

# Important Misc in Probab/Stat<sup>1</sup>

## 1 Important matrix decompositions

- Eigen Decomposition:** Let  $\mathbf{P}$  be a matrix of eigenvectors of a given square matrix  $\mathbf{A}$  and  $\mathbf{D}$  be a diagonal matrix with the corresponding eigenvalues on the diagonal. Then, as long as  $\mathbf{P}$  is a square matrix,  $\mathbf{A}$  can be written as an **eigen decomposition**

$$\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1},$$

where  $\mathbf{D}$  is a diagonal matrix. Furthermore, if  $\mathbf{A}$  is symmetric, then the columns of  $\mathbf{P}$  are **orthogonal** vectors.

If  $\mathbf{P}$  is not a square matrix (for example, the space of eigenvectors of  $\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$  is one-dimensional), then  $\mathbf{P}$  cannot have a matrix inverse and  $\mathbf{A}$  does not have an eigen decomposition. However, if  $\mathbf{P}$  is  $m \times n$  (with  $m > n$ ), then  $\mathbf{A}$  can be written using a so-called **singular value decomposition**.

- QR-decomposition** (Gram-Schmidt orthogonality): For any matrix  $\mathbf{A}_{n \times m}$ , there exists QR-decomposition  $\mathbf{A} = \mathbf{Q}_{n \times m}\mathbf{R}_{m \times m}$ , where  $\mathbf{Q}^T\mathbf{Q} = \mathbf{I}_{m \times m}$  and  $\mathbf{R}$  is an upper triangular.
- QR-decomposition** (Another version) For any matrix  $\mathbf{A}_{n \times m}$  of rank  $k$ , there exists QR-decomposition  $\mathbf{A} = \mathbf{Q}_{n \times m}\mathbf{R}_{m \times m}$ , where  $\mathbf{Q}^T\mathbf{Q}$  is diagonal and  $\mathbf{R}$  is a **unit** upper triangular.
- Left orthogonal decomposition** For any matrix  $\mathbf{A}_{n \times m}$ , there exist non-singular matrix  $\mathbf{P}_{m \times m}$  and orthogonal matrix  $\mathbf{G}_{n \times n}$  such that  $\mathbf{A} = \mathbf{G} \begin{bmatrix} \mathbf{I}_r & 0 \\ 0 & 0 \end{bmatrix} \mathbf{P}$ , where  $r = \text{rank}(\mathbf{A})$ .
- Cholesky Decomposition:** For positive matrix  $\mathbf{A}_{n \times n}$ , there exists Cholesky-decomposition  $\mathbf{A} = \mathbf{T}^T\mathbf{T}$ , where  $\mathbf{T}$  is an upper triangular. Also  $\mathbf{T}$  is unique.
- $\mathbf{A}$  and  $\mathbf{B}$  are real symmetric matrices. Then there exists a orthogonal matrix  $\mathbf{P}$  such that  $\mathbf{P}^T\mathbf{A}\mathbf{P}$  and  $\mathbf{P}^T\mathbf{B}\mathbf{P}$  are both diagonal if and only if  $\mathbf{AB} = \mathbf{BA}$ .
- Spectral decomposition** (A special case of item 1): Let  $\mathbf{A}$  be  $n \times n$  symmetric matrix. There exists an orthogonal matrix  $\mathbf{T} = (t_1, \dots, t_n)$  such that  $\mathbf{T}^T\mathbf{A}\mathbf{T} = \text{diag}(\lambda_1, \dots, \lambda_n) = \mathbf{\Lambda}$ , where  $\lambda_1 \geq \dots \geq \lambda_n$  are the ordered eigenvalues of  $\mathbf{A}$ . With this ordering,  $\mathbf{\Lambda}$  is unique and  $\mathbf{T}$  is unique up to a postfactor.

$$\mathbf{A} = \sum_{i=1}^n \lambda_i t_i t_i^T.$$

## 2 Normal Distribution

- If

$$\begin{pmatrix} \mathbf{Y} \\ \mathbf{X} \end{pmatrix} \sim \text{Normal} \left\{ \begin{pmatrix} \mu_y \\ \mu_x \end{pmatrix}, \begin{pmatrix} \Sigma_{yy} & \Sigma_{yx} \\ \Sigma_{xy} & \Sigma_{xx} \end{pmatrix} \right\}.$$

Then  $\mathbf{Y}|\mathbf{X} \sim \text{Normal}(\mu_{y|x}, \Sigma_{y|x})$  with  $\mu_{y|x} = \mu_y + \Sigma_{yx}\Sigma_{xx}^{-1}(\mathbf{X} - \mu_x)$  and  $\Sigma_{y|x} = \Sigma_{yy} - \Sigma_{yx}\Sigma_{xx}^{-1}\Sigma_{xy}$

<sup>1</sup>1-VIP-misc-pocket-HL.tex

- A useful equality ( $\Phi, \phi$ : normal CDF, pdf)

$$\int \Phi(a + bx)\phi(x)dx = \Phi\left(\frac{a}{\sqrt{1+b^2}}\right).$$

**Proof.** Let  $X, Z$  be iid random variables following standard normal. Then  $\Phi(a + bx) = \Pr(Z \leq a + bx)$ . Note that

$$\begin{aligned} \int \Phi(a + bx)\phi(x)dx &= E_X \Phi(a + bX) \\ &= E_X \{P_Z(Z \leq a + bX)\} \\ &= E_X \{P_Z(Z - bX \leq a)\} \\ &= P_{(Z, X)}(Z - bX \leq a) \\ &= P\left\{N(0, 1) \leq a/\sqrt{1+b^2}\right\} = \Phi\left(\frac{a}{\sqrt{1+b^2}}\right) \end{aligned}$$

- Stein formulae:**  $X$  is a  $N(0, 1)$  random variable and  $g$  is an indefinite integral of the Lebesgue measurable function such that  $E|g'(X)| < \infty$ . Then

$$E\{g'(X)\} = E\{Xg(X)\}.$$

## 3 Linear and Quadratic Forms

- $q = \mathbf{y}^T \mathbf{A} \mathbf{y}$  is called a quadratic form in  $\mathbf{y}$ .  $E(q) = \text{tr}(\mathbf{A}\mathbf{V}) + \mu^T \mathbf{A} \mu$  (**y may not be normal**) and  $\text{cov}(q_1, q_2) = 2\text{tr}(\mathbf{A}_1 \mathbf{V} \mathbf{A}_2 \mathbf{V}) + 4\mu^T \mathbf{A}_1 \mathbf{V} \mathbf{A}_2 \mu$ .
- $\mathbf{y} \sim N(\mu, \mathbf{V})$  then its characteristic function is  $m_{\mathbf{y}}(t) = \exp\{t^T \mu + \frac{1}{2} t^T \mathbf{V} t\}$ . The conditional density of  $\mathbf{y}_2$  given  $\mathbf{y}_1$  is  $N(\mu_2 + \mathbf{V}_{21}\mathbf{V}_{11}^{-1}(\mathbf{y}_1 - \mu_1), \mathbf{V}_{22} - \mathbf{V}_{21}\mathbf{V}_{11}^{-1}\mathbf{V}_{12})$ .
- Craig's Theorem**  $\mathbf{y} \sim N(\mu, \mathbf{V})$ .  $\mathbf{B}\mathbf{y}$  is independent of  $\mathbf{y}^T \mathbf{A} \mathbf{y}$  iff  $\mathbf{B}\mathbf{V}\mathbf{A} = 0$ ;  $\mathbf{y}^T \mathbf{B} \mathbf{y}$  is independent of  $\mathbf{y}^T \mathbf{A} \mathbf{y}$  iff  $\mathbf{B}\mathbf{V}\mathbf{A} = 0$ .  $q = \mathbf{y}^T \mathbf{A} \mathbf{y} \sim \chi^2(r, \lambda)$ , with  $r = r(\mathbf{A})$  and  $\lambda = 1/2\mu^T \mathbf{A} \mu$ , if and only if  $\mathbf{A}\mathbf{V}$  is idempotent.
- Cochran's Theorem**  $r(\mathbf{A}_i) = r_i$ . Let  $\mathbf{A} = \sum_1^k \mathbf{A}_i$  if  $\mathbf{A}\mathbf{V}$  is idempotent and  $r(\mathbf{A}) = \sum_1^k r_i$ , then  $q_i = \mathbf{y}^T \mathbf{A}_i \mathbf{y}$  are mutually independent, noncentral chi-squared variables with  $\chi^2(r_i, 1/2\mu^T \mathbf{A}_i \mu)$ .

## 4 Algebra

algebra.tex

$\mathbf{A} = (a_{ij})_{sn}$ ,  $\mathbf{B} = (b_{ij})_{nm}$ ;  $\text{tr}(\mathbf{A})$ : the trace of  $\mathbf{A}$ ;  $r(\mathbf{A})$ : the rank of  $\mathbf{A}$ ;  $\det(\mathbf{A})$ : the determinant of  $\mathbf{A}$ .

### 4.1 Trace and Eigenvalues

- $\text{tr}(\mathbf{A} + \mathbf{B}) = \text{tr}(\mathbf{A}) + \text{tr}(\mathbf{B})$ .
- $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$ .
- $\mathbf{A}$  and  $\mathbf{B}$  are real symmetric. Then  $\text{tr}(\mathbf{ABAB}) \leq \text{tr}(\mathbf{A}^2 \mathbf{B}^2)$ , the equality holds if and only if  $\mathbf{AB} = \mathbf{BA}$ .

- If  $\mathbf{A}_n$  is a symmetric and  $r(\mathbf{A}) = 1$ , then  $|I_n + \mathbf{A}_n| = 1 + \text{tr}(\mathbf{A})$ .
- For any  $n \times n$  matrix  $\mathbf{A}$  with eigenvalues  $\lambda_1, \dots, \lambda_n$ , we have the following:
  - $\text{tr}(\mathbf{A}) = \sum \lambda_i$
  - $\det(\mathbf{A}) = \prod \lambda_i$
  - $\det(I_n \pm \mathbf{A}) = \prod (1 \pm \lambda_i)$
- For conformable matrices, the nonzero eigenvalues of  $\mathbf{AB}$  are the same as those of  $\mathbf{BA}$ .

## 4.2 Rank

- $r(\mathbf{AB}) \geq r(\mathbf{A}) + r(\mathbf{B}) - n$ .
- $\mathbf{A} = (a_{ij})_{nm}$ . If  $\mathbf{A}^2 = I_n$ . Then  $r(\mathbf{A} + I_n) + r(\mathbf{A} - I_n) = n$ .
- $r(\mathbf{A} + \mathbf{B}) \leq r(\mathbf{A}) + r(\mathbf{B})$ .
- If  $\mathbf{AB} = 0$ .  $r(\mathbf{A}) + r(\mathbf{B}) \leq n$ .
- $r(\mathbf{A}) = r(\mathbf{A}^T \mathbf{A}) = r(\mathbf{A}\mathbf{A}^T)$ .
- If  $\mathbf{A}, \mathbf{B}, \mathbf{C}$  are  $m \times n$ ,  $n \times p$ ,  $p \times q$  matrices, then  $r(\mathbf{AB}) + r(\mathbf{BC}) \leq r(\mathbf{B}) + r(\mathbf{ABC})$ .

## 4.3 Patterned Matrices

- If  $\mathbf{A}$  and  $\mathbf{C}$  are symmetric and all inverses exist,

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{A}^{-1} + \mathbf{F}\mathbf{E}^{-1}\mathbf{F}^T & -\mathbf{F}\mathbf{E}^{-1} \\ -\mathbf{E}^{-1}\mathbf{F}^T & \mathbf{E}^{-1} \end{pmatrix}$$

where  $\mathbf{E} = \mathbf{C} - \mathbf{B}^T \mathbf{A}^{-1} \mathbf{B}$  and  $\mathbf{F} = \mathbf{A}^{-1} \mathbf{B}$ .

- $\begin{vmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{vmatrix} = \begin{cases} |\mathbf{D}||\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}| & \text{if } \mathbf{D}^{-1} \text{ exists,} \\ |\mathbf{A}||\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B}| & \text{if } \mathbf{A}^{-1} \text{ exists.} \end{cases}$

**Proof.**

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{pmatrix} \cdot \begin{pmatrix} \mathbf{I}_r & \mathbf{O} \\ -\mathbf{D}^{-1}\mathbf{C} & \mathbf{I}_s \end{pmatrix} = \begin{pmatrix} \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C} & \mathbf{B} \\ \mathbf{O} & \mathbf{D} \end{pmatrix}.$$

- For matrixes  $\mathbf{B}_{n \times m}$  and  $\mathbf{C}_{m \times n}$ , and non-singular  $\mathbf{A}_{n \times n}$ ,  $|\mathbf{A} + \mathbf{BC}| = |\mathbf{A}||\mathbf{I}_m + \mathbf{CA}^{-1}\mathbf{B}|$ .
- $(\mathbf{I} + \mathbf{AB})^{-1} = \mathbf{I} - \mathbf{A}(\mathbf{I} + \mathbf{BA})^{-1}\mathbf{B}$ ;  $|\mathbf{I} + \mathbf{AB}| = |\mathbf{I} + \mathbf{BA}|$ ;  $|\mathbf{A}| = |\mathbf{A}_{11}||\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}|$ .
- $(\mathbf{A} + \mathbf{UBV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{UB}(\mathbf{B} + \mathbf{BVA}^{-1}\mathbf{UB})^{-1}\mathbf{BVA}^{-1}$ . For the particular case  $\mathbf{B} = \mathbf{I}$ ,  $\mathbf{U} = \mathbf{u}$ , and  $\mathbf{V} = \mathbf{v}^T$ , we have  $(\mathbf{A} + \mathbf{uv}^T)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{uv}^T\mathbf{A}^{-1}(1 + \mathbf{v}^T\mathbf{A}^{-1}\mathbf{u})^{-1}$ . Furthermore  $\mathbf{x}^T(\mathbf{A} + \mathbf{xx}^T)^{-1}\mathbf{x} = \frac{\mathbf{x}^T\mathbf{A}^{-1}\mathbf{x}}{1 + \mathbf{x}^T\mathbf{A}^{-1}\mathbf{x}}$ .

## 4.4 Positive (semi)definite Matrices

1. Denote  $\mathbf{A}_n = (a_{ij})_{mn}$ . If  $\mathbf{A}_n$  is a positive matrix,  $|\mathbf{A}_n| \leq a_{nn}|\mathbf{A}_{n-1}|$ , and so  $|\mathbf{A}_n| \leq \prod_{i=1}^n a_{ii}$ .
2.  $\mathbf{A}$  and  $\mathbf{B}$  are real symmetric matrices. Then there exists a orthogonal matrix  $\mathbf{P}$  such that  $\mathbf{P}^\top \mathbf{A} \mathbf{P}$  and  $\mathbf{P}^\top \mathbf{B} \mathbf{P}$  are both diagonal if and only if  $\mathbf{A}\mathbf{B} = \mathbf{B}\mathbf{A}$ .
3.  $\mathbf{A}$  is symmetric. Then  $\mathbf{A}$  is nonnegative if and only if there exists a matrix  $\mathbf{C} = (c_{ij})_m$  for  $r = \text{rank}(\mathbf{A})$  such that  $\mathbf{A} = \mathbf{C}^\top \mathbf{C} \iff$  there exists  $\mathbf{B} = (b_{ij})_{mn}$  such that  $\mathbf{A} = \mathbf{B}^\top \mathbf{B} \iff$  all of eigenvalues of  $\mathbf{A}$  are nonnegative.
4.  $\mathbf{A}$  is positive and  $\mathbf{B}$  is non-negative. Then  $|\mathbf{A} + \mathbf{B}| \geq |\mathbf{A}|$  and equality holds if and only if  $\mathbf{B} = \mathbf{0}$ .
5. If  $\mathbf{A}$  is positive, then

$$f(Y) = \begin{bmatrix} \mathbf{A} & \mathbf{Y} \\ \mathbf{Y}^\top & \mathbf{0} \end{bmatrix}$$

is negative, where  $\mathbf{Y} = (y_1, \dots, y_n)^\top$ .

Note that  $f(Y) = Y^\top \{(-1)^{2n+2} \mathbf{A}^*\} Y$ , where  $\mathbf{A}^* = |\mathbf{A}| \mathbf{A}^{-1}$ .

