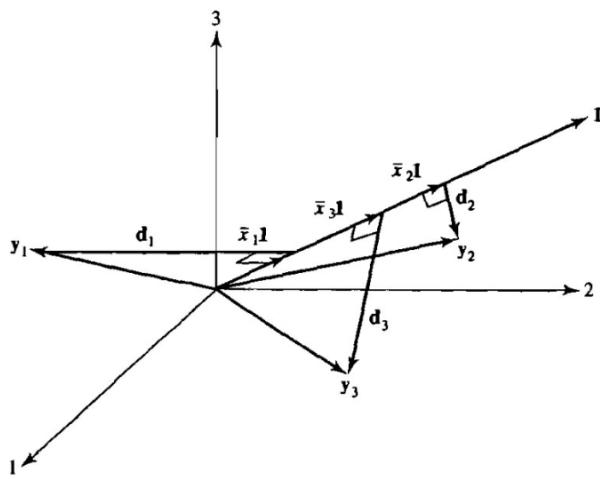


## 3.2

## The Geometry of the Sample

The scatter plot of  $n$  points in  $p$ -dimensional space provides information on the locations and variability of the points.

- \* coefficient of determination ( $R^2$ ) is the determinant of variance-covariance matrix.
- \* mean vector  $\bar{X}$  is the center of the  $n$  points in  $p$ -dimension
- although graphical representation is limited to 3-dimensions, geometrical relationships and the associated statistical concepts remain valid for extended dimensions.



$\bar{x}_i$ 's and corresponding  $d_i$ 's are orthogonal.

Decomposition of  $y_i$  vectors

$$d_i = y_i - \bar{x}_i 1$$

$$y_i = d_i + \bar{x}_i 1$$

Let  $y = [x_{1i}, x_{2i}, \dots, x_{ni}]$ . The projection of  $y_i$  on the unit vector  $\frac{1}{\sqrt{n}} 1$  is

$$y'_i \left( \frac{1}{\sqrt{n}} 1 \right) \frac{1}{\sqrt{n}} 1 = \frac{x_{1i} + x_{2i} + \dots + x_{ni}}{\sqrt{n}} = \bar{x}_{i1}$$

The squared length of the deviation vectors equals the sum of squared deviations

\* the length of deviation vector is proportional to the standard deviation.

Squared Lengths of the deviation vectors

$$L^2_{d_i} = d'_i d_i = \underbrace{\sum_{j=1}^n}_{(\text{Length})^2} (x_{ji} - \bar{x}_i)^2 \quad \underbrace{\sum}_{\text{sum of squared deviations}}$$

\* note that the squared length  $L^2$  is proportional to the variance of the measurements on the  $j^{th}$  variable. Equivalently, the length is proportional to the standard deviation. Longer vectors represent more variability than shorter vectors.

For any two deviation vectors  $d_i$  and  $d_k$ ,

$d_i^* d_k = \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k)$ , and let  $\theta_{ik}$  be the angle between  $d_i$  and  $d_k$ ,

$$d_i^* d_k = L_{d_i} L_{d_k} \cos(\theta_{ik}) = \sqrt{\sum_{j=1}^n (x_{ji} - \bar{x}_i)^2} \sqrt{\sum_{j=1}^n (x_{jk} - \bar{x}_k)^2} \cdot \cos(\theta_{ik}),$$

then we notice that

$$\frac{d_i \cdot d_k}{L_{d_i} \cdot L_{d_k}} = \frac{\sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k)}{\sqrt{\sum_{i=1}^n (x_{ji} - \bar{x}_i)^2} \cdot \sqrt{\sum_{j=1}^n (x_{jk} - \bar{x}_k)^2}} = \cos(\theta_{ik}) = r_{ik} = \text{Cor}(i, k)$$

### 3.3 Random Samples and the Expected Values of the Sample Mean and Covariance Matrix

- Euclidean distance appears appropriate when samples are independent and have the same variances.
  - If not, the influence of individual measurements on location is asymmetric so the use of statistical distances or quadratic forms are appropriate.

\* Refer to the proof at pg 142-143

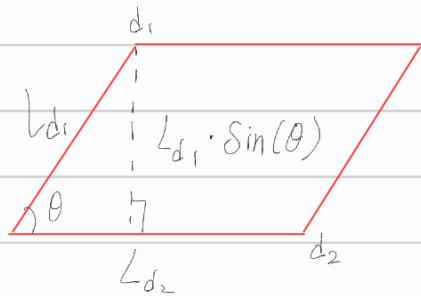


## Unbiased Sample Variance - Covariance Matrix

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

## 3.4 Generalized Variance

- $|S|$ , the determinant of a variance-covariance matrix, is a single numerical value for the variation expressed (**Generalized Variance**)
  - Consider the area of a trapezoid,



$$\text{then, } \text{Area} = |L_{d_1} \sin(\theta) | L_{d_2} \\ = L_{d_1} L_{d_2} \sqrt{1 - \cos^2 \theta}, \text{ and}$$

$$\text{we know) } L_{d_1} = \sqrt{\sum_{j=1}^n (X_{ji} - \bar{X})^2} = \sqrt{(n-1)S_{11}},$$

$$L_{d_2} = \sqrt{\sum_{j=1}^n (X_{j2} - \bar{X})^2} = \sqrt{(n-1)S_{22}},$$

$$\cos(\theta) = r_{12}, \text{ so}$$

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

$$S_{12} = S_{21} = \sqrt{S_{11}} \sqrt{S_{22}} R_{12}, \quad \text{so } b/c \quad |S| = S_{11}S_{22} - S_{11}S_{22}R_{12}^2 = S_{11}S_{22}(1-R_{12}^2),$$

$$|S| = \frac{(\text{Area})^2}{(n-1)^2} = \underbrace{(n-1)^{-p}}_{\text{determinant}} (\text{volume})^2$$

Generalized sample variance

The generalized variance with  $|S|=0$  is indicative of degeneracy.

