

## Contents

<b>1</b>	<b>Foundation</b>	<b>1</b>
1.1	Matrix algebra . . . . .	1
1.2	Matrix decomposition . . . . .	4
<b>2</b>	<b>Linear regression</b>	<b>5</b>
2.1	Geometric approach . . . . .	5
2.2	Algebraic approach . . . . .	5
2.2.1	When X is full rank . . . . .	5
2.2.2	When X is not full rank . . . . .	5
2.3	Projection matrix P . . . . .	6
2.3.1	When X is full rank . . . . .	6
2.3.2	When X is not full rank . . . . .	6
2.4	Residual Sums of Squares (RSS) . . . . .	6
2.5	PROPERTIES OF LEAST SQUARES ESTIMATES . . . . .	7
2.5.1	Best Linear Unbiased Estimator, BLUE . . . . .	7
2.5.2	Unbiased estimation of $\sigma^2$ . . . . .	7
<b>3</b>	<b>Linear regression with distribution assumption</b>	<b>8</b>
3.1	MLE . . . . .	8
3.1.1	review on MLE properties . . . . .	9
3.1.2	Orthogonal columns in the regression matrix . . . . .	9
3.2	Estimation with linear constrain . . . . .	9
3.3	Identifiability . . . . .	11
3.4	Estimability . . . . .	11
<b>4</b>	<b>Generalized least square</b>	<b>12</b>
<b>5</b>	<b>Hypothesis Testing</b>	<b>13</b>
5.1	Likelihood ratio test . . . . .	13
5.2	F test . . . . .	14
<b>6</b>	<b>Some comments about ANOVA</b>	<b>15</b>

## 1 Foundation

### 1.1 Matrix algebra

#### Vectors dependency definition

A set of vectors  $D = \{x_1, x_2, \dots, x_r\}$  is called *linearly dependent* if there is a set of scalar  $\alpha_1, \alpha_2, \dots, \alpha_r$  not all zero such that

$$\sum_{i=1}^r \alpha_i x_i = 0$$

Conversely, if  $\sum_{i=1}^r \alpha_i x_i = 0 \Rightarrow \alpha_i = 0, i = 0, 1, \dots, r$ , then  $D = \{x_1, x_2, \dots, x_r\}$  are *linearly independent*.

#### Column space

Suppose A is an  $n \times p$  matrix. Then each column of A is a vector in  $\mathbb{R}^n$ . We can write  $A = (x_1, \dots, x_n)$ , where each  $x_i \in \mathbb{R}^n, i = 1, \dots, p$ . The space spanned by the columns of A is called the *column space* of A, written  $C(A)$ . That is  $S(A) = C(A)$ , where  $S(A)$  is the space spanned by A.

#### Vector differentiation

Define the vector differentiation as follows

$$\frac{d}{d\beta} = \left( \frac{d}{d\beta_i} \right)$$

where  $\beta$  is a  $n \times 1$  vector. Then we have the following properties

$$\frac{d(\beta' a)}{d\beta} = a$$

$$\frac{d(a' \beta)}{d\beta} = a$$

$$\frac{d(\beta' A \beta)}{d\beta} = 2A\beta$$

### Patterned matrices

If all inverses exist

$$\begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}^{-1} = \begin{pmatrix} A_{11}^{-1} + B_{12}B_{22}^{-1}B_{21} & -B_{12}B_{22}^{-1} \\ -B_{22}^{-1}B_{21} & B_{22}^{-1} \end{pmatrix}$$

$$= \begin{pmatrix} C_{11}^{-1} & -C_{11}^{-1}C_{12} \\ -C_{21}C_{11}^{-1} & A_{22}^{-1} + C_{21}C_{11}^{-1}C_{12} \end{pmatrix}$$

where  $B_{22} = A_{22} - A_{21}A_{11}^{-1}A_{12}$ ,  $B_{12} = A_{11}^{-1}A_{12}$ ,  $B_{21} = A_{21}A_{11}^{-1}$ ,  $C_{11} = A_{11} - A_{12}A_{22}^{-1}A_{21}$ ,  $C_{12} = A_{12}A_{22}^{-1}$ , and  $C_{21} = A_{22}^{-1}A_{21}$

### Nonsingular

Suppose  $A$  is an  $n \times n$  square matrix. Then  $A$  is said to be nonsingular if there exists a matrix  $A^{-1}$  such that

$$A^{-1}A = AA^{-1} = I$$

### Null space

The set of all  $x$  such that  $Ax = 0$  is a vector space and is called the null space of  $A$ , written  $N(A)$ .

**Theorem 1.1.** Suppose  $A$  is  $n \times n$ . If  $r(A) = r$  then  $r(N(A)) = n - r$

### Trace

Suppose  $A$  is an  $n \times n$  square matrix with  $ij$ th element  $a_{ij}$ . The trace of  $A$  is defined as

$$\text{tr}(A) = \sum_{i=1}^n a_{ii}$$

Property 1.  $\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B)$

Property 2. The trace is invariant under cyclic

Property 3. Suppose  $A$ ;  $B$ ;  $C$  are  $n \times n$  square matrices. Then

$$\text{tr}(ABC) = \text{tr}(BCA) = \text{tr}(CAB)$$

Property 4. If  $A$  is an  $n \times n$  matrix with eigenvalues  $\lambda_j$ , then  $\text{tr}(A) = \sum_{i=1}^n \lambda_j$  and  $\det(A) = \prod_{i=1}^n \lambda_j$

Property 5. Assume that  $A$  is symmetric, then  $\text{tr}(A^s) = \sum_{i=1}^n \lambda_i^s$

Property 6. Assume that  $A$  is symmetric and nonsingular, then the eigenvalues of  $A^{-1}$  are  $\lambda_i^{-1}$  ( $i = 1, \dots, n$ ) and hence  $\text{tr}(A^{-1}) = \sum_{i=1}^n \lambda_i^{-1}$

### Rank

Property 1.  $\text{rank}(AB) \leq \text{minimum}(\text{rank } A, \text{rank } B)$

Property 2. If  $A$  is any matrix, and  $P$  and  $Q$  are any conformable nonsingular matrices, then  $\text{rank}(PAQ) = \text{rank}(A)$

Property 3.  $\text{rank}(A) = \text{rank}(A') = \text{rank}(AA')$

Property 4. Let  $A$  be any  $m \times n$  matrix such that  $r = \text{rank}(A)$  and  $s = \text{nullity}(A)$ , [the dimension of  $N(A)$ , the null space or kernel of  $A$ , i.e., the dimension of  $\{x : Ax = 0\}$ ]. Then  $r + s = n$

Property 5. If  $C(A)$  is the column space of  $A$ , then  $C(A'A) = C(A')$

Property 6. If  $A$  is symmetric, then  $\text{rank}(A)$  is equal to the number of nonzero eigenvalues

## Eigenvalues and Eigenvectors of a matrix

Suppose  $A$  is an  $n \times n$  square matrix.

$$Ax = \lambda x, \lambda \in \mathbb{R}^1$$

then  $\lambda$  is called an eigenvalue of  $A$  and  $x$  is called an eigenvector. Note that eigenvectors are not unique.

