

Chapter 3

Sample Geometry and Random Sampling

$$X = \begin{bmatrix} x_{11} & \dots & x_{1p} \\ x_{21} & \dots & x_{2p} \\ \vdots & \vdots & \vdots \\ x_{n1} & \dots & x_{np} \end{bmatrix} = \begin{bmatrix} x_1' \\ x_2' \\ \vdots \\ x_n' \end{bmatrix}$$

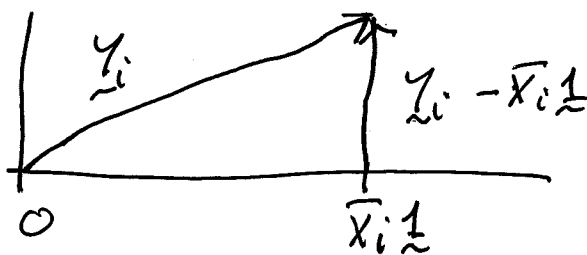
n vectors -
p-dim space
 $x_i' = [x_{i1} \dots x_{ip}]$

$$\bar{x} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix}$$

Computing the mean vector: $\bar{x} = \frac{1}{n} X' \mathbf{1}_n$

$X = (\underline{y}_1; \underline{y}_2; \dots; \underline{y}_p)$ p vectors in n-dim space.

$$y_i' \left(\frac{1}{\sqrt{n}} \mathbf{1} \right) \left(\frac{1}{\sqrt{n}} \mathbf{1} \right) = \frac{\sum_{j=1}^n x_{ji}}{n} \mathbf{1} = \bar{x}_i \mathbf{1}$$



$$\underline{d}_i = \underline{y}_i - \bar{x}_i \mathbf{1} = \begin{bmatrix} x_{i1} - \bar{x}_1 \\ x_{i2} - \bar{x}_2 \\ \vdots \\ x_{ip} - \bar{x}_p \end{bmatrix} \quad (i=1, \dots, p)$$

Note: Book uses different D

Let $D = [d_1; d_2; \dots; d_p]$

$$D'D = \begin{bmatrix} \sum_j (x_{j1} - \bar{x}_1)^2 & \sum_j (x_{j1} - \bar{x}_1)(x_{j2} - \bar{x}_2) & \dots & \sum_j (x_{j1} - \bar{x}_1)(x_{jp} - \bar{x}_p) \\ \vdots & \vdots & \ddots & \vdots \\ \sum_j (x_{j1} - \bar{x}_1)(x_{jp} - \bar{x}_p) & \sum_j (x_{j2} - \bar{x}_2)(x_{jp} - \bar{x}_p) & \dots & \sum_j (x_{jp} - \bar{x}_p)^2 \end{bmatrix}$$

$$S_n = \frac{1}{n} D'D \quad V^{1/2} = \text{sqrt}(\text{diag}(S_n)) \quad \underline{3.2}$$

$$R = V^{-1/2} S_n V^{-1/2}$$

(Equation 2.6)

$$d_i'd_k = \sum_j (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k) = L_{d_i} L_{d_k} \cos \theta_{ik}$$

$$\Rightarrow \sum_j (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k) = \sqrt{\sum_j (x_{ji} - \bar{x}_i)^2} \sqrt{\sum_j (x_{jk} - \bar{x}_k)^2} \cos(\theta_{ik})$$

$$\Rightarrow r_{ik} = \frac{S_{ik}}{\sqrt{S_{ii}} \sqrt{S_{kk}}} = \cos \theta_{ik}$$

θ_{ik} is angle
formed by
vectors $\underline{d}_i, \underline{d}_k$

Geometrical Interpretation of the Sample

(column of X)

1) Projection of y_i onto vector $\underline{1}$ is $\bar{x}_i \underline{1}$

$$\text{Length of } \bar{x}_i \underline{1} = \sqrt{\bar{x}_i \underline{1}' \underline{1} \bar{x}_i} = |\bar{x}_i| \sqrt{n}$$

2) $S_n = \frac{1}{n} D'D \quad D = [d_i; d_k; \dots; d_p]$ (Deviations)
 $d_i = y_i - \bar{x}_i \underline{1}$ $(L_{d_i})^2 = n S_{ii}$ $(L_{d_i} L_{d_k}) = n S_{ik} = d_i'd_k$

3) $r_{ik} \equiv \cos$ of angle between $\underline{d}_i, \underline{d}_k$.

