

(فصل دوازدهم را ملاحظه کنید) باشند. علاوه بر این مؤلفه های اصلی (مقیاس بندی) شده یک «عامل سازی» ماتریس کوواریانس برای الگوی تحلیل عاملی مورد نظر در فصل نهم است.

۲-۸ مؤلفه های اصلی جامعه

مؤلفه های اصلی از نظر جبری ترکیبات خطی ویژه p متغیر تصادفی X_1, X_2, \dots, X_p است. این ترکیبات خطی از نظر هندسی انتخاب یک دستگاه مختصات جدید را نشان می دهد که از دوران دستگاه اولیه با X_1, X_2, \dots, X_p به عنوان محورهای مختصات به دست می آید. محورهای جدید جهتها را با بیشترین تغییرپذیری نشان می دهد و بیان ساده تر و ممسک تری از ساختمان کوواریانس را فراهم می کند.

چنان که ملاحظه خواهیم نمود مؤلفه های اصلی تنها به ماتریس کوواریانس Σ (یا ماتریس همبستگی ρ) X_1, X_2, \dots, X_p مربوط می شود. برای بسط آنها فرض نرمال چندمتغیری لازم نیست. از سوی دیگر مؤلفه های اصلی که برای جامعه های نرمال چندمتغیری به دست می آید تعابیر مفیدی بر حسب بیضویهای چگالی ثابت دارد. علاوه بر این وقتی جامعه نرمال چندمتغیری است (بخش ۵-۸ را ملاحظه نمایید) استنباطهایی را از مؤلفه های نمونه می توان به عمل آورد.

فرض کنید بردار تصادفی $X' = [X_1, X_2, \dots, X_p]$ دارای ماتریس کوواریانس Σ با مقادیر ویژه $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ است.

ترکیبات خطی زیر را در نظر می گیریم:

$$\begin{aligned} Y_1 &= \ell'_1 X = \ell_{11}X_1 + \ell_{21}X_2 + \dots + \ell_{p1}X_p \\ Y_2 &= \ell'_2 X = \ell_{12}X_1 + \ell_{22}X_2 + \dots + \ell_{p2}X_p \\ &\vdots \\ Y_p &= \ell'_p X = \ell_{1p}X_1 + \ell_{2p}X_2 + \dots + \ell_{pp}X_p \end{aligned} \quad (1-8)$$

در این صورت با استفاده از (۲-۴۵)، داریم:

$$\text{Var}(Y_i) = \ell'_i \Sigma \ell_i \quad i = 1, 2, \dots, p \quad (2-8)$$

$$\text{Cov}(Y_i, Y_k) = \ell'_i \Sigma \ell_k \quad i, k = 1, 2, \dots, p \quad (3-8)$$

مؤلفه های اصلی آن ترکیبات خطی ناهمبسته Y_1, Y_2, \dots, Y_p هستند که واریانسهای آنها در (۲-۸) تا جایی که ممکن است، بزرگ باشد.

اولین مؤلفه اصلی یک ترکیب خطی با واریانس ماکزیمم است. یعنی $\text{Var}(Y_1) = \ell_1' \Sigma \ell_1$ را ماکزیمم می کند. واضح است که $\text{Var}(Y_1) = \ell_1' \Sigma \ell_1$ را می توان با ضرب کردن هر ℓ_1 در یک ثابت ماکزیمم کرد. برای از بین بردن این ابهام بهتر است بردارهای ضرایب با طول واحد را مورد توجه قرار دهیم. بنابراین تعریف می کنیم،

ترکیب خطی $l_1' X$ که $\text{Var}(l_1' X) = 1$ را با توجه به $l_1' l_1 = 1$ ماکزیمم کند = اولین مؤلفه اصلی.
 ترکیب خطی $l_2' X$ که $\text{Var}(l_2' X) = 1$ را با توجه به $l_2' l_2 = 1$ و $\text{Cov}(l_1' X, l_2' X) = 0$ = دومین مؤلفه اصلی ماکزیمم کند.

در مرحله i ام

ترکیب خطی $l_i' X$ که $\text{Var}(l_i' X) = 1$ را با توجه به $\ell_i' \ell_i = 1$ و مؤلفه اصلی i ام
 $\text{Cov}(l_i' X, l_k' X) = 0$ ، $k < i$ ماکزیمم کند.

نتیجه ۸-۱. فرض کنید Σ ماتریس کوواریانس بردار تصادفی $X' = [X_1, X_2, \dots, X_p]$ باشد. فرض کنید Σ دارای زوج مقدار ویژه - بردار ویژه $(\lambda_1, e_1), (\lambda_2, e_2), \dots, (\lambda_p, e_p)$ باشد که $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ ، مؤلفه اصلی i ام با

$$Y_i = e_i' X = e_{1i} X_1 + e_{2i} X_2 + \dots + e_{pi} X_p, \quad i = 1, 2, \dots, p \quad (4-8)$$

داده می شود. با این انتخابها، داریم:

$$\text{Var}(Y_i) = e_i' \Sigma e_i = \lambda_i \quad i = 1, 2, \dots, p$$

$$\text{Cov}(Y_i, Y_k) = e_i' \Sigma e_k = 0 \quad i \neq k \quad (5-8)$$

در صورتی که بعضی از λ_i ها برابر باشند، انتخابهای بردارهای ضرایب مربوط e_i و در نتیجه Y_i یکتا نخواهند بود.

اثبات. از (۲-۵۱) با $B = \Sigma$ می دانیم که

$$\max_{\ell \neq 0} \frac{\ell' \Sigma \ell}{\ell' \ell} = \lambda_1 \quad (\ell = e_1 \text{ وقتی به دست می آید که})$$

اما چون بردارهای ویژه نرمال شده اند، لذا $e_i' e_i = 1$. از این رو

$$\max_{\ell \neq 0} \frac{\ell' \Sigma \ell}{\ell' \ell} = \lambda_1 = \frac{e_1' \Sigma e_1}{e_1' e_1} = e_1' \Sigma e_1 = \text{Var}(Y_1)$$

به طور مشابه با استفاده از (۲-۵۲)، داریم:

$$\max_{\ell \perp e_1, e_2, \dots, e_k} \frac{\ell' \Sigma \ell}{\ell' \ell} = \lambda_{k+1} \quad k = 1, 2, \dots, p-1$$

برای انتخاب $l = e_{k+1}$ که $e'_{k+1} e_k = 0$ ، $k = 1, 2, \dots, p-1$ ،

$$e'_{k+1} \Sigma e_{k+1} / e'_{k+1} e_{k+1} = e'_{k+1} \Sigma e_{k+1} = \text{Var}(Y_{k+1})$$

اما $\lambda_{k+1} e'_{k+1} e_{k+1} = \lambda_{k+1} (\Sigma e_{k+1}) = e'_{k+1}$ بنابراین $\text{Var}(Y_{k+1}) = \lambda_{k+1}$. حال باید عمود بودن e_k بر e_i (یعنی $e'_i e_k = 0$ ، $i \neq k$) که $\text{Cov}(Y_i, Y_k) = 0$ را می دهد اثبات کنیم . اکنون بردارهای ویژه Σ در صورتی که مقادیر ویژه $\lambda_1, \lambda_2, \dots, \lambda_p$ متمایز باشند ، متعامدند . اگر تمام مقادیر ویژه متمایز نباشند ، در آن صورت بردارهای ویژه متناظر با مقادیر ویژه مشترک را می توان متعامد انتخاب کرد . بنابراین برای هر دو بردار ویژه e_i و e_k ، $e'_i e_k = 0$ ، $i \neq k$. چون $\Sigma e_k = \lambda_k e_k$ لذا اگر آن را از سمت چپ در e'_i ضرب کنیم ، برای هر $i \neq k$ ، داریم :

$$\text{Cov}(Y_i, Y_k) = e'_i \Sigma e_k = e'_i \lambda_k e_k = \lambda_k e'_i e_k = 0$$

و بدین ترتیب اثبات کامل می شود .

از نتیجه (۸-۱) معلوم می شود که مؤلفه های اصلی ناهمبسته بوده و واریانس آنها برابر مقادیر ویژه Σ است .

نتیجه ۸-۲ . فرض کنید $X' = (X_1, X_2, \dots, X_p)$ دارای ماتریس کوواریانس Σ با زوج مقدار ویژه-بردار ویژه $(\lambda_1, e_1), (\lambda_2, e_2), \dots, (\lambda_p, e_p)$ ، $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ ، است . فرض کنید $Y_1 = e'_1 X, Y_2 = e'_2 X, \dots, Y_p = e'_p X$ مؤلفه های اصلی باشند ، آن گاه

$$\sigma_{11} + \sigma_{22} + \dots + \sigma_{pp} = \sum_{i=1}^p \text{Var}(X_i) = \lambda_1 + \lambda_2 + \dots + \lambda_p = \sum_{i=1}^p \text{Var}(Y_i)$$

اثبات . از تعریف ۲ الف . ۲۸ داریم ، $\sigma_{11} + \sigma_{22} + \dots + \sigma_{pp} = \text{tr}(\Sigma)$ ، از (۲-۲۰) با $\Lambda = \Sigma$ می توان نوشت $\Sigma = P \Lambda P'$ که Λ ماتریس قطری مقادیر ویژه و $P = [e_1, e_2, \dots, e_p]$ چنان است که $PP' = P'P = I$. با استفاده از نتیجه ۲ الف . ۱۲ (ج) ، داریم :

$$\text{tr}(\Sigma) = \text{tr}(P \Lambda P') = \text{tr}(\Lambda P'P) = \text{tr}(\Lambda) = \lambda_1 + \lambda_2 + \dots + \lambda_p$$

از این رو

$$\sum_{i=1}^p \text{Var}(X_i) = \text{tr}(\Sigma) = \text{tr}(\Lambda) = \sum_{i=1}^p \text{Var}(Y_i)$$

نتیجه ۸-۲ بیان می کند که

$$\sigma_{11} + \sigma_{22} + \dots + \sigma_{pp} = \lambda_1 + \lambda_2 + \dots + \lambda_p \quad (۸-۶)$$

کل واریانس جامعه

و در نتیجه نسبت واریانس کل مربوط به (بیان شده با) مؤلفه اصلی k ام ، عبارت است از :

$$\left(\begin{array}{c} \text{سهم کل واریانس} \\ \text{جامعه مربوط} \\ \text{به مؤلفه اصلی} \\ \text{ام } k \end{array} \right) = \frac{\lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_p} \quad k = 1, 2, \dots, p \quad (7-8)$$

اگر برای p بزرگ بیشتر واریانس کل جامعه (مثلاً ۸۰ تا ۹۰ درصد) آن را بتوان به سه مؤلفه اول نسبت داد در آن صورت این مؤلفه ها را بدون این که اطلاعات زیادی را از دست دهیم ، می توان «جایگزین» p متغیر اولیه کرد .

هر مؤلفه بردار ضرایب $e'_i = [e_{1i}, \dots, e_{ki}, \dots, e_{pi}]$ نیز ارزش بررسی را دارد . مقدار e_{ki} اهمیت متغیر k ام را در مؤلفه اصلی i ام صرف نظر از متغیرهای دیگر اندازه می گیرد . به ویژه e_{ki} با ضریب همبستگی Y_i و X_k متناسب است .

نتیجه ۳-۸ . اگر $Y_1 = e'_1 X$, $Y_2 = e'_2 X$, ..., $Y_p = e'_p X$ ، مؤلفه های اصلی به دست آمده از ماتریس کوواریانس Σ باشد آن گاه

$$\rho_{Y_i, X_k} = \frac{e_{ki} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}} \quad i, k = 1, 2, \dots, p \quad (8-8)$$

ضرایب همبستگی بین مؤلفه های Y_i و متغیرهای X_k است . در این جا $(\lambda_1, e_1), (\lambda_2, e_2), \dots, (\lambda_p, e_p)$ زوجهای مقدار ویژه - بردار ویژه Σ هستند .