Property 1. Assume that  $A$  is symmetric, then the eigenvalues of  $(I_n + cA)$  are  $1 + c\lambda_i, (i = 1, \dots, n)$

Property 2. Any  $n \times n$  symmetric matrix  $A$  has a set of  $n$  orthonormal eigenvectors, and  $C(A)$  is the space spanned by those eigenvectors corresponding to nonzero eigenvalues

**Theorem 1.2.** *if  $x_1$  and  $x_2$  are eigenvectors with the same eigenvalue, then any nonzero linear combination of  $x_1$  and  $x_2$  is also an eigenvector with the same eigenvalue.*

**Theorem 1.3.**  $\lambda$  is an eigenvalue of  $A$  if and only if  $A - \lambda I$  is singular.

The eigenvalues of a matrix  $A$  are found by finding the solutions of the equation for  $\lambda$

$$\det(A - \lambda I) = 0$$

**Theorem 1.4.** *Suppose  $A$  is  $n \times n$  with eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_n$*

- $\det(A) = \prod_{i=1}^n \lambda_i$
- if  $A$  is singular, then  $\det(A) = 0$
- if  $A$  is nonsingular then  $A^{-1}$  exists and the eigenvalues are given by  $\lambda_1^{-1}, \dots, \lambda_n^{-1}$
- the eigenvalues of  $A'$  are the same as those of  $A$
- $\text{tr}(A) = \sum_{i=1}^n \lambda_i$  and  $\text{tr}(A^{-1}) = \sum_{i=1}^n \lambda_i^{-1}$
- if  $A$  is symmetric then  $\text{tr}(A^r) = \sum_{i=1}^n \lambda_i^r$  for any integer  $r$

## Orthogonal matrix

A square matrix is orthogonal if  $PP' = P'P = I$

**Theorem 1.5.** *An  $n \times n$  matrix  $P$  is orthogonal if and only if the columns of  $P$  form an orthonormal basis for  $\mathbb{R}^n$ , that is the columns of  $P$  are all unit vectors and orthogonal to each other.*

## Positive-semi definite matrices

A symmetric matrix  $A$  is said to be positive-semidefinite (p.s.d.) if and only if  $x'Ax \geq 0 \quad \forall x$

Property 1. The eigenvalues of a p.s.d. matrix are nonnegative

Property 2. If  $A$  is p.s.d., then  $\text{tr}(A) \geq 0$

Property 3.  $A$  is p.s.d. of rank  $r$  if and only if there exists an  $n \times n$  matrix  $R$  of rank  $r$  such that  $A = RR'$

Property 4. If  $A$  is an  $n \times n$  p.s.d. matrix of rank  $r$ , then there exists an  $n \times r$  matrix  $S$  of rank  $r$  such that  $S'AS = I_r$

Property 5.  $A$  is p.s.d. then  $x'Ax = 0 \implies Ax = 0$

## Positive-definite matrices

A symmetric matrix  $A$  is said to be positive-definite (p.d.) if  $x'Ax > 0 \quad \forall x \neq 0$ . We note that a p.d. matrix is also p.s.d.

Property 1. The eigenvalues of a p.d. matrix  $A$  are all positive; thus  $A$  is also nonsingular

Property 2.  $A$  is p.d. if and only if there exists a nonsingular  $R$  such that  $A = RR'$

Property 3. If  $A$  is p.d. then so is  $A^{-1}$

Property 4. If  $A$  is p.d. then  $\text{rank}(CAC') = \text{rank}(C)$

Property 5. If  $A$  is an  $n \times n$  p.d. matrix and  $C$  is  $p \times n$  of rank  $p$ , then  $CAC'$  is p.d.

Property 6. If  $X$  is  $n \times p$  of rank  $p$ , then  $X'X$  is p.d.

Property 7.  $A$  is p.d. if and only if all the leading minor determinants of  $A$  [including  $\det(A)$  itself] are positive.

Property 8. The diagonal elements of a p.d. matrix are all positive

Property 9. (Cholesky decomposition) If  $A$  is p.d., there exists a unique upper tri-angular matrix  $R$  with positive diagonal elements such that  $A = R'R$

Property 10. (Square root of a positive-definite matrix) If  $A$  is p.d., there exists a p.d. square root  $A^{1/2}$  such that  $(A^{1/2})^2 = A$

### Idempotent matrices

A matrix  $P$  is idempotent if  $P^2 = P$ . A symmetric idempotent matrix is called a **projection matrix**

Property 1. If  $P$  is symmetric, then  $P$  is idempotent and of rank  $r$  if and only if it has  $r$  eigenvalues equal to unity and  $n-r$  eigenvalues equal to zero.

Property 2. If  $P$  is a projection matrix then  $\text{tr}(P) = \text{rank}(P)$

Property 3. If  $P$  is idempotent then so is  $I-P$

Property 4. Projection matrices are positive-semidefinite

Property 5. If  $P_i (i = 1, 2)$  is a projection matrix and  $P_1 - P_2$  is p.s.d. then  $P_1 P_2 = P_2 P_1 = P_2$  and  $P_1 - P_2$  is a projection matrix

### Generalized inverse

Consider the linear transformation  $A : \mathbb{R}^p \rightarrow \mathbb{R}^n$ . A generalized inverse of  $A$  is the linear transformation  $A^-$  such that

$$AA^-y = y \text{ for all } y \in C(A)$$

**Equivalently**, suppose  $A$  is an  $n \times p$  matrix, then  $A_{p \times n}^-$  is a generalized inverse of  $A$  if

$$AA^-A = A$$

from the definition we can get

$$(A^-A)(A^-A) = A^-(AA^-A) = A^-A$$

thus  $A^-A$  is idempotent and hence a projection. The generalized inverse is not unique, but always exists.

### Moore-Penrose generalized inverse

Suppose  $A$  is an  $n \times p$  matrix. If the generalized inverse  $A_{p \times n}^-$  satisfies four conditions

- $AA^-A = A$
- $A^-AA^- = A^-$
- $(AA^-)' = AA^-$
- $(A^-A)' = A^-A$

then  $A_{p \times n}^-$  is called the Moore-Penrose generalized inverse. The Moore-Penrose generalized inverse is unique.

## 1.2 Matrix decomposition

### Theorem 1.6. Spectral decomposition

Suppose  $A$  is an  $n \times n$  symmetric matrix. Then there exists an orthogonal matrix  $P$  such that

$$A = P \Lambda P'$$

where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$  is an  $n \times n$  diagonal matrix of the eigenvalues of  $A$  with  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  and  $P$  is the orthogonal matrix of orthonormal eigenvectors corresponding to the eigenvalues of  $A$ .

### Theorem 1.7. Singular value decomposition

Suppose  $A$  is an  $n \times p$  matrix of rank  $r$ , where  $r \leq \min(n, p)$ . There exists orthogonal matrices  $U_{p \times p}$  and  $V_{n \times n}$  such that

$$V'AU = \begin{pmatrix} \Delta & 0 \\ 0 & 0 \end{pmatrix} \rightarrow A = VDU', \text{ where } D = \begin{pmatrix} \Delta & 0 \\ 0 & 0 \end{pmatrix}$$

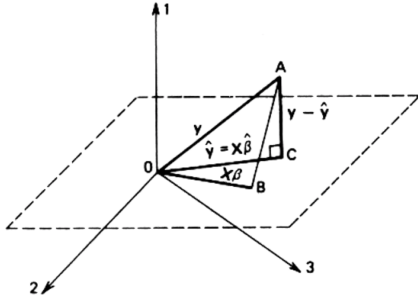
where  $\Delta = \text{diag}(\delta_1, \dots, \delta_r)$  is an  $r \times r$  diagonal matrix with  $\delta_1 \geq \delta_2 \geq \dots \geq \delta_r > 0$ . The  $\delta_i$  are called the singular values of  $A$ .

## 2 Linear regression

We are going to learn about the general linear model in the form of  $E(Y) = X\beta$  and estimation of  $\beta$  using the least squared method and associated distribution theory. The LSE method consists of minimizing  $\sum_i \epsilon_i^2$  with respect to  $\beta$

Let  $\theta = X\beta$ , we minimize  $\epsilon'\epsilon = \|Y - \theta\|^2$  subject to  $\theta \in C(X) = \Omega$ , where  $\Omega$  is the column space of  $X$ .