- when at least one deviation vector lies in the (hyper) plane formed by linear combinations of the others  $\Leftrightarrow$  linearly dependent (Column Degeneracy)

↳ this means that one of the deviation vectors lies in the plane generated by the other residual vectors of the others. (the volume or area is 0)

Conditions for generalized variance = 0 (either all true or all false)

1.  $S\alpha = 0$ ,  $\alpha$  is a scaled eigenvector of  $S$  with eigenvalue 0.

2.  $\alpha'(X_j - \bar{X}) = 0$ , linear combination of the mean corrected data is 0.

3.  $\alpha'X_j = C$ ,  $C = \alpha'\bar{X}$ , linear combinations of the original data is a constant

- In any statistical analysis,  $|S|=0$  means that the measurements on some variables should be removed from the study as far as the mathematical computations are concerned.

Some Pre-alerts of  $|S|=0$  cases

1.  $n \leq p$

? Covariance matrix

2. If the linear combination  $\alpha'X_i$  has positive variance for each constant vector  $\alpha \neq 0$ , then, provided  $p < n$ ,  $S$  has full rank with probability 1 and  $|S| > 0$ .

3. If, with probability 1,  $\alpha'X_i$  is a constant, then  $|S|=0$ .

Generalized Variance Determined by  $|R|$  and its Geometrical Interpretation

- The generalized sample variance is unduly affected by the variability of measurements on a single data

Generalized sample variance of the standardized variables  $= |R| = (n-1)^p (\text{volume})^2$

- the generalized sample variance of the standardized variables will be large when these vectors are nearly perpendicular and will be small when two or more of these vectors are in almost the same direction

$$|S| = (S_{11} S_{22} S_{33} \dots S_{pp}) |R|$$

$$(n-1)^p |S| = (n-1)^p (S_{11} S_{22} S_{33} \dots S_{pp}) |R|$$

Another Generalization of Variance

$$\text{Total Sample Variance} = \sum_{i=1}^n S_{ii} = \text{tr}(S)$$

### 3.5 Sample Mean, Covariance, and Correlation as Matrix Operations

$$\bar{X} = \frac{1}{n} X' 1$$

$$1 \bar{X}' = \frac{1}{n} 1 1' X = \begin{bmatrix} \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_p \\ \bar{x}_1 & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \bar{x}_1 & \cdots & \bar{x}_p \end{bmatrix}$$

$$X - \frac{1}{n} 1 1' X = \begin{bmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \cdots & x_{1p} - \bar{x}_p \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 & \cdots & x_{np} - \bar{x}_p \end{bmatrix} \rightarrow \text{from we can derive}$$

$$(n-1) S = \begin{bmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \cdots & x_{1p} - \bar{x}_p \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 & \cdots & x_{np} - \bar{x}_p \end{bmatrix} \begin{bmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \cdots & x_{1p} - \bar{x}_p \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 & \cdots & x_{np} - \bar{x}_p \end{bmatrix}$$

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

$$= (I - \frac{1}{n} 1 1')$$

$$\therefore \bar{X} = \frac{1}{n} X 1$$

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

Let  $D^{\frac{1}{2}} = \begin{bmatrix} \sqrt{s_{11}} & 0 & \cdots & 0 \\ 0 & \sqrt{s_{22}} & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & \cdots & \sqrt{s_{pp}} \end{bmatrix}$ ,  $D^{-\frac{1}{2}} = \begin{bmatrix} \frac{1}{\sqrt{s_{11}}} & 0 & \cdots & 0 \\ 0 & \frac{1}{\sqrt{s_{22}}} & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \cdots & \cdots & \frac{1}{\sqrt{s_{pp}}} \end{bmatrix}$ , and

$$S = \begin{bmatrix} s_{11} & s_{12} & \cdots & s_{1p} \\ \vdots & \ddots & \ddots & \vdots \\ s_{p1} & \cdots & \cdots & s_{pp} \end{bmatrix} = D^{\frac{1}{2}} R D^{\frac{1}{2}}$$

$$R = \begin{bmatrix} 1 & r_{12} & \cdots & r_{p1} \\ r_{21} & 1 & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ r_{p1} & \cdots & \cdots & 1 \end{bmatrix} = D^{-\frac{1}{2}} S D^{-\frac{1}{2}}$$

### 3.6 Sample Values of Linear Combinations of Variables

If  $C'X = C_1X_1 + C_2X_2 + \dots + C_pX_p$ ,

Sample mean =  $C'\bar{X}$ , and sample variance of  $C'X = C'SC$

If  $b'X = b_1X_1 + b_2X_2 + \dots + b_pX_p$ ,

Covariance of  $C'X$  = covariance  $b'X = b'SC$