### 6. Cholesky Decomposition.

For positive matrix  $\mathbf{A}_{n \times n}$ , there exists Cholesky-decomposition  $\mathbf{A} = \mathbf{T}^\top \mathbf{T}$ , where  $\mathbf{T}$  is an upper triangular. Also  $\mathbf{T}$  is unique.

7. Let  $\mathbf{X}^\top = (x_1, \dots, x_n)$ , where the  $x_i$  are  $n$  independent  $d$ -dimensional vectors of random variables, and let  $\mathbf{A}$  be a positive semidefinite  $n \times n$  matrix of rank  $r(\geq d)$ . Suppose that for each  $x_i$  and all  $\mathbf{b}(\neq \mathbf{0})$  and  $c$ ,  $\text{prob}[\mathbf{b}^\top x_i = c] = 0$ . Then  $\text{prob}(\mathbf{X}^\top \mathbf{A} \mathbf{X} > 0) = 1$ .
8. If  $\mathbf{A}$  is positive and  $\mathbf{B}$  is symmetric. Then there exists a nonsingular matrix  $\mathbf{C}$  such that  $\mathbf{C}^\top \mathbf{A} \mathbf{C} = \text{identity matrix}$  and  $\mathbf{C}^\top \mathbf{B} \mathbf{C} = \mathbf{\Lambda}$ , a diagonal matrix  $\text{diag}(\lambda_1, \dots, \lambda_r, 0, \dots, 0)$ , where  $\lambda_i$  is the eigenvalue of  $\mathbf{A}$ .

**Proof.** There exists a  $\mathbf{D}$  such that  $\mathbf{A} = \mathbf{D}^\top \mathbf{D}$ . Note that  $\mathbf{D}^{-1\top} \mathbf{B} \mathbf{D}^{-1}$  is still symmetric. There is a  $\mathbf{C}$  such that  $\mathbf{C}^\top \mathbf{D}^{-1\top} \mathbf{B} \mathbf{D}^{-1} \mathbf{C} = \mathbf{\Lambda}$ . Taking  $\mathbf{C} = \mathbf{D}^{-1} \mathbf{C}$ , the proof follows.

9. Suppose  $\mathbf{A}$  and  $\mathbf{B}$  are positive.  $\mathbf{A} - \mathbf{B}$  is positive if and only if  $\mathbf{B}^{-1} - \mathbf{A}^{-1}$  is positive.

**Proof.** Denote  $\mathbf{S} = \mathbf{A} - \mathbf{B}$ . If  $\mathbf{S}$  is positive, there exists a nonsingular matrix  $\mathbf{C}$  such that  $\mathbf{C}^\top \mathbf{S} \mathbf{C} = \mathbf{I}$  and  $\mathbf{C}^\top \mathbf{B} \mathbf{C} = \mathbf{\Lambda}$ . It follows that  $\mathbf{A} = \mathbf{C}^{-1\top} (\mathbf{\Lambda} + \mathbf{I}) \mathbf{C}^{-1}$  and  $\mathbf{B} = \mathbf{C}^{-1\top} \mathbf{\Lambda} \mathbf{C}^{-1}$ , and then  $\mathbf{B}^{-1} - \mathbf{A}^{-1} = \mathbf{C}^{-1\top} \{ \mathbf{\Lambda}^{-1} - (\mathbf{\Lambda} + \mathbf{I})^{-1} \} \mathbf{C}^{-1} \geq \mathbf{0}$ . Conversely, the proof for significant condition becomes trivial because the above arguments.

The conclusion can be generalized to **nonnegative** case since we can consider  $\mathbf{A} - \mathbf{B} + 1/n\mathbf{I}$  and finally let  $n \rightarrow \infty$ . Without the assumption of  $\mathbf{A}$  and  $\mathbf{B}$  being positive, the conclusion is false, see for example,  $\mathbf{A} = \mathbf{I}$  and  $\mathbf{B} = -0.5\mathbf{I}$ .

## 4.5 Idempotent Matrices

A matrix  $\mathbf{A}$  is idempotent if  $\mathbf{A}^2 = \mathbf{A}$ . A symmetric idempotent matrix is called a projection matrix.

1. If  $\mathbf{A}$  is a projection matrix of rank  $r$ , then it can be expressed in the form

$$\mathbf{A} = \sum_{i=1}^r \mathbf{t}_i \mathbf{t}_i^\top.$$

where  $\mathbf{t}_1, \dots, \mathbf{t}_r$  form an orthonormal set.

2. If  $\mathbf{A}$  is a projection matrix, then  $r(\mathbf{A}) = \text{tr}(\mathbf{A})$ .
3. If  $\mathbf{A}$  is idempotent, then so is  $\mathbf{I} - \mathbf{A}$ .

## 4.6 Vector and Matrix Differentiation

$$\frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a} \quad \frac{\partial \mathbf{x}^\top \mathbf{x}}{\partial \mathbf{x}} = 2\mathbf{x}$$

$$\frac{\partial \mathbf{x}^\top \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x} \quad \frac{\partial \mathbf{x}^\top \mathbf{A} \mathbf{y}}{\partial \mathbf{x}} = \mathbf{A} \mathbf{y}$$

$$\frac{\partial |\mathbf{X}|}{\partial x_{ij}} = \begin{cases} X_{ij} & \text{if } i = j \\ 2X_{ij} & \text{otherwise} \end{cases} \quad \begin{matrix} \text{if all elements of } \mathbf{X} \text{ are distinct} \\ \mathbf{X} \text{ is symmetric} \end{matrix}$$

$$\frac{\partial \text{tr} \mathbf{X} \mathbf{Y}}{\partial \mathbf{X}} = \begin{cases} \mathbf{Y}^\top & \text{if all elements of } \mathbf{X} \text{ are distinct} \\ \mathbf{Y} + \mathbf{Y}^\top - \text{Diag}(\mathbf{Y}) & \text{if } \mathbf{X} \text{ is symmetric} \end{cases}$$

$$\frac{\partial \mathbf{X}^{-1}}{\partial x_{ij}} = \begin{cases} \mathbf{X}^{-1} \mathbf{J}_{ij} \mathbf{X}^{-1} & \text{if } i = j \\ \mathbf{X}^{-1} \mathbf{J}_{ij} \mathbf{X}^{-1} & \text{if } i = j \\ \mathbf{X}^{-1} (\mathbf{J}_{ij} + \mathbf{J}_{ji}) \mathbf{X}^{-1} & \text{otherwise} \end{cases} \quad \begin{matrix} \text{if all elements of } \mathbf{X} \text{ are distinct} \\ \mathbf{X} \text{ is symmetric} \end{matrix}$$

where  $X_{ij}$  denotes the cofactor of  $x_{ij}$  in  $\mathbf{X}$  and  $\mathbf{J}_{ij}$  denotes a matrix with 1 in the  $(i, j)$ th place and zeros elsewhere.

## 4.7 Basic Concepts and Facts

1. For  $s = m$  and  $m \geq n$ ,  $|\lambda \mathbf{E} - \mathbf{A} \mathbf{B}| = \lambda^{m-n} |\lambda \mathbf{E} - \mathbf{B} \mathbf{A}|$ .
2.  $\mathbf{A}$  is  $n$ -order matrix. If  $|a_{ii}| \geq \sum_{i \neq j} |a_{ij}|$  for  $i = 1, \dots, n$ , then  $|\mathbf{A}| \neq 0$ ; If  $a_{ii} \geq \sum_{i \neq j} |a_{ij}|$  for  $i = 1, \dots, n$ , then  $|\mathbf{A}| > 0$ .
3.  $\mathbf{A}$  is  $n$ -order matrix. There exists a orthogonal matrix  $\mathbf{P}$  such that  $\mathbf{P}^\top \mathbf{A} \mathbf{P}$  is trigonal matrix if and only if all eigenvalues of  $\mathbf{A}$  are real.
4. If  $\mathbf{A}$  is orthogonal and its eigenvalues are all of real, then  $\mathbf{A}$  must be symmetric and  $\mathbf{A}^2 = \mathbf{E}$ .
5. Suppose that  $\{\mathbf{A}_k\}$  is a sequence of real symmetric matrices, and  $\mathbf{A}_i \mathbf{A}_j = \mathbf{A}_j \mathbf{A}_i$  for  $i \neq j$ . Then there exists a orthogonal matrix  $\mathbf{P}$  such that  $\mathbf{P}^\top \mathbf{A}_k \mathbf{P}$  are all diagonal.
6. Any real inverable matrix  $\mathbf{A}$  can be decomposed into  $\mathbf{S} \mathbf{S}_1 \mathbf{S}_2$ , where  $\mathbf{S}$  is positive,  $\mathbf{S}_1$  and  $\mathbf{S}_2$  are real symmetric orthogonal.

1. A square matrix  $\mathbf{A}$  such that  $\mathbf{A}^k = \mathbf{0}$  for some integer  $k$  is called **nilpotent** matrix, and the smallest positive integral exponent  $k$  such that  $\mathbf{A}^k = \mathbf{0}$  is called the **index** of  $\mathbf{A}$ .

2.  $\mathbf{A} = -\mathbf{A}^\top$  is called **skew-symmetric**.

3. A matrix  $\mathbf{A}$  with complex elements is said to be **Hamiltonian** if  $\mathbf{A} = \mathbf{A}^*$ , and **skew-Hermitian** if  $\mathbf{A} = -\mathbf{A}^*$ .

4.  $\mathbf{A}$  is **equivalent** to  $\mathbf{B}$  if  $\mathbf{B}$  can be obtained from  $\mathbf{A}$  by the successive application of finitely many elementary row and column operations, and we write  $\mathbf{A} \stackrel{\text{E}}{=} \mathbf{B}$ .
5.  $\mathbf{A} \stackrel{\text{E}}{=} \mathbf{B}$  iff  $\mathbf{A} = \mathbf{S} \mathbf{A} \mathbf{T}$ , where  $\mathbf{S}$  and  $\mathbf{T}$  are  $m \times m$  and  $n \times n$  non-singular matrices, respectively.

6. Every nonzero  $\mathbf{A}$  is **equivalent** to  $\begin{bmatrix} \mathbf{I}_r & \mathbf{O}_{r, n-r} \\ \mathbf{O}_{m-r, r} & \mathbf{O}_{m-r, m-r} \end{bmatrix}$ .

7. Matrices of the form  $\begin{bmatrix} \mathbf{I}_r & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{bmatrix}$  are called **canonical matrices**.

8. If  $\mathbf{A}$  and  $\mathbf{B}$  are two square matrices of order  $n$  and  $m$ , respectively, the matrix

$$\mathbf{A} \dot{+} \mathbf{B} = \begin{bmatrix} \mathbf{A} & \mathbf{O} \\ \mathbf{O} & \mathbf{B} \end{bmatrix}$$

is called their **direct sum**.

9.  $\mathbf{A}_{ij}$  is called the **cofactor** of element  $a_{ij}$ .

10. (The Hamilton-Cayley Theorem) Let  $\mathbf{A}$  be an  $n \times n$  matrix and let

$$f(\lambda) \equiv (-1)^n \{\lambda^n - p_1 \lambda^{n-1} + \dots + (-1)^n p_n\}$$

be the characteristic function of  $\mathbf{A}$ . Then

$$\mathbf{A}^n - p_1 \mathbf{A}^{n-1} + \dots + (-1)^n p_n \mathbf{I}_n = \mathbf{O}_n.$$

11. Denote  $\mathbf{I}_n$  by  $\mathbf{A}_0$ , one successively computes  $c_1, \mathbf{A}_1, c_2, \mathbf{A}_2, \dots, \mathbf{A}_{n-1}, c_n$  by the two formulas

$$c_k = (1/k) \text{tr}(\mathbf{A} \mathbf{A}_k) \quad \mathbf{A}_k = \mathbf{A} \mathbf{A}_{k-1} - c_k \mathbf{I}$$

Then  $\mathbf{A}^{-1} = \mathbf{A}_{n-1}/c_n$ .

12. If  $\lambda_1, \lambda_2, \dots, \lambda_n$  are the characteristic roots, distinct or not, of an  $n \times n$  matrix  $\mathbf{A}$ , and if  $g(\mathbf{A})$  is any polynomial function of  $\mathbf{A}$ , then the characteristic roots of  $g(\mathbf{A})$  are  $g(\lambda_1), \dots, g(\lambda_n)$ .
13.  $\mathbf{A}$  is **similar** to  $\mathbf{B}$  if and only if there exists a nonsingular matrix  $\mathbf{P}$  such that  $\mathbf{A} = \mathbf{P}^{-1} \mathbf{B} \mathbf{P}$ . We write  $\mathbf{A} \stackrel{\text{S}}{=} \mathbf{B}$ .
14. Any square matrix  $\mathbf{A}$  is **similar** to an upper triangular matrix whose diagonal elements are the eigenvalues of  $\mathbf{A}$ .
15. A square matrix  $\mathbf{B}$  is said to be **congruent** to  $\mathbf{A}$  if and only if there exists a nonsingular matrix  $\mathbf{P}$  such that  $\mathbf{A} = \mathbf{P}^\top \mathbf{B} \mathbf{P}$ . We write  $\mathbf{A} \stackrel{\text{C}}{=} \mathbf{B}$ .
16. An elementary row operation applied to a square matrix, and followed by the corresponding elementary column operation, is called an **elementary cogredient operation** on the matrix.
17. Every symmetric  $\mathbf{R}$  matrix  $\mathbf{A}$  of rank  $r$  is **congruent** to a matrix of the form  $\text{diag}(\mathbf{I}_r, \mathbf{O})$ .
18. Any  $n \times n$  real symmetric matrix  $\mathbf{A}$  is orthogonal similar to a diagonal matrix whose diagonal elements are the eigenvalues of  $\mathbf{A}$ .

- $\mathbf{A}\mathbf{B} = \mathbf{B}\mathbf{A}$  commutative
- $\mathbf{A}\mathbf{B} = -\mathbf{B}\mathbf{A}$  anti-commute
- $\mathbf{A}^2 = \mathbf{I}$  involutory

## 4.8 Optimization and Inequalities

1. Consider the matrix function  $f$ , where

$$f(\mathbf{X}) = -\log |\mathbf{X}| + \text{tr}(\mathbf{X}^{-1}\mathbf{A}).$$

If  $\mathbf{A} > \mathbf{O}$ , then, subject to  $\mathbf{X} > \mathbf{O}$ ,  $f(\mathbf{X})$  is minimized uniquely at  $\mathbf{X} = \mathbf{A}$ .

2. Let  $f: \Theta \rightarrow f(\Theta)$  be a real-valued function with domain  $\Theta$ , and let  $g: \Theta \rightarrow g(\Theta) = \Phi$  be a bijective (one-to-one) function from  $\Theta$  onto  $\Phi$ . Since  $g$  is bijective, it has an inverse,  $g^{-1}$ , say, and we can define  $h(\phi) = f(g^{-1}(\theta))$  for  $\phi \in \Phi$ .

(a) If  $f(\theta)$  attains a maximum at  $\theta = \hat{\theta}$ ,  $h(\phi)$  attains its maximum at  $\hat{\phi} = g(\hat{\theta})$ .

(b) If the maximum of  $f(\theta)$  occurs uniquely at  $\hat{\theta}$ , then the maximum of  $h(\phi)$  occurs uniquely at  $\hat{\phi}$ .

3. *Frobenius norm approximation.* Let  $\mathbf{B}$  be a  $p \times q$  matrix of rank  $r$  with singular value decomposition  $\sum_{i=1}^r \delta_i l_i m_i^\top$ , and let  $\mathbf{C}$  be a  $p \times q$  matrix of ranks  $s$  ( $s < r$ ). Then

$$\|\mathbf{B} - \mathbf{C}\|^2 = \sum_{i=1}^p \sum_{j=1}^q (b_{ij} - c_{ij})^2$$

is minimized when

$$\mathbf{C} = \mathbf{B}_{(s)} = \sum_{i=1}^s \delta_i l_i m_i^\top.$$

The minimum value is  $\sum_{i=s+1}^r \delta_i^2$ .