### 2.1 Geometric approach



From the image we know that  $\|Y - \hat{\theta}\|^2$  minimized when  $\Omega \perp (Y - \hat{\theta})$ , which is  $X'(Y - \hat{\theta}) = 0 \Rightarrow X'\hat{\theta} = X'Y$  thus  $\hat{\theta} = (X')^{-1}X'Y$ . Here  $\hat{\theta}$  is *uniquely determined* being the unique orthogonal projection of  $Y$  onto  $\Omega$ , but  $\beta$  is not necessarily unique.

We have

$$X'\hat{\theta} - X'Y = 0$$

defined as **normal equation**

### 2.2 Algebraic approach

To derive  $\hat{\beta}$  algebraically. Write

$$\begin{aligned} \epsilon'\epsilon &= (Y - X\beta)'(Y - X\beta) \\ &= Y'Y - 2\beta'X'Y + \beta'X'X\beta \end{aligned}$$

By the vector differentiation we have

$$\begin{aligned} -2X'Y + 2X'X\beta &= 0 \\ X'X\beta &= X'Y \end{aligned}$$

To prove that we get the minimal  $\beta$  from this equation, we still need to take the second derivative, from which we get  $2X'X \geq 0$ .

#### 2.2.1 When X is full rank

When the columns of  $X$  are linearly independent i.e.  $X$  is full rank, then there exists a unique vector

$$\hat{\beta} = (X'X)^{-1}X'Y$$

Cause when  $X$  is full rank,  $X'X$  is positive-definite and therefore non-singular

#### 2.2.2 When X is not full rank

When the columns of  $X$  are linearly dependent i.e.  $X$  is not full rank, then the solution is given by

$$\hat{\beta} = (X'X)^- X'Y$$

where  $(X'X)^-$  is any generalized inverse of  $(X'X)$

*Proof:*  $X'X\beta = X'Y$ . Consider a g-inverse  $(X'X)^-$ . We know that  $(X'X)(X'X)^-(X'X) = (X'X)$ . Then we have  $(X'X)(X'X)^-(X'X)\hat{\beta} = (X'X)(X'X)^-X'Y = X'Y$ , by comparing this equation with  $X'X\beta = X'Y$ . We get that  $\hat{\beta} = (X'X)^- X'Y$  is a solution.

## 2.3 Projection matrix P

From the normal equation

$$\hat{\theta} = X\hat{\beta} = X(X'X)^{-}X'Y = PY$$

we define the projection matrix P as

$$P = X(X'X)^{-}X'$$

P is unique and does not depend on the g-inverse used. When the inverse of  $X'$  exists

$$P = X(X'X)^{-1}X'$$

### 2.3.1 When X is full rank

Suppose that  $X$  is  $n \times p$  of rank  $p$ , so that  $P = X(X'X)^{-1}X'$  then following hold.

- (i)  $P$  and  $I_n - P$  are symmetric and idempotent.
- (ii)  $\text{rank}(I_n - P) = \text{tr}(I_n - P) = n - p$ .
- (iii)  $PX = X$

*Proof*

- (i)  $PP = X(X'X)^{-1}X'X(X'X)^{-1}X' = X(X'X)^{-1}X'$
- (ii)  $(I - P)(I - P) = I - P - P + P \cdot P = I - P$
- (iii)  $PX = X(X'X)^{-1}X'X = X$

### 2.3.2 When X is not full rank

If  $X$  has rank  $r < p$ , then the above result still holds but with  $p$  replaced by  $r$

Theorem: Suppose that  $X$  is  $n \times p$  of rank  $r$  so that  $P = X(X'X)^{-}X'$  then following hold.

- (i)  $P$  and  $I_n - P$  are symmetric and idempotent.
- (ii)  $\text{rank}(I_n - P) = \text{tr}(I_n - P) = n - r$
- (iii)  $PX = X$

*Proof*

$X$  has rank  $r$ , let  $X_1$  be the  $n \times r$  matrix with  $r$  linearly independent column then  $C[X_1] = C[X]$ , then

$$P = X_1(X_1'X_1)^{-1}X_1'$$

cause the linear space is the same, then it's easily got that  $P^2 = P$ ,  $(I - P)^2 = (I - P)$ .

Also  $\exists L$  such that  $X = X_1L$

thus  $PX = X_1(X_1'X_1)^{-1}X_1' \cdot X_1L = X_1L = X$

## 2.4 Residual Sums of Squares (RSS)

We denote the fitted values  $X\hat{\beta}$  by  $\hat{Y} = (\hat{Y}_1, \dots, \hat{Y}_n)'$ . The elements of the vector

$$\begin{aligned} Y - \hat{Y} &= Y - X\hat{\beta} \\ &= (I_n - P)Y, \end{aligned}$$

then

$$\begin{aligned} RSS &= [(I - P)Y]'[(I - P)Y] \\ &= Y'(I - P)Y \end{aligned}$$

Another way of doing this is

$$\begin{aligned} e'e &= (Y - X\hat{\beta})'(Y - X\hat{\beta}) \\ &= Y'Y - 2\hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta} - \hat{\beta}'X'X\hat{\beta} + \hat{\beta}'X'X\hat{\beta} \\ &= Y'Y - \hat{\beta}'X'X\hat{\beta} + 2\hat{\beta}'[X'X\hat{\beta} - X'Y] \\ &= Y'Y - \hat{\beta}'X'X\hat{\beta} \end{aligned}$$

thus

$$RSS = Y'Y - \hat{\beta}'X'X\hat{\beta}$$

## 2.5 PROPERTIES OF LEAST SQUARES ESTIMATES

$\hat{\beta}$  is an unbiased estimate of  $\beta$ . That is

$$E(\hat{\beta}) = \beta$$

The variance of the Least Square Estimator of  $\beta$  is given by

$$\text{Var}[\hat{\beta}] = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$$

*Proof*

$$\text{Var}(\hat{\beta}) = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \text{Var}(Y) \mathbf{X} \left[ (\mathbf{X}'\mathbf{X})^{-1} \right]' = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$$

$$E(\hat{\beta}) = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' E[Y] = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbf{X} \beta = \beta$$

Similar result holds for  $\hat{\theta}$

$$E(\hat{\theta}) = P E(Y) = P \underbrace{\mathbf{X}\beta}_{\theta} = \mathbf{X}\beta = \theta$$

$$\text{Var}(\hat{\theta}) = P \text{Var}(Y) P' = \sigma^2 P P' = \sigma^2 P$$

### 2.5.1 Best Linear Unbiased Estimator, BLUE

**THEOREM 3.2:** Let  $\hat{\theta}$  be the least squares estimate of  $\theta = \mathbf{X}\beta$  where  $\theta \in \Omega = C(\mathbf{X})$  and  $\mathbf{X}$  may not have full rank. Then among the class of linear unbiased estimates of  $c'\theta$ ,  $c'\hat{\theta}$  is the unique estimate with minimum variance. We say that  $c'\hat{\theta}$  is the best linear unbiased estimate *BLUE* of  $c'\theta$

*Proof*

$$\hat{\theta} = PY \quad (\text{LSE})$$

$$E(c'\hat{\theta}) = c' E(\hat{\theta}) = c'\theta, \quad \forall \theta \in \Omega = c[x] \Rightarrow \text{unbiasness}$$

Let  $d'Y$  be another estimator which is linear and unbiased then

$$E(d'Y) = d' E(Y) = d' \mathbf{X}\beta = d'\theta \xrightarrow{\text{unbiasness}} (d' - c')\theta = 0 \Rightarrow (d - c) \perp \Omega \Rightarrow P(c - d) = 0 \Rightarrow Pc = Pd$$

then

$$\begin{aligned} \text{Var}(d'Y) - \text{Var}(c'\hat{\theta}) &= d' \text{Var}(Y) d - \text{Var}(c'\hat{\theta}) \leftarrow c'\hat{\theta} = c'PY = (Pc)'Y = (Pd)'Y \\ &= d' (\sigma^2 I) d - \text{Var}[(Pd)'Y] = \sigma^2 d'd - (Pd)'\sigma^2(Pd) \\ &= \sigma^2 d'd - \sigma^2 d'Pd \\ &= \sigma^2 d'(I - P)d \\ &= \sigma^2 d'(I - P)'(I - P)d \\ &= \sigma^2 [(I - P)d]'[(I - P)d] \geq 0 \end{aligned}$$

equality holds when

$$(I - P)d = 0 \Rightarrow d = Pd = Pc$$

If  $\mathbf{X}$  is full rank, then  $a'\hat{\beta}$  is the BLUE of  $a'\beta$  for every vector  $a$ .

