Sample Geometry and Random Sampling

Shyh-Kang Jeng

Department of Electrical Engineering/ Graduate Institute of Communication/ Graduate Institute of Networking and Multimedia

Outline

- The Geometry of the Sample
- Random Samples and the Expected Values of the Sample Mean and Covariance Matrix
- Generalized Variance
- Sample Mean, Covariance, and Correlation as Matrix Operations
- Sample Values of Linear
 Combinations of Variables

Outline

- The Geometry of the Sample
- Random Samples and the Expected Values of the Sample Mean and Covariance Matrix
- Generalized Variance
- Sample Mean, Covariance, and Correlation as Matrix Operations
- Sample Values of Linear Combinations of Variables

Questions

- + How to represent a sample of size *n* from a *p*-variate population?
- What is the geometrical representation of sample mean and deviation?
- How to calculate lengths and angles of deviation vectors?
- What is the geometric meaning of the correlation coefficient?















Lengths and Angles of Deviation Vectors $L_{\mathbf{d}_{i}}^{2} = \mathbf{d}_{i} \mathbf{d}_{i} = \sum_{j=1}^{n} (x_{ji} - \overline{x}_{i})^{2} = ns_{ii}$ $\mathbf{d}_{i} \mathbf{d}_{k} = \sum_{j=1}^{n} (x_{ji} - \overline{x}_{i})(x_{jk} - \overline{x}_{k}) = ns_{ik}$ $= L_{\mathbf{d}_{i}} L_{\mathbf{d}_{k}} \cos \theta_{ik}$ $= \sqrt{\sum_{j=1}^{n} (x_{ji} - \overline{x}_{i})^{2}} \sqrt{\sum_{j=1}^{n} (x_{ji} - \overline{x}_{i})^{2}} \cos \theta_{ik}$ $\cos \theta_{ik} = \frac{s_{ik}}{\sqrt{s_{ii}} \sqrt{s_{kk}}} = r_{ik}$





$$\mathbf{X} = \begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{np} \end{bmatrix} = \begin{bmatrix} \mathbf{X}_{1}^{'} \\ \mathbf{X}_{2}^{'} \\ \vdots \\ \mathbf{X}_{n}^{'} \end{bmatrix}$$



20

Random Sample

- Measurements of a single trial, such as X_j'=[X_{j1},X_{j2},...,X_{jp}], will usually be correlated
- The measurements from different trials must be independent
- The independence of measurements from trial to trial may not hold when the variables are likely to drift over time

19

Geometric Interpretation of Randomness

- Column vector $\mathbf{Y}_{k}' = [X_{1k}, X_{2k}, \dots, X_{nk}]$ regarded as a point in *n* dimensions
- The location is determined by the joint probability distribution f(y_k) = f(x_{1k}, x_{2k},...,x_{nk})
- For a random sample, $f(\mathbf{y}_k) = f_k(x_{1k}) f_k(x_{2k}) \dots f_k(x_{nk})$
- Each coordinate x_{jk} contributes equally to the location through the same marginal distribution $f_k(x_{jk})$







Proof of Result 3.1

$$\begin{aligned}
& (\overline{\mathbf{X}}\overline{\mathbf{X}}') = E((\overline{\mathbf{X}} - \mu + \mu)(\overline{\mathbf{X}} - \mu + \mu)) = \frac{1}{n}\Sigma + \mu\mu' \\
& (E(\mathbf{S}_n) = \frac{1}{n}E(\sum_{j=1}^n \mathbf{X}_j \mathbf{X}_j - n\overline{\mathbf{X}}\overline{\mathbf{X}}') \\
& = \frac{1}{n}\left(n(\Sigma + \mu\mu') - n\left(\frac{1}{n}\Sigma - \mu\mu'\right)\right) = \frac{n-1}{n}\Sigma
\end{aligned}$$

Some Other Estimators

The expectation of the (i,k)th entry of $\frac{n}{n-1}\mathbf{S}_n$

 $E(\frac{n}{n-1}s_{ik}) = E(\frac{1}{n-1}\sum_{j=1}^{n} \left(X_{ji} - \overline{X}_{i}\right) \left(X_{jk} - \overline{X}_{k}\right) = \sigma_{ik}$ $E(\sqrt{s_{ii}}) \neq \sqrt{\sigma_{ii}}, \quad E(r_{ik}) \neq \rho_{ik}$

Biases $E(\sqrt{s_{ii}}) - \sqrt{\sigma_{ii}}$ and $E(r_{ik}) - \rho_{ik}$ can usually be ignored if size *n* is moderately large



Questions

- How to define a generalized sample variance?
- What is the geometric interpretation of a generalized sample variance for bivariate cases?
- What is the geometric interpretaion of a generalized sample variance for multivariate cases?