4. Let  $\mathbf{A}$  be an  $n \times n$  symmetric matrix with eigenvalues  $\lambda_1 \geq \dots \geq \lambda_n$ , and a corresponding set of orthogonal eigenvectors  $\mathbf{t}_1, \dots, \mathbf{t}_n$ . Define  $T_k = (\mathbf{t}_1, \dots, \mathbf{t}_k)$  ( $k = 1, \dots, n-1$ ) and  $T = (\mathbf{t}_1, \dots, \mathbf{t}_n)$ . Then, if we assume that  $\mathbf{x} \neq 0$ , we have the following:

(a)

$$\sup_{\mathbf{x}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \lambda_1,$$

and the supremum is attained if  $\mathbf{x} = \mathbf{t}_1$ .

(b)

$$\sup_{T_k^\top \mathbf{x} = 0} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \lambda_{k+1},$$

and the supremum is attained if  $\mathbf{x} = \mathbf{t}_{k+1}$ .

(c)

$$\inf_{\mathbf{x}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \lambda_n,$$

and the infimum is attained if  $\mathbf{x} = \mathbf{t}_n$ .

(d) If  $T_{n-k} = (\mathbf{t}_{n-k+1}, \dots, \mathbf{t}_n)$

$$\inf_{T_{n-k}^\top \mathbf{x} = 0} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \lambda_{n-k},$$

and the infimum is attained if  $\mathbf{x} = \mathbf{t}_{n-k}$ .

(e) *Courant-Fischer min-max theorem.*

$$\inf_{L \times k} \sup_{L^\top \mathbf{x} = 0} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \lambda_{k+1},$$

and the result is attained if  $L = T_k$  and  $\mathbf{x} = \mathbf{t}_{k+1}$ .

(f)

$$\sup_{L \times k} \inf_{L^\top \mathbf{x} = 0} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \lambda_{n-k},$$

and the result is attained if  $L = T_{n-k}$  in (d) and  $\mathbf{x} = \mathbf{t}_{n-k}$ .

5. Let  $\mathbf{A}$  be an  $n \times n$  symmetric matrix and let  $\mathbf{D}$  be any  $n \times n$  positive definite matrix. Let  $\gamma_1 \geq \dots \geq \gamma_n$  be eigenvalues of  $\mathbf{D}^{-1}\mathbf{A}$  with corresponding eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_n$ . Then

$$\sup_{\mathbf{x}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{D} \mathbf{x}} = \gamma_1, \quad \text{and} \quad \inf_{\mathbf{x}} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{D} \mathbf{x}} = \gamma_n,$$

with the bounds being attained when  $\mathbf{x} = \mathbf{v}_1$  and  $\mathbf{x} = \mathbf{v}_n$ , respectively.

6. If  $\mathbf{D}$  is positive definite, then for any  $\mathbf{a}$

$$\sup_{\mathbf{x}} \frac{(\mathbf{a}^\top \mathbf{x})^2}{\mathbf{x}^\top \mathbf{D} \mathbf{x}} = \mathbf{a}^\top \mathbf{D}^{-1} \mathbf{a}.$$

The supremum occurs when  $\mathbf{x}$  is proportional to  $\mathbf{D}^{-1}\mathbf{a}$ .

7. Let  $\mathbf{M}$  and  $\mathbf{N}$  be positive definite, then

$$\sup_{\mathbf{x}, \mathbf{y}} \frac{\mathbf{x}^\top \mathbf{L} \mathbf{x}}{\mathbf{x}^\top \mathbf{M} \mathbf{x} \cdot \mathbf{y}^\top \mathbf{N} \mathbf{y}} = \theta_{\max},$$

where  $\theta_{\max}$  is the largest eigenvalue of  $\mathbf{M}^{-1}\mathbf{L}^\top \mathbf{N}^{-1}\mathbf{L}$ . The supremum occurs when  $\mathbf{x}$  is an eigenvector of  $\mathbf{M}^{-1}\mathbf{L}^\top \mathbf{N}^{-1}\mathbf{L}$  corresponding to  $\theta_{\max}$ , and  $\mathbf{y}$  is an eigenvector of  $\mathbf{M}^{-1}\mathbf{L}^\top \mathbf{M}^{-1}\mathbf{L}$  corresponding to  $\theta_{\max}$ .

8. Let  $\mathbf{C}$  be  $p \times q$  matrix of rank  $m$  and let  $\rho_1^2 \geq \dots \geq \rho_m^2 > 0$  be the nonzero eigenvalues of  $\mathbf{C}\mathbf{C}^\top$ . Let  $\mathbf{t}_1, \dots, \mathbf{t}_m$  be the corresponding eigenvectors of  $\mathbf{C}\mathbf{C}^\top$  and let  $\mathbf{w}_1, \dots, \mathbf{w}_m$  be the corresponding eigenvectors of  $\mathbf{C}^\top \mathbf{C}$ . If  $T_k = (\mathbf{t}_1, \dots, \mathbf{t}_k)$  and  $W_k = (\mathbf{w}_1, \dots, \mathbf{w}_k)$  ( $k < m$ ), then

$$\sup_{T_k^\top \mathbf{x} = 0, W_k^\top \mathbf{y} = 0} \frac{(\mathbf{x}^\top \mathbf{C} \mathbf{y})^2}{\mathbf{x}^\top \mathbf{x} \cdot \mathbf{y}^\top \mathbf{y}} = \rho_{k+1}^2,$$

and the supremum occurs when  $\mathbf{x} = \mathbf{t}_{k+1}$  and  $\mathbf{y} = \mathbf{w}_{k+1}$ .

9. Let  $\mathbf{A}$  and  $\mathbf{B}$  be an  $n \times n$  symmetric matrices with eigenvalues  $\rho_1(\mathbf{A}) \geq \dots \geq \rho_n(\mathbf{A})$  and  $\rho_1(\mathbf{B}) \geq \dots \geq \rho_n(\mathbf{B})$ , respectively. If  $\mathbf{A} - \mathbf{B} \geq \mathbf{O}$ , then we have the following:

(a)  $\rho_i(\mathbf{A}) \geq \rho_i(\mathbf{B})$  ( $i = 1, \dots, n$ )

(b)  $\text{tr}(\mathbf{A}) \geq \text{tr}(\mathbf{B})$

(c)  $|\mathbf{A}| \geq |\mathbf{B}|$

(d)  $\|\mathbf{A}\| \geq \|\mathbf{B}\|$ , where  $\|\mathbf{A}\| = \{\text{tr}(\mathbf{A}\mathbf{A}^\top)\}^{1/2}$ .

10. Let  $\mathbf{A}$ ,  $\mathbf{B}$  and  $\mathbf{A} - \mathbf{B}$  be an  $n \times n$  positive semidefinite matrices, with  $r(\mathbf{B}) \leq r$ , and let  $\rho_i(\cdot)$  represent the  $i$ th largest eigenvalues. Then

$$\rho_i(\mathbf{A} - \mathbf{B}) \geq \begin{cases} \rho_{r+i}(\mathbf{A}) & i = 1, \dots, n-r \\ 0 & i = n-r+1, \dots, n \end{cases}$$

Equality occurs if

$$\mathbf{B} = \mathbf{B}_0 = \sum_{i=1}^r \rho_i(\mathbf{A}) \mathbf{t}_i \mathbf{t}_i^\top,$$

where  $\mathbf{t}_1, \dots, \mathbf{t}_n$  are orthogonal eigenvectors corresponding to  $\rho_1(\mathbf{A}), \dots, \rho_n(\mathbf{A})$ .

## 4.9 Jacobians and Transformations

1. If the distinct elements of a symmetric  $d \times d$  matrix  $\mathbf{A}$  have a joint density function of the form  $g(\lambda_1, \dots, \lambda_d)$  where  $\lambda_1 \geq \dots \geq \lambda_d$  are the eigenvalues of  $\mathbf{A}$ , then the joint density function of the eigenvalues is

$$\pi^{d^2/2} g(\lambda_1, \dots, \lambda_d) \left\{ \prod_{j < k} (\lambda_j - \lambda_k) \right\} / \Gamma_d(d/2)$$

where  $\Gamma_d(d/2) = \pi^{d(d-1)/4} \prod_{j=1}^d \Gamma(\frac{1}{2}[d+1-j])$ .

2. Let  $\mathbf{X}$  be an  $m \times n$  matrix of distinct random variables and let  $\mathbf{Z} = a(\mathbf{X})$ , where  $\mathbf{Z}$  is  $m \times n$  and  $a$  is a bijective function. Then there exists an inverse function  $b = a^{-1}$ , so that  $\mathbf{X} = b(\mathbf{Z})$ . If  $\mathbf{X}$  has density  $f$  and  $\mathbf{Z}$  has density  $g$ , then

$$g(\mathbf{Z}) = f(b(\mathbf{Z})) \left| \frac{d\mathbf{X}}{d\mathbf{Z}} \right|,$$

where  $d\mathbf{X}/d\mathbf{Z}$  represents the Jacobian of the transformation from  $\mathbf{X}$  to  $\mathbf{Z}$ .

- (a) If  $\mathbf{X} = \mathbf{A}\mathbf{Z}\mathbf{B}$ , where  $\mathbf{A}$  and  $\mathbf{B}$  are  $m \times m$  and  $n \times n$  nonsingular matrices, respectively, then

$$\frac{d\mathbf{X}}{d\mathbf{Z}} = |\mathbf{A}|^m |\mathbf{B}|^n.$$

- (b) If  $\mathbf{X}$  and  $\mathbf{Z}$  are  $n \times n$  symmetric matrices,  $\mathbf{A}$  is nonsingular matrices, and  $\mathbf{X} = \mathbf{A}\mathbf{Z}\mathbf{A}^\top$

$$\frac{d\mathbf{X}}{d\mathbf{Z}} = |\mathbf{A}|^{n+1}$$

- (c) Let  $\mathbf{E}$  and  $\mathbf{H}$  be  $d \times d$  positive definite matrices, and let  $\mathbf{Z} = \mathbf{E} + \mathbf{H}$  and  $\mathbf{V} = (\mathbf{E} + \mathbf{H})^{-1/2} \mathbf{H} (\mathbf{E} + \mathbf{H})^{-1/2}$ .

$$\frac{d(\mathbf{H}, \mathbf{E})}{d(\mathbf{V}, \mathbf{Z})} = |\mathbf{Z}|^{(d+1)/2}.$$

## 4.10 Generalized Inverse

### 4.10.1 $g$ -inverse

If  $\mathbf{AGA} = \mathbf{A}$ ,  $\mathbf{G}$  is called a generalized inverse ( $g$ -inverse).

1. For each matrix  $\mathbf{A} = \mathbf{A}_L \mathbf{A}_R \stackrel{\Delta}{=} \mathbf{B}\mathbf{C}$ ,  $\mathbf{A}^- = \mathbf{C}^\top (\mathbf{C}\mathbf{C}^\top)^{-1} (\mathbf{B}^\top \mathbf{B})^{-1} \mathbf{B}^\top$ .
2.  $r(\mathbf{A}) = r(\mathbf{A}^-) = r(\mathbf{A}\mathbf{A}^-) = r(\mathbf{A}^- \mathbf{A})$ .
3. If  $\mathbf{A}$  is an  $m \times n$  matrix of rank  $m$ , then  $\mathbf{A}^- = \mathbf{A}^\top (\mathbf{A}\mathbf{A}^\top)^{-1}$  (**right-inverse**) and  $\mathbf{A}\mathbf{A}^- = \mathbf{I}_m$ . If  $r(\mathbf{A}) = n$ , then  $\mathbf{A}^- = (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top$  (**left-inverse**) and  $\mathbf{A}^- \mathbf{A} = \mathbf{I}_n$ .
4. It is not always true that  $(\mathbf{G}\mathbf{H})^- = \mathbf{H}^- \mathbf{G}^-$  for all matrices  $\mathbf{H}, \mathbf{G}$ .
5. Let  $\mathbf{B}$  be an  $m \times r$  matrix of rank  $r$  and  $\mathbf{C}$  be an  $r \times m$  matrix of rank  $r$ ; then  $(\mathbf{B}\mathbf{C})^- = \mathbf{C}^- \mathbf{B}^-$ .
6.  $(\mathbf{A}^\top \mathbf{A})^- = \mathbf{A}^- (\mathbf{A}^\top)^-$  for any matrix  $\mathbf{A}$ .
7. Let  $\mathbf{P}$  be an  $m \times m$  orthogonal matrix,  $\mathbf{Q}$  be an  $n \times n$  orthogonal matrix, and  $\mathbf{A}$  is any  $m \times n$  matrix. Then  $(\mathbf{P}\mathbf{A}\mathbf{Q})^- = \mathbf{Q}^- \mathbf{A}^- \mathbf{P}^-$ .

#### 4.10.2 $c$ -inverse

1.  $r(\mathbf{X}^c) \geq r(\mathbf{X}) = r(\mathbf{X}\mathbf{X}^c) = r(\mathbf{X}^c \mathbf{X})$  for any matrix  $\mathbf{X}$ .
2.  $\mathbf{X}^c \mathbf{X}$  and  $\mathbf{X}\mathbf{X}^c$  are idempotent matrices.
3. If  $\mathbf{X}^c$  is any  $c$ -inverse of  $\mathbf{X}$ , then  $(\mathbf{X}^c)^\top$  is a  $c$ -inverse of  $\mathbf{X}^\top$ .
4. For any  $m \times n$  matrix  $\mathbf{X}$  of rank  $r > 0$ , define

$$\mathbf{K} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^c \mathbf{X}^\top.$$

Then  $\mathbf{K}$  is invariant for any  $c$ -inverse of  $\mathbf{X}^\top \mathbf{X}$ .

5.  $\mathbf{X}(\mathbf{X}^\top \mathbf{X})^c \mathbf{X}^\top = \mathbf{X}\mathbf{X}^-$  for any  $c$ -inverse  $(\mathbf{X}^\top \mathbf{X})^c$  of  $\mathbf{X}^\top \mathbf{X}$ .
6.  $r(\mathbf{K}) = r(\mathbf{X})$ .
7.  $\mathbf{K}\mathbf{X} = \mathbf{X}; \mathbf{X}^\top \mathbf{K} = \mathbf{X}^\top$
8.  $(\mathbf{X}^\top \mathbf{X})^c \mathbf{X}^\top$  is a  $c$ -inverse of  $\mathbf{X}$  for any  $c$ -inverse of  $\mathbf{X}^\top \mathbf{X}$ .
9.  $\mathbf{X}(\mathbf{X}^\top \mathbf{X})^c$  is a  $c$ -inverse of  $\mathbf{X}^\top$  for any  $c$ -inverse of  $\mathbf{X}^\top \mathbf{X}$ .

### 4.11 Linear Equations

Let  $\mathbf{A}$  be an  $m \times n$  matrix and  $\mathbf{A}^c$  be any  $c$ -inverse of  $\mathbf{A}$ . Suppose a solution exists to the system  $\mathbf{A}\mathbf{x} = \mathbf{g}$ . For each  $n \times 1$  vector  $\mathbf{h}$ , the vector  $\mathbf{x}_0$  is a solution, where

$$\mathbf{x}_0 = \mathbf{A}^c \mathbf{g} + (\mathbf{I}_n - \mathbf{A}^c \mathbf{A}) \mathbf{h}. \quad (1)$$

Also, every solution to the system can be written in the form of Equation (1) for some  $n \times 1$  vector  $\mathbf{h}$ .