3.3 Random Samples and $E\{\bar{X}\}$ and $E\{S_n^2\}$

$$\textcircled{1} X = \begin{bmatrix} X_{11} & \dots & X_{1p} \\ \vdots & \ddots & \vdots \\ X_{n1} & \dots & X_{np} \end{bmatrix} = \begin{bmatrix} X_1' \\ \vdots \\ X_n' \end{bmatrix} \quad X_i' = [X_{i1} \dots X_{ip}]$$

Row vectors independent observations from common

$f(\underline{x}) = f(x_1, \dots, x_p) \Rightarrow \underline{x}_1, \dots, \underline{x}_n$ are random sample
from $f(\underline{x})$

\Rightarrow joint density is $f(\underline{x}_1) f(\underline{x}_2) \dots f(\underline{x}_n)$ $f(\underline{x}_i) = f(x_{i1}, \dots, x_{ip})$

1) measurements of observation in a single trial

$X_i' = [X_{i1}, \dots, X_{ip}]$ typically correlated. Measurements
across trials must be independent.

2) When measurements are observed over time

(stock prices or economic indicators), independence
may not hold across trials.

Result 3.1 X_1, \dots, X_n = random sample from joint distribution with mean $\underline{\mu}_X = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_p \end{bmatrix}$

and Covariance matrix $\Sigma_X = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{12} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{bmatrix}$

$$\bar{X} = \frac{1}{n} [X_1 + X_2 + \dots + X_n] = \begin{bmatrix} \bar{X}_1 \\ \bar{X}_2 \\ \vdots \\ \bar{X}_p \end{bmatrix}$$

$$E\{\bar{X}\} = \frac{1}{n} [\underline{\mu}_X + \underline{\mu}_X + \dots + \underline{\mu}_X] = \frac{1}{n} n \underline{\mu}_X = \underline{\mu}_X$$

$$\begin{aligned} \text{Cov}\{\bar{X}\} &= \text{Cov}\left\{\frac{1}{n} [X_1 + X_2 + \dots + X_n]\right\} \\ &= \frac{1}{n^2} [V\{X_1\} + \dots + V\{X_n\} + 2\text{Cov}\{X_1, X_2\} + \dots + 2 \end{aligned}$$

$$\text{Cov}\{\bar{X}\} = \frac{1}{n} \Sigma_X$$

$$E\{S_n\} = \frac{n-1}{n} \Sigma_X = \Sigma_X - \frac{1}{n} \Sigma_X$$

$$\Rightarrow E\left\{\frac{n}{n-1} S_n\right\} = \Sigma_X$$

$$\bar{x} - \mu = \frac{1}{n} [(x_1 - \mu) + \dots + (x_n - \mu)] = \frac{1}{n} \sum_{j=1}^n (x_j - \mu)$$

$$\Rightarrow (\bar{x} - \mu)(\bar{x} - \mu)' = \left[\frac{1}{n} \sum_{j=1}^n (x_j - \mu) \right] \left[\frac{1}{n} \sum_{l=1}^n (x_l - \mu) \right]'$$

$$= \frac{1}{n^2} \sum_{j=1}^n \sum_{l=1}^n (x_j - \mu)(x_l - \mu)'$$

Note: if $j \neq l$, $E\{(x_j - \mu)(x_l - \mu)'\} = 0$
(independence across trials)

$$E\{(x_j - \mu)(x_j - \mu)'\} = E\left\{ \begin{pmatrix} x_{j1} - \mu_1 \\ \vdots \\ x_{jp} - \mu_p \end{pmatrix} [x_{j1} - \mu_1 \dots x_{jp} - \mu_p] \right\}$$

$$= \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{12} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{1p} & \sigma_{2p} & \dots & \sigma_{pp} \end{pmatrix}$$

$$\Rightarrow E\{(\bar{x} - \mu)(\bar{x} - \mu)'\} = \frac{1}{n^2} n \Sigma_x = \frac{1}{n} \Sigma_x$$

$$S_n = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})'$$

$$= \frac{1}{n} \left[\sum_{j=1}^n (x_j - \bar{x})x_j' + \sum_{j=1}^n (x_j - \bar{x})(-\bar{x})' \right]$$

$$\Rightarrow S_n = \frac{1}{n} \left[\sum_{j=1}^n \underline{x}_j \underline{x}_j' - \bar{x} \sum_{j=1}^n \underline{x}_j' - \bar{x} \sum_{j=1}^n (\underline{x}_j - \bar{x}) \right]$$

$$= \frac{1}{n} \left[\sum_{j=1}^n \underline{x}_j \underline{x}_j' - n \bar{x} \bar{x}' \right]$$

$$\Rightarrow E\{S_n\} = \frac{1}{n} E\left\{ \sum_{j=1}^n \underline{x}_j \underline{x}_j' - n \bar{x} \bar{x}' \right\}$$

