin{pmatrix} 0.8 & & 0 \\ & 0.15 & \\ 0 & & 0.05 \end{pmatrix}$$



$$y_j = \frac{\lambda_1 + \lambda_2 + \cdots + \lambda_j}{\sum\limits_{j=1}^{p} \lambda_j}$$

# Application of PCA

## ① Data Compression
e.g. image

image [grid] $= 4096 = p$
$64 \times 64$

Data of homework

[diagram: shaded bar] ← human faces

$\Rightarrow \hat{\Sigma} \Rightarrow a_1, a_2 \dots a_q$
$q = 200$

## ②
USE PCA as "summary statistics."

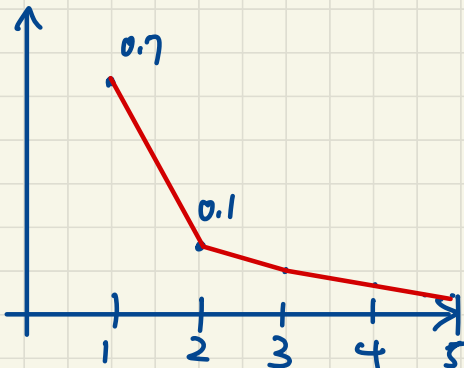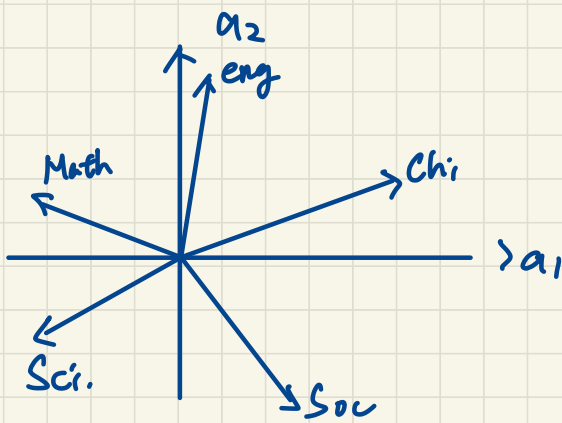e.g. $\begin{pmatrix} Chi \\ Eng \\ Math \\ Social \\ Sci \end{pmatrix}$   $a_1 = ?$   $a_2 = ?$   $a_1 = \begin{pmatrix} 0.47 \\ -0.005 \\ -0.479 \\ 0.468 \\ -0.570 \end{pmatrix}$   $a_2 = \begin{pmatrix} 0.33 \\ 0.79 \\ 0.1 \\ -0.45 \\ -0.19 \end{pmatrix}$

Either you're good at Chi, Soc
or Math.Sci

Second means that
if you're good at Eng.
or not.

[graph with points at 0.7 (x=1), 0.1 (x=2), axis labeled 1 2 3 4 5]

## ③ PCA as feature engineering.

$$X_i \to Z_i \in R^{8}$$

    regression   clustering

## ④ Factor analysis.

$$\underset{p \times 1}{X} = \underset{p \times q}{A} \times \underset{q \times 1}{Z} + \underset{p \times 1}{\varepsilon}$$

What we observed    Latent Variable   error term.

$$\begin{pmatrix} Chi \\ Eng \\ Math \\ Soc. \\ Sci. \end{pmatrix} = \begin{pmatrix} 0.2 & 0.8 \\ 0.5 & 0.5 \\ 0.9 & 0.1 \\ 0.3 & 0.7 \\ 0.8 & 0.2 \end{pmatrix}_{p \times q} \begin{pmatrix} ability \\ in\ Sci \\ ability \\ in \\ literature \end{pmatrix} + \varepsilon_i$$

$$X = AZ + \varepsilon$$

Assum1.
$$\varepsilon \sim N_p(0, \sigma^2 I)_{p \times p}$$

Assum3. $Z \perp\!\!\!\perp \varepsilon$
$$\Rightarrow X = AZ + \varepsilon$$

Assum2
$$Z \sim \underset{q}{N(0, |D|)}, \quad D = \begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \lambda_q \end{pmatrix}$$

Parameters
A. D. $\sigma^2$ ← need to be estimated
(assume $q$ known)

Generative Model

↓

distribution of data is specified.

e.g. In Psycology, Big five traits (大五人格特質)

A person's personality is built by extrovert, friendliness
consciention openess,
—ness
neuroticism

## PCA & factor analysis

need to estimate, $A, D, \sigma^2$

$$A = (a_1, a_2 \dots a_q) \; PCs$$

$$D = \begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \lambda_q \end{pmatrix}$$

when
$\sigma^2 \to 0$

$\mathcal{E} \sim N(0, \sigma^2 I)$

$Z \sim N(0, \underline{D}) \leftarrow$ Only this version would lead t. PCA.

Why factor Model?

① better interpretation.

② statistical inference. on $a_1 \, a_2 \dots a_q$

Hypothesis testing

$H_0: \dfrac{\lambda_1 + \lambda_2}{\sum\limits_{0.1} \lambda_j} > 0.8$   (Benefit of Generative Model)