### 2.5.2 Unbiased estimation of $\sigma^2$

**THEOREM 3.3**  $E[\mathbf{Y}] = \mathbf{X}\beta$  where  $\mathbf{X}$  is an  $n \times p$  matrix of rank  $r$  ( $r \leq p$ ) and  $\text{Var}[\mathbf{Y}] = \sigma^2 \mathbf{I}_n$  then

$$S^2 = \frac{(\mathbf{Y} - \hat{\theta})'(\mathbf{Y} - \hat{\theta})}{n - r} = \frac{RSS}{n - r}$$

is an unbiased estimate of  $\sigma^2$

*Proof*

$$\begin{aligned}
\text{residual} &= (Y - X\hat{\beta}) = (I - P)Y \\
RSS &= [(I - P)Y]'[(I - P)Y] = Y'(1 - P)Y \\
E(RSS) &= E\{Y'(I - P)Y\} = \text{tr}[(I - P) * \sigma^2 I] + \underbrace{(X\beta)'(1 - P)(x\beta)}_{X-PX=0} \\
&\Rightarrow E\left(\frac{RSS}{n - r}\right) = \sigma^2 \quad \text{unbiased estimator}
\end{aligned}$$

### 3 Linear regression with distribution assumption

Until now the only assumptions we have made about the  $\epsilon_i$  are that  $E(\epsilon) = 0$  and  $\text{Var}(\epsilon) = \sigma^2 I_n$ . If we assume that the  $\epsilon_i$  are also normally distributed, i.e.  $\epsilon \sim N_n(0, \sigma^2 I_n)$  and hence  $Y \sim N_n(X\beta, \sigma^2 I_n)$ . A number of distributional results then follow.

THEOREM 3.5 if  $Y \sim N_n(X\beta, \sigma^2 I_n)$ , where  $X$  is  $n \times p$  of rank  $p$  then

- $\hat{\beta} \sim N_p(\beta, \sigma^2 (X'X)^{-1})$
- $(\hat{\beta} - \beta)'X'X(\hat{\beta} - \beta)/\sigma^2 \sim \chi_p^2$
- $\hat{\beta}$  is independent of  $S^2$
- $RSS/\sigma^2 = (n - p)S^2/\sigma^2 \sim \chi_{n-p}^2$

*Proof*

- $\hat{\beta} = \underbrace{(X'X)^{-1}X'Y}_c$  and also  $Y \sim N_n(X\beta, \sigma^2 I_n)$ .  
then  $cY \sim MVN(cX\beta, c\Sigma c')$  where  $\text{rank}(c) = \text{rank}(X) = \text{rank}(X')$   
 $cX\beta = \beta, c\Sigma c' = \sigma^2 cc' = \sigma^2 (X'X)^{-1}X'X \left[ (X'X)^{-1} \right]' = \sigma^2 (X'X)^{-1}$
- $(\hat{\beta} - \beta)'X'X(\hat{\beta} - \beta)/\sigma^2 \sim \chi_p^2$  since  $\hat{\beta} \sim N_p(\beta, \sigma^2 (X'X)^{-1})$
- $\hat{\beta} = \underbrace{(X'X)^{-1}X'Y}_c$  and  $(n - p)S^2 = Y'(I - P)Y = [(I - P)Y]'[(I - P)Y]$  then  $(X'X)^{-1}X'[I - P] = (X'X)^{-1}X'[I - X(X'X)^{-1}X'] = (X'X)^{-1}X' - (X'X)^{-1}X'X(X'X)^{-1}X' = 0$
- $\frac{RSS}{\sigma^2} = \frac{Y'(1-P)Y}{\sigma^2}$  since  $(1 - P)$  is idempotent with  $\text{rank}(1 - P) = n - r$  based on THEOREM 2.7 we have  $RSS/\sigma^2 = (n - p)S^2/\sigma^2 \sim \chi_{n-p}^2$

#### 3.1 MLE

Assuming full rank of  $X$ , the likelihood is

$$L(\beta, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp\left[-\frac{1}{2\sigma^2}\|Y - X\beta\|^2\right]$$

then the log-likelihood is

$$\ell(\beta, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(Y - X\beta)'(Y - X\beta) = -\frac{n}{2} \log(2\pi\mu) - \frac{1}{2\mu}(Y - X\beta)'(Y - X\beta)$$

where  $\sigma^2 = \mu$

then

$$\frac{\partial \ell}{\partial \beta} = -\frac{1}{2\mu}(-2X'Y + 2X'X\beta) \stackrel{\text{set}}{=} 0 \Rightarrow \hat{\beta}_{\text{mle}} = (X'X)^{-1}X'Y \Rightarrow \text{lse} = \text{mle}$$

and

$$\frac{\partial \ell}{\partial \mu} = \frac{-n}{2\mu} + \frac{1}{2\mu^2}(Y - X\hat{\beta})'(Y - X\hat{\beta}) \stackrel{\text{set}}{=} 0 \Rightarrow \hat{\mu} = \frac{(Y - X\hat{\beta})'(Y - X\hat{\beta})}{n} = \frac{RSS}{n} \neq \hat{\mu}_{\text{lse}} = \underbrace{\frac{(Y - X\hat{\beta})'(Y - X\hat{\beta})}{n - p}}_{\text{unbiased estimator for } \sigma^2}$$



then for the distribution

$$\frac{\partial^2 l}{\partial \beta \partial \beta'} = -\frac{1}{\sigma^2} (X'X) \Rightarrow -E \left[ -\frac{1}{\sigma^2} (X'X) \right] = \frac{X'X}{\sigma^2}$$

$$\frac{\partial^2 l}{\partial \beta \partial \mu} = \frac{1}{2\mu^2} (-2X'Y + 2X'X\beta), \text{ since } \hat{\beta} = (X'X)^{-1} X'Y, \frac{\partial^2 l}{\partial \beta \partial \mu} = \frac{\partial^2 l}{\partial \mu \partial \beta} = 0$$

$$\frac{\partial^2 l}{\partial \mu^2} = \frac{n}{2\mu^2} - \frac{1}{\mu^3} (Y-X\beta)'(Y-X\beta) \Rightarrow -E \left[ -\frac{\partial^2 l}{\partial \mu^2} \right] = \frac{-n}{2\mu^2} + \frac{1}{\mu^3} \underbrace{E[(Y-X\beta)'(Y-X\beta)]}_{(Y-X\beta)'(\sigma^2 I)^{-1}(Y-X\beta) \sim \chi_n^2} = \frac{-n}{2\mu^2} + \frac{n\sigma^2}{\mu^3} = \frac{-n}{2\mu^2} + \frac{n\mu}{\mu^3} = \frac{n}{2\mu^2}$$

then

$$I = \begin{bmatrix} \frac{1}{\mu} X'X & 0 \\ 0 & \frac{n}{2\mu^2} \end{bmatrix} \Rightarrow I^{-1} = \begin{bmatrix} \mu (X'X)^{-1} & 0 \\ 0 & \frac{2\mu^2}{n} \end{bmatrix}$$

then we have

$$\begin{pmatrix} \hat{\beta}_{\text{mle}} \\ \hat{\sigma}_{\text{mle}}^2 \end{pmatrix} \text{ asymptotically normal } \left( \begin{pmatrix} \beta \\ \sigma^2 \end{pmatrix}, \begin{pmatrix} (X'X)^{-1} \sigma^2 & 0 \\ 0 & \frac{2\sigma^4}{n} \end{pmatrix} \right)$$