Let V be a vector with mean μ_V , $\text{Cov}\{V\} = \Sigma_V$

$$VV' = \begin{bmatrix} V_1 \\ V_2 \\ \vdots \\ V_k \end{bmatrix} [V_1 \ V_2 \ \dots \ V_k] = \begin{bmatrix} V_1^2 & V_1 V_2 & \dots & V_1 V_k \\ V_1 V_2 & V_2^2 & \dots & V_2 V_k \\ \vdots & \vdots & \ddots & \vdots \\ V_1 V_k & V_2 V_k & \dots & V_k^2 \end{bmatrix}$$

$$E\{V_i^2\} = V\{V_i\} + [E\{V_i\}]^2 = \Sigma_{V_{ii}} + \mu_i^2$$

$$V\{V_i, V_j\} = \text{Cov}\{V_i, V_j\} + \mu_i \mu_j = \Sigma_{V_{ij}} + \mu_i \mu_j$$

$$V\{V_i, V_k\} = \text{Cov}\{V_i, V_k\} + \mu_i \mu_k = \Sigma_{V_{ik}} + \mu_i \mu_k$$

$$\Rightarrow E\{VV'\} = \Sigma_V + \mu_V \mu_V'$$

$$\Rightarrow E\{x_j \underline{x}_j'\} = \Sigma_x + \mu_x \mu_x'$$

$$E\{\bar{x} \bar{x}'\} = \frac{1}{n} \Sigma_x + \mu_x \mu_x'$$

$$\Rightarrow E\{S_n\} = \frac{1}{n} \left[n(\bar{X} + \underline{\mu}_x \underline{\mu}_x') - n\left(\frac{1}{n} \bar{X} + \underline{\mu}_x \underline{\mu}_x'\right) \right]$$

$$= \frac{1}{n} [(n-1) \bar{X}] = \frac{n-1}{n} \bar{X}$$

$$\Rightarrow E\left\{\frac{n}{n-1} S_n\right\} = E\{S\} = \bar{X}$$

Unbiased Sample Variance-Covariance Matrix

$$S = \frac{n}{n-1} S_n = \frac{1}{n} \sum_{j=1}^n (\underline{x}_j - \bar{X})(\underline{x}_j - \bar{X})'$$

3.4 Generalized Variance

$$S = \begin{bmatrix} S_{11} & \dots & S_{1p} \\ \vdots & \ddots & \vdots \\ S_{1p} & \dots & S_{pp} \end{bmatrix} \quad S_{ik} = \frac{1}{n-1} \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k)$$

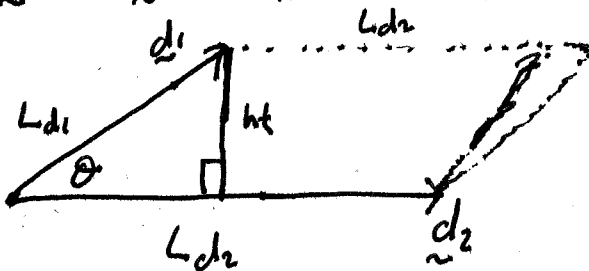
S contains p variances and $\binom{p}{2} = \frac{p(p-1)}{2}$ covariances

Generalized Sample Variance = $|S|$

Geometric Interpretation

$$\underline{d}_1 = \underline{y}_1 - \bar{x}_1 \underline{1}$$

$$\underline{d}_2 = \underline{y}_2 - \bar{x}_2 \underline{1}$$



$$L_{d1} = \sqrt{\sum_j (x_{j1} - \bar{x}_1)^2} = \sqrt{(n-1)S_{11}}$$

$$L_{d2} = \sqrt{\sum_j (x_{j2} - \bar{x}_2)^2} = \sqrt{(n-1)S_{22}}$$

$$\sin \theta = \frac{ht}{L_{d1}} \Rightarrow ht = L_{d1} \sin \theta$$

$$\text{Area of Trapezoid} = |L_{d1} \sin \theta| L_{d2} = L_{d1} L_{d2} |\sin \theta|$$

$$\textcircled{1} \cos^2 \theta + \sin^2 \theta = 1 \Rightarrow |\sin \theta| = \sqrt{1 - \cos^2 \theta}$$

$$\Rightarrow \text{Area} = L_{d1} L_{d2} \sqrt{1 - \cos^2 \theta} \quad \cos \theta = r_{12}$$