اثبات . قرار می دهیم $l'_k = [0, \dots, 0, 1, 0, \dots, 0]$ به طوری که $X_k = l'_k X$ و بنا به (۲-۴۵) ،

$\text{Cov}(X_k, Y_i) = \text{Cov}(l'_k X, e'_i X) = l'_k \Sigma e_i$ چون $\Sigma e_i = \lambda_i e_i$ لذا $\text{Cov}(X_k, Y_i) = \lambda_i e_{ki}$ در این صورت $\text{Var}(Y_i) = \lambda_i$ (۵-۸) را ملاحظه نمایید) و $\text{Var}(X_k) = \sigma_{kk}$ ،

$$\rho_{Y_i, X_k} = \frac{\text{Cov}(Y_i, X_k)}{\sqrt{\text{Var}(Y_i)} \sqrt{\text{Var}(X_k)}} = \frac{\lambda_i e_{ki}}{\sqrt{\lambda_i} \sqrt{\sigma_{kk}}} = \frac{e_{ki} \sqrt{\lambda_i}}{\sqrt{\sigma_{kk}}} \quad i, k = 1, 2, \dots, p$$

را حاصل می کند .

مثال فرضی زیر مفاد نتایج ۱-۸ ، ۲-۸ و ۳-۸ را تشریح می کند .

مثال ۱-۸

فرض کنید متغیرهای تصادفی X_1 ، X_2 و X_3 دارای ماتریس کوواریانس زیر باشند :

$$\Sigma = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

زوجهای مقدار ویژه- بردار ویژه ، عبارتند از :

$$\lambda_1 = 5.83, \quad \mathbf{e}'_1 = [.383, -.924, 0]$$

$$\lambda_2 = 2.00, \quad \mathbf{e}'_2 = [0, 0, 1]$$

$$\lambda_3 = 0.17, \quad \mathbf{e}'_3 = [.924, .383, 0]$$

بنابراین مؤلفه های اصلی به صورت زیر خواهند بود :

$$Y_1 = \mathbf{e}'_1 \mathbf{X} = .383X_1 - .924X_2$$

$$Y_2 = \mathbf{e}'_2 \mathbf{X} = X_3$$

$$Y_3 = \mathbf{e}'_3 \mathbf{X} = .924X_1 + .383X_2$$

متغیر X_3 یکی از مؤلفه های اصلی است ، زیرا با دو متغیر دیگر ناهمبسته است .
معادله (۸-۵) را با توجه به اصول اولیه می توان تشریح کرد . برای مثال

$$\begin{aligned} \text{Var}(Y_1) &= \text{Var}(.383X_1 - .924X_2) \\ &= (.383)^2 \text{Var}(X_1) + (-.924)^2 \text{Var}(X_2) + 2(.383)(-.924) \text{Cov}(X_1, X_2) \\ &= .147(1) + .854(5) - .708(-2) \\ &= 5.83 = \lambda_1 \end{aligned}$$

$$\begin{aligned} \text{Cov}(Y_1, Y_2) &= \text{Cov}(.383X_1 - .924X_2, X_3) \\ &= .383 \text{Cov}(X_1, X_3) - .924 \text{Cov}(X_2, X_3) \\ &= .383(0) - .924(0) = 0 \end{aligned}$$

همچنین به آسانی می توان دید

$$\sigma_{11} + \sigma_{22} + \sigma_{33} = 1 + 5 + 2 = \lambda_1 + \lambda_2 + \lambda_3 = 5.83 + 2.00 + .17$$

که معادله (۸-۶) را برای این مثال تأیید می کند . سهم کل واریانس که به وسیله اولین مؤلفه اصلی به حساب می آید $\frac{\lambda_1}{(\lambda_1 + \lambda_2 + \lambda_3)} = \frac{5.83}{8} = .73$ است . با ارائه این عمل سهم واریانس جامعه که با دو مؤلفه اول به حساب می آید $\frac{(5.83 + 2)}{8} = .98$ است . در این حالت بدون این که کمترین اطلاعی را از دست دهیم می توانیم مؤلفه های Y_2 و Y_3 را جانشین سه متغیر اولیه کنیم .
بالاخره با استفاده از (۸-۸) ، داریم :

$$\rho_{Y_1, X_1} = \frac{e_{11} \sqrt{\lambda_1}}{\sqrt{\sigma_{11}}} = \frac{.383 \sqrt{5.83}}{\sqrt{1}} = .925$$

$$\rho_{Y_1, X_2} = \frac{e_{21} \sqrt{\lambda_1}}{\sqrt{\sigma_{22}}} = \frac{-.924 \sqrt{5.83}}{\sqrt{5}} = -.998$$

نتیجه می گیریم که اهمیت هر یک از متغیرهای X_1 و X_2 در اولین مؤلفه اصلی تقریباً یکسان است. همچنین

$$\rho_{Y_2, X_1} = \rho_{Y_2, X_2} = 0 \quad \text{و} \quad \rho_{Y_2, X_3} = \frac{\sqrt{\lambda_2}}{\sqrt{\sigma_{33}}} = \frac{\sqrt{2}}{\sqrt{2}} = 1$$

(که باید چنین باشد)

چون مؤلفه سوم بی اهمیت است، لذا از بقیه همبستگیها می توان صرف نظر کرد.

بررسی مؤلفه های اصلی به دست آمده از متغیرهای تصادفی نرمال چندمتغیری مفید است.

فرض کنید X دارای توزیع $N_p(\mu, \Sigma)$ باشد. از (۷-۴) می دانیم که بیضویهای متمرکز شده در μ با چگالی ثابت

$$(x - \mu)' \Sigma^{-1} (x - \mu) = c^2$$

دارای محورهای $\pm c \sqrt{\lambda_i} e_i$ ، $i = 1, 2, \dots, p$ بوده که (λ_i, e_i) زوجهای مقدار ویژه-بردار ویژه Σ است. نقطه ای که روی محور i ام بیضوی قرار دارد دارای مختصات متناسب با $e'_i = [e_{i1}, e_{i2}, \dots, e_{ip}]$ در دستگاه مختصات با مبدأ μ و محورهای x_1, x_2, \dots, x_p است. در بحث بعدی^۱ بهتر است $\mu = 0$ قرار دهیم.

با توجه به بحث بخش ۲-۳ که $A = \Sigma^{-1}$ ، می توان نوشت:

$$c^2 = x' \Sigma^{-1} x = \frac{1}{\lambda_1} (e'_1 x)^2 + \frac{1}{\lambda_2} (e'_2 x)^2 + \dots + \frac{1}{\lambda_p} (e'_p x)^2$$

که در آن $e'_1 x$ ، $e'_2 x$ ، \dots ، $e'_p x$ را به عنوان مؤلفه های اصلی x می شناسیم. با قرار دادن $y_1 = e'_1 x$ ، $y_2 = e'_2 x$ ، \dots ، $y_p = e'_p x$ داریم:

$$c^2 = \frac{1}{\lambda_1} y_1^2 + \frac{1}{\lambda_2} y_2^2 + \dots + \frac{1}{\lambda_p} y_p^2$$

که این معادله یک بیضوی (زیرا $\lambda_1, \lambda_2, \dots, \lambda_p$ مثبت اند) را در یک دستگاه مختصات تعریف می کند که محورهای y_1, y_2, \dots, y_p آن به ترتیب در جهات e_1, e_2, \dots, e_p قرار دارد. اگر λ_1 بزرگترین مقدار ویژه باشد، آن گاه محور بزرگ آن در جهت e_1 واقع است و محورهای کوچک باقی مانده در جهاتی که با e_2, \dots, e_p تعریف شده اند، قرار دارند.

به طور خلاصه مؤلفه های اصلی $y_1 = e'_1 x$ ، $y_2 = e'_2 x$ ، \dots ، $y_p = e'_p x$ در جهات محورهای

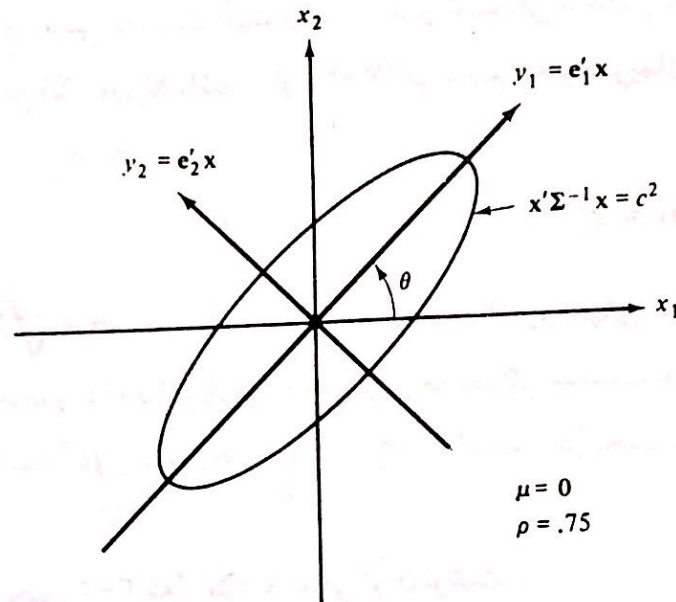
بیضوی چگالی ثابت واقع اند. بنابراین هر نقطه روی محور بیضوی i ام دارای مختصات x متناسب

۱- این را بدون این که به کلیت خللی وارد شود می توان انجام داد، زیرا بردار تصادفی نرمال X را همیشه می توان

به بردار تصادفی نرمال $W = X - \mu$ و $E(W) = 0$ تبدیل کرد. با وجود این $\text{Cov}(W) = \text{Cov}(X)$.

تحلیل آماری چندمتغیری کاربردی

با $e'_i = [e_{1i}, e_{2i}, \dots, e_{pi}]$ و الزاماً مختصات مؤلفه اصلی به شکل $[0, \dots, 0, y_p, 0, \dots, 0]$ است. یک بیضی چگالی ثابت و مؤلفه های اصلی مربوط به بردار تصادفی نرمال دو متغیری با $\mu = 0$ و $\rho = .75$ را در شکل ۸-۱ نشان داده ایم. می بینیم که مؤلفه های اصلی از دوران محورهای مختصات اولیه به اندازه زاویه θ تا این که بر محورهای بیضی چگالی ثابت منطبق شوند، به دست می آیند. این نتیجه برای ابعاد $p > 2$ نیز برقرار است.



شکل ۸-۱ بیضی چگالی ثابت $x'\Sigma^{-1}x = c^2$ و مؤلفه های اصلی y_1, y_2 برای یک بردار تصادفی نرمال دو متغیره X .

محاسبه مؤلفه های اصلی متغیرهای استاندارد شده

مؤلفه های اصلی را می توان برای متغیرهای استاندارد شده نیز به دست آورد:

$$Z_1 = \frac{(X_1 - \mu_1)}{\sqrt{\sigma_{11}}}$$

$$Z_2 = \frac{(X_2 - \mu_2)}{\sqrt{\sigma_{22}}}$$

$$\vdots$$

$$Z_p = \frac{(X_p - \mu_p)}{\sqrt{\sigma_{pp}}}$$

با نماد ماتریسی می توان نوشت:

$$Z = (V^{1/2})^{-1}(X - \mu)$$

(۸-۱۰)

که در آن ماتریس انحراف استاندارد قطری $V^{1/2}$ را در (۲-۳۵) تعریف کردیم. براساس (۲-۳۷) واضح

است که $E(Z) = 0$ و

$$\text{Cov}(Z) = (V^{1/2})^{-1} \Sigma (V^{1/2})^{-1} = \rho$$

مؤلفه های اصلی Z را از بردارهای ویژه ماتریس همبستگی ρ ، X می توان به دست آورد. چون واریانس هر Z_i برابر واحد است، لذا تمام نتایج قبلی را به سهولت می توان به کار برد. ما در آینده نیز از نماد Y_i برای نشان دادن مؤلفه اصلی i ام و (λ_i, e_i) برای زوج مقدار ویژه - بردار ویژه استفاده می کنیم. با این وجود کمیت های به دست آمده از Σ در کل همان کمیت های به دست آمده از ρ نیستند.