Questions

- What is the equation for points within a constant statistical distance c from the sample mean?
- Example 3.8
- Example 3.9
- Examples causing zero generalized variance

Questions

- Example 3.10
- Result 3.3
- Result 3.4
- Generalized Sample Variance of Standardized Variables
- Example 3.11
- Total Sample Variance

Generalized Sample Variance

Generalized Sample Variance = $|\mathbf{S}|$ Example 3.7 : Employees and profits per employee for 16 largest publishing firms in US

$$\mathbf{S} = \begin{bmatrix} 252.04 & -68.43 \\ -68.43 & 123.67 \end{bmatrix}$$

|S| = 26.487









Example 3.8: Sample Mean and
Variance-Covariance Matrices
$$\mathbf{S} = \begin{bmatrix} 5 & 4 \\ 4 & 5 \end{bmatrix}, r = 0.8$$
$$\mathbf{S} = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}, r = 0$$
$$\mathbf{S} = \begin{bmatrix} 5 & -4 \\ -4 & 5 \end{bmatrix}, r = -0.8$$
$$\mathbf{\bar{x}} = [2, 1], |\mathbf{S}| = 9 \text{ for all three cases}$$

Example 3.8:
Eigenvalues and Eigenvectors

$$\begin{bmatrix} 5 & 4 \\ 4 & 5 \end{bmatrix} : \lambda_1 = 9, \lambda_2 = 1$$

$$\mathbf{e}_1 = [1/\sqrt{2}, 1/\sqrt{2}], \mathbf{e}_2 = [1/\sqrt{2}, -1/\sqrt{2}]$$

$$\begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix} : \lambda_1 = 3, \lambda_2 = 3$$

$$\mathbf{e}_1 = [1, 0], \mathbf{e}_2 = [0, 1]$$

$$\begin{bmatrix} 5 & -4 \\ -4 & 5 \end{bmatrix} : \lambda_1 = 9, \lambda_2 = 1$$

$$\mathbf{e}_1 = [1/\sqrt{2}, -1/\sqrt{2}], \mathbf{e}_2 = [1/\sqrt{2}, 1/\sqrt{2}]$$

Example 3.8: Mean-Centered Ellipse $(\mathbf{x} - \overline{\mathbf{x}})\mathbf{S}^{-1}(\mathbf{x} - \overline{\mathbf{x}}) \le c^2$ $(\mathbf{x} - \overline{\mathbf{x}})\mathbf{S}^{-1}(\mathbf{x} - \overline{\mathbf{x}}) = \frac{y_1^2}{2}$

 $\left(\mathbf{x} - \overline{\mathbf{x}}\right)^{*} \mathbf{S}^{-1} (\mathbf{x} - \overline{\mathbf{x}}) = \frac{y_{1}^{2}}{\lambda_{1}} + \frac{y_{2}^{2}}{\lambda_{2}}$ $\mathbf{S}^{-1} : \text{eigenvalues } \frac{1}{\lambda_{1}}, \frac{1}{\lambda_{2}}; \text{ eigenvectors } \mathbf{e}_{1}, \mathbf{e}_{2}$ $\left(\because \mathbf{S} \mathbf{e} = \lambda \mathbf{e}, \quad \mathbf{e} = \lambda \mathbf{S}^{-1} \mathbf{e} \right)$ $\begin{bmatrix} y_{1} \\ y_{2} \end{bmatrix} = \begin{bmatrix} \mathbf{e}_{1} \\ \mathbf{e}_{2} \end{bmatrix} \begin{bmatrix} x_{1} - \overline{x}_{1} \\ x_{2} - \overline{x}_{2} \end{bmatrix}$ Choose $c^{2} = 5.99$ to cover approximately 95% observations

Example 3.8: Semi-major and Semi-minor Axes	
$\mathbf{S} = \begin{bmatrix} 5 & 4 \\ 4 & 5 \end{bmatrix}, a = 3\sqrt{5.99}, b = \sqrt{5.99}$	
$\mathbf{S} = \begin{bmatrix} 0 & 3 \end{bmatrix}, a = \sqrt{3}\sqrt{5.99}, b = \sqrt{3}\sqrt{5.99}$ $\mathbf{S} = \begin{bmatrix} 5 & -4 \\ -4 & 5 \end{bmatrix}, a = 3\sqrt{5.99}, b = \sqrt{5.99}$	
	38