$$\Rightarrow \text{Area} = (n-1) \sqrt{S_{11}} \sqrt{S_{22}} \sqrt{1 - r_{12}^2}$$

$$= (n-1) \sqrt{S_{11} S_{22} (1 - r_{12}^2)}$$

$$S = \begin{bmatrix} S_{11} & S_{12} \\ S_{12} & S_{22} \end{bmatrix} = \begin{bmatrix} S_{11} & r_{12} \sqrt{S_{11}} \sqrt{S_{22}} \\ r_{12} \sqrt{S_{11}} \sqrt{S_{22}} & S_{22} \end{bmatrix}$$

$$|S| = S_{11} S_{22} - \left[r_{12} \sqrt{S_{11}} \sqrt{S_{22}} \right]^2$$

$$= S_{11} S_{22} - r_{12}^2 S_{11} S_{22} = \cancel{S_{11} S_{22}} \sqrt{1 - r_{12}^2}$$

$$= S_{11} S_{22} (1 - r_{12}^2) \Rightarrow |S| = \frac{\text{area}^2}{(n-1)^2}$$

Generalized to $p-1$ deviation vectors in a space

$$\text{Generalized Sample Variance} = |S| = (n-1)^{-p} (\text{volume})^2$$

p -dimensional space interpretation

$$(\underline{x} - \bar{\underline{x}})' S^{-1} (\underline{x} - \bar{\underline{x}}) = c^2 \Rightarrow \text{ellipsoid (hyperellipsoid)} \\ \text{of point } \underline{x} \text{ a distance} \\ \text{of } c \text{ from } \bar{\underline{x}}.$$

$$\text{Volume of } \{ \underline{x} : (\underline{x} - \bar{\underline{x}})' S^{-1} (\underline{x} - \bar{\underline{x}}) \leq c^2 \} = k_p |S|^{1/2} c^p$$

$$\Rightarrow \text{Volume}^2 = k_p^2 c^{2p} |S| = \text{constant} \times \text{Generalized} \\ \text{Sample Variance}$$

where: $k_p = \frac{2\pi^{p/2}}{p\Gamma(\frac{p}{2})}$

Generalized Variance from Correlation Matrix

$$\text{Generalized Sample Variance of Standardized vars} = |R|$$

$$\text{Total Sample Variance} = S_{11} + S_{22} + \dots + S_{pp}$$

3.5 Sample Mean, Covariance, Correlation in Matrix Form

$$\bar{X} = \begin{bmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_p \end{bmatrix} = \begin{bmatrix} \frac{y'_1 \mathbf{1}}{n} \\ \vdots \\ \frac{y'_p \mathbf{1}}{n} \end{bmatrix} = \frac{1}{n} X' \mathbf{1}$$

• $n \times p$ matrix of means: $\begin{bmatrix} \bar{x}_1 \mathbf{1}' \\ \bar{x}_2 \mathbf{1}' \\ \vdots \\ \bar{x}_p \mathbf{1}' \end{bmatrix}$

$$= \frac{1}{n} \bar{X}' = \frac{1}{n} \left(\frac{1}{n} X' \mathbf{1} \right)' = \frac{1}{n} \mathbf{1} \mathbf{1}' X = \frac{1}{n} J_n X$$

• $n \times p$ matrix of deviations (Residuals): $\begin{bmatrix} x_{11} - \bar{x}_1 & \dots & x_{1p} - \bar{x}_p \\ \vdots & \ddots & \vdots \\ x_{n1} - \bar{x}_1 & \dots & x_{np} - \bar{x}_p \end{bmatrix} = E$

• $X - \frac{1}{n} J_n X = (I_n - \frac{1}{n} J_n) X$ ($J_n = \mathbf{1} \mathbf{1}'$)

~~PKP matrix~~ of sums of squares and cross-products:

$$(n-1)S = \begin{bmatrix} \sum_{j=1}^n (x_{j1} - \bar{x}_1)^2 & \dots & \sum_{j=1}^n (x_{j1} - \bar{x}_1)(x_{jp} - \bar{x}_p) \\ \vdots & \ddots & \vdots \\ \sum_{j=1}^n (x_{j1} - \bar{x}_1)(x_{jp} - \bar{x}_p) & \dots & \sum_{j=1}^n (x_{jp} - \bar{x}_p)^2 \end{bmatrix}$$

$$= E'E = \left[(I - \frac{1}{n} J) \right]' (I - \frac{1}{n} J) X = X' (I - \frac{1}{n} J)' (I - \frac{1}{n} J) X$$