نتیجه ۸-۴. مؤلفه اصلی i ام متغیرهای استاندارد شده $Z' = [Z_1, Z_2, \dots, Z_p]$ با $\text{Cov}(Z) = \rho$ با

$$Y_i = e_i' Z = e_i' (V^{1/2})^{-1} (X - \mu), \quad i = 1, 2, \dots, p$$

داده می شود. علاوه بر این

$$\sum_{i=1}^p \text{Var}(Y_i) = \sum_{i=1}^p \text{Var}(Z_i) = p \quad (۱۱-۸)$$

و

$$\rho_{Y_i, Z_k} = e_{ki} \sqrt{\lambda_i}, \quad i, k = 1, 2, \dots, p$$

در این حالت $(\lambda_1, e_1), (\lambda_2, e_2), \dots, (\lambda_p, e_p)$ زوج های مقدار ویژه - بردار ویژه ρ با $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ هستند.

اثبات. اگر Z_1, Z_2, \dots, Z_p را جایگزین X_1, X_2, \dots, X_p کنیم، ρ را به جای Σ قرار دهیم، نتیجه ۸-۴ از نتایج ۸-۱، ۸-۲ و ۸-۳ به دست می آید.

از (۱۱-۸) می بینیم که واریانس کل جامعه (متغیرهای استاندارد شده) برابر p یعنی مجموع اعضای قطری ماتریس ρ است. اگر در (۷-۸) از Z به جای X استفاده کنیم سهم واریانس کل بیان شده با مؤلفه اصلی k ام Z ،

$$\left(\begin{array}{c} \text{سهم واریانس (استاندارد شده)} \\ \text{جامعه مربوط به مؤلفه اصلی} \\ \text{ام } k \end{array} \right) = \frac{\lambda_k}{p}, \quad k = 1, 2, \dots, p \quad (۱۲-۸)$$

است که در آن λ_k ها مقادیر ویژه ρ هستند.

مثال ۸-۲ (مؤلفه های اصلی حاصل از ماتریسهای کوواریانس همبستگی)
ماتریس کوواریانس

$$\Sigma = \begin{bmatrix} 1 & .4 \\ .4 & 100 \end{bmatrix}$$

و ماتریس همبستگی به دست آمده از آن

$$\rho = \begin{bmatrix} 1 & .4 \\ .4 & 1 \end{bmatrix}$$

را در نظر می گیریم . زوجهای مقدار ویژه-بردار ویژه Σ ، عبارتند از :

$$\lambda_1 = 100.16, \quad e'_1 = [.040, .999]$$

$$\lambda_2 = .84, \quad e'_2 = [.999, -.040]$$

به طور مشابه زوجهای مقدار ویژه-بردار ویژه حاصل از ρ ، عبارتند از :

$$\lambda_1 = 1 + \rho = 1.4, \quad e'_1 = [.707, .707]$$

$$\lambda_2 = 1 - \rho = .6, \quad e'_2 = [.707, -.707]$$

مؤلفه های اصلی مربوط به صورت زیر در می آید :

$$\Sigma: \begin{aligned} Y_1 &= .040 X_1 + .999 X_2 \\ Y_2 &= .999 X_1 - .040 X_2 \end{aligned}$$

و

$$\begin{aligned} Y_1 &= .707Z_1 + .707Z_2 = .707\left(\frac{X_1 - \mu_1}{1}\right) + .707\left(\frac{X_2 - \mu_2}{10}\right) \\ &= .707(X_1 - \mu_1) + .0707(X_2 - \mu_2) \end{aligned}$$

ρ :

$$\begin{aligned} Y_2 &= .707Z_1 - .707Z_2 = .707\left(\frac{X_1 - \mu_1}{1}\right) - .707\left(\frac{X_2 - \mu_2}{10}\right) \\ &= .707(X_1 - \mu_1) - .0707(X_2 - \mu_2) \end{aligned}$$

X_2 به خاطر واریانس بزرگی که دارد اولین مؤلفه اصلی تعیین شده از Σ را کاملاً تحت تأثیر قرار می دهد. علاوه بر این اولین مؤلفه اصلی یک نسبت

$$\frac{\lambda_1}{\lambda_1 + \lambda_2} = \frac{100.16}{101} = .992$$

از واریانس کل جامعه را بیان می کند . در عین حال وقتی متغیرهای X_1 و X_2 را استاندارد می کنیم ، متغیرهای حاصل سهمی یکسان در مؤلفه های اصلی تعیین شده از ρ دارند . با استفاده از ۸-۴ ، داریم :

$$\rho_{Y_1, Z_1} = e_{11} \sqrt{\lambda_1} = .707 \sqrt{1.4} = .837$$

و

$$\rho_{Y_1, Z_2} = e_{21} \sqrt{\lambda_1} = .707 \sqrt{1.4} = .837$$

در این حالت اولین مؤلفه اصلی یک نسبت

$$\frac{\lambda_1}{p} = \frac{1.4}{2} = .7$$

از واریانس کل جامعه (استاندارد شده) را بیان می کند .

ملاحظه می کنیم که اهمیت نسبی متغیرها ، مثلاً روی اولین مؤلفه اصلی به طور قابل ملاحظه ای تحت تأثیر استاندارد کردن واقع می شود . هنگامی که مؤلفه های اصلی به دست آمده از p را بر حسب X_1 و X_2 بیان می کنیم ، مقدار نسبی وزنه های $.707$ و $.707$ با مقدار نسبی وزنه های $.040$ و $.999$ منضم به این متغیرها در مؤلفه های اصلی به دست آمده از Σ در جهت مخالف یکدیگرند . مثال قبلی بیان می کند که مؤلفه های اصلی به دست آمده از Σ با مؤلفه های اصلی به دست آمده از p متفاوت اند . علاوه بر این یک مجموعه از مؤلفه های اصلی تابعی ساده از سایر مؤلفه ها نیست . که این بدان معنی است که استاندارد کردن بی ربط نیست .

متغیرها را در صورتی که با مقیاسهایی که به طور وسیعی با هم تفاوت دارند ، واحدهای اندازه گیری آنها متناسب نیست اندازه گیری کنیم ، احتمالاً باید استاندارد شوند . برای مثال اگر X_1 فروش سالانه را در فاصله $10,000$ دلار تا $350,000$ دلار نشان دهد و X_2 نسبت درآمد خالص سالانه / کل دارایی که در فاصله 0.01 تا 0.60 قرار دارد را نشان دهد در آن صورت کل تغییرات تقریباً به خاطر فروش دلار خواهد بود . در این حالت انتظار یک مؤلفه اصلی (مهم) با وزن دار کردن زیاد X_1 را داریم . از طرفی اگر هر دو متغیر را استاندارد کنیم ، مقادیر بعدی آنها به همان ترتیب خواهد بود و X_2 (یا Z_2) نقش مهمتری را در ساختار مؤلفه ها ایفا می کند . که این رفتار را در مثال (۸-۲) مشاهده کردیم .

مؤلفه های اصلی برای ماتریسهای کوواریانس با ساختارهای ویژه

ماتریسهای کوواریانس و همبستگی دارای طرحهای خاصی هستند که مؤلفه های اصلی آنها را به

شکلهای ساده می توان بیان کرد . فرض کنید Σ ماتریس قطری زیر باشد :

$$\Sigma = \begin{bmatrix} \sigma_{11} & 0 & \cdots & 0 \\ 0 & \sigma_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_{pp} \end{bmatrix} \quad (8-13)$$

با قرار دادن $e'_i = [0, \dots, 0, 1, 0, \dots, 0]$ که جمله i ام آن است، مشاهده می کنیم که

$$\begin{bmatrix} \sigma_{11} & 0 & \dots & 0 \\ 0 & \sigma_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_{pp} \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \sigma_{ii} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{یا} \quad \Sigma e_i = \sigma_{ii} e_i$$

و نتیجه می گیریم که (σ_{ii}, e_i) زوج مقدار ویژه- بردار ویژه i ام است. چون ترکیب خطی $X_i = e'_i X$ ، لذا مجموعه مؤلفه های اصلی درست مجموعه اولیه متغیرهای تصادفی ناهمبسته است.

برای یک ماتریس کوواریانس با طرح (۸-۱۳) با خلاصه کردن مؤلفه های اصلی چیزی به دست نمی آید. از نقطه نظر دیگر اگر X دارای توزیع $N_p(\mu, \Sigma)$ باشد مسیرهای چگالی ثابت بیضویهایی هستند که محورهای آن قبلاً در جهات بیشترین تغییر بود. بنابراین نیازی به دوران دستگاه مختصات نیست.

استاندارد کردن وضعیت را برای Σ در (۸-۱۳) به طور اساسی تغییر نمی دهد. در این حالت

$\rho = I$ ، ماتریس همانی $p \times p$ است. واضح است که $\rho e_i = 1 e_i$ ، لذا مقدار ویژه ۱ دارای مضرب ρ است و $e'_i = [0, \dots, 0, 1, 0, \dots, 0]$ ، $i = 1, 2, \dots, p$ انتخابهای مناسب بردارهای ویژه هستند. در نتیجه مؤلفه های اصلی که برای ρ به دست می آیند، مؤلفه های اصلی متغیرهای اولیه Z_1, Z_2, \dots, Z_p نیز هستند. علاوه بر این در این حالت که مقادیر ویژه برابرند، بیضویهای نرمال چندمتغیری چگالی ثابت کروی هستند.

طرح دیگری از ماتریس کوواریانس که اغلب رابطه میان بعضی متغیرهای مربوط به زیست شناسی نظیر میزان حیات اشیاء را بیان می کند، شکل کلی زیر را دارد:

$$\Sigma = \begin{bmatrix} \sigma^2 & \rho\sigma^2 & \dots & \rho\sigma^2 \\ \rho\sigma^2 & \sigma^2 & \dots & \rho\sigma^2 \\ \vdots & \vdots & \ddots & \vdots \\ \rho\sigma^2 & \rho\sigma^2 & \dots & \sigma^2 \end{bmatrix} \quad (۸-۱۴)$$

ماتریس همبستگی که از این ماتریس نتیجه می شود نیز ماتریس کوواریانس متغیرهای استاندارد شده است.

$$\rho = \begin{bmatrix} 1 & \rho & \dots & \rho \\ \rho & 1 & \dots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \dots & 1 \end{bmatrix} \quad (۸-۱۵)$$

از ماتریس (۸-۱۵) نتیجه می شود که متغیر X_1, X_2, \dots, X_p همبستگی یکسانی دارند .
اثبات این مطلب که (تمرین ۸-۵ را ملاحظه کنید) p مقدار ویژه ماتریس همبستگی (۸-۱۵) را
به دو دسته می توان تقسیم کرد مشکل نیست . وقتی p مثبت است بزرگترین مقدار ویژه

$$\lambda_1 = 1 + (p - 1)\rho \quad (۸-۱۶)$$

متناظر با بردار ویژه

$$\mathbf{e}'_1 = \left[\frac{1}{\sqrt{p}}, \frac{1}{\sqrt{p}}, \dots, \frac{1}{\sqrt{p}} \right] \quad (۸-۱۷)$$

است . $p - 1$ مقدار ویژه باقی مانده

$$\lambda_2 = \lambda_3 = \dots = \lambda_p = 1 - \rho$$

بوده و یک انتخاب بردارهای ویژه ، عبارت است از :

$$\mathbf{e}'_2 = \left[\frac{1}{\sqrt{1 \times 2}}, \frac{-1}{\sqrt{1 \times 2}}, 0, \dots, 0 \right]$$

$$\mathbf{e}'_3 = \left[\frac{1}{\sqrt{2 \times 3}}, \frac{1}{\sqrt{2 \times 3}}, \frac{-2}{\sqrt{2 \times 3}}, 0, \dots, 0 \right]$$

⋮

$$\mathbf{e}'_i = \left[\frac{1}{\sqrt{(i-1)i}}, \dots, \frac{1}{\sqrt{(i-1)i}}, \frac{-(i-1)}{\sqrt{(i-1)i}}, 0, \dots, 0 \right]$$

⋮

$$\mathbf{e}'_p = \left[\frac{1}{\sqrt{(p-1)p}}, \dots, \frac{1}{\sqrt{(p-1)p}}, \frac{-(p-1)}{\sqrt{(p-1)p}} \right]$$

اولین مؤلفه اصلی

$$Y_1 = \mathbf{e}'_1 \mathbf{X} = \frac{1}{\sqrt{p}} \sum_{i=1}^p X_i$$

با مجموع p متغیر اولیه متناسب است ، که این را می توان به عنوان یک «شاخص» با وزنهای مساوی
تلقی کرد . این مؤلفه اصلی یک نسبت

$$\frac{\lambda_1}{p} = \frac{1 + (p-1)\rho}{p} = \rho + \frac{1-\rho}{p} \quad (۸-۱۸)$$

از کل تغییرات جامعه را بیان می کند. ملاحظه می کنیم که برای ρ نزدیک به ۱ یا p بزرگ $\frac{\lambda_1}{p} = \rho$.
 برای مثال اگر $\rho = 0.80$ و $p = 5$ باشد، مؤلفه اول 84% کل واریانس را بیان می کند. وقتی ρ نزدیک به ۱ است $p - 1$ مؤلفه آخر جمعاً سهم بسیار کمی از کل واریانس را دارند و اغلب از آنها می توان صرف نظر کرد.