### 3.1.1 review on MLE properties

**Score:** The partial derivative with respect to  $\theta$  of the natural logarithm of the likelihood function is called the score

$$Z = l' = \frac{\partial}{\partial \theta} \log f(X; \theta)$$

$$E(Z) = 0 \text{ and } Z \xrightarrow{d} N(0, I(\theta_0))$$

under  $\theta_0$

**Fisher information:** The variance of the score is defined to be the Fisher information

$$\mathcal{I}(\theta) = E \left[ \left( \frac{\partial}{\partial \theta} \log f(X; \theta) \right)^2 \mid \theta \right] = -E \left[ \frac{\partial^2}{\partial \theta^2} \log f(X; \theta) \mid \theta \right]$$

Property 1. If  $\hat{\theta}$  is the MLE estimate of  $\theta_0$ , then it has the following property:

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} N\left(0, \frac{1}{I(\theta_0)}\right)$$

### 3.1.2 Orthogonal columns in the regression matrix

Suppose that in the full-rank model  $E[\mathbf{Y}] = \mathbf{X}\beta$  the matrix  $\mathbf{X}$  has a column representation where the columns are all mutually orthogonal

$$\mathbf{X} = \left( \mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(p-1)} \right)$$

then we will have

$$\hat{\beta} = (X'X)^{-1} X'Y = \begin{bmatrix} [x^{(0)'}x^{(0)}]^{-1} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & [x^{(p-1)'}x^{(p-1)}]^{-1} \end{bmatrix} \begin{bmatrix} x^{(0)'}Y \\ \vdots \\ x^{(p-1)'}Y \end{bmatrix} = \begin{bmatrix} \frac{x^{(0)'}Y}{x^{(0)'}x^{(0)}} \\ \vdots \\ \frac{x^{(p-1)'}Y}{x^{(p-1)'}x^{(p-1)}} \end{bmatrix}$$

this implies that under the orthogonal condition, if we only want to estimate certain  $\beta_j$  then we don't need to fit the whole model instead we can only fit  $Y$  on  $x^{(j)}$

## 3.2 Estimation with linear constrain

**The conclusion from this part is applicable under both MLE and LSE cause we only estimate  $\hat{\beta}$**

Let  $Y = X\beta + \epsilon$  where  $X$  is  $n \times p$  of full rank  $p$ . Suppose that we wish to find the minimum of  $\epsilon'\epsilon$  subject to the linear restrictions  $A\beta = c$  where  $A$  is a known  $q \times p$  matrix of rank  $q$  and  $c$  is a known  $q \times 1$  vector then with **Lagrange multiplier** we can get

$$\hat{\beta}_H = \hat{\beta} + (X'X)^{-1} A' \left[ A(X'X)^{-1} A' \right]^{-1} (c - A\hat{\beta})$$

where  $\hat{\beta}$  is the estimation without constrain, i.e.  $\hat{\beta} = (X'X)^{-1} X'Y$ .

*Proof* let  $\lambda = \begin{pmatrix} \lambda_1 \\ \vdots \\ \lambda_q \end{pmatrix}$  then we define Lagrange multiplier as

$$\text{Lagrange multiplier} = (\beta' A' - c') \lambda$$

let

$$S = \min_{\beta} (Y - X\beta)'(Y - X\beta) + (\beta' A' - c') \lambda$$

then

$$\frac{\partial S}{\partial \beta} = -2X'Y + 2X'X\beta + A'\lambda \stackrel{\text{set}}{=} 0 \Rightarrow (X'X)\beta = X'Y - \frac{1}{2}A'\lambda \Rightarrow \hat{\beta}_H = (X'X)^{-1} X'Y - \frac{1}{2}(X'X)^{-1} A'\hat{\lambda}_H$$

suppose  $\hat{\beta}_H$  and  $\hat{\lambda}_H$  are the solution under constrain  $H : A\beta = c$  then we have

$$c = A\hat{\beta}_H = A \underbrace{(X'X)^{-1} X'Y}_{\hat{\beta} \text{ under no constrain}} - \frac{1}{2}A(X'X)^{-1} A'\hat{\lambda}_H = A\hat{\beta} - \frac{1}{2}A(X'X)^{-1} A'\hat{\lambda}_H$$

since  $\text{rank}(A) = q$ ,  $A(X'X)^{-1} A'$  is non-singular. Thus

$$-\frac{1}{2}\hat{\lambda}_H = \left[ A(X'X)^{-1} A' \right]^{-1} (c - A\hat{\beta})$$

therefore

$$\hat{\beta}_H = \hat{\beta} + (X'X)^{-1} A' \left[ A(X'X)^{-1} A' \right]^{-1} (c - A\hat{\beta})$$

we also need to prove the  $\hat{\beta}_H$  minimize  $\epsilon'\epsilon$  under constrain

$$\begin{aligned} (Y - X\beta)'(Y - X\beta) &= (Y - X\hat{\beta} + X\hat{\beta} - X\beta)'(Y - X\hat{\beta} + X\hat{\beta} - X\beta) \\ &= (Y - X\hat{\beta})'(Y - X\hat{\beta}) + (\hat{\beta} - \beta)'X'X(\hat{\beta} - \beta) \end{aligned}$$

where

$$\begin{aligned} (\hat{\beta} - \beta)'X'X(\hat{\beta} - \beta) &= (\hat{\beta} - \hat{\beta}_H + \hat{\beta}_H - \beta)'X'X(\hat{\beta} - \hat{\beta}_H + \hat{\beta}_H - \beta) \\ &= (\hat{\beta} - \hat{\beta}_H)'X'X(\hat{\beta} - \hat{\beta}_H) + (\hat{\beta}_H - \beta)'X'X(\hat{\beta}_H - \beta) + 2(\hat{\beta} - \hat{\beta}_H)'X'X(\hat{\beta}_H - \beta) \\ &= (\hat{\beta} - \hat{\beta}_H)'X'X(\hat{\beta} - \hat{\beta}_H) + (\hat{\beta}_H - \beta)'X'X(\hat{\beta}_H - \beta) \end{aligned}$$

since

$$(\hat{\beta} - \hat{\beta}_H)'X'X(\hat{\beta}_H - \beta) = \left[ \frac{1}{2}(X'X)^{-1} A'\hat{\lambda}_H \right]' X'X (\hat{\beta}_H - \beta) = \frac{1}{2}\hat{\lambda}_H' A (\hat{\beta}_H - \beta) = \frac{1}{2}\hat{\lambda}_H'(c - c) = 0$$

we get

$$(Y - X\beta)'(Y - X\beta) = (Y - X\hat{\beta})'(Y - X\hat{\beta}) + \left[ X(\hat{\beta} - \hat{\beta}_H) \right]' \left[ X(\hat{\beta} - \hat{\beta}_H) \right] + \left[ X(\hat{\beta}_H - \beta) \right]' \left[ X(\hat{\beta}_H - \beta) \right]$$

all three terms are positive, thus minimum achieved when

$$\left[ X(\hat{\beta}_H - \beta) \right] = 0$$

i.e.

$$\beta = \hat{\beta}_H$$

then we get

$$\left\| \mathbf{Y} - \hat{\mathbf{Y}}_H \right\|^2 = \left\| \mathbf{Y} - \hat{\mathbf{Y}} \right\|^2 + \left\| \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_H \right\|^2$$

### 3.3 Identifiability

The conclusion from this part is applicable under both MLE and LSE cause we only estimate  $\hat{\beta}$

A model is identifiable if it is theoretically possible to learn the true values of this model's underlying parameters after obtaining a sample. Mathematically, this is equivalent to saying that different values of the parameters must generate different probability distributions of the observable variables.