1. If  $\mathbf{A}$  is an  $m \times m$  symmetric matrix such that  $\mathbf{1}^\top \mathbf{A} = \mathbf{0}$ , then

$$\begin{bmatrix} \mathbf{A} \\ \mathbf{1}^\top \end{bmatrix}^- = \begin{bmatrix} \mathbf{A}^- \\ \frac{1}{m} \mathbf{1} \end{bmatrix}$$

2. If  $\mathbf{A}$  is an  $m \times m$  symmetric matrix of rank  $m - 1$  such that  $\mathbf{1}^\top \mathbf{A} = \mathbf{0}$ , then  $\mathbf{B} = \mathbf{A} + \mathbf{1}\mathbf{1}^\top/n$  is nonsingular and its inverse is  $\mathbf{A}^- + \mathbf{J}/n$ ;  
Meanwhile

$$\begin{bmatrix} \mathbf{A} & \mathbf{1} \\ \mathbf{1}^\top & 0 \end{bmatrix}^- = \begin{bmatrix} \mathbf{A}^- & \frac{1}{n} \mathbf{1} \\ \frac{1}{n} \mathbf{1}^\top & 0 \end{bmatrix}.$$

## 5 Analysis

analysis.tex

1.  $f(x)$  is bound and  $g(x)$  is differentiable on  $[a, b]$ .  $g(\lambda) = 0$  for some  $\lambda \neq 0$ , if  $|g(x)f(x) + \lambda g'(x)| \leq |g(x)|$ , then  $g(x) = 0$  for all  $x \in [a, b]$ .
2.  $f(x)$  is monotone on  $[0, \infty]$  and  $\int_0^\infty f(x)dx$  is well-defined. Then

$$\lim_{h \rightarrow 0^+} h \sum_{n=1}^{\infty} f(nh) = \int_0^\infty f(x)dx.$$

3.  $f(x)$  is  $2n$ -differentiable on  $[a, b]$ , and  $|f^{(2n)}(x)| \leq M$ ,  $f^{(m)}(a) = f^{(m)}(b) = 0$  for  $m = 0, \dots, n-1$ . Then

$$|\int_a^b f(x)dx| \leq \frac{(n!)^2 M}{(2n)!(2n+1)!} (b-a)^{2n+1}.$$

4.  $f(x)$  and  $g(x)$  are bound on any sub-interval of  $[0, +\infty)$ , and satisfy that  $g(x+T) > g(x)$  for some  $T > 0$  and any  $x > 0$ ,  $g(x) \rightarrow +\infty$ . In addition,

$$\lim_{x \rightarrow \infty} \frac{f(x+T) - f(x)}{g(x+T) - g(x)} = l,$$

where  $l$  may be  $+\infty$ . Then

$$\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = l.$$

5.  $f(x)$  and  $g(x)$  are bound on any sub-interval of  $[0, +\infty)$ , and satisfy that  $0 < g(x+T) < g(x)$  for some  $T > 0$  and any  $x > 0$ ,  $\lim_{x \rightarrow \infty} g(x) = 0$ . In addition,

$$\lim_{x \rightarrow \infty} \frac{f(x+T) - f(x)}{g(x+T) - g(x)} = l,$$

where  $l$  may be  $+\infty$ . Then

$$\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = l.$$

6.  $f'(x)$  is absolutely continuous on  $[a, b]$ , then for any  $c \in (a, b)$  and  $p > 1$

$$\int_a^b |f''(x)|^p dx \geq \left\{ \frac{p-1}{2p-1} (b-a) \right\}^{1-p} \left| \frac{f(b)-f(c)}{b-c} - \frac{f(c)-f(a)}{c-a} \right|^p.$$

7. Positive series  $\sum u_n$

- comparing principle;
- integral decision;
- $\lim u_n/u_{n+1} = \rho$ ,  $\rho > 1$  converges and  $\rho < 1$  diverges;
- Cauchy criterion:  $\limsup u_n^{1/n} = \rho$ ,  $\rho > 1$  converges and  $\rho < 1$  diverges;
- $u_n/u_{n+1} = \lambda + \mu/n + o(\theta_n/n^{1+\mu})$ ,  $\lambda > 1$  converges,  $\lambda < 1$  diverges,  $\lambda = 1$   $\mu > 1$  converges and  $\mu < 1$  diverges;
- $\sum |b_n| < \infty$ ,  $u_n/u_{n+1} = \lambda + 1/n + o(b_n)$ , diverges;
- $u_{n+1}/u_n = 1 - \alpha/n + O(1/n^\lambda)$ ,  $\alpha > 1$  converges,  $\alpha < 1$  diverges;

- $u_{n+1}/u_n = 1 - 1/n - \alpha'_n/n \log n$ ,  $\alpha'_n \geq \alpha > 1$  converges,  $\alpha'_n \leq \alpha < 1$  diverges.

8. (Roll Theorem)  $f'(x)$  is bound on finite or infinite interval  $(a, b)$ , and  $\lim_{x \rightarrow a^+} f(x) = \lim_{x \rightarrow b^-} f(x)$ . Then there exists at least one  $c \in (a, b)$  such that  $f'(c) = 0$ .

9. (Dabu Theorem)  $m < |f'(x)| < M$  Then for any  $\mu \in (m, M)$ , there exists  $x_\mu \in (a, b)$  such that  $f'(x_\mu) = \mu$ .

10.  $f(x)$  is differentiable on and  $|f'(x)| \leq M$  on  $(a, b)$ , then

$$\left| \frac{1}{b-a} \int_a^b f(x)dx - \frac{f(b)+f(a)}{2} \right| \leq \frac{M(b-a)}{4} \left\{ 1 - \left( \frac{f(b)-f(a)}{M(b-a)} \right)^2 \right\}.$$

11. From Joseph Edward (1954).

$$\int_0^1 \frac{(-\log x)^p}{(1-x)^2} dx = \int_0^\infty y^p \sum_{n=1}^{\infty} n y^{-ny} dy = p! \sum_{k=1}^{\infty} \frac{1}{k^p}$$

- 12.

$$\int_0^1 \frac{(-\log x)^p}{(1+x)^2} dx = p!(1-2^{1-p}) \sum_{k=1}^{\infty} \frac{1}{k^p}$$

#### The relationship of some important inequalities

Denote  $M_r(a) = (\sum_{i=1}^n p_i a_i^r)^{1/r}$  for  $a = (a_1, \dots, a_n)^\top$  with  $a_i \geq 0$  and  $r \geq 1$ ,  $\{p_i\}$  are weight ( $\sum p_i = 1$ ).  $G(a) = \prod a_i^{p_i}$ .

Cauchy  $\implies M_{2r}(a) > M_r(a) \implies M_r(a) \rightarrow G(a)$  by letting  $r \rightarrow 0 \implies$  Hölder

$$\sum_{k=1}^n a_k^\alpha b_k^\beta \dots l_k^\lambda \leq \left( \sum_{k=1}^n a_k \right)^\alpha \left( \sum_{k=1}^n b_k \right)^\beta \dots \left( \sum_{k=1}^n l_k \right)^\lambda$$

for  $\alpha, \beta, \dots, \lambda > 0$  and  $\alpha + \beta + \dots + \lambda = 1$ .

$$\text{Hölder} \Rightarrow \begin{cases} \text{(Minkowski)} \left\{ \sum_{j=1}^k a_j^r \right\}^{1/r} \leq \left\{ \sum_{j=1}^k a_j^s \right\}^{1/s} & r > 1 \\ \text{(Jensen)} \left( \sum_1^n a_k^s \right)^{1/s} \leq \left( \sum_1^n a_k^r \right)^{1/r} & 0 < r < s \\ \text{(Liap)} \{M_s(a)\}^s \leq [\{M_r(a)\}^r]^{\frac{(s-r)}{r-s}} [\{M_t(a)\}^t]^{\frac{(s-r)}{t-r}} & 0 < r < s \\ \text{(Increasing property of)} f(x) = M_x(a) \end{cases}$$

## 6 Good Examples in Linear Models

linear.tex

1. Consider the linear model

$$\underline{Y} = \underline{X}\underline{\beta} + \underline{Z}\underline{\gamma} + \underline{e} \quad \underline{e} \sim N(\underline{0}, \sigma^2 \mathbf{I}).$$

where  $\underline{Y}$  is  $n \times 1$ ,  $\underline{\gamma}$   $q \times 1$ ,  $\underline{\beta}$  is  $p \times 1$ ,  $\underline{X}$  is  $n \times p$ ,  $\underline{Z}$  is  $n \times q$ ,  $[\underline{X}, \underline{Z}]$  is of rank  $p+q$  and  $n > p+q$ .

- (a) Show that  $\underline{Z}^\top (\mathbf{I} - \underline{X}\underline{X}^\top) \underline{Z}$  is positive

$$\begin{aligned} r[\underline{Z}^\top (\mathbf{I} - \underline{X}(\underline{X}^\top \underline{X})^{-1} \underline{X}^\top)] &= r\{(\underline{X}, \underline{Z})^\top [\mathbf{I} - \underline{X}(\underline{X}^\top \underline{X})^{-1} \underline{X}^\top]\} \geq \\ &= r\{(\underline{X}, \underline{Z})^\top\} + r\{\mathbf{I} - \underline{X}(\underline{X}^\top \underline{X})^{-1} \underline{X}^\top\} - n = \\ &= p+q + (n-p) - n = q \end{aligned}$$

$\underline{Z}^\top [\mathbf{I} - \underline{X}(\underline{X}^\top \underline{X})^{-1} \underline{X}^\top] \underline{Z}$  is full rank and nonnegative definite.

- (b) Show that the MLE of  $\underline{\beta}$  and  $\underline{\gamma}$  are:

$$\hat{\underline{\beta}} = (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top (\underline{Y} - \underline{Z}\hat{\underline{\gamma}}) \text{ and } \hat{\underline{\gamma}} = \{\underline{Z}^\top (\mathbf{I} - \underline{X}\underline{X}^\top) \underline{Z}\}^{-1} \underline{Z}^\top (\mathbf{I} - \underline{X}\underline{X}^\top) \underline{Y}$$

- (c) How would the estimators of  $\underline{\beta}$  and  $\underline{\gamma}$  change if  $\underline{X}$  and  $\underline{Z}$  were orthogonal? Find the joint distribution of the estimators.

2. Consider the linear model

$$\begin{bmatrix} \underline{Y}_1 \\ \underline{Y}_2 \end{bmatrix} = \begin{bmatrix} \underline{X}_1 & 0 \\ 0 & \underline{X}_2 \end{bmatrix} \begin{bmatrix} \underline{\beta}_1 \\ \underline{\beta}_2 \end{bmatrix} + \underline{e} \quad \underline{e} \sim N(\underline{0}, \sigma^2 \mathbf{I}).$$

where  $\underline{Y}_i$  is  $n_i \times 1$ ,  $\underline{\beta}_i$  is  $p \times 1$ ,  $\underline{X}_i$  is  $n_i \times p$  of rank  $p$  and  $n_1 + n_2 = n$ . Now, consider the following three estimators

$$\hat{\underline{\beta}}_1 = (\underline{X}_1^\top \underline{X}_1)^{-1} \underline{X}_1^\top \underline{Y}_1, \quad \hat{\underline{\beta}}_2 = (\underline{X}_2^\top \underline{X}_2)^{-1} \underline{X}_2^\top \underline{Y}_2, \quad \&\hat{\underline{\beta}} = (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \underline{Y}$$

We wish to test  $H_0 : \underline{\beta}_1 = \underline{\beta}_2$  vs  $H_a : \underline{\beta}_1 \neq \underline{\beta}_2$ . Two possible test statistics are

$$F_1 = \frac{(\hat{\underline{\beta}}_1 - \hat{\underline{\beta}}_2)^\top \{(\underline{X}_1^\top \underline{X}_1)^{-1} + (\underline{X}_2^\top \underline{X}_2)^{-1}\}^{-1} (\hat{\underline{\beta}}_1 - \hat{\underline{\beta}}_2)}{SSE_1 + SSE_2} \times \frac{n-2p}{p}$$

and

$$F_2 = \frac{(\hat{\underline{\beta}} - \hat{\underline{\beta}}_1)^\top \{(\underline{X}_1^\top \underline{X}_1)^{-1} - (\underline{X}^\top \underline{X})^{-1}\}^{-1} (\hat{\underline{\beta}} - \hat{\underline{\beta}}_1)}{SSE_1 + SSE_2} \times \frac{n-2p}{p}$$

- (a) Show that

$$\begin{aligned} &(\underline{X}_1^\top \underline{X}_1) \{(\underline{X}_1^\top \underline{X}_1)^{-1} + (\underline{X}_2^\top \underline{X}_2)^{-1}\} (\underline{X}_1^\top \underline{X}_1)^{-1} \\ &= (\underline{X}^\top \underline{X}) (\underline{X}_2^\top \underline{X}_2)^{-1} \{(\underline{X}_1^\top \underline{X}_1)^{-1} + (\underline{X}_2^\top \underline{X}_2)^{-1}\}^{-1} (\underline{X}_2^\top \underline{X}_2)^{-1} \end{aligned} \quad \text{Df1.8}$$

- (b) Show that  $F_1 = F_2$ .

3. Consider the linear model  $\underline{Y} = \underline{X}\underline{\beta} + \underline{e}$ , where  $\underline{X}$  is  $n \times p$  of rank  $p$  and  $\underline{e} \sim N(0, \sigma^2 \mathbf{I})$ . Partition  $\underline{X}$  into  $\underline{X} = [\underline{X}_1 | \underline{X}_2 | \underline{X}_3]$ , where  $\underline{X}_i$  is  $n_i \times p$  of rank  $p_i$ . Equivalently, let  $\underline{\beta}^\top = [\underline{\beta}_1^\top | \underline{\beta}_2^\top | \underline{\beta}_3^\top]$ . Thus  $\underline{Y} = \underline{X}_1 \underline{\beta}_1 + \underline{X}_2 \underline{\beta}_2 + \underline{X}_3 \underline{\beta}_3 + \underline{e}$

- (a) Show that successively fitting the model  $\underline{Y} = \underline{e}$ ,  $\underline{Y} = \underline{X}_1 \underline{\beta}_1 + \underline{e}$ ,  $\underline{Y} = \underline{X}_1 \underline{\beta}_1 + \underline{X}_2 \underline{\beta}_2 + \underline{e}$ , and  $\underline{Y} = \underline{X}\underline{\beta} + \underline{e}$  yields SS's for  $\underline{\beta}_1$ ,  $\underline{\beta}_2$  and  $\underline{\beta}_3$  which are orthogonal.

- (b) Prove that the SSE for the model containing  $\underline{\beta}$  must be at least as small as the SSE for the model with only  $\underline{\beta}_1$ .

- (c) Now, let  $\underline{e} \sim N(0, \sigma^2 \mathbf{V})$ ,  $\hat{\underline{\beta}}_{ols} = (\underline{X}^\top \underline{X})^{-1} (\underline{X}^\top \underline{Y})$  and  $\hat{\underline{\beta}}_{wls} = (\underline{X}^\top \mathbf{V}^{-1} \underline{X})^{-1} (\underline{X}^\top \mathbf{V}^{-1} \underline{Y})$ . Show that  $V(\hat{\underline{\beta}}_{ols}) - V(\hat{\underline{\beta}}_{wls})$  is nonnegative definite.

$$\begin{aligned} &(\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \mathbf{V} \underline{X} (\underline{X}^\top \underline{X})^{-1} \geq (\underline{X}^\top \mathbf{V}^{-1} \underline{X})^{-1} \\ \iff &(\underline{X}^\top \underline{X}) (\underline{X}^\top \mathbf{V} \underline{X})^{-1} (\underline{X}^\top \underline{X}) \leq (\underline{X}^\top \mathbf{V}^{-1} \underline{X}) \\ \iff &\underline{X}^\top \{\mathbf{V}^{-1} - \underline{X} (\underline{X}^\top \mathbf{V} \underline{X})^{-1} \underline{X}^\top\} \underline{X} \geq 0 \\ \iff &\mathbf{V}^{-1} - \underline{X} (\underline{X}^\top \mathbf{V} \underline{X})^{-1} \underline{X}^\top \geq 0 \\ \iff &\mathbf{I} - \mathbf{V}^{1/2} \underline{X} (\underline{X}^\top \mathbf{V} \underline{X})^{-1} \underline{X}^\top \mathbf{V}^{1/2} \geq 0 \end{aligned}$$

Note that  $\mathbf{I} - \mathbf{V}^{1/2} \underline{X} (\underline{X}^\top \mathbf{V} \underline{X})^{-1} \underline{X}^\top \mathbf{V}^{1/2}$  is a symmetric and idempotent matrix. The conclusion follows.

4. Consider two models for  $Y$  where for both models  $e_i$  are iid  $N(0, \sigma^2)$  and  $\underline{X} = [\underline{X}_1 | \underline{X}_2]$  in  $n \times (p_1 + p_2)$  of rank  $(p_1 + p_2)$ .

$$M1 : \underline{Y} = \underline{X}_1 \underline{\beta}_1 + \underline{e}$$

$$M2 : \underline{Y} = \underline{X}_1 \underline{\beta}_1 + \underline{X}_2 \underline{\beta}_2 + \underline{e}$$

- (a) Under what conditions does the estimator of  $\underline{\beta}_1$  using  $M1$  equal the estimator of  $\underline{\beta}_1$  using  $M2$ ?