اگر متغیرهای استاندارد شده Z_1, Z_2, \dots, Z_p دارای یک توزیع نرمال چندمتغیری با ماتریس کوواریانسی که با (۸-۱۵) داده می شود باشد آن گاه بیضویهای چگالی ثابت «سیگاری شکل» بوده و محور بزرگ (اصلی) در امتداد اولین مؤلفه اصلی $Y_1 = (\frac{1}{p})[1, 1, \dots, 1]$ است. محورهای کوچک (فرعی) (و بقیه مؤلفه های اصلی) در جهات شکل کرووی متقارن عمود بر محور بزرگ (و مؤلفه اصلی اول) روی می دهند.

۳-۸ خلاصه کردن تغییرات نمونه به وسیله مؤلفه های اصلی

ما اکنون چارچوبی را که برای مطالعه مسأله خلاصه کردن تغییرات در n راندازه وی p متغیر لازم است را با چند انتخاب درست ترکیبات خطی در دست داریم:

فرض می کنیم داده های x_1, x_2, \dots, x_n استخراج مستقل از جامعه p بعدی با بردار میانگین μ و ماتریس واریانس Σ را نشان می دهد این داده ها بردار میانگین نمونه \bar{x} ، ماتریس واریانس نمونه S و ماتریس همبستگی نمونه R را حاصل می کند.

هدف ما در این بخش این است که ترکیبات خطی ناهمبسته ای از خصیصه های اندازه گیری شده ای را بسازیم که در بیشتر تغییرات نمونه به حساب می آیند ترکیبات ناهمبسته با بزرگترین واریانس را مؤلفه های اصلی نمونه می نامیم.

خاطر نشان می کنیم که n مقدار هر ترکیب خطی

$$l'_1 x_j = l_{11}x_{1j} + l_{21}x_{2j} + \dots + l_{p1}x_{pj}, \quad j = 1, 2, \dots, n$$

دارای میانگین نمونه $l'_1 \bar{x}$ واریانس نمونه $l'_1 S l_1$ است. همچنین زوج مقادیر $(l'_1 x_j, l'_2 x_j)$ برای دو ترکیب خطی دارای ماتریس کوواریانس نمونه $l'_2 S l_2$ است [۳-۳۶] را ملاحظه نمایید.

مؤلفه های اصلی نمونه را به صورت ترکیبات خطی که دارای واریانس نمونه ماکزیمم است، تعریف می کنیم. بردارهای ضرایب l_i مانند کمیت های جامعه باید در $l'_i l_i = 1$ صدق کنند، به ویژه اولین مؤلفه اصلی نمونه = ترکیب خطی $l'_1 x_j$ که واریانس نمونه $l'_1 x_j$ را با شرط $l'_1 l_1 = 1$ ماکزیمم می کند.

دومین مؤلفه اصلی نمونه = ترکیب خطی $l'_2 x_j$ که واریانس نمونه $l'_2 x_j$ را با شرط $l'_1 l_1 = 1$ و کوواریانس صفر برای زوجهای $(l'_1 x_j, l'_2 x_j)$ ماکزیمم کند.

در مرحله i ام

i امین مؤلفه اصلی = ترکیب خطی $l'_i x_j$ که واریانس نمونه $l'_i x_j$ را با شرط $l'_i l_i = 1$ و کوواریانس نمونه صفر برای تمام زوجهای $(l'_i x_j, l'_k x_j)$ ، $k < i$ ماکزیمم کند.

اولین مؤلفه اصلی $l'_1 S l_1$ یا معادل با آن

$$\frac{l'_1 S l_1}{l'_1 l_1} \quad (19-8)$$

را ماکزیمم می کند. بر اساس (۲-۵۱) ماکزیمم آن بزرگترین مقدار ویژه $\hat{\lambda}_i$ است که از انتخاب بردار ویژه $l_1 = e_1$ ، ماتریس S به دست می آید. به شرطی که $0 = l'_1 S e_k = l'_1 \hat{\lambda}_k e_k$ یا l_1 بر e_k عمود باشد انتخابهای متوالی l_i ، (۸-۱۹) را ماکزیمم می کند. بنابراین مانند اثبات نتایج ۸.۱-۸.۳ نتایج زیر که مربوط به مؤلفه های اصلی نمونه است را به دست می آوریم.

اگر $S = \{s_{ik}\}$ ماتریس واریانس نمونه $p \times p$ با زوجهای مقدار ویژه بردار ویژه $(\hat{\lambda}_1, e_1), (\hat{\lambda}_2, e_2), \dots, (\hat{\lambda}_p, e_p)$ باشد، در آن صورت مؤلفه اصلی نمونه i ام با

$$\hat{y}_i = \hat{e}'_i x = \hat{e}_{i1} x_1 + \hat{e}_{i2} x_2 + \dots + \hat{e}_{ip} x_p, \quad i = 1, 2, \dots, p$$

داده می شود که $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$ و هر مشاهده روی متغیرهای X_1, X_2, \dots, X_p است.

همچنین

$$(\hat{y}_k) = \hat{\lambda}_k, \quad k = 1, 2, \dots, p$$

$$(\hat{y}_i, \hat{y}_k) = 0, \quad i \neq k \quad (20-8)$$

علاوه بر این

$$\text{کل واریانس نمونه} = \sum_{i=1}^p s_{ii} = \hat{\lambda}_1 + \hat{\lambda}_2 + \dots + \hat{\lambda}_p$$

و

$$r_{y_i, x_k} = \frac{\hat{e}_{ki} \sqrt{\hat{\lambda}_i}}{\sqrt{s_{kk}}}, \quad i, k = 1, 2, \dots, p$$

ما صرف نظر از این که مؤلفه های اصلی نمونه را از S یا R به دست آورده ایم ، آنها را با $\lambda_1, \lambda_2, \dots, \lambda_p$ نشان می دهیم . مؤلفه هایی که از S یا R حاصل می شوند در کل یکسان نیستند ولی از فرضیه واضح است که از کدام ماتریس استفاده شده و کدام نماد λ_i مناسب است . همچنین بهتر است که مؤلفه بردارهای ضرایب را با \hat{e}_i و مؤلفه واریانسها را برای هر دو وضعیت با $\hat{\lambda}_i$ نشان دهیم . اغلب مشاهدات x_j را با کم کردن از \bar{x} «مرکزی» می کنیم . این موضوع روی ماتریس کوواریانس نمونه S تأثیری ندارد و برای هر بردار مشاهده x ، مؤلفه اصلی زیر را می دهد :

$$\hat{y}_i = \hat{e}_i'(x - \bar{x}), \quad i = 1, 2, \dots, p \quad (21-8)$$

اگر مقادیر مؤلفه i ام

$$\hat{y}_{ij} = \hat{e}_i'(x_j - \bar{x}), \quad i = 1, 2, \dots, p \quad (22-8)$$

که از جایگزین کردن هر مشاهده x_j به جای x دلخواه در (21-8) تولید می شود ، آن گاه :

$$\bar{y}_i = \frac{1}{n} \sum_{j=1}^n \hat{e}_i'(x_j - \bar{x}) = \frac{1}{n} \hat{e}_i' \left(\sum_{j=1}^n (x_j - \bar{x}) \right) = \frac{1}{n} \hat{e}_i' \mathbf{0} = 0 \quad (23-8)$$

یعنی میانگین نمونه هر مؤلفه اصلی صفر است . واریانسهای نمونه ، مانند (20-8) هنوز با $\hat{\lambda}_i$ ها داده می شود .

مثال ۳-۸

آمارگیری سال ۱۹۷۰ اطلاعات ناحیه ای روی ۵ متغیر اقتصادی اجتماعی را برای ناحیه

۱- اگر X_j ها دارای توزیع نرمال باشند (نتیجه ۴-۱۱ را ملاحظه کنید) در آن صورت مؤلفه های اصلی نمونه را می توان از $S_n = \Sigma$ یعنی برآورد درست نمایی ماکزیمم ماتریس کوواریانس Σ نیز به دست آورد . در این حالت به شرط این که مقادیر ویژه Σ متمایز باشند ، مؤلفه های اصلی نمونه را می توان به عنوان برآوردهای درست نمایی ماکزیمم پارامترهای متناظر جامعه مربوط در نظر نخواهیم گرفت ([۱] را ملاحظه کنید) . چون فرض نرمال در این بخش لازم نیست ، لذا Σ را در نظر گرفت . همچنین Σ دارای مقادیر ویژه $[(n-1)/n] \hat{\lambda}_i$ و بردارهای ویژه متناظر \hat{e}_i اند که $(\hat{\lambda}_i, \hat{e}_i)$ زوجهای مقدار ویژه-بردار ویژه S هستند . از این رو S و Σ هر دو مؤلفه های اصلی یکسان x \hat{e}_i' (20-8) را ملاحظه نمایید] و نسبت واریانس بیان شده یکسان $(\hat{\lambda}_1 + \hat{\lambda}_2 + \dots + \hat{\lambda}_p) / \hat{\lambda}_i$ را می دهد . سرانجام S و Σ هر دو ماتریس همبستگی نمونه یکسان R را می دهند ، لذا اگر متغیرها را استاندارد کرده باشیم . انتخاب S یا Σ فرق نمی کند .

مدیسون ویسکانسین فراهم می کند. داده های مربوط به ۱۴ ناحیه را در جدول (۸-۲) در تمرینهای آخرین فصل ثبت کرده ایم. از این داده ها آماره های زیر به دست می آید:

$$\bar{x}' = \begin{bmatrix} 4.32, & 14.01, & 1.95, & 2.17, & 2.45 \\ \text{کل} & \text{میانۀ} & \text{کل استخدای} & \text{خدمات بهداشتی} & \text{میانۀ ارزش} \\ \text{جامعه} & \text{سالهای} & \text{(به هزار)} & \text{استخدای} & \text{منزل} \\ \text{(به هزار)} & \text{مدرسه} & & \text{(به صد)} & \text{(\$10,000s)} \end{bmatrix}$$

و

$$S = \begin{bmatrix} 4.308 & 1.683 & 1.803 & 2.155 & -.253 \\ 1.683 & 1.768 & .588 & .177 & .176 \\ 1.803 & .588 & .801 & 1.065 & -.158 \\ 2.155 & .177 & 1.065 & 1.970 & -.357 \\ -.253 & .176 & -.158 & -.357 & .504 \end{bmatrix}$$

آیا می توان تغییرات را به یک یا دو مؤلفه اصلی خلاصه کرد؟
جدول زیر را پیدا می کنیم.

