

# Statistical Inference and Multivariate Analysis (MA324)

## LECTURE SLIDES Lecture 34

### Principal Component Analysis



Indian Institute of Technology Guwahati

Jan-May 2023

# Principal Components Obtained from Standardized Variables

$$Z_i = \frac{X_i - \mu_i}{\sqrt{\sigma_{ii}}}, i = 1, 2, \dots, p$$

$$\Leftrightarrow \underline{\tilde{Z}} = \begin{pmatrix} Z_1 \\ \vdots \\ Z_p \end{pmatrix} = V^{-1}(\underline{\tilde{X}} - \underline{\tilde{\mu}}), \text{ where } V = \text{diag}(\sqrt{\sigma_{11}}, \sqrt{\sigma_{22}}, \dots, \sqrt{\sigma_{pp}})$$

- $$Cov(\tilde{Z}) = \rho_{p \times p} = \begin{pmatrix} \rho_{11}(=1) & . & . & . & \rho_{1p} \\ . & & & & . \\ . & & & & . \\ . & & & & . \\ \rho_{p1} & . & . & . & \rho_{pp}(=1) \end{pmatrix}$$

- $$\sum_{i=1}^p Var(Y_i) = p$$

- Proportion of standardized population variance due to  $k^{th}$  principal component =  $\frac{\lambda_k}{p}$ ,  $k = 1, 2, \dots, p$ , where  $\lambda_k$ 's are the eigen values pf  $\rho_{p \times p}$ .

# Summarizing Sample Variation using Principal Components

- Let  $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n$  denote  $n$  independent draws from a  $p$ -dimensional population with mean vector  $\mu$  and variance-covariance matrix  $\Sigma$ .
- The sample mean is  $\bar{\tilde{x}}$ , the sample var-cov matrix  $S$  and the sample correlation matrix  $R$ .
- 1<sup>st</sup> sample principal component = Linear combination  $a'_1 \tilde{x}_j$  that maximizes the sample variance of  $a'_1 \tilde{x}_j$  subject to  $a'_1 a_1 = 1$ .

- $i^{th}$  sample principal component = Linear combination  $\underline{a}'_i \underline{x}_j$  that maximizes the sample variance of  $\underline{a}'_i \underline{x}_j$  subject to  $\underline{a}'_i \underline{a}_i = 1$  and zero sample correlation between  $\underline{a}'_i \underline{x}_j$  and  $\underline{a}'_k \underline{x}_j$  for all  $k < i$ .
- Note that sample variance of  $\underline{a}'_1 \underline{x}_j$  is  $\underline{a}'_1 S \underline{a}_1$ . The sample covariance between  $\underline{a}'_i \underline{x}_j$  and  $\underline{a}'_k \underline{x}_j$  is  $\underline{a}'_i S \underline{a}_k$ .

# Theorem

Let  $S$  be  $p \times p$  sample var-cov matrix with eigen-value-eigen-vector pair  $(\hat{\lambda}_1, \hat{e}_1), \dots, (\hat{\lambda}_p, \hat{e}_p)$  with  $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$ . Then  $i^{th}$  sample principal component is given by,

$$\hat{y}_i = \hat{e}_i' \underline{x} = \hat{e}_{i1}x_1 + \dots + \hat{e}_{ip}x_p, \quad i = 1, 2, \dots, p$$

Also, sample variance of  $\hat{y}_i = \hat{\lambda}_i$ ,  $i = 1, 2, \dots, p$  and sample covariance  $\hat{y}_i$  and  $\hat{y}_k$  is 0, for  $i \neq k$ .

# How many components to retain?

- No definite answer
- Scree plot

# Sample principal components based on standardized data

- $S$  matrix needs to be replaced by  $R$ .



# Application of PCA...

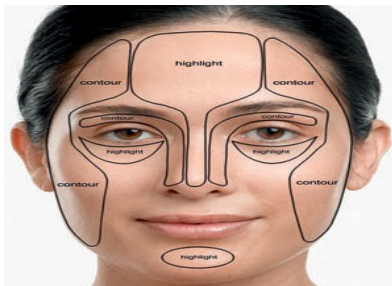


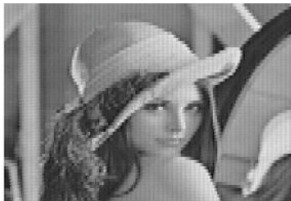
Figure 1. The regions that are either highlighted or contoured when used in PCA.[1]

- Reference: Alkandari, A., & Aljaber, S. J. (2015, April). Principle Component Analysis algorithm (PCA) for image recognition. In 2015 Second International Conference on Computing Technology and Information Management (ICCTIM) (pp. 76-80). IEEE.

# Application of PCA...



a)



b)



c)



d)

- Original and compressed image with two, five, and eight principal components. (a) Original image. (b) Compression with two components. (c) Compression with five components. (d) Compression with eight components.
- Reference: Hernandez, W., Mendez, A., & Göksel, T. (2018). Application of principal component analysis to image compression. Statistics-Growing Data Sets and Growing Demand for Statistics.

# Application of PCA...

Did you spot the bias in AI?