Examples Cause Zero Generalized Variance

Example 1

- -Data are test scores
- Included variables that are sum of others
- -e.g., algebra score and geometry score were combined to total math score
- –e.g., class midterm and final exam scores summed to give total points
- Example 2
 - Total weight of chemicals was included along with that of each component

Example 3.10								
$\begin{bmatrix} 1 & 9 & 10 \end{bmatrix} \qquad \begin{bmatrix} -2 & -1 & -3 \end{bmatrix}$								
4 12 16 1 2 3 2.5 0	2.5							
$\mathbf{X} = \begin{vmatrix} 2 & 10 & 12 \end{vmatrix}, \mathbf{X} - 1\overline{\mathbf{x}}' = \begin{vmatrix} -1 & 0 & -1 \end{vmatrix}, \mathbf{S} = \begin{vmatrix} 0 & 2.5 \end{vmatrix}$	2.5							
5 8 13 2 -2 0 2.5 2.5	5.0							
$ \mathbf{S} = 0 \Longrightarrow \mathbf{S}\mathbf{a} = 0$								
Eigenvector corresponding to zero eigenvalues of S								
\Rightarrow a ' = [1, 1, -1]								
$\therefore 1(x_{j1} - \overline{x}_1) + 1(x_{j2} - \overline{x}_2) - (x_{j3} - \overline{x}_3) = 0$								
	46							

Result 3.3





Proof of Result 3.3

∴ row₁($\mathbf{X} - \mathbf{1}\overline{\mathbf{x}}$ ')' is a linear combination of the remaining row vectors col_k(S) is a linear combination of at most n-1linear independent of transpose of row vectors The rank of **S** is thus less than or equal to n-1, i.e., less than or equal to p-1. Since **S** is a p by p matrix, $|\mathbf{S}| = 0$

Result 3.4 Let the *p* by 1 vectors x_j, x₂, ..., x_n, where x_j' is the *j*th row of the data matrix X, be realizations of the independent random vectors X₁, X₂, ..., X_n. If the linear combination a'X_j has positive variance for each non-zero constant vector a, then, provided that *p* < *n*, S has full rank with probability 1 and |S| > 0 If, with probability 1, a'X_j is a constant *c* for all *j*, then |S| = 0



Generalized Sample Variance of Standardized Variables Generalized sample variance of

the standardized variables = $/\mathbf{R}/$

$\frac{y_i - \overline{x}_i 1}{\sqrt{s_{ii}}} = \begin{bmatrix} \frac{x_{1i} - \overline{x}_i}{\sqrt{s_{ii}}} & \frac{x_{2i} - \overline{x}_i}{\sqrt{s_{ii}}} & \cdots & \frac{x_{ni} - \overline{x}_i}{\sqrt{s_{ii}}} \end{bmatrix}$							
$ \mathbf{R} = (n-1)^{-p} (volume)^2, \mathbf{S} = (s_{11}s_{22}\cdots s_{pp}) \mathbf{R} $							
$ \mathbf{R} $ is large when all r_{ik} are nearly zero, and is small							
when one or more r_{ik} are nearly +1 or -1							





Total Sample Variance

```
Total Sample Variance = s_{11} + s_{22} + \dots + s_{pp}

Pays no attention to the orientation of the residual vectors

Example 3.7 : S = \begin{bmatrix} 252.04 & -68.43 \\ -68.43 & 123.67 \end{bmatrix}

Total sample variance = 375.71

Example 3.9 : S = \begin{bmatrix} 3 & -3/2 & 0 \\ -3/2 & 1 & 1/2 \\ 0 & 1/2 & 1 \end{bmatrix}

Total sample variance = 5
```





Covariance as Matrix Operation $\begin{bmatrix} \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_n \end{bmatrix}$

$$\begin{bmatrix} x_{1} & x_{2} & \cdots & x_{p} \\ \overline{x}_{1} & \overline{x}_{2} & \cdots & \overline{x}_{p} \\ \vdots & \vdots & \ddots & \vdots \\ \overline{x}_{1} & \overline{x}_{2} & \cdots & \overline{x}_{p} \end{bmatrix} = \mathbf{1}\overline{\mathbf{x}}' = \frac{1}{n}\mathbf{1}\mathbf{1}'\mathbf{X}$$