ضرایب مربوط به مؤلفه های اصلی (اعداد داخل پرانتزها ضرایب همبستگی اند)

متغیر	$e_1(r_{1 \cdot x_k})$	$e_2(r_{2 \cdot x_k})$	e_3	e_4	e_5
کل جامعه	۰٫۷۸۱ (۰٫۹۹)	-۰٫۰۷۱ (-۰٫۰۴)	۰٫۰۰۴	۰٫۵۴۲	-۰٫۳۰۲
میانۀ سالهای مدرسه	۰٫۳۰۶ (۰٫۶۱)	-۰٫۷۶۴ (-۰٫۷۶)	-۰٫۱۶۲	-۰٫۵۴۵	-۰٫۰۱۰
کل استخدای	۰٫۳۳۴ (۰٫۹۸)	۰٫۰۸۳ (۰٫۱۲)	۰٫۰۱۵	۰٫۰۵	۰٫۹۳۷
خدمات بهداشتی استخدای	۰٫۴۲۶ (۰٫۸۰)	۰٫۵۷۹ (۰٫۵۵)	۰٫۲۲۰	-۰٫۶۳۶	-۰٫۱۷۳
میانۀ ارزش منزل	-۰٫۰۵۴ (-۰٫۲۰)	-۰٫۲۶۲ (-۰٫۴۹)	۰٫۹۶۲	-۰٫۰۵۱	۰٫۰۲۴
پراش ($\hat{\lambda}_i$)	۶٫۶۳۹	۱٫۷۸۶	۰٫۳۹۰	۰٫۲۳۰	۰٫۰۱۴
درصد تجمعی واریانس کل	۷۴٫۱	۹۳٫۲	۹۷٫۴	۹۹٫۹	۱۰۰

اولین مؤلفه اصلی ۷۴٫۱٪ واریانس کل نمونه را بیان می کند. دو مؤلفه اصلی اول با هم ۹۳٫۲٪ واریانس کل را بیان می کنند. در نتیجه تغییرات نمونه خیلی خوب با دو مؤلفه اصلی خلاصه می شود و یک کاهش در داده ها از چهارده مشاهده روی پنج متغیر به چهارده مشاهده روی دو مؤلفه اصلی معقول است.

با معلوم بودن ضرایب بالا ، به نظر می رسد مؤلفه اصلی اول یک متوسط موزون چهار متغیر اول است . دومین مؤلفه اصلی به نظر می رسد خدمات سلامتی استخدام را با یک متوسط وزن دار میانه سالهای مدرسه و میانه ارزش منزل مقابله می کند .

هنگامی که می خواهیم در باره موضوع تعبیر مؤلفه های اصلی کار کنیم ، شاید همبستگیهای r_{y_i, x_k} راهنماهای قابل اعتمادتری از مؤلفه ضرایب e_{ki} باشند . همبستگیها تفاوتهای در واریانسهای متغیرهای اولیه را مجاز می نمایند که بدین وسیله مشکل تعبیر و تفسیری که مقایسه های متفاوت اندازه گیری موجب آن می شود برطرف می گردد . در مثال (۸-۳) ضرایب همبستگی که در جدول نشان داده شده است ، تعبیری که با مؤلفه ضرایب حاصل می شود را تأیید می کند .

مثال ۸-۴

در مطالعه روابط موجود در اندازه و شکل لاک پشتهای رنگ شده جالیکلور و موسیمان [۱۰] درازا ، عرض و ارتفاع لاک پشت را اندازه گیری نموده اند . این داده ها در تمرین (۶.۱۳) آورده شده و جدول (۶-۵) تحلیلی را بر حسب لگاریتمها پیشنهاد می کند . (جالیکور معمولاً یک تبدیل لگاریتمی را در مطالعه روابط اندازه و شکل پیشنهاد می کند) . یک تحلیل مؤلفه اصلی را انجام می دهیم .

لگاریتمهای طبیعی ابعاد ۲۴ لاک پشت نیز دارای بردار میانگین نمونه

$$\bar{x}' = [4.725, 4.478, 3.703]$$

$$S = 10^{-3} \begin{bmatrix} 11.555 & 8.367 & 8.508 \\ 8.367 & 6.697 & 6.264 \\ 8.508 & 6.264 & 7.061 \end{bmatrix}$$

است . یک تحلیل مؤلفه اصلی خلاصه زیر را به ما می دهد .

ضرایب مربوط به مؤلفه های اصلی (ضرایب همبستگی در داخل پرانتزها هستند)

متغیر	$e_1(r_{y_i, x_k})$	e_2	e_3
لگاریتم (درازا)	۰٫۶۸۳ (۰٫۹۹)	۰٫۱۶۲	۰٫۷۱۲
لگاریتم (عرض)	۰٫۵۱۰ (۰٫۹۷)	۰٫۵۹۱	-۰٫۶۲۴
لگاریتم (ارتفاع)	۰٫۵۲۲ (۰٫۹۷)	-۰٫۷۹۰	-۰٫۳۲۱
واریانس ($\hat{\lambda}_i$)	$۲۴٫۳۱ \times ۱۰^{-۳}$	$۰٫۶۳ \times ۱۰^{-۳}$	$۰٫۳۸ \times ۱۰^{-۳}$
درصد تجمعی واریانس کل	۹۶٫۰	۹۸٫۵	۱۰۰

اولین مؤلفه اصلی که ۹۶٪ واریانس کل را بیان می کند یک تعبیر جالبی دارد . چون

$$\hat{y}_1 = .683 \ln(\text{length}) + .51 \ln(\text{width}) + .522 \ln(\text{height}) \\ = \ln[(\text{ارتفاع})^{.522} (\text{عرض})^{.510} (\text{طول})^{.683}]$$

لذا اولین مؤلفه اصلی را می توان به عنوان (حجم) hn جعبه ای با ابعاد تعدیل شده در نظر گرفت . داده ها را از نظر هندسی می توان به صورت n نقطه در فضای p -بعدی رسم کرد . اگر S معین و

مثبت باشد ، آن گاه تمام بردارهای $1 \times p$ ، x که در

$$(x - \bar{x})' S^{-1} (x - \bar{x}) = c^2 \quad (۲۴-۸)$$

صدق می کند یک ابر بیضوی متمرکز در \bar{x} که محورهای ویژه S^{-1} یا معادل با آن با بردارهای ویژه S داده می شود را تعریف می کند (بخش (۲-۳) و نتیجه (۴-۱) را که S جایگزین Σ می شود ، ملاحظه کنید) طولهای این محورها متناسب با $\sqrt{\hat{\lambda}_i}$ ، $i = 1, 2, \dots, p$ ، بوده که $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$ مقادیر ویژه S هستند .

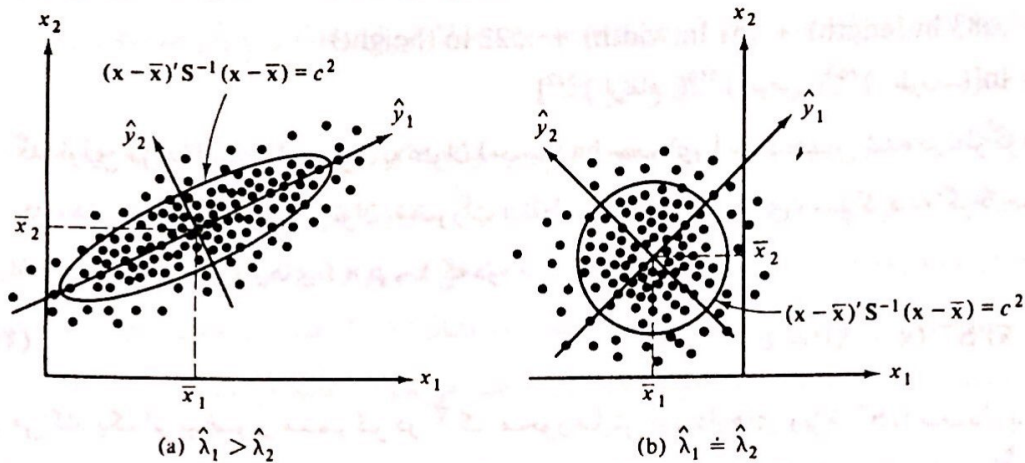
مؤلفه های اصلی نمونه همان ارتباطی را با محورهای بیضویهای با فاصله ثابت در (۲۴-۸) دارند که مؤلفه های جامعه با بیضویهای با چگالی ثابت برای متغیرهای $N_p(\mu, \Sigma)$ دارند . یعنی مؤلفه های نمونه در امتداد محورهای بیضویهای با فاصله ثابت قرار دارد . به این مؤلفه ها می توان به عنوان نتیجه دوران دستگاه مختصات اولیه تا این که محورهای مختصات در جهات واریانس ماکزیمم از نمودار پراکنندگی عبور کند ، نگاه کرد .

قدر مطلق مؤلفه اصلی i ام ، $|e_i(x - \bar{x})| = |e_i| |x - \bar{x}|$ طول تصویر بردار $(x - \bar{x})$ روی بردار واحد e_i را می دهد [(۲-۸) و (۲-۹) را ملاحظه کنید] . این تعبیر هندسی مؤلفه های اصلی نمونه را در شکل ۲-۸ برای $p = 2$ تشریح نموده ایم .

شکل ۲-۸ (الف) یک بیضی با فاصله ثابت متمرکز در \bar{x} را با $\hat{\lambda}_1 > \hat{\lambda}_2$ نشان می دهد . مؤلفه های اصلی نمونه به خوبی تعیین می شوند . این مؤلفه ها در امتداد محورهای بیضی در جهات عمود واریانس نمونه ماکزیمم قرار دارند . شکل ۲-۸ (ب) یک بیضی با فاصله ثابت متمرکز در \bar{x} را با $\hat{\lambda}_1 = \hat{\lambda}_2$ نشان می دهد . در این حالت محورهای بیضی (دایره) با چگالی ثابت به طور منحصر به فردی تعیین نمی شوند و می توانند در هر دو جهت متعامد از جمله در جهات محورهای مختصات اولیه قرار گیرند .

به طور مشابه مؤلفه های اصلی نمونه می تواند در هر دو جهت متعامد از جمله جهات محورهای مختصات اولیه واقع شود . وقتی مسیرهای با فاصله ثابت تقریباً مدورند یا معادل با آن وقتی مقادیر ویژه S تقریباً مساویند ، تغییرات نمونه در تمام جهات همگن است . در این صورت داده ها را نمی توان

در کمتر از p جهت نشان داد .



شکل ۸-۲ مؤلفه‌های اصلی نمونه و بیضی‌های با فاصله ثابت

اگر x_1, x_2, \dots, x_n را بتوان به صورت نمونه‌ای از یک جامعه نرمال تلقی کرد، در آن صورت مؤلفه‌های اصلی نمونه، $\hat{y}_i = \hat{e}_i'(x - \bar{x})$ ، مصادیق مؤلفه‌های اصلی جامعه، $Y_i = e_i'(X - \mu)$ هستند که دارای توزیع $N_p(0, \Lambda)$ هستند. ماتریس قطری Λ دارای درایه‌های $\lambda_1, \lambda_2, \dots, \lambda_p$ است و (λ_i, e_i) زوجهای مقدار ویژه-بردار ویژه Σ هستند. همچنین بیضویهای با فاصله ثابت (۸-۲۴) برآوردهای بیضویهای با چگالی ثابت $(x - \mu)'\Sigma^{-1}(x - \mu) = c^2$ هستند. فرض نرمال برای روشهای استنباطی مورد بحث در (۸-۵) مفید است، ولی برای تعمیم خواص مؤلفه‌های اصلی نمونه خلاصه شده در (۸-۲۰) لازم نیستند.