$(I - \frac{1}{n}J) \equiv \text{Symmetric}$

$$\begin{aligned} (I - \frac{1}{n}J)(I - \frac{1}{n}J) &= II - I\frac{1}{n}J - \frac{1}{n}JI + \frac{1}{n}J\frac{1}{n}J \\ &= I - \frac{1}{n}J - \frac{1}{n}J + \frac{1}{n}J\frac{1}{n}J \end{aligned}$$

$$\frac{1}{n}J = \frac{1}{n} \begin{bmatrix} 1 & \dots & 1 \\ \vdots & & \vdots \\ 1 & \dots & 1 \end{bmatrix} \quad \frac{1}{n}J\frac{1}{n}J = \frac{1}{n^2} \begin{bmatrix} n & \dots & n \\ \vdots & & \vdots \\ n & \dots & n \end{bmatrix}$$

$$= \frac{1}{n}J$$

$$\Rightarrow (I - \frac{1}{n}J)(I - \frac{1}{n}J) = I - \frac{1}{n}J \quad (\text{idempotent})$$

$$\Rightarrow S = \frac{1}{n-1} X'(I - \frac{1}{n}J)X$$

Sample Std. Deviation matrix: $D^{1/2} = \begin{bmatrix} \sqrt{s_{11}} & & 0 \\ & \sqrt{s_{22}} & \\ 0 & & \dots \\ & & & \sqrt{s_{pp}} \end{bmatrix}$

$$D^{-1/2} = \begin{bmatrix} 1/\sqrt{s_{11}} & & 0 \\ & 1/s_{22} & \\ 0 & & \dots \\ & & & 1/\sqrt{s_{pp}} \end{bmatrix}$$

$$R = D^{-1/2} S D^{-1/2}$$

$$S = D^{1/2} R D^{1/2}$$

3.6. Sample values of Linear Combination of vars

$$\underline{\tilde{X}} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix} \quad \underline{c}'\underline{X} = c_1 X_1 + \dots + c_p X_p \quad (\text{General form})$$

observed value on j^{th} trial:

$$\underline{c}'\underline{\tilde{x}}_j = c_1 X_{j1} + \dots + c_p X_{jp} \quad j=1, \dots, n$$

$$\text{Sample mean: } \frac{\underline{c}'\underline{\tilde{x}}_1 + \dots + \underline{c}'\underline{\tilde{x}}_n}{n} = \frac{1}{n} \underline{c}'(\underline{\tilde{x}}_1 + \dots + \underline{\tilde{x}}_n)$$

$$= \underline{c}'\bar{\underline{X}}$$

$$\text{Sample variance: } \frac{1}{n-1} \sum_{j=1}^n (\underline{c}'\underline{\tilde{x}}_j - \underline{c}'\bar{\underline{X}})^2$$

$$= \frac{1}{n-1} \sum_{j=1}^n [c'(x_j - \bar{x})]^2 = \frac{1}{n-1} \sum_{j=1}^n [c'(x_j - \bar{x})(x_j - \bar{x})'c]$$

$$= \underline{c}' \frac{\sum_{j=1}^n (x_j - \bar{x})(x_j - \bar{x})'}{n-1} \underline{c} = \underline{c}' S \underline{c}$$

Second Linear combination: $\underline{b}'\underline{X}$

$$\text{Sample Mean: } \underline{b}'\bar{\underline{X}} \quad \text{Sample variance: } \underline{b}' S \underline{b}$$

Sample Covariance

$$\underline{(b_0' \underline{\tilde{x}}_1 - b_0' \bar{\underline{X}})(c_0' \underline{\tilde{x}}_1 - c_0' \bar{\underline{X}}) + \dots + (b_0' \underline{\tilde{x}}_n - b_0' \bar{\underline{X}})(c_0' \underline{\tilde{x}}_n - c_0' \bar{\underline{X}})}$$

$$n-1$$

$$= \frac{b'(x_1 - \bar{x})(x_1 - \bar{x})'c + \dots + b'(x_n - \bar{x})(x_n - \bar{x})'c}{n-1}$$

$$= b'Sc = \text{Cov}\{b'X, c'X\}$$

Generalized to q linear combinations:

$$a_{i1}X_1 + a_{i2}X_2 + \dots + a_{ip}X_p \quad i=1, \dots, q$$

$$\Rightarrow AX \quad \text{w/} \quad A = \begin{bmatrix} a_{11} & \dots & a_{1p} \\ \vdots & \ddots & \vdots \\ a_{q1} & \dots & a_{qp} \end{bmatrix}$$

sample mean: $A\bar{X}$

sample covariance matrix: ASA'