- (b) Prove that the SSE under  $M2$  is less than or equal to SSE under  $M1$ .

Note that

$$\begin{aligned} &\text{sth wrong—} \underline{X}_1^\top \{ \underline{X}_2 (\underline{X}_2^\top \underline{X}_2)^{-1} \underline{X}_2^\top \underline{X}_1 [\underline{X}_1^\top (\mathbf{I} - \underline{H}_{\underline{X}_2}) \underline{X}_1]^{-1} \} \underline{X}_1^\top + \\ &\underline{X}_2 \{ \underline{X}_2^\top (\mathbf{I} - \underline{H}_{\underline{X}_1}) \underline{X}_2 \}^{-1} \underline{X}_2^\top \text{ is nonnegative definite.} \end{aligned}$$

—HL

$$SSE_2 - SSE_1 = \underline{Y}^\top (\underline{H}_x - \underline{H}_{x_1}) \underline{Y} = \underline{y} \tilde{\underline{X}}_2 (\tilde{\underline{X}}_2^\top \tilde{\underline{X}}_2)^{-1} \tilde{\underline{X}}_2^\top \underline{Y}$$

where  $\tilde{\underline{X}}_2 = (\mathbf{I} - \underline{H}_{x_1}) \underline{X}_2$  (Th7.1 of Ronald, p247).

## 7 Essential Probability

- (i) A  $\pi$ -class  $\mathcal{P}$  is a class of subsets of  $\Omega$  such that  $A, B \in \mathcal{P}$  implies  $A \cap B \in \mathcal{P}$

- (ii) A semi-field  $\mathcal{S}$  is a class of subsets of  $\Omega$  such that  $\mathcal{S}$  is closed under finite intersections, and  $A \in \mathcal{S}$  implies  $A^c = \bigcap_{i=1}^m B_i$  where  $B_i \in \mathcal{S}$  and disjoint and  $m < \infty$ .

- iii)  $\lambda$ -class  $\mathcal{L} \iff \Omega \in \mathcal{L}; A, B \in \mathcal{L}$  and  $A \subset B$  implies  $BA^c \in \mathcal{L}; A_n \in \mathcal{L}$  and  $A_n \uparrow A$  implies  $A \in \mathcal{L}$ .

- Th1.11 (**Dynkin's Class Theorem**) Suppose  $\mathcal{P}$  is a  $\pi$ -class for  $\Omega$ . Then  $\sigma(\mathcal{P}) = \lambda(\mathcal{P})$ .

- Th1.12 Let  $(\Omega, \mathcal{F}_i, \mu_i)$  be measure spaces ( $i=1,2$ ). Suppose  $\mathcal{P}$  is a  $\pi$ -class such that  $\mathcal{P} \subset \mathcal{F}_i$ ,  $\mu_1$  and  $\mu_2$  agree on  $\mathcal{P}$  and there exist  $A_n \uparrow \Omega$  with  $A_n \in \mathcal{P}$  and  $\mu_i(A_n) < \infty$ . Then  $\mu_1$  and  $\mu_2$  agree on  $\sigma(\mathcal{P})$ .

- Th1.14 (**Carathéodory's Extension Theorem**) Suppose  $\mathcal{F}$  is a field of subsets of  $\Omega$  and  $\mu : \mathcal{P} \rightarrow R^+$ . If  $\mu(\emptyset) = 0$ ;  $\mu(A) \geq 0$  for all  $A \in \mathcal{F}$ ; if  $A_i \in \mathcal{F}$  disjoint and  $A = \bigcup A_i$  is in  $\mathcal{F}$  then  $\mu(A) = \sum \mu(A_i)$ . Then there exists a unique extension of  $\mu$  to  $\sigma(\mathcal{F})$ .

- Th1.16  $f$  is measurable iff  $f^{-1}((-\infty, x]) \in \mathcal{F}$  for every  $x \in R$  and  $f^{-1}(\{-\infty\}) \in \mathcal{F}$ ,  $f^{-1}(\{\infty\}) \in \mathcal{F}$ .

- Th1.17  $f$  is measurable iff it is the pointwise limit of simple functions.

- Th1.18 If  $f_n$  are measurable, then  $\lim f_n$  is measurable;  $f_1 + f_2$  is also measurable; Continuous and monotone functions are **Borel** measurable.

- Df1.21 Let  $f : \Omega \rightarrow \bar{R}$  be  $\mathcal{F}$  measurable and  $\mu$  be  $\sigma$ -finite (i)  $f = \sum_{i=1}^m a_i 1_{A_i}$  is simple. Then  $\int f d\mu = \sum a_i \mu(A_i)$ .

- (ii)  $f \geq 0$  and  $f_n \uparrow f$ , where  $f_n \geq 0$  is simple.  $\int f d\mu = \lim \int f_n d\mu$ .

- (iii)  $f$  is measurable.  $\int f d\mu = \int f_+ d\mu - \int f_- d\mu$ .

- Th1.22 (vii) If  $f \geq 0$  a.e. and  $\int f d\mu < \infty$  then  $f < \infty$  a.e.

- (viii) If  $f \geq 0$  a.e. and  $\mu(\{\omega : f(\omega) > 0\}) > 0$  then  $\int f d\mu > 0$ .

- Th1.24 If  $g$  is Riemann integrable on  $[a, b]$  then it is Lebesgue integrable on  $[a, b]$  (it is also bounded and continuous a.e., i.e. let  $A = \{x : x_n \rightarrow x \text{ but } g(x_n) \not\rightarrow g(x)\}$  then  $\mu(A) = 0$ . ( $[a, b]$  has to be bounded, otherwise not true. For example,  $\int_0^\infty \frac{\sin x}{x} dx = \pi$  but  $\int_0^\infty \frac{\sin x}{x} d\mu$  doesn't exist)

- Th1.26 (**Monotone Convergence**) Suppose  $f_n \geq 0$  measurable and  $f_n \uparrow f$  a.s., then  $\int f_n d\mu \uparrow \int f d\mu$ .

- Th1.27 (**Fatou's Theorem**)  $f_n \geq 0$  measurable, then

$$\liminf_{n \rightarrow \infty} \int f_n d\mu \geq \int \liminf_{n \rightarrow \infty} f_n d\mu.$$

$$f_n \leq f \text{ integrable, then } \limsup_{n \rightarrow \infty} \int f_n d\mu \leq \int \limsup_{n \rightarrow \infty} f_n d\mu.$$

- Th1.28 (**Dominated Convergence Theorem**)  $|f_n| \leq g$  and  $g$  is integrable.  $f_n \rightarrow f$  means  $\lim_{n \rightarrow \infty} \int f_n d\mu = \int f d\mu$ .

- Th1.29 (**Extended DCT**)  $|f_n| \leq g_n \rightarrow g$ ,  $f_n \rightarrow f$ .  $g_n$  and  $g$  integrable and  $\lim_{n \rightarrow \infty} \int g_n d\mu = \int g d\mu$ . Then  $\lim_{n \rightarrow \infty} \int f_n d\mu = \int f d\mu$ .

- Th1.33 (**Fubini's Theorem**)  $\mu_1 \times \mu_2$  is a product measure on  $(\Omega_1 \times \Omega_2, \mathcal{F}_1 \times \mathcal{F}_2)$ .  $\mu_i$  is  $\sigma$ -finite.  $f(w_1, w_2)$  is  $\mathcal{F}_1 \times \mathcal{F}_2$  measurable and is either non-negative or  $\mu_1 \times \mu_2$  integrable. Then

$$\begin{aligned} \int_{\Omega_1 \times \Omega_2} f(w_1, w_2) d(\mu_1 \times \mu_2) &= \int_{\Omega_1} \left\{ \int_{\Omega_2} f(w_1, w_2) d\mu_2 \right\} \mu_1 \\ &= \int_{\Omega_2} \left\{ \int_{\Omega_1} f(w_1, w_2) d\mu_1 \right\} \mu_2 \end{aligned}$$

- Th1.35 Let  $(\Omega, \mathcal{F}, \mu)$  be a measure space. If  $f, g \geq 0$  and  $\mu$  is a  $\sigma$ -finite, then  $\int_A f d\mu = \int_A g d\mu$  for all  $A \in \mathcal{F} \iff f = g$  a.e. ( $\mu$ ).

- If  $f, g$  are  $\mu$ -integrable and  $\mathcal{P}$  is  $\pi$ -class generating  $\mathcal{F}$ , then  $\int_A f d\mu = \int_A g d\mu$  for all  $A \in \mathcal{P} \iff f = g$  a.e. ( $\mu$ ).

## 7.1 Random Variables

- Th2.7  $X \sim F$  and  $r > 0$ . Then  $E(|X|^r) = r \int_0^\infty x^{r-1} \{1 - F(x) + F(-x)\} dx$ .

- Th2.10 (**Minkowski Inequality**)

$$(\int |f+g|^r d\mu)^{1/r} \leq (\int |f|^r d\mu)^{1/r} + (\int |g|^r d\mu)^{1/r}$$



Co2.19 Let  $X_1, X_2, \dots$ , be independent random variables (or vectors). If  $g_1, g_2, \dots$  are measurable, then  $g_1(X_1), g_2(X_2), \dots$  are also independent.

$A_1, A_2, \dots$  independent  $\iff B_1, B_2, \dots$  independent if each  $B_n \in \{\emptyset, A_n, A_n^c, \Omega\}$

Th2.20 Let  $F_1, \dots$  be probability d.f.'s. There exists  $(\Omega, \mathcal{F}, P)$  and a sequence of **independent** r.v.'s  $X_1, \dots$  such that  $X_i \sim F_i$  for all  $i \geq 1$ .

## 7.2 Convergence of Random Variables

Ex3.8 (A useful lower bound) Let  $Y \geq 0$  with  $E(Y^2) < \infty$  and let  $a < E(Y)$ . Then we have  $P(Y > a) \geq (EY - a)^2 / EY^2$ . This is often used with  $a = 0$ .

**Proof.** Using Cauchy-Schwarz inequality to  $Y1_{(Y>a)}$ , we obtain that

$$P(Y > a) \geq \frac{(EY1_{(Y>a)})^2}{EY^2} = \frac{(EY - EY1_{(Y \leq a)})^2}{EY^2} \geq \frac{(EY - a)^2}{EY^2}$$

Th3.8

- Suppose  $X_n, X$  are a.s. finite and  $g$  is continuous. Then  $X_n \rightarrow X$  a.s. (pr) implies  $g(X_n) \rightarrow g(X)$  a.s. (pr)
- Suppose  $X_n, Y_n, X, Y$  are a.s. finite and  $X_n \rightarrow X$  a.s. (pr) and  $Y_n \rightarrow Y$  a.s. (pr). Then  $X_n + Y_n \rightarrow X + Y$  a.s. (pr) and  $\max(X_n, Y_n) \rightarrow \max(X, Y)$  a.s. (pr)
- If  $X_n, X, M_n$  are a.s. finite and  $M_n \rightarrow \infty$  and  $X_n \rightarrow X$  a.s. Then  $X_{M_n} \rightarrow X$  a.s. (not true for probability convergence).

Th3.9 (Borel-Cantelli Theorem).

- $\sum_{n=1}^{\infty} P(A_n) < \infty$ , then  $P(\overline{\lim} A_n) = 0$ .
- If  $A_n$  are independent, then  $\sum_{n=1}^{\infty} P(A_n) = \infty$  if and only if  $P\{\overline{\lim} A_n\} = 1$  or  $P\{\underline{\lim} A_n\} = 0$ .

Th3.10  $X_n \xrightarrow{pr} X \iff \forall$  subsequence  $n_k \rightarrow \infty$  there exist a further subsequence  $n_{k_j} \rightarrow \infty$  such that  $X_{n_{k_j}} \rightarrow X$  a.s. (Useful result)

Th3.11  $X_n \xrightarrow{pr} X$  and  $|X_n| \leq Y$  and  $E(Y) < \infty$ . Then  $E(X_n) \rightarrow E(X)$ .

Th3.17 (Glivenko-Cantelli Theorem)  $X_n$  iid  $\sim F$  with empirical distribution  $F_n$ . Then  $P\{\sup_x |F_n(x) - F(x)| \rightarrow 0\} = 1$ .

Df3.18 (Tail events) Let  $X_1, \dots$  be r.v.'s on  $(\Omega, \mathcal{F}, P)$ . An event  $A$  is a tail event for  $\{X_n\}$  if  $A \in \sigma(X_n, X_{n+1}, \dots)$  for every  $n$ .  $\{\lim X_n \text{ exists}\}$ ,  $\{\lim X_n = X\}$ ,  $\{\limsup X_n \leq x\}$ ,  $\{\sum^{\infty} |X_n| < \infty\}$  are all tail events, but  $\{\sum^{\infty} X_n \leq x\}$  not.

Df3.20  $\{X_n\}$  is uniformly integrable if  $\sup_n \int_{|X_n| > a} |X_n| dP \rightarrow 0$  as  $a \rightarrow \infty$ .

Th3.21 If either of following holds,  $\{X_n\}$  is uniformly integrable

- $|X_n| \leq Y$  a.s. for each  $n$  and  $E|Y| < \infty$ ;
- $|X_n| \leq Y_n$  and  $\{Y_n\}$  is uniformly integrable;
- $\sup_n E|X_n|^{1+\epsilon} < \infty$  for some  $\epsilon > 0$ ;
- $\sup_n E(|g(X_n)|) < \infty$  and  $|g(x)|/|x| \rightarrow \infty$  as  $|x| \rightarrow \infty$

Th3.22  $\{X_n\}$  is uniformly integrable  $\iff$  (i)(uniformly bounded)  $\sup_n E(|X_n|) < \infty$ , and (ii)(uniformly continuous) for each  $\epsilon > 0$ , there is  $\delta > 0$  such that  $P(A) < \delta \implies \sup_n \int_A |X_n| dP < \epsilon$ .

Th3.23 Suppose  $0 < p < \infty$  and  $X_n \xrightarrow{pr} X$ . Then the following are equivalent.

- $\{X_n\}$  is uniformly integrable
- $X_n \rightarrow X$  ( $L^p$ ) and either  $X_n \in L^p$  all  $n$  or  $X \in L^p$ .
- $E(|X_n|^p) \rightarrow E(|X|^p) < \infty$

## 7.3 Convergence of Distributions $\rightarrow^L$

Th4.3  $X_n \xrightarrow{pr} X$ , then  $X_n \rightarrow^L X$ . If  $X_n \in (\Omega, \mathcal{F}, P)$ , then  $X_n \xrightarrow{pr} a$  (constant)  $\iff X_n \rightarrow^L a$

Th4.4 Let  $P_n, P$  be probability measures with densities wrt a common measure  $\mu$ .  $f_n = dP_n/d\mu$  and  $f = dP/d\mu$ . If  $f_n \rightarrow f$  a.e.  $(\mu)$ , then  $\sup_{A \in \mathcal{F}} |P_n(A) - P(A)| \leq \int |f_n - f| d\mu \rightarrow 0$

Th4.5 (Slutsky)  $X_n \rightarrow^L X$  and  $Y_n \rightarrow^L C_1, Z_n \rightarrow^L C_2$ . Then  $X_n + Y_n \rightarrow^L X + C_1$  and  $X_n Z_n \rightarrow^L C_2 X$ . (Even  $Y_n \xrightarrow{pr} Y$  may  $X_n + Y_n \not\rightarrow^L X + Y$ )

Th4.9 (Skorohod's Theorem) Suppose  $F_n, F$  are probability distribution functions on  $R$  and  $F_n \rightarrow^L F$ . Then there exists a probability space  $(\Omega, \mathcal{F}, P)$  and r.v.'s  $Y_n, Y$  such that  $Y_n \sim F_n$  and  $Y \sim F$  and  $Y_n(\omega) \rightarrow Y(\omega)$  for every  $\omega \in \Omega$ . (because of this theorem, some topic about expectation involved into  $\rightarrow^L$  can be transferred into that into  $\rightarrow^{a.s.}$ )

Th4.10  $X_n \rightarrow^L X$

- If  $X_n \geq 0$  then  $\liminf_n E(X_n) \geq E(X)$
- If  $P(|X_n| > x) \leq P(|Y| > x)$  and  $E(|Y|) < \infty$ , then  $\lim_n E(X_n) = E(X)$ .
- $\{X_n\}$  is uniformly integrable iff  $E(X_n) \rightarrow E(X)$  as  $n \rightarrow \infty$ .