به طور کلی مؤلفه‌های اصلی نمونه نسبت به تغییرات در مقیاس (تمرین ۲.۸ را ملاحظه نمایید) پایا نیستند. به طوری که در عملکرد مؤلفه‌های جامعه متذکر شدیم، متغیرهایی که با مقیاسهای مختلف یا مقیاس مشترک ولی در فاصله‌های با تغییر زیاد اندازه گیری می‌شوند را اغلب استاندارد می‌کنیم. استاندارد کردن را برای نمونه با ساختن

$$z_j = D^{-1/2}(x_j - \bar{x}) = \begin{bmatrix} \frac{x_{1j} - \bar{x}_1}{\sqrt{s_{11}}} \\ \frac{x_{2j} - \bar{x}_2}{\sqrt{s_{22}}} \\ \vdots \\ \frac{x_{pj} - \bar{x}_p}{\sqrt{s_{pp}}} \end{bmatrix} \quad j = 1, 2, \dots, n \quad (۲۵-۸)$$

انجام می‌دهیم

ماتریس داده های $p \times n$ مشاهدات استاندارد شده

$$\mathbf{Z} = [\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n] = \begin{bmatrix} z_{11} & z_{12} & \dots & z_{1n} \\ z_{21} & z_{22} & \dots & z_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ z_{p1} & z_{p2} & \dots & z_{pn} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{x_{11} - \bar{x}_1}{\sqrt{s_{11}}} & \frac{x_{12} - \bar{x}_1}{\sqrt{s_{11}}} & \dots & \frac{x_{1n} - \bar{x}_1}{\sqrt{s_{11}}} \\ \frac{x_{21} - \bar{x}_2}{\sqrt{s_{22}}} & \frac{x_{22} - \bar{x}_2}{\sqrt{s_{22}}} & \dots & \frac{x_{2n} - \bar{x}_2}{\sqrt{s_{22}}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{x_{p1} - \bar{x}_p}{\sqrt{s_{pp}}} & \frac{x_{p2} - \bar{x}_p}{\sqrt{s_{pp}}} & \dots & \frac{x_{pn} - \bar{x}_p}{\sqrt{s_{pp}}} \end{bmatrix} \quad (26-8)$$

بردار میانگین نمونه $(24-3)$ را ملاحظه کنید،

$$\bar{\mathbf{z}} = \frac{1}{n} \mathbf{Z} \mathbf{1} = \frac{1}{n} \begin{bmatrix} \sum_{j=1}^n \frac{x_{1j} - \bar{x}_1}{\sqrt{s_{11}}} \\ \sum_{j=1}^n \frac{x_{2j} - \bar{x}_2}{\sqrt{s_{22}}} \\ \vdots \\ \sum_{j=1}^n \frac{x_{pj} - \bar{x}_p}{\sqrt{s_{pp}}} \end{bmatrix} = \mathbf{0} \quad (27-8)$$

و ماتریس کوواریانس نمونه $(27-3)$ را ملاحظه نمایید،

$$\mathbf{S}_z = \frac{1}{n-1} \left(\mathbf{Z} - \frac{1}{n} \mathbf{Z} \mathbf{1} \mathbf{1}' \right) \left(\mathbf{Z} - \frac{1}{n} \mathbf{Z} \mathbf{1} \mathbf{1}' \right)' = \frac{1}{n-1} (\mathbf{Z} - \bar{\mathbf{z}} \mathbf{1}') (\mathbf{Z} - \bar{\mathbf{z}} \mathbf{1}')'$$

$$= \frac{1}{n-1} \mathbf{Z} \mathbf{Z}'$$

$$= \frac{1}{n-1} \begin{bmatrix} \frac{(n-1)s_{11}}{s_{11}} & \frac{(n-1)s_{12}}{\sqrt{s_{11}}\sqrt{s_{22}}} & \dots & \frac{(n-1)s_{1p}}{\sqrt{s_{11}}\sqrt{s_{pp}}} \\ \frac{(n-1)s_{12}}{\sqrt{s_{11}}\sqrt{s_{22}}} & \frac{(n-1)s_{22}}{s_{22}} & \dots & \frac{(n-1)s_{2p}}{\sqrt{s_{22}}\sqrt{s_{pp}}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{(n-1)s_{1p}}{\sqrt{s_{11}}\sqrt{s_{pp}}} & \frac{(n-1)s_{2p}}{\sqrt{s_{22}}\sqrt{s_{pp}}} & \dots & \frac{(n-1)s_{pp}}{s_{pp}} \end{bmatrix} = \mathbf{R} \quad (28-8)$$

را می دهد . مؤلفه های اصلی نمونه مشاهدات استاندارد شده به وسیله (۸-۲۰) که در آن ماتریس R جانشین S شده است داده می شود . چون مشاهدات را بیشتر «متمرکز» نموده ایم ، لذا نوشتن مؤلفه ها به صورت (۸-۲۱) لزومی ندارد .

اگر z_1, z_2, \dots, z_p مشاهدات استاندارد شده با ماتریس کوواریانس R باشند ، مؤلفه اصلی نمونه i ام ، عبارت است از :

$$\hat{y}_i = \hat{e}_i' z = \hat{e}_{1i} z_1 + \hat{e}_{2i} z_2 + \dots + \hat{e}_{pi} z_p, \quad i = 1, 2, \dots, p$$

که در آن $(\hat{\lambda}_i, \hat{e}_i)$ زوج مقدار ویژه- بردار ویژه i ام R با $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$ است . همچنین

$$(\hat{y}_i) \text{ واریانس نمونه} = \hat{\lambda}_i, \quad i = 1, 2, \dots, p$$

$$(\hat{y}_i, \hat{y}_k) \text{ کوواریانس نمونه} = 0, \quad i \neq k$$

علاوه بر این

$$\text{کل واریانس نمونه} = \text{tr}(R) = p = \hat{\lambda}_1 + \hat{\lambda}_2 + \dots + \hat{\lambda}_p$$

(استاندارد شده)

و

$$r_{\hat{y}_i, z_k} = \hat{e}_{ki} \sqrt{\hat{\lambda}_i}, \quad i, k = 1, 2, \dots, p$$

با استفاده از (۸-۲۹) نسبت کل واریانس نمونه بیان شده با مؤلفه اصلی نمونه i ام ، عبارت است از :

$$\left(\begin{array}{c} \text{سهم واریانس نمونه (استاندارد شده)} \\ \text{مربوط به مؤلفه اصلی نمونه} \\ \text{م} i \end{array} \right) = \frac{\hat{\lambda}_i}{p} \quad i = 1, 2, \dots, p \quad (۸-۳۰)$$

یک قانون سرانگشتی پیشنهاد می کند که فقط آن مؤلفه هایی را که واریانس $\hat{\lambda}_i$ آنها بزرگتر از واحد است یا معادل با آن فقط مؤلفه هایی را که به تنهایی حداقل نسبت $\frac{1}{p}$ واریانس کل را بیان می کنند ، نگاه داریم . با این وجود این قانون را نظریه زیاد تأیید نمی کند و آن را نباید بدون دلیل به کار برد .

مثال ۸-۵

نرخهای هفتگی سوددهی پنج سهم (الاید کمیکال ، دوپونت ، یونیون کارباید ، اکسون و تکزاکو) در بازار بورس نیویورک که برای دوره ژانویه ۱۹۷۵ تا دسامبر ۱۹۷۶ تعیین شده است را ثبت نموده ایم . نرخهای هفتگی سوددهی را به صورت (قیمت سهام این جمعه - قیمت سهام جمعه

قبل) / (قیمت سهام جمعه قبل) برای خرد کردن و تقسیم سهام تعدیل می شوند .
 داده ها در جدول (۸-۱) مربوط به تمرینها ثبت شده است . مشاهدات درصد هفته متوالی
 به نظر می رسد مستقلاً توزیع شده اند ، ولی نرخهای سوددهی در سرتاسر سهام همبسته اند ، زیرا
 انتظار می رود سهام در پاسخگویی به کل شرایط اقتصادی با یکدیگر تغییر کنند .
 فرض کنید x_1, x_2, \dots, x_5 به ترتیب نرخهای هفتگی سوددهی مشاهده شده شرکت‌های
 الاید کیمیکال ، دوپونت ، یونیدن کارباید ، اکسون و تکزاکو باشد . در این صورت

$$\bar{x}' = [.0054, .0048, .0057, .0063, .0037]$$

و

$$R = \begin{bmatrix} 1.000 & .577 & .509 & .387 & .462 \\ .577 & 1.000 & .599 & .389 & .322 \\ .509 & .599 & 1.000 & .436 & .426 \\ .387 & .389 & .436 & 1.000 & .523 \\ .462 & .322 & .426 & .523 & 1.000 \end{bmatrix}$$

توجه می کنیم که R ماتریس کوواریانس مشاهدات استاندارد شده

$$z_1 = \frac{x_1 - \bar{x}_1}{\sqrt{s_{11}}}, z_2 = \frac{x_2 - \bar{x}_2}{\sqrt{s_{22}}}, \dots, z_5 = \frac{x_5 - \bar{x}_5}{\sqrt{s_{55}}}$$

است . مقادیر ویژه و بردارهای ویژه نرمال شده متناظر R با رایانه محاسبه شده و در زیر داده می شود :

$$\hat{\lambda}_1 = 2.857, \quad \hat{e}'_1 = [.464, .457, .470, .421, .421]$$

$$\hat{\lambda}_2 = .809, \quad \hat{e}'_2 = [.240, .509, .260, -.526, -.582]$$

$$\hat{\lambda}_3 = .540, \quad \hat{e}'_3 = [-.612, .178, .335, .541, -.435]$$

$$\hat{\lambda}_4 = .452, \quad \hat{e}'_4 = [.387, .206, -.662, .472, -.382]$$

$$\hat{\lambda}_5 = .343, \quad \hat{e}'_5 = [-.451, .676, -.400, -.176, .385]$$

با استفاده از متغیرهای استاندارد شده دو مؤلفه اصلی اول نمونه را به دست می آوریم :

$$\hat{y}_1 = \hat{e}'_1 z = .464z_1 + .457z_2 + .470z_3 + .421z_4 + .421z_5$$

$$\hat{y}_2 = \hat{e}'_2 z = .240z_1 + .509z_2 + .260z_3 - .526z_4 - .582z_5$$

این مؤلفه ها که

$$\left(\frac{\hat{\lambda}_1 + \hat{\lambda}_2}{p} \right) 100\% = \left(\frac{2.857 + .809}{5} \right) 100\% = 73\%$$

واریانس نمونه (استاندارد شده) کل را حاصل می‌ند، دارای تعابیر جالبی است. اولین مؤلفه یک مجموع وزن دار شده با وزنهای (تقریباً) مساوی یا «شاخص» پنج سهم است. این مؤلفه را یک مؤلفه بازار-بورس کلی یا به طور خلاصه یک مؤلفه بازار می‌نامند. (در حقیقت این پنج سهم در متوسط صنعتی دو جونز منظور می‌شود).

دومین مؤلفه یک مقایسه بین سهام کمیکال (الاید کمیکال، دوپونت، و یونیون کارباید) و سهام مربوط به نفت (اکسون و تکزاکو) را می‌دهد که آن را یک مؤلفه صنعتی می‌نامند. از این رو می‌بینیم که بیشتر تغییرات در این سوددهی سهام به خاطر فعالیت بازار و فعالیت ناهمبسته صنعت است. این تغییر رفتار قیمت اوراق توسط کینگ [۱۱] نیز پیشنهاد شده است.

تعبیر بقیه مؤلفه‌ها آسان نیست و روی هم رفته تغییری را نشان می‌دهند که احتمالاً مختص هریک از سهام است. به هر تقدیر آنها بخش زیادی از واریانس نمونه کل را بیان نمی‌کنند. این مثال موردی را که نگاهداشتن یک مؤلفه (۲) مربوط به یک مقدار ویژه کمتر از ۱ معقول به نظر می‌رسد را نشان می‌دهد.