Let

$$P = \{P_\theta : \theta \in \Theta\}$$

then

$$P_{\theta_1} = P_{\theta_2} \implies \theta_1 = \theta_2$$

for  $\forall \theta_1, \theta_2 \in \Theta$

We have  $E(Y) = X\beta$ . The design matrix  $X$  is  $n \times p$  and is not full rank. So  $\hat{\beta}$  is not unique. But  $\theta = X\beta$  and  $\hat{\theta}$  is unique. Also  $RSS = Y'(I - P)Y$  is unique. We can always have solution to  $\beta$  based on **g inverse**. The other way is that we impose constraints  $H$  so that models are identifiable.

$$H\beta = 0$$

let

$$G = \begin{bmatrix} X \\ H \end{bmatrix}$$

then we have

$$\hat{\beta} = (G'G)^{-1} X'Y$$

*Proof*

$$G = \begin{pmatrix} X_{n \times p} \\ H_{(p-r) \times p} \end{pmatrix} \Rightarrow G' = (X', H') \Rightarrow G'G = X'X + H'H$$

then  $G^{-1}$  exists cause  $G$  is full rank. We know

$$(X'X)\hat{\beta} = X'Y, \quad H\hat{\beta} = 0$$

thus

$$(G'G - H'H)\hat{\beta} = X'Y \Rightarrow G'G\hat{\beta} = X'Y \Rightarrow \hat{\beta} = (G'G)^{-1} X'Y$$

We have that  $\hat{\beta}$  is unbiased.

*Proof*

$$\begin{aligned} E(\hat{\beta}) &= E \left[ (G'G)^{-1} X'Y \right] = (G'G)^{-1} X'E(Y) = (G'G)^{-1} X'X\beta \\ &= (G'G)^{-1} (G'G - H'H)\beta \\ &= \beta - \underbrace{(G'G)^{-1} H'H\beta}_0 \\ &= \beta \end{aligned}$$

### 3.4 Estimability

Since  $\hat{\beta}$  is not unique,  $\beta$  is not estimable. We consider function of elements of  $\beta$ , i.e.  $a'\beta$

**Definition** The parametric function  $a'\beta$  is said to be estimable if it has a linear unbiased estimate, say  $b'Y$ . This implies that

$$\begin{aligned} E(b'Y) &= b'E(Y) = b'X\beta, \quad \forall \beta \\ b'X\beta &= a'\beta, \quad \forall \beta \\ a' &= b'X \text{ or } a = X'b \end{aligned}$$

thus THEOREM 1

$$\begin{aligned} a'\beta \text{ is estimatable} &\Leftrightarrow a = X'b \\ &\Leftrightarrow a \in C[X'] \text{ or} \end{aligned}$$

**THEOREM 2** if  $a'\beta$  is estimable and  $\hat{\beta}$  is any solution of the normal equation then (1)  $a'\hat{\beta}$  is unique and (2)  $a'\hat{\beta}$  is the BLUE of  $a'\beta$

**THEOREM 3**  $a'\beta$  is estimable if and only if  $a'(X'X)^- X'X = a'$

*Proof*

$\Rightarrow a'\beta$  is estimable then

$$\begin{aligned} \exists c \text{ s.t. } a' &= c'X \\ \Rightarrow a'(X'X)^{-1}X'X &= c'X(X'X)^{-1}X'X = c'PX = c'X = a' \end{aligned}$$

$\Rightarrow$  suppose  $a'(X'X)^{-1}X'X = a'$  then

$$E[a'\hat{\beta}] = E[a'(X'X)^{-1}X'Y] = a'(X'X)^{-1}X'X\beta = a'\beta$$

## 4 Generalized least square

Having developed a least squares theory for the full-rank model  $Y = X\beta + \epsilon$ , where  $E(\epsilon) = 0$  and  $\text{Var}(\epsilon) = \sigma^2I$ , we now consider what modifications are necessary if we allow the  $\epsilon_i$  to be correlated. In particular, we assume the  $\text{Var}[\epsilon] = \sigma^2V$ , where  $V$  is a known  $n \times n$  positive-definite matrix.

**Theorem 4.1.** *Under the above setting, we have*

$$\begin{aligned} \beta^* &= (X'V^{-1}X)^{-1}X'V^{-1}Y \\ \text{Var}[\beta^*] &= \sigma^2(X'V^{-1}X)^{-1} \\ RSS &= (Y - X\beta^*)'V^{-1}(Y - X\beta^*) \end{aligned}$$

*Proof*  $Y = X\beta + \epsilon$ , where  $\epsilon \sim N(0, \sigma^2V)$ . Since  $V$  is a positive-definite

$$\Rightarrow \exists K_{n \times n} \text{ s.t. } V = KK', K^{-1} \text{ exist}$$

then write

$$\underbrace{K^{-1}Y}_Z = \underbrace{K^{-1}X}_B\beta + \underbrace{K^{-1}\epsilon}_\eta \Rightarrow z = B\beta + \eta$$

, where  $E(\eta) = 0, \text{Var}(\eta) = \sigma^2I$  since

$$E(\eta) = E(K^{-1}\epsilon) = 0 \quad ; \quad \text{Var}(\eta) = \text{Var}(K^{-1}\epsilon) = K^{-1}\text{Var}(\epsilon)(K^{-1})' = K^{-1}\sigma^2KK'(K^{-1})' = \sigma^2I$$

then

$$\beta^* = (B'B)^{-1}B'Z = (X'K^{-1}K^{-1}X)^{-1}X'K^{-1}K^{-1}Y = (X'V^{-1}X)^{-1}X'V^{-1}Y$$

then

$$E(\beta^*) = (X'V^{-1}X)^{-1}X'V^{-1}X\beta = \beta \quad \text{unbiased}$$

$$\begin{aligned} \text{Var}(\beta^*) &= (X'V^{-1}X)^{-1}X'V^{-1}\text{Var}(Y)\left[(X'V^{-1}X)^{-1}X'V^{-1}\right]' = (X'V^{-1}X)^{-1}X'V^{-1}\sigma^2V\left[(X'V^{-1}X)^{-1}X'V^{-1}\right]' \\ &= (X'V^{-1}X)^{-1}\sigma^2 \end{aligned}$$

$$E[\beta^*] = (X'V^{-1}X)^{-1}X'V^{-1}X\beta = \beta$$

$$E[(Y - X\beta^*)'V^{-1}(Y - X\beta^*)] = (n - p)\sigma^2$$

$$\begin{aligned} RSS &= (Z - B\beta^*)'(Z - B\beta^*) = (K^{-1}Y - K^{-1}X\beta^*)'(K^{-1}Y - K^{-1}X\beta^*) \\ &= (Y - X\beta^*)'(K^{-1})'K^{-1}(Y - X\beta^*) \\ &= (Y - X\beta^*)'V^{-1}(Y - X\beta^*) \end{aligned}$$

Alternative method for deriving the  $\beta$  is minimizing  $\eta'\eta$  with respect to  $\beta$

$$\begin{aligned} \eta'\eta &= \epsilon'V^{-1}\epsilon \\ &= (Y - X\beta)'V^{-1}(Y - X\beta) \\ &= Y'V^{-1}Y - 2\beta'X'V^{-1}Y + \beta'X'V^{-1}X\beta \end{aligned}$$

$$\frac{\partial \eta'\eta}{\partial \beta} = -2X'V^{-1}Y + 2X'V^{-1}X\beta \stackrel{\text{set}}{=} 0$$

we also get

$$\beta^* = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1} \mathbf{X}'\mathbf{V}^{-1}\mathbf{Y}$$

A special case is that when

$$\mathbf{V} = \text{diag}(w_1^{-1}, w_2^{-1}, \dots, w_n^{-1}) (w_i > 0)$$

we have

$$\beta^* = \frac{\sum_i w_i Y_i x_i}{\sum_i w_i x_i^2}$$

$$(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1} = (x'\mathbf{V}^{-1}x)^{-1} = \left(\sum w_i x_i^2\right)^{-1}$$

**Theorem 4.2.**  $a'\beta^*$  is the best linear unbiased estimate (BLUE) of  $a'\beta$  under the generalized linear model.