Th4.11 (Continuous Mapping Theorem) Let  $X_n \sim F_n$  and  $X \sim F$ . If  $X_n \rightarrow^L X$  and  $h: R \rightarrow R$  measurable and discontinuous only on  $D$  with  $P(x \in D) = 0$ . Then  $h(X_n) \rightarrow^L h(X)$ .

Th4.20 (Taylor expansion of the characteristic function). If  $E|X|^n < \infty$ , then  $E(X^k) = (-i)^k \phi^{(k)}(0)$  for  $k = 1, \dots, n$  and

$$\phi(t) = \sum_{k=0}^{n-1} \frac{\phi^{(k)}(0)}{k!} t^k + O(t^n) = \sum_{k=0}^n \frac{\phi^{(k)}(0)}{k!} t^k + o(t^n) \text{ (as } t \rightarrow \infty)$$

Le4.21  $(X, Y) \sim f(x, y)$ . Then

$$X + Y \sim \int_{R^1} f(x, z - x) dx, \quad X - Y \sim \int_{R^1} f(x + z, x) dx$$

$$X/Y \sim \int_{R^1} f(xz, x) |x| dx.$$

Th4.22 (Inverse formula) For any  $x_1 < x_2$ ,

$$\begin{aligned} & \frac{1}{2} \{F(x_2 + 0) + F(x_2)\} - \frac{1}{2} \{F(x_1 + 0) + F(x_1)\} \\ &= \lim_{T \rightarrow \infty} \frac{1}{2\pi} \int_{-T}^T \frac{e^{-itx_1} - e^{-itx_2}}{it} \phi(t) dt \end{aligned}$$

Th4.26 (Continuity Theorem) Let  $\{F_n\}$  be a sequence of d.f.'s with c.f.'s  $\phi_n$ . Then  $F_n \rightarrow^L F$  for some d.f.  $F \iff \phi_n(t) \rightarrow \phi(t)$  for all  $t$  and  $\phi$  is continuous at  $t = 0$ . In this case,  $\phi$  is the c.f. of  $F$  and the convergence  $\phi_n \rightarrow \phi$  is uniformly in every finite interval.

Th4.28 (Lindeberg-Lévy central limit theorem)  $\{X_i\}$  are i.i.d. with mean zero and finite variance. Then

$$\lim_{n \rightarrow \infty} \left\{ \frac{1}{\sqrt{n}\sigma} \sum_{i=1}^n X_i < x \right\} = N(0, 1)$$

De4.30 (Lindeberg Condition) A triangular array  $\{X_{nk}\}$  with mean zero satisfies the Lindeberg condition if

$$\lim_{n \rightarrow \infty} \frac{1}{B_n^2} \sum_{k=1}^{r_n} \int_{|X_{nk}| > \tau B_n} X_{nk}^2 dP = 0 \text{ for any } \tau > 0. B_n^2 = \sum_{i=1}^{r_n} EX_{ni}^2.$$

Th4.31 (Central limit theorem)  $\{X_i\}$  are independent with finite variance  $\sigma_i^2$ . Then

$$\lim_{n \rightarrow \infty} \left\{ \frac{1}{B_n} \sum_{i=1}^n (X_i - EX_i) < x \right\} = N(0, 1) \text{ and } \lim_{n \rightarrow \infty} \max_{k \leq n} \frac{\sigma_k}{B_n} = 0$$

if and only if Lindeberg condition holds. The second one is called Feller condition or uniformly asymptotically negligible.

Co4.32 (Lyapounov theorem) Suppose  $\{X_i\}$  are independent with satisfying, for some  $\delta > 0$ ,

$$\lim_{n \rightarrow \infty} \sum_{i=1}^n \frac{E|X_i|^{2+\delta}}{\{\sum_{i=1}^n \sigma_i^2\}^{1+\delta/2}} = 0$$

Then  $\{X_n\}$  satisfies a central limit theorem.

## 7.4 Absolute continuity/Conditional Expectation

De5.4 Let  $X$  be a r.v. on  $(\Omega, \mathcal{F}, P)$  such that its mean exists and let  $Q$  be a sub  $\sigma$ -field of  $\mathcal{F}$ . The conditional expectation of  $X$  given  $Q$  is any r.v. (denote  $E(X|Q)$ ) satisfying

- $E(X|Q)$  is measurable wrt  $Q$ , and
- $E\{1_B E(X|Q)\} = E(1_B X)$  for any  $B \in Q$ .

$$E\{g(X)|Y\} = \frac{\int g(x)f(x, Y) dx}{\int f(x, Y) dx}$$

Th5.9  $\mathcal{P}_1 \subset \mathcal{P}_2$ . Then  $E\{E(X|\mathcal{P}_1)|\mathcal{P}_2\} = E\{E(X|\mathcal{P}_2)|\mathcal{P}_1\} = E(X|\mathcal{P}_1)$  a.s. In particular,  $E\{E(X|Y)|Y, Z\} = E\{E(X|Y, Z)|Y\} = E(X|Y)$  a.s.

**An interesting counterexample-Sometimes,**  $E\{E(Y|X)Z\} \neq E(YZ)$ . For example,  $Y = g(X) + \epsilon$ . The error term  $\epsilon$  is independent of  $X$  with zero mean and finite variance  $\sigma^2$ . Let  $Z = \epsilon$ . Then  $E\{E(Y|X)Z\} = E\{g(X)Z\} = 0$ . But  $E(YZ) = E(\epsilon^2) = \sigma^2$

Th5.10  $\mathcal{P}_1$  and  $\mathcal{P}_2$  are all  $\sigma$ -fields and  $E\{E(X|\mathcal{P}_1)|\mathcal{P}_2\} = E(X|\mathcal{P}_2 \cap \mathcal{P}_1)$  a.s.

## 7.5 Related Materials

1. (**Feller-Chung Theorem**) For each integral number  $j$ ,  $A_j A_{j-1}^C \cdots A_0^C$  and  $B_j$  are independent, where  $A_0 = \emptyset$ . Then  $P(\cup_j A_j B_j) \geq \alpha P(\cup_j A_j)$  for  $\alpha = \inf_j P(B_j)$ .

**Proof.** Note that

$$\begin{aligned} \cup_j A_j B_j &= A_1 B_1 + A_2 B_2 (A_1 B_1)^C + \cdots \supset A_1 B_1 + A_2 B_2 A_1^C + \cdots \\ P(\cup_j A_j B_j) &\geq P(A_1 B_1) + P(A_2 A_1^C) P(B_2) + \cdots \geq \alpha \{P(A_1) + P(A_2 A_1^C) + \cdots\} \\ &= \alpha P(\cup_j A_j). \end{aligned}$$

2.  $X_1, \dots, X_n$  are independent and **symmetric** random variables. Write  $S_k = \sum_{i=1}^k X_i$  for  $k = 1, \dots, n$ . Then

$$P(\max_{1 \leq k \leq n} S_k > a) \leq 2P(S_n < a).$$

**Proof.** Set  $A_k = \{S_1 \leq a, \dots, S_{k-1} \leq a, S_k > a\}$  and  $B_k = \{S_n - S_k \geq 0\}$ . Using the results of the previous item, we can prove this conclusion. *This inequality is often used in Wiener processes.*

3.  $X$  and  $Y$  are independent with means zero. Then  $E|X+Y|^r \geq \max(E|X|^r, E|Y|^r)$  for any  $r \geq 1$ .

**Proof.** For any  $x$ ,  $|x|^r = |E(x+Y)|^r \leq E|x+Y|^r$ . It follows that

$$\begin{aligned} E|X+Y|^r &= \int |x+y|^r dF_X(x) dF_Y(y) \\ &= \int dF_X(x) \left\{ \int |x+y|^r dF_Y(y) \right\} \\ &= \int E|x+Y|^r dF_X(x) \geq \int |x|^r dF_X(x) = E|X|^r. \end{aligned}$$

4. (**Generalized Kolmogorov inequality**)  $X_1, \dots, X_n$  are independent random variables with mean zeros. Write  $S_k = \sum_{i=1}^k X_i$  for  $k = 1, \dots, n$  and  $A = \{\sup_{k \leq n} |S_k| \geq C\}$  for some positive constant  $C$ . Then

$$C^r P(A) \leq E(|S_n|^r I_A) \leq E|S_n|^r \quad \text{for } r \geq 1.$$

**Proof.** Set  $A_0 = \emptyset$  and  $A_k = \{\sup_{j < k} |S_j| < C, |S_k| \geq C\}$ . Then  $A_1, \dots, A_n$  are disconnect each other and  $A = \sum A_k$ . It follows from the result of the former item that

$$E|S_n|^r I_A = \sum_{k=1}^n E|S_n|^r I_{A_k} \geq \sum_{k=1}^n E|S_k|^r I_{A_k} \geq \sum_{k=1}^n C^r E I_{A_k} = C^r P(A),$$

and we complete the proof of the first assertion. The second one is trivial.

5. If random variable  $X$  is integral, then  $|\text{median}(X) - E(X)| \leq \{2\text{var}(X)\}^{1/2}$ .
6.  $E|X|^p \leq \infty$  ( $p \geq 1$ ) if and only if  $\sum_{n=1}^{\infty} \int_{|x| \geq n} |x|^{p-1} dF(x) < \infty$ .
7. (**Borel law of large number**)  $\mu_n \sim \text{Bernoulli}(n, p)$ , then

$$P\left\{\lim_{n \rightarrow \infty} \frac{\mu_n}{n} = p\right\} = 1 \iff P\left\{\left|\frac{\mu_n}{n} - p\right| > \varepsilon\right\} \leq \frac{E|\mu_n - np|^4}{n^4 \varepsilon^4}$$

8. (**Hajek-Renyi inequality**)  $X_i$  are independent each other with finite variances.

$$P\left\{\max_{m \leq j \leq n} \left|C_j \sum_{i=m}^j (X_i - EX_i)\right| > \varepsilon\right\} \leq \frac{1}{\varepsilon^2} \left(C_2^m \sum_{i=1}^m \sigma_i^2 + \sum_{i=1+m}^n C_i^2 \sigma_i^2\right)$$

9. (**Kolmogorov strong law of large number**)  $X_i$  are independent each other with finite variances. If  $\sum_{n=1}^{\infty} D(X_n)/b_n^2 < \infty$  for  $b_n \uparrow \infty$ , then  $\sum_{i=1}^n (X_i - EX_i)/b_n \rightarrow 0$  a.s.

10. (**Kolmogorov strong law of large number**)  $X_i$  are i.i.d.  $\sum_{i=1}^n X_i/n$  a.s. converges to  $a$  if and only if  $EX_i < \infty$  and  $a = E(X_1)$ .

11. Chebyshev inequality, Markov inequality.

12. (**Khinchine law of large number**)  $X_i$  i.i.d. with finite mean. Then  $\{X_n\}$  is satisfied with law of large numbers.

13.  $X_i \sim F_i(x)$ . If  $\lim_{A \rightarrow \infty} \sup_{1 \leq n < \infty} \int_{|x| > A} |x| dF_n(x) = 0$ , Then  $\{X_n\}$  is satisfied with law of large numbers. (Use **Kolmogorov three series theorem** to prove)

14.  $X_i$  are independent. If there exists  $\alpha > 1$  and  $\beta > 0$  such that  $E|X_n|^\alpha \leq \beta$ , Then  $\{X_n\}$  is satisfied with law of large numbers.

15. (**Markov law of large number**)  $X_i$  are independent with mean zero. There exists a  $0 < \delta \leq 1$  such that

$$\frac{1}{n^{1+\delta}} \sum_{i=1}^n E|X_i|^{1+\delta} \rightarrow 0$$

Then  $\{X_n\}$  is satisfied with law of large numbers.

16. (**Elementary Inequality**)  $f(x)$  is a **non-decreasing continuous** function. Define

$$a.e \sup f(\xi) = \inf\{c : P(f(\xi) > c) = 0\}.$$

Then we have the following conclusion:

$$\frac{Ef(|\xi|) - f(\varepsilon)}{a.e \sup f(\xi)} \leq P(|\xi| > \varepsilon) \leq \frac{Ef(|\xi|)}{f(\varepsilon)}$$

In this case, assume  $f(0) = 0$ , then  $\xi_n \rightarrow 0$  in probability if and only if  $Ef(|\xi_n|) \rightarrow 0$ .

17. (**Gnedenko law of large number**) Taking  $f(x) = x^2/(1+x^2)$  in **elementary inequality**, we know that  $\{X_n\}$  is satisfied with law of large numbers if and only if

$$E \left[ \frac{\{\sum_{i=1}^n (X_i - EX_i)\}^2}{n^2 + \{\sum_{i=1}^n (X_i - EX_i)\}^2} \right] \rightarrow 0 \iff \sum_{i=1}^n E \left\{ \frac{(X_i - EX_i)^2}{1 + (X_i - EX_i)^2} \right\} \rightarrow 0 \text{ if independent}$$

18.  $\{X_i\}$  are independent. There exist constants  $k_n$  such that  $\max_{1 \leq j \leq n} |X_j| \leq k_n$  and  $k_n/B_n \rightarrow 0$ . Then  $\{X_i\}$  obey central limit theorem.

19.  $\{X_i\}$  are independent and obey central limit theorem.  $\{X_n\}$  are satisfied with law of large numbers if and only if  $B_n^2 = o(n^2)$ .

20.  $\{X_i\}$  are independent.  $X_1 \sim U[-1, 1]$  and  $X_k \sim N(0, 4^{k-1})$  for  $k = 2, \dots$ . Then  $\{X_k\}$  satisfy central limit theorem (using c.f. to prove), but not Lindeberg condition and Feller condition because  $b_n/B_n^2 \rightarrow 1/2$ .

21. (**The 1st Helly Theorem**) Suppose  $f(x)$  is a continuous function on  $[a, b]$ ,  $F_n(x)$  uniformly bound nondecreasing and converge to  $F(x)$  on  $[a, b]$ .  $a, b$  are the continuous points of  $F(x)$ , then

$$\lim_{n \rightarrow \infty} \int_a^b f(x) dF_n(x) = \int_a^b f(x) dF(x).$$

22. (**The 2nd Helly Theorem**) Suppose  $f(x)$  is a continuous bounded function on  $R^1$ ,  $F_n(x)$  uniformly bound nondecreasing and converge to  $F(x)$  on  $R^1$ . In addition  $F_n(-\infty) \rightarrow F(-\infty)$  and  $F_n(\infty) \rightarrow F(\infty)$ . Then

$$\lim_{n \rightarrow \infty} \int_{-\infty}^{\infty} f(x) dF_n(x) = \int_{-\infty}^{\infty} f(x) dF(x).$$

- 147 Suppose  $X_n$  are independent, then  $X_n \rightarrow^{a.s.} 0 \iff \forall \varepsilon, \sum_{n=1}^{\infty} P\{|X_n| \geq \varepsilon\} < \infty$ . The proof can be completed using **Borel-Cantelli Theorem** by setting  $A_n = \{|X_n| \geq \varepsilon\}$ .

- 179 Suppose  $\{X_n, n \geq 1\}$  is iid. Then there exists a sequence of constants  $C_n$  such that

$$\frac{1}{n} \left( \sum_{j=1}^n X_j - C_n \right) \rightarrow^{pr} 0 \iff \lim_{n \rightarrow \infty} nP(|X_1| \geq n) = 0$$

and  $C_n = nP(|X_1| \geq n)$ . If  $E|X_1| < \infty$ ,  $C_n$  can be taken  $nE(X_1)$ .