مثال ۸-۶

علمای ژنتیک اغلب با خصیصه‌های ارثی که آنها را در طول عمر حیوانات می‌توان چندین بار اندازه‌گیری نمود، سروکار دارند. وزن $n = 150$ موش ماده را (به گرم) بلافاصله بعد از تولد چهار بچه موش اول آنها به دست آوردیم. بردار میانگین نمونه و ماتریس همبستگی نمونه، عبارت است از:

$$\bar{x}' = [39.88, 45.08, 48.11, 49.95]$$

$$R = \begin{bmatrix} 1.000 & .7501 & .6329 & .6363 \\ .7501 & 1.000 & .6925 & .7386 \\ .6329 & .6925 & 1.000 & .6625 \\ .6363 & .7386 & .6625 & 1.000 \end{bmatrix}$$

مقادیر ویژه این ماتریس، عبارتند از:

$$\hat{\lambda}_1 = 3.058, \quad \hat{\lambda}_2 = .382, \quad \hat{\lambda}_3 = .342, \quad \text{و} \quad \hat{\lambda}_4 = .217$$

توجه می‌کنیم که اولین مقدار ویژه تقریباً مساوی $1 + (p-1)\bar{r} = 1 + (4-1)(.6854) = 3.056$ است. گرچه $\hat{\lambda}_4$ تا اندازه‌ای کمتر از $\hat{\lambda}_2$ و $\hat{\lambda}_3$ است ولی مقادیر ویژه باقی مانده کوچک و تقریباً مساویند. از این رو گواهی داریم که ماتریس همبستگی جامعه متناظر p ممکن است از نظر «برابری همبستگیها» شکلی مانند (۱۸-۱۵) را داشته باشد. در مورد این نماد در مثال (۸-۹) بیشتر تفحص می‌کنیم.

اولین مؤلفه اصلی

$$\hat{y}_1 = \hat{e}'_1 z = .49z_1 + .52z_2 + .49z_3 + .50z_4$$

تولدهای بعدی در طول زمان افزایش می یابد ولی تغییر در وزنها تقریباً به وسیله مؤلفه اصلی اول با ضرایب (تقریباً) مساوی به خوبی بیان می شود .

تذکر . یک مقدار کم غیر معمول برای مقدار ویژه آخر ماتریس کوواریانس یا ماتریس همبستگی نمونه می تواند یک وابستگی خطی در مجموعه داده ها که ما متوجه آن نبوده ایم را نشان دهد . اگر این امر روی می دهد یکی (یا بیش از یکی) از متغیرها زائد است و باید حذف شود . وضعیتی را که x_1 ، x_2 و x_3 نمرات زیر آزمون و نمره کل x_4 مجموع $x_1 + x_2 + x_3$ است را در نظر می گیریم . در این صورت گرچه ترکیب خطی $e'x = [1, 1, 1, -1]x = x_1 + x_2 + x_3 - x_4$ همیشه صفر است ولی خطای حاصل از گرد کردن محاسبه مقادیر ویژه ، ممکن است منجر به یک مقدار مخالف صفر کوچک گردد . اگر عبارت خطی که x_4 را به (x_1, x_2, x_3) مربوط می کند ، از اول نادیده گرفته می شد در آن صورت کوچکترین زوج مقدار ویژه - بردار ویژه باید سرنخی از وجودش را به ما می داد .

از این رو مقادیر ویژه «بزرگ» و بردارهای ویژه مربوط به آن در تحلیل یک مؤلفه اصلی با اهمیت اند و مقادیر ویژه بسیار نزدیک به صفر را نباید معمولاً نادیده گرفت . بردارهای ویژه مربوط به این مقادیر ویژه اخیر ممکن است به نابسنجیهای خطی در مجموعه داده ها که می تواند مسائل تعبیری و محاسبه ای را در یک تحلیل بعدی به وجود آورد ، اشاره کند .

۴-۸ نمودار مؤلفه های اصلی

نمودار مؤلفه های اصلی می تواند مشاهدات مورد شک را آشکار نماید و بررسیهای فرض نرمال را فراهم کند . چون مؤلفه های اصلی ترکیبات خطی متغیرهای اولیه هستند ، لذا تقریباً نرمال بودن آنها نامعقول نیست . وقتی از مؤلفه های اصلی به عنوان داده های ورودی برای تحلیلهای دیگر استفاده می کنیم ، اغلب لازم است توزیع چند مؤلفه اصلی اول را نرمال در نظر بگیریم .

مؤلفه های اصلی آخر می توانند در مورد مشاهدات مشکوک به ما کمک کنند . هر مشاهده x_j را

می توان به صورت یک ترکیب خطی

$$\begin{aligned} x_j &= (x'_j \hat{e}_1) \hat{e}_1 + (x'_j \hat{e}_2) \hat{e}_2 + \dots + (x'_j \hat{e}_p) \hat{e}_p \\ &= \hat{y}_{1j} \hat{e}_1 + \hat{y}_{2j} \hat{e}_2 + \dots + \hat{y}_{pj} \hat{e}_p \end{aligned}$$

روشهای تشخیصی مؤلفه های اصلی را به خوبی می توان برای بررسی فرضهای مربوط به الگوی رگرسیون چندگانه چندمتغیری به کار برد. در حقیقت اگر هر الگو را با هر روش برآورد برآزش کنیم در نظر گرفتن

(بردار مقادیر تخمین زده شده (برآورد شده)) - (بردار مشاهده) = بردار باقی مانده

یا

$$\hat{\varepsilon}_j = \mathbf{y}_j - \mathbf{z}'_j \hat{\boldsymbol{\beta}}, \quad j = 1, 2, \dots, n \quad (31-8)$$

برای الگوی خطی چندمتغیری عاقلانه است. مؤلفه های اصلی به دست آمده از ماتریس باقی مانده

$$\frac{\sum_{j=1}^n (\hat{\varepsilon}_j - \bar{\hat{\varepsilon}}_j)(\hat{\varepsilon}_j - \bar{\hat{\varepsilon}}_j)'}{n - p} \quad (32-8)$$

را به همان شکلی که از یک نمونه تصادفی به دست آمدند، مورد بررسی قرار داد. شما باید آگاه باشید که وابستگیهای خطی در میان باقی مانده های حاصل از یک تحلیل رگرسیون خطی وجود دارند، لذا مقادیر ویژه آخر در بین خطاهای ناشی از گرد کردن صفر خواهد بود.

۵-۸ استنباطهای مبتنی بر نمونه های بزرگ

دیدیم که مقادیر ویژه و بردارهای ویژه ماتریس کوواریانس (همبستگی) جوهره یک تحلیل مؤلفه های اصلی است. بردارهای ویژه جهات بیشترین تغییرات و مقادیر ویژه واریانسها را مشخص می کند. هنگامی که چند مقدار ویژه اول خیلی بزرگتر از بقیه هستند، بیشتر واریانس کل را می توان در کمتر از «بعد» بیان کرد.

در عمل تصمیمهایی در مورد کیفیت تقریب مؤلفه اصلی بایستی بر پایه زوجهای مقدار ویژه. بردار ویژه به دست آمده از S یا R ساخته شود. این مقادیر ویژه و بردارهای ویژه به علت تغییرات نمونه گیری با مقادیر ویژه و بردارهای ویژه جامعه مورد بررسی تفاوت می کنند. به دست آوردن توزیعهای نمونه گیری $\hat{\lambda}_i$ و $\hat{\varepsilon}_i$ مشکل و خارج از بحث این کتاب است. اگر علاقه مند باشید می توانید برخی از اثباتها را برای توزیعهای نرمال چندمتغیری در [۱] و [۲] و [۴] ببینید. نتایج نمونه های بزرگ مناسب را به آسانی خلاصه می کنیم.

خواص نمونه های بزرگ $\hat{\lambda}_i$ و $\hat{\varepsilon}_i$

در نتایج مربوط به فواصل اطمینان نمونه های بزرگ برای $\hat{\lambda}_i$ و $\hat{\varepsilon}_i$ که در حال حاضر

در دسترس هستند، X_1, X_2, \dots, X_n یک نمونه تصادفی از یک جامعه نرمال در نظر گرفته می شود. همچنین باید فرض کنیم که مقادیر ویژه (نامعلوم) Σ متمایز و مثبت اند، به طوری که $\lambda_1 > \lambda_2 > \dots > \lambda_p > 0$. حالتی که تعداد مقادیر ویژه مساوی معلوم اند یک استثناست. معمولاً نتایج مربوط به مقادیر ویژه متمایز به کار برده می شود، مگر این که یک دلیل قوی وجود داشته باشد که باور کنیم Σ برای حصول مقادیر ویژه مساوی دارای ساختمان ویژه ای است. حتی وقتی از فرض نرمال عدول می شود، فواصل اطمینانی که به این طریق به دست می آید، هنوز دلیلی بر عدم قطعیت $\hat{\lambda}_i$ و \hat{e}_i فراهم می کند.

اندرسن [۲] و گیرشیک [۴] نظریه توزیع نمونه های بزرگ زیر را برای مقادیر ویژه $\hat{\lambda}' = [\hat{\lambda}_1, \dots, \hat{\lambda}_p]$ و بردارهای ویژه $\hat{e}_1, \dots, \hat{e}_p$ ماتریس S ثابت می کنند.

۱- فرض کنید Λ ماتریس قطری مقادیر ویژه $\lambda_1, \dots, \lambda_p$ ، Σ باشد در آن صورت $\sqrt{n}(\hat{\lambda} - \lambda)$ تقریباً $N_p(0, 2\Lambda^2)$ است.

۲- فرض کنید

$$E_i = \lambda_i \sum_{\substack{k=1 \\ k \neq i}}^p \frac{\lambda_k}{(\lambda_k - \lambda_i)^2} e_k e_k'$$

در این صورت $\sqrt{n}(\hat{e}_i - e_i)$ تقریباً $N_p(0, E_i)$ است.

۳- هر $\hat{\lambda}_i$ مستقل از اعضای مربوط به \hat{e}_i توزیع می شود.

از نتیجه ۱ معلوم می شود که برای n بزرگ، $\hat{\lambda}_i$ ها مستقلاً توزیع می شوند. علاوه بر این $\hat{\lambda}_i$

دارای یک توزیع تقریبی $N(\lambda_i, \frac{2\lambda_i^2}{n})$ است. با استفاده از این توزیع نرمال، داریم:

$P[|\hat{\lambda}_i - \lambda_i| \leq z(\alpha/2)\lambda_i\sqrt{2/n}] = 1 - \alpha$ به این ترتیب یک فاصله اطمینان $(1 - \alpha)100\%$ نمونه بزرگ برای λ_i ، عبارت است از:

$$\frac{\hat{\lambda}_i}{(1 + z(\alpha/2)\sqrt{2/n})} \leq \lambda_i \leq \frac{\hat{\lambda}_i}{(1 - z(\alpha/2)\sqrt{2/n})} \quad (۳۳-۸)$$

که در آن $z(\alpha/2)$ صدک $(\alpha/2)$ ۱۰۰ ام یک توزیع نرمال استاندارد است. فواصل $(1 - \alpha)100\%$ هم زمان از نوع بونفرونی برای m تا λ_i با جایگزینی $z(\alpha/2)$ با $z(\alpha/2m)$ به دست می آید (بخش ۴-۵ را ملاحظه نمایید). از نتیجه ۲ معلوم می شود که برای نمونه های بزرگ \hat{e}_i ها حول e_i های مربوطه به طور نرمال توزیع می شوند. اعضای هر \hat{e}_i همبسته بوده و همبستگی تا حد زیادی به جدایی مقادیر ویژه $\lambda_1, \lambda_2, \dots, \lambda_p$ (که نامعلوم اند) و حجم نمونه n بستگی دارد. خطاهای معیار تقریبی ضرایب \hat{e}_{ki}

با اعضای قطری $(1/n)\hat{E}_i$ داده می شود که \hat{E}_i با قرار دادن $\hat{\lambda}_i$ ها به جای λ_i ها از E_i به دست می آید.

مثال ۸-۸

با استفاده از داده های قیمت سهام در جدول (۸-۱) یک فاصله اطمینان ۹۵٪ برای λ_1 یعنی واریانس اولین مؤلفه اصلی جامعه، به دست می آوریم.

