NTHU STAT 5191, 2010



NTHU STAT 5191, 2010

Lecture Notes

| • $\begin{pmatrix} Proportion of total \\ population variance \\ due to kth principal \\ component \end{pmatrix} = \frac{\lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_p} k = 1, 2, \dots, p$ If most (for instance, 80 to 90%) of the total population variance, for large p , can be attributed to the first one, two, or three components, then these components can $\lambda_k = HW2$ "replace" the original p variables without much loss of information. $\lambda_3 = HW3$ Each component of the coefficient vector $\mathbf{e}'_i = [e_{i1}, \dots, e_{i_k}, \dots, e_{i_p}]$ also merits $y_{i=0.63}$ (humpspection. The magnitude of e_{i_k} measures the importance of the kth variable to the $\mathbf{e} \circ \mathbf{y} = \mathbf{e} \mathbf{y} \mathbf{x}$ are the principal components, irrespective of the other variables. In particular, e_{i_k} is pro- $\mathbf{e} \circ \mathbf{y} = 0.52$ (huppertunction to the correlation coefficient between Y_i and X_k . $y_{2=0.51}(\mathbf{a}, \mathbf{h})$ such that $\mathbf{e} = \mathbf{e}_1 \mathbf{x}, Y_2 = \mathbf{e}_2 \mathbf{x}, \dots, Y_p = \mathbf{e}_p' \mathbf{x}$ are the principal components $\mathbf{e} = \mathbf{e} \circ \mathbf{x} + \mathbf{e} \circ \mathbf{x} + \mathbf{e} \circ \mathbf{x} + \mathbf{e} \circ \mathbf{x} + \mathbf{x} + \mathbf{e} \circ \mathbf{x} + \mathbf{e} + e$ |
|--|
| • the <i>i</i> th principal component scores $PX = [\mathbf{e}_{1} \cdots \mathbf{e}_{p}] = PX - X^{2} P'$ • $Y_{i} = \mathbf{e}_{i}' \mathbf{X}$ • $Y_{i} = \mathbf{e}_{i}' (\mathbf{X} - \boldsymbol{\mu})$ |
| Principal Components Obtained from Standardized Variables Principal Components Obtained from Standardized Variables Q: What if one variables is measured in the millions whereas the others are measured in tens? or one variable has much larger scales than other variable? ⇒ The 1 st PC will essentially be just that variable Principal components may also be obtained for the standardized variables Z ₁ = (X ₁ - µ ₁) Z ₂ = (X ₂ - µ ₂) √σ ₂₂ ,, Z _p = (X _p - µ _p) (o ⁶ → 0 ⁶ → 0 ⁶) Cov(Z) = (V ^{1/2}) ⁻¹ Σ(V ^{1/2}) ⁻¹ = ρ Z ₂ • Result 8.4. The <i>i</i> th principal component of the standardized variables Z' = [Z ₁ , Z ₂ ,, Z _p] with Cov(Z) = ρ, is given by Y _i = e ⁱ (Z) = e ⁱ (V ^{1/2}) ⁻¹ (X - µ), <i>i</i> = 1, 2,, p Moreover, P(C scores) and p(Var(Y _i) = $\sum_{i=1}^{p} Var(Z_i) = p$ ⇒ all variables equally important p(Note, i, e ₁), (λ ₂ , e ₂),, (λ _p , e _p) are the eigenvalue-eigenvector pairs for ρ, with $\lambda_1 \ge \lambda_2 \ge \cdots \ge \lambda_p \ge 0$. |

made by S.-W. Cheng (NTHU, Taiwan)



made by S.-W. Cheng (NTHU, Taiwan)

NTHU STAT 5191, 2010



p. 4-10 2. Let $\mathbf{E}_i = \lambda_i \sum_{k=1}^p \frac{\lambda_k}{(\lambda_k - \lambda_i)^2} \, \mathbf{e}_k \mathbf{e}'_k$ then $\sqrt{n} (\hat{\mathbf{e}}_i - \mathbf{e}_i)$ is approximately $N_p(\mathbf{0}, \mathbf{E}_i)$. 3. Each $\hat{\lambda}_i$ is distributed independently of the elements of the associated $\hat{\mathbf{e}}_i$. - A large sample $100(1 - \alpha)$ % confidence interval for λ_i is thus provided by $\frac{\hat{\lambda}_{i}}{(1+z(\alpha/2)\sqrt{2/n})} \leq \lambda_{i} \leq \frac{\hat{\lambda}_{i}}{(1-z(\alpha/2)\sqrt{2/n})} \xrightarrow{\lambda_{i} - \lambda_{i} \wedge N(o, \frac{2}{n}\lambda_{i})}{\frac{\hat{\lambda}_{i} - \lambda_{i} \wedge N(o, \frac{2}{n}\lambda_{i})}$ • Result 2 implies that the $\hat{\mathbf{e}}_i$'s are normally distributed about the corresponding \mathbf{e}_i 's for large samples. The elements of each $\hat{\mathbf{e}}_i$ are correlated, and the correlation depends to a large extent on the separation of the eigenvalues $\lambda_1, \lambda_2, \ldots, \lambda_p$ li≈p 22 - 2p Small. Some important issues > Analyses of principal components are more of a means to an end rather than an end in themselves, because they frequently serve as intermediate steps in much larger investigation. Example: regression, factor analysis, cluster analysis, > PCs with variance almost zero (i.e., $\hat{\lambda}_i \approx 0$) $\hat{y}_i = \hat{\mathbf{e}}_i \mathbf{x} = \hat{e}_{i1} x_1 + \hat{e}_{i2} x_2 + \dots + \hat{e}_{ip} x_p \approx c, c:$ a constant \Rightarrow an unusually small value for eigenvalue can indicate a linear dependency in the data set p. 4-11 \triangleright equal eigenvalues (i.e., $\lambda_i = \lambda_{i+1} = \cdots = \lambda_i$) population principal component sample principal component > selecting a subset of variables $X_1, ..., X_p$ • PCA does not quite reduce the data since PCs are linear combinations of all variables X_1, \ldots, X_p so all variables are in some sense still needed. The technique, however, can help us discard some variables by keeping only those with the highest factor loading connection to the singular value decomposition (SVD) • Let us apply a SVD to the matrix $\mathbf{X} - \mathbf{1}\mathbf{\bar{x}}'$ $(\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}')_{n \times p} = U_{n \times p} L_{p \times p} V'_{n \times p}$ where $U'U = I_m$, $V'V = I_m$, and L is a diagonal matrix. Then, $(n-1)\mathbf{S} = (\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}')'(\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}') = VLU'ULV' = V(LL)V' = VL^2V'$ $\mathbf{S} = V\left(\frac{1}{n-1}L^2\right)V' = V\left(\frac{1}{\sqrt{n-1}}L\right)^2 V' \equiv VL^{*2}V'$ which means \bullet columns of V: eigenvectors of S • $\frac{l_i^2}{n-1}$: eigenvalues of **S**, where l_i : singular value • $\mathbf{Y} = (\mathbf{X} - \mathbf{1}\mathbf{\bar{x}}')V = UL$: PC scores