- 177 Let  $p > 0$  and  $F$  be the cdf of r.v.  $X$ . (i) If  $E|X|^p < \infty$ , then for any  $\alpha > -1, \beta > 0$  and  $\gamma \geq 0$  satisfying  $\frac{\alpha+1}{\beta} + \gamma = p$ , we have

$$\sum_{n=1}^{\infty} n^\alpha \int_{|x| > n^\beta} |x|^\gamma dF(x) < \infty. \quad (2)$$

Conversely, if there is a set of  $(\alpha, \beta, \gamma)$  satisfying (2), then  $E|X|^p < \infty$ . (ii) If “ $\alpha > -1$ ” was changed into “ $\alpha < -1$ ”, the integral in (2) should also changed as  $\int_{|x| \leq n^\beta} |x|^\gamma dF(x)$ . (This is a very useful result. For example, prove 183. If  $p < 1$ , taking  $\alpha = 0, \beta = 1/p$  and  $\gamma = 0$  we have  $E|X|^p < \infty \iff \sum_{n=1}^{\infty} P(|X| \geq n^{1/p}) < \infty$ ; If  $p > 1$ , taking  $\alpha = 0, \beta = 1$  and  $\gamma = p-1$  we have  $E|X|^p < \infty \iff \sum_{n=1}^{\infty} E\{|X|^{p-1} I_{|X| \geq n}\} < \infty$ )

- 183 (**Marcinkiewicz-Zygmund Theorem**) Suppose  $\{X_n, n \geq 1\}$  is iid and  $p \in (0, 2)$ . Then there exists a sequence of constants  $C_n$  such that

$$n^{-1/p} \left( \sum_{j=1}^n X_j - C_n \right) \rightarrow^{a.s.} 0 \iff E|X_1|^{1/p} < \infty,$$

and  $C_n = 0$  if  $0 < p < 1$  and  $nE(X_1)$  otherwise.

- 211 Suppose  $\{X_n, n \geq 1\}$  is iid. Then

$$\max_{1 \leq i \leq n} X_i \rightarrow^{pr} 0 \iff nP(|X_1| > n) = o(1) \text{ and } \frac{1}{n} \max_{1 \leq i \leq n} X_i \rightarrow^{a.s.} 0 \iff E|X| < \infty$$

- (**Doob Inequality**) For independent sequence  $\{X_n\}$  with mean zero and  $p > 1$ ,

$$E\left(\max_{1 \leq k \leq n} \left| \sum_{j=1}^k X_j \right|^p\right) \leq \left(\frac{p}{p-1}\right)^p E\left(\left| \sum_{j=1}^n X_j \right|^p\right)$$

- 65, 81, 138 in the big notebook

- Characteristic Functions

Dis.	density	c.f.	additivity
$B(n, p)$	$\binom{n}{k} p^k (1-p)^{n-k}$	$(pe^{it} + 1-p)^n$	$B(n_1, p) * B(n_2, p) = B(n_1 + n_2, p)$
$P(\lambda)$	$\frac{\lambda^k}{k!} \exp(-\lambda)$	$\exp(\lambda(e^{it} - 1))$	$P(\lambda_1) * P(\lambda_2) = P(\lambda_1 + \lambda_2)$
$N(\mu, \sigma^2)$	$(\sqrt{2\pi}\sigma)^{-1} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$	$\exp(i\mu t - \frac{\sigma^2 t^2}{2})$	$N(\mu_1, \sigma_1^2) * N(\mu_2, \sigma_2^2) = N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$
$\Gamma(\alpha, \lambda)$	$\frac{\lambda^\alpha}{\Gamma(\lambda)} x^{\lambda-1} e^{-\alpha x}$	$\left(1 - \frac{it}{\lambda}\right)^{-\lambda}$	$\Gamma(\alpha, \lambda_1) * \Gamma(\alpha, \lambda_2) = \Gamma(\alpha, \lambda_1 + \lambda_2)$
$C(\alpha, \mu)$	$\frac{\alpha}{\pi} \frac{\pi(\alpha^2 + (x-\mu)^2)^{-1}}{2\mu}$	$e^{i\mu t - \alpha t }$	$C(\alpha_1, \mu_1) * C(\alpha_2, \mu_2) = C(\alpha_1 + \alpha_2, \mu_1 + \mu_2)$
$\chi^2(n)$	$\frac{1}{2^{n/2} \Gamma(n/2)} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}$	$(1-2it)^{-\frac{n}{2}}$	$\chi^2(n_1) * \chi^2(n_2) = \chi^2(n_1 + n_2)$

## 8 Basics of Statistical Inference

**Definition 1** Suppose  $X$  is an observation from an unknown distribution  $P \in \mathcal{P}$  where  $\mathcal{P}$  is a family of distributions. A statistic  $T = T(X)$  is said to be sufficient for  $P \in \mathcal{P}$  (or  $\theta \in \Theta$  where  $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$  is a parametric family) if the conditional distribution of  $X$  given  $T$  does not depend on  $P$  (or  $\theta$  in parametric family).

**Theorem 1 (N-F factorization Theorem).** Suppose  $X$  is an observation from  $P \in \mathcal{P} = \{f(\cdot, \theta), \theta \in \Theta\}$ . A statistic  $T = T(X)$  is sufficient for  $\mathcal{P}$  iff  $f(x, \theta) = g(t(x), \theta)h(x)$  where  $g$  and  $h$  are two functions and  $h$  is free of  $\theta$ .

**Theorem 2** Consider the model  $\mathcal{P} = \{f(\cdot, \theta); \theta \in \Theta\}$ . Let  $T$  be a minimal sufficient statistic ( $T$  is a function of any sufficient statistic). Then for any  $x, y$  in the sample space.

$$T(x) = T(y) \iff L(\theta, x) \propto L(\theta, y)$$

- Remark 1**
- Theorem 2 shows that the likelihood function is equivalent to the following: Sufficiency principle inference for a given model which admits a minimal sufficient statistic  $T$  should be identical for any  $x$  and  $y$  such that  $T(x) = T(y)$ .
  - Theorem 2 also provides methods for identifying minimal sufficient statistic. If  $T$  is a statistic such that for any  $x, y$

$$T(x) = T(y) \iff L(\theta, x) \propto L(\theta, y).$$

Then  $T$  is minimal sufficient.

**Definition 2** The family  $\{f(\cdot, \theta); \theta \in \Theta\}$  constitutes an exponential family if

$$f(x, \theta) = h(x) \exp \left\{ \sum_{i=1}^k \phi_i(\theta) T_i(x) - \tau(\theta) \right\}$$

where  $T_1, \dots, T_k$  and  $h$  are functions of  $x$  not depending on  $\theta$ , and  $\phi_1, \dots, \phi_k$  and  $\tau$  are functions of  $\theta$  not depending on  $x$ .

**Theorem 3** If the exponential family is in reduced form, then  $(T_1(x), \dots, T_k(x))$  is **minimal sufficient statistic**.

**Theorem 4** If  $k = 1$  and  $f(x, \theta) = h(x) \exp\{\phi(\theta)T(x) - \tau(\theta)\}$ , then the moments of  $T$  can be represented in terms of  $\phi$  and  $\tau$

$$E_\theta\{T(X)\} = \frac{\tau'(\theta)}{\phi'(\theta)} \quad \text{and} \quad \text{Var}_\theta\{T(X)\} = \frac{\phi'(\theta)\tau''(\theta) - \tau'(\theta)\phi''(\theta)}{\phi'^3(\theta)}.$$

and so on.

**Definition 3** A statistic  $T(x)$  is an ancillary statistic if its distribution does not depend on the parameter.

**Definition 4** Suppose  $\mathbf{x}$  is an observation from an unknown distribution  $\mathcal{P} = \{P_\theta; \theta \in \Theta\}$ . A statistic  $T = T(X)$  is said to be complete if for any measurable function  $g$ ,  $Eg(T) = 0$  for all  $P \in \mathcal{P}$  means that  $g(t) = 0$  a.e.  $P$

We say that  $T = T(X)$  is boundedly complete if the previous statement holds for and bounded measurable function  $g$ .

**Example 1** (a) Suppose  $X \sim B(n, p)$ . Suppose  $g$  is such that  $\sum_{x=0}^n g(x) \binom{n}{x} p^x (1-p)^{n-x} = 0$  for all  $p$ . This means  $g(x) = 0$  for  $x = 0, 1, \dots, n$ , and  $X$  is complete.

(b)  $X_1, \dots, X_n$  are iid  $U(0, \theta)$  for  $\theta \in (0, \infty)$ . Then  $X_{(n)}$  is sufficient. Suppose  $g$  is such that  $0 = Eg(X_{(n)}) \propto \int g(t)t^{n-1}dt$  for all  $\theta$ . This means that  $g(t)t^{n-1} = 0$  a.e. from measure theory. Thus  $X_{(n)}$  is complete.

(c)  $X_1, \dots, X_m$  are iid with pdf  $N(\mu, \sigma_1^2)$  and  $Y_1, \dots, Y_n$  are iid with pdf  $N(\mu, \sigma_2^2)$  and they are independent.  
 $T = (\sum_{i=1}^m X_i^2, \sum_{i=1}^n Y_i^2, \sum_{i=1}^m Y_i, \sum_{i=1}^n X_i)$  is sufficient. But  $T$  is not complete since  $E \sum X_i - E \sum Y_i = 0$  for all  $\mu, \sigma_1^2, \sigma_2^2$ .

(d)  $X_1, \dots, X_n$  are iid  $U(\theta, \theta + 1)$  for  $\theta \in R^1$ . Then  $T = (X_{(1)}, X_{(n)})$  is minimal sufficient but not complete (HW).

(e) Suppose  $X$  is observation from

$$P \in \left\{ f(x, \phi_1, \dots, \phi_k) = h(x) \exp \left\{ \sum_{i=1}^k \phi_i T_i(x) - \xi(\phi) \right\}, \phi \in \Theta \right\}$$

where  $\Theta$  contains an open set (the family is said to be full-rank). Then  $(T_1, \dots, T_k)$  is **complete**.

**Theorem 5** Let  $T$  be a one-dimensional complete and sufficient statistic. Then it is minimal sufficient.

**Lemma 1** Let  $X, Y$  be r.v.'s where  $Y$  has finite variance, then (i)  $E(E(Y|X)) = E(Y)$  and (ii)  $\text{Var}\{E(Y|X)\} \leq \text{Var}(Y)$ .

**Theorem 6 (Basu Theorem)** Let  $V$  and  $T$  be two statistics based on an observation  $X$  from  $P_\theta \in \mathcal{P}$ . If  $T$  is boundedly complete and sufficient and the distribution of  $V$  doesn't depend on  $\theta$ . Then  $V$  and  $T$  are independent for any  $\theta$ .

### 8.1 Point Estimation

**Theorem 7 (Blackwell-Lehmann-Rao-Scheffé Theorem)(B-L-R-S)** Let  $X$  be an observation from a distribution in a family  $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$ . Assume that  $g(\theta)$  is  $U$ -estimable and  $U$  is an unbiased estimator of  $g(\theta)$ .

- If  $T$  is sufficient for  $g(\theta)$ ,  $E(U|T)$  is also unbiased for  $g(\theta)$  and  $\text{Var}_\theta[E(U|T)] \leq \text{Var}_\theta(U)$  for all  $\theta$ .
- If  $T$  is complete and sufficient for  $\theta$ , Then  $E(U|T)$  is the unique UMVUE of  $g(\theta)$ . (Here unique means if there exists another estimator  $V$  which is UMVUE, then  $V(x) = U(X)$  a.e.  $P_\theta$ )
- Therefore if a unbiased estimator is not one function of a complete and sufficient statistic, the estimator must not be UMVUE. For example the sample variance  $S_n^2$  for  $\sigma^2$  in  $N(0, \sigma^2)$ .

The following are typical approaches for deriving UMVUE when a complete and sufficient statistic  $T$  is available.

- We happen to know that  $\phi(T)$  is unbiased for  $g(\theta)$ , then  $\phi(T)$  is UMVUE of  $g(\theta)$ .
- We first identify an unbiased estimator  $U$  of  $g(\theta)$ , and then calculate  $E(U|T)$ , which is UMVUE.
- In some case, one can solve  $E_\theta \phi(T) = g(\theta)$  for  $\phi$ .

**Theorem 8** Let  $\mathcal{U}$  be the set of all unbiased estimators of 0 with finite variance and  $T$  an unbiased estimator of  $g(\theta)$ . A necessary and sufficient condition for  $T$  to be UMVUE is that  $\text{Cov}(U, T) = 0$  for all  $U \in \mathcal{U}$ .

**Theorem 9 (The Cramér-Rao Lower Bound).** Let  $X$  be an observation from  $P \in \mathcal{P} = \{P_\theta, \theta \in \Theta\}$ , where  $\Theta$  is an open set in  $R^k$ . Suppose that  $T = T(X)$  is an unbiased estimator of  $g(\theta)$ , where  $g$  is differentiable at all  $\theta \in \Theta$ . Further, suppose that  $P_\theta$  has a density function  $f(x, \theta)$  w.r.t some measure  $\nu$  for all  $\theta \in \Theta$ , and  $f(x, \theta)$  is differentiable in  $\theta$  and satisfies that

$$\frac{\partial}{\partial \theta} \int h(x) f(x, \theta) d\nu = \int h(x) \frac{\partial}{\partial \theta} f(x, \theta) d\nu \quad (3)$$

for all  $\theta \in \Theta$ , and  $h(x) = 1$  and  $h(x) = T(x)$ . Then

$$\text{Var}_\theta\{T(X)\} \geq \left( \frac{\partial}{\partial \theta} g(\theta) \right) I^{-1}(\theta) \left( \frac{\partial}{\partial \theta} g(\theta) \right)^T$$

where  $I(\theta) = E_\theta \left\{ \left( \frac{\partial}{\partial \theta} \log f(X, \theta) \right) \left( \frac{\partial}{\partial \theta} \log f(X, \theta) \right)^T \right\}$

The r.v.  $\frac{\partial}{\partial \theta} f(x, \theta)$  is called the efficient score of  $\theta$ .  $I(\theta)$  is called Fisher Information Matrix.

**Series Expansion Method** Often we wish to estimate  $g(\theta)$  when we have an unbiased estimator  $T$  of  $\theta$ . We are attempted to use  $g(T)$  as the estimator of  $g(\theta)$ , but this is typically biased. We can express  $g(T)$  about  $\theta$  using Taylor series

$$g(T) \approx g(\theta) + g'(\theta)(T - \theta) + \frac{1}{2} g''(\theta)(T - \theta)^2$$

Taking expectation both side, we get

$$Eg(T) \approx g(\theta) + \frac{1}{2} g''(\theta) \text{Var}(T).$$

Often  $\text{Var}(T) = O(1/n)$ . This means that the bias has order  $1/n$ . In some cases, we can estimate  $g''(\theta) \text{Var}(T)$  and modify  $g(T)$  accordingly. So that it will have smaller bias.

**Jackknife.**  $\bar{T}_{n-1, \cdot} = \frac{1}{n} \sum_{j=1}^n T_{n-1, j}$ . Define  $T_n^J = nT_n - (n-1)\bar{T}_{n-1, \cdot}$ .

### 8.2 Maximum Likelihood Estimation (MLE)

**Definition 5** Suppose  $\mathbf{X}$  is a sample from  $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$  where  $P_\theta$  is assumed to have a density  $f(x, \theta)$ . Let  $L(\theta, x)$  be the likelihood function. A statistic  $\hat{\theta} \in \Theta$  satisfying

$$L(\hat{\theta}, \mathbf{X}) = \max_{\theta \in \Theta} L(\theta, \mathbf{X})$$

is called maximum likelihood estimate of  $\theta$ .  $\hat{\theta}$  viewed as an estimator is called maximum likelihood estimator.

**Newton-Raphson** Let  $\theta_0$  be a fixed point. Write

$$Dl(\theta, x) = \left[ \frac{\partial}{\partial \theta_1} l(\theta, x), \dots, \frac{\partial}{\partial \theta_k} l(\theta, x) \right]^T \quad \text{and} \\ D^2 l(\theta, x) = \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} l(\theta, x) \right]_{i, j=1 \sim k}. \quad \hat{\theta} \approx \theta_0 - [D^2 l(\theta_0, x)]^{-1} Dl(\theta_0, x).$$

**EM-algorithm**

**E-step** (Estimation step). Compute  $Q(\theta|\theta_k) = E_{\theta_k}[\log L(\theta, X)|Y]$ ;

**M-step** (Maximization step). Select  $\theta_{k+1}$  as the maximization of  $Q(\theta|\theta_k)$ . Apply these steps iteratively until "convergence".