فرض می کنیم نرخهای سوددهی سهام استخراجهای مستقل از یک $N_s(\mu, \Sigma)$ که Σ معین و مثبت با مقادیر ویژه متمایز $0 < \lambda_1 < \lambda_2 < \dots < \lambda_p$ است را نشان دهد. چون $n = 100$ بزرگ است، می توانیم از (۸-۳۳) با $i = 1$ برای ساختن یک فاصله اطمینان ۹۵٪ برای λ_1 استفاده کنیم از تمرین ۱۰.۸ داریم $\hat{\lambda}_1 = .0036$ و علاوه بر این $z(.025) = 1.96$. بنابراین با اطمینان ۹۵٪ داریم:

$$.0028 \leq \lambda_1 \leq .0050 \quad \text{یا} \quad \frac{.0036}{(1 + 1.96\sqrt{1/100})} \leq \lambda_1 \leq \frac{.0036}{(1 - 1.96\sqrt{1/100})}$$

هرگاه یک مقدار ویژه مانند ۱۰۰ یا حتی ۱۰۰۰ بزرگ باشد فواصلی که به وسیله (۸-۳۳) تولید می شوند برای سطوح اطمینان معقولی می توانند کاملاً عریض باشند و لو این که n نسبتاً بزرگ باشد. به طور کلی فاصله اطمینان به همان نحی که $\hat{\lambda}_i$ بزرگتر می شود عریض تر می گردد. بنابراین باید بر پایه یک امتحان $\hat{\lambda}_i$ ها در حذف یا نگهداری مؤلفه های اصلی دقت کنیم.

آزمون مربوط به ساختار همبستگیهای مساوی

ساختار ویژه همبستگی $\text{Cov}(X_i, X_k) = \sqrt{\sigma_{ii}\sigma_{kk}}\rho$ یا $\text{Corr}(X_i, X_k) = \rho$ برای هر $i \neq k$ ساختار مهمی است که مقادیر ویژه Σ متمایز نیست و نتایج قبلی کاربرد نداشته باشند. برای آزمون این ساختار، فرض کنید:

$$H_0: \rho = \rho_0 = \begin{bmatrix} 1 & \rho & \dots & \rho \\ \rho & 1 & \dots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \dots & 1 \end{bmatrix}$$

و

$$H_1: \rho \neq \rho_0$$

یک آزمون H_{11} در مقابل H_1 می تواند بر پایه یک آماره نسبت درست نمایی قرار داشته باشد، ولی لاولی [۱۲] نشان می دهد که می توان یک روش آزمون معادلی را از اعضای غیر قطری R ساخت. در روش لاولی کمیتهای زیر لازم است:

$$\bar{r}_k = \frac{1}{p-1} \sum_{\substack{i=1 \\ i \neq k}}^p r_{ik} \quad k = 1, 2, \dots, p; \quad \bar{r} = \frac{2}{p(p-1)} \sum_{i < k} r_{ik} \quad (34-8)$$

$$\hat{\gamma} = \frac{(p-1)^2 [1 - (1 - \bar{r})^2]}{p - (p-2)(1 - \bar{r})^2}$$

واضح است که \bar{r}_k متوسط اعضای غیرقطری در ستون (یا سطر) k ام R و \bar{r} متوسط تمام اعضای غیرقطری است.

آزمون تقریبی نمونه های بزرگ در سطح α دارای شکل زیر است: H_0 را به نفع H_1 رد می کنیم هرگاه

$$T = \frac{(n-1)}{(1-\bar{r})^2} \left[\sum_{i < k} (r_{ik} - \bar{r})^2 - \hat{\gamma} \sum_{k=1}^p (\bar{r}_k - \bar{r})^2 \right] > \chi_{(p+1)(p-2)/2}^2(\alpha) \quad (35-8)$$

که در آن $\chi_{(p+1)(p-2)/2}^2(\alpha)$ صدک (100α) ام بالایی یک توزیع کی دو با $\frac{(p+1)(p-2)}{2}$ درجه آزادی است.

مثال ۸-۹

ماتریس همبستگی نمونه ساخته شده از وزنهای تولد بعدی موشهای ماده در مثال (۸-۶) به شرح زیر است:

$$R = \begin{bmatrix} 1.0 & .7501 & .6329 & .6363 \\ .7501 & 1.0 & .6925 & .7386 \\ .6329 & .6925 & 1.0 & .6625 \\ .6363 & .7386 & .6625 & 1.0 \end{bmatrix}$$

ما از این ماتریس همبستگی برای تشریح آزمون نمونه بزرگ در (۳۵-۸) استفاده خواهیم کرد.

در این جا $p = 4$ بوده و قرار می دهیم:

$$H_0: \rho = \rho_0 = \begin{bmatrix} 1 & \rho & \rho & \rho \\ \rho & 1 & \rho & \rho \\ \rho & \rho & 1 & \rho \\ \rho & \rho & \rho & 1 \end{bmatrix}$$

$$H_1: \rho \neq \rho_0$$

با استفاده از (۳۴-۸) و (۳۵-۸)، داریم:

تحلیل آماری چندمتغیری کاربردی

$$\bar{r}_1 = \frac{1}{3}(.7501 + .6329 + .6363) = .6731, \quad \bar{r}_2 = .7271,$$

$$\bar{r}_3 = .6626, \quad \bar{r}_4 = .6791$$

$$\bar{r} = \frac{2}{4(3)}(.7501 + .6329 + .6363 + .6925 + .7386 + .6625) = .6855$$

$$\sum_{i < k} \sum (r_{ik} - \bar{r})^2 = (.7501 - .6855)^2 + (.6329 - .6855)^2 + \dots + (.6625 - .6855)^2 = .01277$$

$$\sum_{k=1}^4 (\bar{r}_k - \bar{r})^2 = (.6731 - .6855)^2 + \dots + (.6791 - .6855)^2 = .00245$$

$$\hat{\gamma} = \frac{(4 - 1)^2 [1 - (1 - .6855)^2]}{4 - (4 - 2)(1 - .6855)^2} = 2.1329$$

و

$$T = \frac{(150 - 1)}{(1 - .6855)^2} [.01277 - (2.1329)(.00245)] = 11.4$$

چون $5 = 5(2)/2 = (p+1)(p-2)/2$ لذا مقدار بحرانی 5% برای آزمون در (۸-۳۵) $\chi^2_5(0.05) = 11.07$ است. مقدار آماره آزمون ما تقریباً مساوی نقطه بحرانی 5% نمونه بزرگ است، لذا شواهد علیه H_0 (همبستگیهای مساوی) قوی بوده ولی بیش از حد اندازه نیست. چنان که در مثال (۸-۶) دیدیم، کوچکترین مقادیر ویژه $\hat{\lambda}_2$ ، $\hat{\lambda}_3$ و $\hat{\lambda}_4$ قدری با هم تفاوت دارند و $\hat{\lambda}_4$ تا اندازه ای کمتر از دوتای دیگر است. در نتیجه با حجم نمونه بزرگی که در این مسأله است تفاوت های کم از ساختار همبستگیهای مساوی یک معنی داری آماری را نشان می دهند.

تمرینها

۱.۸ مؤلفه های اصلی Y_1 و Y_2 جامعه را برای ماتریس کوواریانس

$$\Sigma = \begin{bmatrix} 5 & 2 \\ 2 & 2 \end{bmatrix} \quad \text{SAS}$$

تعیین کنید همچنین نسبت واریانس کل جامعه بیان شده با اولین مؤلفه اصلی را محاسبه کنید.

۲.۸ ماتریس کوواریانس در تمرین (۱.۸) را به ماتریس همبستگی ρ تبدیل کنید.

(الف) مؤلفه های اصلی Y_1 و Y_2 را از ρ تعیین کرده و نسبت واریانس کل جامعه بیان شده با Y_1 را محاسبه کنید.

(ب) مؤلفه ها را با مؤلفه های اصلی به دست آمده در تمرین (۱.۸) مقایسه کنید. آیا

این دو یکسانند؟ آیا باید یکسان باشند؟

(ج) همبستگیهای ρ_{Y_1, Z_1} ، ρ_{Y_1, Z_2} و ρ_{Y_2, Z_1} را محاسبه کنید.

۳.۸ فرض کنید:

$$\Sigma = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 4 \end{bmatrix}$$

مؤلفه های اصلی Y_1 و Y_2 و Y_3 را تعیین کنید. نظر شما در مورد بردارهای ویژه (و مؤلفه های

اصلی) مربوط به مقادیر ویژه ای که متمایز نیستند چیست؟

۴.۸ مؤلفه های اصلی و نسبت واریانس کل جامعه بیان شده توسط هر یک را وقتی ماتریس کوواریانس

$$\Sigma = \begin{bmatrix} \sigma^2 & \sigma^2 \rho & 0 \\ \sigma^2 \rho & \sigma^2 & \sigma^2 \rho \\ 0 & \sigma^2 \rho & \sigma^2 \end{bmatrix}, \quad -\frac{1}{\sqrt{2}} < \rho < \frac{1}{\sqrt{2}}$$

پیدا کنید.

۵.۸ (الف) مقادیر ویژه ماتریس همبستگی زیر را پیدا کنید،

$$\rho = \begin{bmatrix} 1 & \rho & \rho \\ \rho & 1 & \rho \\ \rho & \rho & 1 \end{bmatrix}$$

آیا نتایج شما با (۸-۱۶) و (۸-۱۷) سازگار است؟

(ب) زوجهای مقدار ویژه-بردار ویژه را برای ماتریس $p \times p$ ، ρ داده شده در (۸-۱۵) به دست آورید.

۶.۸ داده های مربوط به دارایی x_1 و درآمد خالص x_2 برای ده تا از بزرگترین شرکتهای صنعتی ایالات متحده در تمرین (۴.۱) فصل ۱ ثبت شده است .
از مثال (۴-۱۱) ، داریم :

$$\bar{x} = \begin{bmatrix} 19.32 \\ 1.51 \end{bmatrix}, \quad S = \begin{bmatrix} 70.41 & 5.87 \\ 5.87 & .97 \end{bmatrix}$$

(الف) مؤلفه های اصلی نمونه و واریانس آنها را برای این داده ها تعیین کنید (شما ممکن است برای به دست آوردن مقادیر ویژه S یک فرمول درجه دوم لازم داشته باشید) .

(ب) نسبت واریانس کل نمونه بیان شده با λ_1 را پیدا کنید .

(ج) بیضی با چگالی ثابت $(x - \bar{x})' S^{-1} (x - \bar{x}) = 1.4$ را رسم کرده و مؤلفه های اصلی λ_1 و λ_2 را روی نمودار نشان دهید .

(د) ضرایب همبستگی r_{λ_1, x_k} ، $k = 1, 2$ را محاسبه کنید . چه تعبیری از اولین مؤلفه اصلی ارائه می کنید ؟

۷.۸ ماتریس کوواریانس S را در تمرین (۶.۸) به یک ماتریس همبستگی نمونه R تبدیل کنید .

(الف) مؤلفه های اصلی نمونه λ_1 و λ_2 و واریانسهایشان را پیدا کنید .

(ب) نسبت واریانس کل نمونه بیان شده با λ_1 را محاسبه کنید .

(ج) ضرایب همبستگی r_{λ_1, x_k} ، $k = 1, 2$ را محاسبه نموده و λ_1 را تعبیر کنید .

(د) مؤلفه های به دست آمده در بخش الف را با مؤلفه های به دست آمده در تمرین (۶.۸) الف

مقایسه کنید . داده های اولیه نمایش داده شده در تمرین (۴.۱) معلوم اند، آیا احساس

می کنید بهتر است مؤلفه های اصلی را از ماتریس کوواریانس نمونه یا ماتریس همبستگی

نمونه تعیین کنیم ؟ بیان کنید .

۸.۸ با استفاده از نتایج مثال (۵-۸)

(الف) همبستگیهای r_{λ_i, z_k} ، را برای $i = 1, 2$ ، $k = 1, 2, \dots, 5$ محاسبه کنید . آیا این

همبستگیها تعابیر ارائه شده برای دو مؤلفه اول را تقویت می کند ؟

(ب) فرض

$$H_0: \rho = \rho_0 = \begin{bmatrix} 1 & \rho & \rho & \rho & \rho \\ \rho & 1 & \rho & \rho & \rho \\ \rho & \rho & 1 & \rho & \rho \\ \rho & \rho & \rho & 1 & \rho \\ \rho & \rho & \rho & \rho & 1 \end{bmatrix}$$