*Proof*

$$a'\beta^* = a' \underbrace{(\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1} \mathbf{X}'\mathbf{V}^{-1}}_{b'} \mathbf{Y} = b'\mathbf{Y}$$

is linear in  $\mathbf{Y}$  and also unbiased. Let  $b'_1\mathbf{Y}$  be another unbiased linear estimator of  $a'\beta$

$$a'\beta^* = \alpha' (\mathbf{B}'\mathbf{B})^{-1} \mathbf{B}'\mathbf{Z}$$

$$b'_1\mathbf{Y} = b'_1 \mathbf{K}\mathbf{K}^{-1}\mathbf{Y} = (\mathbf{K}'b_1)' \mathbf{Z}$$

by previous theorem BLUE:

$$\text{Var}(a'\beta^*) \leq \text{Var}((\mathbf{K}'b_1)' \mathbf{Z}) = \text{Var}(b'_1\mathbf{Y})$$

equality holds if and only if

$$(\mathbf{K}'b_1)' = a' (\mathbf{B}'\mathbf{B})^{-1} \mathbf{B}' \Rightarrow b'_1 \mathbf{K} = a' (\mathbf{B}'\mathbf{B})^{-1} \mathbf{B}' \Rightarrow b'_1 = a' (\mathbf{B}'\mathbf{B})^{-1} \mathbf{B}' \mathbf{K}^{-1} = a' (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1} \mathbf{X}'\mathbf{V}^{-1} = b'$$

## 5 Hypothesis Testing

We are interested in the form of  $H_0 = \mathbf{A}\beta = 0$

### 5.1 Likelihood ratio test

$$G : Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{p-1} x_{i,p-1} + \epsilon_i, \text{ full model}$$

$$H : Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{r-1} x_{i,r-1} + \epsilon_i, \text{ reduced model}$$

$$\text{Hypothesis of interest is } H_0 : \beta_r = \beta_{r+1} = \dots = \beta_{p-1} = 0$$

Given the linear model  $G : \mathbf{Y} = \mathbf{X}\beta + \epsilon$ , where  $\mathbf{X}$  is  $n \times p$  of rank  $p$  and  $\epsilon \sim N_n(0, \sigma^2 I_n)$ , we wish to test the hypothesis  $H_0 : \mathbf{A}\beta = 0$ , where  $\mathbf{A}$  is  $q \times p$  of rank  $q$ . the likelihood function for  $G$  is

$$L(\beta, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp\left[-\frac{1}{2\sigma^2} \|\mathbf{Y} - \mathbf{X}\beta\|^2\right]$$

under the MLE

$$\hat{\sigma}^2 = \|\mathbf{Y} - \mathbf{X}\hat{\beta}\|^2/n$$

then the likelihood becomes

$$L(\hat{\beta}, \hat{\sigma}^2) = (2\pi\hat{\sigma}^2)^{-n/2} e^{-n/2}$$

let  $\hat{\beta}_H$  and  $\hat{\sigma}_H^2$  be the estimator of  $\mathbf{Y} = \mathbf{X}\beta + \epsilon$  when  $\mathbf{A}\beta = 0$ , then under the hypothesis

$$L(\hat{\beta}_H, \hat{\sigma}_H^2) = (2\pi\hat{\sigma}_H^2)^{-n/2} e^{-n/2}$$

then the likelihood ratio test of  $H$  is given by

$$\Lambda = \frac{L(\hat{\beta}_H, \hat{\sigma}_H^2)}{L(\hat{\beta}, \hat{\sigma}^2)} = \left(\frac{\hat{\sigma}^2}{\hat{\sigma}_H^2}\right)^{n/2}$$

We have learned that  $-2\log\Lambda$  has a chi-squared distribution.

## 5.2 F test

In summary

$$F = \frac{n-p}{q} (\Lambda^{-2/n} - 1)$$

has an  $F_{q, n-p}$  distribution when  $H_0$  is true. And we also have

•

$$RSS_H - RSS = \|\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_H\|^2 = (\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{c})' \left[ \mathbf{A} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{A}' \right]^{-1} (\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{c})$$

•

$$\begin{aligned} E[RSS_H - RSS] &= \sigma^2 q + (\mathbf{A}\boldsymbol{\beta} - \mathbf{c})' \left[ \mathbf{A} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{A}' \right]^{-1} (\mathbf{A}\boldsymbol{\beta} - \mathbf{c}) \\ &= \sigma^2 q + (RSS_H - RSS)_{Y=E[\mathbf{Y}]} \end{aligned}$$

•

$$F = \frac{\overbrace{(RSS_{\tilde{H}} - RSS)}^{\text{improvement}} / q}{RSS / (n-p)} = \frac{(\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{c})' \left[ \mathbf{A} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{A}' \right]^{-1} (\mathbf{A}\hat{\boldsymbol{\beta}} - \mathbf{c})}{qS^2}$$

is distributed as  $F_{q, n-p}$  (the F-distribution with  $q$  and  $n-p$  degrees of freedom, respectively)

• when  $c = 0$ ,  $F$  can be expressed in the form

$$F = \frac{n-p}{q} \frac{\mathbf{Y}'(\mathbf{P} - \mathbf{P}_H)\mathbf{Y}}{\mathbf{Y}'(\mathbf{I}_n - \mathbf{P})\mathbf{Y}}$$

where  $\mathbf{P}_H$  is symmetric and idempotent and  $\mathbf{P}_H\mathbf{P} = \mathbf{P}\mathbf{P}_H = \mathbf{P}_H$

*Proof*  $y = X\beta + \varepsilon$ , where  $\text{rank}(X) = p$ ,  $\text{rank}(A) = q$  and  $H_0 : A\beta = c$ . Under the constrain,  $H_0 : A\beta = c$  we know

$$\hat{\boldsymbol{\beta}}_H = \hat{\boldsymbol{\beta}} + (X'X)^{-1} A' \left[ A(X'X)^{-1} A' \right]^{-1} (c - A\hat{\boldsymbol{\beta}}) \quad (1)$$

also since we have  $\|\mathbf{Y} - \hat{\mathbf{Y}}_H\|^2 = \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 + \|\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_H\|^2$  we get

$$RSS_H = RSS + (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_H)' X'X (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_H)$$

from 1 we know that

$$\hat{\boldsymbol{\beta}}_H - \hat{\boldsymbol{\beta}} = (X'X)^{-1} A' \left[ A(X'X)^{-1} A' \right]^{-1} (c - A\hat{\boldsymbol{\beta}})$$

and

$$\begin{aligned} RSS_H - RSS &= (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_H)' X'X (\hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}_H) \\ &= (c - A\hat{\boldsymbol{\beta}})' \left[ A(X'X)^{-1} A' \right]^{-1} A(X'X)^{-1} (X'X) (X'X)^{-1} A' \left[ A(X'X)^{-1} A' \right]^{-1} (c - A\hat{\boldsymbol{\beta}}) \\ &= (c - A\hat{\boldsymbol{\beta}})' \left[ A(X'X)^{-1} A' \right]^{-1} (c - A\hat{\boldsymbol{\beta}}) \\ &= (A\hat{\boldsymbol{\beta}} - c)' \left[ A(X'X)^{-1} A' \right]^{-1} (A\hat{\boldsymbol{\beta}} - c) \end{aligned}$$