**Theorem 10** Let  $X_1, \dots, X_n$  be iid with a common density  $f(x, \theta)$  w.r.t a  $\sigma$ -finite measure (focus on pdf and pmf) where  $\theta$  is real-valued. Assume the following conditions.

- (a) The parameter space  $\theta$  is an open interval (finite or infinite)
- (b) The distribution  $P_\theta$  of  $X_i$  have common support so that  $A = \{x, f(x, \theta) > 0\}$  is independent of  $\theta$ .
- (c) For any  $x \in A$ , the density  $f(x, \theta)$  is three times differentiable in  $\theta$ ,
- (d)

$$E_\theta \left[ \frac{\partial}{\partial \theta} \log f(X, \theta) \right] = 0.$$

and

$$E_\theta \left[ \frac{\partial^2}{\partial \theta_i \partial \theta_j} \log f(X, \theta) \right] = -I(\theta).$$

- (e) There exists a finite neighbor  $c(\theta_0 - \varepsilon, \theta_0 + \varepsilon)$  and a function  $M(x)$  such that

$$\left| \frac{\partial^3}{\partial \theta^3} \log f(x, \theta) \right| \leq M(x)$$

for all  $x \in A$  and  $\theta \in c(\theta_0 - \varepsilon, \theta_0 + \varepsilon)$  with  $E_{\theta_0} M(X) < \infty$ .

Then any consistent sequence  $\hat{\theta}_n$  of roots of the likelihood equation satisfies

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \longrightarrow N(0, I^{-1}(\theta_0)).$$

**Definition 6** Let  $X_1, \dots, X_n$  be iid with a common density  $f(x, \theta)$ . Let  $I(\theta)$  be the Fisher information of  $X$ , which is assumed to be well-defined and finite. Let  $T_n$  be an estimator of  $g(\theta)$  where  $g(\theta)$  is differentiable with  $g'(\theta) > 0$ . We say  $T_n$  is **asymptotically efficient** if

$$\sqrt{n}(T_n - g(\theta)) \longrightarrow N(0, v(\theta))$$

where  $v(\theta) = [g'(\theta)]^2 / I(\theta)$ .

**Theorem 11** Suppose that the conditions of Theorem 10 hold and that  $\tilde{\theta}_n$  is any root-n consistent estimators, i.e.  $\sqrt{n}(\tilde{\theta}_n - \theta_0) = O_p(1)$ . Then the estimator sequence

$$T_n = \tilde{\theta}_n - \frac{l'(\tilde{\theta}_n, x)}{l''(\tilde{\theta}_n, x)}$$

is asymptotically efficient. (In fact  $n^{1/4}(\tilde{\theta}_n - \theta_0) = o_p(1)$  is enough)

## 8.3 Robustness

**Stielties Integral**  $F(x) = pF_c(x) + qF_d(x)$ , where  $F_c(x)$  has a derivative  $f_c(x)$  and  $F_d(x)$  is a step function with discontinuous points  $x_1, \dots, x_n$ , the size of jump at the  $x_i$  is  $p_i$ . Then

$$\int g(x) dF(x) = p \int g(x) dF_c(x) + q \sum_{i=1}^n g(x_i) p_i.$$

**Definition 7** The influence function (or curse) of  $T$  at  $F$  is defined for each  $x$  by

$$IF(x, T, F) = \lim_{\varepsilon \downarrow 0} \frac{T((1 - \varepsilon)F + \varepsilon \delta_x) - T(F)}{\varepsilon}$$

when the limit exists.

It is usually true that  $\sqrt{n}\{T(F_n) - T(F)\} \longrightarrow N\{0, \int IF^2(x, T, F) dF(x)\}$ .

## 8.4 Decision Theoretic Estimation

- Risk function:  
 $R_T(\theta) = \int L(T(x), g(\theta)) p_\theta(x) d\mu(x) = E_\theta[L(T(x), g(\theta))]$ . An estimator  $T$  is said to be admissible if there is no other estimator  $T^*$  such that  $R_{T^*}(\theta) \leq R_T(\theta)$  for all  $\theta \in \Theta$  and  $R_{T^*}(\theta) < R_T(\theta)$  for some  $\theta$ .
- Bayes Risk:  $\bar{R}_T = \int R_T(\theta) \pi(\theta) d\theta$ . An estimator  $T_\pi$  is said to be **Bayes** estimator with respect to  $\pi$  if  $\bar{R}_{T_\pi} \leq \bar{R}_T$  for all estimator  $T$ .  
Special Cases:

- $L(T(x), g(\theta)) = \{T(x) - g(\theta)\}^2$ .  
 $\rho_T(x) = \int \{T(x) - g(\theta)\}^2 \pi(\theta|x) d\theta$ . Thus  $\rho_T(x)$  is minimized by  $T(x) = E\{g(\theta)|X\}$ .
- $L(T(x), g(\theta)) = |T(x) - g(\theta)|$ .  
 $\rho_T(x) = \int |T(x) - g(\theta)| \pi(\theta|x) d\theta$ . Thus  $\rho_T(x)$  is minimized by  $T(x)$  = the **median** of  $g(\theta)$  given  $X$ , which is just Bayes estimator.
- $L(T(x), g(\theta)) = w(\theta)\{T(x) - g(\theta)\}^2$ .  
 $\rho_T(x) = \int w(\theta)\{T(x) - g(\theta)\}^2 \pi(\theta|x) d\theta$ . Thus  $\rho_T(x)$  is minimized by

$$T(x) = \frac{E\{w(\theta)g(\theta)|X\}}{E\{w(\theta)|X\}}.$$

- An estimator  $T_\pi$  is said to be **Minimax** estimator if  $\sup_{\theta \in \Theta} R_{T_\pi}(\theta) \leq \sup_{\theta \in \Theta} R_{T^*}$  for any other estimator  $T^*$  of  $\theta$ .  
Suppose  $T_\pi$  is Bayes with respect to  $\pi$  and  $T_\pi$  has constant risk. Then  $T_\pi$  is minimax.

Remark: The MLE or UMVUE may be inadmissible!

**Proposition 8.1** If  $T_\pi$  is unique Bayes with respect to prior  $\pi$ , then  $T_\pi$  is admissible.

## 8.5 Hypothesis Test

- parameter space  $\Theta$
- Hypothesis  $H_0 : \theta \in \Theta_0 \subset \Theta \iff H_1 : \theta \in \Theta_0^c$
- Reject  $H_0$  if  $\delta(x) = 1$  and accept (do not reject)  $H_0$  if  $\delta(x) = 0$ .
- $\gamma(\theta) = P_\theta(\delta(x) = 1)$  is called **power function** of  $\delta$  at  $\theta$ .
- $\alpha(\theta) = \gamma(\theta)$  = probability of type I error for  $\theta \in \Theta_0$ ;
- $\beta(\theta) = 1 - \gamma(\theta)$  = probability of type II error for  $\theta \in \Theta_0^c$ .

1. (UMT) A test  $\phi$  of size  $\alpha$  is a uniformly most powerful test if  $\gamma_\phi(\theta) \geq \gamma_{\bar{\phi}}(\theta)$  for all  $\theta \in \Theta_0^c$  and size of  $\alpha$ .

2. (**N-P Lemma**) Any type of the form  $\phi(x) = \begin{cases} 1 & f_1/f_0 > c \\ \xi(x) & f_1/f_0 = c \\ 0 & o.w \end{cases}$

for some  $c \geq 0$  and  $0 \leq \xi(x) \leq 1$  satisfying  $E_{\theta_0}\{\phi(X)\} = \alpha$  is MP-level  $-\alpha$  test.

3. Let  $\{f_\theta; \theta \in \Theta\}$  be a family with **MLR** in  $T(x)$ . (for all  $\theta < \theta'$ ,  $f_{\theta'}(x)/f_\theta(x)$  is a non-decreasing function of  $T(x)$ )  
(i) For testing  $H_0 : \theta \leq \theta_0$  vs  $H_1 : \theta > \theta_0$ , there exists a UMP test of level  $\alpha$  given by

$$\phi(x) = \begin{cases} 1 & T(x) > c \\ \xi & T(x) = c \\ 0 & o.w \end{cases} \quad (4)$$

where  $c$  and  $\xi$  are determined by

$$E_{\theta_0} \phi(X) = \alpha \quad (5)$$

- (ii) The power function  $\gamma(\theta) = E_\theta \phi(X)$  of the test (4) is strictly increasing for all  $\theta$ .
- (iii) For all  $\theta'$ , the test determined by (4) and (5) is UMP for testing  $H_0 : \theta \leq \theta'$  vs  $H_1 : \theta > \theta'$  at level  $\alpha' = \gamma(\theta')$ .

4. Suppose  $X = (X_1, \dots, X_n)$  is a random variable from the **one-dimensional exponential** family

$$f_\theta(x) = h(x) \exp\{T(x)\theta - \tau(\theta)\}$$

then the **UMPU** test for  $H_0 : \theta = \theta_0 \iff H_1 : \theta \neq \theta_0$  is given by

$$\phi(x) = \begin{cases} 1 & T(x) < C_1 \text{ or } T(x) > C_2 \\ \xi_i & T(x) = C_i \text{ for } i = 1, 2 \\ 0 & C_1 < T(x) < C_2 \end{cases} \quad (6)$$

where  $C_i$  and  $\xi_i$  are determined by the following two equations

$$E_{\theta_0}\{\phi(X)\} = \alpha \quad E_{\theta_0}\{T(X)\phi(X)\} = \alpha E_{\theta_0}\{T(X)\}.$$

5. Suppose  $X$  has pdf  $f_{\theta, \eta}(x) = c(\theta, \eta)h(x) \exp\{\theta u(x) + \sum_{i=1}^k \eta_i T_i(x)\}$  where  $(\theta, \eta) \in R^{k+1}$ . Define

$$\phi(u, t) = \begin{cases} 1 & u > c(t) \\ \xi(t) & u = c(t) \\ 0 & o.w \end{cases}$$

where  $c(t)$  and  $\xi(t)$  are determined by

$$E_{\theta_0}[\phi\{U(X), T(X)\} | T(X) = t] = \alpha$$

for all  $t$  with  $T = (T_1, \dots, T_k)$ . Then  $\phi$  is a **UMPU** level  $1 - \alpha$  test for  $H_0 : \theta \leq \theta_0 \iff H_1 : \theta > \theta_0$ .

6. Assume the setting as in item 5 and consider testing  $H_0 : \theta = \theta_0 \iff H_1 : \theta \neq \theta_0$ . Then the **UMPU** test of level  $\alpha$  is given by

$$\phi(u, t) = \begin{cases} 1 & u > c_1(t) \text{ or } u < c_2(t) \\ \xi_i(t) & u = c_i(t) \\ 0 & o.w \end{cases}$$

where  $c_i(t)$  and  $\xi_i(t)$  are determined by  $E_{\theta_0}\{\phi\{U, T\} | T = t\} = \alpha$  and  $E_{\theta_0}[U\phi\{U, T\} | T = t] = \alpha E_{\theta_0}\{U | T = t\}$  for all  $t$ .

7. Suppose  $X$  has pdf  $f_{\theta, \eta}(x) = c(\theta, \eta)h(x) \exp\{\theta u(x) + \sum_{i=1}^k \eta_i T_i(x)\}$  and that  $V = V(u, T)$  is independent of  $T$  when  $\theta = \theta_0$ .

- (a) Assume further that  $V(u, t)$  is increasing in  $u$  for each fixed  $t$ .  
Then the **UMPU** test of  $H_0 : \theta \leq \theta_0 \iff H_1 : \theta > \theta_0$  is given

$$\varphi(u, t) = \begin{cases} 1 & V > C \\ \xi & V = C \\ 0 & o.w \end{cases}$$

where  $C$  and  $\xi$  are determined by  $E_{\theta_0} \{\varphi(X)\} = \alpha$ .

- (b) Assume further that  $V(u, t) = a(t)u + b(t)$  where  $a(t) > 0$  for all  $t$ . Then the **UMPU** test of  $H_0 : \theta = \theta_0 \iff H_1 : \theta \neq \theta_0$  is given

$$\varphi(V) = \begin{cases} 1 & V < C_1 \text{ or } V > C_2 \\ \xi_i & V = C_i \\ 0 & o.w \end{cases}$$

where  $C_i$  and  $\xi_i$  are determined by  $E_{\theta_0} \{\varphi(V)\} = \alpha$  and  $E_{\theta_0} \{\varphi(V)V\} = \alpha E_{\theta_0}(V)$ .

## 8. Bse( $\hat{\theta}$ ): Bootstrap Se( $\hat{\theta}$ )

$$Bse(\hat{\theta}) = \left[ \frac{1}{B} \sum_{b=1}^B \{\hat{\theta}_b - be(\hat{\theta})\}^2 \right]^{1/2} \quad \text{where } be(\hat{\theta}) = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b$$

8. **(LRT)**  $\lambda(x) = \frac{\sup_{\theta \in \Theta} f_{\hat{\theta}}(x)}{\sup_{\theta \in \Theta_0} f_{\theta}(x)}$ .
9. **Wald** test statistic  

$$W = n \{R(\hat{\theta})\}^T \left[ \left\{ \frac{\partial}{\partial \theta} R(\hat{\theta}) \right\}^T I^{-1}(\hat{\theta}) \left\{ \frac{\partial}{\partial \theta} R(\hat{\theta}) \right\} \right]^{-1} R(\hat{\theta})$$

$$(H_0 : R(\theta) = 0 \iff H_1 : R(\theta) \neq 0)$$
10. **Rao's** score statistic ( $H_0 : \theta = \theta_0 \iff H_1 : \theta \neq \theta_0$ )  

$$S = \frac{1}{n} \left\{ \frac{\partial l(\theta_0)}{\partial \theta} \right\}^T I^{-1}(\theta_0) \left\{ \frac{\partial l(\theta_0)}{\partial \theta} \right\}.$$
11. For each  $\theta_0 \in \Theta$ , let  $A(\theta_0)$  be the *acceptance* region of a level  $\alpha$  test of  $H_0 : \theta = \theta_0$ . For each  $x \in \mathcal{X}$ , define  $C(x) = \{\theta_0 : x \in A(\theta_0)\}$ . Then random set  $C(X)$  is a  $1 - \alpha$  *confidence set*. Vice versa.
12. A  $1 - \alpha$  **highest posterior density** (HPD) credible set for  $\theta$  is a subset  $C$  of  $\theta$ , of the form  $C_\alpha = \{\theta \in \Theta : \pi(\theta|x) > K(\alpha)\}$  where  $K(\alpha)$  is the largest const s.t.  $P(C_\alpha|x) \geq 1 - \alpha$ .

## 8.6 Standard Errors

1.  $Se(\hat{\theta}) = \{E(\hat{\theta} - \theta)^2\}^{1/2}$
2.  $Rmse(\hat{\theta}) = \left\{ \frac{1}{r} \sum_{j=1}^r (\hat{\theta}_j - \theta)^2 \right\}^{1/2}$ , where  $\{\hat{\theta}_j\}$  are the estimators  $\hat{\theta}$  for  $r$  replications
3.  $Ese(\hat{\theta})$ : estimate  $Se(\hat{\theta})$ , when it depends another parameters, e.g.  
 $Se(\bar{X}) = \sigma/\sqrt{n}$ , then  $Ese(\hat{\theta}) = S/\sqrt{n}$ .
4.  $Ase(\hat{\theta})$
5.  $Ease(\hat{\theta})$ : estimate  $Ase(\hat{\theta})$
6. All of above assume that one knows the PDF of what one is sampling from
7.  $Jse(\hat{\theta})$ : Jackknife  $Se(\hat{\theta})$

$$Jse(\hat{\theta}) = \left[ \frac{n-1}{n} \sum_{i=1}^n \{\hat{\theta}_{(i)} - Je(\hat{\theta})\}^2 \right]^{1/2} \quad \text{where } Je(\hat{\theta}) = \bar{\bar{\theta}}_{(i)}$$

$\bar{\bar{\theta}}_{(i)} = \hat{\theta}$  with the  $i$ -th obs left over.