$$\begin{bmatrix} x_{11} - \overline{x}_{1} & x_{12} - \overline{x}_{2} & \cdots & x_{1p} - \overline{x}_{p} \\ x_{21} - \overline{x}_{1} & x_{22} - \overline{x}_{2} & \cdots & x_{2p} - \overline{x}_{p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} - \overline{x}_{1} & x_{n2} - \overline{x}_{2} & \cdots & x_{np} - \overline{x}_{p} \end{bmatrix} = \mathbf{X} - \mathbf{1}\mathbf{1}'\mathbf{X}$$

$$\begin{aligned} \textbf{Covariance as Matrix Operation} \\ \textbf{(n-1)s} &= \begin{bmatrix} x_{11} - \overline{x}_1 & x_{21} - \overline{x}_1 & \cdots & x_{n1} - \overline{x}_1 \\ x_{12} - \overline{x}_2 & x_{22} - \overline{x}_2 & \cdots & x_{n2} - \overline{x}_2 \\ \vdots & \vdots & \ddots & \vdots \\ x_{1p} - \overline{x}_p & x_{2p} - \overline{x}_p & \cdots & x_{np} - \overline{x}_p \end{bmatrix} \\ \textbf{x} \\ \begin{bmatrix} x_{11} - \overline{x}_1 & x_{12} - \overline{x}_2 & \cdots & x_{1p} - \overline{x}_p \\ x_{21} - \overline{x}_1 & x_{22} - \overline{x}_2 & \cdots & x_{2p} - \overline{x}_p \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} - \overline{x}_1 & x_{n2} - \overline{x}_2 & \cdots & x_{np} - \overline{x}_p \end{bmatrix} \\ = \begin{pmatrix} \textbf{x} - \frac{1}{n} \textbf{1}^* \textbf{X} \end{pmatrix} \begin{pmatrix} \textbf{x} - \frac{1}{n} \textbf{1}^* \textbf{X} \end{pmatrix} \end{aligned}$$

Covariance as Matrix Operation

$$\left(\mathbf{X} - \frac{1}{n}\mathbf{11'X}\right) \left(\mathbf{X} - \frac{1}{n}\mathbf{11'X}\right)$$

= $\mathbf{X'}(\mathbf{I} - \frac{1}{n}\mathbf{11'})'(\mathbf{I} - \frac{1}{n}\mathbf{11'})\mathbf{X}$
 $(\mathbf{I} - \frac{1}{n}\mathbf{11'})'(\mathbf{I} - \frac{1}{n}\mathbf{11'}) = \mathbf{I} - \frac{1}{n}\mathbf{11'} - \frac{1}{n}\mathbf{11'} + \frac{1}{n^2}\mathbf{11'11'}$
= $\mathbf{I} - \frac{1}{n}\mathbf{11'}$ (:: $\mathbf{1'1} = n$)
 $\mathbf{S} = \frac{1}{n-1}\mathbf{X'}(\mathbf{I} - \frac{1}{n}\mathbf{11'})\mathbf{X}$





61

Result 3.5

b'**X** = $b_1X_1 + b_2X_2 + \dots + b_pX_p$ **c**'**X** = $c_1X_1 + c_2X_2 + \dots + c_pX_p$ Sample mean of **b**'**X** = **b**'**X** Sample variance of **b**'**X** = **b**'**Sb** Sample covariance of **b**'**X** and **c**'**X** = **b**'**Sc**



Proof of Result 3.5

Sample covariance =
$$\frac{1}{n-1} \sum_{j=1}^{n} (\mathbf{b}' \mathbf{x}_{j} - \mathbf{b}' \overline{\mathbf{x}}) (\mathbf{c}' \mathbf{x}_{j} - \mathbf{c}' \overline{\mathbf{x}})$$

= $\frac{1}{n-1} \mathbf{b}' \sum_{j=1}^{n} (\mathbf{x}_{j} - \overline{\mathbf{x}}) (\mathbf{x}_{j} - \overline{\mathbf{x}})' \mathbf{c} = \mathbf{b}' \mathbf{S} \mathbf{c}$

67

	F	Resu	ılt 3	.6				
AX =	$egin{array}{c} a_{11} \\ a_{21} \\ \vdots \\ a_{q1} \end{array}$	$egin{array}{c} a_{12} \ a_{22} \ dots \ a_{q2} \ dots \ a_{q2} \end{array}$	···· ··· ··.	$egin{aligned} a_{1p} \ a_{2p} \ dots \ a_{qp} \end{bmatrix}$	$\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{bmatrix}$			
Sample mean of $\mathbf{A}\mathbf{X} = \mathbf{A}\overline{\mathbf{x}}$								
Sample covariance matrix $= ASA'$								
						68		