the first equation is proved

We also know that

$$A\hat{\boldsymbol{\beta}} \sim N(A\boldsymbol{\beta}, \underbrace{A(X'X)^{-1} A'}_B \sigma^2)$$

let  $Z = A\hat{\boldsymbol{\beta}} - c$  then  $\text{Var}(Z) = B\sigma^2$  we have

$$\begin{aligned} E(RSS_H - RSS) &= E \left[ (A\hat{\boldsymbol{\beta}} - c)' \left[ A(X'X)^{-1} A' \right]^{-1} (A\hat{\boldsymbol{\beta}} - c) \right] = E \left[ Z' B^{-1} Z \right] \\ &= \text{tr}(\sigma^2 B^{-1} B) + (A\boldsymbol{\beta} - c)' B^{-1} (A\boldsymbol{\beta} - c) \\ &= \text{tr}(\sigma^2 I_{q \times q}) + (A\boldsymbol{\beta} - c)' B^{-1} (A\boldsymbol{\beta} - c) \\ &= \sigma^2 q + (A\boldsymbol{\beta} - c)' \left[ A(X'X)^{-1} A' \right]^{-1} (A\boldsymbol{\beta} - c) \end{aligned}$$

the second equation is proved

$H_0 : A\beta = c$ , under  $H_0$ ,  $A\hat{\beta} \sim N(c = A\beta, \sigma^2 A(X'X)^{-1}A')$  thus

$$(A\hat{\beta} - c)' \left[ \sigma^2 A(X'X)^{-1}A' \right]^{-1} (A\hat{\beta} - c) = \frac{RSS_H - RSS}{\sigma^2} \sim \chi_q^2$$

we also have

$$\frac{RSS}{\sigma^2} \sim \chi_{n-p}^2$$

and previously we learned that  $RSS_H - RSS \perp RSS$ . Therefore,

$$\frac{RSS_H - RSS / (\sigma^2/q) \sim \chi_q^2}{RSS / (\sigma^2/(n-p)) \sim \chi_{n-p}^2} > \text{independence} \implies \sim F_{q, n-p}$$

the third equation is proved

when  $c = 0$ ,

$$\begin{aligned} \hat{Y}_H &= X\hat{\beta}_H = X \left[ \hat{\beta} + (X'X)^{-1}A' \left[ A(X'X)^{-1}A' \right]^{-1} (c - A\hat{\beta}) \right] \\ &= PY - X \underbrace{(X'X)^{-1}A' \left[ A(X'X)^{-1}A' \right]^{-1} A(X'X)^{-1}X'Y}_{P_1} \\ &= \underbrace{(P - P_1)}_{P_H} Y \end{aligned}$$

here  $P_H$  is symmetric since we know  $P_1$  is symmetric and idempotent. We can also show  $P_1P = PP_1 = P_1$ . Further,  $P_H$  is idempotent

$$\begin{aligned} P_1^2 &= (P - P_1)(P - P_1) = P^2 - P_1P - PP_1 + P_1^2 \\ &= P - 2P_1 + P_1 = P - P_1 = P_H \end{aligned}$$

also  $PP_H = P_H$  and

$$\begin{aligned} RSS_H &= \|Y - X\hat{\beta}_H\|^2 = \|Y - P_H Y\|^2 = Y'(I - P_H)Y \\ RSS &= Y'(I - P)Y \end{aligned}$$

## 6 Some comments about ANOVA

Consider one-way ANOVA

		simple mean
trt 1	$Y_{11}, Y_{12}, \dots, Y_{1, J_1}$	$\bar{y}_1$
trt 2	$Y_{21}, Y_{22}, \dots, Y_{2, J_2}$	$\bar{y}_2$
$\vdots$	$\vdots$	$\vdots$
trt I	$Y_{I1}, Y_{I2}, \dots, Y_{I, J_I}$	$\bar{y}_I$

the model is  $Y_{ij} = \mu_i + \varepsilon_{ij}$ ,  $(i = 1, 2, \dots, I, j = 1, \dots, J_i)$ ,  $\beta = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_I \end{pmatrix}$ , assume  $\varepsilon_{ij} \sim N(0, \sigma^2)$ , then

we can write as

$$\begin{pmatrix} y_{11} \\ \vdots \\ y_{1, J_1} \\ y_{21} \\ \vdots \\ y_{2, J_2} \\ \vdots \\ y_{I1} \\ \vdots \\ y_{I, J_I} \end{pmatrix} = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \\ & & & & \vdots \\ & 1 & 0 & 0 & \cdots & 0 \\ & 0 & 1 & 0 & \cdots & 0 \\ & & & & & \vdots \\ & 0 & 1 & 0 & \cdots & 0 \\ & 0 & 0 & 0 & \cdots & 1 \\ & & & & & \vdots \\ & 0 & 0 & 0 & \cdots & 1 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_2 \\ \vdots \\ \mu_I \end{bmatrix}, \implies X = \begin{bmatrix} \mathbb{1}_{J_1} & \mathbb{0}_{J_1} & \cdots & \mathbb{0}_{J_1} \\ \mathbb{0}_{J_2} & \mathbb{1}_{J_2} & \cdots & \mathbb{0}_{J_2} \\ & & & \vdots \\ \mathbb{0}_{J_I} & \mathbb{0}_{J_I} & \cdots & \mathbb{1}_{J_I} \end{bmatrix}. \text{ To test hypoth-}$$

esis,  $H_0 : \mu_1 = \mu_2 = \dots = \mu_I$  (testable since  $X$  is full rank)  $A\beta = 0 \Leftrightarrow \begin{pmatrix} 1 & -1 & 0 & \dots & 0 \\ 0 & 1 & -1 & \dots & 0 \\ 0 & 0 & 1 & -1 & \dots \\ 0 & & & & \\ & & & \vdots & \\ 0 & 0 & 0 & 1 & -1 \end{pmatrix} \beta = 0.$

Then proceed with usual F-test  $F = \frac{(RSS_H - RSS)/(I-1)}{RSS/(n-I)}$ . **ALTERNATIVELY**, we can write  $S = \sum_i \sum_j (y_{ij} - \mu_i)^2$ , where  $S$  is the error sum of squared  $\epsilon_{ij}^2$ ,  $\frac{\partial S}{\partial \mu_i} = 0 \Rightarrow \sum_j 2(y_{ij} - \mu_i) = 0 \Rightarrow \mu_i = \frac{\sum_j y_{ij}}{J_i} = \bar{y}_i$ ,  $RSS = \sum_i \sum_j (y_{ij} - \hat{\mu}_i)^2 = \sum_i \sum_j (y_{ij} - \bar{y}_i)^2$  and under  $H_0$   $S = \sum_i \sum_j (y_{ij} - \mu)^2$  - error sum of squared,  $\frac{\partial S}{\partial \mu} = 0 \Rightarrow \sum_i \sum_j (y_{ij} - \mu) = 0 \Rightarrow \hat{\mu} = \frac{\sum_i \sum_j y_{ij}}{\sum_i J_i} = \bar{Y}$  ← overall mean. Thus  $RSS_H = \sum_i \sum_j (y_{ij} - \hat{\mu})^2 = \sum_i \sum_j (y_{ij} - \bar{Y})^2$ . Then  $RSS_H - RSS = \sum_i \sum_j (\hat{y}_{ij} - \hat{y}_H)^2 = \sum_i \sum_j (\bar{y}_i - \bar{Y})^2 = \sum_i J_i (\bar{y}_i - \bar{Y})^2$ . Therefore, using the F-test  $F = \frac{\sum_i J_i (\bar{y}_i - \bar{Y})^2 / (I-1)}{\sum_i \sum_j (y_{ij} - \bar{y}_i)^2 / (n-I)}$ ,  $p = I, q = I - 1$